

# **Identification of prostate cancer diagnostic and prognostic biomarkers in urine expression data with a focus on extracellular vesicles**



**Helen Marie Curley**

This thesis is submitted for the degree of  
*Doctor of Philosophy*

University of East Anglia  
School of Biological Sciences  
September 2018

© This copy of the thesis has been supplied on condition that anyone who consults it is understood to recognise that its copyright rests with the author and that use of any information derived there-from must be in accordance with current UK Copyright Law. In addition, any quotation or extract must include full attribution.

I would like to dedicate this thesis to Andy Ripley.

## **Declaration**

No portion of the work referred to in this thesis has been submitted in support of an application for another degree or qualification at this or any other university or other institute of learning.

Helen Marie Curley  
September 2018

## **Acknowledgements**

Mostly I would like to thank my supervision team, Professor Colin Cooper, Dr Daniel Brewer and Dr Jeremy Clark. Dr Daniel Brewer was supportive throughout the whole PhD particularly throughout the challenging task of writing the thesis. Dr Jeremy Clark was very helpful and gave me a great amount of insight into the field of prostate cancer research. I would also like to thank the other members of the lab who made this research possible; Dr Rachel Hurst and Marcelino Yazbeck-Hanna.

A huge thank you goes to the friends and family of Andy Ripley, who funded the studentship. I would also like to thank my family and friends who have been there for me throughout.



## Abstract

Prostate Cancer (PCa) is a major clinical problem worldwide with considerable variability in clinical outcome of patients. PCa diagnostics and prognostics currently lack specific and sensitive clinical biomarkers and treatment is not well individualised. The *PCA3* test, amongst others, highlights the utility of urine in PCa diagnostics and prognostics. Urine contains cells and extracellular vesicles (EV) that originate in the prostate. There are many areas of the PCa clinical process that could be aided with an expression based urine test, including diagnosis, prognosis and response to therapy.

NanoString data (167 transcripts) from 485 EV RNA samples were collected from PCa patients and used to build models that would aid in PCa diagnosis and prognosis i.e. i) PCa (low- (L), intermediate-(I), and high-risk(H)) vs CB (Clinically Benign/No evidence for cancer), ii) high-risk PCa vs CB, and iii) trend in expression across CB>L>I>H. These models were validated in 235 samples, with AUCs of i) 0.851 ii) 0.897 and iii) 0.709, respectively.

The potential of using urine EVs to predict patient response to treatments was also investigated. In a pilot data set a signature of seven transcripts was identified that could optimally predict progression of patients on hormone therapy ( $p = 2.3 \times 10^{-05}$ ; HR = 0.04288). Models were also built using NanoString data from 92 cell RNA samples. Intercomparing expression data from matched cell and EV fractions of urine showed that transcripts significantly higher in the EV samples were associated with the prostate, PCa and cancer in general, supporting them as a viable source of biomarkers in the clinical management of PCa.

In conclusion my analyses have demonstrated the utility of examining urine RNA for the diagnosis and prognosis of PCa. My studies have formed the basis of the production of a Prostate Urine Risk test that is currently under development at UEA.

# Table of Contents

<b>Introduction .....</b>	<b>25</b>
<b>1.1 The Research Gap .....</b>	<b>25</b>
<b>1.2 Biomarkers.....</b>	<b>26</b>
1.2.1 Biomarkers in Cancer.....	26
1.2.2 Problems with the use of current and new biomarkers in clinical diagnostics .....	28
1.2.3 Biomarkers pave the way for stratified treatment of cancers.....	29
<b>1.3 Biomarkers in Prostate Cancer.....</b>	<b>30</b>
1.3.1 Prostate Cancer .....	30
1.3.2 Factors influencing PCa risk of incidence and progression .....	31
1.3.3 Current clinical practice for the diagnosis of PCa .....	33
1.3.4 Current process for the clinical treatment of PCa.....	36
<b>1.4 Known PCa Biomarkers .....</b>	<b>41</b>
1.4.1 Prostate Specific Antigen (PSA).....	41
1.4.2 PCA3.....	46
1.4.3 AMACR.....	49
1.4.4 AR.....	49
1.4.5 SPOP .....	50
1.4.6 TMPRSS2:ERG.....	51
1.4.7 Biomarkers for pre-disposition to PCa .....	53
<b>1.5 Urine and Exosomes.....</b>	<b>55</b>
<b>1.6 Methods in Biomarker Discovery .....</b>	<b>57</b>
1.6.1 Nanostring.....	58
1.6.2 Sequencing.....	60
1.6.3 Polymerase Chain Reaction (PCR) .....	67
1.6.4 Microarrays.....	69
1.6.5 Mass Spectrometry.....	74
1.6.6 Fluorescent In-situ Hybridisation (FISH).....	75
1.6.7 Immunohistochemistry (IHC) .....	76
1.6.8 Enzyme-linked Immunosorbent Assay (ELISA).....	77
1.6.9 Methylation Assays .....	77
1.6.10 Supervised and Unsupervised Analyses.....	78
<b>1.7 Summary and Aim .....</b>	<b>80</b>
1.7.1 Summary.....	80
1.7.2 Aims & Objectives .....	81
<b>Materials and Methods.....</b>	<b>84</b>
<b>2.1 Sample Collection and Processing .....</b>	<b>84</b>
2.1.1 Sample Collection.....	84
2.1.2 Micro-filtration harvesting of Urine Extracellular Vesicles .....	85
2.1.3 Qiagen RNA Extraction.....	86
2.1.4 Nugen Amplification of RNA as cDNA .....	87
2.1.5 NanoString.....	89
2.1.6 PCR (Polymerase Chain Reaction) .....	89
<b>2.2 Clinical Data Collection .....</b>	<b>90</b>

<b>2.3 NanoString Pre-processing</b>	<b>91</b>
2.3.1 Normalisation	91
2.1.1 Normalisation by <i>KLK2</i> and <i>KLK3</i>	91
2.3.2 Normalisation by housekeeping genes	92
2.3.2 NanoStringNorm and NanoString QC Pro	93
2.3.3 Log and Square-root Transformation	94
2.3.4 ComBat	94
<b>2.4 Basic Statistical Tests</b>	<b>95</b>
2.4.1 Mann-Whitney U test (Wilcoxon Rank Sum test)	95
2.4.2 Spearman's Correlation	96
2.4.3 Pearson's Correlation	96
2.4.4 Pearson's Chi-Squared	97
2.4.5 Welch <i>t</i> -test	97
2.4.6 ANOVA – Analysis of Variance	98
2.4.7 Tukey test	98
2.4.8 Kruskal-Wallis	99
2.4.9 Variance and IQR	99
2.4.10 Log rank Test	100
2.4.11 Shapiro-Wilk	100
2.4.12 Brent	101
2.4.13 Benjamini – Hochberg Multiple Testing Correction	101
2.4.14 Receiver Operator Characteristics (ROC)	101
<b>2.5 Clustering</b>	<b>102</b>
2.5.1 Principal Component Analysis (PCA)	102
2.5.2 Hierarchical Clustering	103
2.5.3 <i>k</i> -means Clustering	103
2.5.4 Silhouette, Dunn and Davies-Bouldin Indices	104
2.5.5 Latent Process Decomposition	105
<b>2.6 Model Optimisation</b>	<b>106</b>
2.6.1 GLM: Generalised Linear Model	106
2.6.2 Lasso	107
2.6.3 Random Forest	107
2.6.4 Step for feature selection	109
<b>2.7 Pathway Analysis</b>	<b>109</b>
2.7.1 DAVID	109
<b>2.8 Survival Analysis Tools</b>	<b>109</b>
2.8.1 Kaplan Meier (KM) Curves	110
2.8.2 Cox Proportional Hazard	110
<b>NanoString Data Analysis 1: The Pilot Study</b>	<b>111</b>
<b>3.1 Summary</b>	<b>111</b>
<b>3.2 Introduction</b>	<b>113</b>
3.2.1 The Research Gap	113
3.2.2 The Pilot Study Aims	113
3.2.3 The Probe Targets	114
3.2.4 Risk classification of prostate cancer patients	114
<b>3.3 Data Pre-processing and Technical Variation</b>	<b>115</b>
3.3.1 Normalisation and Background correction	115
3.3.2 NanoStringNorm – Quality of Data and its Normalisation	116
3.3.3 Experimental and Technical Investigation	117
3.3.4 Transforming data to a normal distribution and the Shapiro-Wilk test	119
3.3.5 Housekeeping Probes	123
3.3.6 Removal of Outliers	126
3.3.7 Correlating Gene Probes	127

<b>3.4 Identification of Prostate and Cancer Specific Transcripts and DRE relevance.....</b>	<b>132</b>
3.4.1 Kallikrein identification.....	132
3.4.2 <i>TMPRSS2:ERG</i> Identification.....	133
3.4.3 PCA3 Test.....	138
3.4.4 RNA yield, clinical group and DRE.....	138
<b>3.5 Clustering.....</b>	<b>141</b>
3.5.1 Principal Component Analysis (PCA) and <i>k</i> -means Clustering.....	141
3.5.2 Hierarchical Clustering.....	141
3.5.3 Cluster A.....	143
3.5.4 Cluster B.....	150
3.5.5 Latent Process Decomposition (LPD).....	150
<b>3.6 Significantly varying genes.....</b>	<b>155</b>
<b>3.7 Low-risk, intermediate-risk and high-risk trend.....</b>	<b>160</b>
<b>3.8 Clinical Prediction models.....</b>	<b>161</b>
3.8.1 Logistic regression models using step wise variable selection.....	161
3.8.2 Lasso logistic regression models.....	164
3.8.3 Random Forest.....	166
3.8.4 Random Forest applied to all clinical categories.....	170
3.8.5 Comparing the Models.....	173
<b>3.9 Transcripts that show high-importance.....</b>	<b>174</b>
<b>3.10 Conclusions.....</b>	<b>176</b>
<b>Response to Hormone Therapy.....</b>	<b>179</b>
<b>4.1 Summary.....</b>	<b>179</b>
<b>4.2 Introduction.....</b>	<b>180</b>
4.2.1 The Research Gap.....	180
4.2.2 Aim.....	181
4.2.3 Summary of the HT patient cohort.....	181
<b>4.3 Hormone Therapy Predictor constructed using Nanostring 1 data.....</b>	<b>182</b>
4.3.1 Differentially expressed genes based on initial response, 12 month relapse and 24 month relapse.....	182
4.3.2 Survival analyses of time to progression after HT.....	183
4.3.3 Determining the optimal predictor of progression after HT.....	186
4.3.4 Validation of the seven-transcript signature using NanoString 2 data....	195
<b>4.4 Identifying novel progression related transcripts in the NanoString 2 data.....</b>	<b>197</b>
<b>4.5 Hormone Therapy Predictor using <i>KLK2</i> ratio data on Nanostring 1.....</b>	<b>199</b>
4.5.1 Validation of the final model on <i>KLK2</i> ratio NanoString 2 data.....	206
<b>4.6 Conclusion.....</b>	<b>206</b>
<b>NanoString Data Analysis 2.....</b>	<b>209</b>
<b>5.1 Summary.....</b>	<b>209</b>
<b>5.2 Introduction.....</b>	<b>210</b>
5.2.1 The Research Gap.....	210
5.2.2 Aims.....	211
5.2.3 The Probe Targets.....	211
5.2.4 Classification of prostate cancer patient samples.....	211
<b>5.3 Data Preprocessing and Technical Variation.....</b>	<b>213</b>
5.3.1 Normalisation and Background correction.....	213
5.3.2 Quality of Normalisation.....	215
5.3.3 Experimental and Technical Investigations.....	217
5.3.4 ComBat – Removing collection-centre based significance.....	221
5.3.5 Correlating Gene Probes.....	222
5.3.6 Comparison of NanoString2 with NanoString1.....	224

<b>5.4 Identification of Prostate and Cancer Specific Transcripts and DRE</b>	
<b>relevance</b> .....	<b>224</b>
5.4.1 Kallikrein identification.....	224
5.4.2 <i>TMPRSS2:ERG</i> Identification.....	225
5.4.3 PCA3 Test.....	227
<b>5.5 Clustering</b> .....	<b>230</b>
5.5.1 Principal Component Analysis.....	230
5.5.2 Latent Process Decomposition (LPD).....	230
<b>5.6 Further processing techniques</b> .....	<b>238</b>
<b>5.7 Clinical Prediction models</b> .....	<b>239</b>
5.7.1 Models predicting presence of cancer CB and cancer (L, I, H) samples..	240
5.7.2 Models to distinguish the extreme categories i.e. CB and high-risk cancer	
samples .....	247
5.7.3 Models to predict risk categories using trends in expression.....	251
5.7.4 Models to predict patient type using trends in expression .....	255
5.7.5 Conclusions.....	260
<b>Expression Profile of the Cell Sediment Urine Fraction</b> .....	<b>263</b>
<b>6.1 Summary</b> .....	<b>263</b>
<b>6.2 Introduction</b> .....	<b>264</b>
6.2.1 The Research Gap.....	264
6.2.2 The Aims.....	265
6.2.3 The Data.....	265
<b>6.3 Models predicting presence of cancer CB and cancer (L, I, H) samples using</b>	
<b>cell sediment data</b> .....	<b>267</b>
6.3.2 CB vs High risk cancer patients.....	274
6.3.3 Trend CBN-L-I-H .....	279
<b>6.4 Summary of Predictive Models</b> .....	<b>287</b>
<b>6.5 Comparison of the urine expression profiles of Extracellular vesicle and</b>	
<b>Cell fractions in Prostate Cancer</b> .....	<b>291</b>
6.5.1 Microarray comparison of the global expression profile of Extracellular	
vesicle and Cell fractions.....	291
6.5.2 NanoString comparison of the global expression profile of Extracellular	
vesicle and Cell fractions.....	294
<b>6.6 Discussion</b> .....	<b>297</b>
<b>Discussion</b> .....	<b>301</b>
<b>6.7 Summary</b> .....	<b>301</b>
7.1.1. Chapter 3: NanoString Data Analysis 1: The Pilot Study.....	301
<b>6.8 Chapter 4: NanoString2 Analysis: The Movember GAP1 Project</b> .....	<b>302</b>
6.8.1 Chapter 5: Response to treatment.....	303
6.8.2 Chapter 6: Analysis of Cell Fraction and comparison with EV fraction..	304
<b>6.9 Discussion</b> .....	<b>305</b>
<b>6.10 Future Work</b> .....	<b>308</b>
<b>6.11 Conclusions</b> .....	<b>310</b>
<b>References</b> .....	<b>312</b>
<b>Appendices</b> .....	<b>329</b>
6.12 Binomial Testing between CB and Ca.....	350
6.13 Binomial Testing between CB and Ca (Random Sampling) .....	365
6.14 Binomial Testing between CB and High-risk Ca.....	373
6.15 Multinomial CBLIH Trend .....	386
6.16 Multinomial CBCaM Trend.....	399
6.17 Looking for Housekeepers.....	411

<b>6.18 Cancer Vs CB .....</b>	<b>413</b>
<b>6.19 High Risk Vs CB .....</b>	<b>440</b>
<b>6.20 CB- L-I-H Trend .....</b>	<b>465</b>
<b>6.21 Cell vs EV fraction.....</b>	<b>492</b>

## List of Figures

Figure 0.1 The different zones of the Prostate. 75-85% PCas originate in the peripheral zone, whereas, ~25% originate in the transitional zone. Adapted from Akin O., et al 2006 <sup>2</sup> .....	31
Figure 0.2 The Gleason grading standard drawing. Shows the histopathological pattern of prostate cancers, starting at normal looking prostate cells with normal cellular architecture to fully differentiated PCa cells with no formal cellular architecture. Adopted from Humphrey, P et al., 2004 <sup>39</sup> .....	35
Figure 0.3 The NCI website breaks down the results of PSA screening of 1,000 men between the ages of 55-69. Taken from the National Cancer Institute 2015 <sup>69</sup> .....	45
Figure 0.4: SPOP frequency of substitutions and substrate binding cleft <sup>87</sup> . A) the frequency of substitutions in SPOP across four PCa cohorts from Weill Cornell Medical College (WCMC), University of Michigan (UM), UroPath and University of Washington (UW). B) the substrate-binding cleft of SPOP with the positions of all eight residues that can be possibly mutated. Adopted from Barbieri, C. E. et al. 2012 <sup>87</sup> .....	51
Figure 0.5 <i>ETS</i> family partners for <i>TMPRSS2</i> fusion and their splice variant diversity. Adopted from Clark, J <i>et al.</i> , 2009 <sup>92</sup> .....	52
Figure 0.6 Anatomy of the prostate. Adapted from Drake <i>et al.</i> , 2015 <sup>101</sup> .....	55
Figure 0.7 Tumour cells send signals to distant cells through exosomes. A) Production of exosomes and how they can be sent to recipient cells. B) The different materials that can be found inside exosomes. Adopted from Bátiz, L.F., 2016 <sup>103</sup> .....	56
Figure 0.8 NanoString Ncounter system. A) The set up of the two probes (capture and reporter), one target system. B) The elongation and fixing of probes using a current for imaging. C) Imaging of the uniquely labelled reporter probes. Adapted from Geiss, G et al., 2008 <sup>115</sup> .....	59
Figure 0.9 Sequencing cost per genome from 2001 to 2015. Sudden drops seen in ~2008 and again in 2015. Adapted from National Human Genome Research Institute (NHI) 2016 <sup>119</sup> .....	60
Figure 1.10 Solexa's sequencing methodology using bridge amplification. DNA strands bound by complimentary oligonucleotides to a flow cell arch over to prime the next round of polymerization. This creates clusters of clonal populations via PCR. Fluorophores that can be cleaved between steps show the incorporation of the next dNTP. Adapted from Voelkerdig et al., 2009 <sup>126</sup> .....	63
Figure 0.11 A schematic for oligonucleotide and two-channel microarrays. Both show RNA isolation from the cells of interest, followed by reverse transcriptase labeling to create cDNA from RNA and then hybridisation to array. In two channel arrays, cDNA from the normal cells and the "condition" cells are combined prior to hybridisation. Adapted from Vermeeren et al., 2011 <sup>142</sup> .....	73
Figure 3.1 To ensure there were no batch issues PCA plots were produced of NanoString loading batches and RNA extraction protocol. A) PCA did not identify any clustering associated with NanoString cartridge or scanner used. Along with the Kruskal-Wallis rank sum results also (Cartridge: $p = 0.17$ , Scanner $p = 0.71$ ), it was deemed there was no batch effect produced by NanoString loading. B) PCA does not identify any clustering associated with	

RNA extraction protocol used and the Kruskal-Wallis rank sum test was also insignificant ( $p = 0.16$ ). Thus it was deemed that using no filter, a 45 $\mu$ m filter, and a 45 $\mu$ m filter with a 30-minute wait along side the Qiagen micro RNA RNeasy kit using manufactures' protocols made no difference. ....118

Figure 3.2 A) Amplification cDNA yield shows mild clustering (cDNA yields were grouped: group 1 = 1-2 $\mu$ g, group 2 = 2-3 $\mu$ g, group 3 = 3-4 $\mu$ g, group 4 = 4-5 $\mu$ g, and >5 $\mu$ g in group 5). B) Amplification cDNA yield shows no influence on sample mean expression C) Amplification cDNA yield shows dependence on clinical category.....119

Figure 3.3 Density plots showing the distribution of a) the non-transformed data. B) the log<sub>2</sub> transformed data. C) the square root transformed data. ....120

Figure 3.4 Tukey test comparisons of clinical category for housekeeping probes. When the bar does not cross the mid-point of the x-axis then the comparison is significant. The Tukey test takes each of the five probes (*ALAS1*, *B2M*, *GAPDH*, *HPRT*, and *TBP*) and detects significant expression differences between the six clinical categories. The significant comparisons with S (high PSA/negative Bx samples) is treated cautiously as there were only  $n = 4$  samples within this group. This leaves only one group comparison (CB with Advanced samples in *HPRT*) that showed any significant difference. A good housekeeping probe would be expected to not differ between clinical categories.....124

Figure 3.5 Correlation plots between the housekeeper transcripts: *ALAS1*, *B2M*, *GAPDH*, *HPRT*, and *TBP*. ....125

Figure 3.6 PCA plot of all log<sub>2</sub> normalised data identifies an outlier samples M\_19\_5. ....127

Figure 3.7 Heatmap showing correlation between all NanoString probe data. The colours reflect the  $R$  value of the correlation, where 1 is perfect correlation (represented by yellow) and 0 is uncorrelated (represented by red), with orange in between.....129

Figure 3.8 Correlation plots between data for probes: *M.genitalium* RplA, *M.genitalium* RplB, *HOXC6* and *ERG 5'*.....131

Figure 3.9 Correlation plots for a second group of probes that correlate: *SPINK1*, *SLC12A1* and *UPK2*. All correlate with  $p < 2.26 \times 10^{-16}$  and  $R < 0.6$ . ....132

Figure 3.10 A) Kallikreins are observed at higher expression levels than the blood, kidney and bladder specific markers in the NanoString data. B) Correlation between the two *KLK3* probes is strong.....133

Figure 3.11 A) Density plot for *TMPRSS2:ERG* expression coloured by clinical category. Generally, two peaks are seen suggesting an on/off pattern of expression. B) Density plot for *ERG 3'* expression coloured by clinical category. Again, two bumps are generally seen suggesting an on/off pattern. C) Density plot for *ERG 5'* expression coloured by clinical category. No observable on/off pattern can be seen. D) Box plot showing spread of *TMPRSS2:ERG* expression across clinical categories. Higher expression is observed in cancer than benign. E and F] Box plots showing expression of *ERG 3'* and *ERG5'* respectively across clinical categories. Median expression is Higher in cancer than benign.....135

Figure 3.12 Detection of *TMPRSS2:ERG* by NanoString probes for *TMPRSS2:ERG* (upper) and *ERG3'* (lower) versus PCR detection of *TMPRSS2:ERG* transcripts. T1/E4 indicates a *TMPRSS2* ex1/*ERG*ex4 fusion transcript, 'Other' indicates a different fusion transcript, 'Plus' indicates a mixture of T1/E4 and other



transcripts. The dotted lines are the optimal thresholds (4.93 for *TMPRSS2:ERG* and 7.28 for *ERG3'*) calculated using the Brent method, similarly the solid line is the min curve of a density plot (6.78 and 4.58 for *TMPRSS2:ERG* and *ERG3'* respectively) containing all of the *TMPRSS2:ERG* and *ERG* data..... 137

Figure 3.13 Nanostring PCA3 score calculation (*PCA3* divided by *KLK3* multiplied by 1000 as per the usual *PCA3* score (section 1.4.2). The *PCA3* score is significantly increased in PCa samples compared to those with no clinical evidence of PCa (CB). However, there is no significant difference between the intra-clinical categories of PCa. The uPM3™ assay has shown to be able to detect PCa from non-PCa samples. The NanoString probes have shown to follow this same pattern. .... 138

Figure 3.14 Most of the transcripts detected are from the prostate; DRE boosts transcript level detection and post radical prostatectomy patients offer very low signals in their samples. Samples  $n = 389$ . The advanced (A), high-risk (H), intermediate risk (I), low-risk (L) and no evidence of clinical PCa (CB) samples were taken post-DRE. Pre-DRE and post-RP urine samples have been taken without DRE..... 140

Figure 3.15 The NanoString probe expression distribution of four patient paired samples (pre- and post-DRE). .... 140

Figure 3.16 PCA plots coloured by A) clinical category and B) k-means to identify cluster cut-offs. Cluster A shown by red circle. Cluster B shown by orange circle..... 141

Figure 3.17 Hierarchical clustering provides further evidence for Cluster A and B identification. A) Samples belonging to Cluster A and B are shown in red and yellow boxes, respectively. B) Clusters with significant AU  $p$ -values are encapsulated within a red box. Both Cluster A and B are not included within this main. .... 149

Figure 3.18 A) LPD group bar charts B,C,D,E) Clinical distribution, PSA, Gleason score and age without LPD group, respectively. F,G) PCA plots for  $k$ -means and LPD clustering comparison. .... 152

Figure 3.19 Violin plots showing distribution of each probe across clinical category..... 157

Figure 3.20 Boxplots showing the expression levels in significantly differentially expressed genes between cancer and non-cancer samples found by Mann Whitney U test..... 158

Figure 3.21 Boxplots showing differential expression between aggressive cancer and not aggressive PCa samples for those deemed significant by Mann Whitney U test..... 159

Figure 3.22 Boxplots showing differential expression between advanced cancer and non-cancer samples for those deemed significant via Mann Whitney U testing..... 160

Figure 4.1 Kaplan Meier plots for each of the candidate probes (section 4.3.1). Expression for each probe is grouped into high and low expression using  $K$ -means clustering..... 185

Figure 4.2 Kaplan Meier showing the seven-transcript signature (*AR* exons 4-8 \* *AGR2* \* *DLX1* \* *KLK2* \* *NAALADL2* \* *PPAP2A* / *AMACR*) separated into low and high expression using  $k$ -means. The significance was measured using the cox model (Table 4.17),  $p = 2.3 \times 10^{-05}$ . .... 193

Figure 4.3 Kaplan Meier plot showing the seven transcript signature on NanoString 2 data. The signature was separated using <i>k</i> -means.....	196
Figure 4.4 Individual Kaplan Meier plots for the seven transcripts involved in the signature .....	196
Figure 4.5 Kaplan Meier plots (expression separated via <i>k</i> -means) for the fourteen transcripts identified via Mann Whitney U, Cox and log-rank tests for early HT relapse. ....	203
Figure 4.6 Kaplan Meier plots (expression separated via median) for the fourteen transcripts identified via Mann Whitney U, Cox and log-rank tests for early HT relapse. ....	204
Figure 4.7 Kaplan Meier plots for the three transcripts in the model for KLK2 adjusted hormone therapy data: <i>CAMKK2</i> , <i>PSGR</i> and <i>UPK2</i> .....	206
Figure 5.1 A) Positive control normalised data. B) Positive control normalised and Log <sub>2</sub> transformed data. The data shows a more normal distribution after Log <sub>2</sub> transformation. ....	214
Figure 5.2 Median Vs. IQR of samples on the second NanoString study. Six samples were identified with low medians and/or IQRs, which could be problematic to further analyses. ....	218
Figure 5.3 DNA extracted from EVs was collected from four different centres (Dublin, ICR, UEA, and the USA). DNA extracted from the cell pellet was only collected at UEA (UEA_Cell). PCA plot clearly identifies cell sediment derived samples as a separate cluster from EV derived samples.....	219
Figure 5.4 PCA plot of only EV derived DNA shows evidence of collection-centre of origin based clustering. ....	219
Figure 5.5 Average Log <sub>2</sub> expression across centres shows similar expression levels. ....	220
Figure 5.6 Boxplots showing average expression across cartridge and position on cartridge are similar and are showing no batch effects.....	220
Figure 5.7 PCA plot of EV derived samples, showing a lack of clustering by cDNA yield.....	221
Figure 5.8 Boxplots show the log <sub>2</sub> expression across each sample, coloured by location before and after the application of ComBat.....	222
Figure 5.9 PCA plots of post-ComBat data, shows no clustering by location of origin.....	222
Figure 5.10 Heatmap showing correlation between NanoString Probes in post-ComBat data. <i>R</i> -values between 0 (darker) and 1 (lighter). Correlations with <i>R</i> > 0.8 have been highlighted.....	223
Figure 5.11 KLK2, KLK3 and KLK4 expression is higher than the tissue specific controls for blood, kidney and bladder. The two KLK3 probes are highly correlated (Pearson's correlation: <i>R</i> = 0.84, <i>p</i> < 2.2x10 <sup>-16</sup> ). ....	225
Figure 5.12 Density plots and Boxplots showing the expression changes of <i>TMPRSS2:ERG</i> , two <i>ERG</i> 3' probes, and <i>ERG</i> 5' across clinical categories. ....	227
Figure 5.13 PCA3 Test on post-ComBat NanoString2 data (PCA3 transcript expression/average KLK3 transcript expression * 1000) .....	228
Figure 5.14 PCA plot of post-ComBat data, shows no clustering by clinical category.....	230
Figure 5.15 LPD of post-ComBat data separated into five processes and coloured by clinical category. ....	232

Figure 5.16 Clinical breakdown of each LPD group. Chi-square test: $p$ -value = $7.46 \times 10^{-14}$ , $X^2 = 115$ (ignoring samples from unknown LPD group).....	233
Figure 5.17 LPD of post-ComBat data separated into five processes and coloured by location of origin. ....	234
Figure 5.18 Location of origin breakdown of each LPD group. Chi-square test: $p$ -value = 0.095, $X^2 = 18.7$ (ignoring unknown LPD group samples).....	235
Figure 5.19 ROC curve of top performing model for the prediction of CB vs. Cancer (Low-, Intermediate- and High-risk).....	247
Figure 5.20 ROC curve of the training set for the GAPDH and RPLP2 normalised model built using the 5 significant probes post multiple testing correction....	248
Figure 5.21 Top Significant Probe for CB, low-risk, intermediate-risk and high-risk cancer trend in all four data normalisations.....	251
Figure 5.22 Boxplot showing the expression level of each transcript featured in the CB-L-I-H model built using the multiple tested correction significant probes from the GAPDH and RPLP2 normalised data. This model showed the best test data AUC (0.7008).....	252
Figure 5.23 Top Significant Probe for CB, Cancer, Metastatic trend in all four data normalisations.....	255
Figure 5.24 Boxplot showing the expression level of each transcript featured in the CB-Cancer-Metastatic cancer model built using the significant probes from the GAPDH and RPLP2 normalised data. This model showed the best test data AUC (0.6469).....	257
Figure 6.1 PCA plot of the expression levels for samples taken from the cell sediment and the extracellular vesicle fraction of urine.....	295

## List of Tables

Table 0.1: PCa risk stratification table. Proposed risk categorization from NICE Guidelines 175 <sup>41</sup> .....	37
Table 0.2: Cost of different technologies available for biomarker discovery.....	58
Table 2.1 PCR product sizes for <i>TMPRSS2_exon1</i> (T1) and <i>ERG_ex6</i> (E6) PCR primers (nests are 139bp smaller than primaries).....	90
Table 3.1 Classification and Frequency of the sample types based on NICE criteria <sup>40</sup> . The quantity of samples for each clinical group can be seen as well as the clinical description of the group in terms of Gleason score, PSA level and T stage. ....	114
Table 3.2 Median age and PSA at diagnosis for each clinical category, of samples that are used in subsequent analysis. ....	115
Table 3.3 Sample numbers used in i) 'Cancer', ii) 'Aggression' and iii) 'Extreme' computational analyses.....	115
Table 3.4 Three samples were flagged by NanoStringNorm.....	116
Table 3.5 Three probes were flagged by NanoStringNorm.....	117
Table 3.6 Shapiro-Wilk test results on the first 70 and last 70 samples (all probes) for the non-transformed, log <sub>2</sub> transformed and square root transformed datasets.....	121
Table 3.7 Shapiro-Wilk test results for 10 randomly selected probes for the non-transformed, log <sub>2</sub> transformed and square root transformed datasets....	121
Table 3.8 Shapiro-Wilk test results for the three probes identified by NanoStringNorm as having potential quality issues in the three datasets: non-transformed, log <sub>2</sub> transformed and square root transformed.....	121
Table 3.9 Housekeeper probe Pearson's correlation results, looking for correlating housekeeping probes.....	126
Table 3.10 Four Clusters of probes that correlate with each other (Pearson's correlation).....	130
Table 3.11 Median expression values for kallikreins (prostate specific transcripts) and other tissue markers.....	133
Table 3.12 Testing for Cluster A association to clinical and technical variables. ....	145
Table 3.13 Transcripts significantly associated ( $p < 0.05$ ) with Cluster A via Mann-Whitney U test after using Hochberg multiple testing correction. ....	146
Table 3.14 Gene Ontology (GO) over-represented biological processes in Cluster A's significantly associated transcript list via DAVID. ....	147
Table 3.15 Gene Ontology (GO) over-represented biological processes in all of the transcripts used on NanoString via DAVID. ....	147
Table 3.16 Composition of sample type in each LPD cluster (Cluster B samples and bacterial probes removed). Chi-square test: $p = 2.8 \times 10^{-08}$ , $X = 65.47$ .....	151
Table 3.17 Transcripts significantly different between each LPD group members and those that are not. These are the transcripts that define each LPD cluster.....	153
Table 3.18 Kruskal-Wallis identified 16 probes that significantly differ across clinical category.....	155
Table 3.19 Transcripts differentially expressed between cancer (A, H, I, L) and non-cancer samples (Mann Whitney U test). ....	156

Table 3.20 Transcripts differentially expressed between aggressive cancer and non-aggressive samples (Mann Whitney U test). .....	159
Table 3.21 differentially expressed transcripts when comparing advanced samples with benign (no evidence of cancer) samples (Mann Whitney U test). .....	160
Table 3.22 Spearman's correlation results comparing expression with ordered clinical categories: Low-, Intermediate- and High-risk.....	161
Table 3.23 Transcripts in the Step derived model for comparing cancer to non-cancer.....	162
Table 3.24 Category predictions using the cancer vs. non-cancer step model. ....	162
Table 3.25 Transcripts in the Step derived model for comparing aggressive cancers (A, H) to non-aggressive cancers (I, L).....	163
Table 3.26 Category Predictions when using the aggressive cancer model derived from Step.....	164
Table 3.27 Transcripts in the extreme model (A Vs. CB) derived from Step....	164
Table 3.28 Category predictions using the extreme model derived from Step. ....	164
Table 3.29 Lasso coefficients for three models A) cancer Vs. Non-cancer B) Aggressive cancer Vs. Non-aggressive cancer C) extreme model (A Vs. CB)....	165
Table 3.30 Category predictions using the Lasso cancer Vs. non-cancer model .....	165
Table 3.31 Category predictions using the Lasso aggressive cancer Vs. non-aggressive cancer model.....	165
Table 3.32 Category prediction for the Lasso extreme model (A Vs. CB) .....	166
Table 3.33 Confusion matrix for random forest modelling samples on cancer vs. non-cancer. OOB error estimate of 18.52%.....	166
Table 3.34 Gini values for the random forest model to categorise the samples into cancer and non-cancer. ....	167
Table 3.35 Confusion matrix for random forest modelling samples on aggressive cancer vs. non-aggressive cancer. OOB error estimation of 22.82%. ....	167
Table 3.36 Gini values for the random forest model to categorise the samples into aggressive cancer and non-aggressive cancer.....	168
Table 3.37 Confusion matrix for random forest modelling the samples belonging to the extreme clinical categories (A vs. CB). OOB error estimate of 15.79%.....	169
Table 3.38 Gini values for the random forest model to categorise the extreme samples (A vs. CB).....	170
Table 3.39 Confusion matrix for random forest on all 5 clinical categories. OOB error estimate of 45.5%. ....	172
Table 3.40 Sensitivity, Specificity and PPV for each category after categorising samples into five clinical categories using random forest.....	172
Table 3.41 Confusion matrix for random forest on all 5 categories with random sampling to equalise categorical sample sizes. OOB error estimate of 53.44%. ....	172
Table 3.42 Transcripts identified to distinguish between PCa and non-cancer using Mann Whitney U and Lasso. Random Forest rank is also shown.....	174

Table 3.43 Transcripts repeatedly shown to be differentially expressed between aggressive PCa and non-aggressive PCa. ....	175
Table 3.44 Transcripts commonly found to be differentially expressed by the Mann Whitney U test, GLM and Lasso and Random Forest between advanced and benign samples. ....	175
Table 4.1 Clinical summary of the hormone therapy cohort ( $n=32$ ). ....	181
Table 4.2 Mann-Whitney U test results for comparing samples that respond to HT and those that don't at different time points. ....	182
Table 4.3 Cox results for relapse to hormone therapy ....	183
Table 4.4 Significant probes using log rank test applied to data separated by $k$ -means. ....	183
Table 4.5 The probes included the in the glm after LASSO shrinkage and variable selection, (of the Mann-Whitney U selected probes) with the corresponding beta coefficients ....	186
Table 4.6 The probes included in the cox model after step variable selection (of the Mann-Whitney U selected probes) with the hazard values and $p$ -values. The overall performance of the model to predict progression on HT is $p = 0.00024$ . ....	187
Table 4.7 The importance of each probe in the random forest predictor for HT relapse (of the Mann-Whitney U selected probes). ....	187
Table 4.8 Overall performance of the models (created from the probes originally identified by Mann-Whitney U) tested by cox. ....	188
Table 4.9 The probes included the in the glm after LASSO shrinkage and variable selection, (of the cox selected probes) with the corresponding beta coefficients ....	188
Table 4.10 The probes included in the cox model after step variable selection (of the cox selected probes) with the hazard values and $p$ -values. The overall performance of the model to predict progression on HT is $p = 0.00323$ . ....	189
Table 4.11 The importance of each probe in the random forest predictor for HT relapse (of the cox selected probes). ....	189
Table 4.12 Overall performance of the models (created from the probes originally identified by cox) tested by cox. ....	190
Table 4.13 The probes included the in the glm after LASSO shrinkage and variable selection, (of the log-rank selected probes) with the corresponding beta coefficients. ....	190
Table 4.14 The probes included in the cox model after step variable selection (of the log-rank selected probes) with the hazard values and $p$ -values. The overall performance of the model to predict progression on HT is $p = 0.0012$ . ....	191
Table 4.15 The importance of each probe in the random forest predictor for HT relapse (of the log-rank selected probes). ....	191
Table 4.16 Overall performance of the models (created from the probes originally identified by log rank) tested by cox. ....	192
Table 4.17 Comparing the Cox regression models of various linear combination scores producing from combining gene selection lists. Mann-Whitney U = candidate probes identified as differentially expressed at initial response, 12 month relapse or 24 month relapse; cox = candidate probes identified by step applied to cox regression models; Log rank = candidate probes identified by the log rank test on expression dichotomised into low and high expression. ....	193

Table 4.18 Univariate cox models showing the significance of clinical variables, LPD group and the seven-transcript signature on predicting HT relapse.....	194
Table 4.19 Multivariate cox model for predicting early relapse on HT. ....	195
Table 4.20 Clinical breakdown of the 27 HT patients unique to NanoString 2. ....	195
Table 4.21 Cox regression modelling identified ten probes that were predictors of progression after HT. None were significant after multiple testing correction. ....	197
Table 4.22 Log-rank test identified probes that could significantly predict progression on HT. K-means was used to separate into high- and low-expression of each probe.....	198
Table 4.23 Log-rank test identified probes that could significantly predict progression on HT. Median was used to separate into high- and low-expression of each probe. ....	198
Table 4.24 Optimising models using the four probes common to log-rank and cox tests. The cox model had an overall $p$ -value: $p = 0.0013$ . ....	199
Table 4.25 Mann-Whitney U test identifies probes differentially expressed between those that have relapsed and those that are still responding to HT at different time periods (initial response relapse, within 6 month relapse, with 12 month relapse and within 24 month relapse.....	200
Table 4.26 Cox identified probes that are differentially expressed in NanoString 1 data normalised by <i>KLK2</i> ratio.....	200
Table 4.27 Log rank (using median for separating high and low expression) identified probes that differ between response to HT.....	200
Table 4.28 Lasso selects three transcripts for HT progression prediction in <i>KLK2</i> adjusted data.....	201
Table 4.29 Stepwise regression selects six probes for early HT relapse prediction in <i>KLK2</i> adjusted data. ....	201
Table 4.30 Random forest shows the importance of each transcript in distinguishing early HT relapse in <i>KLK2</i> adjusted data.....	202
Table 4.31 Lasso (with glm) selects three transcripts from the five shown to be differential from Kaplan Meier plots using $k$ -means for separation. An overall Cox model using these three probes proves to be significant ( $p = 0.007$ ).....	205
Table 4.32 Step (with Cox) selects four transcripts from the five shown to be differential from Kaplan Meier plots using $k$ -means for separation. An overall Cox model using these four probes is not significant ( $p = 0.07$ ). ....	205
Table 4.33 Random forest shows the importance of each of the five transcripts identified via Kaplan Meier plots using $k$ -means for separation. The top three important transcripts are identical to the Lasso output. ....	205
Table 5.1 Classification and Frequency of the sample types based on NICE criteria <sup>40</sup> . The quantity of samples for each clinical group are provided as well as the clinical description of the group in terms of Gleason score, PSA level and T stage.....	212
Table 5.2 Sample collection-site breakdown of the EV samples from NanoString2. ....	212
Table 5.3 Median age and PSA of each clinical category within the training and test datasets. ....	212
Table 5.4 Shapiro-Wilk tests show that $\text{Log}_2$ data is not normally distributed. ....	215

Table 5.5 Expression values from different collection-centres of origin compared by Mann Whitney U tests show that all centres are significantly different.....	219
Table 5.6 Pearson's Correlation between the 49 common probes and 131 common samples between NanoString1 and NanoString2.....	224
Table 5.7 Mann Whitney U test of PCA3 Test scores between the different clinical categories. ....	229
Table 5.8 Location of origin breakdown of LPD groups.....	235
Table 5.9 Top ten significantly associated transcripts involved in the separation of samples into LPD groups. The p-value shown is adjusted using Benjamin Hochberg multiple testing correction.....	236
Table 5.10 The different normalisations of the data that the predictive models were built using (separately).....	238
Table 5.11 Clinical predictive models built using the training set and tested using the test set.....	240
Table 5.12 Training model outcomes comparing CB with Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.....	243
Table 5.13 Test model outcomes comparing CB with Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.....	244
Table 5.14 Training model outcomes comparing CB with randomly selected Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.....	245
Table 5.15 Test model outcomes comparing CB with randomly selected Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction. ....	246
Table 5.16 Training model outcomes comparing CB with high-risk Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction. ....	249
Table 5.17 Test model outcomes comparing CB with high -risk Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.....	250
Table 5.18 Training model outcomes comparing CB, low-, intermediate- and high- risk cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction .....	253
Table 5.19 Test model outcomes comparing CB, low-, intermediate- and high-risk cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.....	254
Table 5.20 Training model outcomes comparing CB, Cancer (low-, intermediate- and high- risk) and metastatic (A) cancer samples for the four	



different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.....	258
Table 5.21 Test model outcomes comparing CB, Cancer (low-, intermediate- and high- risk) and metastatic (A) cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.....	259
Table 6.1 Clinical breakdown of cell sediment fraction samples subjected to NanoString (within the second NanoString set). Twelve samples were CB (no evidence of cancer). Thirty raised PSA samples were negative for PCa on biopsy, but other abnormalities were found such as, HGPIN, prostatitis and atypia. Forty-six had localised cancer on TRUS biopsy of which four were D'Amico graded as Low risk, twenty-eight Intermediate risk and fourteen High-risk. Four samples had shown signs of metastasis. ....	266
Table 6.2 Clinical predictive models built using the cell dataset. ....	266
Table 6.3 Top ten transcripts with biggest log <sub>2</sub> fold change in the baseline normalised data. ....	269
Table 6.4 Top ten transcripts with biggest log <sub>2</sub> fold change in the <i>KLK2</i> ratio data.....	269
Table 6.5 Top ten transcripts with biggest log <sub>2</sub> fold change in the HK normalised data. ....	269
Table 6.6 AUC, Sensitivity and Specificity of models to predict CB vs. Cancer (L, I, H) in different data normalisations of cell NanoString data.....	271
Table 6.7 Beta values of individual transcripts within models suggested by Lasso using different input transcripts for the baseline normalised data.....	272
Table 6.8 Beta values of individual transcripts within models suggested by Lasso using different input transcripts for the <i>KLK2</i> ratio data.....	272
Table 6.9 Beta values of individual transcripts within models suggested by Lasso using different input transcripts for the HK normalised data.....	273
Table 6.10 Frequency of transcripts in top 5 for random forests.....	273
Table 6.11 Mean square of residuals for random forest models for predicting CB vs. cancer (L, I and H) samples using different input probes across three different normalisations.....	274
Table 6.12 Top 10 transcripts with biggest log <sub>2</sub> fold change between CB and HR-cancer in the baseline data. ....	275
Table 6.13 Top 10 transcripts with biggest log <sub>2</sub> fold change between CB and HR-cancer in the <i>KLK2</i> ratio data. ....	275
Table 6.14 Top 10 transcripts with biggest log <sub>2</sub> fold change between CB and HR-cancer in the HK normalised data.....	275
Table 6.15 AUC, Sensitivity and Specificity of models to predict CB vs. high-risk cancer (H) in different data normalisations of cell NanoString data. ....	277
Table 6.16 Beta values of individual transcripts within HR cancer and CB models suggested by Lasso using different input transcripts for the baseline normalised data. ....	278
Table 6.17 Beta values of individual transcripts within HR cancer and CB models suggested by Lasso using different input transcripts for the <i>KLK2</i> ratio data.....	278

Table 6.18 Beta values of individual transcripts within HR cancer and CB models suggested by Lasso using different input transcripts for the HK normalised data.....	278
Table 6.19 Frequency of transcripts in top 5 for random forests (CB vs high-risk cancer models).....	279
Table 6.20 Mean Square of residuals error for each random forest model produced using different input probes in three different normalisations.....	279
Table 6.21 Top 15 significant transcripts identified by polr to have trend across CB - L - I -H clinical categories in the baseline normalised cell data.....	281
Table 6.22 Top 15 significant transcripts identified by polr to have trend across CB - L - I -H clinical categories in the <i>KLK2</i> ratio cell data. ....	281
Table 6.23 Top 15 significant transcripts identified by polr to have trend across CB - L - I -H clinical categories in the HK normalised cell data.....	281
Table 6.24 Optimal multinomial models for predicting clinical category (CB, low-risk, intermediate-risk, and high-risk cancer) with different subsets of input transcripts (from preliminary ordered glm and polr tests) in the baseline normalised cell data. ....	283
Table 6.25 Optimal multinomial models for predicting clinical category (CB, low-risk, intermediate-risk, and high-risk cancer) with different subsets of input transcripts (from preliminary ordered glm and polr tests) in <i>KLK2</i> ratio cell data.....	283
Table 6.26 Optimal multinomial models for predicting clinical category (CB, low-risk, intermediate-risk, and high-risk cancer) with different subsets of input transcripts (from preliminary ordered glm and polr tests) in HK normalised cell data. ....	284
Table 6.27 AUC, Sensitivity and Specificity of models to predict trend across clinical categories: CB > L- > I- > H-risk cancer in different data normalisations of cell NanoString data. ....	285
Table 6.28 OOB error rates for random forest models built to predict trend over clinical categories: CB > L > I >H.....	286
Table 6.29 Frequency of transcripts in top 5 for random forests (CB > L > I >H trend models).....	286
Table 6.30 Comparison of AUCs from models using cell and EV data.....	287
Table 6.31 Transcripts identified by all selection models for the different clinical category tests across the different normalisations on the cell NanoString data.....	288
Table 6.32 Transcripts selected for models in EV data. ....	289
Table 6.33 A list of the top 20 microarray detected transcripts out of 98 that were found to be significantly more abundant in extracellular vesicles compared with sediment from the same urine.....	292
Table 6.34 A list of the top 20 microarray detected gene-transcripts out of 116 that were found to be significantly more abundant in the cell sediment compared with extracellular vesicles from the same urine.....	294
Table 6.35 NanoString top twenty transcripts that were up-regulated in extracellular vesicle fractions compared to cell sediment fractions. ....	296
Table 6.36 NanoString top twenty transcripts that were up-regulated in cell sediment fractions compared to extracellular vesicle fractions.....	297

## Abbreviations

A – Advanced	FFPE - Formalin-Fixed Paraffin-Embedded
ADT - Androgen deprivation therapy	FOV - field of view
AED – androstenedione	GAP1 – Global Action Plan 1
AIC – Akaike Information Criteria	GLM - Generalised Linear Model
AMACR - Alpha-methylacyl-CoA racemase	GWAS - Genome-wide association studies
ANOVA - Analysis of variance	H - High risk
AR – Androgen Receptor	HG-PIN – High grade prostatic intraepithelial neoplasia
AUC – Area under the Curve	HR – Hazard Ratio
BCR - Biochemical recurrence	HT – Hormone Therapy
BPH - Benign prostatic hyperplasia	IHC – Immunohistochemistry
Bx - Biopsy	I - Intermediate risk
CB – No evidence for cancer	IQR – Interquartile range
CRPC – Castration Resistant Prostate Cancer	KM - Kaplan Meier
CTCs - circulating tumour cells	L – Low risk
DHT – Dihydrotestosterone	LASSO - Least absolute shrinkage and selection operator
DNA - Deoxyribonucleic acid	LHRH – Luteinizing Hormone Releasing Hormone
DRE – Digital Rectal Examination	lincRNA - long noncoding RNA
ELISA - Enzyme-linked Immunosorbent Assay	LPD – Latent Process Decomposition
EV - Extracellular vesicle	MFISH - multi-fluorochrome assays
FDA - U S Food and Drug Administration	MIP - molecular inversion probes
FISH - Fluorescent In-situ Hybridisation	MMPs - Matrix metalloproteinases

MRI - Magnetic resonance imaging	RF – Random forest
MS - Mass spectrometry	RP - Radical prostatectomy
NICE - National Institute for Health and Care Excellence	RT - Reverse transcriptase
NF - normalisation factor	RTPCR - Reverse transcription polymerase chain reaction
NGS – Next generation sequencing	SD – Standard deviation
NNUH - Norfolk and Norwich University Hospital	SKY - Spectral Karyotyping
OOB - Out-of-bag	SNP - Single nucleotide polymorphism
PCA – Principal Component Analysis	SNV - Single nucleotide variant
PCa – Prostate Cancer	TIMPs - Tissue inhibitors of metalloproteinase
PCA3 – Prostate Cancer Antigen 3	TNBC - Triple negative breast cancer
PCR - Polymerase chain reaction	TNM – Tumour, Lymph Nodes, Metastasis
PPV – Positive Prediction Value	TRUS - Transrectal ultrasound guided
PR- - progesterone receptor negative	UPGMA - Unweighted pair-group method using arithmetic averages
PSA – Prostate Specific Antigen	WHO – World Health Organisation
ROC – Receiver Operator Characteristic	
RMH - Royal Marsden Hospital NHS Foundation Trust	
RNA - Ribonucleic acid	

# 1

## Introduction

### **1.1 The Research Gap**

Prostate cancer (PCa) is the second most common male cancer worldwide and the most common in the UK<sup>1</sup>. The current available biomarkers for PCa lack specificity and/or sensitivity to detect the disease and are unable to distinguish indolent from aggressive disease or predict treatment response. PCa is generally slow-growing, the vast majority requiring no therapeutic intervention at all whilst some of these cancers progress to fatal disease. There is no genetic stratification for treatment unlike many other cancer types, PCa is instead treated with a risk-adjusted patient specific method<sup>2</sup> that aims to improve the control of the cancer whilst reducing risk of complications from treatment. Biopsies

## CHAPTER 1: INTRODUCTION

are commonly performed at diagnosis, but can miss the cancerous area of the prostate and thus lead to a misdiagnosis of “no cancer”. There are limitations to biomarkers capable of predicting positive subsequent biopsy results. There is an urgent clinical need for biomarkers to determine which patients have PCa, which patients have disease that will progress rapidly, and individualise treatment to optimise response.

### **1.2 Biomarkers**

Biomarkers have become widely used in clinical and basic research. The National Institute of Health defines biomarkers as “characteristics that are objectively measured and evaluated as indicators of normal biological processes, pathogenic processes, or pharmacological responses to therapeutic intervention”<sup>3</sup>. Whilst the WHO (World Health Organisation) have a much broader definition that also includes measurable effects of exposure to chemicals or nutrients that allow for risk assessment<sup>4</sup>. Clinically they are used for diagnosis (identification of disease), prognosis (predicting the likely course/outcome of the disease), treatment response stratification and monitoring treatment response in patients. Examples range from blood pressure to more complex genetic screens of tissues, blood, urine and other samples<sup>5</sup>.

#### **1.2.1 Biomarkers in Cancer**

Within the field of cancer management, biomarkers are used for risk assessment, diagnostics, prognostics, treatment stratification and monitoring the effects of treatments. Tumour biomarkers are any measurable molecule that is either produced by the tumour itself or through the host’s response to the tumour that indicates the presence of cancerous processes. Tumour biomarkers can be proteins, glycoproteins, antigens, hormones, receptors, metabolites, and genetic markers; including DNA and RNAs and their epigenetic changes<sup>6</sup>.

Examples of biomarkers in risk assessment include hereditary germ line mutations that increase a person’s risk of developing a certain type of cancer, for example, presence of

## CHAPTER 1: INTRODUCTION

germ line *BRCA1* or *BRCA2* mutations increases the crude life time rate (number of incidences within a population during a specific time period, not considering subgroups within the population) risk of breast cancer in women from 12.5% to 65% and 45%, respectively. Likewise in ovarian cancer, crude rate risk increases from 0.02% to 39% and 11%, respectively<sup>7</sup>. *BRCA* mutation screens are offered to people with known family history of these cancers and positive results can lead to optional preventive measures (e.g. a mastectomy). Other risk assessment biomarkers include p53 but its mutant occurrence in such a range of cancers (50% of all cancers) makes it unusable for screening and diagnosis purposes. As, it could be detected but you would not know where the cancer was or if both alleles were mutated. Also, p53 mutation levels differ between cancer types also, for example, only 3-20% of PCas have a p53 mutation detected at diagnosis<sup>8</sup>.

An example of a biomarker in use in cancer diagnostics is prostate specific antigen (PSA). Serum PSA is currently the first test for PCa diagnosis in the clinic, as elevated levels can suggest the presence of malignancy. PSA, however, does not have great specificity as discussed later: Section 1.4.1

Tissue inhibitors of metalloproteinase (*TIMPs*) are examples of prognostic biomarkers in cancer. *TIMPs* are glycoproteins able to promote proliferation and block apoptosis by inhibiting matrix metalloproteinases (*MMPs*). Increased levels of *TIMPs* have been shown to correlate with poorer survival in many cancers including multiple myeloma, melanoma, breast, lung, colorectal, gastric and head & neck cancers<sup>9</sup>.

Examples of biomarkers in treatment stratification include Human Epidermal Growth Factor Receptor 2 (*HER2*) and Estrogen Receptor  $\alpha$  (*ER $\alpha$* ) in breast cancer. *HER2* and *ER $\alpha$*  receptors may be over-expressed in the breast cancer cells and a simple molecular test (Immunohistochemistry (IHC)) can determine this. This allows treatments to be applied to target the expression profiles of different biomarkers. Herceptin is a drug that specifically targets *HER2*, whereas Tamoxifen is an *ER $\alpha$*  antagonist.

## CHAPTER 1: INTRODUCTION

A biomarker for treatment stratification does not necessarily have to be the drug target. The monoclonal antibody therapies Cetuximab and Panitumumab, which target *EGFR* in colorectal cancers, can only be administered to a cohort of patients who have wild-type *KRAS*. *KRAS* is a signal mediator (extracellular ligand binding and intracellular transduction) between *EGFR* and the nucleus<sup>10</sup>. *KRAS* mutants provide a resistance to these monoclonal antibody therapies. *KRAS* mutations can also occur in response to these treatments and has been shown to be (non-invasively) detectable as early as 10 months prior to radiographic detection of disease progression, allowing administration of *MEK* inhibitors to delay or reverse the resistance<sup>11</sup>.

For treatment resistance monitoring in lung cancer patients, a second *EGFR* mutation, Thr790Met, which can be acquired as a result of treatment or can be pre-existing, provides resistance to *EGFR* tyrosine kinase inhibitors and has been associated with a shorter progression-free survival. Therefore could be used to eliminate people out of the *EGFR* tyrosine kinase inhibitor treatment cohort<sup>10</sup>.

### **1.2.2 Problems with the use of current and new biomarkers in clinical diagnostics**

There is a striking discrepancy between the efforts made to discover cancer biomarkers and the number of biomarkers that actually make it into clinical practice<sup>6</sup>. Major investments have been made to identify and validate novel cancer biomarkers. Using the search terms novel biomarker cancer and new biomarker cancer, a literature search yields 5,358 hits in 2016 alone. Over the past 5 years (2012-2016), 29,775 papers were published using the same search criteria.

However, very few major diagnostic biomarkers have been put into clinical use in the last 25 years<sup>12</sup>. Clinical programs have promised to revolutionize the diagnosis of cancer and the management of its patients. Considerable improvements to how tumours are characterized at a molecular level have shifted treatments towards the use of



## CHAPTER 1: INTRODUCTION

targeted therapies<sup>13</sup>. New PCa tests that have been developed recently include OncotypeDx<sup>14</sup> (section 1.6.3.1), Decipher<sup>15</sup> (section 1.6.4.1) and Prolaris<sup>16</sup> (section 1.6.3.2). However, there is a gap in the number of patients having these tests in clinic to help determine which therapies are suitable for them, and the number of patients that could benefit from these tests. In 2014, the NHS provided 39,298 molecular diagnostic tests for lung, colorectal and melanoma patients in England. Yet the demand was 59,294, leaving 15,929 patients without testing. If this demand was met, it is estimated that 3,552 patients would have been eligible for targeted therapies<sup>17</sup>.

Effective cancer biomarkers need to produce a reliable, reproducible clinically useful assay that is cost effective<sup>6</sup>. The process between biomarker identification to a clinical assay used in practice is lengthy, expensive and convoluted; many researchers working on identifying biomarkers are unaware of clinical practice<sup>6</sup>. Even if a useful tumour biomarker is discovered in the lab there still must be commercialisation incentives in place to develop the assays. Before widespread clinical use the biomarker must be tested in many large datasets and trials carried out by pharmaceutical companies in partnership with academics and also optimised to increase predictive power. Therefore, it can be complicated to determine at which point patenting for the biomarker should be awarded. Regulatory authorities also play a crucial role in validation and quantification of biomarker assays to justify the test to health care providers<sup>13</sup>.

### **1.2.3 Biomarkers pave the way for stratified treatment of cancers**

The current goal of biomarker research is personalized medicine. It aims to provide targeted therapy for individual patients, given their specific clinical, genetic and environmental state. Cancer treatment success is often limited by the heterogeneity among patients; giving patients with genetically different cancers the same treatments can often lead to failure of response with toxic side effects<sup>18</sup>.

Stratified medicine is considered the first step towards personalized medicine. It works by grouping patients via tumour mutations for targeted therapy, using omics

## CHAPTER 1: INTRODUCTION

technologies. It has shown good results within breast cancer patients<sup>19</sup>, amongst other cancers. Breast cancer patients are often stratified between *HER2+* and *HER2-*, *ER+* and *ER-*, *PR+* and *PR-* and triple-negative groups. *HER2+* and *ER+* breast cancers can receive Herceptin and Tamoxifen, respectively: Biological therapies, which are targeted towards those specific receptors.

There is a subset of breast cancers known as triple negative breast cancer (TNBC), where the cancerous cells are *HER2-*, *ER-* and progesterone receptor negative (*PR-*). These cancers have proven to be difficult to treat in the past especially when in their late stages, but promising results have been seen using targeted treatments such as *EGFR* inhibitors and *VEGF* inhibitors that have been previously used for other cancers of different tissues<sup>20,21,22</sup>.

In order for stratified medicine to be effective, biomarker assays that can be routinely applied are needed to accurately stratify patients into treatment cohorts. These assays need to be easily performed with minimal risk to the patient and include immediate or rapid return of the results to ensure early initiation of treatment<sup>23</sup>.

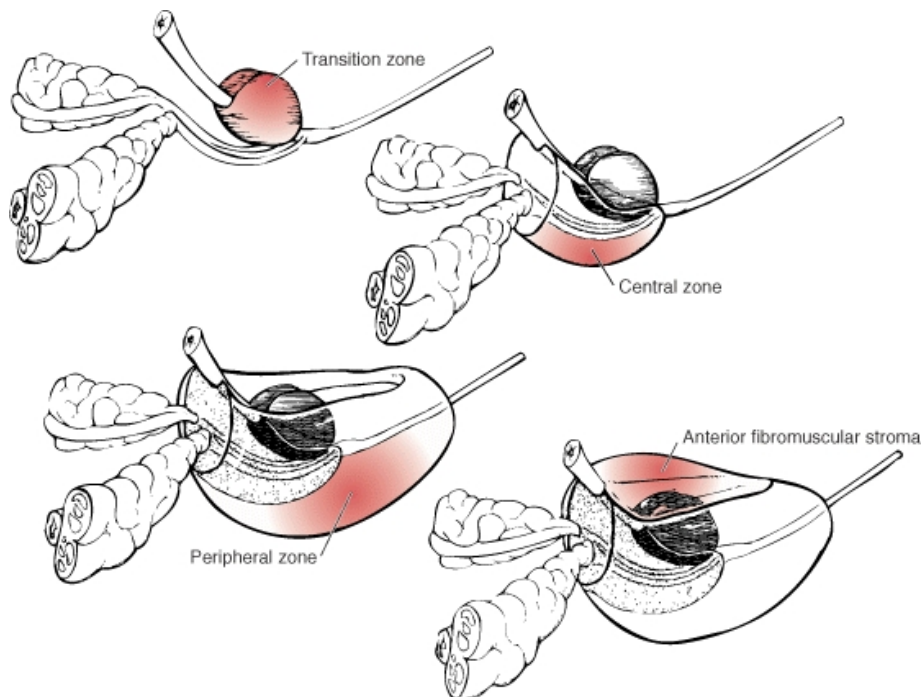
### **1.3 Biomarkers in Prostate Cancer**

#### **1.3.1 Prostate Cancer**

PCa is the second most common male cancer worldwide<sup>24</sup> and the most commonly diagnosed cancer in the UK<sup>25</sup>. In 2010 it accounted for 25% of all cancers diagnosed in men, with 40,975 cases. In 2012 an estimated 307,000 men died from PCa worldwide<sup>24</sup> whilst in the UK 10,721 males died of PCa in 2010; PCa is the second most common cause of cancer death in males. Detected incidence increased by 22% in the last decade and is the fifth fastest increasing cancer in males. Mortality rate, however, has fallen by 11% over the same period<sup>26</sup>; 81.4% of PCa patients survived for five or more years in the UK during 2005 – 2009. Both, the increased incidence rate and the decreased

## CHAPTER 1: INTRODUCTION

mortality rate are associated with the use of the PSA test (section 1.4.1). Changes to classification of PCa deaths and improvements in treatment are also likely to have affected mortality rates. 90% of PCas are acinar adenocarcinomas that originate in the gland cells of the prostate<sup>27</sup>. In approximately 75-85% of PCas<sup>2</sup>, the cancer originates in the peripheral zone rather than the transitional zone (Figure 0.1). The other 10% of PCas fall into different types: signet ring carcinoma, ductal adenocarcinoma, transitional cell (urothelial cancer), squamous cell cancer, carcinoid of the prostate, small cell cancer and sarcoma/sarcomatoid cancer<sup>28</sup>. These will not be considered in the rest of this thesis.



**Figure 0.1** The different zones of the Prostate. 75-85% PCas originate in the peripheral zone, whereas, ~25% originate in the transitional zone. Adapted from Akin O., et al 2006<sup>2</sup>.

### 1.3.2 Factors influencing PCa risk of incidence and progression

There are many factors influencing PCa risk including age, race and family history. PCa is primarily found in older men and risk of developing PCa increases with age. Between 2009 and 2011 36% of UK diagnosed cases of PCa were in men aged above 75, whilst

## CHAPTER 1: INTRODUCTION

only 1% were in those younger than 50<sup>29</sup>. Men aged over 70 also had a statistically significant association with higher clinical stage and Gleason score<sup>30</sup>.

African Americans have a 60% higher risk of developing PCa and mortality is approximately double that of white Americans<sup>31</sup>, and a more aggressive form of the disease can be seen in African Americans<sup>32</sup>. In comparison, native Asian men show a much lower frequency of developing PCa; African American men show a 60-fold higher risk than those in Shanghai, China<sup>31</sup>, although the incidence in Asian populations is increasing<sup>33</sup>. This extraordinary variation of occurrence across the world is boiled down to genetic and environmental factors, which is thought to largely include a Western diet. American-Japanese men have higher incidence rates of PCa than their counterparts in Japan, and this is independent of if they migrated early or late in life, suggesting that life style can accelerate progression of PCa<sup>31</sup>. Asian-American cohorts still hold a lower rate of incidence than white American men<sup>34</sup>.

Evidence of familial risk of PCa has been seen from epidemiological studies, which suggest a two- to three-fold risk increase when there has been a first degree relative diagnosed. Familial clustering patterns have been seen in segregation studies that show high penetrance genetic mutations (including those at the putative susceptibility loci)<sup>31</sup>. PCa aggregates with other familial cancer types (like breast and ovarian). The genes that infer increased susceptibility to these cancers have also shown to increase susceptibility to PCa, e.g. *BRCA1*, *BRCA2*, *CHEK2* and *BRIP1*<sup>1</sup>. Leongamornlert et al., discovered frequent germline mutations in DNA repair genes that were associated with familial PCa as well as a more aggressive phenotype; the cancers were more likely to have nodule involvement, metastasis and be stage 4<sup>1</sup>.

Genome wide association studies (GWAS) identified 76 susceptibility loci associated with PCa risk largely within the European population<sup>35</sup>. These occur commonly but with low penetrance and act multiplicatively to substantially increase risk. GWAS are where genetic variants across whole genomes of different individuals are examined to identify if any variants are associated with specific traits. Investigation of >10 million

## CHAPTER 1: INTRODUCTION

SNPs in a more diverse ancestry population (European, African, Japanese and Latino) in ~43,000 PCa cases and ~43,000 controls revealed 23 novel susceptibility loci<sup>36</sup>. Combining these 23 novel variants with already known variants, we can now explain 33% of the familial risk of PCa in populations of European ancestry. The per allele effects of the 23 variants ranged from 1.06-1.14 and were consistent with log-additive effects of the 23 variants, 15 were exclusive to the European ancestry population, 7 were multi-ethnic, 17 were associated with earlier onset (<55 years compared to >55 years) and 1 was associated with disease severity<sup>37</sup>.

### **1.3.3 Current clinical practice for the diagnosis of PCa**

The current clinical process uses a risk-adjusted patient specific method<sup>2</sup> that aims to improve control of the cancer whilst reducing risk of complications from treatment. The initial step is for a PSA blood test (section 1.4.1.2) to be performed at a GP after a patient has shown symptoms or has other factors increasing their risk such as family history and/or ethnicity. A PSA test is an antibody-based test that measures the concentration of the prostate specific antigen (PSA) in the peripheral blood. A digital rectal examination (DRE) is then performed by a clinician, during which they feel the prostate for any abnormalities. DRE tests have about a 59% overall accuracy<sup>32</sup>. PSA testing is a better predictor of PCa than DRE. In a multicentre trial ( $n = 6$ ) with a total of 6,630 men, 1,167 underwent TRUS biopsies due to PSA>4ng/ml or suspicious DRE result. PSA detected 82% of tumours, whilst DRE only detected 55%, PSA was significantly superior at detecting PCa ( $p = 0.001$ , PPV for PSA: 32% and PPV for DRE: 21%)<sup>38</sup>. However, a DRE is useful because it can often detect cancers missed by the other tests; especially those with normal PSA levels<sup>32</sup>. It can also be used to investigate other abnormal prostatic conditions such as BPH.

If the PSA test (section 1.4.1.2) result is above normal but below 100ng/ml, then a transrectal ultrasound-guided (TRUS) biopsy of the prostate is performed. Using an ultrasound probe, sound waves are reflected off of tissues and organs providing a black

## CHAPTER 1: INTRODUCTION

and white image of the prostate. The probe and biopsy needle gain access to the prostate via the rectum. At the histopathology department, the collected material is examined for cancerous cells and given a Gleason score. In the case of a PSA of greater than 100ng/ml no TRUS is performed, an advanced diagnosis of metastasis is made usually alongside an MRI and/or Bone scan.

The Gleason scoring system (Figure 0.2) is a histopathology score for staging PCa based on how differentiated the cellular structure is in the prostate. This helps evaluate the patient's prognosis, the higher the score the worse the prognosis. It is obtained by combining the scores of the two most common non-normal patterns of histopathology found in the biopsy. The patterns are scored as such: Grade 1 and grade 2 patterns means the tissue is mostly normal; glands are small, well formed and compactly packed, grade 2 has more intracellular space between. Pattern of grade 3 shows recognisable gland units and darker cells that have begun to decrease in size and invade surrounding tissue, the invasion is the most defining feature. Grade 3 is the most common identified, followed by grade 4. A grade 4 pattern has few recognisable gland units with many cells invading surrounding tissue, this can be achieved in many ways resulting in this being the most difficult grade to identify. The fifth grade has no recognisable glands with many cells within the surrounding tissue, there are sheets of cells that lack any nuclear arrangement and a complete loss of gland architecture is observed. In common practice no lower than a 3+3 is seen (giving an overall Gleason grade of 6) and this offers a good prognosis. A Gleason score of 4+3 offers a worse prognosis than that of a 3+4<sup>39</sup>.

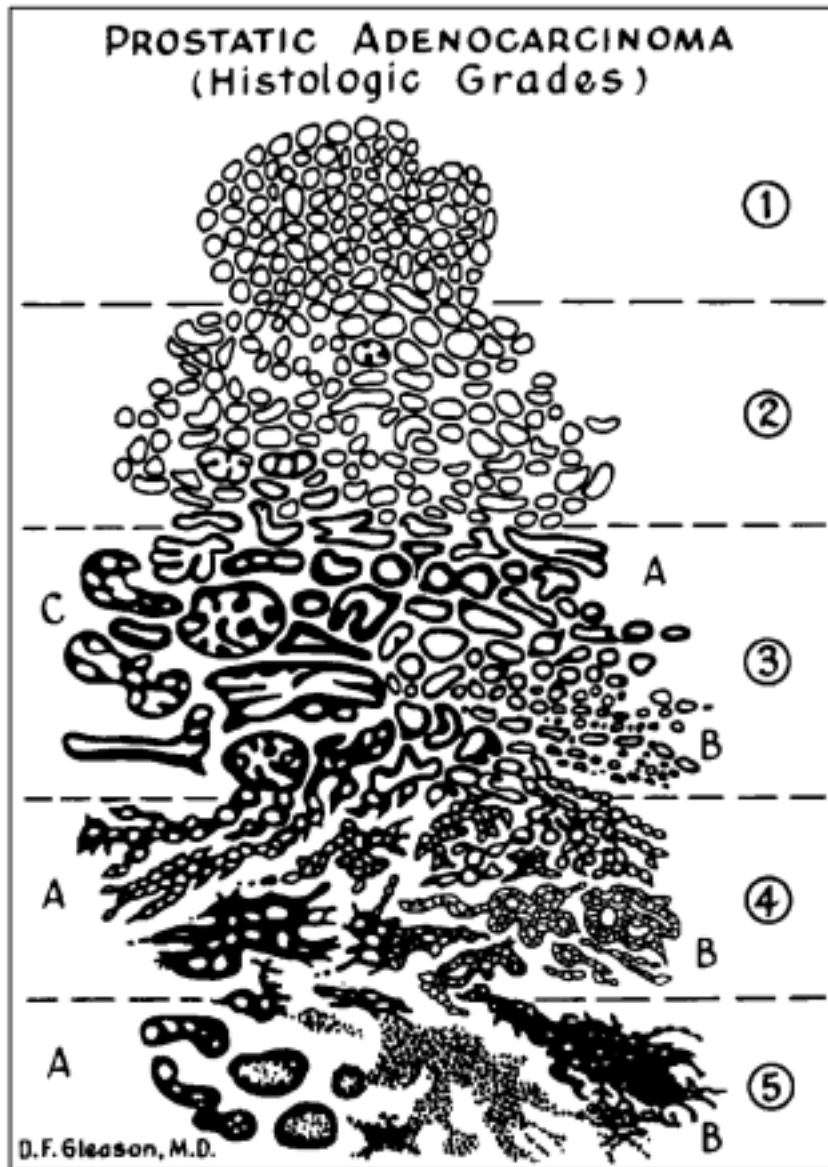


Figure 0.2 The Gleason grading standard drawing. Shows the histopathological pattern of prostate cancers, starting at normal looking prostate cells with normal cellular architecture to fully differentiated PCa cells with no formal cellular architecture. Adopted from Humphrey, P et al., 200439.

Following a negative TRUS biopsy result, if the PSA maintains a high value, a template biopsy can be performed. This differs from the TRUS biopsy as it uses a template or grid over the perineum, which the biopsy needle is entered through to the prostate. However, an ultrasound probe is still used to help guide the needle to the prostate tissue. Generally, more cores are obtained during a template biopsy.

PSA testing lacks specificity and so many men undergo unnecessary TRUS biopsies. TRUS biopsies have risks including serious infection, bleeding, urine retention as well as extra medical costs<sup>40</sup>. Therefore, it is important to identify molecular biomarkers for

## CHAPTER 1: INTRODUCTION

PCa detection that are specific and reliable from a non-invasive source such as blood or urine (section 1.5).

### **1.3.4 Current process for the clinical treatment of PCa**

There are many treatment regimes open to patients with PCa. However, there is a lack of specific and accurate biomarker to stratify patients between the different treatments. For many clinical pathways in PCa there is variability in how long the treatment lasts or whether there is any response at all. For example, resistance to hormone therapies (section 1.3.4.2.1) are inevitable but patients will remain responsive for different lengths of time; from no initial response at all to anywhere between 6 months and 10 years. Another example is how long a patient will last on active surveillance (section 0) before requiring treatment. No biomarkers currently exist that are able to detect which patients will have long term response and which patients response will be short lived and therefore, could benefit from receiving a different/ more aggressive treatment more rapidly. This would offer each patient a more effective treatment first time around.

There are many factors taken into consideration when deciding which treatment is best for a specific PCa patient including general health, age, Gleason score, TNM stage, PSA and whether it is metastatic or not. However, there are not any molecular tests currently available.

#### ***1.3.4.1 Localised Prostate Cancer***

Localised PCa is stratified by their risk of metastasis using the NICE risk categories (Table 0.1) which incorporates PSA level, Gleason score and clinical staging in order to decide treatment style for each patient. Each level of risk is offered a different course of therapy.



## CHAPTER 1: INTRODUCTION

**Table 0.1: PCa risk stratification table. Proposed risk categorization from NICE Guidelines 175<sup>41</sup>**

<i>Level of Risk</i>	<i>PSA</i>	<i>Gleason Score</i>	<i>Clinical Stage</i>
<i>Low risk</i>	<i>&lt;10 ng/ml</i>	<i>and ≤6</i>	<i>and T1-T2a</i>
<i>Intermediate risk</i>	<i>10-20 ng/ml</i>	<i>or 7</i>	<i>or T2b</i>
<i>High Risk<sup>l</sup></i>	<i>≥20ng/ml</i>	<i>or 8-10</i>	<i>or ≥T2c</i>
<i>High-risk localised PCa is also included in the definition of locally advanced PCa.</i>			

### 1.3.4.1.1 Surgery as a treatment for PCa

Radical Prostatectomy (removal of the whole prostate gland) is a treatment considered for men with T1 or T2 PCa (localised to the prostate gland without spread). Side effects can include urinary incontinence, impotence and loss of fertility. Transurethral resection of the prostate (TURP) is considered for men with benign prostate growth (BPH) and for advanced cancer to alleviate symptoms; the inner area of the prostate (that surrounding the urethra) is removed.

### 1.3.4.1.2 Radiotherapy

Radiation therapy is the provided course of treatment for low-grade, localised PCAs (with similar cure rates as those who receive radical prostatectomy). It can also be provided alongside hormone therapy for cancers that have spread out of the gland to nearby tissues, for recurring tumours (post-surgery), and also to advanced patients to reduce tumour size (offering some relief from symptoms). Side effects can include urinary incontinence, impotence cystitis and radiation proctitis.

### 1.3.4.1.3 Biochemical Recurrence

Men treated with either radiotherapy or radical prostatectomy (RP) can develop biochemical recurrence (BCR), which is characterised by a state of elevating PSA level post treatment and indicates growing tumour or metastases<sup>42</sup>. Within 10 years, of those patients treated with radiotherapy ~30-50% and ~20-40% of patients post RP will develop BCR<sup>43</sup>. Increase in PSA does not necessarily mean imminent death or threat

## CHAPTER 1: INTRODUCTION

and can often be treated with hormone therapy. There has been much research in treatment options for these patients, which includes when to administer hormone treatment as well as non-hormonal alternatives including targeted agents and immunotherapies<sup>43</sup> due to the morbidities associated with hormone treatment.

UHRF1 expression in tissue samples has been identified as a potential biomarker for predicting BCR post RP. UHRF1 expression negatively correlates with mean months of BCR-free survival ( $p < 0.001$ ). However, UHRF1 expression was less significant than pre-operation PSA levels and Gleason score<sup>44</sup>. Other studies have identified biomarkers that are linked to BCR; Prx6 (an oxidative stress marker) expression is associated with shortened biochemical recurrence free survival and overall survival in 240 post RP patients ( $p = 0.02$  and  $p = 0.033$ , respectively)<sup>45</sup>. PTEN deletion has been associated with an increased risk of BCR ( $p < 0.01$ , HR: 3.58)<sup>46</sup>. Metallothionein-2A (MT-2A), E-cadherin, and cyclin-E were investigated for BCR association by microarray immunostaining. Positive MT-2A and cyclin E expression along with negative E-cadherin expression showed a decrease in biochemical recurrence-free survival ( $p = 0.009$  (HR = 2.15, 95% CI = 1.14 - 3.08),  $p = 0.037$  (HR = 1.45, 95% CI = 1.02 – 1.92), and  $p = 0.047$  (HR = 1.31, 95% CI = 1.03 – 2.21), respectively)<sup>47</sup>. In a multivariate analysis all three were deemed to independently predict BCR<sup>47</sup>. Still, the promise of these biomarkers have not been translated into use in the clinic.

Other clinical features such as tumour volume and percentage tumour volume have also been reported to predict BCR post RP in a meta-analysis of multicentre data ( $p = 0.03$ , HR: 1.04 and  $p = 0.02$ , HR: 1.01, respectively)<sup>48</sup>.

### Active Surveillance, Watchful Waiting and PSA monitoring

To attempt to reduce the number of over-treated patients, programs like active surveillance, watchful waiting and PSA monitoring have been implemented.

A high proportion of PCa are localised and non-aggressive and are unlikely to cause any problems in the patient at all, whereas others progress into more problematic cancers that require more aggressive treatments. Active surveillance is offered to

## CHAPTER 1: INTRODUCTION

patients with low-risk localised PCa whom are suitable for radical prostatectomy or radiotherapy as treatment<sup>40</sup>. They monitor the patients looking for indications that their less aggressive cancers are becoming more aggressive problematic cancers. Active surveillance is a close monitoring of the patients and usually involves frequent tests, such as PSA blood tests, DREs, ultrasounds and biopsies.

Watchful waiting is offered to asymptomatic PCa patients for whom there is no curative treatment options or intent. Watchful waiting however is implemented with more aggressive cancers, where treatment would cause problems due to the patients' age or general health. These patients are monitored for disease progression (a rapidly rising PSA or bone pain). Compared to active surveillance, less frequent tests and more reliance on patient symptoms for indication of change are implemented in watchful waiting.

PSA monitoring exists to identify patients who have continual raised PSA in the "grey zone" (PSA between 4 and 10ng/ml) rather than just an intermittently raised PSA on one test. Patients can receive multiple PSA tests to monitor them prior to biopsy. This can help to eliminate the number of unnecessary biopsies if there is a continual PSA>4ng/ml then it is more likely to be due to PCa and thus these patients require biopsies.

### ***1.3.4.2 Metastatic Prostate Cancer***

Metastatic PCa is detected in 21% of men at their time of diagnosis<sup>29</sup>. It is usually identified by a PSA>100ng/ml<sup>49</sup> and/or a positive bone scan. Those with metastasis are primarily prescribed hormone therapy agents that block androgen signaling<sup>50</sup>.

#### **1.3.4.2.1 Hormone Therapy**

Androgens are male hormones, which include testosterone and dihydrotestosterone (DHT), and aid in the signalling for prostate cell growth. Androgen deprivation therapy (ADT) lowers levels of these androgens and/or prevents them from reaching the

## CHAPTER 1: INTRODUCTION

prostate cells, resulting in shrinking and slower growth of the cancer. ADT is not a cure but can prolong life.

Luteinizing hormone-releasing hormone (LHRH) analogs and antagonists reduce levels of testosterone released from the testicles by blocking the feedback loop to the hypothalamus. Anti-androgens bind androgen receptors, preventing cell growth signalling, though these are usually added to LHRH treatments when patients begin to become resistant. However, it is a controversial question of when anti-androgens should be added to LHRH treatment to gain full androgen blockage, it is thought in some cases initial hormone therapy should include both LHRH treatments and anti-androgens<sup>51</sup>.

Patients receiving ADT develop resistance leading to castration resistant PCa (CRPC), with a median survival of 1-2 years<sup>52</sup>. It is likely that the high level of heterogeneity within the prostate tumour contributes to this resistance<sup>53</sup>. CRPC develops when cells become hypersensitive to the residual levels of testosterone that are left during chemical castration. Castration does not remove all testosterone; the maintenance of intratumoral androgens is due (at least partly) to the intratumoral or intracrine biosynthesis of steroid hormones (adrenal androgens) or potentially de novo steroidogenesis, from cholesterol or progesterone precursors within the tumour<sup>54</sup>. Hypersensitivity to these residual levels of testosterone are believed to be due to androgen receptor (*AR*)- mutations that alter ligand binding, alterations in *AR* co-regulators or *AR* over-expression (considered to be the main driver of CRPC progression)<sup>54</sup>. *AR* over-expression has also shown to convert anti-androgen treatments (like bicalutamide, flutamide and enzalutamide) from *AR* antagonists to *AR* agonists<sup>55,56</sup>.

Abiraterone was the first drug in clinical practice to target the production of androgens by the tumour. It irreversibly and selectively inhibits CYP17A activity. CYP17A is a critical enzyme; it facilitates the hydroxylase and lyase activity required in the production of adrenal androgens, DHEA and androstenedione (AED), from cholesterol<sup>54</sup>. Although, abiraterone has had impressive responses in clinical trials, not all men respond and resistance occurs (seen by a rising PSA), the mechanisms for

## CHAPTER 1: INTRODUCTION

which are currently unknown. Abiraterone's place in the treatment of PCa is so far undetermined and many clinical trials are in place to investigate this.

### 1.3.4.2.2 Castrate Resistant Prostate Cancer (CRPC)

Once PCa becomes castrate resistant, there are other treatment options available such as chemotherapy and vaccine therapy.

Chemotherapies are usually given to PCa patients who have metastasis but are not, or no longer responding to hormone therapies. It is generally not given to patients with early PCa, although studies are currently investigating its use following surgery. The first chemotherapy agent of choice for PCa is Docetaxel (administered alongside the steroid prednisone) and if this doesn't work or stops working, Cabazitaxel is often a second drug choice<sup>57</sup>. Chemotherapy is used again with the focus on increasing life expectancy and/or quality of life for PCa patients (by slowing the growth of the cancer) but is considered unlikely to result in a cure to the disease.

## **1.4 Known PCa Biomarkers**

Biomarkers in PCa fall into different categories: Biomarkers to predict the presence of PCa (screening and diagnosis), biomarkers to stratify patients (into those requiring active surveillance and those requiring more radical treatments), biomarkers for identifying those whom can be treated with biological targeted therapies and predisposition biomarkers for those who are more likely to develop PCa in their lifetime.

### **1.4.1 Prostate Specific Antigen (PSA)**

PSA, a kallikrein like serine protease (coded for by the gene *KLK3*), is a molecular biomarker currently and routinely used for the diagnosis of PCa, as well as roles in prognosis and treatment response. In normal prostate glands, PSA is highly compartmentalized and found at levels 1 million times fold higher within the prostate compared to that in blood serum. However, in prostatic disease it is thought that this

## CHAPTER 1: INTRODUCTION

compartmentalization is disrupted resulting in increased levels of escaped circulating PSA<sup>58</sup>.

PSA is prostate specific but not cancer specific; elevated serum PSA can be the result of benign prostatic hyperplasia (BPH), chronic inflammation, and infection. Normal and diseased prostatic epithelial cells produce PSA, therefore, weakening its specificity as a cancer biomarker.

Research into men with a PSA less than 4ng/ml has shown that there are many men with low PSA (0.6-1ng/ml) that have PCa (10.1%) and even high-grade (Gleason 7+) PCa (10%)<sup>59</sup>. Evidence suggested there was no PSA threshold for which a man can be assured he has no risk of PCa, but men with <0.5ng/ml PSA do have a decreased risk of developing PCa. Risk of PCa in men with PSA <0.5ng/ml was 6.6%, this increased to 26.9% in men with PSA 3.1-4ng/ml<sup>59</sup>. PSA level effect on the risk of PCa was significant,  $p < 0.001$  (odds ratio 1.66 per unit increase in PSA, 95% CI 1.50 – 1.85)<sup>59</sup>.

PSA levels are affected by both age and race; when deciding on a reference range for diagnosis and deciding which men will undergo TRUS biopsies, it is important to consider these factors. A study on 77,700 men showed that not only does the PSA level rise but also that the range increases with increased age (ages 40-49; mean PSA: 0.83, SD: 0.79, ages 50-59; mean PSA: 1.23, SD: 1.33, ages 60-69; mean PSA: 1.83, SD: 1.94, and ages 70-79; mean PSA: 2.31, SD: 2.35). The differences between the age groups and their variances were significant,  $p < 0.0001$  and  $p = 0.0001$ , respectively<sup>60</sup>. Significant differences in PSA levels were observed between different races also; pairwise differences were seen between white and black people, white and Latino people, black and Asian people, and Asian and Latino people ( $p < 0.0001$ ). Black people have the highest mean PSA values in each age cohort<sup>60</sup>.

### ***1.4.1.1 PSA - Screening***

Due to the lack of specificity that PSA holds, using it for screening purposes has led to over diagnosis and over-treatment as well as downgrading and down staging at

## CHAPTER 1: INTRODUCTION

diagnosis and fewer PCa related deaths<sup>61</sup>. A cohort of men diagnosed with PCa, have a form of the cancer that grows so slowly that it is unlikely to pose a threat to the patient. Treating of these cancers is known as over-treatment. PSA's lack of specificity for PCa means it is not recommended for a screening biomarker due to the over-detection and overtreatment costs it would lead to<sup>61</sup>.

The National Cancer Institute estimate that screening 1,000 men between 55 and 69 every 1-4 years would result in 100-120 men getting a false positive diagnosis (Figure 0.3). False positive diagnoses lead to anxiety and stress for the patient and his family, as well as extra medical costs in further diagnostic procedures. Procedures include TRUS biopsies, which also add further risk to patients; serious infections are not uncommon. Of the 1,000 men screened, and the 110 patients to receive a true positive result, it is estimated that only 1 man would be saved due to screening, compared to the 4-5 men that would die without screening<sup>62</sup>.

### ***1.4.1.2 PSA – Diagnosis***

Similarly to its use in screening, PSA makes a weak diagnostic biomarker due to its lack of specificity to cancer. However, it is the current first diagnostic test for PCa. The sensitivity and predictive value of PSA as a biomarker for PCa decreases greatly for patients in the “grey zone”. PSA levels in the approximate range 2-10ng/ml is known as the “grey zone” as it is difficult to distinguish which elevations are due to cancer and which are associated to other factors including age and BMI, or due to conditions such as benign prostatic hyperplasia (BPH). Investigations into the PSA grey zone generally use cutoffs between 2/4 to 10ng/ml to define it. For every 5 patients, whose PSA level resides between 2.5-10ng/ml, 4 will have a negative biopsy result, and the predictive value of PSA in the grey zone drops from >90% to <25%<sup>63</sup>.

As an individual variable, PSA is a much better PCa predictor than a digital rectal examination (DRE) or transrectal ultrasound<sup>61</sup>, but its modest diagnostic accuracy has led to other PSA forms being investigated.

## CHAPTER 1: INTRODUCTION

### ***1.4.1.3 Free PSA and Pro-PSA***

To improve abilities in distinguishing BPH from PCa in patients who fall in the “grey zone”, investigations into the percent free PSA (or ratio of free to complex PSA) and its most significant cut-off for biopsy, and different isoforms of pro-PSA were performed.

Antibodies were developed that could distinguish between and measure the amounts of tPSA and fPSA, a higher ratio of fPSA:tPSA correlates with a lower risk of PCa. This comparison allowed a small yet significant improvement in the ability of PSA to distinguish PCa from BPH (and other benign diseases that raise PSA levels)<sup>64</sup>.

A study of 773 men with PSA levels between 4-10ng/ml with confirmed histological diagnosis (379 with PCa and 394 with BPH) resulted with a suggested 25% free PSA cut-off. The 25% free PSA cut-off was able to detect 95% of patients with PCa and was also able to avoid 20% of unnecessary biopsies<sup>64</sup>.

PSA is secreted as the inactive enzyme pro-PSA, this can be cleaved at different locations resulting in the mature/active form of PSA. Some remain uncleaved and pro-PSA can have many isoforms. The [-2]proPSA consistently correlates with PCa<sup>65</sup>; it is observed in greater abundance if the prostate is neoplastic (25-95% of free PSA compared to only 6-19% in men without PCa<sup>66</sup>).

Guazzoni et al., showed that the use of %[-2]pro-PSA alone was better at discriminating between PCa and BPH (in patients with PSA ranges 2—10ng/ml) compared to that of total PSA and percentage free PSA, with AUCs of 75.7%, 53%, and 58%, respectively<sup>67</sup>. Using an artificial neural network, Stephan et al., showed that the combination of %[-2]proPSA, %free-PSA, total PSA and age (but not prostate volume) offered highest accuracy (AUC 0.85). It was also shown that %[-2]proPSA was better at discriminating between T2 and T3 tumours as well as Gleason <7 and Gleason >7 cancers<sup>68</sup>.



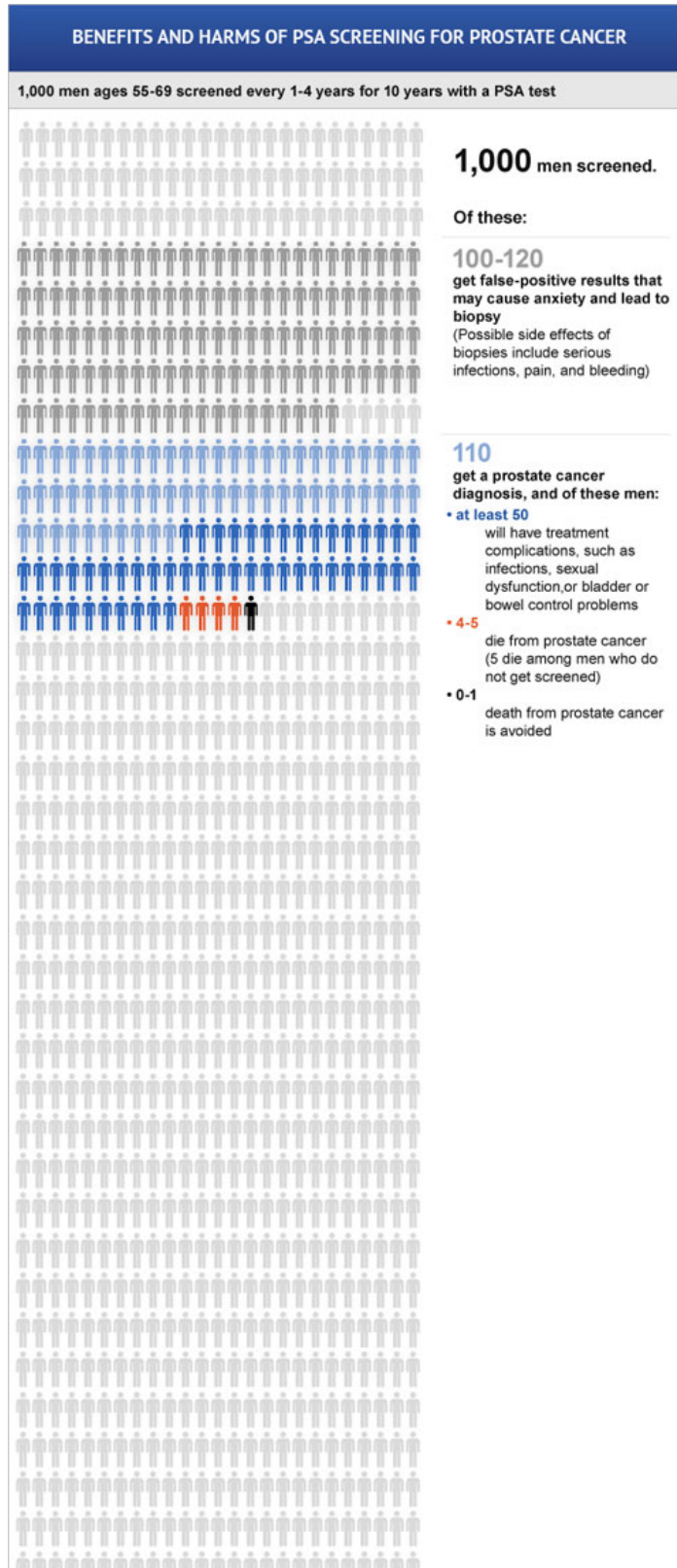


Figure 0.3 The NCI website breaks down the results of PSA screening of 1,000 men between the ages of 55-69. Taken from the National Cancer Institute 2015<sup>69</sup>.

## CHAPTER 1: INTRODUCTION

### ***1.4.1.4 PSA – Treatment***

PSA is commonly used within treatment plans available for PCa, it is a good indicator of progression and drug resistance. PSA levels are routinely and frequently checked in PCa patients; looking for progression in AS patients (section 0), resistance in HT patients (section 1.3.4.2.1) and BCR in radiotherapy or post-radical prostatectomy patients (section 1.3.4.1.3).

PSA is one of the key factors in determining treatment options for patients. A PSA above 100 is indication of metastasis and so hormone therapy is usually provided. PSA also is involved in determining treatment of lower grade localised PCa (Table 0.1).

Investigations into [-2]pro-PSA combined with percentage fPSA identified a correlation for the need of more radical treatment rather than active surveillance<sup>70</sup>. Also, other proPSA isoforms ([-5] and [-7]pro-PSA) correlate with a need for more radical treatments in active surveillance patients, when found in the tissue surrounding the tumour in biopsies.

### ***1.4.1.5 Concluding PSA***

PSA is not a specific PCa biomarker, yet it is the first clinical diagnostic test given to patients and is also a determining factor in treatment options and changes. PSA remains a very useful biomarker in following patients with PCa to look for resistance to treatment, further progression and recurrence. Though other biomarkers are unlikely to replace PSA, they are required to improve the sensitivity and specificity of PSA as a PCa biomarker.

## **1.4.2 PCA3**

PCa gene 3 (*PCA3*) is a PCa specific long noncoding RNA (lincRNA), also known as DD3 on chromosome 9q21-22 that is over-expressed in PCa tissue<sup>71</sup>. *PCA3* is not expressed in normal prostate tissue and expression is seen at low rates for hyperplastic prostate tissue, making it the most specific PCa biomarker identified so far. The non-

## CHAPTER 1: INTRODUCTION

coding *PCA3* mRNA functions as a polyadenylated RNA transcript, which does not result in a cytoplasmic protein<sup>63</sup>.

### ***1.4.2.1 PCA3 – Diagnosis***

*PCA3* can be found in urine, but only at sufficient levels, after a DRE is performed<sup>72</sup>, and that comparing the ratio of *PCA3* mRNA quantities with *KLK3* mRNA (which is the transcript for PSA) quantities (very slightly over-expressed in prostate cells in urine) gave high sensitivity and specificity rates, 67% and 83% respectively<sup>73</sup>. The comparison of *PCA3* and *KLK3* mRNA quantities found in prostate cells in urine is known as the *PCA3* score. An assay was generated to simultaneously detect *PCA3* mRNA as well as *KLK3* mRNA in urine: the uPM3<sup>TM</sup> assay. The assay was tested on 158 patients with elevated PSA and/or an abnormal DRE, whom provided a sample with a sufficient amount of prostate cells in the urine. The assay identified PCa in 62 of the 158 patients (39%), with sensitivity and specificity rates of 82% and 76%, respectively. The positive and negative predictive values for the assay were 67% and 87%, respectively. Comparably, PSA had sensitivity and specificity rates of 98% and 5%, with positive and negative predictive values of 40% and 83%<sup>63</sup>.

The performance of the uPM3<sup>TM</sup> assay at different PSA levels (<4ng/ml, 4-10ng/ml and >10ng/ml) was examined, with outcome sensitivity levels of 73%, 84% and 84%, respectively, and specificity levels of 61%, 80, and 70%, respectively<sup>63</sup>. A more stable re-designed assay was later designed and evaluated in a multicenter assessment: The assay had between 94%-100% discriminatory rates in samples after a DRE with at least 3 strokes<sup>72</sup>. This test was then applied to 72 men with known biopsy outcomes, of which 17 were positive for and 55 were negative for PCa, at two centers. Taking the *PCA3* score as a continuous variable, a ROC analysis was performed and both sites were able to correctly classify 49/72 (68.1%) of patients, and the AUC were not significantly different ( $p = 0.9289$ ), this demonstrates significant accuracy between the

## CHAPTER 1: INTRODUCTION

sites ( $p = 0.0085$ )<sup>72</sup>, highlighting the *PCA3* assay as an accurate, reproducible test for the diagnosis of PCa.

Another multicentre saw improvement of *PCA3* on PSA in the “grey zone” (PSA 3-15ng/ml); AUC increased from 0.57 to 0.66 and specificity increased from 47% to 66% for PSA and *PCA3*, respectively<sup>74</sup>. A study looking at multi-gene expression profiling of prostatectomy tissues yielded an AUC for *PCA3* of 0.85 individually but increased with the addition of *EZH2*, prostein and *TRPM8* to 0.90<sup>75</sup>.

### **1.4.2.2 Repeat Biopsies**

The *PCA3* test is effective at identifying patients who were likely to have a positive second biopsy result, after receiving a negative first. A multicentre clinical study of 466 men evaluated the clinical usefulness of the *PCA3* assay for the prediction of repeat biopsy outcome. The study resulted in a suggested *PCA3* cutoff of 25, with patients with a *PCA3* score lower than 25 were 4.56 times as likely to have a second negative result for their repeat biopsy<sup>76</sup>. The *PCA3* test is FDA approved but generally only used in private healthcare in the UK.

### **1.4.2.3 *PCA3* Conclusions**

Although the *PCA3* assay shows significant improvements in specificity and sensitivity compared to PSA, it is significantly more expensive: A *PCA3* test costs between approximately £300 and £400, (whereas a PSA test costs approximately £7) and this cost will increase with the use of gene panels. In comparison a TRUS biopsy costs £312<sup>50</sup>, as you can see the *PCA3* test can be more expensive than just doing the repeat biopsy. The literature and improved sensitivity show that the *PCA3* test is clearly useful but where it fits into PCa diagnostics is unclear at this time. The *PCA3* test is currently available privately but not on the public health care system/NHS in the UK.

## CHAPTER 1: INTRODUCTION

### 1.4.3 AMACR

Alpha-methylacyl-CoA racemase (*AMACR*) is used as an immunohistochemical (section 1.6.7) diagnostic biomarker for PCa. Needle biopsy specimens are stained for *AMACR* during diagnosis of PCa patients<sup>77</sup>, as *AMACR* expression is increased in PCa but may decrease with progression<sup>78</sup>. *AMACR* expression alone was not informative for the prediction of metastatic or lethal PCa; age, Gleason score and stage were also indicative<sup>78</sup> and out of 64 prostate adenocarcinomas no significant correlation was seen between *AMACR* expression levels and histopathological grade<sup>79</sup>.

*AMACR* is an enzyme that regulates the metabolism of branched-chain lipids and drugs and is often overexpressed in PCa tissues<sup>80,81</sup>. It is thought that the synthesis of fatty acids and increased use of branched chain fatty acids plays a role in PCa progression. It is essential for optimal growth of PCa cells in vitro and offers a potential treatment target complementary to hormone therapy. *AMACR* is also frequently seen in tumours of patients with hereditary links to PCa<sup>79</sup>.

### 1.4.4 AR

The Androgen receptor (*AR*) binds androgens leading to the development and survival of prostate epithelial cells. In PCa it allows survival and growth of the tumour and is a known contributor to its progression. Whilst PCas show great heterogeneity, it is obvious that *AR* plays an important role in the survival of the bulk of prostate tumour cells<sup>82</sup>.

Hormone therapies work by blocking androgen-*AR* signalling, inhibiting growth and survival of the tumour. *AR* transcriptional reactivation/rearrangements are fundamental to the inevitable resistance of PCa to hormone therapies and androgen-independent activation of the *AR* pathway. One resistance mechanism is the production of *AR* variants that lack the canonical ligand-binding domain<sup>82</sup>, allowing the transcription of *AR* target genes without the initiating signal of androgen binding. 17 of these *AR*

## CHAPTER 1: INTRODUCTION

variants have been identified, all containing a common core of the DNA binding domain and then NH2 terminal domain and lacking the ligand-binding domain. There are several mechanisms for the production of *AR* variants including: proteolytic cleavage, genomic alterations, and altered exon splicing.

Levels of specific *AR* variants observed in clinical samples are highly variable and not all variants are equivalent at predicting progression and resistance. As well, clinical studies using *AR* variants are limited by the lack of clinically validated assays for the detection of the individual variants. However, these limitations are currently being addressed<sup>82</sup>, suggesting potential clinical use of *AR* variants as biomarkers.

### 1.4.5 SPOP

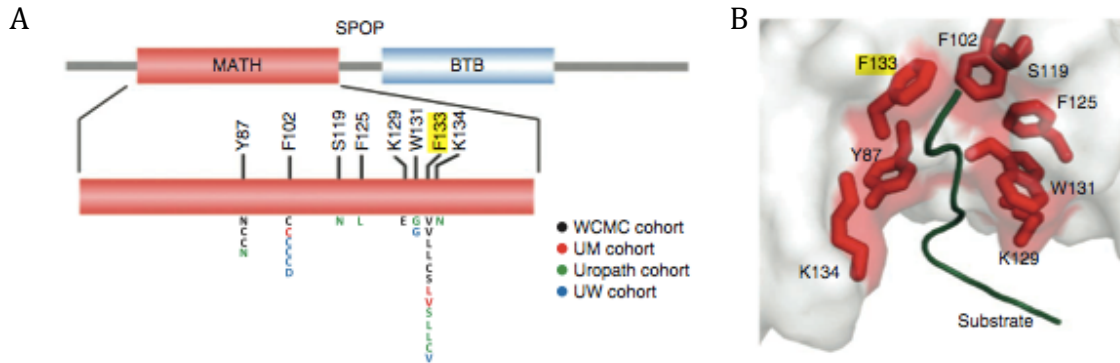
*SPOP*, otherwise known as E3 ubiquitin ligase adaptor speckle-type poxvirus and zinc finger (POZ) domain protein, interacts directly with and regulates SRC-3 (p160 steroid receptor coactivator-3). The p160 SRCs play fundamental roles in the cell proliferations and *AR* transcriptional activity as well as resistance to androgen deprivation therapy<sup>83</sup>. *SPOP* binds wild-type *AR* leading to its degradation; this is promoted by anti-androgens but antagonized by androgens. Whereas, *SPOP* mutants and *AR* alternative splicing leads to *AR* stabilization suggesting a key role in acquiring ADT resistance<sup>84</sup>.

A new molecular subtype of PCa can be defined by mutations in *SPOP*; *SPOP* mutations are found in PCas that lack ETS family rearrangements<sup>85,86</sup>. *SPOP* missense mutations within the substrate-binding cleft were identified in 13% PCas and were the most common mutations in 111 prostate tumours that underwent exome sequencing<sup>87</sup>. This substrate-binding cleft harbours many residues that can be mutated in PCas (Figure 0.4B). The cleft central F133 is the most common site of mutations (Figure 0.4A).

Exome sequencing of 50 lethal heavily pre-treated CRPCs and 11 treatment naïve high-grade localized PCas', showed that four CRPCs had *SPOP* oncogene mutations; 2 point mutations, 1 frame-preserving indel and 1 copy-number call increase<sup>88</sup>. *SPOP* mutations correlate with somatic deletions at chromosome 5q21 and 6q21. *CHD1*, *FOXO3* and

## CHAPTER 1: INTRODUCTION

*PRDM1* are found at these chromosomal regions and are also correlated with *SPOP* mutated PCas<sup>89</sup>. As well as *TMPRSS2:ERG* fusions, *SPOP* does not appear to be mutated in cancers with *Tp53*, *PTEN* and *PIK3CA* mutations<sup>87</sup>.



**Figure 0.4: SPOP frequency of substitutions and substrate binding cleft<sup>87</sup>.** A) the frequency of substitutions in SPOP across four PCa cohorts from Weill Cornell Medical College (WCMC), University of Michigan (UM), Uropath and University of Washington (UW). B) the substrate-binding cleft of SPOP with the positions of all eight residues that can be possibly mutated. Adopted from Barbieri, C. E. et al. 2012<sup>87</sup>

*SPOP* associations with *AR* highlight the need for examining *SPOP* mutation frequencies in men whom do not initially respond to, or very quickly acquire resistance to PCa; *SPOP* mutation detection could potentially be used to stratify patients out of hormone therapy as a treatment.

### 1.4.6 TMPRSS2:ERG

*TMPRSS2:ERG* is a fusion gene that is formed as a result of structural chromosomal rearrangements. *TMPRSS2* is an androgen responsive, prostate specific gene and *ERG* is a transcription factor oncogene belonging to the *ETS* family, both located on chromosome 21. *ETS* family genes are involved in proliferation, differentiation, angiogenesis, inflammation and apoptosis. The fusion occurs via a translocation of sequences that can involve deletion of the intervening sequences between *TMPRSS2* and *ERG*<sup>90</sup>.

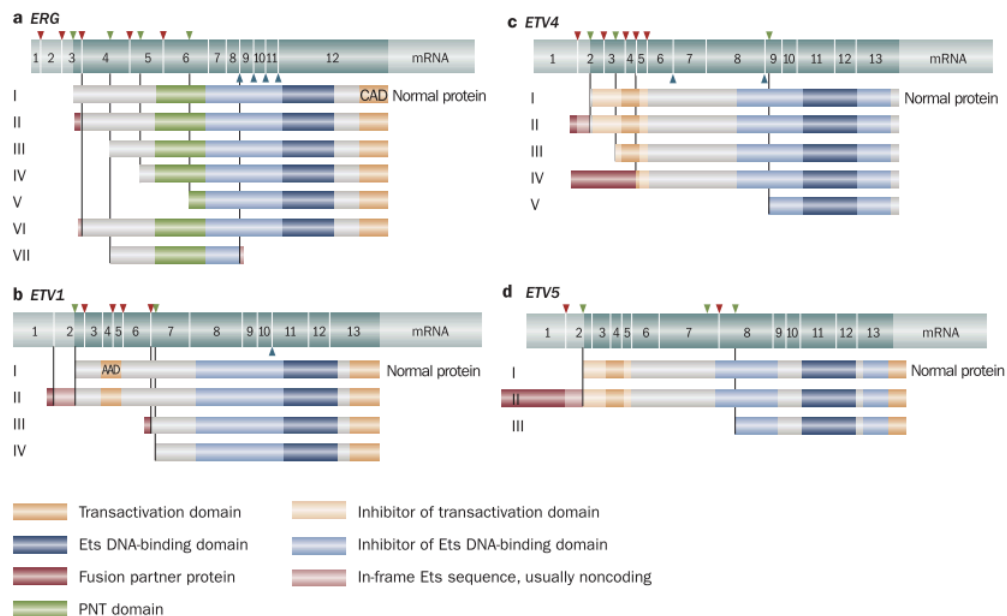
*ERG* has been identified in fusion genes in other cancers; leukaemia and Ewing's sarcoma. *ERG* knockdown inhibits cell growth and invasion and oppositely over-expression leads to invasion and the induction of PCa like lesions on *in vivo* models.

## CHAPTER 1: INTRODUCTION

*ERG* has also been identified to work with mutated members of the *PI3K* pathways leading to the progression of PCa in animal models.

*TMPRSS2:ERG* fusions are seen in ~50% of PCas<sup>91</sup>. *TMPRSS2* also fuses with other members of the *ETS* family (*ETV1*, *ETV4* and *ETV5*) in PCas but at much lower frequency (Figure 0.5). Diversity is also observed in the splice variants of *TMPRSS2:ERG* (Figure 0.5) not only between PCas but also within an individual PCa. The most commonly identified *TMPRSS2:ERG* fusion is *TMPRSS* exon 1 fused with *ERG* exon 4, this is described as T1/E4, the second most commonly found is T1/E5<sup>92</sup>.

It remains controversial for if *TMPRSS2:ERG* fusions are implicated in a poor clinical outcome. A number of studies now suggest it is not the major factor of clinical outcome, but that in a combination of copy number gain and other genetic aberrations (like *PTEN* loss) it can offer prognostic information<sup>92</sup>. Yet many papers still suggest that *TMPRSS2:ERG* fusions are implicated in mediating advanced PCas<sup>93</sup>. However, it has also been shown that early cancers and HG-PIN can also harbour *TMPRSS2:ERG* fusions.



**Figure 0.5** *ETS* family partners for *TMPRSS2* fusion and their splice variant diversity. Adopted from Clark, J *et al.*, 2009<sup>92</sup>.



## CHAPTER 1: INTRODUCTION

### ***1.4.6.1 TMPRSS2:ERG as a Therapeutic Target***

During the 1990's, a leukaemia fusion; BCR-Abl (the Philadelphia chromosome) emerged as a target of treatment (Imatinib) in Philadelphia chromosome positive (Ph+) myeloid leukaemia<sup>94</sup>. *TMPRSS2:ERG* has a prevalence of approximately 50% and is one of the commonest of all cancer fusion genes in solid tumours, making it a good potential therapeutic target. However, studies have shown that *TMPRSS2:ERG* does not increase cellular proliferation or anchorage-independent growth, but instead induces a transcriptional program associated with invasion<sup>95</sup>. Knockdown of ERG transcriptional programming in ETS-positive cancers lead to an inhibition of invasion in the VCaP cell line. Direct over expression of *ERG* in both VCaP and benign prostate cells mediate cellular invasion through engagement with plasminogen activation pathway components, potentially showing a downstream target that could be used as a drug target<sup>96</sup>. *TMPRSS2:ERG* fusions have also been implicated in signalling pathways and ion channel genes creating further opportunities for therapeutic targeting of these fusion positive cancers<sup>92</sup>.

Shao *et al.*, have shown that targeting the most common and clinically significant alternatively spliced isoforms of the *TMPRSS2:ERG* fusion using siRNAs delivered by liposomal nanovectors resulted in the inhibition of tumour growth *in vivo*<sup>97</sup>. The mice with orthotopic or subcutaneous xenograft tumours (with the target fusions) also showed no sign of toxicity. Therefore, *TMPRSS2:ERG* targeting could be a potential future therapy for PCa.

### **1.4.7 Biomarkers for pre-disposition to PCa**

Family history has been significantly associated with a higher risk of PCa ( $p = 0.01$ , odds ratio, 1.39; 95 percent confidence interval, 1.07 to 1.79;) in a study of 2,950 men, all with an initial PSA of less than 4ng/ml. Of the 2,950 men, 477 were family history positive and 2,473 were family history negative. After a seven-year follow up, 449 men were diagnosed with PCa; 94/477 (19.7%) that were family history positive and

## CHAPTER 1: INTRODUCTION

355/2,473 (14.4%) that were family history negative<sup>59</sup>. Family history in a first-degree relative (brother, father, or son) is said to double a man's risk of developing PCa, with increasing risk as the number of affected relatives rises<sup>98</sup>.

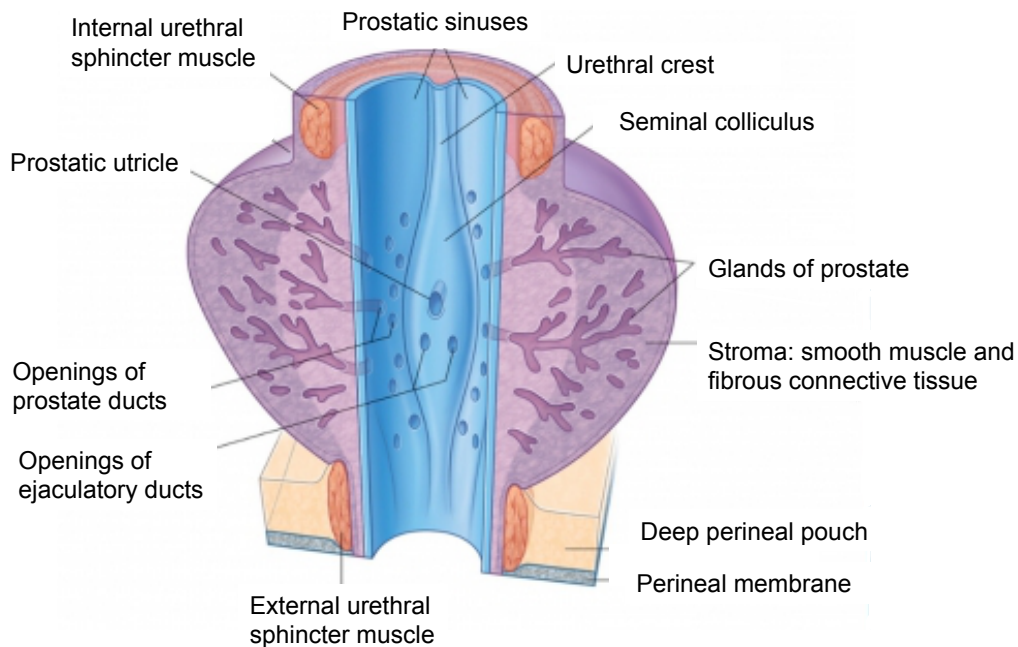
*BRCA2* mutations increase relative risk by 5-23 fold in men above 60 years of age, however, the frequency of *BRCA2* mutations is low and can only account for a small number of PCa susceptibility cases<sup>99</sup>. *BRCA2* mutation carriers are in higher risk of developing PCa than *BRCA1* mutation carriers and studies into *BRCA1* mutations suggest they have limited contribution to PCa risk<sup>100</sup>. Breast cancer linkage consortium studies (BCLC) found that *BRCA2* carriers risk was also based largely on age and the mutation location<sup>100</sup>.

Genome-wide association studies (GWAS) have led to the identification of more than 46 single nucleotide polymorphisms that have low penetrance in PCa<sup>99</sup>. As discussed by Goh et al., these include SNPs at loci or close to loci known to be involved in PCa such as *KLK3*, *AR*, and *AR* transporter genes<sup>99</sup>. Low penetrance genes were investigated because evidence has suggested that the risk of developing PCa is likely related to a combination of loci conferring low to moderate risk of the disease and, not so commonly, alleles with higher risk such as *BRCA2*<sup>99</sup>.

As targeted therapies and screening for PCa becomes more widely used, the use for pre-disposition biomarkers will become increasingly important<sup>1</sup>.

## 1.5 Urine and Exosomes

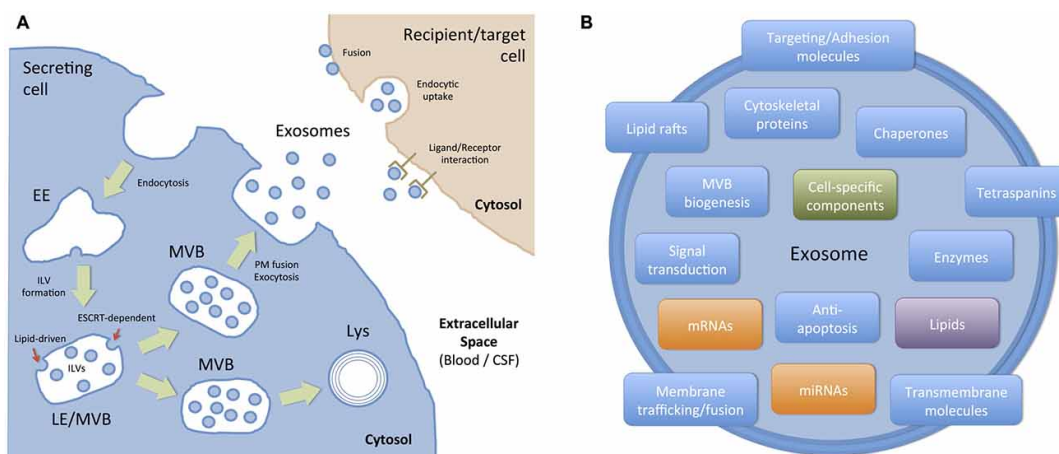
The PCA3 test (section 1.4.2), as previously discussed, proves that urine contains PCa specific biomarkers. The anatomy and location of the prostate make urine a viable source of prostate biomarkers; urine from the bladder passes through the middle of the prostate, where secretions from the prostate glands can enter the urine (Figure 0.6). DRE manipulates a more abundant release from these glands allowing prostate and PCa specific markers to be detectable in urine (such as *PCA3*, *KLK3* and *TMPRSS2:ERG*)<sup>72</sup>.



**Figure 0.6 Anatomy of the prostate.** Adapted from Drake *et al.*, 2015<sup>101</sup>.

Urine holds an advantage over tissue biopsies in that it potentially allows an overview of all foci of cancer in one go. More than ~80% of cancerous prostates have more than one tumour focus<sup>102</sup>, and each cancer focus will have a number of variant tumour clones with divergent genetic and epigenetic changes. Biopsy sampling is incapable of capturing the diversity of cancer within a prostate.

## CHAPTER 1: INTRODUCTION



**Figure 0.7 Tumour cells send signals to distant cells through exosomes. A) Production of exosomes and how they can be sent to recipient cells. B) The different materials that can be found inside exosomes. Adopted from Bätz, L.F., 2016<sup>103</sup>.**

Exosomes are endocytic membrane derived microvesicles 30-120nm in size. They can be found in many biological fluids including those that are easily attainable like blood and urine, which also see elevated exosome secretions during malignancy<sup>104</sup>. Exosomes are a key component of biological trafficking across membranes and play a key role in cell homeostasis. In cancers, aberrant exportation of proteins and RNAs via exosomes can lead to miss-expression in cells that take up the exosome. Exosomes contain proteins, lipids and nucleic acids that can be involved in cell-to-cell communication (Figure 0.7), through their release into surrounding cells. Exosomes derived from tumour cells have roles involved in tumourigenesis, metastasis, and response to therapy by transferring mRNA, miRNA and proteins between cancer cells and the tumour microenvironment<sup>105</sup>. Also ligand binding can trigger a signalling cascade in the target cell. Exosomes have the ability to cross talk/influence key tumour-related pathways (such as those involved in the hallmarks of cancer<sup>106</sup>) including hypoxia driven EMT, evading immune responses, angiogenesis and metastasis<sup>107</sup>. The content of exosomes (miRNA, proteins and mRNA) have been shown to cause changes in a) neighbouring cells, b) the tumour microenvironment and c) in distant cells. “Exosomal shuttle RNA” can be transferred via exosomes from the cell of origin to a recipient cell where it can be translated<sup>107</sup>. Exosomes originating from tumours have been shown to educate non-transformed cells in host tissues to create a pro-metastatic phenotype pre-metastasis.

## CHAPTER 1: INTRODUCTION

Hoshino *et al.*, showed that treatment of organ-specific cells with lung-tropic model derived exosomes can redirect metastasis of bone-tropic tumour cells<sup>108</sup>. Specific exosomal integrins are associated with organ-specific metastasis and so could be useful in predicting which organs metastasis will occur in. Costa-Silva *et al.*, showed that exosomes derived from pancreatic ductal adenocarcinomas was able to create a pre-metastatic niche in livers of naïve mice and also increased the metastatic burden within the liver<sup>109</sup>.

Thus, it could be said that looking for biomarkers in exosomes is like raiding cancers' letterbox. The molecular composition of exosomes vary with cell and tissue of origin<sup>107</sup> and can also be altered by pathophysiological changes in the cell of origin, meaning exosomes have great potential for cancer biomarkers.

Some RNAs are enriched within the exosomes at several 100-fold compared to cells, and transcripts that may have very low copy numbers in tumour cells could be detected at much higher relative levels within exosomes<sup>110</sup>. Nilsson *et al.*, were able to show that exosomes in urine contained genetic information that is directly from PCa cells<sup>111</sup>. Both *PCA3* and *TMPRSS2:ERG* transcripts were detected in the exosomes. Dijkstra *et al.*, showed that the genetic content of exosomes differs from that of the cell sediment<sup>112</sup>. Exosome membranes can resist ribonuclease and DNase digestion of their contents allowing a better-protected RNA inside in comparison to cell RNA. Exosomal RNA will be similar on harvest as when it left the cell, in contrast to cellular RNA which will be altered on loss of cell:cell contact and entry into the non-life sustaining environment of urine. These points make exosomes a stable, viable, and more promising source of PCa biomarkers than cells harvested from urine.

### **1.6 Methods in Biomarker Discovery**

Over the past two decades extensive investigations have proven that cancer is a heterogeneous disease with diverse genomic aberrations<sup>113</sup>. These genomic aberrations

## CHAPTER 1: INTRODUCTION

consist of gains, losses and rearrangements of chromosomal fragments, specific gene mutations and epigenetic alterations including methylation. These can lead to aberrant transcript expression and incorrect protein production at differing levels between disease and benign states.

Many cytogenetic and molecular tests have been developed to detect such aberrations. As technologies advance, more effective, less time consuming and cheaper methods are available for biomarker discovery and their validation (Table 0.2).

**Table 0.2: Cost of different technologies available for biomarker discovery.**

<i>Technique example</i>	<i>Number of transcripts</i>	<i>Batches of Samples</i>	<i>Amplification Required</i>	<i>RNA usage</i>	<i>~Cost/Sample</i>
<i>NanoString</i>	<i>&lt;800</i>	<i>12</i>	<i>Y</i>	<i>20ng</i>	<i>£50/sample</i>
<i>Microarray</i>	<i>30,000</i>	<i>1</i>	<i>Y</i>	<i>20ng</i>	<i>£400/sample</i>
<i>Sequencing</i>	<i>All</i>	<i>1</i>	<i>N</i>	<i>100ng*</i>	<i>£1,000/sample</i>
<i>qRTPCR</i>	<i>1+</i>	<i>1</i>	<i>Y</i>	<i>20ng</i>	<i>\$35/sample</i>
<i>Targeted Sequencing</i>	<i>250</i>	<i>1</i>	<i>N</i>	<i>1ng*</i>	<i>\$50/sample</i>

\*RNA not amplified and used directly in technology.

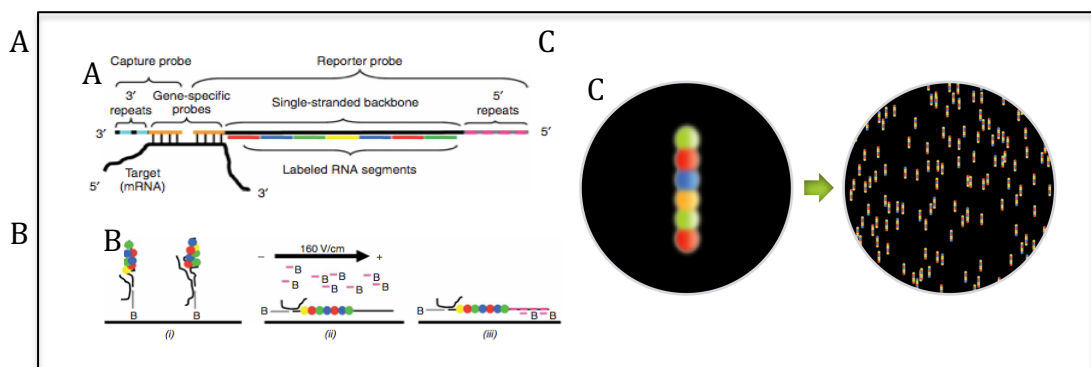
### 1.6.1 Nanostring

The Nanostring nCounter gene expression system was made available in 2008 and is capable of capturing and counting individual mRNA transcripts. It provides direct count data for each of the target genes via a two-probe system: A capture probe and a reporter probe. Both probes are hybridised to the mRNA, the reporter probe hybridising to sequence adjacent to the capture probe (Figure 0.8A). The reporter probes are specifically labelled with a series of fluorescent ‘beads’ that are unique for each gene. The capture probe is biotinylated and the mRNA/probe combination is captured by binding to a streptavidin coated slide. The DNA on the slide is then subjected to a voltage which stretches out the molecules on the slide (Figure 0.8B). The slide is then washed to remove excess probes, and the slide is photographed. The bead codes are counted to give the frequency of each mRNA in the sample (Figure 0.8C)<sup>114</sup>.

## CHAPTER 1: INTRODUCTION

In comparison to microarrays (section 1.6.4), NanoString technologies allow quantification of small amounts of starting materials (100ng), and mRNA levels can be measured without the need for amplification. By allowing the customer to choose specific targets, use of NanoString over array can work out cheaper per sample. Microarrays will provide >34,000 targets and cost ~£400-500, however, if you want a select cohort of genes (maximum 800 per analysis), NanoString can allow a cheaper overall experiment. NanoString is also more specific and has a better dynamic range than microarrays. The reaction is performed in solution and not fixed to a solid surface allowing the reaction to be driven to completion and so boasts higher sensitivity. The Nanostring nCounter system also allows a pure digital readout of transcript counts that claim to have less background noise, and be less ambiguous in downstream analyses than those that use analog signals, like microarrays<sup>115</sup>.

A disadvantage is that due to the barcode system it utilises, there is a limited number of probes (capped at 800 for a custom codeset)<sup>116</sup>. Again, like microarrays unknown mutations are not identified via Nanostring, and so for the identification of these, sequencing is still preferred and similarly to microarrays and PCR, the quality of the data is dependent on the quality of the probe.



**Figure 0.8 NanoString Ncounter system. A) The set up of the two probes (capture and reporter), one target system. B) The elongation and fixing of probes using a current for imaging. C) Imaging of the uniquely labelled reporter probes. Adapted from Geiss, G et al., 2008<sup>115</sup>.**

## 1.6.2 Sequencing

In 1977 Frederick Sanger published Sanger sequencing, a method using the incorporation of chain-terminating dideoxynucleotides by DNA polymerase, which cause base-specific termination during DNA synthesis<sup>117</sup>. This was a fundamental breakthrough for science and allowed a monumental accomplishment: the finished grade human genome sequence in 2001. Since then, sequencing technologies have advanced and become considerably cheaper: In 2001 it cost \$100 million to sequence a genome and since late 2014 it is ~\$1,000. The biggest price drop occurred in 2008 (Figure 0.9) and was a consequence of the introduction of commercialised next-generation sequencing (NGS) technologies. In 2015 the production of Illumina’s HiSeq X Ten allowed the first \$1,000 sequenced genome<sup>118</sup>.

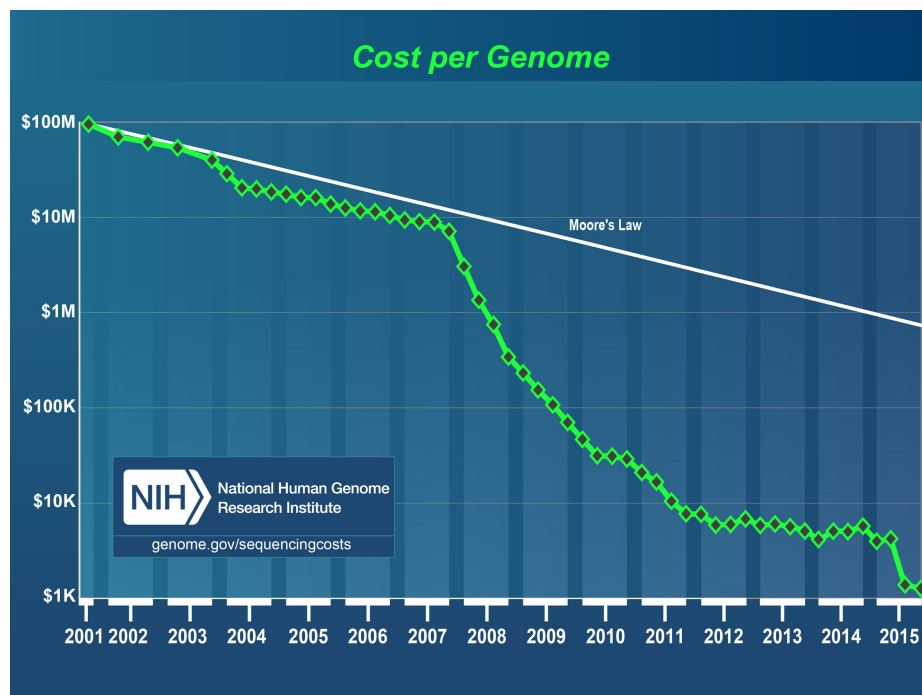


Figure 0.9 Sequencing cost per genome from 2001 to 2015. Sudden drops seen in ~2008 and again in 2015. Adapted from National Human Genome Research Institute (NHI) 2016<sup>119</sup>.

### 1.6.2.1 Next Generation Sequencing (NGS)

Next generation sequencing began with the discovery of the pyrosequencing method using luminescent for measuring pyrophosphate synthesis. This was a two-enzyme



## CHAPTER 1: INTRODUCTION

process whereby ATP sulfurylase converts pyrophosphate into ATP. ATP is the substrate for luciferase, which produces a proportional amount of light to the amount of pyrophosphate produced as each nucleotide is washed over template DNA that is fixed to a solid phase. This method is still, similarly to Sanger sequencing, a sequence by synthesis method. Benefits included using natural dNTPs, and being observed in real time without the need for electrophoreses. A disadvantage of this was that identification of more than 4-5 identical nucleotides proved to be difficult. Further improvements in methodology including using beads for DNA attachment and enzymes for degraded unused dNTPs (removing the lengthy wash step), led to the first commercial NGS technology by 454 Life Sciences. This allowed massive parallelisation of sequencing reactions, meaning the amount of DNA sequenced in one run was significantly increased<sup>120</sup>.

Following the success of 454's high throughput sequencing machines, a number of new techniques were developed, including the Solexa method of sequencing, which was later acquired by Illumina. The Solexa method used bridge amplification, where DNA molecules were run across complementary oligonucleotides bound to a flowcell. Here, the original flow-cell binding DNA strands arch over to prime the next round of polymerisation for neighbouring oligonucleotides to create clusters of clonal populations by solid phase PCR. This is another example of sequence by synthesis, although here modified dNTPs with a fluorescent 'reversible-terminator' occupies the 3' hydroxyl position. These fluorophores needs to be cleaved prior to the next polymerisation step, allowing sequencing in a synchronous manner (

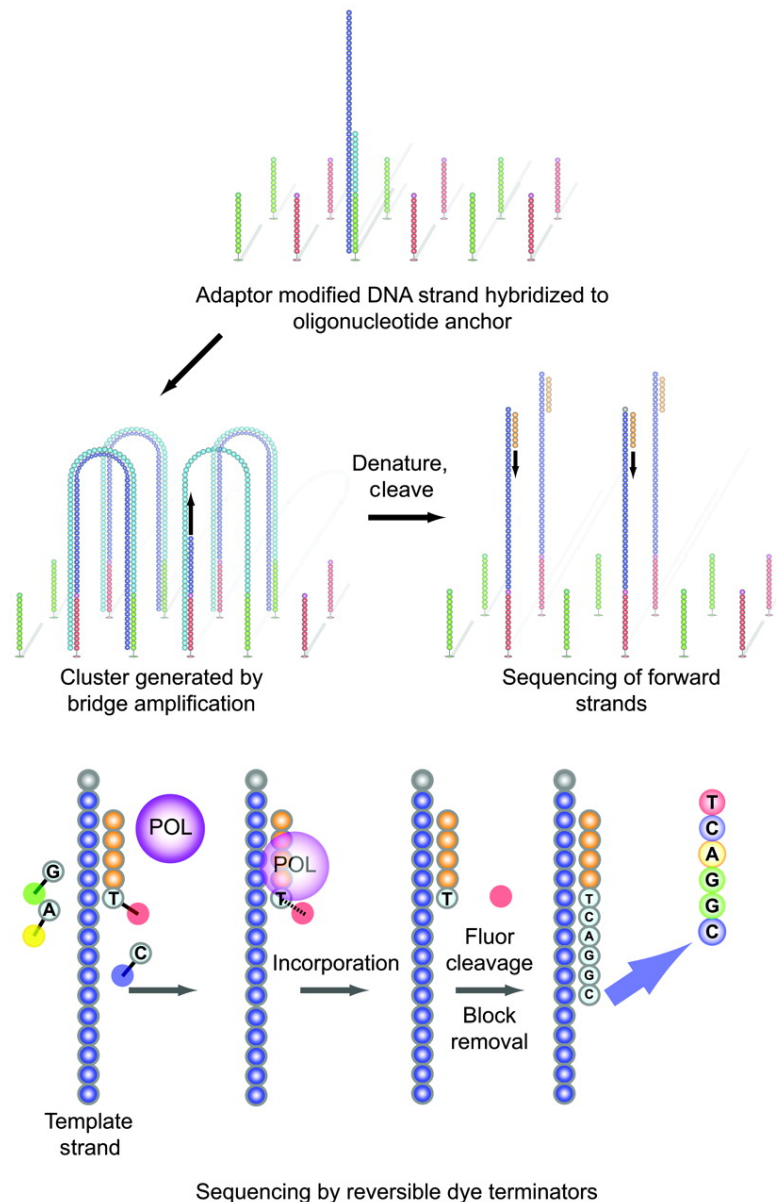
## CHAPTER 1: INTRODUCTION

Figure 1.10).

Illumina created the first Paired end sequencing, improving efficiency and accuracy when aligning to a reference genome by providing positional information<sup>121,122</sup> and decreased sequencing costs per template<sup>123</sup>. Paired-end sequencing enables improved biological applications, allowing genome-wide identification of gene fusions, insertions, deletions and translocations and spliced exons because it retains information on the distance and relationship between two ends of DNA fragments<sup>121,123</sup>.

Illumina's HiSeq series then used a further improved method to allow longer read length and depths. Disadvantages include substitution errors (commonly after "G" incorporation), under-representation of AT-rich and GC-rich regions (due to amplification bias) and a 2.5% false positive rate for novel single nucleotide variants (SNVs)<sup>124</sup>. Illumina is the most commonly used sequencing platform: The HiSeq series is still used commonly for genome sequencing, whilst Illumina's other machines are used for other applications. MiSeq is used for experiments that require lower-throughput and longer read lengths with a faster turn around<sup>121</sup>. NextSeq machines are desktop sequencing tools with fast turn around time used for transcriptome and targeted re-sequencing and thus is commonly used for clinical settings.

Although there are many NGS platforms (Roche/454, Illumina, and Pacific Biosciences, etc.), all use spatially separated, amplified or single DNA molecules, in a flow cell that are massively parallel sequenced<sup>125</sup>. NGS technologies have provided us with an ability to produce enormous amounts of data at a relatively cheap cost. The ever-increasing amounts of DNA sequenced, longer reads and faster turn around times are constantly improving the sequencing technologies.



**Figure 1.10** Solexa's sequencing methodology using bridge amplification. DNA strands bound by complimentary oligonucleotides to a flow cell arch over to prime the next round of polymerization. This creates clusters of clonal populations via PCR. Fluorophores that can be cleaved between steps show the incorporation of the next dNTP. Adapted from Voelkerdig et al., 2009<sup>126</sup>

### 1.6.2.2 Third Generation Sequencing

The Oxford Nanopore, Pacific Biosciences' (PacBio) Single Molecule Real Time (SMRT) and Illumina's Tru-seq Synthetic Long-Read are the three commercially available third generation sequencing technologies<sup>127</sup>. Third generation sequencers can be considered as those that are capable of sequencing single molecules (SMS), which negates the need for DNA amplification<sup>121</sup>, and can produce much longer reads (generally between 5,000-15,000 bp)<sup>128</sup>. Pac Bio's SMRT was the first, released in

## CHAPTER 1: INTRODUCTION

2010, and the reads generally had a raw error rate of 10-15%. However algorithmic techniques and a 50x long read coverage (for *de novo* genome assembly) can allow correction. The main limitation is cost compared to second generation technologies<sup>128</sup>.

Illumina's Tru-seq Synthetic Long-read was released in 2012. Long DNA molecules are clonally amplified and barcoded prior to sequencing using a short read instrument this results in synthetically produced long reads from the short read sequences. This technology boasts a high accuracy without the need for correction but the standard illumina shortcomings are the same; high GC content and tandem repeats remain troublesome. For *de novo* genome assembly, cost can be even greater than that of PacBio's SMRT because for 30x long read coverage you need 900x – 1500x short read coverage<sup>128</sup>.

Oxford Nanopore's MinION is the newest, released in 2014 and is a handheld device. It works by measuring the small disruptions to an electric current as DNA molecules flow through a nanopore. The MinION has low accuracy and throughput compared to the other third generation technologies. Accuracy can be improved with correction algorithms like those used for the PacBio SMRT. A major benefit of the MinION is its size, cost, and speed, allowing its use in remote areas and for breakout classification<sup>128</sup>. Further improvements on accuracy can make the MinION a powerful tool for the future.

### **1.6.2.3 Exome sequencing**

It is estimated that 85% of mutations that cause disease can be found in coding and functional regions of the genome, and therefore, can be identified through exome sequencing rather than whole genome sequencing. Sequencing only the exomes provides a lower cost per genome/exome and whole exome sequencing provides coverage of more than 95% of the exons<sup>129</sup>. Therefore, exome sequencing can be used to identify the majority of cancer biomarkers at a much lower cost. Exome sequencing can also be used to target non-coding elements such as microRNA and lincRNA<sup>130</sup>. The exome can be captured using either solution based or array based technologies. Solution

## CHAPTER 1: INTRODUCTION

based exome capture is most commonly used. Biotinylated oligonucleotide probes to target regions in the genome are used to capture fragmented DNA. Streptavidin beads then bind the probes and untargeted DNA is washed away. PCR is used to amplify the captured target DNA and this is then sequenced<sup>131</sup>. Solution based capture is most commonly used even though array based capture was the first to be used, this is likely due to less input DNA requirements. However, array based capture has proven to be useful in low GC content regions and SNP detection<sup>131</sup>.

### ***1.6.2.4 RNASeq***

RNASeq, first used in 2008, is when next generation sequencing approaches are used to sequence total cDNA, allowing quantitative expression scores (similar to microarrays). However, the entire transcriptome can be observed (without prior knowledge requirements for probe production), including novel transcribed regions and transcript structures, such as alternatively spliced isoforms, can also be identified<sup>132</sup>. Due to the desire to determine differential splicing activity, antisense transcription and novel transcriptional regions in eukaryotes, RNASeq has been key milestone for biological experiments in these organisms. The resolution and sensitivity that can be achieved and the range of different changes that can be observed give RNASeq advantages over microarrays. However, there is a significant extra cost, bioinformatics requirements and data storage required for RNASeq experiments<sup>133</sup>. Due to the role of NGS in RNASeq experiments, the limitations of NGS technologies are still present (section 1.6.2.1).

RNASeq experiments have allowed a better understanding of transcription initiation sites, improved detection of alternatively spliced variants, and fusion genes as well as a better identification of sense and antisense transcripts. All of these things are key to cancer research<sup>134</sup>. Developments in RNASeq methods to allow low-input (cDNA pre-amplification) and the use of unique molecular identifiers (UMIs) have allowed single cell RNA sequencing experiments that can identify transcriptomic variation between genetically homogenous cells. This is very important in cancer research where cancer

## CHAPTER 1: INTRODUCTION

cells are known to have subpopulations with heterogeneous mutations and transcriptomes<sup>135</sup>.

### **1.6.2.5 *CHiPSeq***

ChiPSeq, chromatin immuno precipitation sequencing, is the sequencing of DNA fragments that co-precipitate with a DNA binding protein. The most common of the DNA binding proteins investigated with CHiPSeq are transcription factors, chromatin modifying enzymes or modified histones that interact with the DNA. DNA segments that are associated with a specific DNA-binding protein can be identified with ChiPSeq in an unbiased manner, without existing knowledge of precise DNA binding sites<sup>136</sup>. ChiPSeq allows experiments to study gene regulation.

### **1.6.2.6 *Targeted Sequencing***

The decreasing cost and improvements to second-generation sequencing technologies mean sequencing of complex organisms will eventually become routine. Currently, sequencing large numbers of whole genomes of Eukaryotes routinely is not yet feasible and thus enrichment for areas of interest can reduce time and cost<sup>137</sup>. There are a number of methods to selectively “capture” genomic regions for sequencing, known as target-enrichment; each has their own advantages and drawbacks. These include PCR (Section 2.1.6), molecular inversion probes (MIP), on-array hybrid capture and in-solution hybrid capture<sup>137</sup>.

PCR has been widely used prior to sequencing in experiments. It boasts high sensitivity, good specificity, uniformity and robustness. However, there are issues such as cost, difficulty to multiplex (with the simultaneous use of multiple primers, high levels of nonspecific amplification are observed due to interaction between primer pairs), and an upper limit to the generated amplicon size. Also, in practice not all amplification reactions yield products, which is a key problem when working with clinical samples<sup>137</sup>.

## CHAPTER 1: INTRODUCTION

MIP uses the enzyme ligase to circularize single stranded oligonucleotides formed of a common linker flanked by target-specific sequences. Exonucleases are then used to digest uncircularised species, leaving only the circularised oligonucleotides to be amplified via PCR, using primers targeting the linker. DNA polymerase is used to “gap fill” between target specific MIP sequences. Gap fill and PCR can occur in small volume, aqueous solution, meaning they are easy to scale to large numbers via a 96-well plate. Another advantage is that barcodes for identifying purposes can be incorporated into the primers allowing pooling of multiple samples and input requirements can be as low as 200ng<sup>137</sup>. Issues include capture uniformity, which have been improved modestly but remain this technique’s biggest downfall.

Hybrid capture is performed using immobilised specific probes that hybridise the shotgun fragment library and the un-targeted DNA strands are washed away whilst those captured are eluted. Arrays can hold 2.1 million probes per array with the ability to capture 34Mb<sup>137</sup>. Compared to PCR based approaches, array techniques are quicker and less laborious. Hybrid capture also has its drawbacks including expensive hardware, high starting material requirements (10-15µg) and limits to a) the number that can be performed in a day and b) the number of samples in a study (large numbers aren’t feasible).

In solution capture is similar to array capture, with an excess of probes allowing less starting material. Again this technique can be used in 96-well plates meaning it is readily scalable without the need for specialist equipment<sup>137</sup>.

### **1.6.3 Polymerase Chain Reaction (PCR)**

PCR is an important laboratory technique that is capable of amplifying a single DNA sequence to make thousands/millions of copies. The PCR procedure has multiple heating and cooling steps. The reaction mix (DNA, dNTPS, DNA polymerase, buffer) is heated to 94-98°C to denature double stranded DNA and then cooled to enable the

## CHAPTER 1: INTRODUCTION

sequence-specific hybridisation of the primers to the single stranded DNA. DNA polymerase then makes a complementary DNA strand extending from the 3' end of the hybridised primer. These heating and cooling steps can then be repeated to create more and more copies of the DNA.

PCR can be used to detect presence/absence of a specific target as well as to quantify the amount of target present. Presence/absence can be observed via gel electrophoresis, using a ladder of known sizes to obtain product size. Quantification is generally performed using fluorescent dyes.

There are multiple uses for PCR, for example real-time PCR can monitor the amplification of target nucleotide sequences in real time by either using fluorescent dyes that intercalate between dsDNA in a non-specific manner or by using target specific probes that are fluorescently labelled. The number of cycles required for the product to exceed a predetermined fluorescence threshold is measured (as a cycle threshold- or ct-value) to infer the amount of starting target material. Quantification can be also be performed post-PCR. Nested PCR uses using two sets of primer pairs in sequential reactions. It is used to reduce non-specific probe binding, and increase sensitivity: PCR product from a first PCR is used to seed a subsequent PCR containing a second set of 'nested' primers that hybridise to sequences 3' to the first round primers in the amplified product. This improves specificity as it is unlikely that DNA other than the intended target sequence would hybridise to both primer pairs.

PCR has been used to detect mutations and biomarkers, and to diagnose cancer. Leading up to the development and cost reduction of NGS (section 1.6.2.1), many scientists were using PCR-based investigations into cancer biomarker discovery: A reverse transcription-PCR assay of 761 transcripts was used for the discovery of colon cancer biomarkers<sup>138</sup>. Comparisons with targeted NGS have shown that real-time PCR and NGS have significant concordance (96.3 to 100%) for detecting *EGFR*, *KRAS* and *BRAF* mutations in FFPE materials. However, NGS was capable of identifying seven non-synonymous SNVs and an indel in EGFR that was not detected by the real-time



## CHAPTER 1: INTRODUCTION

PCR method<sup>139</sup>. PCR is also commonly used for target enrichment for targeted sequencing of genes or specific transcript splice variants (section 1.6.2.6).

### ***1.6.3.1 OncotypeDx***

OncotypeDx<sup>14</sup> is a multi-gene expression array that uses quantitative reverse transcription polymerase chain reaction-based assay. It is used clinically to give prognostic and predictive value in early stage ER+ breast cancers, to predict the benefit of chemotherapy with adjuvant hormone therapy<sup>14</sup>.

### ***1.6.3.2 Prolaris***

Prolaris is another quantitative reverse transcription polymerase chain reaction-based assay. It can be used (alongside patient and tumour information) to predict the aggressiveness of PCa. It utilises thirty-one cell cycle progression genes and fifteen housekeeper control genes for PCa tissue<sup>16</sup>. The expression of the thirty-one cell cycle progression genes are correlated with PCa proliferation to serve as a risk-stratification tool: a lower score means lower risk and these men may be prime candidates for AS and a higher score represents those needing treatment<sup>140</sup>. It can also be used to predict ten-year PCa specific mortality and ten-year PCa BCR<sup>16</sup>. Cell cycle progression genes have also been used in the prognosis of other cancers<sup>16</sup>.

## **1.6.4 Microarrays**

Prior to the affordability of sequencing as a method to identify biomarkers microarrays were frequently used, and are still incredibly valuable due to cost and availability of standard pipelines for analysis. Microarrays are significantly cheaper than sequencing and so are still often used today. Microarray technology allows the user to assess DNA copy number or RNA expression levels in cells or tissues in different disease states. They are relatively cheap, not considerably time consuming and the array data can be re-investigated for many different questions. Microarrays have been utilised in gene

## CHAPTER 1: INTRODUCTION

discovery and regulation involved in physiological, developmental and pathological processes, in diagnosis and drug discovery<sup>141</sup>.

Microarrays are an array of specific DNA sequence 'probes'. Fluorescently labelled DNA/cDNA samples are hybridised to the probes, the excess DNA washed off, and the quantity of each DNA sequence is assessed by the strength of the fluorescent signal that remains attached to each probe. on the array, However, cross-hybridisation is an issue in microarray experiments, leading to false positives, and masking of eg down-regulated transcript signals. Analyses can use either single-channel (one sample hybridised) or two-channel microarrays (two differentially labelled samples hybridised at the same time).

Two-channel microarrays have been used to directly compare gene expression between two different conditions, e.g. cancer cells with normal cells to identify genes that have expression changes in cancer. The two samples (one cancer and one normal) are labelled with two different fluorophores (often cy3 (green) and cy5 (far-red)) the two samples are then mixed and simultaneously applied to the microarray for hybridisation (Figure 0.11).

Single-channel microarrays provide intensity data for each specific target DNA/cDNA that hybridises to its matching probe. It provides data on the relative abundance of each probe sequence in a sample and can be compared to data from multiple samples A downside to single-channel microarrays is that unless great care is taken in consistency of sample preparation, microarray hybridisation and washing conditions etc., then error rates can be higher than those achieved from two-channel microarrays.

Oligonucleotide microarrays, like single channel microarrays use one fluorescent label for all of the samples (Figure 0.11). They use short genomic ssDNA fragments that allow sequence coverage of an entire genome, and therefore, can be used for extensive genetic profiling and mutational analysis by providing absolute yield values for each specific target gene. They are capable of providing a presence or absence call for each gene, but two separate arrays are required to allow the comparison of healthy controls to

## CHAPTER 1: INTRODUCTION

cancer patient samples<sup>142</sup>. Affymetrix is the major producer of microarrays: They provide standard arrays for many species, for example, the human genome U133 array, which contains 45,000 probe sets for 39,000 transcripts from 33,000 well-substantiated human genes. Affymetrix also produce of custom arrays for a wide variety of different uses. A standard affymetrix array contains oligonucleotide probes, 25 bases long, specific to targets are fixed to a glass wafer, in set locations. Each oligo is present in millions of copies to allow accurate interpretation of expression levels, from measured intensities of fluorescence given by the tagged hybridised nucleotide sequences.

### ***1.6.4.1 Decipher***

Microarrays have been used in cancer biomarker platforms such as Decipher<sup>15</sup>. Decipher is a classifier score calculated from a gene expression microarray analysis of 22 coding and non-coding RNA probes<sup>15</sup>, that predicts metastatic PCa progression/high risk of recurrence and PCa related mortality within 5 years of RP. High-risk of recurrence is defined by extra-prostatic extension, seminal vesicle invasion, positive margins or biochemical recurrence. Whilst the 22 specific probes are unknown, the panel represents known pathways involved in aggressive PCa, including cell proliferation, cell structure, immune system modulation, cell cycle progression and androgen signalling<sup>143</sup>.

The score ranges from 0-1, with every 0.1 increase representing a 10% increase in risk of metastatic progression. The score is then more generalised into three categories; low-risk 0-0.44, intermediate risk 0.45-0.59 and high risk 0.6-1.

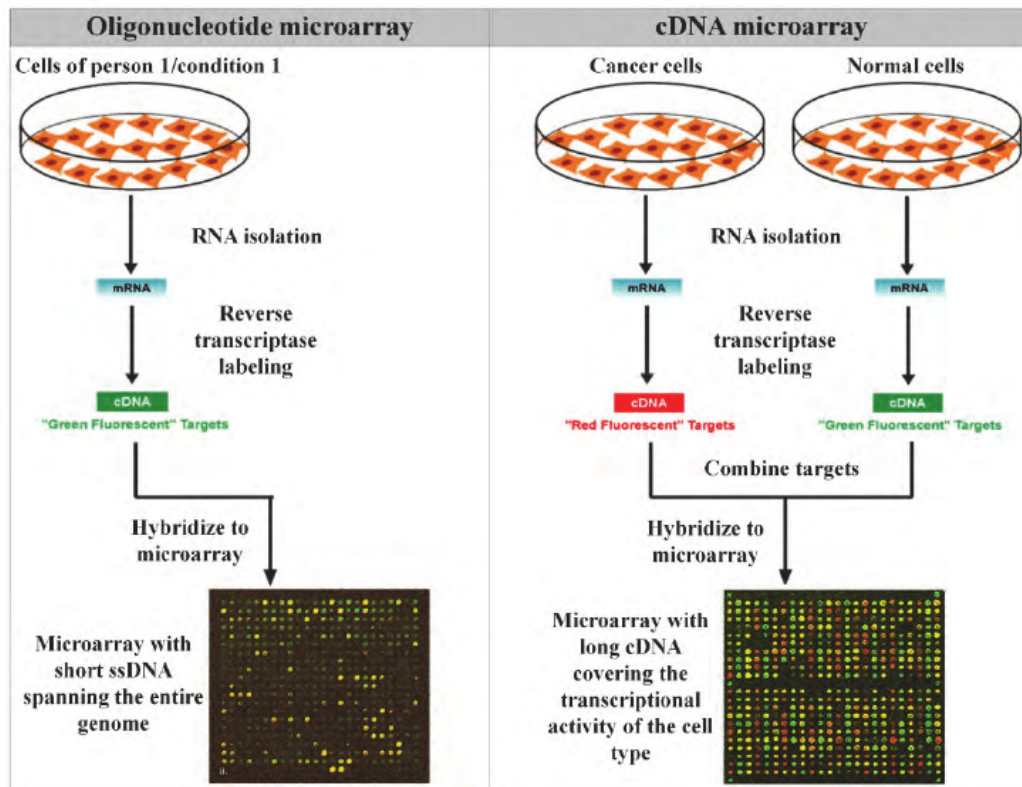
Whilst Decipher was originally established as a predictor for metastatic progression post-RP, there have been further applications since. Decipher has been evaluated for its ability to ease decision making between adjuvant and salvage radiation therapy (second-line treatments); Dalela et al., showed that any two or more risk factors from pT3b-T4, G8-10, lymph node invasion or >0.6 decipher score showed a 4-fold reduction in metastatic progression at 10 years with adjuvant therapy. A score >0.6 had up to 80%

## CHAPTER 1: INTRODUCTION

reduction in metastatic progression if adjuvant radiation therapy was received<sup>144</sup>. The PRO\_IMPACT was a multi-institutional study that showed decipher could significantly decrease decision conflict and patient anxiety. Decisions on adjuvant and salvage therapy were altered with the addition of a decipher score in 18% and 32% of cases, respectively<sup>145</sup>.

Another key finding was that decipher could also be performed on small amounts of genetic material like that obtained from biopsy and also including FFPE tissues. Decipher was tested on the biopsy material of 219 men who then went on to have RP to validate the findings, this gave HR = 7.3 and HR = 11 when moving from low-risk to high-risk on multivariate analysis<sup>145</sup>. In a second study, decipher was applied to the biopsy material of 57 men, who proceeded to undergo RP and also had long term follow up. Here decipher was capable of predicting metastatic progression as an individual predictor with AUC = 0.72<sup>146</sup>. This highlights its potential use in aiding decision making for primary treatment also, helping to identify those who are safe for active surveillance and those who should receive treatment more swiftly.

A limitation of decipher in the aid of primary treatment decision making, is that it relies on biopsies, and so carries the same limitations of a biopsy: PCa is often a multifocal disease and the lower-grade or lower decipher scoring foci could be picked up by biopsy, whilst the higher-scoring foci is missed. Leading to a less severe prediction occurring. Decipher has shown great promise as a second line treatment informer and clearly has a role here in PCa management.



**Figure 0.11** A schematic for oligonucleotide and two-channel microarrays. Both show RNA isolation from the cells of interest, followed by reverse transcriptase labeling to create cDNA from RNA and then hybridisation to array. In two channel arrays, cDNA from the normal cells and the “condition” cells are combined prior to hybridisation. Adapted from Vermeeren et al., 2011<sup>142</sup>.

Sequencing trumps microarrays with its ability to provide further information about specific unknown mutations. Mutations can be detected via microarrays; however, the probes must be designed to hybridise that specific mutation as the target gene, meaning the mutation must first be known. Cross-hybridisation problems in microarrays also mean that SNVs will be unable to be detected. Sequencing can also detect novel gene fusions. However, if you only require count data microarrays hold some advantages in comparison to sequencing: They are cheaper and the analysis of the data produced is easier. There are well-known analysis pathways to take, whereas, the best method for sequencing data analysis is still being investigated. Also, there is a lot of data available to the public from standardized platforms, which can be utilised for comparison with your own samples.

### 1.6.5 Mass Spectrometry

Mass spectrometry (MS) has shown great promise in proteomics and the identification of protein biomarkers. Proteomics is the large scale analysis of proteins including their structure and function; it provides information about the complex end products of a gene<sup>147</sup>. MS has positioned itself as one of the key technologies for the unbiased identification of cancer biomarkers<sup>148</sup>. Combining MS with liquid chromatography allows easy profiling of bodily fluids (samples which generally are less invasively obtained from patients) whilst MALDI-MS (matrix-assisted laser desorption/ionization) is useful for identifying biomarkers in FFPE tissues<sup>149</sup>.

MS works by bombarding molecules with electrons (ionizing) to create charged molecules and measuring their mass-to-charge ratio, by accelerating them and applying an electric or magnetic field. Ions of the same mass-to-charge ratio are deflected at similar amounts and can be detected via an electron multiplier. Results are available as a “spectra”, which can be correlated to previously known masses to identify atoms or molecules present in the sample.

Proteomics has an advantage over genomics, as it will be clear if a mutation is making a big difference to the protein, which can never truly be proven with genomics, just inferred<sup>148,150</sup>. However, MS suffers from insufficient sensitivity when detecting low-concentration biomarkers in a sample with a high-abundance of proteins, making depletion of abundant protein fractions and enrichment of biomarkers imperative to improving MS sensitivities<sup>150</sup>. MS also suffers from accuracy and reproducibility problems caused by software issues, meaning samples need to be run, typically, 10 times, increasing the amount of material required. These issues need to be addressed in order for Mass Spectrometry to develop as an efficient tool for biomarker discovery. Also, unlike other biomarker detection technologies, MS is limited to providing information of presence/absence and levels of the proteins cannot be determined unless targeted MS is performed. Targeted MS means the experiment must be focused on a

## CHAPTER 1: INTRODUCTION

small subset of protein targets to achieve their quantification, thus reducing the scope of proteins that can be quantified within the experiment<sup>151</sup>.

### **1.6.6 Fluorescent In-situ Hybridisation (FISH)**

FISH, a cytogenetic technique, can be used to detect chromosomal abnormalities; changes in chromosomal structure and numbers (including genomic deletion and fusion genes) can be observed when viewing cells or chromosome preparations upon a slide. Chromosomal abnormalities are common in many tumours, and Some of these abnormalities can be used for diagnostic and prognostic purposes<sup>152</sup>. FISH was being used to identify specific chromosomal regions and loci (chromosomal mapping) by the late 1980s. It works by labelling DNA with fluorophores, which emit light detectable by microscopy. FISH probes are capable of hybridising to DNA and RNA of circulating tumour cells (CTCs) and FFPE tissue sections that are fixed. This allows FISH to be useful for solid tumours as well as hematological cancers. Probes are designed for specific target sequences, and usually consist of cloned DNA sequences in the form of BACs, PACs, fosmids and cosmids, but can also be PCR products. The DNA probes can be tagged directly with fluorophores, or with biotin or DIG that can be bound post-hybridisation with streptavidin linked fluorophore or anti-DIG antibodies bound to fluorophore. Short DNA fragments are added to block repetitive DNA sequences and then the probes are applied to the cell preparations on a glass slide. Hybridisation requires approximately 12 hours followed by several wash steps in order to remove non-bound or partly bound probes. After which a microscope can be used to excite the dye and record the images for location and quantification of aberrations.

Improvements in fluorescent dyes and advances in microscopy and imaging allowed for MFISH (multi-fluorochrome assays), particularly SKY (Spectral Karyotyping). This new method allowed for entire metaphase spreads to be investigated using 24-colours, which showed the chromosomal origins of structural rearrangements<sup>153</sup>.

## CHAPTER 1: INTRODUCTION

FISH is a reliable, simple and specific assay for biomarker detection, and because of this, even though it is a low throughput method, it remains to be a cornerstone in genetic labs and even in clinical practice for the diagnosis, treatment stratification and prognosis of cancers. Whilst high-resolution molecular profiling techniques (microarrays and sequencing) are advanced in identifying novel chromosomal abnormalities, FISH remains a reliable validation method for any potential biomarkers identified<sup>154</sup>.

### **1.6.7 Immunohistochemistry (IHC)**

IHC is still commonly used in cancer diagnosis, and can validate biomarkers identified from other methods. Now that molecular, quantitative, global methods exist for novel biomarker identification, it is used much less to identify these but more to locate where in the cell the biomarker is and to validate its presence/absence in cancer tissues<sup>155</sup>.

For IHC, first tissue needs to be collected, fixed (commonly with paraformaldehyde) and often embedded in paraffin wax. The tissue is then sliced (4-40 $\mu$ m), mounted on slides and dehydrated with alcohol washes and cleaned with xylene before imaging via microscope. Blocking buffers are often used to reduce background staining. Positive and negative controls are required, a tissue known to express and a tissue known not to express the specific protein. Antibodies specific to the target antigen must be extracted from animals; the protein of interest is injected into the animal to elicit an immunological response producing the desired antibody. Therefore, this can make IHC time-consuming. Monoclonal or polyclonal antibodies can be used, targeting one epitope or multiple, respectively. Antibodies are often linked (using biotin) to reporter molecules. Reporter molecules can either be fluorophores or enzymes allowing fluorescence or chromogenic detection. There are two methods of antibody detection: Direct and Indirect. The direct method is where the labelled antibody directly binds the antigen, and the indirect uses an intermediate antibody to bind the antigen to which the labelled antibody binds.



## CHAPTER 1: INTRODUCTION

In order to compete with the new molecular methods, IHC will need to be quantitative<sup>156</sup>. IHC is specifically useful for the validation of protein biomarkers, similar to ELISA (section 1.6.8).

### **1.6.8 Enzyme-linked Immunosorbent Assay (ELISA)**

Similarly to IHC (section 1.6.7), ELISA is used for protein detection. ELISA works by using enzyme linked antibodies to capture antigens, and colour changes from the enzyme binding its substrate provide detectable signals, which are proportional to the amount of antibodies bound to the antigens present. IHC and ELISA have their own advantages; ELISA is fully quantifiable and easily standardised with quality assured measurements obtained. IHC is at best semi-quantitative but allows insight into tissue heterogeneity and can be performed on both frozen and paraffin-embedded tissues (section 1.6.7)<sup>157</sup>.

There are three main types of ELISA: Indirect, Sandwich and Competitive, all use 96-well microtitre plates as the immobilising surface, allowing moderately high-throughput investigations. ELISA is a versatile and robust tool and so ELISA is often used for validation of biomarkers<sup>150</sup>.

ELISA only allows the detection of a single antigen and often requires a large amount of sample. This along with the narrow dynamic detection range means it not useful in biomarker discovery. ELISA is costly and its quality is dependent upon antibody quality, users skill and experience, and shows problems with accuracy and reproducibility. ELISA's downfalls mean it is less useful at biomarker validation when there are multiple proteins, which is often the case.

### **1.6.9 Methylation Assays**

Epigenetic gene regulation, such as methylation and histone modification, is an important factor in normal development and disease states like cancer<sup>158</sup>. These are modifications to our gene expression that is not encoded for by the DNA, but inherited

## CHAPTER 1: INTRODUCTION

mitotically<sup>159, 160</sup>. Hyper-methylation of promoter regions is commonly seen in cancers to knock-down the expression tumour suppressor genes<sup>161</sup>. Methylation is the addition of a methyl or hydroxymethyl group to the C5 position of cytosine, which occurs at or around CG dinucleotide regions (known as CpG islands and shores)<sup>162</sup>. Methylation is known to aid in cell cycle regulation and cellular differentiation processes<sup>163,158</sup>. The role of DNA methylation has been well established in many cancers including PCa<sup>163, 164, 165</sup>. Hyper-methylation of several genes, including GSTP1, is commonly observed during the transition between intraepithelial neoplasia to carcinoma<sup>166</sup>. Hyper-methylation detection has shown promise as biomarkers for the diagnosis and prognosis of cancers.

There are many methods for the identification of methylated sites, which method you chose can be based on many things, but importantly is the biological question you are asking: There are different methods available for whole genome methylation profiling, identifying regions of differential methylation status, or for determining the methylation status of specific genes of interest<sup>167</sup>. Other factors to include when choosing a method are the amount and quality of the sample, the sensitivity and specificity requirements for the experiment, robustness and simplicity of method, and its bioinformatics analysis, as well as the availability of specialist equipment and overall experiment cost<sup>167</sup>.

### **1.6.10 Supervised and Unsupervised Analyses**

Due to the developments in genomic technology more and more biological data is being developed that needs to be analysed; to identify patterns and trends and understand what the data means to the biological question. The application of statistics on such data can be called statistical learning or machine learning, and can roughly be separated into two categories: supervised and unsupervised analyses<sup>168</sup>. These can be referred to as classification and clustering, respectively<sup>169</sup>.

For supervised learning, or “classification” of observation  $x$ , an observation with multivariate  $p$  dimensions (also called features) and associated with class  $c$ . The

## CHAPTER 1: INTRODUCTION

purpose is to “learn” a mathematical function that when evaluated with the input  $x$  provides a prediction of its class  $c$ . In general practice data is subset into training and test datasets. The training set is used to “set” the mathematical function to correctly predict the class for each observation provided. This function, with the parameters set from the training data is then applied to the test data to observe its ability to correctly classify the data without bias<sup>169</sup>. Examples of supervised machine learning are generalized linear models (glm) (section 2.6.1), probit regression (polr) for ordered multivariate models, random forest (section 2.6.3). These methods can be accompanied by a shrinkage method to reduce over fitting and thus improve predictability; examples of such methods are Lasso (section 2.6.2) and Step (section 2.6.4).

Unsupervised analysis, or “clustering” can also be referred to as class discovery. A key difference between unsupervised analysis and supervised analysis is a lack of training set for the former and thus no cross-validation. A second important difference is that clustering algorithms are set using optimality criterion and there is a lack of guarantee that the global optimal solution is found, and therefore a heuristic approach is often taken. A choice of a) features to be used, b) similarity metric, and c) algorithm needs to be made for many methods<sup>169</sup>. Unsupervised learning can be further partitioned into hierarchical clustering (section 2.5.2), (which can then be subdivided into agglomerative and divisive) and partitioning. Hierarchical will cluster data into a tree like feature and then to achieve a desired number of clusters one can cut the dendrogram at a desired height. However, partitioning generally requires the user to specify the number of clusters prior to clustering<sup>169</sup>. Examples of partitioning are k-means clustering (section 2.5.3), principal component analysis (section 2.5.1), and latent process decomposition (section 2.5.5).

### **1.7 Summary and Aim**

#### **1.7.1 Summary**

Whilst the introduction of the PSA test has decreased mortality from PCa, the increased incidence rate that can also be attributed to it comes with problems of over-diagnosis and over-treatment. Highlighting the need for additional biomarkers for the diagnosis of PCa. A need for biomarkers for hormone therapy response prediction, BCR prediction, further treatment stratification, and prognosis were also highlighted.

The heterogeneity of PCa means that there have been a lot of potential biomarkers discovered, but also that they are not always consistent in the tumours. Meaning a limited number are capable of being used for the diagnosis and prognosis of PCa. However, combinations of biomarkers in a panel could be of great clinical use. The utility of urine in PCa biomarkers is well established via the PCA3 test (section 1.4.2) and the role of exosomes in cancer development and metastasis (section 1.5) has highlighted a resource to be investigated.

The development of NGS technologies and the continuous advancements in sequencing technologies are making it possible to investigate a large number of genes across a large number of samples, at continuously decreasing costs. Sequencing is an important technology for the discovery of novel biomarkers, as it is capable of identifying expression changes and mutations at high-throughput. The reducing costs of sequencing are closing the gap between data production costs and data processing costs, it is said that there may come a time when processing the data will become more costly. Bioinformatic analysis of the data is still under on-going development to identify the optimal pathways for analysis. Currently, cheaper methods for high-throughput expression analysis (microarrays and Nanostring) still hold a firm place within biomarker discovery. Whilst, mass Spectrometry for proteomic biomarker discovery holds massive potential, however there are issues still to overcome.

## CHAPTER 1: INTRODUCTION

Older techniques of biomarker discovery hold great sensitivity but are at considerably low-throughput, making them very good for validation and clinical detection after a few potential biomarkers are selected from higher-throughput methods. These include: FISH for gains/losses, rearrangements and chromosomal instability investigations, IHC and ELISA for the validation of particularly proteins and to see where in the cell these biomarkers are gained to or are lost from and PCR-based methods for confirming mutations in the biomarkers.

Knuutila *et al.*, compared NGS, aCGH, FISH, PCR and IHC methods for specific biomarker analysis of FFPE tumour tissues. Their conclusions suggest that NGS has the potential to replace all other methods tested for the analysis of tumour biomarkers, especially as the reducing costs and required sample material decreases to that near of FISH or PCR. NGS allows the investigation of mutations, gene fusions and copy number changes in one single analysis<sup>170</sup>. However, NGS has not currently reached the position where it is commonly used in clinical practice.

### 1.7.2 Aims & Objectives

PCa diagnostics and prognostics currently lack specific and sensitive clinical biomarkers and treatment is not well individualised. The PCA3 test highlights the utility of urine in PCa diagnostics and prognostics. The aim of our work is to interrogate PCa patient's urine samples, mostly the exosomal fraction to identify novel biomarkers or sets of biomarkers to aid in PCa management. My objectives are as follows:

**O1:** To determine whether RNA expression from urine extracellular vesicles in prostate cancer patients are a viable target for the development of biomarkers through the use of Nanostring technology.

**O2:** To determine an optimal combination of probes to predict cancer presence and aggression in prostate cancer patients.

**O3:** To determine whether an optimal combination of probes can predict response to hormone therapy treatment.

## CHAPTER 1: INTRODUCTION

**O4:** To evaluate the differences between urine fractions (extracellular vesicles and cell sediment) and determine whether cell sediment can be used to predict cancer presence and aggression in prostate cancer patients.

Below are described more detailed aims for each chapter.

### ***3.1.1.1 Chapter 3: NanoString Data Analysis 1: The Pilot Study***

This chapter encompasses the analysis of the pilot study of samples sent to NanoString to investigate exosomal RNA expression level changes of 57 target sequences. The RNA was extracted from the EV fraction of urinary samples collected at the NNUH as part of the Movember study. The aims were to primarily determine if the transcript content of urinary exosomes contained any PCa derived transcripts and if transcript level could be utilised for risk stratification. Also, it was important to investigate if NanoString was a suitable method for obtaining expression data from these cDNA-amplified samples and to determine suitable methods for analysis.

### ***3.1.1.2 Chapter 4: NanoString 2 Analysis: The Movember GAP1 Project***

A second analysis for the Movember study. RNA was extracted from the EV fraction of urinary samples that were collected from multiple centres (NNUH, Norwich, St James' Hospital, Dublin, Royal Marsden Hospital, London, and Emory Healthcare, Atlanta). 864 samples were sent to NanoString for the quantification of 167 transcripts. The aims were to primarily identify optimal models capable of predicting PCa and to risk-stratify PCa without the need for biopsy. Models were built to answer four important clinical questions: 1) determine which samples were from PCa and which were from samples with no evidence of Ca 2) determine which samples were from high-risk PCa only and which were from samples with no evidence of cancer 3) determine if there was a trend in expression that corresponds to a trend in risk category (CB>L>I>H) and 4) determine if there was a trend in expression that corresponds to a trend in patient type (CB>Ca>Metastatic cancer).

### ***3.1.1.3 Chapter 5: Response to treatment***

Many cancers have benefitted from treatment stratification due to expression of certain genes, however not yet PCa. With hormone therapy (HT) it is known that patients will inevitably progress to castration resistant prostate cancer (CRPC). How long each patient will last on HT varies widely from months to years. It is our aim to use the NanoString data of the advanced patients in the pilot study (n = 32) to see if a significant predictor of early progression in patients on HT can be built and whether this predictor improves on current clinical information collected (e.g. PSA, Gleason score and bone scan). The NanoString 2 data can then be used for validation of this predictor.

### ***3.1.1.4 Chapter 6: Analysis of Cell Fraction and comparison with exosomal fraction***

The use of RNA extracted from EV fractions and cell sediment fractions were used to compare the transcriptome profiles from PCa patients and controls (taken from patients with no evidence of cancer (CB)). The aim was to identify if both fractions contained similar expression profiles of genes and if either contained higher levels of prostate or PCa associated transcripts. The fraction with the highest level of these transcripts is likely to be a better source of material for PCa diagnosis and risk stratification. Data from microarray of samples collected from NNUH, Norwich and Royal Marsden Hospital, London was used.

Secondly, I am to use NanoString data from cell sediment fraction derived transcripts (collected only from NNUH, Norwich) to identify optimal models to answer the four important clinical questions asked of the EV derived data (Chapter 4).

# 2

## Materials and Methods

### **2.1 Sample Collection and Processing**

**Overview:** Urine samples were collected from patients attending hospital clinics. Extracellular vesicle (EV) RNA was harvested by urine microfiltration (Section 2.1.2). EV and cell pellet RNA was extracted (Section 2.1.3), converted to cDNA and amplified as cDNA (Section 2.1.4), ready for NanoString expression analysis (Section 2.1.5).

Not all the procedures in this section were performed by me but were included in this thesis as essential information relative to the study.

#### **2.1.1 Sample Collection**

Note: The procedures in this section were not performed by me.



## CHAPTER 2: MATERIALS AND METHODS

Urine samples were collected in a 30ml Universal tube (Sterilin) from urology clinics at the Norfolk and Norwich University Hospital (NNUH, Norwich, UK), St. James Hospital (Dublin, Republic of Ireland) and from primary care and urology clinics of Emory Healthcare (Atlanta, USA), between 2012 and 2015. Most samples were collected as first void post-DRE but a few matching pre-DRE samples were collected for comparison (these were labelled as such). All samples were collected from treatment naïve patients. Control samples were collected at a micro-haematuria clinic at the NNUH, again first void post-DRE in a 30ml Universal tube. Microvesicular RNA was harvested by ultracentrifugation (section 2.1.2), extracted (section 2.1.3), converted to cDNA and amplified (section 2.1.4). RNA from the cell pellet was also processed, using either the Qiagen Allprep DNA/RNA mini kit cat no: 80204 or RNeasy micro kit cat no: 74004 according to manufacturer's instructions).

The lab also had access to urine samples collected as part of the active surveillance prospective study at the Royal Marsden Hospital NHS Foundation Trust (RMH) between 2009 and 2012. The active surveillance prospective study collected samples, first void post-DRE, specifically from men with untreated, low-risk prostate cancer. Low-risk PCa defined as having clinical stage T1/T2a, Gleason 3+3 (or 3+4 of older than 65), PSA<15 and <50% positive cores. Three of these samples were collected pre-DRE from post-radical prostatectomy patients for comparison. Microvesicular RNA was harvested as above.

The study was given favourable ethical opinion by the NRES Committee East of England – Norfolk on 21<sup>st</sup> August 2014 under the study title “Urine biomarkers for detecting prostate cancer”. Ethics was approved to Dr Marcelino Yazbek Hanna of NNUH with REC reference: 12/EE/0058 and IRAS project ID: 96199.

### **2.1.2 Micro-filtration harvesting of Urine Extracellular Vesicles**

Note: The procedures in this section were generally not performed by me (I performed these procedures on ~20 samples).

## CHAPTER 2: MATERIALS AND METHODS

Urine samples were processed within four hours of samples collection. Urine was centrifuged at 1200g for 5 mins, and then the supernatant was transferred to a 50ml tube and re-centrifuged at 2000g for 5 mins. Supernatant was decanted and then filtered through a 0.8µM filter (Millipore), transferred to an Amicon Ultra-15 100KDa MWCO microfiltration unit and spun at 3400g r/t for 15 mins or until the volume was reduced to below 500µL. PBS (15ml) was added to the sample followed by further centrifugation until the volume containing the EVs was reduced to 200µL. Transmission electron microscopy (TEM) was performed to confirm the presence of EVs.

### **2.1.3 Qiagen RNA Extraction**

Note: This section was generally not performed by me (I performed this step on ~20 samples).

The Qiagen Micro RNA RNeasy kit was used for RNA extraction from EVs and cell pellet as per the manufacture's manual. 700µL of buffer RLT was added to the cell pellet or EV samples. The cell pellet/RLT mix then had an extra step, which was to disrupt the cells using a QIAshredder spin column for 2 mins at full speed (~12,000 rpm) in a microfuge. From this point onwards the cell pellet and EVs were treated the same. 70% ethanol was added and the mixture pipetted into a MinElute spin column and centrifuged in a microfuge (15 seconds, >10,000rpm). 350µL of buffer RW1 was added to the MinElute spin column before re-spinning (15 seconds, >10,000rpm). Then 80µL of Qiagen DNase mix I was directly applied to the membrane and left to stand at room temperature for 15 mins to complete DNA digestion. The wash step with RW1 was then repeated followed by the addition of 500µL of buffer RPE and re-spun (15 seconds, >10,000rpm). 500µL 80% ethanol was added and then spun (2 mins, >10,000rpm). To dry the membrane the column was spun for a further 5 mins with an open lid. The column was transferred to a fresh collection tube, and the RNA was eluted with 14µL of RNase free water and centrifuged (1

## CHAPTER 2: MATERIALS AND METHODS

minute, 12,000 rpm) in a microfuge. Nanodrop and Bioanalyzer were used to confirm that RNA was of a good quality.

### **2.1.4 Nugen Amplification of RNA as cDNA**

Note: I performed the amplification step for 286 samples.

Amplification was performed using the Nugen Ovation picoSL WTA V2 kit as per the manufacture's instructions<sup>171</sup>. The kit works via the following mechanisms: Firstly, the first strand of cDNA was generated using a DNA/RNA chimeric primer mix (containing a mix of random and oligo dT primers such that priming occurs throughout the whole transcript) and reverse transcriptase (RT). The RT extends the 3' end of the DNA for each primer resulting in a cDNA/mRNA hybrid containing a unique RNA tag sequence known as the SPIA tag at the 5' end of each cDNA strand. The SPIA tag was used for a priming site for the SPIA process.

Secondly, fragmentation of this cDNA/mRNA complex was required to provide priming sites for RNA polymerase to synthesise a second cDNA strand. This includes DNA complementary to the 5' SPIA tag and results in a double stranded cDNA with a DNA/RNA heteroduplex, which corresponds to the SPIA tag. Finally, strand displacement amplification occurs that uses a DNA/RNA chimeric primer (SPIA primer), DNA polymerase and RNase H in an isothermic assay. RNase H removes the RNA part of the heteroduplex SPIA tag allowing the SPIA primer to bind. DNA polymerase can then synthesise from the 3' end of the primer displacing the existing forward strand with new cDNA. Priming with the chimeric SPIA primer then in turn makes a new heteroduplex SPIA tag, which becomes the new substrate for RNase H and can initiate the next round of cDNA synthesis. These last few steps were repeated in a highly processive manner allowing rapid accumulation of  $\mu\text{g}$  of amplified cDNA from  $\text{ng}$  of total RNA.

The actual process was as follows: samples were diluted with RNase free water to ensure all contain 20ng of total RNA in a PCR tube. The first strand synthesis primers were added

## CHAPTER 2: MATERIALS AND METHODS

(2 $\mu$ L) to each sample and they were heated to 65°C for 2 mins. The first strand buffer and enzyme were pre-mixed, 2.5 $\mu$ L and 0.5 $\mu$ L per sample, respectively. 2.9 $\mu$ L of this mix was then added to the sample tubes and samples were returned to the thermal cycler for program 2: 4°C for 2 mins, 25°C for 30 mins, 42°C for 15 mins, 70°C for 15 mins and hold at 4°C. Second strand synthesis required second strand buffer and enzyme to be pre-mixed; 9.7 $\mu$ L and 0.3 $\mu$ L per sample, respectively. 9.5 $\mu$ L of this mix was added to each sample, mixed via pipetting (5x) and returned to the thermal cycler for program 3: 4°C for 1 minute, 25°C for 10 mins, 50°C for 30 mins, 80°C for 20 mins and hold at 4°C. The cDNA must then be purified using the magnetic beads provided in the NuGEN kit; 32 $\mu$ L of the beads were added to each sample and mixed via pipetting (10x) and were incubated at room temperature for 10 mins.

Following this the samples were placed into the 96 well magnet and were incubated at room temperature for a further 5 mins. Using long thin pipette tips, 45 $\mu$ L of buffer was removed as the beads (that were attached to the cDNA complexes) were pulled to the tube side via the magnet. The tubes were then washed with 70% ethanol (200 $\mu$ L) three times and left to air dry at room temperature (roughly 25 mins but until there was no liquid left in the tubes). The last step was then the SPIA amplification; where the SPIA buffer, the SPIA primer mix and the SPIA enzyme were pre-mixed in order; 20 $\mu$ L, 10 $\mu$ L and 10 $\mu$ L per sample, respectively. 38 $\mu$ L of the SPIA mix was added to each sample and the samples were returned the thermal cycler for program 4: 4°C for 1 minute, 47°C for 75 mins and 95°C for 5 mins. At a different bench, the PCR tubes were returned to the magnet for 5 mins, and the liquid that contained the amplified cDNA was collected.

TE was added to a final concentration of 0.2xTE and yields determined via Nanodrop or Qubit.

## CHAPTER 2: MATERIALS AND METHODS

### 2.1.5 NanoString

The NanoString nCounter gene expression system technology uses 2 probes; a capture and a reporter probe<sup>115</sup>. The probes were designed to have a complementary sequence to the specific transcripts we wished to study. Each transcript reporter-probe had a distinct string of fluorescent coloured beads attached<sup>115</sup>. Up to ~800 different bead string combinations are available, and so up to 800 different transcripts can be detected in one analysis.

The probes were hybridised to the complementary nucleic acid sequences in each sample, forming a tripartite complex of the 2 probes and target mRNA or cDNA. The complexes were then pulled down and immobilised onto a capture surface, unbound sequences were washed away and an electric field was passed across the surface to stretch out the nucleotide and bead complexes. The bead-complexes were then imaged and the number and type of each string of coloured beads counted. This provided a direct measure of RNA or cDNA counts per transcript<sup>115</sup>. Twelve samples loaded onto the NanoString machine at a time (in each cartridge).

### 2.1.6 PCR (Polymerase Chain Reaction)

Note: I performed PCR detection of *TMPRSS2:ERG* fusions for 113 samples.

*TMPRSS2:ERG* fusions were detected by primary PCR and confirmed with a secondary PCR that used nested primers. A master mix was made using the following components (volumes provided for a single PCR): 2.5µL 10x PCR buffer, 1µL 50mM MgSO<sub>4</sub>, 0.5µL 10mM dNTP mixture, 0.5µL Primer 1 (10µM), 0.5µL Primer 2 (10µM), 0.1µL Platinum Taq (Thermo Fisher) and 19µL HPLC H<sub>2</sub>O. Primer 1: CAGGAGCGGAGGCGGA (*TMPRSS2* exon 1 Forward). Primer 2: GGCGTTGTAGCTGGGGGTGAG (*ERG* exon 6 Reverse). The master mix was pipetted into a clean 0.25ml tube and 1µL of the template cDNA was added. PCR conditions were as follows: 94°C for 30 seconds, followed by 35 cycles of 94°C for 20 seconds to denature and 68°C for 60 seconds to extend. A second master mix was created using the same reagents but using 0.5µL of

## CHAPTER 2: MATERIALS AND METHODS

the nested primer 1 and 0.5µL of the nested primer 2. Nested primer 1: GGAGCGCCGCTGGAG (*TMPRSS2* exon 1 nest Forward) and nested primer 2: CCATATTCTTTCACCGCCCACTCC (*ERG* exon 6 nest Reverse). This master mix was aliquoted and 0.25µL of the primary PCR was added. PCR conditions were as above but with a 66°C annealing temperature instead of 68°C. The resulting amplification products of the primary PCR and the nest PCR were run in adjacent wells on a 2% agarose gel with a 100bp DNA ladder (New England Biolabs (N3231L)) to determine product sizes (Table 2.1) and thus infer which of the *TMPRSS2:ERG* fusions were present in each sample.

**Table 2.1 PCR product sizes for *TMPRSS2*\_exon1 (T1) and *ERG*\_ex6 (E6) PCR primers (nests are 139bp smaller than primaries)**

	<i>Primary</i>	<i>Nest</i>
<i>T1/E4</i>	<b>596</b>	<b>457</b>
<i>T1/E5</i>	<b>379</b>	<b>240</b>
<i>T1/E6</i>	<b>227</b>	<b>88</b>
<i>T2/E2</i>	<b>856</b>	<b>717</b>
<i>T1/E3, -, 5, 6</i>	<b>465</b>	<b>326</b>
<i>T1/E2, 3, 4, -, 6</i>	<b>661</b>	<b>522</b>
<i>T2/E5</i>	<b>450</b>	<b>311</b>
<i>T3/E4</i>	<b>891</b>	<b>752</b>
<i>T4/E5</i>	<b>760</b>	<b>621</b>
<i>T5/E4</i>	<b>1098</b>	<b>959</b>

### **2.2 Clinical Data Collection**

Note: I completed part of the clinical data collection.

Clinical data was collected for NNUH samples from a number of different NHS databases such as ICE (the NNUH database), Somerset (the NNUH cancer database) and also from the patient's forms completed for the study within the clinic. Information from the patient's forms were manually typed into an Excel sheet and uploaded to a pseudo-anonymised online database for the Movember project. A clinical NHS colleague and I updated and checked over clinical information for the majority of the samples, this included information

## CHAPTER 2: MATERIALS AND METHODS

such as age, initial PSA reading, Gleason score and further biopsy information, scan conclusions, prostate volume, family history, health altering habits, general health and current medications as well as subsequent information (ensuing PSA readings for example). Clinical data for samples from other centres were provided by them and uploaded into the Movember database.

### **2.3 NanoString Pre-processing**

#### **2.3.1 Normalisation**

The NanoString output data file provides the nCount data for 6 spiked non-human positive control probes and 8 non-human negative control probes for each of the samples being analysed. The six positive control probes matched to spiked-in RNAs and was used to calculate a normalisation factor (NF): the average nCount for each samples' positive controls were calculated and this number was divided by the sum of all samples' averages. Each nCount value was then multiplied by the sample-specific NF. This results in a shift of all samples so that the means of the positive controls was identical across samples. Background correction and background subtraction using the negative controls was found to be inappropriate for this data.

##### **2.1.1 Normalisation by *KLK2* and *KLK3***

Normalisation using *KLK3* and *KLK2*, separately, was conducted as follows. For *KLK2*, a ratio was determined (Equation 2.1) and then applied to the data, this data was referred to as *KLK2* ratio data.

**Equation 2.1 *KLK2* ratio normalisation, similar to the normalisation of *PCA3* by *KLK3* in the *PCA3* test**

$$\left( \frac{(x_{ij})}{(\bar{x}_{KLK2})} \right) * 1000$$

## CHAPTER 2: MATERIALS AND METHODS

Additionally, for *KLK2* and *KLK3* an adjustment normalisation was conducted using the median and IQR (Equation 2.2). For this data, any samples observing low *KLK2* or *KLK3* levels, respectively were removed from the data set prior to adjustment. The threshold for “low” expression was determined using a density plot and the Brent method to find the minima of the curve. For *KLK2*, and *KLK3* the same nineteen samples were identified and removed for low expression. As well as removing low Kallikrein expression samples, six CB samples that had high *TMRPSS2:ERG* expression were also identified and removed. Samples with high *TMRPSS2:ERG* expression were again identified through density plots and the Brent method.

**Equation 2.2 Kallikrein adjustment of data using median and IQR. Where  $i$  is the sample and  $j$  is the transcript.**

$$\left( x_{ij} - \frac{\text{median}(x_j)}{\text{IQR}(x_j)} \right) * \text{IQR}(KLK) + \text{median}(KLK)$$

### 2.3.2 Normalisation by housekeeping genes

Five previously identified housekeeping transcripts were included in the NanoString1 pilot study: *ALAS1*, *B2M*, *HPRT*, *GAPDH*, and *TBP*. *RPLP2* was added in NanoString2. Tukey tests (section 2.4.7) were used to identify transcripts that were not significantly different between any clinical category ( $p < 0.05$ ). ANOVA (section 2.4.6), variance and IQR (section 2.4.8), and Pearson’s correlation (section 2.4.3) were also utilised, to identify novel transcripts to use for housekeeping purposes. In NanoString2 EV data, *RPLP2* and *GAPDH* were selected to normalise the data, whereas for the NanoString2 cell data, *RPLP2* and *TWWAST1* were selected.

For each sample, the mean of the two transcripts was calculated, as well as the mean of those means across samples. Each sample was then multiplied by a normalisation factor (ratio of the mean of means with the individual sample mean).



## CHAPTER 2: MATERIALS AND METHODS

### 2.3.2 NanoStringNorm and NanoString QC Pro

The quality of the normalisation was evaluated using the R packages NanoStringNorm<sup>172</sup> and NanoString QC Pro<sup>173</sup>.

#### 2.3.2.1 NanoStringNorm

NanoStringNorm<sup>174</sup> investigates the normalisation of the data as well as identification of samples and transcripts that were outliers. The first test performed by NanoStringNorm was to plot the mean verses the standard deviation (SD) with a Loess curve of best fit. Positive controls and potential housekeeping probes should have high means and low SD, whilst negative controls should have low means and low SD. Batch effects and potential confounding were also tested for using sample summary features, including mean, SD, proportion of missing (0 counts) or positive/negative control counts. These features were plotted independently by NanoStringNorm, where the location of the point relative to the horizontal line shows how different it was from the others and the size of the (green) dot was proportional to the level of its significance. Orange dots were not significant. Potential influencing outliers were identified by looking into the normalisation factors: if the normalisation parameters extended beyond 100% difference from the mean, it was flagged as a potential outlier.

#### 2.3.2.2 NanoString QC Pro

NanoStringQCPro<sup>173</sup> (an R library) was conducted to check the quality control of the NanoString data, specifically looking at the control probe metrics and count probe metrics (similar to NanoStringNorm) but additionally looks at other metrics. The field of view (FOV) was a discrete area of each lane being imaged by the ncounter® digital analyser. Within the FOVs, bubbles and insufficient oiling can make unsuccessful imaging attempts. A low ratio between successful and unsuccessful attempts can be indicative of low imaging performance. NanoStringQCPro highlights any samples with less than 80% successful imaging attempts. If the binding density was too high in a sample, there can be overlapping

## CHAPTER 2: MATERIALS AND METHODS

of the barcodes, which leads to errors in correctly imaging the number of probes. According to NanoString a binding density of less than 0.5 and higher than 2.25 can lead to these errors. NanoStringQCPro flags samples that have binding densities outside of the recommended thresholds.

Positive controls were spiked into each NanoString experiment, they should show linearity with positive control A having highest values, down to positive control F with the lowest. Control range and Interquartile range (IQR) were also examined. The counts were also examined; any samples with unusually low counts were flagged using cutoffByMMAD to identify the threshold. This was based on the median of the data and the upper and lower thresholds were counted using  $\text{median}(x) - d * \text{mad}(x)$  and  $\text{median}(x) + d * \text{mad}(x)$ , respectively (where  $d$  was a scalar).

### 2.3.3 Log and Square-root Transformation

Sometimes, to obtain a more normal distribution of the data, it can be useful to transform the data. Many inferential statistical tests assume that the data was of normal distribution and violating these assumptions can cause an increase in both type 1 and type 2 errors. For regression-based models, the relationship between input and output variables should be approximately linear (so the input variables have a normal distribution and the output has constant variance, thus the variance of output variables was independent of input variables). Two transformations that have been used in this project were log transforming the data and square root transformation. Square root transformation has been shown to be appropriate for transforming count data<sup>175</sup>. However, square-root transformation of data has its drawbacks; if your data contains both values greater than 1 and values between 0 and 1, these two types of values will be treated differently.

### 2.3.4 ComBat

Batch effects occur in many high-throughput experiments, they can be caused due to laboratory conditions, reagent lots and personnel differences. ComBat was determined to be

## CHAPTER 2: MATERIALS AND METHODS

the best performing of six methods for removing batch effects in microarray data<sup>176</sup>. The ComBat function is an empirical Bayes method, where location and scale model adjustments are made as follows:

$$Y_{ijg} = \alpha_g + X\beta_g + \gamma_{ig} + \delta_{ig}\varepsilon_{ijg}$$

Where  $\alpha_g$  is the overall gene expression, X is a design matrix for sample conditions,  $\beta_g$  is the vector of regression coefficients corresponding to X,  $\varepsilon_{ijg}$  is the error terms (which are assumed to follow a normal distribution) with expected value of 0 and variance  $\sigma_g^2$ . The  $\gamma_{ig}$  and  $\delta_{ig}$  represent the additive and multiplicative batch effects of batch i for gene g.

The adjusted data is then given by:

$$Y_{ijg}^* = \frac{Y_{ijg} - \hat{\alpha}_g - X\hat{\beta}_g - \hat{\gamma}_{ig}}{\hat{\delta}_{ig}} + \hat{\alpha}_g + X\hat{\beta}_g$$

Where  $\hat{\alpha}_g, \hat{\beta}_g, \hat{\gamma}_{ig}$  and  $\hat{\delta}_{ig}$  are estimators for the parameters  $\alpha_g, \beta_g, \gamma_{ig}$  and  $\delta_{ig}$  based on the above model. The ComBat function of the sva R package was used with R version 3.2.1.

### **2.4 Basic Statistical Tests**

Basic statistical functions used and described below were part of the R stats package and were used with default settings, under R version 3.2.1.

#### **2.4.1 Mann-Whitney U test (Wilcoxon Rank Sum test)**

The Mann-Whitney U test was a non-parametric log-rank test capable of identifying differential expression of genes between two different states, for example, cancer vs. non-cancer. The test works by assigning a rank to each individual value from 1 to  $n$  (where  $n$  was the number of samples) and 1 was assigned to the smallest value. It then compares the sum of the ranks in the first group ( $R_1$ ) to the expected sum of the ranks given the sample size of group 1 and then the sum of the ranks in the second group ( $R_2$ ) was compared to the expected sum of the ranks given the sample size of group 2 (these values were considered

## CHAPTER 2: MATERIALS AND METHODS

$U_1$  and  $U_2$ , respectively, see Equation 2.3). The smallest of these numbers was then used to calculate the significance.

**Equation 2.3 Mann-Whitney U test**

$$U_1 = R_1 - \frac{n_1(n_1 + 1)}{2}$$

One advantage was that more accurate results than a  $t$ -test were obtained when used on data with a non-normal distribution<sup>177</sup>.

### 2.4.2 Spearman's Correlation

Spearman's rank correlation coefficient was a non-parametric test of statistical dependence between two sets of data, most commonly two variables. It measures the relationship between these variables providing a value between -1 and 1, where 1 or -1 means complete dependence, whilst 0 means that no dependence was observed. Spearman's correlation uses the rank of the variables rather than exact values (as used in Pearson's correlation). The covariance of these ranks was divided by the standard deviation of the ranks also (Equation 2.4). Here,  $d_i$  was the difference in ranks for variables  $x$  and  $y$ ,  $r_s$  was the notation for the coefficient for a sample statistic and  $n$  was the number of samples. Spearman's correlation was preferred over Pearson's correlation typically when one of the variables was ordinal and the other was continuous or if the relationship was non-linear<sup>178</sup>.

**Equation 2.4 Spearman's Correlation**

$$r_s = 1 - \frac{\sum_{i=1}^n d_i^2}{n(n^2 - 1)}$$

### 2.4.3 Pearson's Correlation

Pearson's product moment correlation coefficient was calculated in a very similar method to that of Spearman's correlation in that the covariance of the two variables was divided by the standard deviation of those variables. The key difference was that the exact values were used instead of their ranks (Equation 2.5). Here the correlation coefficient was noted by  $r$  and  $x_i$  and  $y_i$  were the  $i$ th individuals of  $x$  and  $y$  variables. Pearson's correlation was

## CHAPTER 2: MATERIALS AND METHODS

typically used when both variables were continuous, normally distributed (as extreme values can bias the strength of a relationship), and the tested relationship was linear<sup>178</sup>.

### Equation 2.5 Pearson's Correlation

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{[\sum_{i=1}^n (x_i - \bar{x})^2][\sum_{i=1}^n (y_i - \bar{y})^2]}}$$

### 2.4.4 Pearson's Chi-Squared

This statistical procedure was typically used to identify if the frequency distribution of events was independent from the labels assigned to the event. It can be used to suggest if two groups of variables were related or not, for example in clustering, to see if the clusters were significantly related to the clinical category, a frequency distribution table can be produced. To calculate what frequencies were likely to occur from chance, the number of observations ( $O_{ij}$ ) was divided by the number of cells in the table, this gives what was known as the theoretical frequency ( $E_{ij}$ ). This can then be used to calculate the test statistic (Equation 2.6) and with  $n-1$  degrees of freedom, the  $p$ -value can also be determined<sup>179</sup>.

### Equation 2.6 Pearson's Chi Square test

$$X^2 = \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

### 2.4.5 Welch $t$ -test

The Welch  $t$ -test was a parametric test to measure how the means and variance of two groups differ in normally distributed data, where the variances of the two populations were assumed to be non-equal. The mean of the data points in-group A and B, along with the squared sums ( $\sum x$ )<sup>2</sup> and also the sum of the squares  $\sum(x^2)$  were used to calculate a  $t$  value (Equation 2.7). This provides a  $t$  value, which was comparable to values designated using different degrees of freedom (dependent on the number of samples in your two groups of data). Combining the  $t$ -value with the relevant degrees of freedom (sum of the variables in each group minus 2) yields a  $p$ -value.

## CHAPTER 2: MATERIALS AND METHODS

**Equation 2.7** The Welch  $t$ -test, comparing the mean and standard deviation between two sets of data to conclude if they were significantly different from each other

$$t = \frac{\bar{X}_A - \bar{X}_B}{\sqrt{\left[ \frac{\left( \sum A^2 - \frac{(\sum A)^2}{n_A} \right) + \left( \sum B^2 - \frac{(\sum B)^2}{n_B} \right)}{n_A + n_B - 2} \right] \cdot \left[ \frac{1}{n_A} + \frac{1}{n_B} \right]}}$$

### 2.4.6 ANOVA – Analysis of Variance

Analysis of Variance (ANOVA) was a procedure used to analyse the differences among group means. In this work, it has a similar function to the  $t$ -test but allows the analysis of more than two subgroups. Firstly, the mean sum of squares within each group,  $MSS_w$  (Equation 2.8) and the mean sum of squares between the groups,  $MSS_B$  (Equation 2.9), was calculated. The ratio of these then provides the test statistic,  $F$  ( $F = \frac{MSS_B}{MSS_w}$ ). Combining the  $F$  statistic with the degrees of freedom allows a  $p$ -value of significance to be determined. There were two degrees of freedom to calculate in ANOVA,  $df_B = k-1$  and  $df_w = n-k$ , where  $n$  was the total number of samples and  $k$  was the total number of groups.

**Equation 2.8** Mean sum of squares within each group of data, where  $n$  was the total number of samples,  $k$  was the total number of groups,  $g$  was the value and  $G$  was all of the values across all groups.

$$MSS_w = \frac{\sum_{g \in G} (x - \bar{x}_g)^2}{n - k}$$

**Equation 2.9** Mean sum of squares between each group of data, where  $n$  was the total number of samples,  $k$  was the total number of groups,  $g$  was the value and  $G$  was all of the values across all groups and  $n_g$  was the number in each group.

$$MSS_B = \frac{\sum_{g \in G} n_g (\bar{x}_g - \bar{x}_G)^2}{k - n}$$

### 2.4.7 Tukey test

The Tukey test allows us to make multiple mean comparisons within the data with just a one step procedure (Equation 2.10). It was essentially a  $t$ -test that takes into consideration multiple testing. By assigning known groups to the data one can infer if these groups have significantly different means from all other groups within the data. Pairwise comparisons of all the possible groups' means were made and the difference was compared to the standard error.

**Equation 2.10** The Tukey test, where  $Y_A$  was the greater of the two means,  $Y_b$  was the smaller of the two means and SE was the standard error of the sum of the means.

$$q_s = \frac{Y_A - Y_B}{SE}$$

### 2.4.8 Kruskal-Wallis

The Kruskal-Wallis test is a one-way ANOVA on ranks, it is essentially an extension to the Mann Whitney U test. Similarly to Mann Whitney U, the Kruskal Wallis test uses ranks, and therefore, is a non-parametric test useful for non-normally distributed data. Additionally, similarly to ANOVA, the Kruskal-Wallis test can allow testing of >2 categories of data.

$$H = (N - 1) \frac{\sum_{i=1}^g n_i (\bar{r}_i - \bar{r})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (r_{ij} - \bar{r})^2}$$

where  $n_i$  is the number of observations in group  $i$ ,  $g$  is the number of groups,  $r_{ij}$  is the rank of observation  $j$  from group  $i$ ,  $N$  is the total number of observations across all groups,  $\bar{r}_i$  is the average rank of all observations in group  $i$  and  $\bar{r}$  is the average of all the  $r_{ij}$ .

A p-value can then be approximated from  $H$  from the table of  $X^2$  distributions and the degrees of freedom ( $g-1$ ). The function `kruskal.test` from the stats R package was used in R version 3.2.1.

### 2.4.9 Variance and IQR

Variance of a dataset can be measured as the sum of the squared distance of the data points from their mean. The IQR of the data was the lower quartile (the data point at 25%) subtracted from the upper quartile (the data point at 75%). The IQR was useful when data was not normally distributed.

### 2.4.10 Log rank Test

The log rank test is used to compare the survival experience of two different experimental statuses. It tests for the null hypothesis that there is no difference for the populations for the probability of an event at any time period, unlike survival curves, where a comparison at arbitrary time points are given.

For each time the number of events in each group are calculated and compared to the number expected if the null hypothesis were to be true. For each group the test statistic is calculated using  $(O-E)^2/E$ , where  $O$  is the number of observed events and  $E$  is the number of expected events. The comparison is completed using a  $X^2$  test (Section 2.4.4) and from the  $X^2$  distribution tables, a  $p$ -value can be provided allowing acceptance or rejection of the null hypothesis<sup>180</sup>.

The log rank test has advantages such that the whole follow up period is utilised, and no information about the shape of the survival curve or distribution of survival times is required.

The log rank test was completed using the `survdiff` function of the survival R package<sup>26</sup>.

### 2.4.11 Shapiro-Wilk

The Shapiro-Wilk test was used to determine if a sample came from a normally distributed population. The null hypothesis was that the data was from a normally distributed population and so was rejected if the  $p$ -value was less than the chosen alpha value (typically 0.05). Equation 2.11 was utilised to determine the  $W$  statistic, where  $x_{(i)}$  was the  $i$ th smallest number in the sample (the  $i$ th order statistic) and  $a_i$  were the constants derived from the covariance matrix of the order statistics<sup>181</sup>. The algorithm used in R also has the ability to calculate a  $p$ -value from  $W$ <sup>182</sup>. This was used with standard settings, under R version 3.2.1 for all Shapiro-Wilk testing.

**Equation 2.11** The Shapiro-Wilk test, where  $x_{(i)}$  was the  $i$ th smallest number in the sample (the  $i$ th order statistic) and  $a_i$  were the constants derived from the covariance matrix of the order statistics



$$W = \frac{(\sum_{i=1}^n a_i x_{(i)})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

### 2.4.12 Brent

The Brent method is an algorithm that combines three root-finding algorithms. It is as quick and more reliable than most other bisection methods. It is an iterative method that moves inwards from two points known to be on a quadratic curve until the root that provides the optimal bisection is discovered. The function `optim` from the `stats` package in R was used, with the argument `method=Brent`, for the optimal bisection of density plots.

### 2.4.13 Benjamini – Hochberg Multiple Testing Correction

In order to limit false discovery rates when completing multiple tests, multiple testing correction is completed. This particularly is useful for removing false positive hits, but has the trade-off of creating false negatives. The Benjamini-Hochberg method is a widely used procedure when completing multiple statistical tests, like testing each gene or probe between two groups. The correction starts by assigning a rank 1 to N, where 1 is assigned to the smallest p-value. Each p-value is then given a Benjamini-Hochberg critical value, using the formula  $(i/m)Q$ , where  $i$  is the assigned rank,  $m$  is the total number of tests and  $Q$  is the false discovery rate (chosen by the user). A comparison between the p-value and its critical value is then made by finding the largest p-value that is smaller than its critical value. Any p-value above this is then considered significant by the Benjamini-Hochberg method, and a new p-value is assigned.

The function `p.adjust` from the R package `stats` was used the Benjamini-Hochberg method passed as an argument.

### 2.4.14 Receiver Operator Characteristics (ROC)

ROC curves were a graphical plot to show the diagnostic ability of a classifier system as its discrimination threshold was varied. There was a trade-off between true predicted positives (sensitivity) and true predicted negatives (specificity) in the outcomes as this threshold was

## CHAPTER 2: MATERIALS AND METHODS

varied. ROC aims to identify the best threshold to give the best balance between the specificity and sensitivity. Each ROC also gives an AUC (area under the curve), which was a value between 0 and 1. Where 1 was a perfect model with all positives classed as positive and all negatives classed as negatives, and 0 shows there was no predictive value of the model at all. Generally, an AUC above 0.8 was valued as good. Two packages were used to produce ROC curves. For the HT chapter, ROC was performed using ROC, also part of the ROC bioconductor package<sup>183</sup>. Alternatively, to analyse the performance of the models built in the NanoString2 chapter the ROC function of the epi package<sup>184</sup> was used, this calculated the sensitivity and specificity as well as the AUC.

### **2.5 Clustering**

#### **2.5.1 Principal Component Analysis (PCA)**

PCA allows the visualisation of the maximum variability of a data set in two-dimensions. For a simple explanation, imagine there were ten samples and five genes and a graph was drawn with five axes, with each of the ten samples placed at the point that represents their value along each axis. Then identify a line that goes through as many of the samples as possible with the highest variation, that imaginary line was the first principal component. The second principal component was the line with the second highest variation and so on. Therefore, the majority of the variation of the data was found in the first two principal components and a 2D plot of these was enough to identify the biggest differences in samples.

This unsupervised mathematical procedure aiming to reduce dimensions of data works using a coordination transformation from the original data space to “eigenspace” using eigenvectors and eigenvalues of a matrix<sup>185</sup>. The first step was to calculate a covariance matrix of the data, with the aim to reduce redundancy and maximise variance. From the covariance matrix, which was used to measure how much the dimensions vary from the

## CHAPTER 2: MATERIALS AND METHODS

mean, the eigenvalues and eigenvectors can be determined. The covariance of two variables tells was a measure of how they vary together (Equation 2.12). Once the eigenvectors and the eigenvalues have been determined the eigenvalues can be sorted in descending size order.

**Equation 2.12 PCA covariance equation**

$$cov(X, Y) = \sum_{i=1}^N \frac{(x_i - \bar{x})(y_i - \bar{y})}{N}$$

### 2.5.2 Hierarchical Clustering

In this work, the commonly used UPGMA (“unweighted pair-group method using arithmetic averages”) method of hierarchical clustering was used. The highest similarity (or smallest distance) was used to identify the next two clusters to be merged. The distance of each sample to members of a cluster were computed with equal weights and the similarity or distance matrix was produced. This was updated and reduced at each computation, as samples/clusters were combined, allowing clustering to proceed by agglomeration as the similarity criterion was relaxed<sup>186</sup>.

#### 2.5.2.1 Pvclust

Pvclust<sup>187</sup> was a bootstrapping method that calculates the  $p$ -value for each cluster in a hierarchical clustering dendrogram object through the application for bootstrap resampling; clusters with significant AU  $p$ -values were shown with a red box.

### 2.5.3 $k$ -means Clustering

$k$ -means clustering aims to separate points into  $k$ -clusters so that the within clusters sum of squares was minimalized by seeking local optima so that moving of a point from one cluster to another will not reduce the sum of squares (Equation 2.13)<sup>188</sup>:

**Equation 2.13 Optimal local within cluster sum of squares.**  $x_i^{(j)}$  was the data point and  $c_j$  was the centroid, where  $i$  was a data point in cluster  $j$ .  $k$  was the number of clusters and  $n_j$  was the number of samples in cluster  $j$ .

$$J = \sum_{j=1}^k \sum_{i=1}^{n_j} \|x_i^{(j)} - c_j\|^2$$

Initially, centroids were arbitrarily picked each point was assigned to the closest centroid in Euclidean distance. Then the centroid was adjusted to the new cluster mean, the samples were reassigned to the closest centroid and this was repeated until convergence was reached. Convergence was when no observations can change the clusters when added and centroids were subsequently redefined<sup>189</sup>. The advantages of  $k$ -means clustering include speed and simplicity, whilst disadvantages include differing results per run due to the random starting centroid points and an unknown input value for  $k$ <sup>190</sup>. In this project, the optimal number of clusters was determined using the Bioconductor function NbClust<sup>191</sup>, which uses 30 metrics including the Silhouette, Dunn and Davies-Bouldin Indices (section 2.5.4).

#### 2.5.4 Silhouette, Dunn and Davies-Bouldin Indices

Three of these main techniques used for comparing how well data was clustering were the Silhouette, Dunn and Davies-Bouldin Indices.

The silhouette index compares the mean distance of a point to the others in it's cluster and then other clusters. It provides an index value between -1 and 1, where 1 was an indication that the point belongs to the correct cluster and -1 means it does not.

The Dunn index was the minimum distance of points between two different clusters divided by the maximum distance of points within a cluster for each cluster. Here, a larger value was representative of good clustering.

The Davies-Bouldin index takes the mean distance of the points within a cluster from their Barycentre and then divides this by the distance between the Barycentres and so a smaller value was an indication of good clustering.

### 2.5.5 Latent Process Decomposition

Latent Process Decomposition (LPD)<sup>192</sup> is a hierarchical Bayesian (probabilistic) model that was designed for the clustering of microarray data and thus can also be used with other forms of count data. It estimates the most probable/optimal number of clusters, and determines the probability of a sample belonging to each cluster, rather than membership of a cluster being assigned. This was important as samples were often heterogeneous made up of cells from different clones of cancer. Also different biological processes often work together to influence expression levels.

LPD makes the assumption that a sample's expression was determined by a series of processes. Each process has an associated expression profile which was determined by the algorithm. A sample's expression profile was then de-convoluted in to these process expression profiles. For example, Gene A has expression of  $n$  genes similar to the expression of the genes that make up the signature of process 1, and expression of  $m$  genes was similar to the expression of genes that make up the signature of process 2.  $n$  genes were of a higher similarity to process 1 than  $m$  genes were to process 2 and so max likelihood was higher for process 1; Gene A has 0.78 for process 1 but still 0.22 for process 2, etc.. So it has some similarity through some genes to process 2 but has more similarity through more genes to process 1.

The first step of LPD was to estimate the most probable number of clusters or "processes" using the maximum likelihood solution and a uniform prior. A uniform prior was a probability assumption with limited knowledge. E.g. a ball under 3 cups A, B, C has probability prior of  $p(A) = p(B) = p(C) = 1/3$ , where changing the order of the probabilities of the cups makes no change to the prediction. In the final model, a prior was defined to avoid over fitting by penalizing over complex. The parameter (sigma) in this prior was estimated next through cross-validation.

After these parameters were defined, the final solution was obtained by iteratively updating various parameter values of the Dirichlet distribution (a collection of multivariate,

## CHAPTER 2: MATERIALS AND METHODS

continuous probability distributions that was a more generalized version of beta distributions) modelling expression. Process mean values were initialized to the mean expression across the data set for each gene, whilst variances were set to the variance of their respective genes<sup>192</sup>.

### **2.6 Model Optimisation**

Predictive models were a supervised learning method, which has been applied to both the NanoString1 and NanoString2 data. For the NanoString2 data, it was divided into training, and test sets for a more robust and accurate model evaluation. A number of different models and modelling techniques were applied.

#### **2.6.1 GLM: Generalised Linear Model**

There are two important aspects of GLM<sup>193</sup>: General and Linear. Linear because the underlying equation was that of a straight line:  $Y = \beta_0 + \beta_1 X_1$ . In this example  $Y$  was the predicted or response variable, whilst  $X$  was a single predictor or explanatory variable.  $\beta_0$  was the y-intercept and was constant, whilst  $\beta_1$  was the slope or weight of variable  $X_1$ . General because the equation was able to handle multiple explanatory ( $X$ ) variables e.g.  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ . Any control variables may be included and if so, should precede the explanatory variable of interest within the equation, in general practice.

The explanatory variables may be numerical and continuous or binomial/factorial with levels. The GLM generalised linear regression by allowing the linear model to be related to the response variable via a link function allowing the modelling of binary response variables through logistic regression and ordinal variables through proportional odds models.

GLM was performed as an initial step in identifying probes that were significant for predicting clinical category (CB vs. Ca, CB vs HR-Ca, CB-L-I-H trend and CB-Ca-Advanced Ca trend) within all of the data (NanoString1) and within the training data (NanoString2). Significant probe lists were then shrunk and selected using techniques such

## CHAPTER 2: MATERIALS AND METHODS

as Lasso, Step and Random Forest (section 2.6.2, section 2.6.4, section 2.6.3, respectively) and then these condensed lists were then used to build a final model. These models were then tested upon the test data (NanoString2). The R function `glm` from the `stats` package was used for logistic regression models and the `polr` function from the `MASS` package for proportional odds models. These were used with the R version 3.2.1.

### **2.6.2 Lasso**

Least absolute shrinkage and selection operator (LASSO) was a regression method that was capable of performing selection and regularisation in order to improve both interpretability and prediction accuracy of statistical models, respectively. A constraint was applied to which the sum of the absolute value of the regression coefficients must be less than. This forces some of the coefficients to be set to zero, allowing these covariates to be disregarded from the optimal statistical model. Thus allowing both subset selection and shrinking large regression coefficients so as to reduce over fitting<sup>194</sup>. Over fitting of models can be problematic because these models tend to have poor predictability and can be over responsive to minor fluctuations within the test data set. Lasso can be easily applied to a variety of statistical models including generalised linear models and proportional hazard models, amongst others.

### **2.6.3 Random Forest**

Random forest<sup>195</sup> (RF) was an ensemble method that was a combination of tree predictors (weak learners) such that each tree was built using a sample set constructed by random selection replacement (bootstrapping). Once built the result of the model was the combination of the results of all trees (votes for binary outcomes and mean for continuous outcomes). The random forest function (from the random forest package) was used for classification and for regression models.

## CHAPTER 2: MATERIALS AND METHODS

Each decision tree was built by taking the bootstrap data and repeatedly separating it at nodes. At each node a small subset of,  $m$ , variables were chosen at random, and the combination that optimises the split, according to some objective function, was found. At the next node another  $m$  variables were chosen and the same method was performed.  $m$  was generally set at  $\sqrt{p}$  or  $\frac{p}{3}$ , where  $p$  was the number of variables. As the number of trees increases, the generalisation error of the forest converges.

The importance of variables in the model were assessed in two ways. Internal out-of-bag (OOB) estimates were used to judge the quality of the model. OOB was the average error calculated for each variable from the trees that do not contain that specific variable in their respective bootstrap sample. The error was calculated using the misclassification rate of the subjects. These estimates were produced using a single run of a forest with 1,000 trees and no test set. Variables with large mean decrease in accuracy or OOB were more important for classification of the data. Additionally, a Gini coefficient was also used to assess importance. This was a measure of how each variable contributes to the homogeneity of the nodes and leaves in the RF. Each time a variable was used to split a node, the Gini Coefficient for the child nodes were calculated and compared to the original nodes coefficient. The coefficient can be between 0 (homogenous) and 1 (heterogeneous). These changes in Gini were summed and normalised for each variable. Again, variables that were more important have a higher mean decrease in Gini.

Random forest was applicable to regression. Mean squares error was usually used to determine error rate when using random forest with regression. MSE was the mean (divided by  $n$  (number of data points)) of the squares of the errors<sup>196</sup>.

Random forest (from the random forest package) was used for classification and for regression.



## CHAPTER 2: MATERIALS AND METHODS

### 2.6.4 Step for feature selection

StepAIC was a function of the MASS package in R<sup>197</sup>. It was an automated model selection technique that takes a model and inserts or removes each variable and assess the model quality using the AIC – Akaike Information Criteria (Equation 2.14). The model with a smallest AIC was selected as the optimal and then this model was fused in the next step, this was repeated until no further improvements in AIC were observed. StepAIC can be run forwards (where you begin with all variables and remove them), backwards (where you begin with a small number of variables and add them) or both (where variables were added or removed as required)<sup>198</sup>.

**Equation 2.14 shows how to calculate AIC. Where the model with the lowest AIC was deemed optimal. Where k was the number of parameters and L was the maximum value of the likelihood function for the model.**

$$AIC = 2k - 2 \log (\hat{L})$$

### 2.7 Pathway Analysis

#### 2.7.1 DAVID

DAVID was the Database for Annotation, Visualization and Integrated Discovery, it was a gene functional classification tool. It was a web-based tool whereby you submit a list of transcripts of interest and DAVID classifies the list into functional related gene groups, ranks the importance of the discovered gene groups (dgg) and summarises the major biology of the dgg<sup>199</sup>. DAVID was used to identify if there were any interesting biological functions of the transcripts identified as significant.

### 2.8 Survival Analysis Tools

Survival analysis was the analysis of data where the response variable was the time to an event, for example to death or as in our case time to failure. Individuals that fail after the end of the study at some point in the future were known to be censored. Survival analysis

tools were used to identify if there were transcripts capable of predicting relapse to hormone therapy (HT) prior to a two-year period. The non-failures were said to be censored as after the last follow up date you don't know if they have failed or not<sup>200</sup>.

### 2.8.1 Kaplan Meier (KM) Curves

The KM survival distribution was a discrete stepped survivorship curve, which gains information as each event (failure) occurs. There were two variables at any time point on the KM (Equation 2.15); those that have failed,  $d(t_i)$  and those at risk of failing,  $r(t_i)$  and this produces a step at each failure.

**Equation 2.15 The KM function.**

$$\hat{S}_{KM} = \prod_{t_i < t} \frac{r(t_i) - d(t_i)}{r(t_i)}$$

Censored points were denoted by a + on KM plots. Kaplan Meier plots were created using the `survfit()` function, specifying `type="kaplan-meier"` from the `survival` package<sup>201</sup> and `ggsurv()` of `GGally` package<sup>202</sup>, on R version 3.2.1. Dichotomised high/low expression levels were determined for each probe using *k*-means clustering and *k*=2 (section 2.5.3).

### 2.8.2 Cox Proportional Hazard

Cox Proportional hazard model was the most commonly used regression model for survival data. It assumes the hazard was of the form  $\lambda(t; Z_i) = \lambda_0(t)r_i(t)$ , where  $Z_i(t)$  was the set of explanatory variables for individual *i* at time *t*. The risk score for individual *i* was  $r_i(t) = e^{\beta Z_i(t)}$ , where  $\beta$  was a vector of parameters from the linear predictor  $\lambda_0(t)$ , which was an unspecified baseline hazard function that will cancel in due course. It guarantees that  $\lambda$  was positive for any regression model. Hazard was the instantaneous risk of failure, or instantaneous rate of change in the log number of survivors per unit time. Coxph was part of the `survival` package<sup>201</sup>.

# 3

## **NanoString Data Analysis 1: The Pilot Study**

### **3.1 Summary**

The Movember GAP1 global PCa biomarker initiative has multiple collaborators working on the identification of urinary biomarkers for the risk-stratification of PCa. Our laboratory is specifically interested in the RNA expression changes in PCa that are detectable within urinary cell sediments and extracellular vesicles (EVs). The EV RNA expression pilot study described here had the following aims:

1. Identify if PCa specific transcripts can be detected in urinary EVs
2. Assess whether transcript levels within urinary EVs were able to i) identify PCa per se, ii) distinguish aggressive from indolent PCa

## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

3. Identify if the NanoString system could be applied to Nugen Ovation amplified cDNA (Nanostring probes are strand specific and designed to be applied to mRNA).
4. Identify suitable methods for the analysis of the NanoString data

In the pilot study, expression levels of 57 transcripts were measured in 194 samples using NanoString technology (section 1.6.1). The NanoString technology was able to detect PCa specific markers (section 1.4.6), such as *TMPRSS2:ERG* which was detected in 58% of all PCa samples and in 19% of samples from men with no clinical evidence of PCa (CB). This result, confirmed by RTPCR demonstrated that i) NanoString technology was capable of capturing cDNA amplified by the Nugen Ovation kit; ii) EV mRNA contains PCa-specific transcripts, and iii) the methodology was sensitive enough to identify PCa in men with undiagnosed cancer or HGPIN.

Latent Process Decomposition unsupervised analysis (section 2.5.5), clustered the EV expression data into four groups: LPD groups 1 and 4 were saturated with high-risk and advanced cancers, whilst LPD groups 2 and 3 showed clinical diversity. The majority of the intermediate-risk samples resided within LPD group 2 and most of the CB were in LPD group 3 (section 3.5.5).

Supervised statistical approaches (Mann Whitney U test) determined nine probes significantly differently expressed between PCa (advanced, high-, intermediate- and low-risk) and non-PCa samples (Table 3.19), eleven probes significantly different between high-risk PCa and non-PCa samples (Table 3.20) and six probes between advanced PCa and non-PCa samples (Table 3.21).

Supervised modelling of the data (using generalised linear models (glm) and Lasso for shrinkage (section 2.6.2)) identified three models that distinguished; i) PCa vs. non-PCa with an AUC of 0.937, ii) aggressive PCa vs. non-aggressive PCa with an AUC of 0.852 and iii) advanced PCa vs. benign with an AUC of 0.983.

Twenty-three transcripts were significantly differentially expressed between PCa and non-PCa (Table 3.42), however, only seven were consistently differentially expressed

between the various data-analytic methods used (*DLX1*, *ERG3'*, *TMPRSS2:ERG*, *HOXC4*, *ERG5'*, *PCA3* and *HPN*). Four transcripts were consistently differentially expressed between aggressive PCa and non-aggressive PCa in various tests (Table 3.43) and two transcripts between advanced and non-PCa (Table 3.44).

These findings highlight that the transcript data collected from urinary EVs in PCa patients comes, at least in part, from the prostate and holds clinically relevant structure.

### **3.2 Introduction**

#### **3.2.1 The Research Gap**

Risk stratification is currently based on PSA, Gleason score and T stage. MRI is being phased in, but has been shown to have only 41% specificity in a recent study of low risk patients<sup>203</sup>. Patient clinical pathways would benefit from additional information on their PCa diagnostic and prognostic status. We propose that urine EV mRNA data could provide useful clinical information that could help tailor patients to treatment pathways based on their genetic composition and potentially improve uncertainty over which treatment pathway each patient should be assigned to. The PCA3 test has shown to provide minor improvements to risk stratification but importantly shows the utility of urine in PCa diagnostics and prognostics.

#### **3.2.2 The Pilot Study Aims**

The pilot study used NanoString technology to investigate the RNA expression level changes of 57 target transcript sequences within EVs extracted from urinary samples collected at the NNUH as part of the Movember study. The aims of this pilot study were:

- a) To identify if the transcript content of urinary EVs contained PCa derived material
- b) To identify if transcript levels within urinary EVs are linked to PCa risk stratification

c) To identify if NanoString is a suitable method for obtaining transcript level data from our cDNA samples

d) To identify appropriate methods for analysing NanoString data.

### 3.2.3 The Probe Targets

The 57 target transcript sequences were selected for the following reasons: i) prostate specific transcripts, ii) transcripts overexpressed in advanced PCa tissue (literature search), iii) suspected housekeeping genes, iv) tissue-specific controls for kidney, bladder and blood.

### 3.2.4 Risk classification of prostate cancer patients

Patients were placed into clinical risk categories based on D'Amico and NICE criteria: Prostate Cancer Diagnosis and Treatment 2014 guidelines<sup>40</sup>. In addition the intermediate risk patients were subdivided on Gleason (G3+4 Vs. G4+3), as progression rates between these two groups are very different (Table 3.1). The median age and PSA at diagnosis for each clinical category have been recorded (Table 3.2). For some computational analyses, specific risk groups were combined (Table 3.3).

**Table 3.1 Classification and Frequency of the sample types based on NICE criteria<sup>40</sup>. The quantity of samples for each clinical group can be seen as well as the clinical description of the group in terms of Gleason score, PSA level and T stage.**

<i>Classification: NICE Groupings</i>		
<i>Sample Class</i>	<i>Description</i>	<i>Number of Samples</i>
<i>Advanced (A)</i>	<i>Metastatic , PSA&gt;100, and G&gt;8</i>	<i>17</i>
<i>High-risk (H)</i>	<i>G7 PSA&gt;20</i>	<i>50</i>
<i>Upper Intermediate-risk (UI)</i>	<i>G4+3 PSA&lt;20</i>	<i>19</i>
<i>Intermediate-risk (I)</i>	<i>G3+4 PSA&lt;20 and IL= G6 PSA&gt;10</i>	<i>53</i>
<i>Low-risk (L)</i>	<i>Low G6 PSA&lt;10</i>	<i>10</i>
<i>Abnormal (S)</i>	<i>High PSA no Bx</i>	<i>4</i>
<i>CB&lt;1*</i>	<i>No evidence of Ca and PSA&lt;1</i>	<i>18</i>
<i>CBn*</i>	<i>No evidence of Ca and PSA normal to age</i>	<i>22</i>

<i>M</i> 19 5	<i>Removed for technical failure</i>	<i>I</i>
<i>Total</i>		<i>194</i>

\*CBN and CB<1 were combined to CB (as there was no significant difference between their expression levels  $p > 0.05$ : Two sample *t*-test) and UI and I were combined to I (as there was no significant difference between their expression levels  $p > 0.05$ : Two sample *t*-test).

Table 3.2 Median age and PSA at diagnosis for each clinical category, of samples that are used in subsequent analysis.

Sample Class	Number of Samples	Median Age	Median PSA at Dx
Advanced (A)	17	78	110
High-risk (H)	50	73	27
Upper Intermediate-risk (UI)	19	74	9.55
Intermediate-risk (I)	53	67.5	8.35
Low-risk (L)	10	68	5.95
CBN and CB<1	40	68	1.1

Table 3.3 Sample numbers used in i) 'Cancer', ii) 'Aggression' and iii) 'Extreme' computational analyses.

Group	Number of Samples
<i>Cancerous (A, H, I and L) and No Evidence of Cancer (CB)</i>	<i>Cancerous =149 / CB =40</i>
<i>Aggressive (A, H) and Non-Aggressive (I, L)</i>	<i>Aggressive = 67/ Non-Aggressive = 82</i>
<i>Extremes (A Vs. CB)</i>	<i>A=17 / CB= 40</i>

### 3.3 Data Pre-processing and Technical Variation

#### 3.3.1 Normalisation and Background correction

The NanoString analyses provided data for 57 test probes, and 14 non-human system control probes (6 positive-control probes and 8 negative-control probes) in 194 Nugen Ovation amplified cDNA samples. The 6 positive control probes detected spiked-in control sequences that were used to assess the overall NanoString assay transcript detection efficiency for each sample, and generated a normalisation factor (NF) in the following way: The average nCount for each samples' positive controls was calculated and this number was divided by the sum of all samples averages. Each nCount value

was then multiplied by the sample-specific NF. Background correction was not applied and background subtraction was found not to be appropriate for this data (data not shown).

### 3.3.2 NanoStringNorm – Quality of Data and its Normalisation

The quality of the normalisation was evaluated using the NanoStringNorm R package (section 2.3.2.1). Other than a few flagged samples (M\_14\_7, M\_19\_5, M\_36\_7, (Table 3.4)), and a few flagged probes (*KLK4*, *GAPDH* and *FOLH1*, (Table 3.5)), the data was of overall good quality. The three probes were flagged due to high mean and/or standard deviation or for *FOLH1* not following the Loess curve of best fit. For *GAPDH*, we predicted similar housekeeping properties as in cell RNA, however that is not what has been observed (section 3.3.5). For *KLK4*, it suggests high expression in the samples with a wide range of signals (considering we have samples across different clinical categories this is expected). For *FOLH1*, the Loess curve of best fit is a non-parametric regression derived curve that is similar to a line of best fit through all of the data. To not follow it simply suggests that this probe is expressed rather differently to the other probes in these samples (again could be due to the range of clinical categories used).

Some cartridges (each cartridge is loaded with twelve samples and then run on the NanoString machine) showed significantly different means and standard deviations in comparison to others in the raw data. The flagged outliers were considered with caution and reviewed further in subsequent analyses

**Table 3.4 Three samples were flagged by NanoStringNorm.**

<i>Samples</i>	<i>Issues</i>
<i>M_14_7</i>	<i>Low sample mean</i>
<i>M_19_5</i>	<i>Low sample mean</i>
<i>M_36_7</i>	<i>Low sample mean Normalisation factor flagged as influential outlier</i>



**Table 3.5 Three probes were flagged by NanoStringNorm.**

<i>Probes</i>	<i>Issues</i>
<i>KLK4</i>	<i>High mean and SD</i>
<i>GAPDH</i>	<i>High mean and SD</i>
<i>FOLH1</i>	<i>Doesn't follow Loess curve of best fit</i>

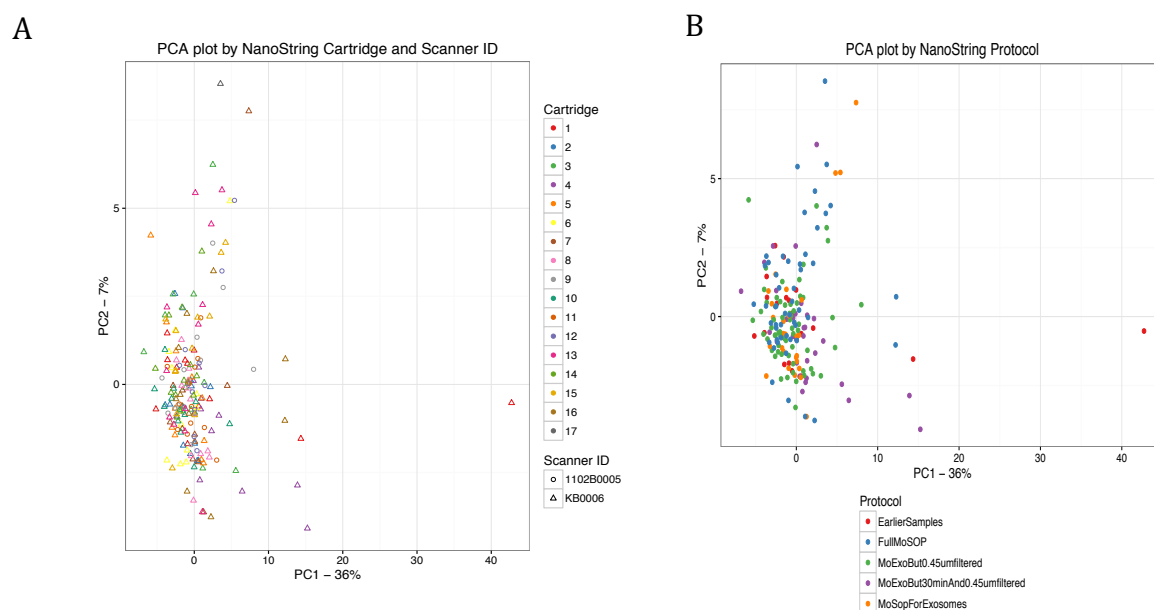
### 3.3.3 Experimental and Technical Investigation

#### 3.3.3.1 NanoString Scanners and Cartridges

NanoStringNorm showed significant differences between the mean and standard deviation of the normalised data between cartridges; indicating there might be batch effects on the scanner and cartridge-dependent variables. Scanner and cartridge-dependent variations were therefore examined using Principal component analyses (PCA) (section 2.5.1). PCA did not detect any clustering based on technical artefacts (Figure 3.1A), and there was no significant association between mean expression per sample and either cartridge (Kruskal-Wallis rank sum test:  $p = 0.17$ ,  $\chi = 21.21$ ), or Scanner (Kruskal-Wallis rank sum test:  $p = 0.71$ ,  $\chi = 0.14$ ).

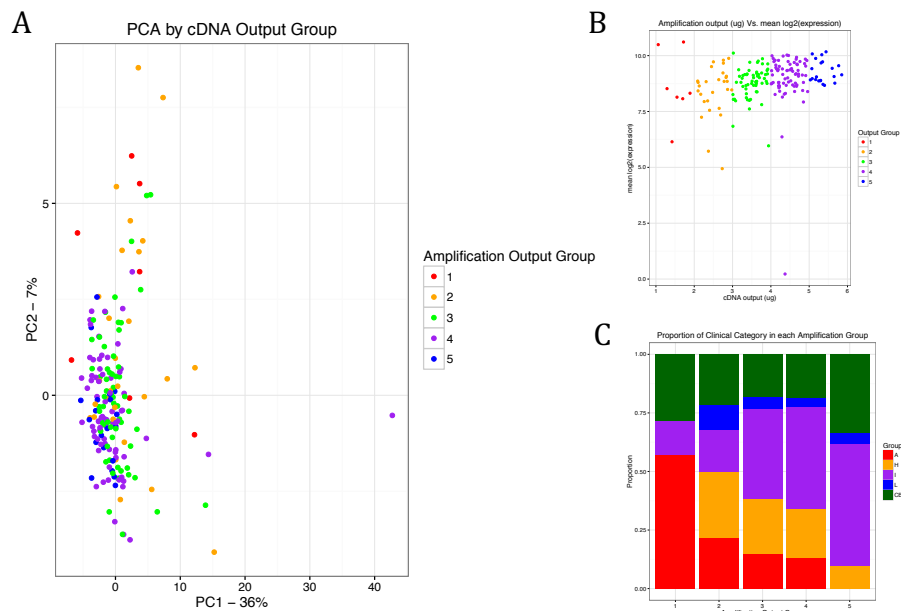
#### 3.3.3.2 RNA Extraction and Amplification

At the beginning of the urine-collection study, protocol optimisation for RNA yield from RNA extractions was conducted (by Marcel Yazbek-Hanna and Rachel Hurst, section 2.1.3), which led to samples from multiple variant protocols being included in the pilot study set. PCA (section 2.5.1) was applied and no clustering was observed due to RNA extraction protocol (Figure 3.1B). There was no significant association between the median expression for each sample and the RNA extraction protocol used (Kruskal-Wallis rank sum test:  $p = 0.16$ ,  $\chi^2 = 6.5$ ).



**Figure 3.1** To ensure there were no batch issues PCA plots were produced of NanoString loading batches and RNA extraction protocol. A) PCA did not identify any clustering associated with NanoString cartridge or scanner used. Along with the Kruskal-Wallis rank sum results also (Cartridge:  $p = 0.17$ , Scanner  $p = 0.71$ ), it was deemed there was no batch effect produced by NanoString loading. B) PCA does not identify any clustering associated with RNA extraction protocol used and the Kruskal-Wallis rank sum test was also insignificant ( $p = 0.16$ ). Thus it was deemed that using no filter, a 45 $\mu$ m filter, and a 45 $\mu$ m filter with a 30-minute wait along side the Qiagen micro RNA RNeasy kit using manufactures' protocols made no difference.

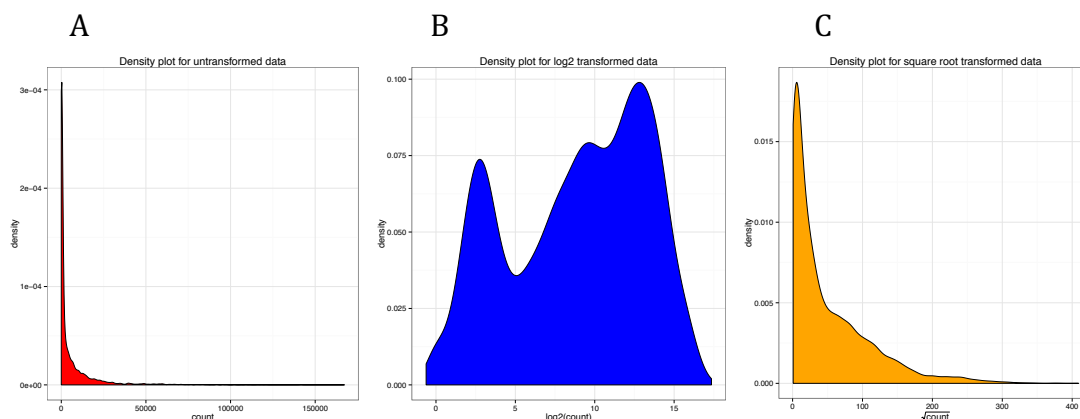
Due to the limited amounts of EV RNA harvestable from urine, 15-20ng RNA from each sample was amplified using a Nugen Ovation WTA2 cDNA amplification kit. The amount of cDNA obtained after amplification (in  $\mu$ g) was investigated for clustering affects using PCA (section 2.5.1) and correlation (section 2.4.3). cDNA yields were split into groups; group 1 = 1-2 $\mu$ g, group 2 = 2-3 $\mu$ g, group 3 = 3-4 $\mu$ g, group 4 = 4-5 $\mu$ g, and group 5: >5 $\mu$ g. Mild clustering affects were observed; samples with lower Ovation output had a lot more spread than higher amounts of output (Figure 3.2A) but no significant correlation was found between cDNA yield and median  $\log_2$  expression per sample ( $p = 0.09$ ,  $r = 0.12$ , Pearson's correlation, Figure 3.2B). The distribution of clinical categories within each Amplification yield group was not statistically significant; ( $\chi = 26.2$ ,  $p > 0.05$ ,  $\chi^2$  test (section 2.4.4), Figure 3.2C).



**Figure 3.2** A) Amplification cDNA yield shows mild clustering (cDNA yields were grouped: group 1 = 1-2 $\mu$ g, group 2 = 2-3 $\mu$ g, group 3 = 3-4 $\mu$ g, group 4 = 4-5 $\mu$ g, and >5 $\mu$ g in group 5). B) Amplification cDNA yield shows no influence on sample mean expression C) Amplification cDNA yield shows dependence on clinical category.

### 3.3.4 Transforming data to a normal distribution and the Shapiro-Wilk test

log<sub>2</sub> and square root transformation (section 2.3.3) was applied to attempt to get the dataset closer to a normal distribution (Figure 3.3). Neither the log<sub>2</sub>-transformed, nor the square root transformed, nor the non-transformed data are normally distribution according to the Shapiro-Wilk test (section 2.4.11, Table 3.7, Table 3.8). However, for the majority of the samples (the first 70 and last 70), the *W* statistic is higher for the log<sub>2</sub>-transformed data, indicating that the data is closer to a normal distribution than for the other transformations (Table 3.6).



**Figure 3.3** Density plots showing the distribution of a) the non-transformed data. B) the log<sub>2</sub> transformed data. C) the square root transformed data.

The Shapiro-Wilk test was also applied to ten randomly selected probes in each of the datasets (un-transformed, log<sub>2</sub> transformed and square root transformed) to see how the distribution of some individual probes varied; the majority were not normally distributed. The NanoStringNorm flagged probes (*KLK4*, *GAPDH* and *FOLH1*) had similar results to the other probes. These results led to the use of non-parametric tests wherever possible during analysis. A log<sub>2</sub> transformation was applied so that probe data was closer to a normal distribution, as is standard practice for NanoString data<sup>204</sup>.

CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

Table 3.6 Shapiro-Wilk test results on the first 70 and last 70 samples (all probes) for the non-transformed, log<sub>2</sub> transformed and square root transformed datasets.

	<i>Un-transformed</i>			<i>Log<sub>2</sub> transformation</i>			<i>Square root transformation</i>		
	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>
<i>The first 70 samples (1-70)</i>	<b>0.5143</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.9479</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.7937</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The last 70 samples (124-194)</i>	<b>0.525</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.9496</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.8105</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>

Table 3.7 Shapiro-Wilk test results for 10 randomly selected probes for the non-transformed, log<sub>2</sub> transformed and square root transformed datasets.

	<i>Un-transformed</i>			<i>Log<sub>2</sub> transformation</i>			<i>Square root transformation</i>		
	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>
<i>Probe 24</i>	<b>0.7704</b>	<b>3.886x10<sup>-16</sup></b>	<b>No</b>	<b>0.6369</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.9421</b>	<b>4.929x10<sup>-07</sup></b>	<b>No</b>
<i>Probe 2</i>	<b>0.8414</b>	<b>2.847x10<sup>-13</sup></b>	<b>No</b>	<b>0.5119</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.942</b>	<b>4.883x10<sup>-07</sup></b>	<b>No</b>
<i>Probe 22</i>	<b>0.7613</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.8579</b>	<b>1.762x10<sup>-12</sup></b>	<b>No</b>	<b>0.9548</b>	<b>7.663x10<sup>-06</sup></b>	<b>No</b>
<i>Probe 17</i>	<b>0.9394</b>	<b>2.906x10<sup>-07</sup></b>	<b>No</b>	<b>0.7784</b>	<b>7.562x10<sup>-16</sup></b>	<b>No</b>	<b>0.9952</b>	<b>0.7955</b>	<b>Yes</b>
<i>Probe 47</i>	<b>0.9888</b>	<b>0.1301</b>	<b>Yes</b>	<b>0.6097</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.9694</b>	<b>0.000306</b>	<b>No</b>
<i>Probe 34</i>	<b>0.8222</b>	<b>4.011x10<sup>-14</sup></b>	<b>No</b>	<b>0.7533</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.9652</b>	<b>9.965x10<sup>-05</sup></b>	<b>No</b>
<i>Probe 13</i>	<b>0.6355</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.9869</b>	<b>0.0708</b>	<b>Yes</b>	<b>0.8741</b>	<b>1.227x10<sup>-11</sup></b>	<b>No</b>
<i>Probe 43</i>	<b>0.1609</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.885</b>	<b>4.967x10<sup>-11</sup></b>	<b>No</b>	<b>0.4477</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>
<i>Probe 29</i>	<b>0.9918</b>	<b>0.3421</b>	<b>Yes</b>	<b>0.5435</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.9626</b>	<b>5.107x10<sup>-05</sup></b>	<b>No</b>
<i>Probe 26</i>	<b>0.9817</b>	<b>0.01245</b>	<b>No</b>	<b>0.5934</b>	<b>2.2x10<sup>-16</sup></b>	<b>No</b>	<b>0.9573</b>	<b>1.362x10<sup>-05</sup></b>	<b>No</b>

Table 3.8 Shapiro-Wilk test results for the three probes identified by NanoStringNorm as having potential quality issues in the three datasets: non-transformed, log<sub>2</sub> transformed and square root transformed.

<i>Un-transformed</i>	<i>Log<sub>2</sub> transformation</i>	<i>Square root transformation</i>
-----------------------	---------------------------------------	-----------------------------------

CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

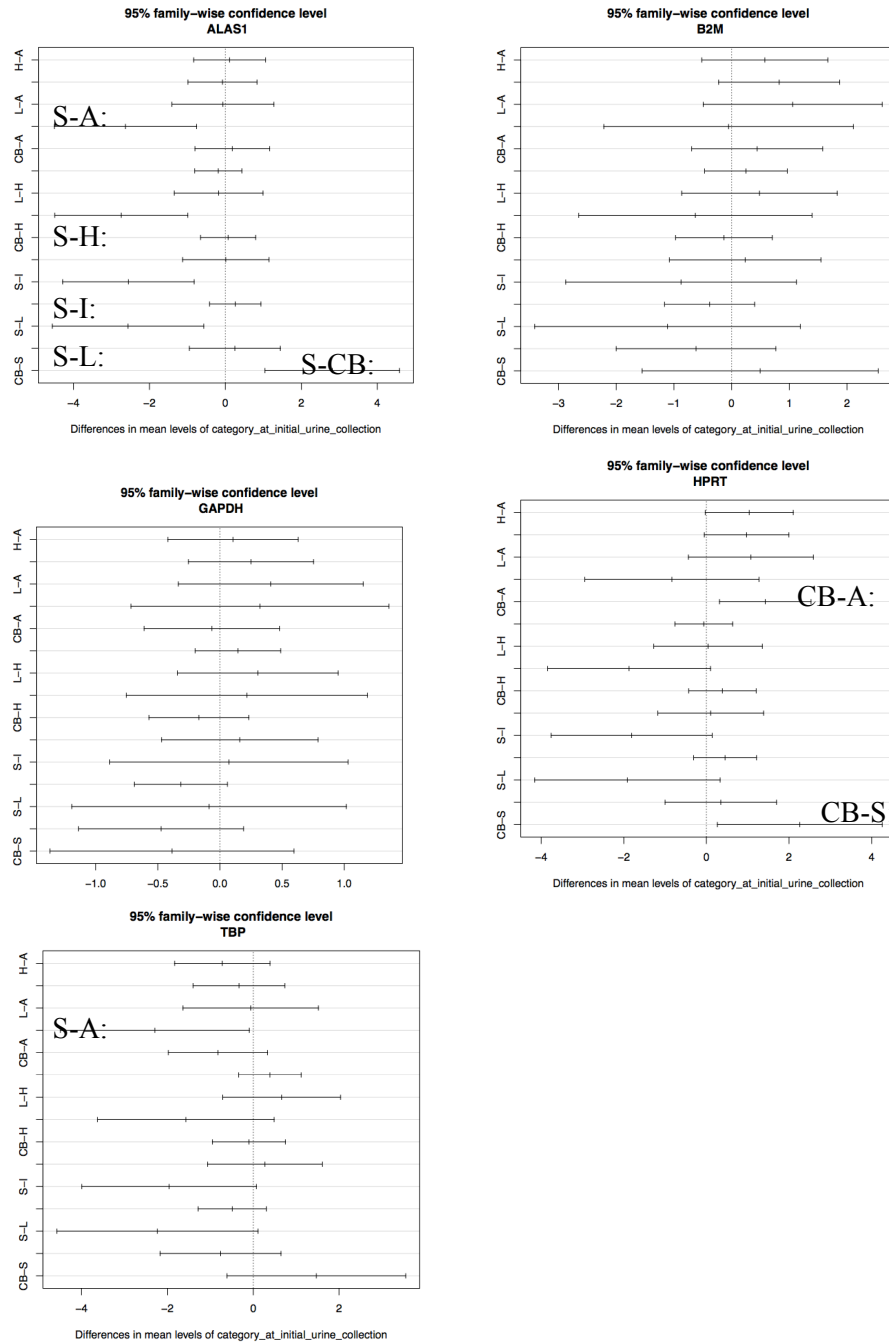
	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>
<i>KLK4</i>	<b>0.9918</b>	<b>0.3421</b>	<b>Yes</b>	<b>0.5435</b>	<b><math>2.2 \times 10^{-16}</math></b>	<b>No</b>	<b>0.9626</b>	<b><math>5.107 \times 10^{-05}</math></b>	<b>No</b>
<i>GAPDH</i>	<b>0.9713</b>	<b>0.0005136</b>	<b>No</b>	<b>0.4575</b>	<b><math>2.2 \times 10^{-16}</math></b>	<b>No</b>	<b>0.9635</b>	<b><math>6.352 \times 10^{-05}</math></b>	<b>No</b>
<i>FOLH1</i>	<b>0.9394</b>	<b><math>2.906 \times 10^{-07}</math></b>	<b>No</b>	<b>0.7784</b>	<b><math>7.562 \times 10^{-16}</math></b>	<b>No</b>	<b>0.9952</b>	<b>0.7955</b>	<b>Yes</b>

### 3.3.5 Housekeeping Probes

Five probes (*ALASI*, *B2M*, *GAPDH*, *HPRT* and *TBP*) were added to the NanoString project to identify housekeeping transcripts (transcripts that remain relatively consistent between samples of different clinical category). Housekeeping transcripts are added so that comparisons between the samples within an expression analysis may be performed accurately. The five transcripts are known housekeeping transcripts in cell mRNA, but there is very little known about EV RNA housekeeping transcripts at present.

There is very little correlation between the six clinical categories (Adv, H, I, L, S, CB) within each housekeeper expression profile (Tukey-ANOVA test, Table 3.1, Figure 3.4); the S clinical group (those with a high PSA but no Bx,  $n = 4$ ) has the most significant differences compared to the other clinical categories; for *ALASI* comparisons with all other clinical groups and the S group were significant. For *HPRT* two comparisons were significant, one between CB and S and the other between CB and Adv. For *TBP* only one comparison was significant, (between the S group and advanced group). However, there were only four samples in the S group and so the results of the significance test for this group were treated cautiously. Ignoring significant comparisons that included the S group, left only one significant comparison (for the *HPRT* probe between CB and Adv).

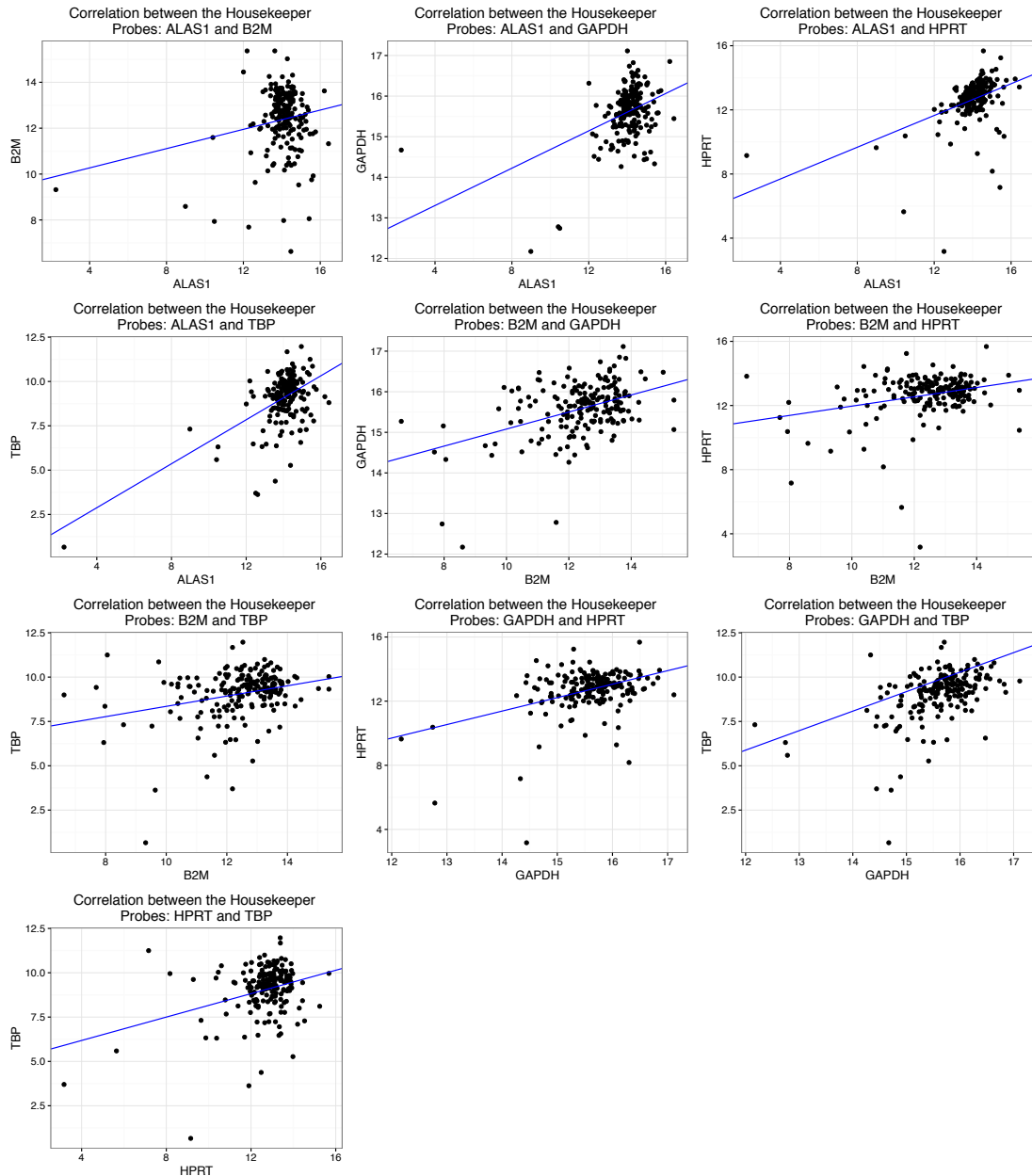
## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY



**Figure 3.4** Tukey test comparisons of clinical category for housekeeping probes. When the bar does not cross the mid-point of the x-axis then the comparison is significant. The Tukey test takes each of the five probes (*ALAS1*, *B2M*, *GAPDH*, *HPRT*, and *TBP*) and detects significant expression differences between the six clinical categories. The significant comparisons with S (high PSA/negative Bx samples) is treated cautiously as there were only  $n = 4$  samples within this group. This leaves only one group comparison (CB with Advanced samples in *HPRT*) that showed any significant difference. A good housekeeping probe would be expected to not differ between clinical categories.



## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY



**Figure 3.5** Correlation plots between the housekeeper transcripts: ALAS1, B2M, GAPDH, HPRT, and TBP.

Pearson's correlation coefficients ( $R$ ) between housekeeping probes was below 0.5 in 9/10 comparisons (0.53 being the highest correlation), which suggests they are not well correlated (Table 3.9, Figure 3.5). This makes the choice of which housekeeping probes to normalise the data with difficult. So, for this reason it was decided to go ahead without using housekeeper style normalisation for these data.

**Table 3.9 Housekeeper probe Pearson's correlation results, looking for correlating housekeeping probes.**

<b>R</b>	<b>ALAS1</b>	<b>B2M</b>	<b>GAPDH</b>	<b>HPRT</b>	<b>TBP</b>
<b>ALAS1</b>	-	<b>0.19</b> ( <i>p</i> = 0.008)	<b>0.43</b> ( <i>p</i> = 7.8x10 <sup>-10</sup> )	<b>0.44</b> ( <i>p</i> = 1.8x10 <sup>-10</sup> )	<b>0.53</b> ( <i>p</i> = 1.6x10 <sup>-15</sup> )
<b>B2M</b>		-	<b>0.44</b> ( <i>p</i> = 2.5x10 <sup>-10</sup> )	<b>0.29</b> ( <i>p</i> = 5.2x10 <sup>-05</sup> )	<b>0.28</b> ( <i>p</i> = 7.7x10 <sup>-05</sup> )
<b>GAPDH</b>			-	<b>0.4</b> ( <i>p</i> = 1.1x10 <sup>-08</sup> )	<b>0.49</b> ( <i>p</i> = 7.8x10 <sup>-13</sup> )
<b>HPRT</b>				-	<b>0.32</b> ( <i>p</i> = 4.8x10 <sup>-06</sup> )
<b>TBP</b>					-

An alternative to using housekeeping transcripts could be to use a similar method to the PCA3 test, which uses *KLK3* (PSA) to enhance the expression of other probes in the data. *KLK3* adjusted data was produced but the resulted data showed much weaker, plateaued, signal strength and therefore, was not used for any subsequent analysis (data not shown).

### 3.3.6 Removal of Outliers

M\_19\_5 was identified via PCA (Figure 3.6) and NanoStringNorm (Table 3.4) as being an outlier that may hinder further analyses. Further investigation into this sample highlighted that 44 out of 57 probes for sample M\_19\_5 failed; in the (positive control normalised, log<sub>2</sub> transformed) data all 44 probes had a value of “-0.07400058”, indicating that they were undetectable. The other samples of this cartridge and scanner appear to have worked. Therefore, this sample alone will be removed for all subsequent analyses.

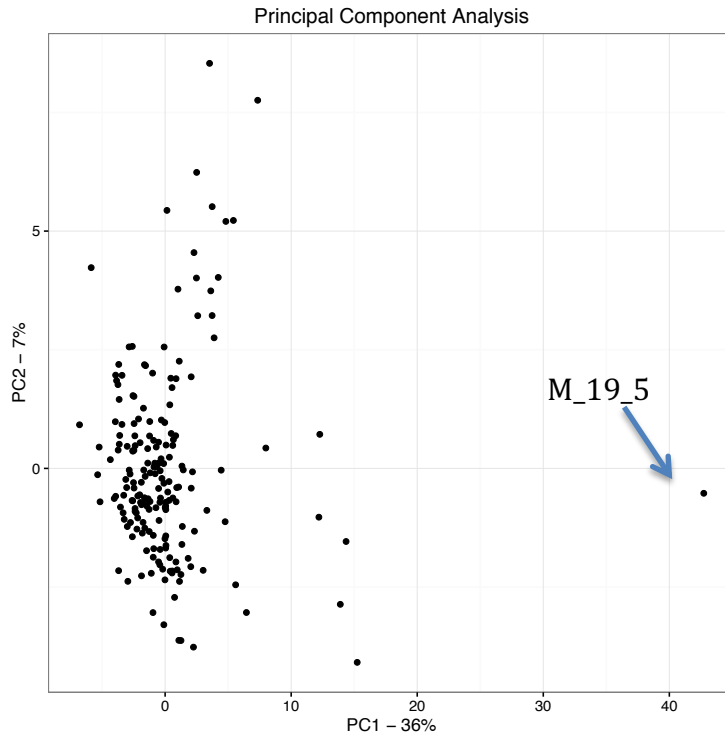


Figure 3.6 PCA plot of all log<sub>2</sub> normalised data identifies an outlier samples M\_19\_5.

### 3.3.7 Correlating Gene Probes

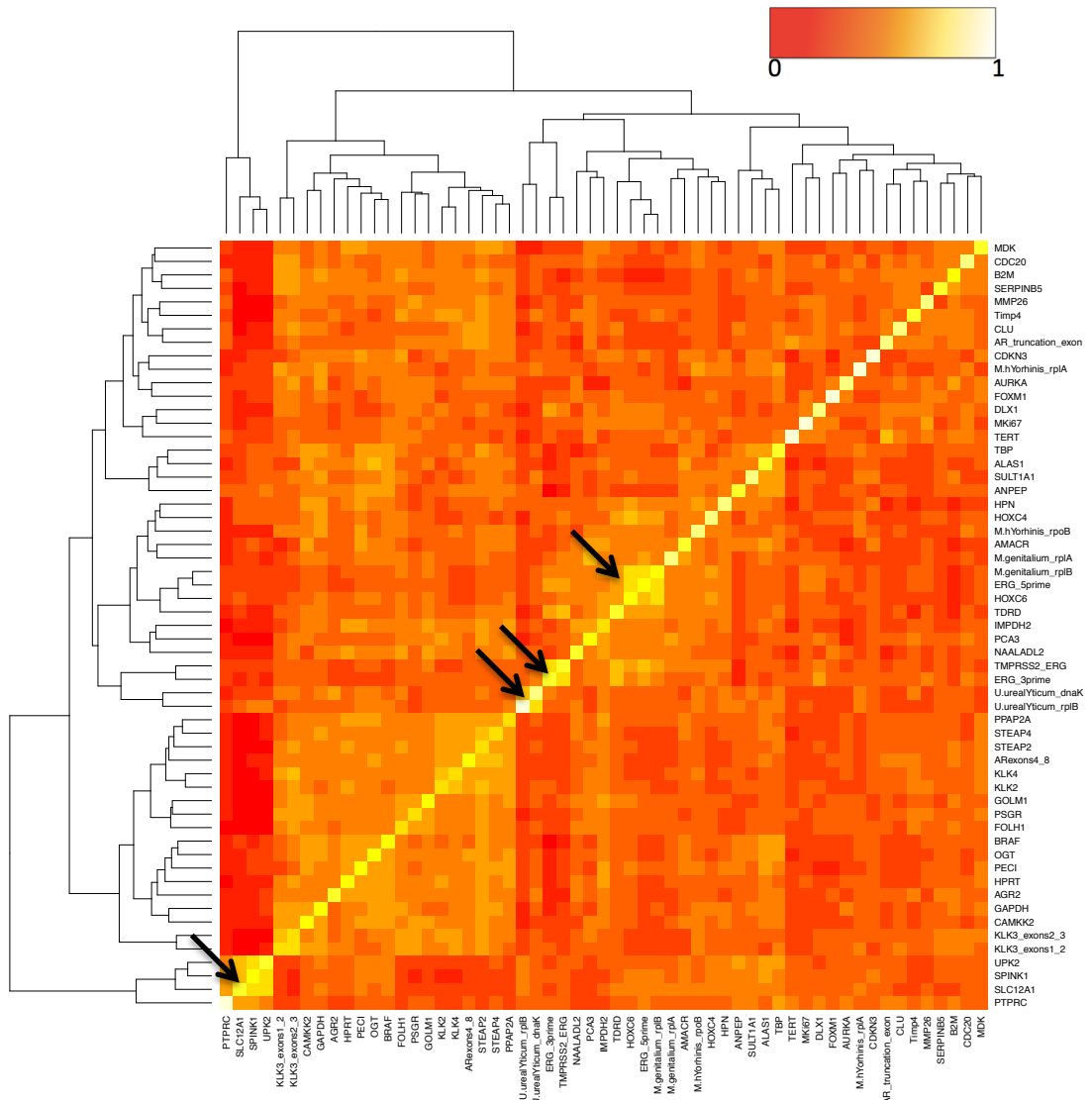
Pearson's correlation (section 2.4.3) between data from all of the probes identifies four clusters of probes that showed strong inter-correlation (Figure 3.7, Table 3.10). Cluster 1: probes for *ERG* 3' and *TMPRSS2:ERG*; Cluster 2: the two probes for the bacteria *U.urealyticum*; Cluster 3: two *M.genitalium* probes, *HOXC6* and *ERG* 5' and Cluster 4: *SLC12A1*, *SPINK1* and *UPK2*.

The data for probes in Clusters 1 and 2 were biologically expected to correlate, as were Cluster 3's two bacterial probes (*M.genitalium* RplA and RplB: Pearson's correlation:  $p = 1.23 \times 10^{-05}$ ,  $R = 0.31$ ). However, Cluster 3's other correlations were not expected and were even more pronounced i) *M.genitalium* RplB and *HOXC6* (Pearson's correlation:  $p < 2.26 \times 10^{-16}$ ,  $R = 0.88$ ) ii) *M.genitalium* RplB and *ERG* 5' (Pearson's correlation:  $p < 2.26 \times 10^{-16}$ ,  $R = 0.83$ ) and iii) between *HOXC6* and *ERG* 5' (Pearson's correlation:  $p < 2.26 \times 10^{-16}$ ,  $R = 0.73$ ) (Figure 3.8). Also in Cluster 3, the two *M.genitalium* probes would be expected to have similar signal strength, which is not the case (Figure 3.8).

## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

*M.genitalium* RplA had signal range ~-1 to 4, whilst *RplB* had a signal range of 0-12 with most samples above 6. *M.genitalium* RplB signal strength was actually more similar to *HOXC6* (~5-16) and *ERG3'* (~2-12).

Figure 3.7 Heatmap showing correlation between all NanoString probe data. The colours reflect the

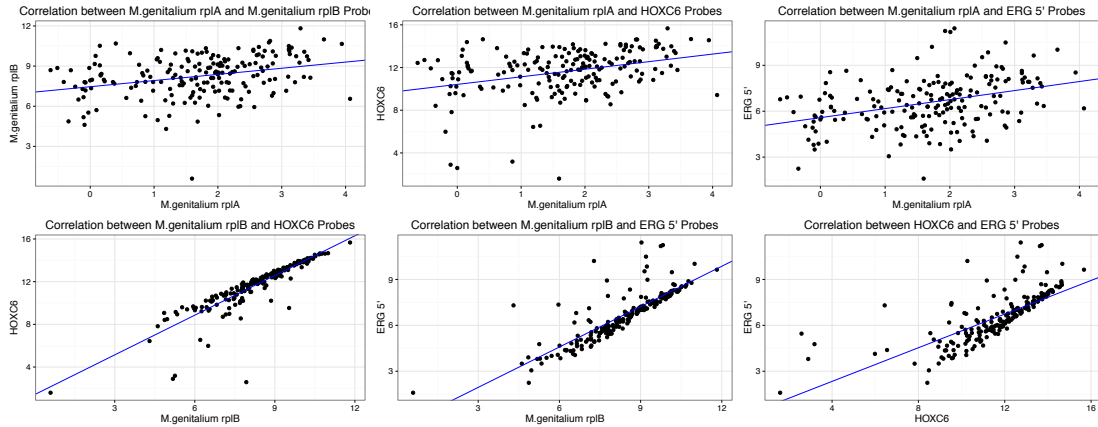


$R$  value of the correlation, where 1 is perfect correlation (represented by yellow) and 0 is uncorrelated (represented by red), with orange in between.

**Table 3.10 Four Clusters of probes that correlate with each other (Pearson's correlation).**

<i>Correlating probes r values</i>			
<b>ERG 3'</b>			
<i>TMPRSS2:ERG</i>	$p < 2.2 \times 10^{-16}$ , $R = 0.74$		
<b><i>U.urealyticum dnaK</i></b>			
<i>U.urealyticum RplB</i>	$p < 2.2 \times 10^{-16}$ , $R = 0.56$		
	<b><i>M.genitalium RplB</i></b>	<b><i>HOXC6</i></b>	<b><i>ERG 5'</i></b>
<i>M.genitalium RplA</i>	$p = 1.23 \times 10^{-05}$ , $R = 0.31$	$p = 4.08 \times 10^{-07}$ , $R = 0.36$	$p = 3.65 \times 10^{-08}$ , $R = 0.38$
<i>M.genitalium RplB</i>		$p < 2.26 \times 10^{-16}$ , $R = 0.88$	$p < 2.26 \times 10^{-16}$ , $R = 0.83$
<i>HOXC6</i>			$p < 2.26 \times 10^{-16}$ , $R = 0.73$
	<b><i>SLC12A1</i></b>	<b><i>UPK2</i></b>	
<i>SPINK1</i>	$p < 2.26 \times 10^{-16}$ , $R = 0.64$	$p < 2.26 \times 10^{-16}$ , $R = 0.78$	
<i>SLC12A1</i>		$p < 2.26 \times 10^{-16}$ , $R = 0.62$	

Needleman-Wunsch alignment of the capture and reporter probes for *HOXC6*, *ERG 5'*, *M.genitalium* RplA and *M.genitalium* RplB gave low percentage alignments and scores with each other. These scores were similar to alignments with three randomly selected NanoString probes (which were selected for a control comparison) that showed no expression correlation; *ALAS1*, *KLK2* and *KLK3*. BLAT analysis detected some homology between *M.genitalium* RplA reporter probe sequence and non-coding sequences on human ChrX, whilst *M.genitalium* RplB capture probe hits non-coding sequence on Chr10. Both *HOXC6* and *ERG 5'* capture and reporter probes only had sequence homologies with their own encoding gene sequences and nowhere else in the genome. These analyses suggest that cross-hybridisation is not likely to be the cause of their correlation.



**Figure 3.8** Correlation plots between data for probes: *M.genitalium* RplA, *M.genitalium* RplB, *HOXC6* and *ERG 5'*.

In Cluster 4, one transcript (*Spink1*) is known to be associated with PCa while the other two are tissue specific controls; *UPK2* is a bladder specific marker and *SLC12A1* is a kidney specific marker. It is understandable to see some correlation between the non-prostate tissue specific markers, as the proportion of these would result from the proportion of cells that are not from the prostate. The correlation between *UPK2* and *SLC12A1* data, whilst significant is not strong enough to suggest that they are cross hybridising ( $p < 2.26 \times 10^{-16}$ ,  $R = 0.62$ ) (Figure 3.9). *UPK2* and *SPINK1* correlate strongly ( $p < 2.26 \times 10^{-16}$ ,  $R = 0.78$ ), whereas *SLC12A1* correlation with *SPINK1* is weaker ( $p < 2.26 \times 10^{-16}$ ,  $R = 0.64$ ) (Figure 3.9). All three probes have similar signal strength also, ranging ~0 to ~15 (Figure 3.9). Needleman-Wunsch alignment of the capture and reporter probes for *SPINK1*, *SLC12A1* and *UPK2* gave low percentage alignments and scores with each other. These scores were similar to those of three randomly selected NanoString control probes that showed no expression correlation; *ALAS1*, *KLK2* and *KLK3*. Furthermore, BLAT analysis detected no other sites of homology in the human genome for *SPINK1* probe sequences, whilst both *UPK2* and *SLC12A1* reporter probes had one partial match each: *CTNNA3* (Chr 10) and *FLRT2* (Chr 14), respectively. The capture probes for *UPK2* and *SLC12A1* also had no alternative sites of homology in the human genome. This suggests that the probes are not cross-hybridising to each others target probes.

## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

It is possible that some of these probes in the two clusters are cross-hybridising (and then it is of course possible that at least one is a true representation for that probe) or that there is a clinical reasoning for their correlation. For this reason, I have included all of these probes in the subsequent analyses but any identification of their significance in clinical comparisons or clustering should be taken with caution.

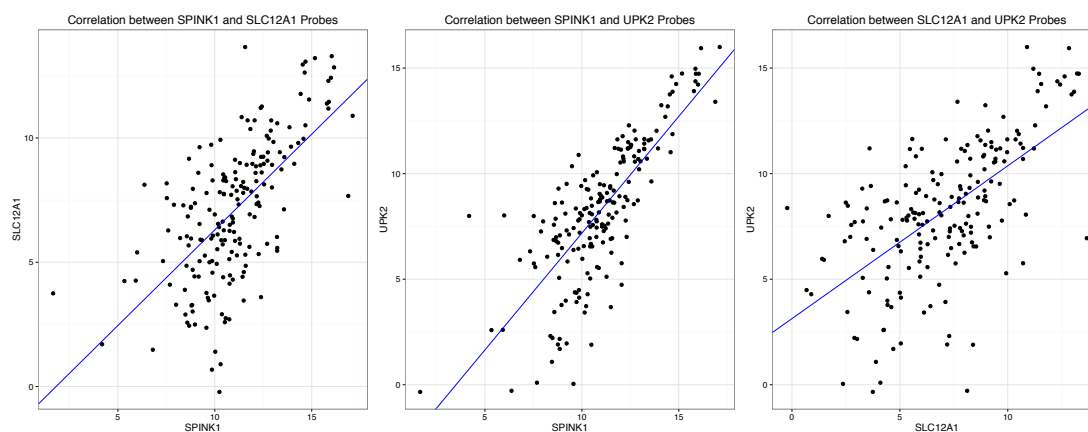


Figure 3.9 Correlation plots for a second group of probes that correlate: *SPINK1*, *SLC12A1* and *UPK2*. All correlate with  $p < 2.26 \times 10^{-16}$  and  $R < 0.6$ .

### **3.4 Identification of Prostate and Cancer Specific Transcripts and**

#### **DRE relevance**

##### **3.4.1 Kallikrein identification**

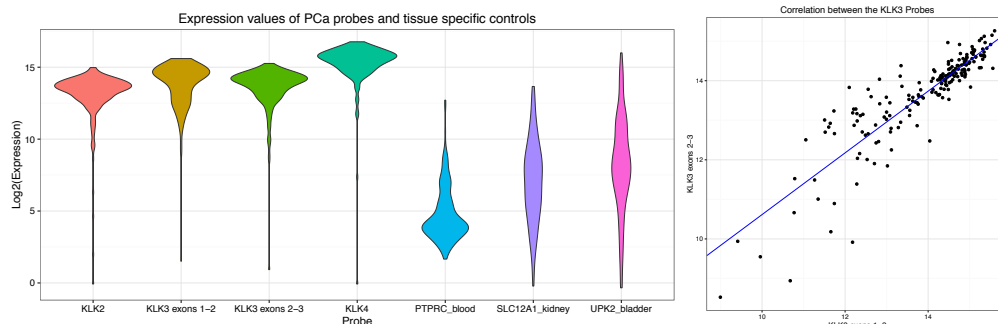
NanoString median signals for the *KLK2*, *KLK3* exons 1-2, *KLK3* exons 2-3 and *KLK4* probes were at significantly higher levels than those for the control tissue probes for blood, kidney and bladder (*PTPRC*, *SLC12A1* and *UPK2* respectively) (Mann Whitney U test:  $p < 2.2 \times 10^{-16}$  in each case, Table 3.11, Figure 3.10).



**Table 3.11 Median expression values for kallikreins (prostate specific transcripts) and other tissue markers.**

<i>Probe</i>	<i>Tissue</i>	<i>Log2 Median expression</i>
<i>KLK2</i>	<i>Prostate</i>	<i>13.49</i>
<i>KLK3 exons 1-2</i>	<i>Prostate</i>	<i>14.35</i>
<i>KLK3 exons 2-3</i>	<i>Prostate</i>	<i>14.02</i>
<i>KLK4</i>	<i>Prostate</i>	<i>15.59</i>
<i>PTPRC</i>	<i>Blood</i>	<i>4.08</i>
<i>SLC12A1</i>	<i>Kidney</i>	<i>7.24</i>
<i>UPK2</i>	<i>Bladder</i>	<i>8.15</i>

The kallikreins are prostate specific transcripts<sup>205</sup>; identification of *KLK2*, *KLK3* and *KLK4* at higher levels in the blood, kidney and bladder specific markers along with the RNA yield of post radical prostatectomy samples (section 3.4.4) suggest that a good proportion of the material captured is in fact from the prostate. Additionally, both the *KLK3* probes (exons 1-2 and exons 2-3) have a strong correlation ( $p < 2.2 \times 10^{-16}$ ,  $R = 0.89$ , Figure 3.10B).



**Figure 3.10 A) Kallikreins are observed at higher expression levels than the blood, kidney and bladder specific markers in the NanoString data. B) Correlation between the two KLK3 probes is strong.**

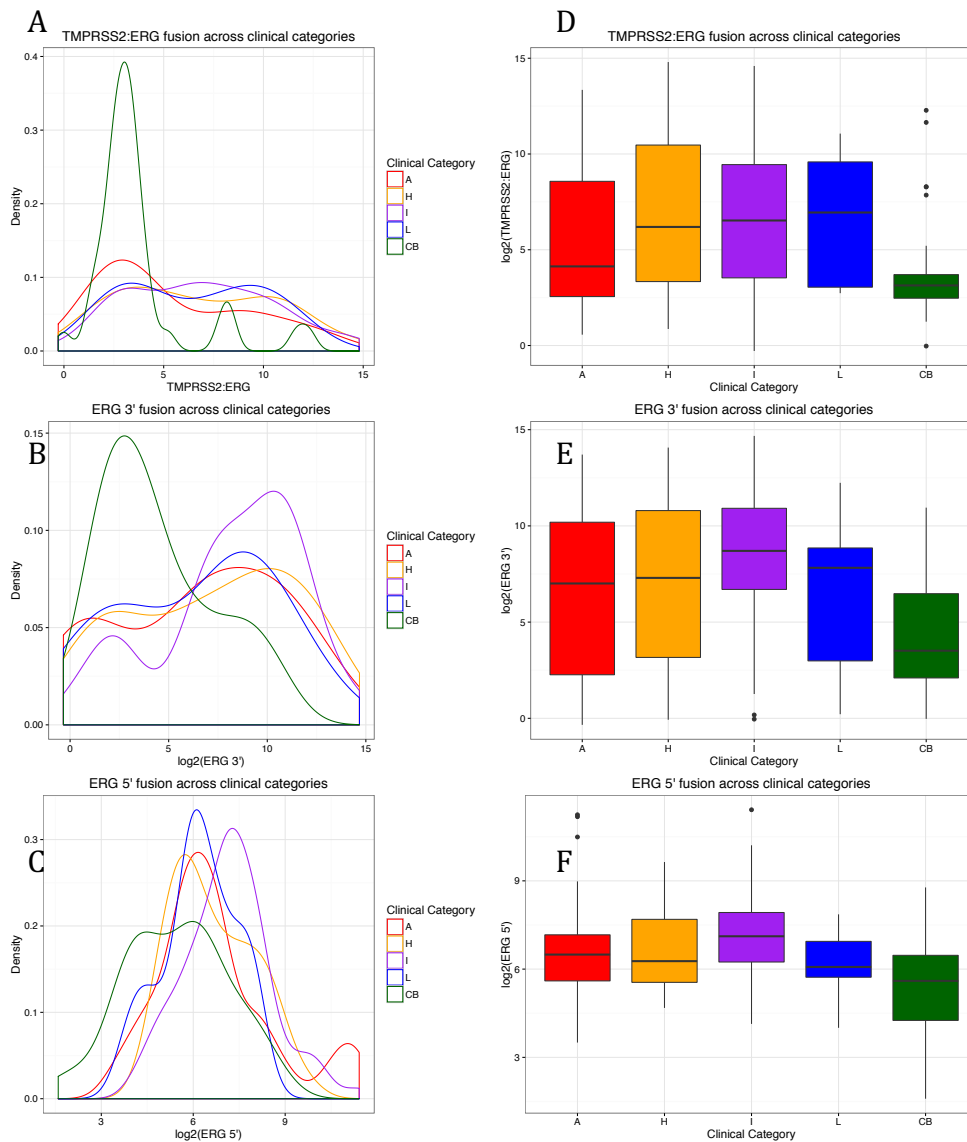
### 3.4.2 *TMPRSS2:ERG* Identification

*TMPRSS2:ERG* fusions, and alleviated *ERG* 3' and *ERG* 5' expression are found in PCa, and this is observed in the Nanostring data where a significant difference is 133

## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

observed between cancer vs CB (Mann Whitney U; *TMPRSS2:ERG*:  $p < 2.2 \times 10^{-16}$ ;  $W = 6179$ , *ERG 3'*:  $p < 2.2 \times 10^{-16}$ ;  $W = 6105$ , and *ERG 5'*:  $p < 2.2 \times 10^{-16}$ ;  $W = 6253$ ; Error! Reference source not found.). The density plots for *TMPRSS2:ERG* and *ERG3'* have two peaks which would be compatible with an on/off pattern of a gene fusion (Error! Reference source not found.). Approximately 50% of the samples from men with cancer have detectable *TMPRSS2:ERG* fusions which is in agreement with the literature (section 1.4.6). The *ERG5'* probe, which is not part of the *TMPRSS2:ERG* fusion transcript, does not follow this pattern. The *ERG 5'* probe was also identified as having potential cross hybridisation (section 3.3.7).

When dicotomised (using the optimal threshold 4.93 identified by the Brent method), *TMPRSS2:ERG* expression had a significant association with clinical category (chi-square test,  $\chi^2 = 37.82$ ,  $p = 4.1 \times 10^{-07}$ ).



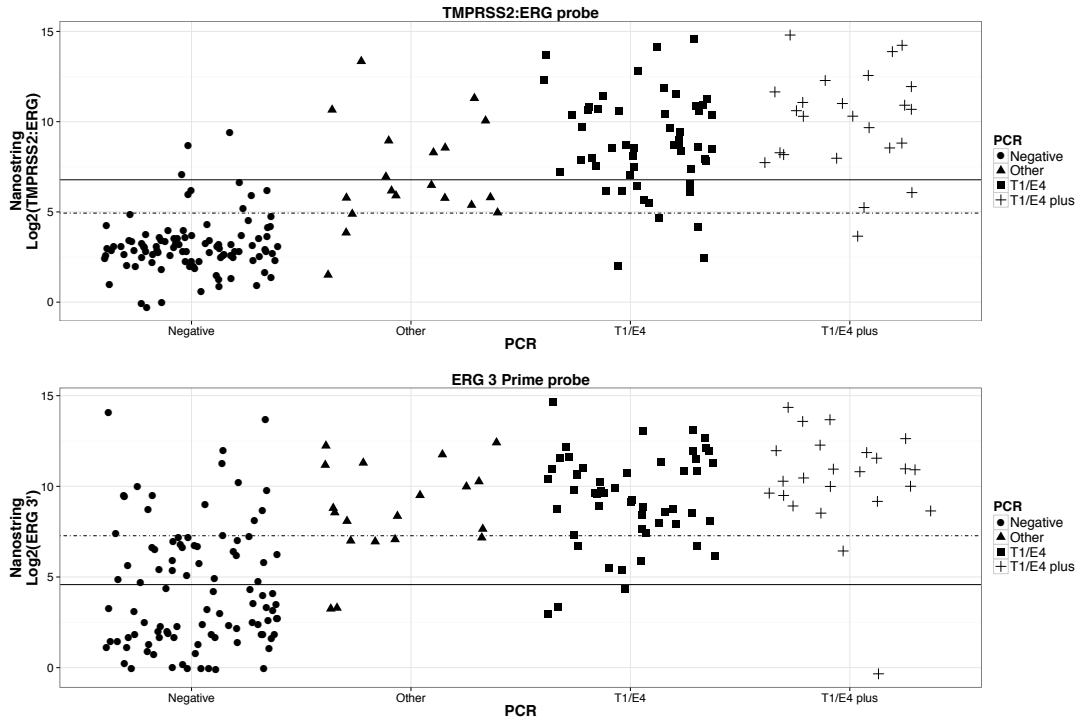
**Figure 3.11** A) Density plot for *TMPRSS2:ERG* expression coloured by clinical category. Generally, two peaks are seen suggesting an on/off pattern of expression. B) Density plot for *ERG 3'* expression coloured by clinical category. Again, two bumps are generally seen suggesting an on/off pattern. C) Density plot for *ERG 5'* expression coloured by clinical category. No observable on/off pattern can be seen. D) Box plot showing spread of *TMPRSS2:ERG* expression across clinical categories. Higher expression is observed in cancer than benign. E and F) Box plots showing expression of *ERG 3'* and *ERGS'* respectively across clinical categories. Median expression is Higher in cancer than benign.

## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

*TMPRSS2:ERG* fusions were identified by NanoString and nested RT-PCR (section 2.1.5 and section 2.1.6) (Figure 3.12). PCR was capable of identifying not only T1/E4 fusion transcripts but also fusions involving other *TMPRSS2* and *ERG* exons (section 1.4.6), including T1/E5, T1/E6 and T2/E2 amongst others. *TMPRSS2:ERG* PCR data were divided into four groups: i) T1/E4 fusions (“T1/E4”), ii) those with T1/E4 plus other fusion types (“T1/E4 plus”), iii) those with only non-T1/E4 products (“other”) and iv) those where no fusions were identified (“negative”).

The minimum curve threshold was calculated from NanoString expression density plots. A cut off of 6.78 for the *TMPRSS2:ERG* probe, showed 97% correlation for the PCR negatives (95/98 are classed as negative in both), 79% accuracy for T1/E4 only fusions (41/52 are classed as positive), 88% accuracy for T1/E4 plus other fusions (21/24 are classed as positive), and 42% accuracy for other fusions (8/19 are classed as positive). The NanoString *TMPRSS2:ERG* probe had been designed to specifically pick up the T1/E4 fusion, and so the poor accuracy for detecting other fusions was expected. Using the optimal threshold identified by the Brent method for the *TMPRSS2:ERG* probe (4.93), the on/off pattern compared with the PCR results, showed an improved and significant association (chi-square test  $\chi^2 = 131.6$ ,  $p < 2.2 \times 10^{-16}$ ).

The *ERG3'* NanoString signal correlates well to the *TMPRSS2:ERG* PCR positive samples for both T1/E4 fusion and non-T1/E4 PCR products. However there are a proportion of the PCR negative samples that also have high *ERG3'* NanoString signals; this would appear to indicate that *ERG3'* has been overexpressed via an alternate mechanism (section 1.4.6).

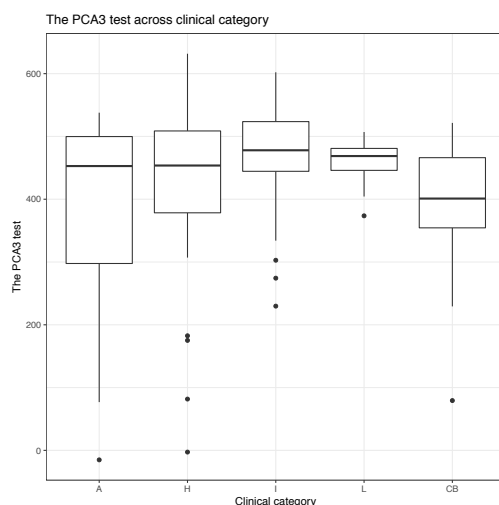


**Figure 3.12** Detection of *TMPRSS2:ERG* by NanoString probes for *TMPRSS2:ERG* (upper) and *ERG3'* (lower) versus PCR detection of *TMPRSS2:ERG* transcripts. T1/E4 indicates a *TMPRSS2* ex1/*ERG*ex4 fusion transcript, 'Other' indicates a different fusion transcript, 'Plus' indicates a mixture of T1/E4 and other transcripts. The dotted lines are the optimal thresholds (4.93 for *TMPRSS2:ERG* and 7.28 for *ERG3'*) calculated using the Brent method, similarly the solid line is the min curve of a density plot (6.78 and 4.58 for *TMPRSS2:ERG* and *ERG3'* respectively) containing all of the *TMPRSS2:ERG* and *ERG* data.

These results suggest that a) NanoString is a sensitive and flexible method for detecting transcripts and b) that a proportion of the genetic material identified is coming from prostate cancer or HG-PIN.

### 3.4.3 PCA3 Test

The *PCA3* test (section 1.4.2) is the ratio of *PCA3* expression with *KLK3* expression in whole urine, and is approved clinically to predict whether a second biopsy will be cancer positive after an initial negative biopsy. The *PCA3* score calculated from the NanoString data shows a significantly increased expression in PCa compared with non-PCa samples (Mann-Whitney U test:  $p < 2.2 \times 10^{-16}$ , Figure 3.13), but was no evidence for a significantly difference between different clinical categories of PCa ( $p < 0.05$ ; Kruskal-Wallis rank sum).



**Figure 3.13** Nanostring *PCA3* score calculation (*PCA3* divided by *KLK3* multiplied by 1000 as per the usual *PCA3* score (section 1.4.2)). The *PCA3* score is significantly increased in PCa samples compared to those with no clinical evidence of PCa (CB). However, there is no significant difference between the intra-clinical categories of PCa. The uPM3™ assay has shown to be able to detect PCa from non-PCa samples. The NanoString probes have shown to follow this same pattern.

### 3.4.4 RNA yield, clinical group and DRE

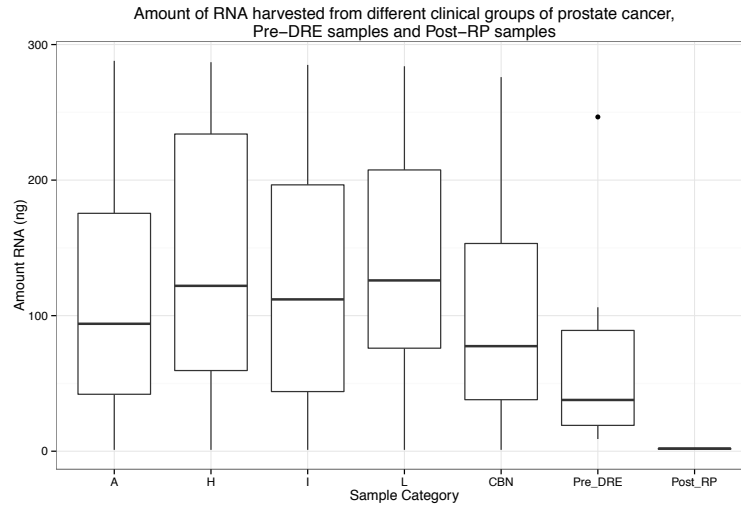
Digital rectal examination (DRE, section 1.3.3) has proven to increase the efficacy of the *PCA3* test (section 1.4.2). It is hypothesised that digital compression on the prostate encourages secreted biomarkers in the gland to flow towards the urethra. Four patient

pairs for pre- and post-DRE urine samples were added to NanoString to see how the transcript levels varied within patients (Figure 3.15). First-void urine post-DRE had higher median RNA yields than non-DRE samples (Figure 3.14). RNA yield is significantly higher in post-DRE collection of localised PCa samples compared to pre-DRE samples ( $p = 0.04$ , Mann Whitney U test) and prostatectomy samples ( $p = 0.01$ , Mann Whitney U test). As seen previously there were also increased numbers of prostate derived transcripts (section 3.4.1) and PCa derived transcripts (section 3.4.2) on post-DRE samples. Overall, the post-DRE samples had 0.178 log<sub>2</sub> fold increased expression of all transcripts compared to the pre-DRE collected samples ( $p = 1.854 \times 10^{-10}$ , paired Man Whitney U test). The median of the sample pairs individually varied with the pre- or post-DRE, however, the post-DRE sample always showed a lower IQR (Figure 3.15).

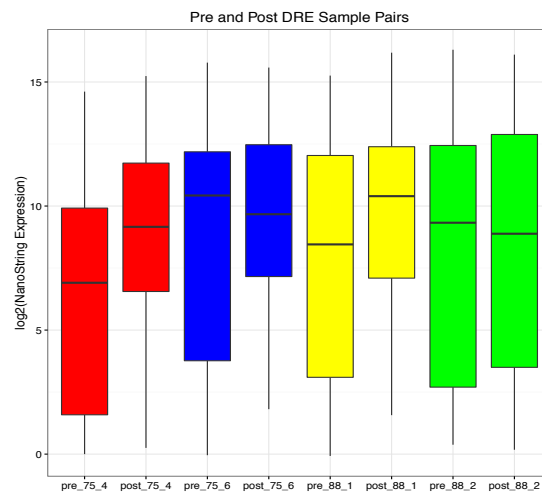
The urine taken from three patients who had previously undergone radical prostatectomy (post-RP) had very low amounts of RNA collected (0.8-2ng) from their urine samples. This suggests that the majority of the EV RNA is likely to have originated in the prostate (Figure 3.14).

The median RNA yields for advanced PCa patients are not significantly lower than for localised-PCa patients ( $p < 0.05$ , Mann Whitney U test, Figure 3.14). The RNA yields for benign samples are observably (Figure 3.14) and significantly lower compared to localised PCa patient samples ( $p = 0.02$ , Mann Whitney U test).

## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY



**Figure 3.14** Most of the transcripts detected are from the prostate; DRE boosts transcript level detection and post radical prostatectomy patients offer very low signals in their samples. Samples  $n = 389$ . The advanced (A), high-risk (H), intermediate risk (I), low-risk (L) and no evidence of clinical PCa (CB) samples were taken post-DRE. Pre-DRE and post-RP urine samples have been taken without DRE.



**Figure 3.15** The NanoString probe expression distribution of four patient paired samples (pre- and post-DRE).

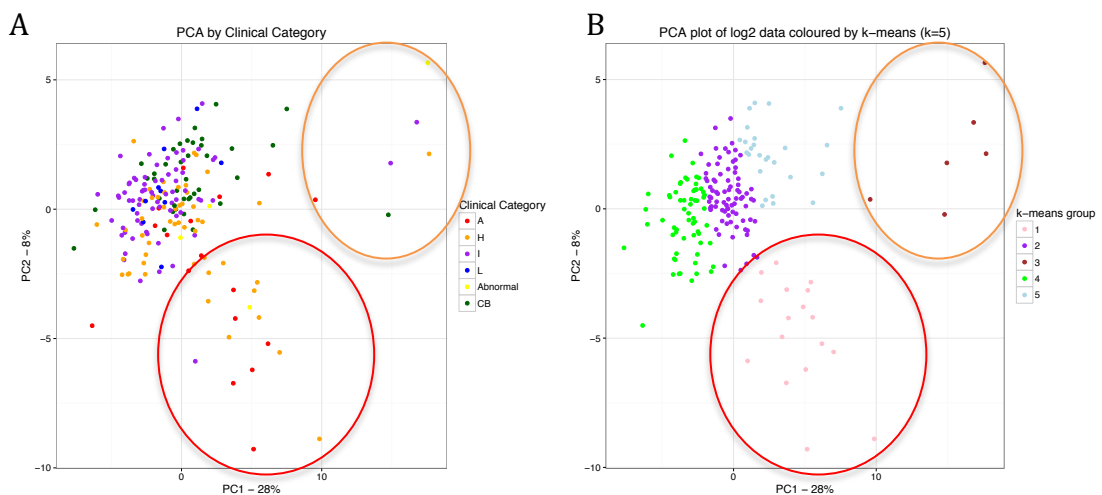


### 3.5 Clustering

#### 3.5.1 Principal Component Analysis (PCA) and *k*-means Clustering

PCA (section 2.5.1) can be utilised to visualise groups in the data; colouring data by clinical category can allow clusters of biological interest to be identified (Figure 3.16A). PCA analysis identified two outlying clusters, A and B, (Figure 3.16). Cluster A had 17 samples consisting mostly of advanced and higher risk samples (6 advanced, 9 high-risk, 1 intermediate risk and 1 abnormal sample). In contrast Cluster B consisted of 6 samples of varying clinical groups (1 advanced, 1 high risk, 2 intermediate risk, 1 abnormal and 1 CB).

**Figure 3.16** PCA plots coloured by A) clinical category and B) *k*-means to identify cluster cut-offs.



**Cluster A** shown by red circle. **Cluster B** shown by orange circle.

#### 3.5.2 Hierarchical Clustering

Hierarchical clustering was performed with an agglomerative approach (section 2.5.2). This showed that samples in Clusters A and B belonged to separate trees to the majority of other samples (Figure 3.17A). Fifteen of the samples belonging to Cluster A form a separate tree, whilst 5 of the 6 samples belonging to Cluster B also form a separate tree with 2 other samples. There was one significant cluster identified by Pvcust (section

2.5.2.1), which contains the bulk of the samples, but does not include the majority of Cluster A or Cluster B samples (Figure 3.17B).

### 3.5.3 Cluster A

Cluster A (identified by PCA and *k*-means clustering (Section 3.5.1) and supported by hierarchical clustering (Section 3.5.2) is predominantly made up by advanced and high-risk samples (6/17 and 9/17, respectively). It has significant over-representation of advanced and high-risk samples (Table 3.12) and there are twenty-three significant differentially expressed transcripts between cluster A and all other samples (Table 3.13). Analysis of the differential expressed transcript list with DAVID (section 2.7.1) identified PCa as an over represented KEGG pathway. This was due to the significantly lowered expression of *AR* and *KLK3* in Cluster A, however the over-representation was not significant at a 95% confidence level ( $p = 8.5 \times 10^{-02}$ ). Ten Gene Ontology (GO) biological processes were associated with the Cluster A defining transcripts (Table 3.14). As expected due to probe selection for involvement in PCa these biological processes were associated with cancer. However, different GO biological processes were identified using all of the transcripts applied to NanoString (Table 3.15). Thus suggesting there is a difference in biological processes involved specifically within cluster A.

RNA amount (ng) extracted is significantly lower in Cluster A compared to all other samples (not including Cluster B), (Table 3.12). Cluster A also has a significantly lower amplification yield, as well as a lower median probe value (Table 3.12). The cartridge number is also significant between members of Cluster A (Table 3.12). However, Scanner ID is not significant (Table 3.12).

Further investigation into the cartridges involved, showed there was no significant differences between the median probe values of these cartridges compared to others, or between the Cluster A samples and non-Cluster A samples on these cartridges (Table 3.12). This suggests that the cartridge is not a factor to why Cluster A may be presenting itself. However, RNA extraction amount, amplification yield and median probe value all seem important in the clustering. Especially as 21/23 significant

## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

differentially expressed transcripts between Cluster A and all other samples (Table 3.13) have lower expression in Cluster A.

Lower RNA yields are observed in a fraction of advanced patients' samples, This is hypothesised to be due to a reduction in tumour microvesicles being harvested in the urine due to: a) efficiency of DRE: the surface of a normal prostate can be depressed by 1cm, however prostates containing advanced/higher grade tumours are commonly firm and not depressible. Samples from patients with advanced tumours are therefore more akin to non-DRE samples. b) Advanced tumours can also have fused glands, poorly formed lumen, and blind-ended lumen that no longer drain into the urethra<sup>206,207</sup>. The position of the advanced tumour within the prostate may also block access of tumour biomarkers from less advanced PCa foci from entering the urine. Thus, the percentage of the tumour that is advanced and its positioning within the prostate can affect the amount of RNA extracted, and the amounts of PCa associated transcripts identified.

**Table 3.12 Testing for Cluster A association to clinical and technical variables.**

<i>Variable</i>	<i>Test and metric</i>	<i>p - value</i>
<i>Clinical Category</i>	<b>Chi square: <math>\chi^2 = 20.29</math></b>	<b><math>p = 6.67 \times 10^{-6}</math></b>
<i>Amount RNA extracted (ng)</i>	<b>Mann-Whitney U test: <math>R = 2443</math></b>	<b><math>p = 4.9 \times 10^{-7}</math></b>
<i>Amplification yield (<math>\mu\text{g}</math>)</i>	<b>Mann-Whitney U test: <math>R = 2410</math></b>	<b><math>p = 3.1 \times 10^{-6}</math></b>
<i>Median probe value</i>	<b>Mann-Whitney U test: <math>R = 1904</math></b>	<b><math>p = 0.02</math></b>
<i>Cartridge</i>	<b>Chi-square test: <math>\chi^2 = 31.9</math></b>	<b><math>p = 0.01</math></b>
<i>Scanner ID</i>	<b>Chi-square test: <math>\chi^2 = 0.03</math></b>	<b><math>p = 0.9</math></b>
<i>Median probe value of Cluster A samples on cartridge 13 (n = 4) compared to other samples on cartridge 13 (n = 8)</i>	<b>Mann-Whitney U test: <math>R = 16</math></b>	<b><math>p = 1</math></b>
<i>Median probe value of samples on cartridge 13 (n = 12) compared to samples on all other cartridges (n = 168)</i>	<b>Mann-Whitney U test: <math>R = 1237</math></b>	<b><math>p = 0.3</math></b>
<i>Median probe value of Cluster A samples on cartridge 15 (n = 3) compared to other samples on cartridge 15 (n = 9)</i>	<b>Mann-Whitney U test: <math>R = 4</math></b>	<b><math>p = 0.1</math></b>
<i>Median probe value of samples on cartridge 15 (n = 12) compared to samples on all other cartridges (n = 168)</i>	<b>Mann-Whitney U test: <math>R = 1074</math></b>	<b><math>p = 0.8</math></b>

It should also be remembered that the vast majority of the NanoString probes were selected due to overexpression in tumour tissue. Thus, it is significant that the expression patterns for Cluster A are more than a general loss of tumour biomarkers as my analyses mark them as a group distinct from the other prostate samples. The only two probes not showing a significant up-regulation in the Cluster A samples are the kidney and bladder controls (Table 3.13).

It is hypothesised that the factors identified as technical issues (RNA amount extracted, amplification yield and median probe value) associated with Cluster A are due to these biological reasons and thus it is important to keep Cluster A's samples within future analyses.

**Table 3.13** Transcripts significantly associated ( $p < 0.05$ ) with Cluster A via Mann-Whitney U test after using Hochberg multiple testing correction.

<i>Transcript</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>Log2 Fold Change</i>
<i>DLX1</i>	$6.6 \times 10^{-04}$	$2.3 \times 10^{-02}$	-1.504
<i>Timp4</i>	$1.2 \times 10^{-04}$	$4.7 \times 10^{-03}$	-1.292
<i>AR exon 9</i>	$2.3 \times 10^{-06}$	$3.2 \times 10^{-04}$	-1.263
<i>MMP26</i>	$7.4 \times 10^{-06}$	$3.2 \times 10^{-04}$	-1.126
<i>CLU</i>	$6.4 \times 10^{-05}$	$2.6 \times 10^{-03}$	-1.017
<i>UPK2</i>	$5.1 \times 10^{-10}$	$2.8 \times 10^{-08}$	0.798
<i>SLC12A1</i>	$2.5 \times 10^{-09}$	$1.4 \times 10^{-07}$	0.736
<i>PSGR</i>	$2.6 \times 10^{-08}$	$1.3 \times 10^{-06}$	-0.555
<i>CDC20</i>	$3 \times 10^{-04}$	$1.1 \times 10^{-02}$	-0.543
<i>SPINK1</i>	$1.5 \times 10^{-10}$	$8.3 \times 10^{-09}$	-0.497
<i>GOLM1</i>	$1.4 \times 10^{-05}$	$6 \times 10^{-04}$	-0.485
<i>PCA3</i>	$1.8 \times 10^{-04}$	$6.9 \times 10^{-03}$	-0.456
<i>SERPINB5</i>	$3 \times 10^{-04}$	$1.1 \times 10^{-02}$	-0.287
<i>KLK3 exons 2-3</i>	$4.2 \times 10^{-10}$	$2.4 \times 10^{-08}$	-0.251
<i>KLK3 exons 1-2</i>	$4.9 \times 10^{-08}$	$2.4 \times 10^{-06}$	-0.235
<i>FOLH1</i>	$6 \times 10^{-08}$	$2.9 \times 10^{-06}$	-0.214
<i>B2M</i>	$1.1 \times 10^{-04}$	$4.2 \times 10^{-03}$	-0.207
<i>AR exons 4-8</i>	$7.9 \times 10^{-07}$	$3.6 \times 10^{-05}$	-0.186
<i>STEAP2</i>	$1.2 \times 10^{-08}$	$6.2 \times 10^{-07}$	-0.183
<i>KLK2</i>	$1.2 \times 10^{-08}$	$6.2 \times 10^{-07}$	-0.174
<i>KLK4</i>	$6.1 \times 10^{-09}$	$3.2 \times 10^{-07}$	-0.132
<i>STEAP4</i>	$8.5 \times 10^{-06}$	$3.6 \times 10^{-04}$	-0.129
<i>PPAP2A</i>	$3.5 \times 10^{-07}$	$1.7 \times 10^{-05}$	-0.128

**Table 3.14 Gene Ontology (GO) over-represented biological processes in Cluster A's significantly associated transcript list via DAVID.**

<i>Term</i>	<i>Count (%)</i>	<i>Transcripts</i>	<i>p - value</i>	<i>Adjusted p - value</i>
<i>Proteolysis</i>	<b>7 (3.9)</b>	<b><i>FOLH1, KLK2, KLK3, KLK4, CLU, MMP26, CDC20</i></b>	<b><math>8.9 \times 10^{-04}</math></b>	<b><math>2.5 \times 10^{-01}</math></b>
<i>Iron ion transport</i>	<b>2 (1.1)</b>	<b><i>STEAP4, STEAP2</i></b>	<b><math>3.4 \times 10^{-02}</math></b>	<b>1</b>
<i>Androgen receptor signalling pathway</i>	<b>2 (1.1)</b>	<b><i>AR, PPAP2A</i></b>	<b><math>4.2 \times 10^{-02}</math></b>	<b><math>9.9 \times 10^{-01}</math></b>
<i>Response to organic substance</i>	<b>4 (2.2)</b>	<b><i>AR, TIMP4, STEAP2, B2M</i></b>	<b><math>5.0 \times 10^{-02}</math></b>	<b><math>9.8 \times 10^{-01}</math></b>
<i>Steroid hormone receptor signalling pathway</i>	<b>2 (1.1)</b>	<b><i>AR, PPAP2A</i></b>	<b><math>6.6 \times 10^{-02}</math></b>	<b><math>9.9 \times 10^{-01}</math></b>
<i>Response to hormone stimulus</i>	<b>3 (1.7)</b>	<b><i>AR, TIMP4, STEAP2</i></b>	<b><math>6.9 \times 10^{-02}</math></b>	<b><math>9.8 \times 10^{-01}</math></b>
<i>Transition metal ion transport</i>	<b>2 (1.1)</b>	<b><i>STEAP4, STEAP2</i></b>	<b><math>8.1 \times 10^{-02}</math></b>	<b><math>9.8 \times 10^{-01}</math></b>
<i>Response to endogenous stimulus</i>	<b>3 (1.7)</b>	<b><i>AR, TIMP4, STEAP2</i></b>	<b><math>8.1 \times 10^{-02}</math></b>	<b><math>9.7 \times 10^{-01}</math></b>
<i>Intracellular receptor-mediated signalling pathway</i>	<b>2 (1.1)</b>	<b><i>AR, PPAP2A</i></b>	<b><math>8.5 \times 10^{-02}</math></b>	<b><math>9.6 \times 10^{-01}</math></b>
<i>Response to molecule of bacterial origin</i>	<b>2 (1.1)</b>	<b><i>TIMP4, B2M</i></b>	<b><math>9.7 \times 10^{-02}</math></b>	<b><math>9.6 \times 10^{-01}</math></b>

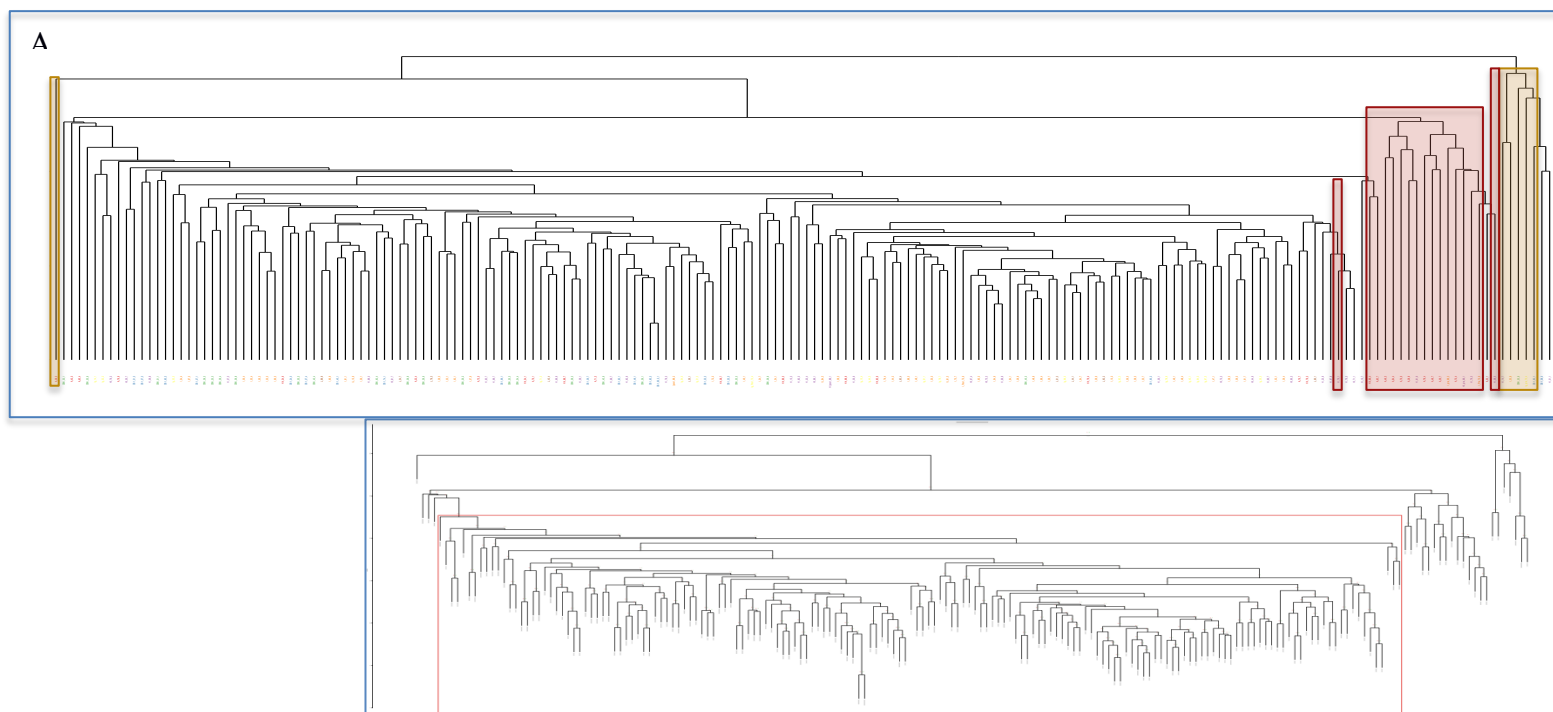
**Table 3.15 Gene Ontology (GO) over-represented biological processes in all of the transcripts used on NanoString via DAVID.**

<i>Term</i>	<i>Count (%)</i>	<i>Transcripts</i>	<i>p - value</i>	<i>Adjusted p - value</i>
-------------	------------------	--------------------	------------------	---------------------------

CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

<i>Proteolysis</i>	7 (0.1)	<i>ANPEP, HPN, KLK2, KLK3, KLK4, MMP26, TMPRSS2</i>	$9.6 \times 10^{-04}$	$3.5 \times 10^{-01}$
<i>Iron ion transport</i>	3 (0)	<i>STEAP2, STEAP4, B2M</i>	$2.2 \times 10^{-03}$	$3.9 \times 10^{-01}$
<b>Negative regulation of neuron apoptotic process</b>	4 (0.1)	<i>BRAF, DLX1, MDK, TERT</i>	$3.5 \times 10^{-03}$	$4.1 \times 10^{-01}$
<b>Ferric iron import into cell</b>	2 (0)	<i>STEAP2, STEAP4</i>	$1.2 \times 10^{-02}$	$7.3 \times 10^{-01}$
<b>Cell cycle</b>	4 (0.1)	<i>AURKA, CDC20, CDKN3, FOXM1</i>	$1.4 \times 10^{-02}$	$7.1 \times 10^{-01}$
<b>Copper ion import</b>	2 (0)	<i>STEAP2, STEAP4</i>	$1.6 \times 10^{-02}$	$7 \times 10^{-01}$
<b>Positive regulation of gene expression</b>	4 (0.1)	<i>BRAF, AR, AGR2, HPN</i>	$2.3 \times 10^{-02}$	$7.7 \times 10^{-01}$
<b>Positive regulation of stem cell proliferation</b>	2 (0)	<i>PTPRC, TERT</i>	$2.8 \times 10^{-02}$	$7.9 \times 10^{-01}$
<b>Positive regulation of transcription, DNA-templated</b>	5 (0.1)	<i>TBP, AR, CAMKK2, FOXM1, MDK</i>	$3.1 \times 10^{-02}$	$7.9 \times 10^{-01}$
<b>Response to cadmium ion</b>	2 (0)	<i>B2M, TERT</i>	$5.7 \times 10^{-02}$	$9.2 \times 10^{-01}$
<b>Negative regulation of endothelial cell apoptotic process</b>	2 (0)	<i>BRAF, TERT</i>	$6.3 \times 10^{-02}$	$9.3 \times 10^{-01}$
<b>Embryonic skeletal system development</b>	2 (0)	<i>DLX1, HOXC6</i>	$6.7 \times 10^{-02}$	$9.2 \times 10^{-01}$
<b>Protein phosphorylation</b>	4 (0.1)	<i>BRAF, ERG, AURKA, CAMKK2</i>	$8.9 \times 10^{-02}$	$9.6 \times 10^{-01}$
<b>Cell differentiation</b>	4 (0.1)	<i>ERG, ANPEP, AGR2, MDK</i>	$9.1 \times 10^{-02}$	$9.5 \times 10^{-01}$
<b>Response to peptide hormone</b>	2 (0)	<i>BRAF, Timp4</i>	$9.7 \times 10^{-02}$	$9.5 \times 10^{-01}$





**Figure 3.17 Hierarchical clustering provides further evidence for Cluster A and B identification. A) Samples belonging to Cluster A and B are shown in red and yellow boxes, respectively. B) Clusters with significant AU  $p$ -values are encapsulated within a red box. Both Cluster A and B are not included within this main.**

### 3.5.4 Cluster B

Cluster B (identified by PCA and *k*-means clustering section 3.5.1 and supported by hierarchical clustering section 3.5.2) contains six samples and are not associated with clinical factors (Chi square:  $\chi^2 = 9.7$ ,  $p = 0.08$ ), suggesting that Cluster B could be associated with technical artefacts: The amount of RNA extracted was lower in Cluster B (Mann-Whitney U test:  $R = 811$ ,  $p = 0.008$ ); the total amount of cDNA from amplification was also lower (Mann-Whitney U test:  $R = 808$ ,  $p = 0.01$ ) and thus was the median probe value (Mann-Whitney U test:  $R = 14196$ ,  $p < 2.2 \times 10^{-16}$ ). Cartridge and Scanner ID were both not significantly associated with Cluster B (Chi square:  $\chi^2 = 0.67$ ,  $p = 0.88$ , and  $\chi^2 = 2.67$ ,  $p = 0.1$ , respectively). It is therefore, unlikely that there is biological reasoning to this cluster.

### 3.5.5 Latent Process Decomposition (LPD)

LPD (section 2.5.5) was performed on 187 of the samples (with M\_19\_5, LNCAP, and the five samples in cluster B removed) and 51 of the transcripts (with *FOXMI* and the six bacterial genes removed) to identify the optimal number of groups and an assign a probability of membership for each group for each sample.

The modelling and estimation stage suggested that there were four clusters, with a sigma parameter of -1. LPD analysis was performed 100 times with these parameters and samples were associated with a probability to each group (Table 3.16, Figure 3.18A, Figure 3.18B). LPD 1 consisted mainly of high-risk samples ( $\chi = 16.5$ ,  $p = 0.01$ ), whilst LPD 4 consisted mostly of advanced and high-risk samples ( $\chi = 29.44$ ,  $p = 5 \times 10^{-05}$ ). Both LPD 2 and 3 contain a mixed representation of clinical category. LPD group 2 consists mostly of the intermediate risk samples ( $\chi = 29.44$ ,  $p = 5 \times 10^{-05}$ ), it holds 66% of the intermediate and low-risk samples. LPD group 3 holds 57% of the

## CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

benign samples though they are not significantly represented within the group ( $\chi = 11.82, p = 0.07$ ).

Clinical category is significantly associated with LPD process, ( $\chi = 65.47, p = 2.83 \times 10^{-08}$ , Table 3.16). PSA level is significantly associated with LPD process, with higher values in LPD groups 1 and 4 (ANOVA,  $p = 7.53 \times 10^{-06}$ , Figure 3.18C), as well as Gleason score ( $\chi = 85.38$  and  $p = 9.98 \times 10^{-10}$ ); a higher Gleason score appears to be associated with LPD process 4, whilst processes 2 and 3 have much lower (Figure 3.18D). Age is also significantly associated with LPD process (ANOVA,  $p = 0.002$ ), with a higher age present in LPD process 4 (Figure 3.18E).

Alternative analysis was performed using NbClust (section 2.5.3), which identified three clusters as the optimal number of clusters in the data, and *k*-means with PCA (section 2.5.3) was used to identify which samples belonged to which cluster (Figure 3.18E). These clusters showed high overlap with the four clusters identified by LPD (Figure 3.18F), providing further evidence that this clustering is reliable.

**Table 3.16 Composition of sample type in each LPD cluster (Cluster B samples and bacterial probes removed). Chi-square test:  $p = 2.8 \times 10^{-08}$ ,  $X = 65.47$ .**

	<i>Total Number of Samples</i>	<i>Number of aggressive samples (A and H risk)</i>	<i>Number of lower-risk cancer samples (L and I risk)</i>	<i>Number of Abnormal samples (S)</i>	<i>Number of CB samples</i>
<i>LPD1</i>	<b>8</b>	<b>6</b>	<b>1</b>	<b>0</b>	<b>1</b>
<i>LPD2</i>	<b>79</b>	<b>20</b>	<b>53</b>	<b>0</b>	<b>6</b>
<i>LPD3</i>	<b>55</b>	<b>14</b>	<b>18</b>	<b>2</b>	<b>21</b>
<i>LPD4</i>	<b>17</b>	<b>13</b>	<b>1</b>	<b>1</b>	<b>2</b>
<i>LPD NA</i>	<b>26</b>	<b>12</b>	<b>7</b>	<b>0</b>	<b>7</b>

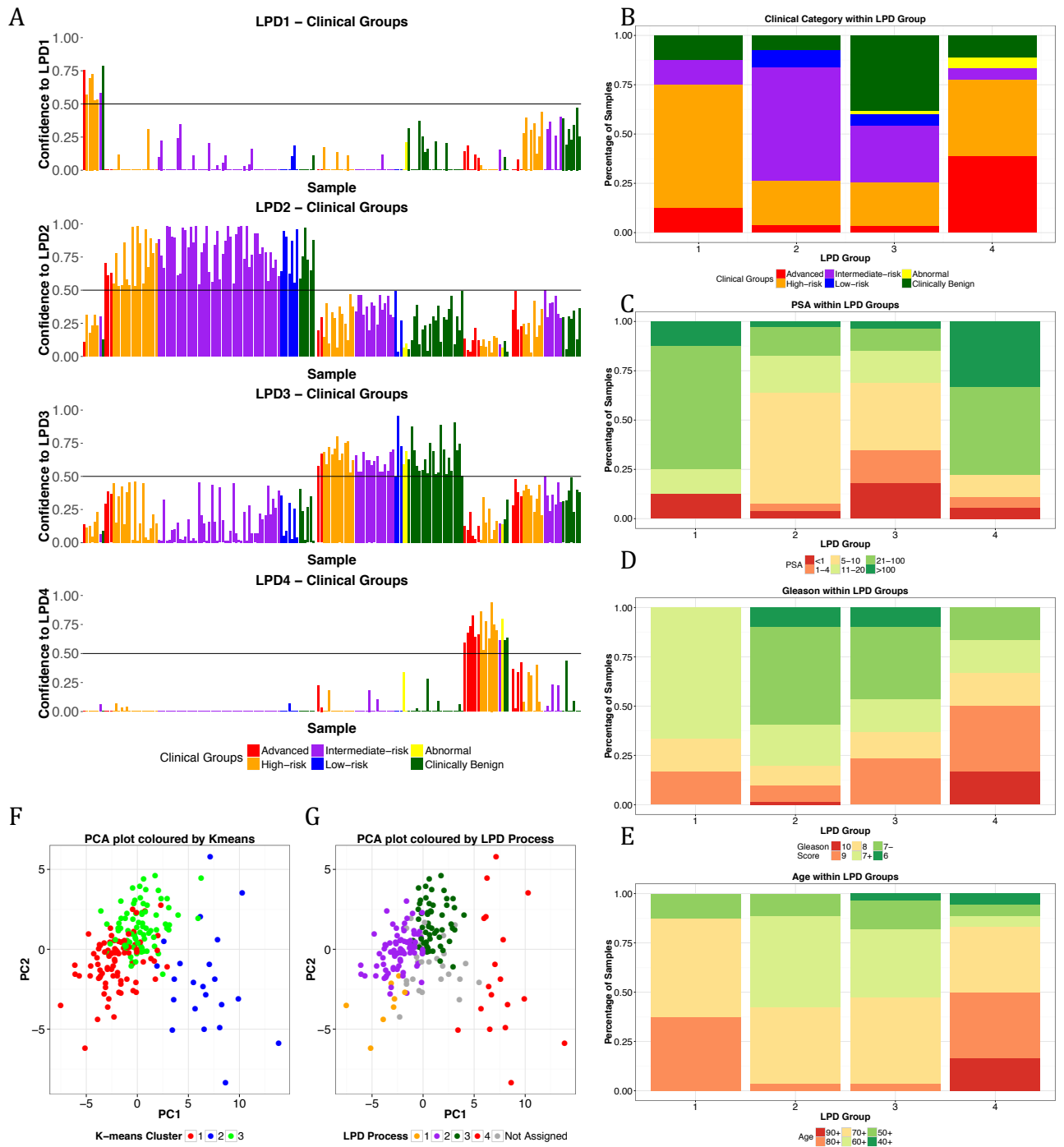


Figure 3.18 A) LPD group bar charts B,C,D,E) Clinical distribution, PSA, Gleason score and age without LPD group, respectively. F,G) PCA plots for *k*-means and LPD clustering comparison.

CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

**Table 3.17** Transcripts significantly different between each LPD group members and those that are not. These are the transcripts that define each LPD cluster.

<i>LPD Process 1</i>			<i>LPD Process 2</i>			<i>LPD Process 3</i>			<i>LPD Process 4</i>		
<i>Gene</i>	<i>p-value</i>	<i>Fold Change</i>	<i>Gene</i>	<i>p-value</i>	<i>Fold Change</i>	<i>Gene</i>	<i>p-value</i>	<i>Fold Change</i>	<i>Gene</i>	<i>p-value</i>	<i>Fold Change</i>
<i>HOXC6</i>	<b>0.000</b>	<b>0.29</b>	<i>KLK3_PSA_exons</i>	<b>9.71E-14</b>	<b>0.08</b>	<i>TMPRSS2_</i>	<b>4.40E-07</b>	<b>-1.19</b>	<i>KLK2</i>	<b>4.60E-09</b>	<b>-0.31</b>
<i>AMACR</i>	<b>0.001</b>	<b>0.27</b>	<i>KLK3_PSA_exons</i>	<b>5.49E-11</b>	<b>0.11</b>	<i>TDRD</i>	<b>9.02E-07</b>	<b>-0.83</b>	<i>KLK3_PSA_exon</i>	<b>4.65E-09</b>	<b>-0.25</b>
<i>ERG_5prime</i>	<b>0.002</b>	<b>0.44</b>	<i>STEAP2</i>	<b>7.09E-11</b>	<b>0.06</b>	<i>ERG_5prime</i>	<b>9.30E-07</b>	<b>-0.28</b>	<i>KLK4</i>	<b>8.32E-09</b>	<b>-0.14</b>
<i>NAALADL2</i>	<b>0.005</b>	<b>0.12</b>	<i>CAMKK2</i>	<b>6.02E-09</b>	<b>0.15</b>	<i>HOXC6</i>	<b>2.37E-05</b>	<b>-0.15</b>	<i>STEAP2</i>	<b>1.13E-08</b>	<b>-0.19</b>
<i>TDRD</i>	<b>0.012</b>	<b>0.86</b>	<i>MMP26</i>	<b>4.07E-08</b>	<b>0.69</b>	<i>AMACR</i>	<b>2.59E-05</b>	<b>-0.14</b>	<i>PPAP2A</i>	<b>4.37E-08</b>	<b>-0.15</b>
<i>PECI</i>	<b>0.012</b>	<b>0.12</b>	<i>KLK4</i>	<b>2.28E-07</b>	<b>0.05</b>	<i>ERG_3prime</i>	<b>8.38E-05</b>	<b>-1.28</b>	<i>FOLH1_PSMA</i>	<b>9.65E-07</b>	<b>-0.21</b>
<i>FOLH1_PSMA</i>	<b>0.016</b>	<b>0.14</b>	<i>GAPDH</i>	<b>2.94E-07</b>	<b>0.03</b>	<i>HOXC4</i>	<b>7.48E-04</b>	<b>-0.70</b>	<i>ARexons4_8</i>	<b>1.55E-06</b>	<b>-0.19</b>
<i>IMPDH2</i>	<b>0.018</b>	<b>0.10</b>	<i>FOLH1_PSMA</i>	<b>3.17E-07</b>	<b>0.09</b>	<i>CAMKK2</i>	<b>2.28E-03</b>	<b>-0.09</b>	<i>STEAP4</i>	<b>2.82E-06</b>	<b>-0.20</b>
<i>DLX1</i>	<b>0.020</b>	<b>0.91</b>	<i>OR52A2_PSGR</i>	<b>2.03E-06</b>	<b>0.13</b>	<i>DLX1</i>	<b>5.33E-03</b>	<b>-1.21</b>	<i>OR52A2_PSGR</i>	<b>2.97E-06</b>	<b>-0.48</b>
			<i>ARexons4_8</i>	<b>3.70E-06</b>	<b>0.07</b>	<i>GAPDH</i>	<b>5.33E-03</b>	<b>-0.03</b>	<i>KLK3_PSA_exon</i>	<b>4.03E-06</b>	<b>-0.24</b>
			<i>KLK2</i>	<b>6.96E-06</b>	<b>0.07</b>	<i>CDKN3</i>	<b>0.007</b>	<b>-0.42</b>	<i>MMP26</i>	<b>3.23E-05</b>	<b>-1.22</b>
			<i>CDC20</i>	<b>1.05E-05</b>	<b>0.50</b>	<i>PECI</i>	<b>0.009</b>	<b>-0.04</b>	<i>PCA3</i>	<b>4.85E-05</b>	<b>-0.57</b>
			<i>TBP</i>	<b>3.98E-05</b>	<b>0.08</b>	<i>MMP26</i>	<b>0.010</b>	<b>-0.46</b>	<i>AR_truncation_exon</i>	<b>5.61E-05</b>	<b>-1.38</b>
			<i>CLU</i>	<b>5.89E-05</b>	<b>0.42</b>	<i>TERT</i>	<b>0.010</b>	<b>-0.17</b>	<i>UPK2</i>	<b>1.62E-05</b>	<b>0.75</b>

CHAPTER 3: NANOSTRING DATA ANALYSIS 1: THE PILOT STUDY

	<i>05</i>						<i>04</i>	
<i>SERPINB5_Maspin</i>	<i>7.52E-05</i>	<i>0.18</i>	<i>HPN</i>	<i>0.042</i>	<i>-0.08</i>	<i>SLC12A1</i>	<i>2.10E-04</i>	<i>0.67</i>
<i>GOLM1</i>	<i>8.48E-05</i>	<i>0.16</i>	<i>BRAF</i>	<i>0.048</i>	<i>-0.03</i>	<i>SPINK1</i>	<i>3.20E-04</i>	<i>0.44</i>
<i>DLX1</i>	<i>1.89E-04</i>	<i>1.24</i>				<i>HPRT</i>	<i>3.25E-04</i>	<i>-0.16</i>
<i>BRAF</i>	<i>1.96E-04</i>	<i>0.06</i>				<i>GOLM1</i>	<i>4.14E-04</i>	<i>-0.65</i>
<i>TERT</i>	<i>2.09E-04</i>	<i>0.25</i>				<i>Timp4</i>	<i>5.49E-04</i>	<i>-1.56</i>
<i>TDRD</i>	<i>2.57E-04</i>	<i>0.61</i>				<i>CLU</i>	<i>6.21E-04</i>	<i>-1.07</i>
<i>TMPRSS2_ERG</i>	<i>4.27E-04</i>	<i>1.11</i>				<i>CDC20</i>	<i>1.20E-03</i>	<i>-0.85</i>
<i>ERG_3prime</i>	<i>6.23E-04</i>	<i>0.83</i>				<i>DLX1</i>	<i>6.39E-03</i>	<i>-1.58</i>
<i>PCA3</i>	<i>6.23E-04</i>	<i>0.16</i>				<i>B2M</i>	<i>0.009</i>	<i>-0.09</i>
<i>AR_truncation_exon</i>	<i>1.71E-03</i>	<i>0.71</i>				<i>IMPDH2</i>	<i>0.010</i>	<i>-0.10</i>
<i>MDK</i>	<i>0.003</i>	<i>0.05</i>				<i>AGR2</i>	<i>0.021</i>	<i>-0.48</i>
<i>MKi67</i>	<i>0.003</i>	<i>0.72</i>				<i>MDK</i>	<i>0.048</i>	<i>-0.21</i>
<i>B2M</i>	<i>0.003</i>	<i>0.08</i>				<i>NAALADL2</i>	<i>0.048</i>	<i>-0.10</i>
<i>PPAP2A</i>	<i>0.006</i>	<i>0.04</i>						
<i>Timp4</i>	<i>0.006</i>	<i>0.21</i>						
<i>UPK2</i>	<i>0.011</i>	<i>-0.23</i>						
<i>OGT</i>	<i>0.011</i>	<i>0.04</i>						
<i>CDKN3</i>	<i>0.012</i>	<i>0.35</i>						
<i>SPINK1</i>	<i>0.012</i>	<i>-0.11</i>						
<i>ERG_5prime</i>	<i>0.016</i>	<i>0.17</i>						
<i>AURKA</i>	<i>0.036</i>	<i>0.15</i>						

### 3.6 Significantly varying genes

Expression distribution of each transcript was fairly even between clinical categories for most probes (Figure 3.19), with only 16 of the 57 probes found to be significantly different between the clinical categories (Kruskal-Wallis rank sum test, adjusted  $p < 0.05$ , Table 3.18).

Mann Whitney U tests (section 2.4.1) were applied to three separate data comparisons; i) cancer vs. non-cancer, ii) aggressive cancer (Advanced and high risk) vs. Non-aggressive cancer (I, L) and iii) the two extremes (Advanced vs. CB), (Table 3.3, page 115). Nine probes were significantly differentially expressed ( $p < 0.05$ , Mann-Whitney U test) between cancer and non-cancer samples after multiple testing correction via the Hochberg method (Table 3.19). All of these transcripts were up-regulated in the cancer (Figure 3.20) and included many well established PCa-associated transcripts such as *ERG*, *TMPRSS2:ERG* and *PCA3*.

**Table 3.18** Kruskal-Wallis identified 16 probes that significantly differ across clinical category.

<i>Probe</i>	<i>p-value</i>	<i>Adjusted p-value</i>	$\chi$
<i>SPINK1</i>	$3.9 \times 10^{-09}$	$2.2 \times 10^{-07}$	47.79
<i>SLC12A1</i>	$2.5 \times 10^{-08}$	$1.4 \times 10^{-06}$	43.78
<i>KLK3 exons 2-3</i>	$1.8 \times 10^{-06}$	$9.9 \times 10^{-05}$	34.60
<i>KLK3 exons 1-2</i>	$2.3 \times 10^{-06}$	$1.3 \times 10^{-04}$	34.04
<i>TMPRSS2:ERG</i>	$1.7 \times 10^{-05}$	$8.8 \times 10^{-04}$	29.69
<i>UPK2</i>	$1.7 \times 10^{-05}$	$8.8 \times 10^{-04}$	29.70
<i>ERG 3'</i>	$2.2 \times 10^{-05}$	0.001	29.09
<i>STEAP2</i>	$2.9 \times 10^{-05}$	0.001	28.52
<i>DLX1</i>	$3.1 \times 10^{-05}$	0.002	28.35
<i>KLK4</i>	$3.6 \times 10^{-05}$	0.002	28.00
<i>HPN</i>	$8.5 \times 10^{-05}$	0.004	26.12
<i>ERG 5'</i>	$1 \times 10^{-04}$	0.005	25.73
<i>PSGR</i>	$1.3 \times 10^{-04}$	0.01	25.08
<i>PCA3</i>	$3.6 \times 10^{-04}$	0.02	22.84
<i>KLK2</i>	$4.1 \times 10^{-04}$	0.02	22.56
<i>CAMKK2</i>	$6.5 \times 10^{-04}$	0.03	21.49

**Table 3.19** Transcripts differentially expressed between cancer (A, H, I, L) and non-cancer samples (Mann Whitney U test).

<i>Transcript</i>	<i>p - value</i>	<i>Adjusted p - value</i>	<i>Log2 Fold Change</i>
<i>DLX1</i>	$3.2 \times 10^{-06}$	<b>0.0002</b>	<b>1.33</b>
<i>ERG 3'</i>	$4.25 \times 10^{-05}$	<b>0.002</b>	<b>1.25</b>
<i>TMPRSS2:ERG</i>	$1.19 \times 10^{-04}$	<b>0.006</b>	<b>0.93</b>
<i>HOXC4</i>	$2.6 \times 10^{-04}$	<b>0.013</b>	<b>0.635</b>
<i>ERG 5'</i>	$1.73 \times 10^{-05}$	<b>0.001</b>	<b>0.281</b>
<i>HOXC6</i>	$4.97 \times 10^{-05}$	<b>0.002</b>	<b>0.242</b>
<i>PCA3</i>	$2.02 \times 10^{-04}$	<b>0.01</b>	<b>0.225</b>
<i>M.genitalium RplB</i>	$4.48 \times 10^{-04}$	<b>0.022</b>	<b>0.144</b>
<i>HPN</i>	$9.02 \times 10^{-06}$	<b>0.0005</b>	<b>0.127</b>

Eleven transcripts were significantly differentially expressed between aggressive and non-aggressive cancers ( $p < 0.05$ , Mann-Whitney U test, Table 3.20). Three of these transcripts were up regulated in the aggressive cancer; *SLC12A1*, *UPK2* and *SPINK1* (Figure 3.21). *SLC12A1* and *UPK2* are tissue specific controls for kidney and bladder, respectively. Advanced tumours often become more solidified and firm which might cause the release of cells and EVs from these prostates to be inhibited. This would cause a relative increase in detection of transcripts from other sources such as the kidney and bladder. Note that *SLC12A1*, *UPK2* and *SPINK1* were heavily correlated across all of the samples (section 3.3.7) and so this result should be taken with some caution. Eight transcripts were down-regulated in the aggressive cancers, again I hypothesise that this is due to a decreased level of cells and EVs emerging from the prostate and it's cancer via DRE, as these transcripts are mostly either prostate or cancer specific.



# CHAPTER 3 – NANOSTRING DATA ANALYSIS 1: The Pilot Study

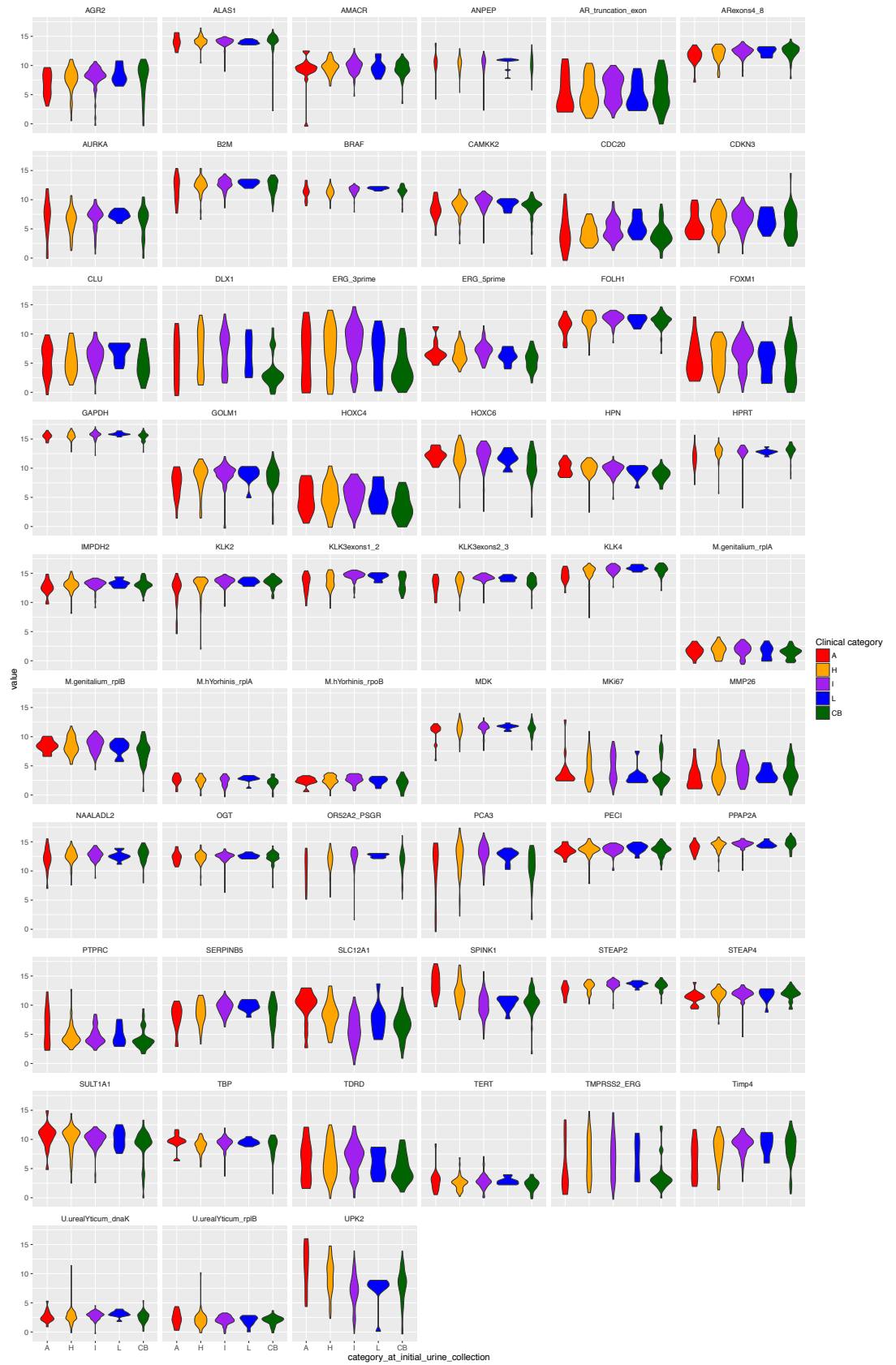
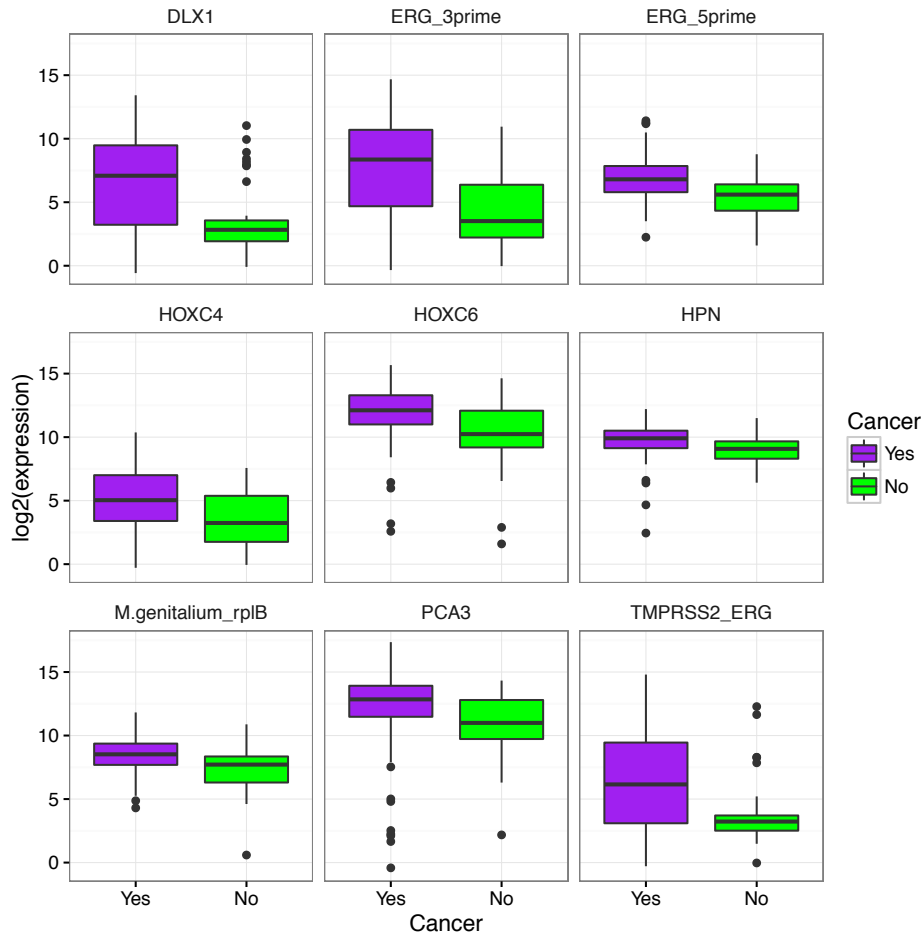


Figure 3.19 Violin plots showing distribution of each probe across clinical category.



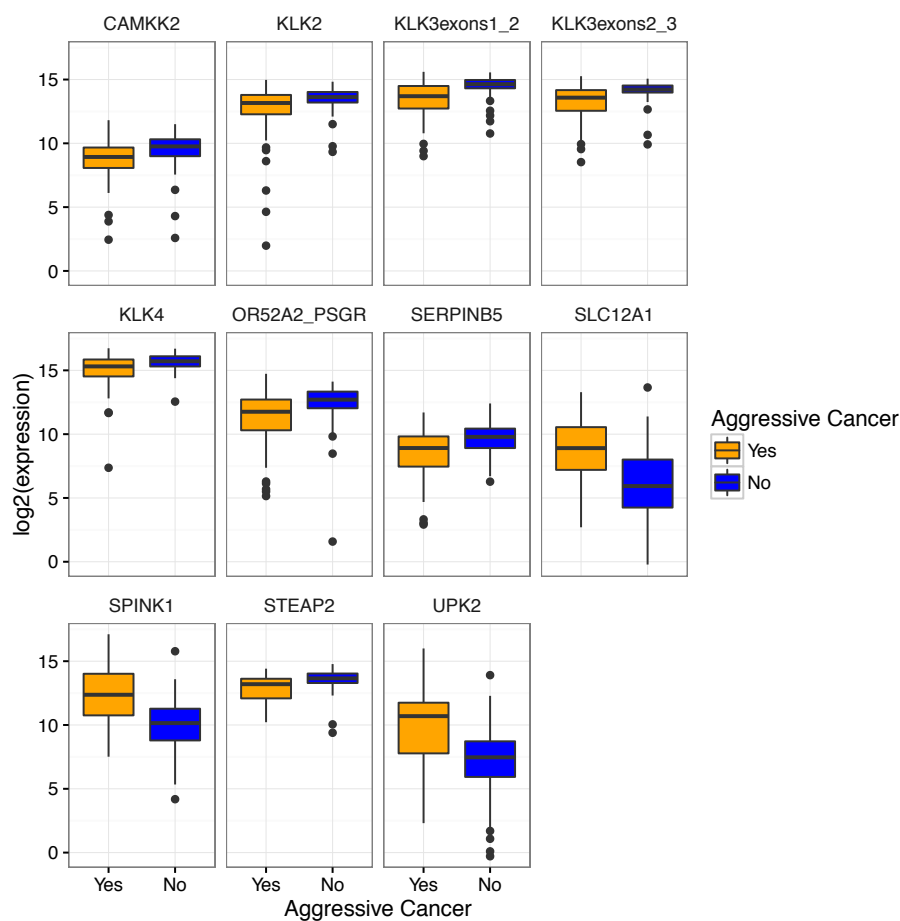
**Figure 3.20** Boxplots showing the expression levels in significantly differentially expressed genes between cancer and non-cancer samples found by Mann Whitney U test.

Six transcripts were significantly differentially expressed between Advanced and CB ( $p < 0.05$ , Mann-Whitney U test, after multiple testing correction, Table 3.21). *SLC12A1* and *SPINK1* are up-regulated as has been previously discussed. The other four transcripts were down-regulated in the advanced samples, and again these include prostate specific transcripts such as *KLK4* and cancer related transcripts such as *PPAP2A* and *STEAP2* (Figure 3.22).

CHAPTER 3 – NANOSTRING DATA ANALYSIS 1: The Pilot Study

**Table 3.20** Transcripts differentially expressed between aggressive cancer and non-aggressive samples (Mann Whitney U test).

<i>Transcript</i>	<i>p - value</i>	<i>Adjusted p - value</i>	<i>Log2 Fold Change</i>
<i>SLC12A1</i>	$2.86 \times 10^{-08}$	$1.6 \times 10^{-06}$	0.59
<i>UPK2</i>	$2.29 \times 10^{-07}$	$1.26 \times 10^{-05}$	0.52
<i>SPINK1</i>	$3.94 \times 10^{-10}$	$2.25 \times 10^{-08}$	0.29
<i>SERPINB5</i>	$2.78 \times 10^{-04}$	$1.36 \times 10^{-02}$	-0.13
<i>CAMKK2</i>	$4.44 \times 10^{-04}$	$2.13 \times 10^{-02}$	-0.13
<i>PSGR</i>	$9.55 \times 10^{-05}$	$4.77 \times 10^{-03}$	-0.11
<i>KLK3 exons 1-2</i>	$6.18 \times 10^{-07}$	$3.34 \times 10^{-05}$	-0.1
<i>KLK3 exons 2-3</i>	$9.15 \times 10^{-07}$	$4.85 \times 10^{-05}$	-0.07
<i>KLK2</i>	$4.77 \times 10^{-04}$	$2.24 \times 10^{-02}$	-0.05
<i>STEAP2</i>	$3.38 \times 10^{-05}$	$1.76 \times 10^{-03}$	-0.05
<i>KLK4</i>	$8.02 \times 10^{-05}$	$4.09 \times 10^{-03}$	-0.04



**Figure 3.21** Boxplots showing differential expression between aggressive cancer and not aggressive PCa samples for those deemed significant by Mann Whitney U test.

Table 3.21 differentially expressed transcripts when comparing advanced samples with benign (no evidence of cancer) samples (Mann Whitney U test).

Transcript	<i>p</i> - value	Adjusted <i>p</i> - value	Log2 Fold Change
<i>SLC12A1</i>	$6.24 \times 10^{-06}$	$3.49 \times 10^{-04}$	0.68
<i>SPINK1</i>	$1.08 \times 10^{-06}$	$6.14 \times 10^{-05}$	0.35
<i>HPRT</i>	$1.29 \times 10^{-04}$	$7.1 \times 10^{-03}$	-0.17
<i>KLK4</i>	$1.53 \times 10^{-04}$	$8.29 \times 10^{-03}$	-0.12
<i>STEAP2</i>	$6.52 \times 10^{-04}$	$3.39 \times 10^{-02}$	-0.09
<i>PPAP2A</i>	$5.6 \times 10^{-04}$	$2.97 \times 10^{-02}$	-0.07

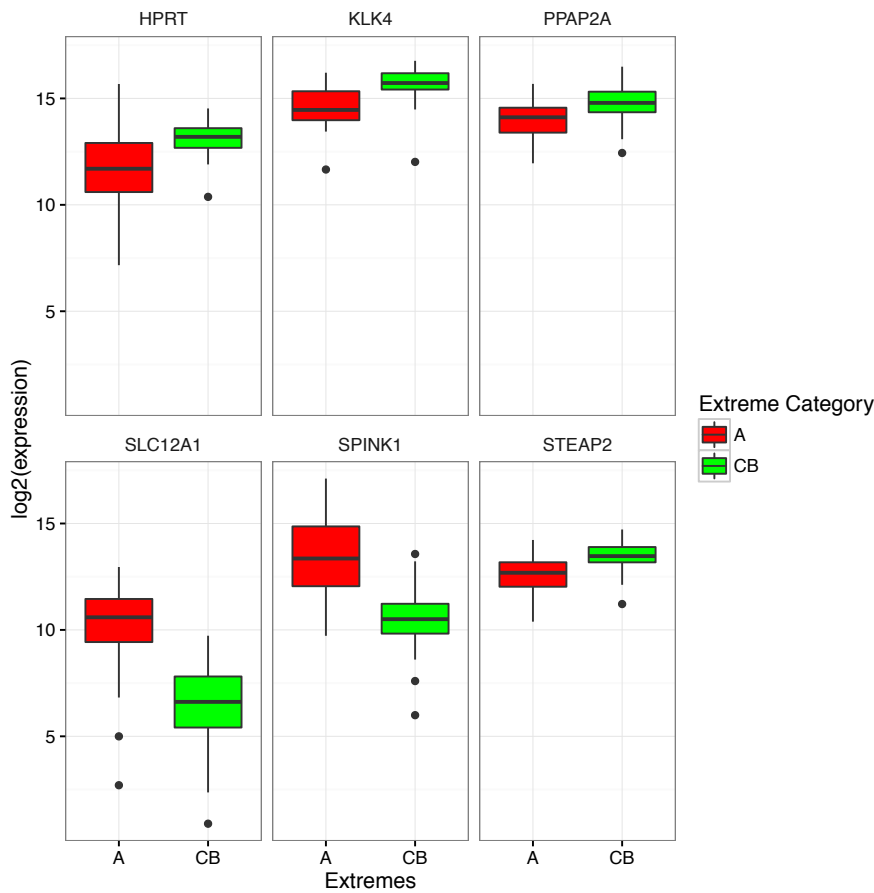


Figure 3.22 Boxplots showing differential expression between advanced cancer and non-cancer samples for those deemed significant via Mann Whitney U testing.

### 3.7 Low-risk, intermediate-risk and high-risk trend

Five probes showed significant increasing or decreasing expression trend with increasing D’Amico risk category (Spearman’s correlation,  $p < 0.05$  after multiple testing correction, Table 3.22). Three of these probes were identified to be highly correlated in general (*SPINK1*, *UPK2* and *SLC12A1*: section 3.3.7). The other two probes are from the same transcript (*KLK3*). The decrease

in *KLK3* with increasing cancer risk has been reported in previous prostate tissue studies (section 1.4.1) and urine.

**Table 3.22 Spearman's correlation results comparing expression with ordered clinical categories: Low-, Intermediate- and High-risk.**

<i>Transcript</i>	<i>p – value</i>	<i>Adjusted p - value</i>	<i>R</i>
<i>SPINK1</i>	$1.27 \times 10^{-06}$	$7.24 \times 10^{-05}$	<b>0.41</b>
<i>UPK2</i>	$2.74 \times 10^{-05}$	$1.53 \times 10^{-03}$	<b>0.36</b>
<i>SLC12A1</i>	$1.08 \times 10^{-04}$	$5.96 \times 10^{-03}$	<b>0.33</b>
<i>KLK3 exons 2-3</i>	$1.12 \times 10^{-04}$	$6.04 \times 10^{-03}$	<b>-0.33</b>
<i>KLK3 exons 1-2</i>	$1.14 \times 10^{-04}$	$6.04 \times 10^{-03}$	<b>-0.33</b>

### **3.8 Clinical Prediction models**

To test the ability of NanoString data derived from urine EVs to predict the presence of cancer and/or it's aggressiveness various models were produced to distinguish between a) PCa and benign samples, b) aggressive PCa and non-aggressive PCa and c) advanced and benign samples (Table 3.3). All samples were used in the training set due to the pilot nature of this work. The modelling techniques applied here are logistic regression models using step wise variable selection (section 2.6.4), Lasso logistic regression models for shrinkage and variable selection (section 2.6.2), and random forest (section 2.6.3).

#### **3.8.1 Logistic regression models using step wise variable selection**

The optimal output cancer vs. non-cancer model contained 33 transcripts and had an AIC score of 68 (Table 3.23), the optimal aggressive cancer vs. non-aggressive cancer model contained 37 transcripts and had an AIC score of 76 (

Table 3.25), and the optimal model for distinguishing Advanced cancer from CB contained 9 transcripts and had an AIC score of 18 (Table 3.27). In each model the sample category was predicted with 100% sensitivity, 100% specificity, and 100% PPV (Table 3.24, Table 3.26, Table 3.28). This may mean the models are over-fitting the data and caution should be taken.

Table 3.23 Transcripts in the Step derived model for comparing cancer to non-cancer.

<i>Transcript</i>	<i>p - value</i>	<i>Coefficient</i>
<i>CDKN3</i>	<b>0.97</b>	<b>374.2</b>
<i>FOLH1</i>	<b>0.97</b>	<b>-838.3</b>
<i>FOXMI</i>	<b>0.97</b>	<b>138.8</b>
<i>HPN</i>	<b>0.97</b>	<b>743.9</b>
<i>IMPDH2</i>	<b>0.97</b>	<b>691.2</b>
<i>KLK3 exons 2-3</i>	<b>0.97</b>	<b>1844</b>
<i>M.genitalium RplB</i>	<b>0.97</b>	<b>-491.5</b>
<i>NAALADL2</i>	<b>0.97</b>	<b>-549.9</b>
<i>AURKA</i>	<b>0.98</b>	<b>-236.7</b>
<i>BRAF</i>	<b>0.98</b>	<b>562.1</b>
<i>KLK2</i>	<b>0.98</b>	<b>-1534.6</b>
<i>M.hyorhinis RplA</i>	<b>0.98</b>	<b>-1415.5</b>
<i>PSGR</i>	<b>0.98</b>	<b>364.4</b>
<i>SULT1A1</i>	<b>0.98</b>	<b>686</b>
<i>TMPRSS2:ERG</i>	<b>0.98</b>	<b>283.5</b>
<i>ANPEP</i>	<b>0.99</b>	<b>876.6</b>
<i>AR truncation exon</i>	<b>0.99</b>	<b>-232.3</b>
<i>AR exons 4-8</i>	<b>0.99</b>	<b>-462.4</b>
<i>B2M</i>	<b>0.99</b>	<b>1219.1</b>
<i>CAMKK2</i>	<b>0.99</b>	<b>-432.1</b>
<i>DLX1</i>	<b>0.99</b>	<b>220.4</b>
<i>ERG 3'</i>	<b>0.99</b>	<b>240.6</b>
<i>ERG 5'</i>	<b>0.99</b>	<b>985.7</b>
<i>KLK4</i>	<b>0.99</b>	<b>-1187.3</b>
<i>MDK</i>	<b>0.99</b>	<b>-1265.6</b>
<i>MMP26</i>	<b>0.99</b>	<b>-730</b>
<i>OGT</i>	<b>0.99</b>	<b>-1185.7</b>
<i>PCA3</i>	<b>0.99</b>	<b>399.8</b>
<i>SERPINB5</i>	<b>0.99</b>	<b>-234.1</b>
<i>TBP</i>	<b>0.99</b>	<b>497</b>
<i>U.urealyticum dnaK</i>	<b>0.99</b>	<b>-919.6</b>
<i>U.urealyticum RplB</i>	<b>0.99</b>	<b>837.9</b>
<i>UPK2</i>	<b>0.99</b>	<b>-132.5</b>

Table 3.24 Category predictions using the cancer vs. non-cancer step model.

<i>Test</i>	<i>Actual Category</i>	
	<i>Disease Present</i>	<i>No evidence of disease</i>
<i>Positive</i>	<b>148</b>	<b>0</b>
<i>Negative</i>	<b>0</b>	<b>40</b>

**Table 3.25** Transcripts in the Step derived model for comparing aggressive cancers (A, H) to non-aggressive cancers (I, L).

<i>Transcript</i>	<i>p - value</i>	<i>Coefficient</i>
<i>AGR2</i>	<b>0.98</b>	<b>-233.03</b>
<i>AMACR</i>	<b>0.98</b>	<b>-1066.97</b>
<i>AURKA</i>	<b>0.98</b>	<b>-395.93</b>
<i>BRAF</i>	<b>0.98</b>	<b>-1106.78</b>
<i>ERG 5'</i>	<b>0.98</b>	<b>164.37</b>
<i>FOXMI</i>	<b>0.98</b>	<b>149.23</b>
<i>HPRT</i>	<b>0.98</b>	<b>317.38</b>
<i>IMPDH2</i>	<b>0.98</b>	<b>913.69</b>
<i>KLK3 exons 1-2</i>	<b>0.98</b>	<b>-607.92</b>
<i>M.genitalium RplA</i>	<b>0.98</b>	<b>-369.55</b>
<i>M.hyorinis RplA</i>	<b>0.98</b>	<b>1236.48</b>
<i>MKi67</i>	<b>0.98</b>	<b>-200.18</b>
<i>NAALADL2</i>	<b>0.98</b>	<b>204.06</b>
<i>PSGR</i>	<b>0.98</b>	<b>-826.31</b>
<i>PCA3</i>	<b>0.98</b>	<b>190.98</b>
<i>PPAP2A</i>	<b>0.98</b>	<b>546.96</b>
<i>SERPINB5</i>	<b>0.98</b>	<b>-357.41</b>
<i>SPINK1</i>	<b>0.98</b>	<b>708.99</b>
<i>STEAP4</i>	<b>0.98</b>	<b>585.87</b>
<i>SULT1A1</i>	<b>0.98</b>	<b>144.98</b>
<i>TMPRSS2:ERG</i>	<b>0.98</b>	<b>-109.31</b>
<i>Timp4</i>	<b>0.98</b>	<b>-304.36</b>
<i>U.urealyticum dnaK</i>	<b>0.98</b>	<b>390.71</b>
<i>U.urealyticum RplB</i>	<b>0.98</b>	<b>-456.65</b>
<i>AR exons 4-8</i>	<b>0.99</b>	<b>148.66</b>
<i>CDC20</i>	<b>0.99</b>	<b>-180.92</b>
<i>DLX1</i>	<b>0.99</b>	<b>47.8</b>
<i>FOLH1</i>	<b>0.99</b>	<b>492.55</b>
<i>GAPDH</i>	<b>0.99</b>	<b>-366.81</b>
<i>GOLM1</i>	<b>0.99</b>	<b>499.32</b>
<i>M.hyorinis rpoB</i>	<b>0.99</b>	<b>218.58</b>
<i>PTPRC</i>	<b>0.99</b>	<b>-126.41</b>
<i>TBP</i>	<b>0.99</b>	<b>-224.34</b>
<i>TDRD</i>	<b>0.99</b>	<b>-160.03</b>
<i>TERT</i>	<b>0.99</b>	<b>156.96</b>
<i>UPK2</i>	<b>0.99</b>	<b>74.34</b>
<i>AGR2</i>	<b>0.98</b>	<b>-233.03</b>

**Table 3.26 Category Predictions when using the aggressive cancer model derived from Step.**

	<b>Aggressive Disease Present</b>	<b>Aggressive undetectable</b>
Positive	<b>68</b>	<b>0</b>
Negative	<b>0</b>	<b>80</b>

**Table 3.27 Transcripts in the extreme model (A Vs. CB) derived from Step.**

<i>Transcript</i>	<i>p - value</i>	<i>Coefficient</i>
<i>ALAS1</i>	<b>0.995</b>	<b>-600.84</b>
<i>KLK4</i>	<b>0.995</b>	<b>-465.6</b>
<i>KLK3 exons 2-3</i>	<b>0.995</b>	<b>330.63</b>
<i>BRAF</i>	<b>0.995</b>	<b>422.54</b>
<i>M.genitalium RplB</i>	<b>0.995</b>	<b>302.06</b>
<i>HPN</i>	<b>0.995</b>	<b>236.54</b>
<i>Timp4</i>	<b>0.995</b>	<b>-158.99</b>
<i>AR truncation exon</i>	<b>0.995</b>	<b>31.34</b>
<i>ALAS1</i>	<b>0.995</b>	<b>-600.84</b>

**Table 3.28 Category predictions using the extreme model derived from Step.**

<i>Test</i>	<i>Actual Category</i>	
	<i>Advanced Cancer Present</i>	<i>CB</i>
<i>Positive</i>	<b>17</b>	<b>0</b>
<i>Negative</i>	<b>0</b>	<b>40</b>

### 3.8.2 Lasso logistic regression models

The cancer vs. non-cancer model had 16 transcripts (Table 3.29), an AUC of 0.937 and 99.32% sensitivity, 52.5% specificity and 88.55% PPV (Table 3.30). The aggressive cancer (A) vs. non-aggressive cancer model had four transcripts (Table 3.29), an AUC of 0.852, and 61.76% sensitivity, 86.25% specificity and 79.25% PPV (Table 3.31). The extreme Lasso model (A Vs. CB samples) had 12 transcripts (Table 3.29), an AUC of 0.983 and 82.35% sensitivity, 100% specificity and 100% PPV (Table 3.32).



CHAPTER 3 – NANOSTRING DATA ANALYSIS 1: The Pilot Study

**Table 3.29** Lasso coefficients for three models A) cancer Vs. Non-cancer B) Aggressive cancer Vs. Non-aggressive cancer C) extreme model (A Vs. CB)

<i>Cancer Model</i>		<i>Aggressive Model</i>		<i>Extreme Model</i>	
<i>Gene name</i>	<i>Coefficient</i>	<i>Gene name</i>	<i>Coefficient</i>	<i>Gene name</i>	<i>Coefficient</i>
<i>ALAS1</i>	<i>-0.042</i>	<i>KLK3 exons 1-2</i>	<i>-0.268</i>	<i>AURKA</i>	<i>0.035</i>
<i>AR exons4-8</i>	<i>-0.278</i>	<i>SERPINB5</i>	<i>-0.087</i>	<i>DLX1</i>	<i>0.043</i>
<i>CLU</i>	<i>0.037</i>	<i>SLC12A1</i>	<i>0.069</i>	<i>ERG 5'</i>	<i>0.235</i>
<i>DLX1</i>	<i>0.164</i>	<i>SPINK1</i>	<i>0.278</i>	<i>HOXC4</i>	<i>0.003</i>
<i>ERG 3'</i>	<i>0.082</i>			<i>HPN</i>	<i>0.188</i>
<i>ERG 5'</i>	<i>0.169</i>			<i>HPRT</i>	<i>-0.554</i>
<i>HOXC4</i>	<i>0.111</i>			<i>PPAP2A</i>	<i>-0.223</i>
<i>HPN</i>	<i>0.283</i>			<i>SLC12A1</i>	<i>0.104</i>
<i>IMPDH2</i>	<i>-0.131</i>			<i>SPINK1</i>	<i>0.321</i>
<i>KLK2</i>	<i>-0.046</i>			<i>STEAP2</i>	<i>-0.297</i>
<i>KLK3 exons 1-2</i>	<i>0.05</i>			<i>STEAP4</i>	<i>-0.113</i>
<i>M.hyorhinis rpoB</i>	<i>0.399</i>			<i>SULT1A1</i>	<i>0.118</i>
<i>MMP26</i>	<i>-0.085</i>				
<i>NAALADL2</i>	<i>-0.003</i>				
<i>PCA3</i>	<i>0.101</i>				
<i>PPAP2A</i>	<i>-0.62</i>				
<i>SULT1A1</i>	<i>0.057</i>				
<i>TMPRSS2:ERG</i>	<i>0.037</i>				
<i>Timp4</i>	<i>-0.033</i>				

**Table 3.30** Category predictions using the Lasso cancer Vs. non-cancer model

<i>Test</i>	<i>Actual Category</i>	
	<i>Disease Present</i>	<i>Disease Absent</i>
<i>Positive</i>	<i>147</i>	<i>19</i>
<i>Negative</i>	<i>1</i>	<i>21</i>

**Table 3.31** Category predictions using the Lasso aggressive cancer Vs. non-aggressive cancer model

<i>Test</i>	<i>Actual Category</i>	
	<i>Aggressive Disease Present</i>	<i>Aggressive Disease Absent</i>
<i>Positive</i>	<i>42</i>	<i>11</i>
<i>Negative</i>	<i>26</i>	<i>69</i>

Table 3.32 Category prediction for the Lasso extreme model (A Vs. CB)

<i>Test</i>	<i>Actual Category</i>	
	<i>Advanced Cancer Present</i>	<i>CB</i>
<i>Positive</i>	<i>14</i>	<i>0</i>
<i>Negative</i>	<i>3</i>	<i>40</i>

### 3.8.3 Random Forest

The cancer vs. non-cancer model had an OOB error estimate of 18.52%, with 87.25% sensitivity, 60% specificity and 89.04% PPV (Table 3.33, the ranked transcript importance provided in Table 3.34). The aggressive vs. non-aggressive cancer model had an OOB error estimate of 22.82%, with 71.64% sensitivity, 81.71% specificity and 76.19% PPV (Table 3.35, Table 3.36). The extremes model had an OOB error estimate of 15.79%, with 70.59% sensitivity, 90% specificity and 75% PPV (

Table 3.37, Table 3.38).

Table 3.33 Confusion matrix for random forest modelling samples on cancer vs. non-cancer. OOB error estimate of 18.52%.

	<i>Cancer not predicted</i>	<i>Cancer predicted</i>	<i>Class error</i>	<i>Sum</i>
<i>CB</i>	<i>24</i>	<i>16</i>	<i>0.4</i>	<i>40</i>
<i>Cancer</i>	<i>19</i>	<i>130</i>	<i>0.13</i>	<i>149</i>

**Table 3.34** Gini values for the random forest model to categorise the samples into cancer and non-cancer.

<i>Transcript</i>	<i>Gini</i>	<i>Rank</i>	<i>Transcript</i>	<i>Gini</i>	<i>Rank</i>	<i>Transcript</i>	<i>Gini</i>	<i>Rank</i>
<i>DLX1</i>	2.2	1	<i>CLU</i>	0.61	20	<i>UPK2</i>	0.4	3
	7						6	9
<i>ERG 3'</i>	2.0	2	<i>SPINK1</i>	0.61	21	<i>GAPDH</i>	0.4	4
	8						5	0
<i>TMPRSS2:ERG</i>	2.0	3	<i>B2M</i>	0.60	22	<i>Timp4</i>	0.4	4
	6						4	1
<i>HOXC6</i>	1.8	4	<i>FOXM1</i>	0.58	23	<i>KLK3 exons 2-3</i>	0.4	4
	6						4	2
<i>HPN</i>	1.8	5	<i>AR exons 4-8</i>	0.58	24	<i>SERPINB5</i>	0.4	4
	3						4	3
<i>PCA3</i>	1.2	6	<i>CAMKK2</i>	0.56	25	<i>CDC20</i>	0.4	4
	1						2	4
<i>PPAP2A</i>	1.1	7	<i>SULT1A1</i>	0.56	26	<i>PECI</i>	0.4	4
	6						2	5
<i>ERG 5'</i>	1.1	8	<i>M.genitalium</i>	0.54	27	<i>AR truncation</i>	0.4	4
	4		<i>RplA</i>			<i>exon</i>	1	6
<i>HOXC4</i>	1.1	9	<i>STEAP2</i>	0.54	28	<i>U.urealyticum</i>	0.4	4
	2					<i>RplB</i>	1	7
<i>M.hYorhini rpoB</i>	1.0	10	<i>MKi67</i>	0.54	29	<i>GOLM1</i>	0.4	4
	4						1	8
<i>ALAS1</i>	0.8	11	<i>MMP26</i>	0.53	30	<i>AURKA</i>	0.4	4
	5						0	9
<i>PTPRC</i>	0.7	12	<i>KLK2</i>	0.52	31	<i>ANPEP</i>	0.4	5
	7						0	0
<i>M.genitalium</i>	0.7	13	<i>AGR2</i>	0.50	32	<i>MDK</i>	0.3	5
<i>RplB</i>	7						9	1
<i>SLC12A1</i>	0.7	14	<i>AMACR</i>	0.49	33	<i>U.urealyticum</i>	0.3	5
	5					<i>dnaK</i>	8	2
<i>HPRT</i>	0.7	15	<i>FOLH1</i>	0.49	34	<i>CDKN3</i>	0.3	5
	5						7	3
<i>KLK3 exons 1-2</i>	0.7	16	<i>TBP</i>	0.48	35	<i>BRAF</i>	0.3	5
	1						6	4
<i>NAALADL2</i>	0.6	17	<i>TERT</i>	0.48	36	<i>KLK4</i>	0.3	5
	7						6	5
<i>TDRD</i>	0.6	18	<i>IMPDH2</i>	0.47	37	<i>PSGR</i>	0.3	5
	3						3	6
<i>STEAP4</i>	0.6	19	<i>M.hYorhinis</i>	0.46	38	<i>OGT</i>	0.3	5
	1		<i>RplA</i>				0	7

**Table 3.35** Confusion matrix for random forest modelling samples on aggressive cancer vs. non-aggressive cancer. OOB error estimation of 22.82%.

	<i>Aggressive Cancer not predicted</i>	<i>Aggressive Cancer predicted</i>	<i>Class error</i>	<i>Sum</i>
<i>Non-aggressive cancer</i>	67	15	0.18	82

<i>Aggressive Cancer</i>	<i>19</i>	<i>48</i>	<i>0.28</i>	<i>67</i>
--------------------------	-----------	-----------	-------------	-----------

Table 3.36 Gini values for the random forest model to categorise the samples into aggressive cancer and non-aggressive cancer.

<i>Transcript</i>	<i>Gini score</i>	<i>Rank</i>	<i>Transcript</i>	<i>Gini score</i>	<i>Rank</i>	<i>Transcript</i>	<i>Gini score</i>	<i>Rank</i>
<i>SPINK1</i>	<i>5.2</i>	<i>1</i>	<i>ERG 3'</i>	<i>1.00</i>	<i>20</i>	<i>HOXC6</i>	<i>0.7</i>	<i>3</i>
<i>KLK3 exons 1-2</i>	<i>4.0</i>	<i>2</i>	<i>MDK</i>	<i>0.99</i>	<i>21</i>	<i>TERT</i>	<i>0.7</i>	<i>4</i>
<i>KLK3 exons 2-3</i>	<i>3.6</i>	<i>3</i>	<i>CAMKK2</i>	<i>0.97</i>	<i>22</i>	<i>NAALADL2</i>	<i>0.7</i>	<i>4</i>
<i>UPK2</i>	<i>3.5</i>	<i>4</i>	<i>CLU</i>	<i>0.95</i>	<i>23</i>	<i>ERG 5'</i>	<i>0.6</i>	<i>4</i>
<i>SLC12A1</i>	<i>3.4</i>	<i>5</i>	<i>GAPDH</i>	<i>0.93</i>	<i>24</i>	<i>AR exons 4-8</i>	<i>0.6</i>	<i>4</i>
<i>SERPINB5</i>	<i>1.8</i>	<i>6</i>	<i>AGR2</i>	<i>0.87</i>	<i>25</i>	<i>M.genitalium RplA</i>	<i>0.6</i>	<i>4</i>
<i>SULT1A1</i>	<i>1.8</i>	<i>7</i>	<i>TMPRSS2:ER G</i>	<i>0.86</i>	<i>26</i>	<i>PCA3</i>	<i>0.6</i>	<i>4</i>
<i>KLK4</i>	<i>1.8</i>	<i>8</i>	<i>CDKN3</i>	<i>0.86</i>	<i>27</i>	<i>DLX1</i>	<i>0.6</i>	<i>4</i>
<i>BRAF</i>	<i>1.5</i>	<i>9</i>	<i>B2M</i>	<i>0.85</i>	<i>28</i>	<i>TBP</i>	<i>0.6</i>	<i>4</i>
<i>PSGR</i>	<i>1.4</i>	<i>10</i>	<i>TDRD</i>	<i>0.83</i>	<i>29</i>	<i>PECI</i>	<i>0.6</i>	<i>4</i>
<i>HPN</i>	<i>1.3</i>	<i>11</i>	<i>ALAS1</i>	<i>0.82</i>	<i>30</i>	<i>MMP26</i>	<i>0.5</i>	<i>4</i>
<i>U.urealyticum dnaK</i>	<i>1.2</i>	<i>12</i>	<i>AR truncation exon</i>	<i>0.82</i>	<i>31</i>	<i>M.hYorhinish rpoB</i>	<i>0.5</i>	<i>5</i>
<i>Timp4</i>	<i>1.2</i>	<i>13</i>	<i>IMPDH2</i>	<i>0.81</i>	<i>32</i>	<i>OGT</i>	<i>0.5</i>	<i>5</i>
<i>STEAP2</i>	<i>1.2</i>	<i>14</i>	<i>GOLM1</i>	<i>0.79</i>	<i>33</i>	<i>HOXC4</i>	<i>0.5</i>	<i>5</i>
<i>U.urealyticum RplB</i>	<i>1.0</i>	<i>15</i>	<i>FOLH1</i>	<i>0.79</i>	<i>34</i>	<i>AMACR</i>	<i>0.5</i>	<i>5</i>
<i>M.hYorhinish RplA</i>	<i>1.0</i>	<i>16</i>	<i>FOXMI</i>	<i>0.79</i>	<i>35</i>	<i>ANPEP</i>	<i>0.5</i>	<i>5</i>
<i>KLK2</i>	<i>1.0</i>	<i>17</i>	<i>MKi67</i>	<i>0.77</i>	<i>36</i>	<i>HPRT</i>	<i>0.5</i>	<i>5</i>
<i>PPAP2A</i>	<i>1.0</i>	<i>18</i>	<i>M.genitalium RplB</i>	<i>0.73</i>	<i>37</i>	<i>PTPRC</i>	<i>0.5</i>	<i>5</i>
<i>CDC20</i>	<i>1.0</i>	<i>19</i>	<i>AURKA</i>	<i>0.73</i>	<i>38</i>	<i>STEAP4</i>	<i>0.5</i>	<i>5</i>

**Table 3.37 Confusion matrix for random forest modelling the samples belonging to the extreme clinical categories (A vs. CB). OOB error estimate of 15.79%.**

	<i>CB predicted</i>	<i>Advanced predicted</i>	<i>Class error</i>	<i>Sum</i>
<i>CB</i>	<b>36</b>	<b>4</b>	<b>0.1</b>	<b>40</b>
<i>Advanced</i>	<b>5</b>	<b>12</b>	<b>0.29</b>	<b>17</b>

Table 3.38 Gini values for the random forest model to categorise the extreme samples (A vs. CB).

<i>Transcript</i>	<i>Gini score</i>	<i>Rank</i>	<i>Transcript</i>	<i>Gini score</i>	<i>Rank</i>	<i>Transcript</i>	<i>Gini score</i>	<i>Rank</i>
<i>SPINK1</i>	5.2 1	1	<i>ERG 3'</i>	1.00	20	<i>HOXC6</i>	0.7 2	3 9
<i>KLK3 exons 1-2</i>	4.0 7	2	<i>MDK</i>	0.99	21	<i>TERT</i>	0.7 2	4 0
<i>KLK3 exons 2-3</i>	3.6 3	3	<i>CAMKK2</i>	0.97	22	<i>NAALADL2</i>	0.7 1	4 1
<i>UPK2</i>	3.5 5	4	<i>CLU</i>	0.95	23	<i>ERG 5'</i>	0.6 9	4 2
<i>SLC12A1</i>	3.4 8	5	<i>GAPDH</i>	0.93	24	<i>AR exons 4-8</i>	0.6 8	4 3
<i>SERPINB5</i>	1.8 9	6	<i>AGR2</i>	0.87	25	<i>M.genitalium RplA</i>	0.6 6	4 4
<i>SULT1A1</i>	1.8 8	7	<i>TMPRSS2:ERG</i>	0.86	26	<i>PCA3</i>	0.6 3	4 5
<i>KLK4</i>	1.8 6	8	<i>CDKN3</i>	0.86	27	<i>DLX1</i>	0.6 2	4 6
<i>BRAF</i>	1.5 7	9	<i>B2M</i>	0.85	28	<i>TBP</i>	0.6 2	4 7
<i>PSGR</i>	1.4 9	10	<i>TDRD</i>	0.83	29	<i>PECI</i>	0.6 2	4 8
<i>HPN</i>	1.3 9	11	<i>ALAS1</i>	0.82	30	<i>MMP26</i>	0.5 9	4 9
<i>U.urealyticum dnaK</i>	1.2 8	12	<i>AR truncation exon</i>	0.82	31	<i>M.hYorhinish rpoB</i>	0.5 9	5 0
<i>Timp4</i>	1.2 5	13	<i>IMPDH2</i>	0.81	32	<i>OGT</i>	0.5 9	5 1
<i>STEAP2</i>	1.2 4	14	<i>GOLM1</i>	0.79	33	<i>HOXC4</i>	0.5 8	5 2
<i>U.urealyticum RplB</i>	1.0 9	15	<i>FOLH1</i>	0.79	34	<i>AMACR</i>	0.5 7	5 3
<i>M.hYorhinish RplA</i>	1.0 8	16	<i>FOXM1</i>	0.79	35	<i>ANPEP</i>	0.5 6	5 4
<i>KLK2</i>	1.0 8	17	<i>MKi67</i>	0.77	36	<i>HPRT</i>	0.5 4	5 5
<i>PPAP2A</i>	1.0 7	18	<i>M.genitalium RplB</i>	0.73	37	<i>PTPRC</i>	0.5 2	5 6
<i>CDC20</i>	1.0 1	19	<i>AURKA</i>	0.73	38	<i>STEAP4</i>	0.5 0	5 7

### 3.8.4 Random Forest applied to all clinical categories

A random forest model was also constructed to classify the samples into their five main types of category (advanced, high-, intermediate-, low-risk and clinically benign), but the results were poor (OOB error estimate of 45.5%, Table 3.39, Table 3.40). The OOB was only modestly improved when the data was adjusted to have equal numbers of samples per category in each tree (53.44%, 170

## CHAPTER 3 – NANOSTRING DATA ANALYSIS 1: The Pilot Study

Table 3.41). This poor performance could be down to the low number of samples per category using the current methods.

CHAPTER 3 – NANOSTRING DATA ANALYSIS 1: The Pilot Study

Table 3.39 Confusion matrix for random forest on all 5 clinical categories. OOB error estimate of 45.5%.

	<i>A</i>	<i>H</i>	<i>I</i>	<i>L</i>	<i>CB</i>	<i>Class error</i>	<i>Sum</i>
<i>A</i>	1	12	4	0	0	0.94	17
<i>H</i>	1	20	21	0	8	0.6	50
<i>I</i>	1	8	58	0	5	0.19	72
<i>L</i>	0	0	6	0	4	1	10
<i>CB</i>	1	6	9	0	24	0.4	40

Table 3.40 Sensitivity, Specificity and PPV for each category after categorising samples into five clinical categories using random forest.

<i>Advanced (A)</i>			
	<i>True A</i>	<i>True Not A</i>	<i>Sensitivity: 25%</i>
<i>Outcome A</i>	1	16	<i>Specificity: 91.35%</i>
<i>Outcome Not A</i>	3	169	<i>PPV: 5.88%</i>

<i>High Risk (H)</i>			
	<i>True H</i>	<i>True Not H</i>	<i>Sensitivity: 40%</i>
<i>Outcome H</i>	20	26	<i>Specificity: 81.29%</i>
<i>Outcome Not H</i>	30	113	<i>PPV: 43.48%</i>

<i>Intermediate Risk (I)</i>			
	<i>True I</i>	<i>True Not I</i>	<i>Sensitivity: 80.56%</i>
<i>Outcome I</i>	58	40	<i>Specificity: 65.81%</i>
<i>Outcome Not I</i>	14	77	<i>PPV: 59.18%</i>

<i>Low Risk (L)</i>			
	<i>True L</i>	<i>True Not L</i>	<i>Sensitivity: 0%</i>
<i>Outcome L</i>	0	0	<i>Specificity: 100%</i>
<i>Outcome Not L</i>	10	179	<i>PPV: *%</i>

<i>Clinically Benign (CB)</i>			
	<i>True CBN</i>	<i>True Not CBN</i>	<i>Sensitivity: 60%</i>
<i>Outcome CBN</i>	24	17	<i>Specificity: 88.59%</i>
<i>Outcome Not CBN</i>	16	132	<i>PPV: 58.54%</i>

Table 3.41 Confusion matrix for random forest on all 5 categories with random sampling to equalise categorical sample sizes. OOB error estimate of 53.44%.

	<i>A</i>	<i>H</i>	<i>I</i>	<i>L</i>	<i>CB</i>	<i>Class error</i>	<i>Sum</i>
<i>A</i>	10	4	2	0	1	0.41	17
<i>H</i>	14	7	15	5	9	0.86	50
<i>I</i>	3	8	46	4	11	0.36	72
<i>L</i>	1	0	4	1	4	0.9	10



<i>CB</i>	<i>2</i>	<i>5</i>	<i>3</i>	<i>6</i>	<i>24</i>	<i>0.4</i>	<i>40</i>
-----------	----------	----------	----------	----------	-----------	------------	-----------

### 3.8.5 Comparing the Models

The OOB error was found to be lowest for modelling the extremes (CB v Aggressive PCa), this was expected as they are the samples that should be least alike in their expression and so should be the easiest categories to separate. The model does give good sensitivity and specificity; however, this error is still fairly high at 15.79%, meaning there could still be improvements. Similarly, the Lasso model (high AUC) and Step model for the extremes comparison both had high sensitivity and specificity, though the step models are likely to be over fitted. From the top fifteen most important transcripts via random forest, five were in common (four uniquely) with the Lasso selected transcripts and five were in common (four uniquely) with the Step selected transcripts. The only transcript common to all three models was HPN, which interestingly only appeared to have mid level importance in each model.

The OOB error for comparing cancer vs. non-cancer was also fairly high, even though it was the second lowest (18.52%). This model showed high sensitivity but was not so specific to identifying cancer in the samples. The Lasso model had a good AUC (0.937) and also showed high sensitivity but not so good specificity to detecting cancer, unlike the Step model, which showed high sensitivity and specificity. However, Step is the least robust of the methods for modelling data. From the top fifteen most important transcripts via random forest, ten were common in the Lasso selected transcripts and seven were common with the Step selected transcripts. There were six transcripts common to all three models: *DLX1*, *ERG 3'*, *TMPRSS2:ERG*, *HPN*, *PCA3* and *ERG 5'*, all of which are transcripts known to be involved or associated to PCa.

The OOB error for comparing aggressive cancers to non-aggressive cancers was higher at 22.82%, though this model had good sensitivity and specificity ratios for the random forest model, 72% and 82%, respectively. From the top fifteen most important transcripts selected via random forest, all four Lasso identified transcripts were in common and ten Step selected transcripts were common. Three of the four Lasso identified transcripts were common to all models: *SPINK1*, *KLK3* exons 1-2 and *SERPINB5*.

This highlights that there is structure in the data that could likely be further improved with data from more samples and more probes.

### 3.9 Transcripts that show high-importance

Seven transcripts were identified by three different methods (Table 3.19, Table 3.29, Table 3.34) to be differentially expressed between cancer and non-cancer samples: *DLX1*, *ERG 3'*, *TMPRSS2:ERG*, *HOXC4*, *ERG 5'*, *PCA3* and *HPN* (Table 3.40). These transcripts all have published associations with PCa. Interestingly, *ERG 5'*, *HOXC6* and *M.genitalium RplB* were significant in the Mann Whitney U test and had been ranked highly by random forest, but were not present in the Lasso model. This is likely due to the inter-correlation of their NanoString signals, as Lasso penalises correlating variables and keeps those it deems to hold the most information (section 3.3.7).

**Table 3.42** Transcripts identified to distinguish between PCa and non-cancer using Mann Whitney U and Lasso. Random Forest rank is also shown.

<i>Transcript</i>	<i>Mann Whitney U</i>	<i>Lasso</i>	<i>Random Forest rank</i>
<i>DLX1</i>	<i>Y</i>	<i>Y</i>	<i>1</i>
<i>ERG 3'</i>	<i>Y</i>	<i>Y</i>	<i>2</i>
<i>TMPRSS2:ERG</i>	<i>Y</i>	<i>Y</i>	<i>3</i>
<i>HOXC4</i>	<i>Y</i>	<i>Y</i>	<i>9</i>
<i>ERG 5'</i>	<i>Y</i>	<i>Y</i>	<i>8</i>
<i>HOXC6</i>	<i>Y</i>	<i>N</i>	<i>4</i>
<i>PCA3</i>	<i>Y</i>	<i>Y</i>	<i>6</i>
<i>M.genitalium RplB</i>	<i>Y</i>	<i>N</i>	<i>13</i>
<i>HPN</i>	<i>Y</i>	<i>Y</i>	<i>5</i>
<i>PPAP2A</i>	<i>N</i>	<i>Y</i>	<i>7</i>
<i>M.hYorhini rpoB</i>	<i>N</i>	<i>Y</i>	<i>10</i>
<i>ALAS1</i>	<i>N</i>	<i>Y</i>	<i>11</i>
<i>KLK3 exons 1-2</i>	<i>N</i>	<i>Y</i>	<i>16</i>
<i>NAALADL2</i>	<i>N</i>	<i>Y</i>	<i>17</i>
<i>CLU</i>	<i>N</i>	<i>Y</i>	<i>20</i>
<i>SULTA1</i>	<i>N</i>	<i>Y</i>	<i>26</i>
<i>MMP26</i>	<i>N</i>	<i>Y</i>	<i>30</i>
<i>KLK2</i>	<i>N</i>	<i>Y</i>	<i>31</i>
<i>IMPDH2</i>	<i>N</i>	<i>Y</i>	<i>37</i>
<i>Timp4</i>	<i>N</i>	<i>Y</i>	<i>41</i>
<i>PTPRC</i>	<i>N</i>	<i>N</i>	<i>12</i>
<i>SLC12A1</i>	<i>N</i>	<i>N</i>	<i>14</i>
<i>HPRT</i>	<i>N</i>	<i>N</i>	<i>15</i>

CHAPTER 3 – NANOSTRING DATA ANALYSIS 1: The Pilot Study

Four transcripts identified by three different methods (Table 3.20, Table 3.31, Table 3.39) were differentially expressed between aggressive cancer and non-aggressive cancer samples: *SLC12A1*, *SPINK1*, *SERPINB5* and *KLK3* exons 1-2.

**Table 3.43 Transcripts repeatedly shown to be differentially expressed between aggressive PCa and non-aggressive PCa.**

<i>Transcript</i>	<i>Mann Whitney U</i>	<i>Lasso</i>	<i>Random Forest rank</i>
<i>SLC12A1</i>	<b>Y</b>	<b>Y</b>	<b>5</b>
<i>UPK2</i>	<b>Y</b>	<b>N</b>	<b>4</b>
<i>SPINK1</i>	<b>Y</b>	<b>Y</b>	<b>1</b>
<i>SERPINB5</i>	<b>Y</b>	<b>Y</b>	<b>6</b>
<i>CAMKK2</i>	<b>Y</b>	<b>N</b>	<b>22</b>
<i>PSGR</i>	<b>Y</b>	<b>N</b>	<b>10</b>
<i>KLK3 exons 1-2</i>	<b>Y</b>	<b>Y</b>	<b>2</b>
<i>KLK3 exons 2-3</i>	<b>Y</b>	<b>N</b>	<b>3</b>
<i>KLK2</i>	<b>Y</b>	<b>N</b>	<b>17</b>
<i>STEAP2</i>	<b>Y</b>	<b>N</b>	<b>14</b>
<i>KLK4</i>	<b>Y</b>	<b>N</b>	<b>8</b>
<i>SULT1A1</i>	<b>N</b>	<b>N</b>	<b>7</b>
<i>BRAF</i>	<b>N</b>	<b>N</b>	<b>9</b>

Two of these transcripts were also identified by three different methods (Table 3.21, Table 3.32, Table 3.38) to be differentially expressed between advanced cancer and clinically benign cancer samples: *SLC12A1* and *SPINK1*. It is perplexing that less transcripts are selected for this extreme comparison, but this may be due to a lack of material coming from the solid advanced cancers.

**Table 3.44 Transcripts commonly found to be differentially expressed by the Mann Whitney U test, GLM and Lasso and Random Forest between advanced and benign samples.**

<i>Transcript</i>	<i>Mann Whitney U</i>	<i>Lasso</i>	<i>Random Forest rank</i>
<i>SLC12A1</i>	<b>Y</b>	<b>Y</b>	<b>5</b>
<i>SPINK1</i>	<b>Y</b>	<b>Y</b>	<b>1</b>
<i>HPRT</i>	<b>Y</b>	<b>Y</b>	<b>55</b>
<i>KLK4</i>	<b>Y</b>	<b>N</b>	<b>8</b>
<i>STEAP2</i>	<b>Y</b>	<b>N</b>	<b>14</b>
<i>PPAP2A</i>	<b>Y</b>	<b>N</b>	<b>18</b>
<i>DLX1</i>	<b>N</b>	<b>Y</b>	<b>46</b>
<i>ERG 5'</i>	<b>N</b>	<b>Y</b>	<b>42</b>
<i>HOXC4</i>	<b>N</b>	<b>Y</b>	<b>52</b>
<i>HPN</i>	<b>N</b>	<b>Y</b>	<b>11</b>
<i>STEAP4</i>	<b>N</b>	<b>Y</b>	<b>57</b>
<i>SULT1A1</i>	<b>N</b>	<b>Y</b>	<b>7</b>

### 3.10 Conclusions

Detection of prostate-specific (*KLK2* and *KLK3*) and PCa-specific (*TMPRSS2:ERG*) transcripts demonstrates that these are present in urine EVs harvested and analysed by our methods. RNA yields post-radical prostatectomy suggests that the vast majority of the urinary RNA originates in the prostate. The identification of differential transcripts between non-aggressive and aggressive cancers demonstrates NanoString's potential ability to distinguish these clinical categories using transcripts from urinary EVs.

It is vital to emphasise that the clinical categories in this study are based on current, and not perfect clinical tests. Hence the current need for novel biomarkers to distinguish accurately between them. Particularly true of CB samples, where it is expected that ~20% of the men that show no clinical evidence of cancer will in fact have PCa. Therefore, it is notable that 12% of CB samples were found to have a *TMPRSS2:ERG* fusion in this study. As *TMPRSS2:ERG* is expected to be in ~50% of PCa, this would suggest clinically undetected PCa in 24% of our CB samples. In LPD, 21 of the 37 CB samples are clustered together, leaving 16 spread amongst the other groups, five (14%) of which are associated to a group where overall *TMPRSS2:ERG* is significantly up-regulated. Seven of the CB samples are left un-grouped, showing no distinct underlying signature. The detection of *TMPRSS2:ERG* by NanoString and confirmation of this find by RTPCR demonstrated the sensitivity of our methods for detection of PCa.

Some Nanostring probes performed much better than others in models throughout the analyses (section 3.9). However, transcripts were identified that were differentially expressed in samples from different clinical categories (PCa present, increasing PCa aggressiveness etc). Due to the nature of the probe set, most of these were known PCa markers, but some were not. The latter demonstrating that it can be difficult to predict what type of transcripts should be targeted in our analyses. As a result of probe selection for advanced PCa associated transcripts, it was expected that we observe unusual distribution and medians for our data.

## CHAPTER 3 – NANOSTRING DATA ANALYSIS 1: The Pilot Study

Housekeeping probes were selected, as they are known to be useful when investigating PCa tissue; their use in urine and urinary EVs has not been studied in detail. Therefore, it may be of interest to investigate further options for urinary microvesicle house keeping transcripts, as the probes selected did not show great correlation.

The analyses have revealed that there is structure in the data, as demonstrated for example, by the detection of differentially expressed transcripts, LPD groups and linear model analysis. Lasso logistic regression predictive models were able to categorise cancer from non-cancer samples and aggressive from non-aggressive cancer samples fairly robustly (AUCs of 0.94 and 0.85, respectively). However, sensitivity and specificity, even on the training set could be improved. For this it could be suggested that we are not using the optimal starting probes and thus more probes should be included to identify the clearest signature available. Another complication is the complexity of cancer within individual prostates, with multifocal tumours being detectable in the majority of cancerous prostates, each with the potential to have a different path and rate of progression. The Mann Whitney U test and random forest results were similar to each other and that of the Lasso models, suggesting that these results are accurate, but also highlighting that the methods were suitable for analysing the NanoString data. The LPD showed some clinical separation of the samples, though again a better selection of probes could provide further discrimination between the lower and intermediate samples with the benign samples. The inclusion of known PCa transcripts in our differential expression and predictive model for cancer results and the inclusion of known prognostic PCa transcripts in our differential expression and predictive model for aggressive cancer results provide evidence of accuracy.

These analyses form the ground work for expansion of the urine biomarker study to include a larger number of probes, and samples which should provide much improved power to dissect the complexities of this disease within individual prostates. The probes that provided no information were determined. These probes were reviewed to unveil if they should be replaced or removed in the larger study (for example the *PCA3* probe didn't work very well and was redesigned for the

## CHAPTER 3 – NANOSTRING DATA ANALYSIS 1: The Pilot Study

large study). *FOXMI* showed no clinical association and was not identified in any clustering or prediction models and so was removed from subsequent studies.

# 4

## Response to Hormone Therapy

### 4.1 Summary

Stratification of treatment by gene expression levels has shown benefits in many cancers, such as breast cancer<sup>208</sup> and lung cancer<sup>209,210</sup>, but it is yet to be utilised successfully in prostate cancer (PCa) treatment. Areas where stratification could benefit PCa patients include: deciding between treatment vs. active surveillance, identifying which radical prostatectomy (RP) and radiotherapy patients will succumb to biochemical recurrence (BCR), and which hormone therapy (HT) patients would benefit from additional treatment (i.e. those patients that are predicted to progress early to castration resistant prostate cancer (CRPC)). In this chapter we focus on men in our

## CHAPTER 4: RESPONSE TO HORMONE THERAPY

cohort treated by hormone therapy and examine whether expression profiles of urinary microvesicles can be used to predict time to CRPC. Unfortunately, the cohort does not have a long enough follow up at this time to examine time to BCR after RP or radiotherapy.

In the normalised NanoString 1 dataset, a signature of seven transcripts was identified that could optimally predict progression of patients on hormone therapy (section 1.3.4.2.1) (cox-regression model;  $p = 2.3 \times 10^{-05}$ ; HR = 0.04288). The transcripts in the predictor are *AGR2*, *DLX1*, *KLK2*, *NAALADL2*, *AR* exons 4-8, *PPAP2A* and *AMACR*. This model was an independent predictor of progression when established clinical variables initial PSA, age, Gleason score and initial bone scan result were taken into account (cox-regression model;  $p = 0.003$ ; HR = 0.03).

When the data was adjusted to *KLK2* levels, similar to *KLK3* adjustment used in the PCA3 test, an optimal model of three transcripts (*CAMKK2*, *PSGR* and *UPK*) was identified (cox-regression model;  $p = 0.007$ , HR = 1.0028). This model was not a significant predictor independent of established clinical factors (cox-regression model;  $p = 0.14$ ; HR = 1.009).

Both of these models were applied to the second NanoString dataset but were not validated.

### **4.2 Introduction**

#### **4.2.1 The Research Gap**

Hormone Therapy (HT) is the primary treatment of men with advanced prostate cancer, that is those diagnosed with a PSA > 100 or with evidence of metastatic spread (generally via bone scan). Response to the treatment is highly variable with some men failing to respond at all, whereas others take years to progress. All men will eventually progress to CRPC (section 1.3.4.2.2). Identification of men that are likely to relapse early could lead to more aggressive first line treatment being used, such as full



## CHAPTER 4: RESPONSE TO HORMONE THERAPY

androgen blockage (section 1.3.4.2.1), combination with chemotherapy, or novel strategies such as combination with Abiraterone. Currently, there is no clinically available test to stratify advanced patients into those who will do well on HT and those that will quickly require further or alternative treatments.

### 4.2.2 Aim

I am to use the NanoString 1 data set (Chapter 3) from advanced patients ( $n = 32$ ), to see if a significant predictor of early progression in patients on HT can be built and whether this predictor improves on current clinical information collected (e.g. PSA, Gleason score and bone scan). I will also attempt to validate these signatures in a second independent cohort (using the second NanoString data).

### 4.2.3 Summary of the HT patient cohort

The breakdown of the clinical data for the 32 patients on HT can be seen in Table 4.1. Many of the advanced patients are diagnosed as being advanced by a PSA  $> 100$  and no biopsy is performed in these circumstances. Other patients with lower PSAs are determined to be advanced, by either a biopsy or a positive bone scan.

**Table 4.1 Clinical summary of the hormone therapy cohort ( $n=32$ ).**

<i>Clinical Variable</i>	<i>Number of patients</i>
<i>Gleason Score</i>	
10	0
9	8
8	4
7 (4+3)	4
<i>No Biopsy: Advanced</i>	<b>16</b>
<i>Bone Scan</i>	
Positive	18
Negative	13
Unknown	1
<i>PSA</i>	<b>Median 98.7 (range: 9.6 – 2508)</b>
<i>Age</i>	<b>Median 78 (range: 55 - 98)</b>

**4.3 Hormone Therapy Predictor constructed using Nanostring 1 data****4.3.1 Differentially expressed genes based on initial response, 12 month relapse and 24 month relapse**

Five transcripts were significantly up regulated in those that had an initial response to HT ( $n = 28$ ) compared to those that did not ( $n = 4$ ) ( $p < 0.05$ ; not adjusted for multiple testing; Mann-Whitney U test; Table 4.2). Three of these five transcripts, were capable of distinguishing patients that relapsed within 12 months ( $n = 6$ ) (Table 4.2;  $p < 0.05$ ; not adjusted for multiple testing): *STEAP4*, *AMACR*, *BRAF*. By 24 months, 14 patients were still responding to HT and 18 had progressed. Four different transcripts were significantly up regulated in patients still responding to treatment (Table 4.2;  $p < 0.05$ ; not adjusted for multiple testing). These results need to be treated with caution due to the low numbers and the lack of significance after multiple testing correction.

**Table 4.2 Mann-Whitney U test results for comparing samples that respond to HT and those that don't at different time points.**

<i>Initial Response</i>			
<i>Transcript</i>	<i>p - value</i>	<i>Adjusted p - value</i>	<i>Log2 Fold change</i>
<i>AGR2</i>	<b>0.047</b>	<b>1</b>	<b>-0.57</b>
<i>STEAP4</i>	<b>0.024</b>	<b>1</b>	<b>-0.26</b>
<i>HPRT</i>	<b>0.029</b>	<b>1</b>	<b>-0.21</b>
<i>AMACR</i>	<b>0.034</b>	<b>1</b>	<b>-0.14</b>
<i>BRAF</i>	<b>0.04</b>	<b>1</b>	<b>-0.13</b>
<i>12 Month Response</i>			
<i>Transcript</i>	<i>p - value</i>	<i>Adjusted p - value</i>	<i>Log2 Fold change</i>
<i>STEAP4</i>	<b>0.019</b>	<b>0.98</b>	<b>-0.18</b>
<i>AMACR</i>	<b>0.01</b>	<b>0.57</b>	<b>-0.14</b>
<i>BRAF</i>	<b>0.033</b>	<b>0.98</b>	<b>-0.08</b>
<i>24 Month Response</i>			
<i>Transcript</i>	<i>p - value</i>	<i>Adjusted p - value</i>	<i>Log2 Fold change</i>
<i>DLX1</i>	<b>0.045</b>	<b>1</b>	<b>-1.28</b>
<i>AR (truncation) exon 9</i>	<b>0.025</b>	<b>1</b>	<b>-1.16</b>
<i>AR exons 4-8</i>	<b>0.018</b>	<b>1</b>	<b>-0.09</b>
<i>TBP</i>	<b>0.034</b>	<b>1</b>	<b>-0.08</b>

**4.3.2 Survival analyses of time to progression after HT**

Using the Mann-Whitney U test as described above lacks statistical power as time to progression is not taken into account, therefore I applied Cox's proportional hazards model and other survival analysis tools (section 2.8). Twelve probes were significant predictors of progression individually (Table 4.3; Cox regression model;  $p < 0.05$ ; multiple testing correction not applied). There were no significant probes after multiple testing correction.

Expression for each gene was divided into two groups, low expression and high expression, using  $k$ -means to determine the threshold (section 2.5.3). Using these grouped data, ten transcripts were identified as having significant different times to progression ( $p < 0.05$ ; log-rank test; Table 4.4), of which only one was significant after multiple testing correction: *NAALADL2*.

**Table 4.3 Cox results for relapse to hormone therapy**

<i>Transcript</i>	<i>p - value</i>	<i>Adjusted p - value</i>	<i>Hazard ratio</i>
<i>KLK2</i>	<b>0.011</b>	<b>0.62</b>	<b>0.74</b>
<i>AMACR</i>	<b>0.011</b>	<b>0.62</b>	<b>0.68</b>
<i>DLX1</i>	<b>0.011</b>	<b>0.62</b>	<b>0.87</b>
<i>PPAP2A</i>	<b>0.014</b>	<b>0.76</b>	<b>0.51</b>
<i>STEAP4</i>	<b>0.017</b>	<b>0.88</b>	<b>0.63</b>
<i>PCA3</i>	<b>0.034</b>	<b>1.00</b>	<b>0.87</b>
<i>CDC20</i>	<b>0.037</b>	<b>1.00</b>	<b>0.81</b>
<i>KLK4</i>	<b>0.039</b>	<b>1.00</b>	<b>0.63</b>
<i>TDRD</i>	<b>0.042</b>	<b>1.00</b>	<b>0.86</b>
<i>STEAP2</i>	<b>0.043</b>	<b>1.00</b>	<b>0.66</b>
<i>NAALADL2</i>	<b>0.045</b>	<b>1.00</b>	<b>0.79</b>
<i>Timp4</i>	<b>0.049</b>	<b>1.00</b>	<b>0.86</b>

**Table 4.4 Significant probes using log rank test applied to data separated by  $k$ -means.**

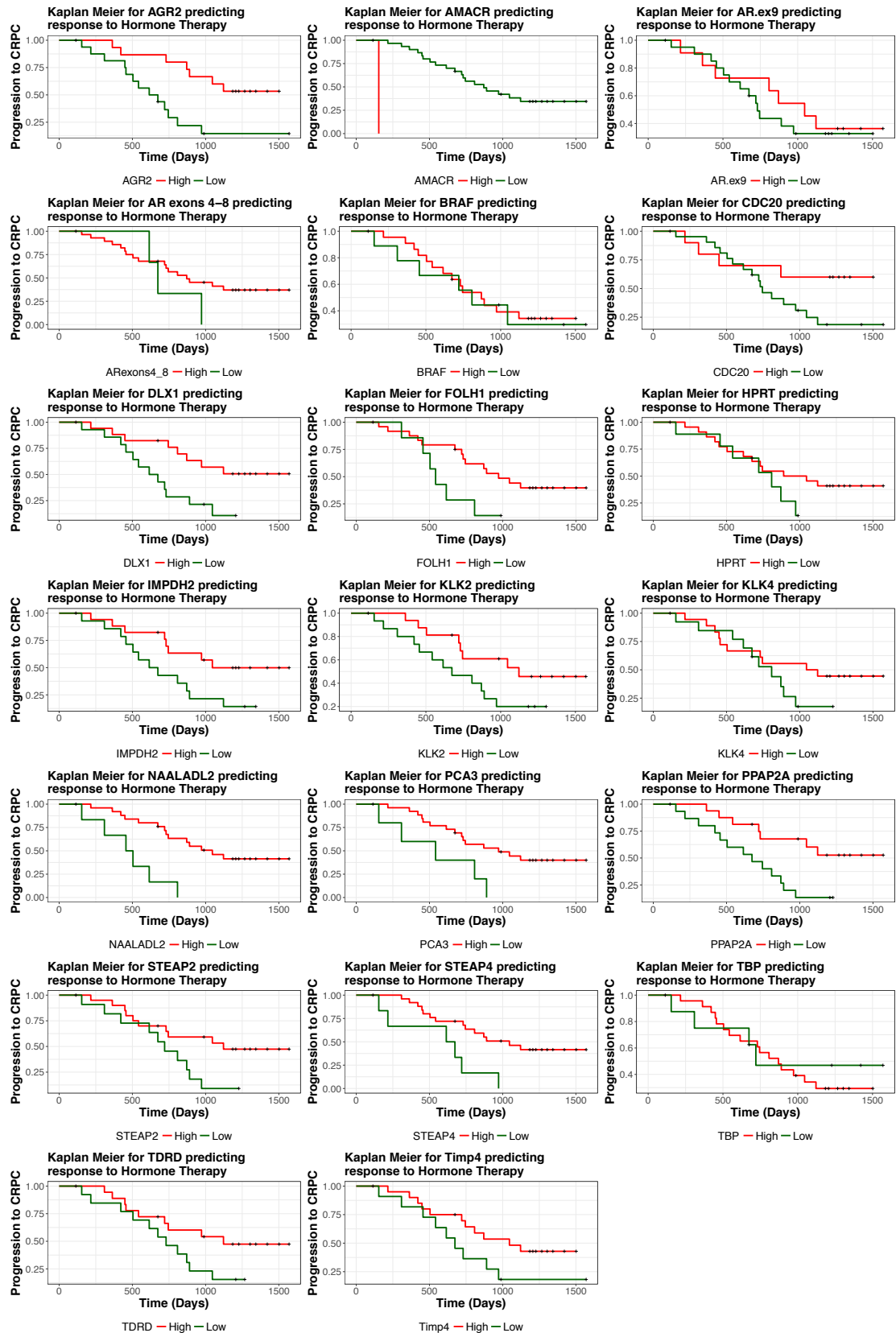
<i>Transcript</i>	<i>p - value</i>	<i>Adjusted p - value</i>	<i>Coefficient</i>
<i>NAALADL2</i>	<b>0.0004</b>	<b>0.03</b>	<b>12.36</b>
<i>PPAP2A</i>	<b>0.005</b>	<b>0.27</b>	<b>7.97</b>
<i>KLK2</i>	<b>0.006</b>	<b>0.31</b>	<b>7.66</b>
<i>STEAP4</i>	<b>0.007</b>	<b>0.4</b>	<b>7.19</b>
<i>DLX1</i>	<b>0.01</b>	<b>0.56</b>	<b>6.55</b>
<i>AGR2</i>	<b>0.01</b>	<b>0.64</b>	<b>6.28</b>
<i>PCA3</i>	<b>0.02</b>	<b>0.93</b>	<b>5.57</b>
<i>IMPDH2</i>	<b>0.03</b>	<b>0.98</b>	<b>4.68</b>
<i>STEAP2</i>	<b>0.03</b>	<b>0.98</b>	<b>4.63</b>

## CHAPTER 4: RESPONSE TO HORMONE THERAPY

<i>FOLHI</i>	<i>0.04</i>	<i>0.98</i>	<i>4.09</i>
--------------	-------------	-------------	-------------

Twenty transcripts have been identified as candidate predictors of progression after HT (Table 4.3 and Table 4.4). For the majority of probes a clear difference in time to progression is seen (Figure 4.1).

## CHAPTER 4: RESPONSE TO HORMONE THERAPY



**Figure 4.1** Kaplan Meier plots for each of the candidate probes (section 4.3.1). Expression for each probe is grouped into high and low expression using K-means clustering.

### 4.3.3 Determining the optimal predictor of progression after HT

The optimal model to detect time to progression after HT is likely to be formed from a combination of the expression from multiple probes. There are various methods for identifying the best combination of probes (variable selection) and here I will investigate three i.e. LASSO (section 2.6.2), stepwise regression (section 2.6.4) and random forest (section 2.6.3). Different starting sets of probes will be used based on results from the previous section.

#### 4.3.3.1 Model built using differentially expressed transcripts based on initial response, 12 month relapse and 24 month relapse

Gene selection and three proposed optimal models were produced based on the nine transcripts identified as differentially expressed at initial response, 12 month relapse or 24 month relapse (Table 4.2): a Cox general linear model with shrinkage and variable selection using LASSO (section 2.6.2) (Table 4.5), Stepwise regression on a Cox model (Table 4.6), and a Random Forest model (section 2.6.3) (Table 4.7). The five transcripts selected by LASSO and step are identical, showing reliability in these results. Four out of these five transcripts most important in the random forest model are also similar (*DLX1*, *AR* exons 4-8, *AMACR* and *AGR2* have high importance), though *STEAP4* appears to have increased importance and *BRAF* has lost importance in the Random Forest model.

**Table 4.5** The probes included the in the glm after LASSO shrinkage and variable selection, (of the Mann-Whitney U selected probes) with the corresponding beta coefficients

<i>Transcript</i>	<i>Beta Coefficient</i>
<i>BRAF</i>	<i>0.27</i>
<i>DLX1</i>	<i>-0.14</i>
<i>AGR2</i>	<i>-0.19</i>
<i>AR exons 4-8</i>	<i>-0.26</i>
<i>AMACR</i>	<i>-0.38</i>

## CHAPTER 4: RESPONSE TO HORMONE THERAPY

**Table 4.6** The probes included in the cox model after step variable selection (of the Mann-Whitney U selected probes) with the hazard values and *p*-values. The overall performance of the model to predict progression on HT is  $p = 0.00024$ .

<i>Transcript</i>	<i>HR</i>	<i>p-value</i>
<i>DLXI</i>	<b>0.83</b>	<b>0.008</b>
<i>AMACR</i>	<b>0.56</b>	<b>0.013</b>
<i>AGR2</i>	<b>0.76</b>	<b>0.038</b>
<i>AR exons 4-8</i>	<b>0.72</b>	<b>0.047</b>
<i>BRAF</i>	<b>1.67</b>	<b>0.054</b>

**Table 4.7** The importance of each probe in the random forest predictor for HT relapse (of the Mann-Whitney U selected probes).

<i>Transcript</i>	<i>Importance</i>	<i>Relative Importance</i>
<i>DLXI</i>	<b>0.042</b>	<b>1</b>
<i>AR exons 4-8</i>	<b>0.015</b>	<b>0.35</b>
<i>STEAP4</i>	<b>0.013</b>	<b>0.32</b>
<i>AMACR</i>	<b>0.011</b>	<b>0.27</b>
<i>AGR2</i>	<b>0.009</b>	<b>0.21</b>
<i>HPRT</i>	<b>0.003</b>	<b>0.08</b>
<i>BRAF</i>	<b>-0.001</b>	<b>-0.01</b>
<i>TBP</i>	<b>-0.003</b>	<b>-0.08</b>
<i>AR exon 9</i>	<b>-0.006</b>	<b>-0.13</b>

Using the selected gene sets determined above, a single score was derived for each gene set and a Cox regression model was constructed (Table 4.8). The top four important transcripts of the random forest model performed best (HR=0.0573;  $p = 3.29 \times 10^{-05}$ ) but all of the models were highly significant ( $p < 1.0 \times 10^{-03}$ ) in discriminating samples from patient's that progressed on HT and those that did not.

**Table 4.8 Overall performance of the models (created from the probes originally identified by Mann-Whitney U) tested by cox.**

<i>Model</i>	<i>p-value</i>	<i>HR (95% confidence intervals)</i>
<i>LASSO genes – (DLX1, AGR2, BRAF, AR exon 9 /AMACR)</i>	$7.5 \times 10^{-04}$	<b>0.106 (0.01761 - 0.6394)</b>
<i>Step genes – (DLX1, AGR2, BRAF, AR exons 4-8 /AMACR)</i>	$5.31 \times 10^{-05}$	<b>0.0786 (0.01299 - 0.4752)</b>
<i>LASSO and Step genes (DLX1, AGR2, BRAF, AR exons 4-8, AR exon 9 /AMACR)</i>	$6.79 \times 10^{-04}$	<b>0.0983 (0.01738 - 0.5567)</b>
<i>Random Forest top 5 genes (DLX1, AGR2, AR exon 4-8, STEAP4 /AMACR)</i>	$3.29 \times 10^{-05}$	<b>0.0573 (0.008039 - 0.4079)</b>

#### 4.3.3.2 Model built using Cox selected transcripts

Using the twelve transcripts identified using the Cox regression model on individual probes (Table 4.3), variable selection was performed. LASSO identified seven transcripts (Table 4.9), stepwise regression identified six transcripts (Table 4.10), and Random Forest identified the relative importance (Table 4.11). The transcripts selected by LASSO and stepwise regression have three common transcripts (*KLK2*, *CDC20* and *STEAP2*) but the importance of these probes was not high in the random forest model.

**Table 4.9 The probes included the in the glm after LASSO shrinkage and variable selection, (of the cox selected probes) with the corresponding beta coefficients**

<i>Transcript</i>	<i>Beta Coefficient</i>
<i>KLK2</i>	<b>-0.009</b>
<i>CDC20</i>	<b>-0.012</b>
<i>PPAP2A</i>	<b>-0.025</b>
<i>STEAP4</i>	<b>-0.032</b>
<i>DLX1</i>	<b>-0.059</b>
<i>NAALADL2</i>	<b>-0.072</b>
<i>STEAP2</i>	<b>-0.076</b>



CHAPTER 4: RESPONSE TO HORMONE THERAPY

**Table 4.10** The probes included in the cox model after step variable selection (of the cox selected probes) with the hazard values and *p*-values. The overall performance of the model to predict progression on HT is  $p = 0.00323$ .

<i>Transcript</i>	<i>HR</i>	<i>p-value</i>
<i>AMACR</i>	<b>0.493</b>	<b>0.003</b>
<i>KLK2</i>	<b>0.496</b>	<b>0.008</b>
<i>STEAP2</i>	<b>0.446</b>	<b>0.057</b>
<i>PCA3</i>	<b>1.395</b>	<b>0.065</b>
<i>KLK4</i>	<b>1.863</b>	<b>0.094</b>
<i>CDC20</i>	<b>0.799</b>	<b>0.097</b>

**Table 4.11** The importance of each probe in the random forest predictor for HT relapse (of the cox selected probes).

<i>Transcript</i>	<i>Importance</i>	<i>Relative Importance</i>
<i>NAALADL2</i>	<b>0.0215</b>	<b>1</b>
<i>AMACR</i>	<b>0.0199</b>	<b>0.928</b>
<i>DLX1</i>	<b>0.0178</b>	<b>0.829</b>
<i>STEAP4</i>	<b>0.0067</b>	<b>0.313</b>
<i>PPAP2A</i>	<b>0.0051</b>	<b>0.239</b>
<i>STEAP2</i>	<b>0.001</b>	<b>0.046</b>
<i>TDRD</i>	<b>0.001</b>	<b>0.046</b>
<i>Timp4</i>	<b>0.0008</b>	<b>0.039</b>
<i>PCA3</i>	<b>0.0002</b>	<b>0.008</b>
<i>CDC20</i>	<b>-0.0013</b>	<b>-0.062</b>
<i>KLK2</i>	<b>-0.0033</b>	<b>-0.154</b>
<i>KLK4</i>	<b>-0.0049</b>	<b>-0.226</b>

The combined score Cox regressions (Table 4.12) showed that the top four important transcripts selected by the random forest model performed best (HR = 0.103;  $p = 7.97 \times 10^{-05}$ ) but all were statistically significant in predicting patient's that progressed on HT.

**Table 4.12 Overall performance of the models (created from the probes originally identified by cox) tested by cox.**

<i>Model</i>	<i>p-value</i>	<i>HR (95% confidence intervals)</i>
<i>LASSO genes – (KLK2, DLX1, NAALADL2, PPAP2A, STEAP2, STEAP4, CDC20)</i>	$7.8 \times 10^{-04}$	<b>0.00038</b> ( $1.168 \times 10^{-06}$ - 0.1239)
<i>Step genes – (KLK2, STEAP2, PCA3, STEAP4, CDC20 /AMACR)</i>	<b>0.01</b>	<b>0.048 (0.004 - 0.641)</b>
<i>LASSO and Step genes (DLX1, NAALADL2, STEAP4, KLK2, STEAP2, PPAP2A, CDC20, PCA3, KLK4 /AMACR)</i>	$1.22 \times 10^{-03}$	<b>0.001</b> ( $5.4 \times 10^{-06}$ - 0.192)
<i>Common to LASSO and Step genes (KLK2, CDC20, STEAP2)</i>	$8.22 \times 10^{-03}$	<b>0.026 (0.001 - 0.536)</b>
<i>Random Forest top 4 genes (NAALADL2, DLX1, STEAP4 /AMACR)</i>	$7.97 \times 10^{-05}$	<b>0.103 (0.021 - 0.504)</b>

### 4.3.3.3 Model built from Log-rank selected transcripts

Variable selection was performed using the ten significant probes identified by the Log rank test (expression divided into two groups using *k*-means) (Table 4.4). LASSO identified seven transcripts (Table 4.13), stepwise regression identified five transcripts (Table 4.14), and Random forest identified the relative importance of each of the ten transcripts (Table 4.15). The transcripts selected by LASSO and step have four in common (*KLK2*, *AGR2*, *DLX1* and *STEAP4*). These four transcripts are also the most important in the random forest model.

**Table 4.13 The probes included the in the glm after LASSO shrinkage and variable selection, (of the log-rank selected probes) with the corresponding beta coefficients**

<i>Transcript</i>	<i>Beta Coefficient</i>
<i>KLK2</i>	<b>-0.003</b>
<i>AGR2</i>	<b>-0.02</b>
<i>PPAP2A</i>	<b>-0.03</b>
<i>STEAP4</i>	<b>-0.03</b>
<i>DLX1</i>	<b>-0.06</b>
<i>NAALADL2</i>	<b>-0.07</b>
<i>STEAP2</i>	<b>-0.08</b>

CHAPTER 4: RESPONSE TO HORMONE THERAPY

**Table 4.14** The probes included in the cox model after step variable selection (of the log-rank selected probes) with the hazard values and *p*-values. The overall performance of the model to predict progression on HT is  $p = 0.0012$ .

<i>Transcript</i>	<i>HR</i>	<i>p-value</i>
<i>DLX1</i>	<b>0.78</b>	<b>0.001</b>
<i>NAALADL2</i>	<b>0.66</b>	<b>0.01</b>
<i>FOLH1</i>	<b>1.44</b>	<b>0.03</b>
<i>AGR2</i>	<b>0.75</b>	<b>0.04</b>
<i>STEAP4</i>	<b>0.69</b>	<b>0.08</b>

**Table 4.15** The importance of each probe in the random forest predictor for HT relapse (of the log-rank selected probes).

<i>Transcript</i>	<i>Importance</i>	<i>Relative Importance</i>
<i>NAALADL2</i>	<b>0.0258</b>	<b>1</b>
<i>DLX1</i>	<b>0.0181</b>	<b>0.701</b>
<i>AGR2</i>	<b>0.0151</b>	<b>0.583</b>
<i>STEAP4</i>	<b>0.0138</b>	<b>0.536</b>
<i>PPAP2A</i>	<b>0.0088</b>	<b>0.34</b>
<i>KLK2</i>	<b>0.0064</b>	<b>0.246</b>
<i>STEAP2</i>	<b>0.0031</b>	<b>0.119</b>
<i>FOLH1</i>	<b>-0.0039</b>	<b>-0.153</b>
<i>PCA3</i>	<b>0.0052</b>	<b>-0.202</b>
<i>IMPDH2</i>	<b>-0.096</b>	<b>-0.373</b>

In the combined score Cox regressions (Table 4.12) showed that the LASSO selected transcripts performed marginally better ( $HR = 0.01$ ;  $p = 4.7 \times 10^{-04}$ ) but all were statistically significant in predicting patient's that progressed on HT.

**Table 4.16 Overall performance of the models (created from the probes originally identified by log rank) tested by cox.**

<i>Model</i>	<i>p-value</i>	<i>HR (95% confidence intervals)</i>
<i>LASSO – (KLK, DLX1, NAALADL2, PPAP2A, STEAP2, STEAP4, AGR2)</i>	<i>4.7x10<sup>-04</sup></i>	<i>0.01 (0.0004 - 0.2654)</i>
<i>Step – (DLX1, NAALADL2, STEAP4, AGR2, FOLH1)</i>	<i>7x10<sup>-04</sup></i>	<i>0.0212 (0.0013 - 0.3362)</i>
<i>LASSO and Step (KLK2, DLX1, NAALADL2, PPAP2A, STEAP2, STEAP4, AGR2, FOLH1)</i>	<i>6x10<sup>-04</sup></i>	<i>0.0072 (0.0002- 0.2865)</i>
<i>Common to LASSO, Step and Random Forest (KLK2, AGR2, DLX1, STEAP4)</i>	<i>7x10<sup>-04</sup></i>	<i>0.0275 (0.0023 - 0.326)</i>

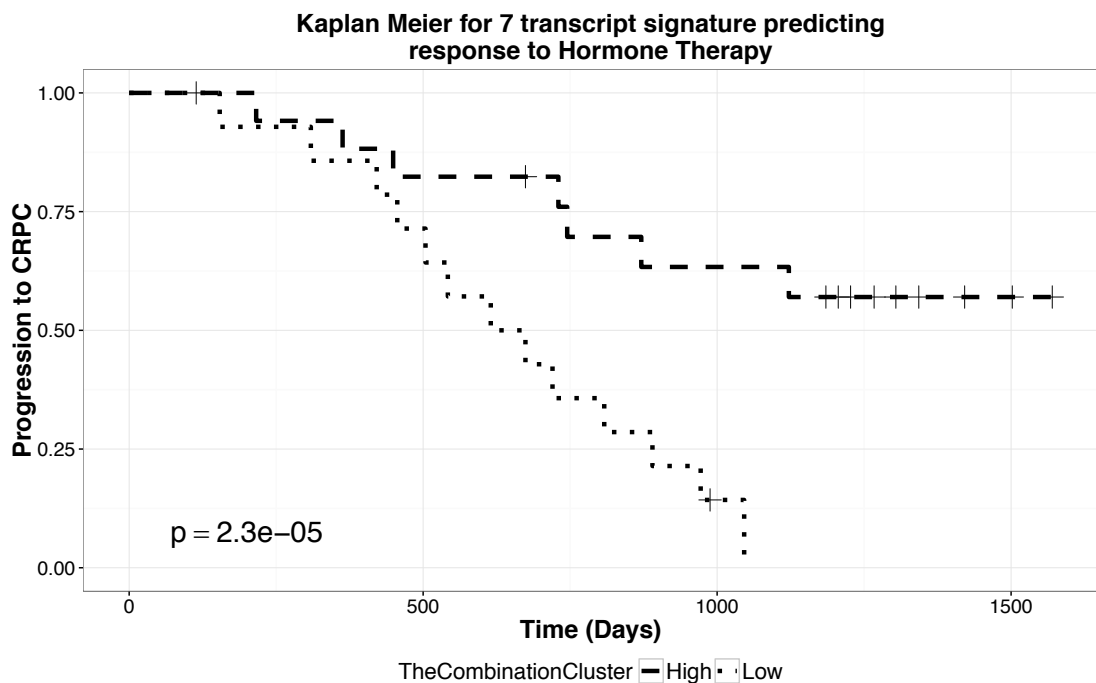
#### ***4.3.3.4 Model built using combining probe selection lists to produce final model***

Combining the candidate probe lists identified using the different gene selection model may produce a better predictor of HT progression. Using each combination of candidate probe lists (Table 4.2, Table 4.3, & Table 4.4), LASSO was applied for variable selection (as it is clear and is designed to avoid over-fitting) and a linear combination score with Cox regression model was produced (Table 4.17). Initiating the variable selection with a combination of the candidate probes identified as differentially expressed at initial response, 12 month relapse or 24 month relapse (Mann Whitney U) and by the log rank test produced the best model ( $p = 1.3 \times 10^{-04}$ ; HR = 0.0091). This model includes *AGR2*, *AR* exons 4-8, *DLX1*, *KLK2*, *NAALADL2*, *PPAP2A* and *AMACR*. It has an AUC of 0.783 (Figure 4.1). This is the best performing model constructed using the Nanostring 1 data. The Kaplan Meier plot for the seven-transcripts combined was also produced (Figure 4.2).

CHAPTER 4: RESPONSE TO HORMONE THERAPY

**Table 4.17 Comparing the Cox regression models of various linear combination scores producing from combining gene selection lists. Mann-Whitney U = candidate probes identified as differentially expressed at initial response, 12 month relapse or 24 month relapse; cox = candidate probes identified by step applied to cox regression models; Log rank = candidate probes identified by the log rank test on expression dichotomised into low and high expression.**

<i>Combination</i>	<i>Method for variable selection</i>	<i>Resulting probes in model</i>	<i>(cox) p-value</i>	<i>(cox) HR</i>
<i>Mann-Whitney U and cox</i>	<b>LASSO</b>	<i>AGR2, AR exons 4-8, DLX1, KLK2, NAALADL2, TDRD/AMACR</i>	$1.3 \times 10^{-04}$	<b>0.0091</b> (0.0004 - 0.214)
<i>Mann-Whitney U and Log rank</i>	<b>LASSO</b>	<i>AGR2, AR exons 4-8, DLX1, KLK2, NAALADL2, PPAP2A/AMACR</i>	$2.3 \times 10^{-05}$	<b>0.04288</b> (0.005 - 0.345)
<i>Cox and Log rank</i>	<b>LASSO</b>	<i>DLX1</i>	<b>0.01</b>	<b>0.871 (0.781 - 0.972)</b>



**Figure 4.2 Kaplan Meier showing the seven-transcript signature (*AR* exons 4-8 \* *AGR2* \* *DLX1* \* *KLK2* \* *NAALADL2* \* *PPAP2A* / *AMACR*) separated into low and high expression using *k*-means. The significance was measured using the cox model (Table 4.17),  $p = 2.3 \times 10^{-05}$ .**

## CHAPTER 4: RESPONSE TO HORMONE THERAPY

The seven-transcript signature was a better predictor of HT relapse than other clinical variables (including bone scan outcome, Gleason score, initial PSA value and age) when treated individually (Table 4.18). LPD group (identified from chapter 3, section 3.5.5) was also tested. The seven-transcript score was a statistically significant independent predictor of HT progression when combined with covariate clinical variables ( $p = 0.003$ ; HR = 0.03; Table 4.19).

**Table 4.18 Univariate cox models showing the significance of clinical variables, LPD group and the seven-transcript signature on predicting HT relapse.**

<b>UNIVARIATE MODELS</b>		
<i>Model</i>	<i>p-value</i>	<i>HR (95% CI)</i>
<i>Age</i>	<b>0.84</b>	<b>1.006 (0.949 - 1.067)</b>
<i>PSA</i>	<b>0.05</b>	<b>1.001 (1 - 1.002)</b>
<i>Gleason Scores</i>	<b>0.27</b>	
<i>Gleason 7</i>		
<i>Gleason 8</i>	<b>0.26</b>	<b>0.252 (0.022 - 2.844)</b>
<i>Gleason 9</i>	<b>0.91</b>	<b>1.101 (0.22 - 5.54)</b>
<i>Gleason Category</i>		
<i>Gleason 7+8</i>		
<i>Gleason 9+NA</i>	<b>0.15</b>	<b>2.291(0.666-7.883)</b>
<i>Bone Scan</i>	<b>0.19</b>	
<i>Negative</i>		
<i>Positive</i>	<b>0.2</b>	<b>1.854 (0.719 - 4.785)</b>
<i>LPD group</i>	<b>0.09</b>	
<i>LPD1</i>		
<i>LPD2</i>	<b>0.78</b>	<b>1.413 (0.128 - 15.59)</b>
<i>LPD3</i>	<b>0.24</b>	<b>3.716 (0.414 - 33.36)</b>
<i>LPD4</i>	<b>0.07</b>	<b>7.043 (0.844 - 58.77)</b>
<i>LPD NA</i>	<b>0.12</b>	<b>5.589 (0.637 - 49.04)</b>
<i>DLX1 * AGR2 * KLK2 * NAALADL2 * AR exons 4-8 * PPAP2A / AMACR</i>	<b><math>2.3 \times 10^{-05}</math></b>	<b>0.04288 (0.005 - 0.345)</b>

**Table 4.19 Multivariate cox model for predicting early relapse on HT.**

<b>MULTIVARIATE MODEL with the seven transcript signature, <i>p</i>-value = <math>7.7 \times 10^{-04}</math></b>		
<i>Variable</i>	<i>p-value</i>	<i>HR (95% CI)</i>
<i>DLX1 * AGR2 * KLK2 * NAALADL2 * AR exons 4-8 * PPAP2A / AMACR</i>	<b>0.003</b>	<b>0.03 (0.003 – 0.313)</b>
<i>Age</i>	<b>0.996</b>	<b>1 (0.949 - 1.054)</b>
<i>PSA</i>	<b>0.176</b>	<b>1 (0.997 – 1.001)</b>
<i>Gleason Category</i>		
<i>Gleason 7 + 8</i>		
<i>Gleason 9 + NA</i>	<b>0.167</b>	<b>2.61 (0.67 – 10.143)</b>
<i>Bone Scan</i>		
<i>Negative</i>		
<i>Positive</i>	<b>0.276</b>	<b>1.85 (0.612 – 5.59)</b>

#### 4.3.4 Validation of the seven-transcript signature using NanoString 2 data

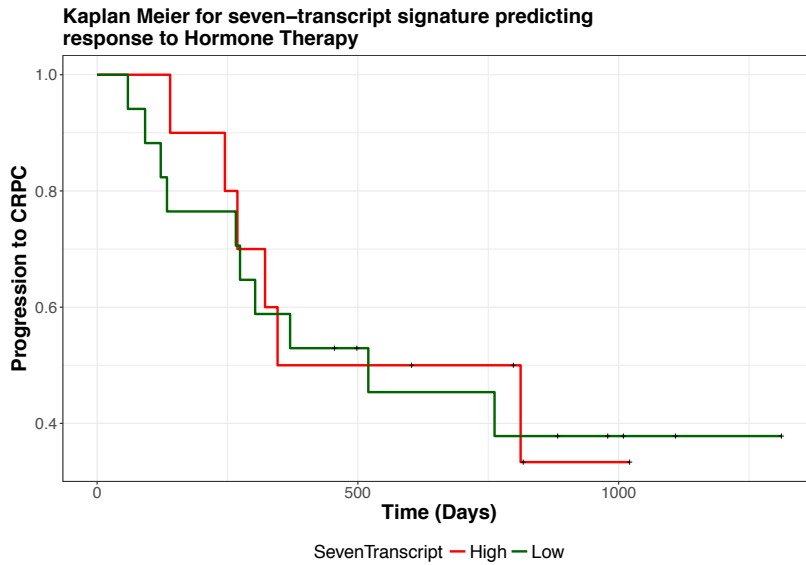
The second set of NanoString data had 43 patients on HT (chapter 5), of which 27 samples were unique to NanoString 2 (Table 4.20).

**Table 4.20 Clinical breakdown of the 27 HT patients unique to NanoString 2.**

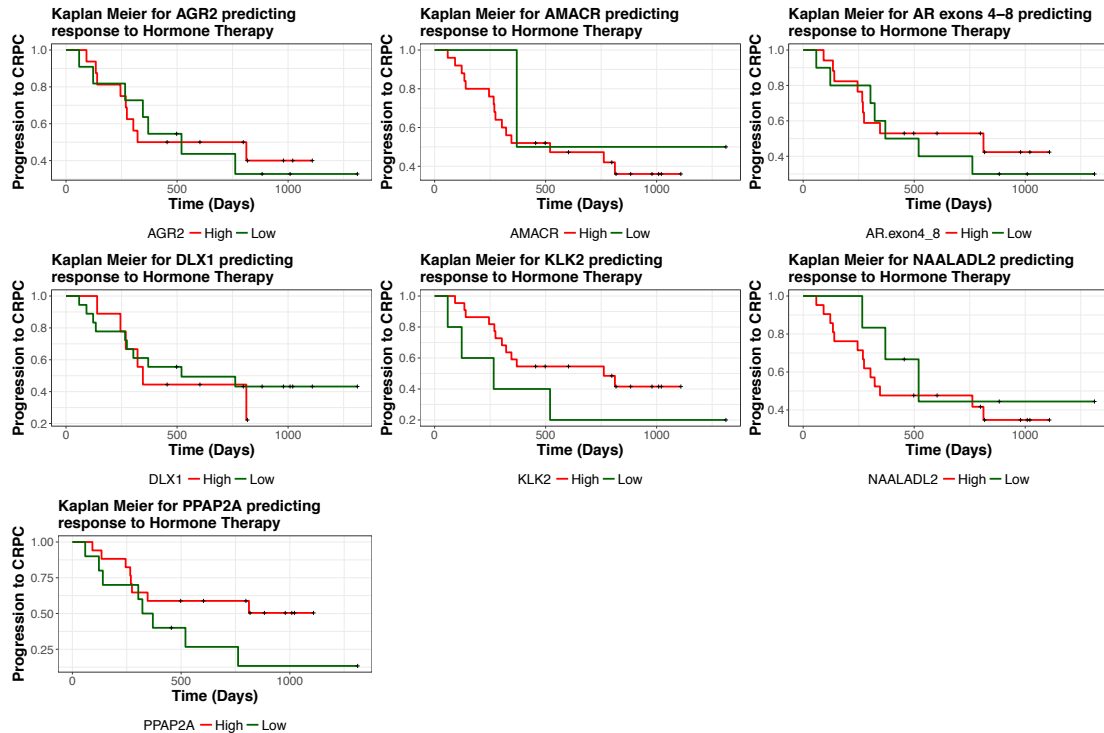
<i>Clinical Variable</i>	<i>Number of patients</i>
<i>Gleason Score</i>	
<i>10</i>	<b>1</b>
<i>9</i>	<b>13</b>
<i>8</i>	<b>3</b>
<i>7 (4+3)</i>	<b>0</b>
<i>No Biopsy: Advanced</i>	<b>10</b>
<i>Bone Scan</i>	
<i>Positive</i>	<b>10</b>
<i>Negative</i>	<b>14</b>
<i>Unknown</i>	<b>3</b>
<i>PSA</i>	<b>Median 63 (7.6 - 9604)</b>
<i>Age</i>	<b>Median 77 (61 - 93)</b>

The seven-transcript signature was not detected as a significant predictor of progression in Nanostring 2 (Cox-regression model;  $p = 0.612$ ,  $HR = 0.640$  (95%  $CI$ : 0.1118 - 3.669). This is confirmed by looking at Kaplan Meier plots of the combined signature (Figure 4.3) and the individual probes (Figure 4.4).

## CHAPTER 4: RESPONSE TO HORMONE THERAPY



**Figure 4.3** Kaplan Meier plot showing the seven transcript signature on NanoString 2 data. The signature was separated using *k*-means.



**Figure 4.4** Individual Kaplan Meier plots for the seven transcripts involved in the signature

To remove any potential batch effect, ComBat was used to normalise the second NanoString data to the pilot study data. Similar results were obtained ( $p = 0.62$ ,  $HR = 0.68$  (95%  $CI$ : 0.15 – 3.18)). Overall Nanostring 1 and 2 are similar.



#### **4.4 Identifying novel progression related transcripts in the NanoString 2 data**

NanoString 2 data contained expression levels from 110 more probes than NanoString 1 data. Therefore, I identified novel progression related transcripts in the NanoString 2 data. Expression of eleven transcripts were identified as significant predictors of progression using Cox regression models ( $p < 0.05$ ), but none were significant after multiple testing correction (Table 4.21). Grouping expression into high and low using  $k$ -means (section 2.5.3), found *MSMB* to be significant even after multiple testing correction ( $p = 0.22 \times 10^{-09}$ , adjusted  $p = 1.54 \times 10^{-06}$ ). Ten other transcripts were significant using this method prior to multiple testing correction ( $p < 0.05$ ; Table 4.22). Log rank test was also performed using the median as a separation cut off for high and low-expression, ten transcripts were observed to be significant prior to multiple testing correction ( $p < 0.05$ ; Table 4.23).

**Table 4.21 Cox regression modelling identified ten probes that were predictors of progression after HT. None were significant after multiple testing correction.**

<i>Probe</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>Coefficient</i>
<i>MSMB</i>	<b>0.0037</b>	<b>0.59</b>	<b>0.7</b>
<i>MIR4435 IHG</i>	<b>0.009</b>	<b>0.98</b>	<b>1.28</b>
<i>BTG</i>	<b>0.017</b>	<b>0.98</b>	<b>1.34</b>
<i>PCSK6</i>	<b>0.021</b>	<b>0.98</b>	<b>0.48</b>
<i>MCTP1</i>	<b>0.028</b>	<b>0.98</b>	<b>1.18</b>
<i>IGFBP3</i>	<b>0.032</b>	<b>0.98</b>	<b>1.2</b>
<i>PCA3</i>	<b>0.036</b>	<b>0.98</b>	<b>0.82</b>
<i>SEC61A1</i>	<b>0.039</b>	<b>0.98</b>	<b>1.19</b>
<i>CLIC2</i>	<b>0.04</b>	<b>0.98</b>	<b>1.65</b>
<i>STOM</i>	<b>0.048</b>	<b>0.98</b>	<b>1.15</b>

**Table 4.22** Log-rank test identified probes that could significantly predict progression on HT. *K*-means was used to separate into high- and low-expression of each probe.

<i>Probe</i>	<i>p-value</i>	<i>Adjusted p-value</i>	$\chi^2$
<i>MSMB</i>	$9.22 \times 10^{-09}$	$1.54 \times 10^{-06}$	33
<i>BTG</i>	0.001	0.2	10.4
<i>CLIC2</i>	0.003	0.5	8.6
<i>MKi67</i>	0.006	0.9	7.7
<i>IGFBP3</i>	0.02	0.99	5.8
<i>PCSK6</i>	0.02	0.99	5.7
<i>APOC1</i>	0.02	0.99	5.5
<i>COL10A1</i>	0.02	0.99	5.4
<i>KLK4</i>	0.02	0.99	5.2
<i>MIC1</i>	0.03	0.99	4.6
<i>SSPO</i>	0.04	0.99	4

**Table 4.23** Log-rank test identified probes that could significantly predict progression on HT. Median was used to separate into high- and low-expression of each probe.

<i>Probe</i>	<i>p-value</i>	<i>Adjusted p-value</i>	$\chi^2$
<i>CLIC2</i>	0.005	0.87	7.8
<i>PCA3</i>	0.008	0.99	7.1
<i>PPAP2A</i>	0.008	0.99	7
<i>SEC61A1</i>	0.012	0.99	6.3
<i>IGFBP3</i>	0.015	0.99	5.9
<i>HIST1H2BG</i>	0.016	0.99	5.8
<i>TBP</i>	0.02	0.99	5.4
<i>PCSK6</i>	0.022	0.99	5.3
<i>BTG2</i>	0.031	0.99	4.6
<i>STOM</i>	0.033	0.99	4.6

There were common transcripts identified in all three different methods: *BTG2*, *CLIC2*, *IGFBP3*, and *PCSK6*. Variable selection using LASSO and stepwise regression identified an optimal model of *BTG2*, *CLIC2* and *PCSK6* (Table 4.24). These three transcripts were also identified as having the greatest importance in a Random Forest model (Table 4.24).

**Table 4.24 Optimising models using the four probes common to log-rank and cox tests. The cox model had an overall  $p$ -value:  $p = 0.0013$ .**

	<i>Lasso Beta value</i>	<i>Cox p-value</i>	<i>Cox HR</i>	<i>Random Forest – Relative importance</i>
<i>BTG2</i>	<b>0.1</b>	<b>0.2</b>	<b>1.17</b>	<b>1</b>
<i>CLIC2</i>	<b>0.43</b>	<b>0.02</b>	<b>1.88</b>	<b>0.49</b>
<i>PCSK6</i>	<b>-0.7</b>	<b>0.003</b>	<b>0.35</b>	<b>0.21</b>
<i>IGFBP3</i>	-	-	-	<b>-0.24</b>

#### **4.5 Hormone Therapy Predictor using *KLK2* ratio data on Nanostring**

##### **1**

For NanoString 2 data I found that refactoring the data using *KLK2* ratio improved the ability to distinguish clinical subtypes (section 5.7.5). Therefore, here I will develop an optimal predictor of progression after HT in the Nanostring 1 data after refactoring using *KLK2*. Differential expression was assessed using the Mann-Whitney U test at three time points: initial non-responders, relapse within 6 months, within 12 months and within 24 months (Table 4.25). Cox regression models (section 2.8.2) identified nine transcripts whose expression were significantly predictors of progression (Table 4.26;  $p < 0.05$ ; multiple testing correction not applied). Log-rank test on expression levels classified as high or low (threshold determined using  $k$ -means, found *AURKA* to be a significant predictor of progression prior to multiple testing correction only ( $p = 0.034$ , Benjamin-Hochberg adjusted  $p = 0.99$ ). Log-rank test when using median for separation into high and low expression, found four transcripts to be differentially expressed between those that relapsed and those that continue to respond to HT (Table 4.27;  $p < 0.05$ ; no multiple testing correction applied; Figure 4.7).

CHAPTER 4: RESPONSE TO HORMONE THERAPY

Table 4.25 Mann-Whitney U test identifies probes differentially expressed between those that have relapsed and those that are still responding to HT at different time periods (initial response relapse, within 6 month relapse, with 12 month relapse and within 24 month relapse.

Initial Response and After 6 months:

<i>Probe</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>Log2(fold change)</i>
<i>KLK3 exons 2-3</i>	<b>0.016</b>	<b>0.92</b>	<b>-0.05</b>
<i>PSGR</i>	<b>0.029</b>	<b>1</b>	<b>-0.09</b>
<i>B2M</i>	<b>0.034</b>	<b>1</b>	<b>-0.06</b>
<i>AURKA</i>	<b>0.047</b>	<b>1</b>	<b>-0.11</b>

After 12 months:

<i>Probe</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>Log2(fold change)</i>
<i>PSGR</i>	<b>0.008</b>	<b>0.47</b>	<b>0.08</b>
<i>FOLH1</i>	<b>0.022</b>	<b>0.98</b>	<b>0.04</b>
<i>KLK3 exons 2-3</i>	<b>0.028</b>	<b>0.98</b>	<b>0.04</b>
<i>B2M</i>	<b>0.038</b>	<b>0.98</b>	<b>0.06</b>

After 24 months:

<i>Probe</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>Log2(fold change)</i>
<i>PECI</i>	<b>0.01</b>	<b>0.56</b>	<b>0.03</b>
<i>PSGR</i>	<b>0.016</b>	<b>0.88</b>	<b>0.06</b>
<i>DLX1</i>	<b>0.018</b>	<b>0.97</b>	<b>-0.05</b>
<i>ALAS1</i>	<b>0.045</b>	<b>1</b>	<b>0.05</b>

Table 4.26 Cox identified probes that are differentially expressed in NanoString 1 data normalised by KLK2 ratio.

<i>Probe</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>Coefficient</i>
<i>CAMKK2</i>	<b>0.015</b>	<b>0.86</b>	<b>2.22</b>
<i>UPK2</i>	<b>0.027</b>	<b>0.98</b>	<b>1.49</b>
<i>KLK3 exons 2-3</i>	<b>0.031</b>	<b>0.98</b>	<b>3.15</b>
<i>PECI</i>	<b>0.031</b>	<b>0.98</b>	<b>2.27</b>
<i>HPN</i>	<b>0.031</b>	<b>0.98</b>	<b>2.48</b>
<i>KLK4</i>	<b>0.034</b>	<b>0.98</b>	<b>3.6</b>
<i>GAPDH</i>	<b>0.036</b>	<b>0.98</b>	<b>2.4</b>
<i>ALAS1</i>	<b>0.038</b>	<b>0.98</b>	<b>2.15</b>
<i>KLK3 exons 1-2</i>	<b>0.048</b>	<b>0.98</b>	<b>2.39</b>

Table 4.27 Log rank (using median for separating high and low expression) identified probes that differ between response to HT.

<i>Probe</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>X<sup>2</sup></i>
<i>STEAP4</i>	<b>0.007</b>	<b>0.32</b>	<b>7.7</b>
<i>PECI</i>	<b>0.009</b>	<b>0.49</b>	<b>6.8</b>
<i>SERPINB5</i>	<b>0.013</b>	<b>0.72</b>	<b>6.1</b>
<i>TBP</i>	<b>0.037</b>	<b>0.99</b>	<b>4.4</b>

## CHAPTER 4: RESPONSE TO HORMONE THERAPY

*PECI* was identified in all comparisons; Mann-Whitney U, Cox and Log-rank (using median for separation). Variable selection on all fourteen transcripts that were identified as candidate predictors (Table 4.25, Table 4.26, Table 4.27) were performed. Lasso identified three transcripts: *CAMKK2*, *DLX1* and *UPK2* (Table 4.28); stepwise regression identified six transcripts of which only *UPK2* was common to Lasso (Table 4.29); and using Random Forest three of the top five important genes were found in either the Lasso or Stepwise regression results (Table 4.30).

**Table 4.28 Lasso selects three transcripts for HT progression prediction in *KLK2* adjusted data.**

<i>Transcript</i>	<i>Beta coefficient</i>
<i>CAMKK2</i>	<b>0.232</b>
<i>DLX1</i>	<b>-0.028</b>
<i>UPK2</i>	<b>0.099</b>
<i>Cox model: p – value = 0.2, HR = 0.999, 95% CI = 0.998 - 1</i>	

**Table 4.29 Stepwise regression selects six probes for early HT relapse prediction in *KLK2* adjusted data.**

<i>Transcript</i>	<i>p-value</i>	<i>HR</i>
<i>B2M</i>	<b>0.096</b>	<b>2.713</b>
<i>FOLH1</i>	<b>0.042</b>	<b>0.246</b>
<i>GAPDH</i>	<b>0.122</b>	<b>0.102</b>
<i>HPN</i>	<b>0.083</b>	<b>5.504</b>
<i>PSGR</i>	<b>0.027</b>	<b>2.773</b>
<i>UPK2</i>	<b>0.039</b>	<b>1.955</b>
<i>Cox model: p – value = 0.039, HR = 86.54, 95% CI = 2.435 - 3076</i>		

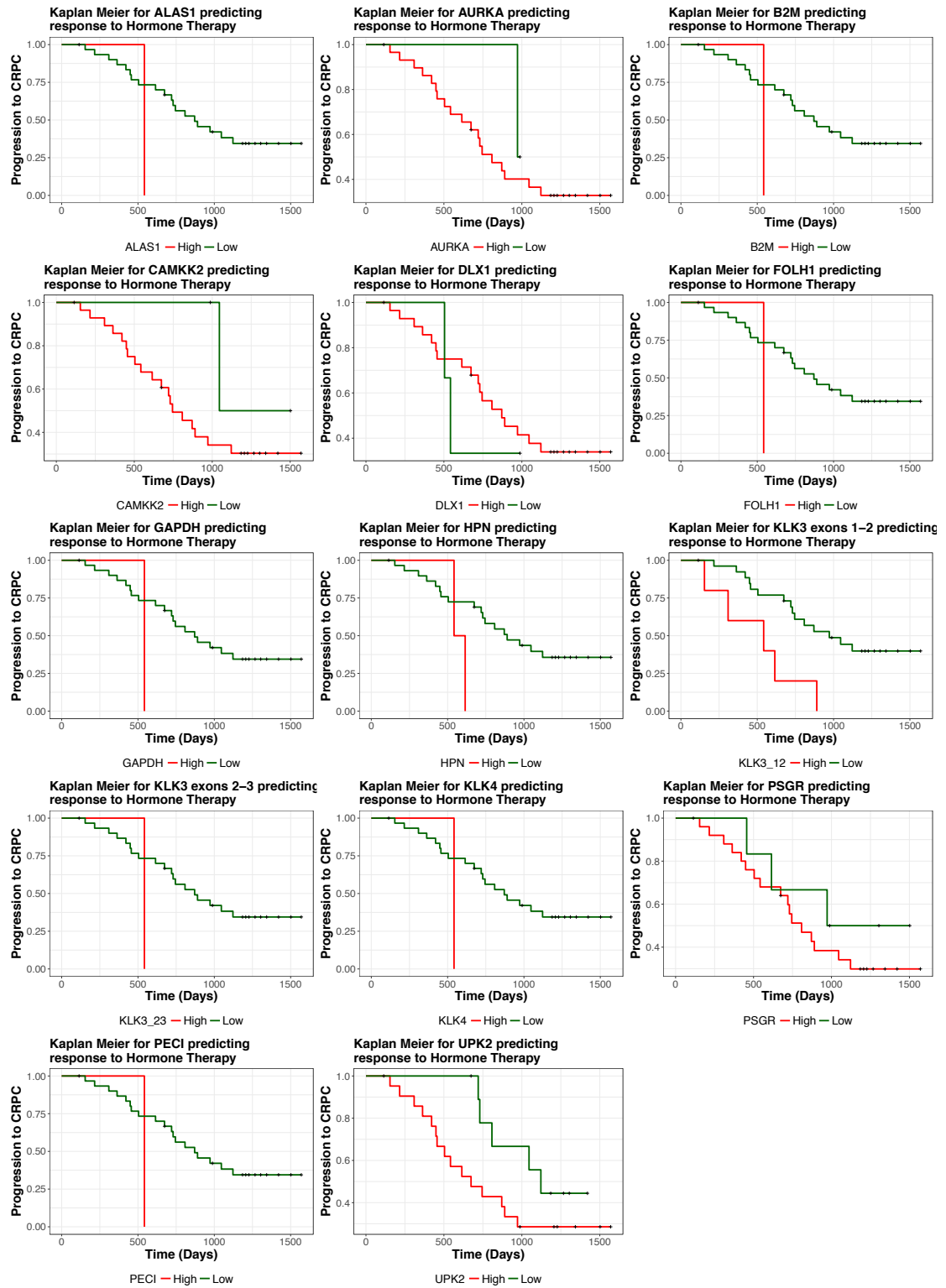
## CHAPTER 4: RESPONSE TO HORMONE THERAPY

**Table 4.30** Random forest shows the importance of each transcript in distinguishing early HT relapse in *KLK2* adjusted data.

<i>Transcript</i>	<i>Importance</i>	<i>Relative Importance</i>
<i>PSGR</i>	<b>0.035</b>	<b>1.00</b>
<i>PECI</i>	<b>0.035</b>	<b>0.99</b>
<i>HPN</i>	<b>0.007</b>	<b>0.20</b>
<i>KLK3 exons 1-2</i>	<b>0.006</b>	<b>0.18</b>
<i>CAMKK2</i>	<b>0.006</b>	<b>0.16</b>
<i>ALAS1</i>	<b>0.006</b>	<b>0.16</b>
<i>FOLH1</i>	<b>0.005</b>	<b>0.13</b>
<i>DLX1</i>	<b>0.003</b>	<b>0.09</b>
<i>AURKA</i>	<b>0.0002</b>	<b>0.01</b>
<i>KLK3 exons 2-3</i>	<b>-0.0002</b>	<b>-0.01</b>
<i>UPK2</i>	<b>-0.0007</b>	<b>-0.02</b>
<i>GAPDH</i>	<b>-0.001</b>	<b>-0.03</b>
<i>B2M</i>	<b>-0.006</b>	<b>-0.02</b>
<i>KLK4</i>	<b>-0.007</b>	<b>-0.19</b>

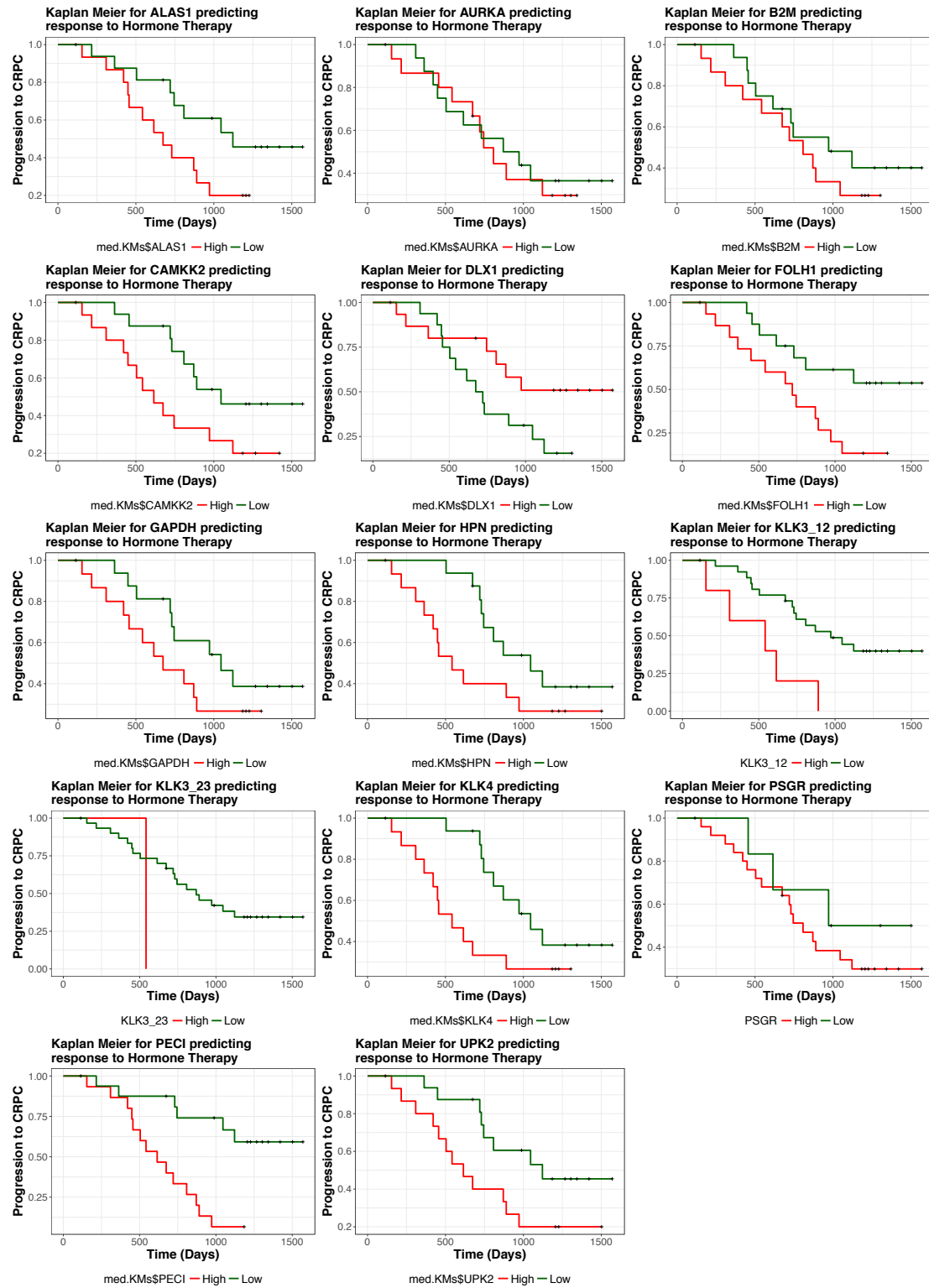
Kaplan Meier plots (section 2.8.1) were produced using a *k*-means determined threshold between high and low expression (Figure 4.5). Applying variable selection on this dichotomised data, Lasso identifies *CAMKK2*, *PSGR* and *UPK2* (Table 4.31); stepwise regression selects *AURKA*, *CAMKK2*, *KLK3* exons 1-2 and *UPK2* (Table 4.32); whilst random forest suggests that *PSGR* is of most importance, followed by *UPK2* and *CAMKK2* (the same three transcripts selected via Lasso) (Table 4.33). *CAMKK2*, *PSGR* and *UPK2*, which were selected by Lasso and also the three most important transcripts according to Random forest, produce a significant cox model ( $p = 0.007$ ,  $HR = 1.0028$ , 95% CI = 1.001 - 1.005).

## CHAPTER 4: RESPONSE TO HORMONE THERAPY



**Figure 4.5** Kaplan Meier plots (expression separated via *k*-means) for the fourteen transcripts identified via Mann Whitney U, Cox and log-rank tests for early HT relapse.

## CHAPTER 4: RESPONSE TO HORMONE THERAPY



**Figure 4.6** Kaplan Meier plots (expression separated via median) for the fourteen transcripts identified via Mann Whitney U, Cox and log-rank tests for early HT relapse.



## CHAPTER 4: RESPONSE TO HORMONE THERAPY

**Table 4.31** Lasso (with glm) selects three transcripts from the five shown to be differential from Kaplan Meier plots using *k*-means for separation. An overall Cox model using these three probes proves to be significant ( $p = 0.007$ ).

<i>Transcript</i>	<i>Beta coefficient</i>
<i>CAMKK2</i>	<b>0.31</b>
<i>PSGR</i>	<b>0.06</b>
<i>UPK2</i>	<b>0.18</b>
<i>Cox model: p – value = 0.007, HR = 1.0028, 95% CI = 1.001 - 1.005</i>	

**Table 4.32** Step (with Cox) selects four transcripts from the five shown to be differential from Kaplan Meier plots using *k*-means for separation. An overall Cox model using these four probes is not significant ( $p = 0.07$ ).

<i>Transcript</i>	<i>p-value</i>	<i>HR</i>
<i>AURKA</i>	<b>0.13</b>	<b>1.25</b>
<i>CAMKK2</i>	<b>0.04</b>	<b>3.51</b>
<i>KLK3 exons 1-2</i>	<b>0.14</b>	<b>0.13</b>
<i>UPK2</i>	<b>0.06</b>	<b>2.10</b>
<i>Cox model: p – value = 0.07, HR = 1, 95% CI = 1-1</i>		

**Table 4.33** Random forest shows the importance of each of the five transcripts identified via Kaplan Meier plots using *k*-means for separation. The top three important transcripts are identical to the Lasso output.

<i>Transcript</i>	<i>Importance</i>	<i>Relative Importance</i>
<i>PSGR</i>	<b>0.098</b>	<b>1</b>
<i>UPK2</i>	<b>0.018</b>	<b>0.179</b>
<i>CAMKK2</i>	<b>0.014</b>	<b>0.137</b>
<i>AURKA</i>	<b>0.003</b>	<b>0.031</b>
<i>KLK3 exons 1-2</i>	<b>-0.002</b>	<b>-0.017</b>

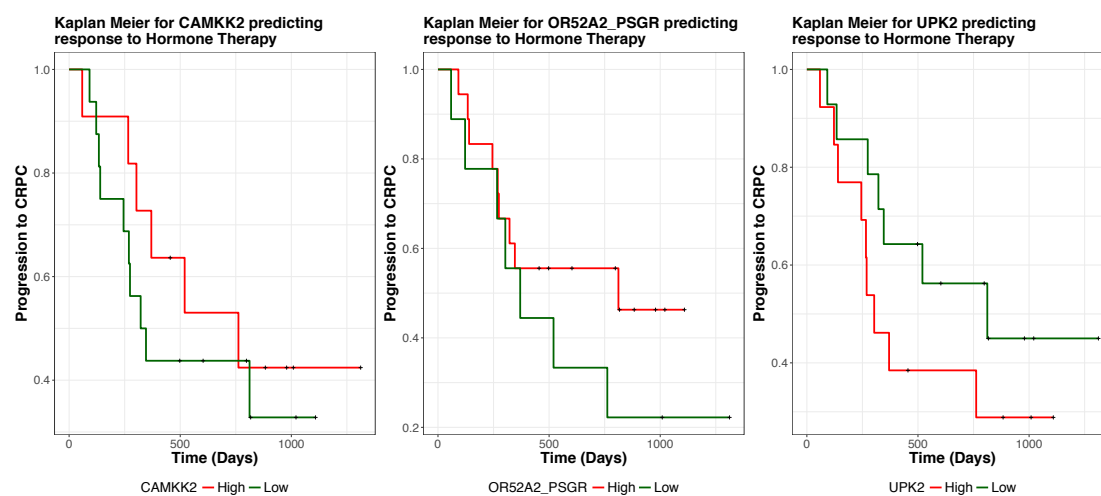
Lasso has consistently selected *CAMKK2* and *UPK2* along with one other transcript (*DLX1* when all that showed significance were used, *PSGR* when only those that appeared to be significant in *k*-means separated Kaplan Meier plots). This consistency of selecting *CAMKK2* and *UPK2* when the input variables are altered shows reproducibility. Though the model including *DLX1* was not significant, the model including *PSGR* was the most significant model identified ( $p = 0.0023$ ,  $HR = 1.0028$ ,  $95\% CI = 1.001 - 1.005$ ). *CAMKK2* was always identified as important by Random forest. Step (with Cox) and Random Forest was not very consistent in creating models

## CHAPTER 4: RESPONSE TO HORMONE THERAPY

with similar variables. Therefore, the most consistent and significant cox regression model identified contained *CAMKK2*, *PSGR* and *UPK2* ( $p = 0.007$ ,  $HR = 1.0028$ , 95%  $CI = 1.001 - 1.005$ ).

### 4.5.1 Validation of the final model on *KLK2* ratio NanoString 2 data

The second set of NanoString data also refactorised using the *KLK2* ratio method was used to test the *CAMKK2*, *PSGR* and *UPK2* Cox regression model identified in NanoString 1 data. The model did not reach statistical significance as a predictor of progression ( $p = 0.4$ ,  $HR = 1.000774$  (95%  $CI: 0.999 - 1.003$ ). Looking at the Kaplan Meier plots of the transcripts individually (Figure 4.7), *CAMKK2*, *PSGR* and *UPK2* showed better survival with low expression in the pilot study, yet in the second set of data, both *CAMKK2* and *PSGR* now show better survival with higher expression.



**Figure 4.7** Kaplan Meier plots for the three transcripts in the model for *KLK2* adjusted hormone therapy data: *CAMKK2*, *PSGR* and *UPK2*.

## 4.6 Conclusion

Stratified medicine enables the optimal treatment for cancer patients to be selected and improve overall survival. There have been successes in breast<sup>208</sup> and lung cancer<sup>209,210</sup>, for example. However, no such robust testing and stratification exists for men with PCa and the route of treatment is not always clear. In particular, it is hard to predict the

## CHAPTER 4: RESPONSE TO HORMONE THERAPY

response to treatments such as radiotherapy, hormone therapy and prostatectomy and determine whether active surveillance is a better option than treatment. These issues are key research areas for the clinical management of PCa patients.

In this chapter, I investigated if expression profiles of urinary microvesicles could be used to estimate how long a patient responded to HT. I successfully built a number of different models based on the normalisation method and dataset used. To produce a non-invasive test for the identification of those who will relapse early on HT, and thus could benefit from additional treatment, would be ground breaking for PCa patients. Depending on how the NanoString data was normalised, we saw correlation between two signatures and HT relapse.

Under the Nanostring 1 dataset with standard normalisation (via NanoString's positive probes, section 2.3.1), the optimal predictor of progression in HT patients included the expression of probes *AGR2*, *DLX1*, *KLK2*, *NAALADL2*, *AR* exons 4-8, *PPAP2A* and *AMACR* (Cox-regression model,  $p = 2.3 \times 10^{-05}$ ,  $HR = 0.043$ ). This seven-probe signature was also a significant independent predictor of progression improving on other clinical factors, initial PSA, Gleason score, initial bone scan results and age (Table 4.18, Table 4.19).

After *KLK2* adjustment of the NanoString 1 data (section 2.1.1), we selected model including *CAMKK2*, *PSGR* and *UPK2* that could significantly separate those that progressed to CPRC and those that continued to respond to HT. This model was significant alone ( $p = 0.0023$ ,  $HR = 1.0028$ , 95% CI = 1.001 - 1.005), which was again more significant than other clinical factors including initial PSA, Gleason score, initial bone scan results and age (Table 4.18).

Unfortunately, both of these models were not validated in the NanoString 2 dataset. There are many possible reasons why the models were not validated but one factor is that there are a relatively small number of patients in this cohort and with a relatively short feedback. This means that the models are very sensitive to outliers in the data. There are also differences between the Nanostring datasets: different centres ran the

## CHAPTER 4: RESPONSE TO HORMONE THERAPY

experiment, there were different probesets, newer samples could have been collected slightly differently, and the cohorts could be somewhat different. It appears the data is very sensitive with no candidate probes being common before and after factorisation. This is common to many expression-based biomarker studies. There is a lack of robustness with proposed tests very rarely being validated in different cohorts<sup>211</sup>. It should also be remembered that this is a targeted based assay, the optimal probes to distinguish treatment outcome may not be included.

In probes that were unique to the second NanoString data set, the optimal model for determining time to progress for HT patients contained *BTG2*, *CLIC2* and *PCSK6*.

In this chapter I have shown the utility of urine derived microvesicle expression profiles for the prediction of outcome after treatment. This is a proof of concept that would require a much larger series with longer feedback to find the best combination of transcripts and become a usable test.

# 5

## NanoString Data

### Analysis 2

#### 5.1 Summary

The Movember GAP1 Urine Biomarker Consortium had multiple collaborators working on the identification of urinary biomarkers for the risk-stratification of PCa. Our laboratory is specifically interested in the RNA expression changes in PCa that are detectable within urinary cell sediments and extracellular vesicles (EVs) from samples collected at multiple centres. The aims of my study were to see if I could identify robust models of expression profiles using data obtained from NanoString that could answer important clinical questions in PCa management: can I detect PCa from non-PCa samples and can I risk stratify PCa, both without the need for biopsy. I therefore,

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

investigated different methods for normalisation of this urinary EV derived data with the aim to build optimal models from the expression of 167 markers for risk stratification and detection of cancer.

I identified robust models for the detection of PCa from non-PCa samples (AUC = 0.851) and of high-risk PCa from non-PCa samples (AUC = 0.897). Models to predict risk stratification between samples with no evidence of cancer (CB) and cancer in order of severity (CB->L->I->H) were also produced (AUC = 0.709). My models used many of the already published transcripts used in whole urine assays but also included novel transcripts that may be present in EV fractions.

### **5.2 Introduction**

NanoString expression analysis of 167 gene-probes was applied to cell and extracellular vesicle (EV) fractions of urine from prostate cancer patients to form the NanoString 2 data set. In this chapter, quality control and technical trouble shooting (section 2.3) is applied to the whole data set, before performing exploratory analysis using just the EV samples. Investigation of the cell fraction samples can be found in chapter 6.

#### **5.2.1 The Research Gap**

Risk stratification is currently based on PSA, Gleason score and T stage but has the potential to be improved by using a novel biomarker panel. This could help tailor patients to treatment pathways and determine, at diagnosis, the aggressiveness of disease. The PCA3 test is an established biomarker that is capable of predicting PCa on a second biopsy. Therefore, showing the utility of the use of urine in PCa diagnostics and prognostics, and has shown some minor improvements to risk stratification. In chapter 3, I performed a pilot project exploring the use of NanoString applied to genetic material obtained from urinary EVs and showed that it was capable of capturing clinically relevant expression profiles.

### **5.2.2 Aims**

In this chapter I used NanoString technology to investigate the RNA expression level changes of 167 target sequences within EVs extracted from urinary samples collected at multiple centres world-wide as part of the Movember study. The aims of this study are:

1. To identify better processing techniques for the EV NanoString data
2. To determine whether EV expression profiles are robust across variable sample cohorts collected from different centres.
3. To identify optimal models built from the expression of 167 markers for risk stratification and detection of cancer.

### **5.2.3 The Probe Targets**

A panel of experts selected the 167 sequence targets used as probes. The probes were primarily selected from publications that highlighting genes overexpressed in prostate tumour tissue. 28 gene probes were selected from Next Generation Sequencing data of 20 urine EV RNA samples from the NNUH. Additionally, some prostate tissue specific controls and controls for kidney, bladder and blood were also included. See Supplementary Table 1 for further details.

### **5.2.4 Classification of prostate cancer patient samples**

NanoString data from 864 samples was obtained, 95 samples were from the cell fraction. 756 samples remained after quality control (Section 5.3.2). Samples were divided in to a training set and a test set based on a 2:1 ratio while maintaining the proportions of each PCa risk category (Table 5.1) and sample collection centre (Table 5.2). The median age and PSA at diagnosis have been recorded for each clinical category within the training and tests, respectively (Table 5.2.3).

**Table 5.1 Classification and Frequency of the sample types based on NICE criteria<sup>40</sup>.** The quantity of samples for each clinical group are provided as well as the clinical description of the group in terms of Gleason score, PSA level and T stage.

<i>Classification: NICE Groupings</i>				
<i>Sample Class</i>	<i>Description</i>	<i>Number of Samples</i>	<i>Number of Training Samples</i>	<i>Number of Test Samples</i>
<i>Advanced</i>	<i>Advanced and Hhh (G8-10 PSA&gt;100) and Hh (G8-10 PSA&lt;100)</i>	<i>31</i>	<i>21</i>	<i>10</i>
<i>High-risk</i>	<i>HL= G7 PSA&gt;20</i>	<i>107</i>	<i>72</i>	<i>35</i>
<i>Intermediate-risk</i>	<i>I= G3+4 PSA&lt;20 and IL= G6 PSA&gt;10</i>	<i>214</i>	<i>142</i>	<i>72</i>
<i>Low-risk</i>	<i>L= Low G6 PSA&lt;10</i>	<i>156</i>	<i>104</i>	<i>52</i>
<i>Abnormal</i>	<i>High PSA no Bx, Prostatitis, Raised PSA negative Bx, HGPIN</i>	<i>137</i>	<i>92</i>	<i>45</i>
<i>CB</i>	<i>CB – no evidence of cancer</i>	<i>111</i>	<i>73</i>	<i>38</i>
<i>Total</i>		<i>756</i>	<i>504</i>	<i>252</i>

**Table 5.2 Sample collection-site breakdown of the EV samples from NanoString2.**

<i>Location</i>	<i>Training Set</i>	<i>Test Set</i>	<i>Number of Samples</i>
<i>Dublin</i>	<i>16</i>	<i>8</i>	<i>27</i>
<i>ICR</i>	<i>84</i>	<i>41</i>	<i>130</i>
<i>UEA</i>	<i>323</i>	<i>163</i>	<i>496</i>
<i>USA</i>	<i>62</i>	<i>23</i>	<i>103</i>
<i>Total</i>			<i>756</i>

**Table 5.3 Median age and PSA of each clinical category within the training and test datasets.**

	<b>Training Set</b>		<b>Test Set</b>	
	Median age	Median PSA	Median age	Median PSA
<b>Advanced</b>	78	273.5	82	285.75
<b>High-risk</b>	69	22.35	73.5	23.7
<b>Intermediate-risk</b>	69	9.2	67	8.45
<b>Low-risk</b>	64.5	6.1	64	5.5
<b>Abnormal</b>	67	8.19	66	7.7
<b>CB</b>	63	1.4	64.5	1.235

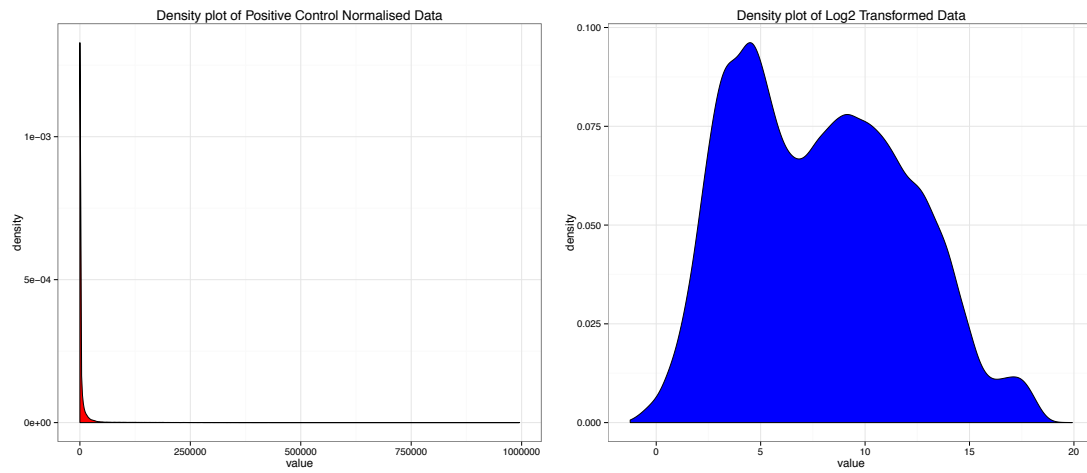


### **5.3 Data Preprocessing and Technical Variation**

#### **5.3.1 Normalisation and Background correction**

There were six positive-control non-human ERCC probes included in the NanoString series and these were used to normalise the data for all samples as per the NanoString manual. As for the pilot data set, a large proportion (33%) of data points were less than zero after negative control correction. Therefore negative control correction was not used in this analysis. As shown in NanoString 1 (section 3.3.4)  $\text{Log}_2$  transformation (section 2.3.3) was used to obtain a more normal distribution in the data (Figure 5.1). The  $\text{Log}_2$  data did not follow a normal distribution using the Shapiro-Wilk test (Table 5.4), this suggests we should use non-parametric methods for analysis.

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2



**Figure 5.1 A) Positive control normalised data. B) Positive control normalised and  $\text{Log}_2$  transformed data. The data shows a more normal distribution after  $\text{Log}_2$  transformation.**

**Table 5.4** Shapiro-Wilk tests show that  $\text{Log}_2$  data is not normally distributed.

	<i>Log<sub>2</sub> transformed</i>		
	<i>W</i>	<i>p-value</i>	<i>Normally Distributed</i>
<i>The first set of randomly selected 29 samples, all probes</i>	<b>0.96</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The second set of randomly selected 29 samples, all probes</i>	<b>0.97</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The third set of randomly selected 29 samples, all probes</i>	<b>0.98</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The fourth set of randomly selected 29 samples, all probes</i>	<b>0.97</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The fifth set of randomly selected 29 samples, all probes</i>	<b>0.98</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The sixth set of randomly selected 29 samples, all probes</i>	<b>0.97</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The first set of randomly selected probes, all samples</i>	<b>0.99</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The second set of randomly selected probes, all samples</i>	<b>0.99</b>	<b>2.205x10<sup>-15</sup></b>	<b>No</b>
<i>The third set of randomly selected probes, all samples</i>	<b>0.94</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The fourth set of randomly selected probes, all samples</i>	<b>0.96</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The fifth set of randomly selected probes, all samples</i>	<b>0.94</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>
<i>The sixth set of randomly selected probes, all samples</i>	<b>0.94</b>	<b>&lt; 2.2x10<sup>-16</sup></b>	<b>No</b>

### 5.3.2 Quality of Normalisation

The quality of the data, and its normalisation and transformation, was assessed using NanoStringNorm (section 2.3.2.1) and NanoStringQCPro (section 2.3.2.2). Overall the quality was good but a few samples and a few probes need to be treated with caution. The samples identified by the IQR/median plot were removed (A210, A216, A517, C147\_1, M\_97\_5, M\_138\_7, M\_149\_7) along with some CBN samples, which were identified through NGS analysis (not shown as this was not performed by me).

#### 5.3.2.1 NanoStringNorm

The negative controls had both low means and standard deviations and the positive controls showed low standard deviation, as expected. The majority of the probes

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

clustered around the loess curve of best fit (96 %) but a few probes were highlighted due to high means and standard deviation: *KLK4*, *RPS10*, *RPLP2*, *M5MB*, and *RPS11*. Whilst *AR* exons 4-8 and *ITPRI* were highlighted due to low mean and standard deviation.

If a sample has many missing values this could be caused by a technical failure or as a result of too little input material. There were a few samples that seemed to have missing values in the normalised data (A216, A210, A196, M\_138\_7, M\_140\_6, M\_147\_3, M\_92\_5, M\_97\_5 and C147\_1). These were watched carefully throughout further analyses.

Each NanoString cartridge holds twelve samples. NanoStringNorm uses a *t*-test to identify cartridges that have a significantly different means, standard deviation and levels of positive controls in comparison to the other cartridges. Cartridges 22, 23, 58, 59, 60, 61, 62, 63 and 64 had higher means and standard deviation, whilst cartridges 15, 29, 36, 37, 43 and 65 through 72 had lower detection levels of positive controls.

Looking into the normalisation factors using NanoStringNorm, a number of samples had normalisation parameters that extended beyond 100% difference from the mean and could be influential outliers: (Supplementary Table 2).

### 5.3.2.2 *NanoStringQCPro*

NanoStringQCPro provided information on the binding density, field of view (FOV) and the positive controls used for initial normalisation. NanoString is only capable of reading un-overlapped barcodes when digitally scanning the image produced. Twenty-eight samples were identified as having overlapping barcodes typically caused by excess RNA input (Supplementary Table 2). No samples were identified as having less than 80% FOV, meaning there were no technical issues due to loading of cartridges (e.g. bubbles, or insufficient oiling).

The slope in the positive control data shows how well an increase in input is reflected by an increase in counts, measured using a linear model ( $\log(\text{counts}) \sim \log(\text{input})$ ).

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

Three samples were highlighted as outliers from the model: M\_122\_2, M\_127\_6 and M\_131\_4. Two of these samples also showed high IQR of positive controls: M\_127\_6 and M\_131\_4. NanoString recommends a positive scaling factor between 0.3 and 3. A scaling factor above this range indicates low performance of that lane during the NanoString counting protocol. Six samples' lanes were flagged as such (M\_95\_1, M\_97\_6, M\_140\_6, M\_144\_1, M\_75\_3 and M\_147\_3) and thus were considered with caution.

### **5.3.3 Experimental and Technical Investigations**

#### ***5.3.3.1 Sample and Centre Investigations***

Comparing the median with the IQR can unveil samples with low medians and/or IQRs (both of which can be problematic). Some samples were identified as such: A210, A216, A517, C147\_1, M\_97\_5, M\_138\_7, M\_149\_7 (Figure 5.2). These samples were removed from the analysis. PCA identified a clear clustering of the cell sediment derived samples compared to EV derived samples from multiple centres (Figure 5.3), further highlighting their need to be analysed separately (Chapter 6). PCA on EV derived samples showed some clustering based on location of origin (Figure 5.4). There is evidence of significant differences in overall expression between some origin centres (Mann Whitney U tests;  $p < 0.05$ ; Table 5.5). However, the average  $\text{Log}_2$  expression appears to be fairly uniform across the centres (Figure 5.5).

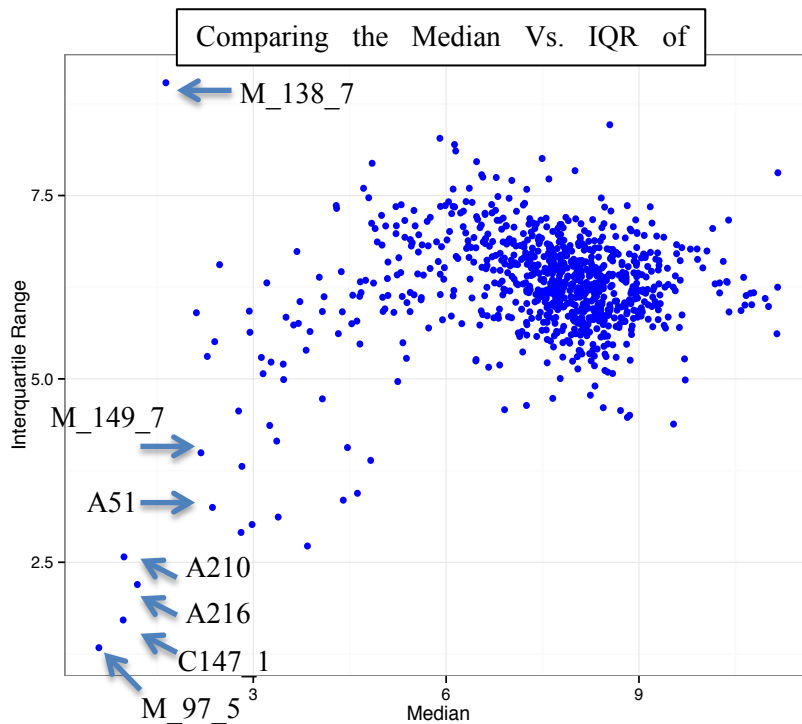
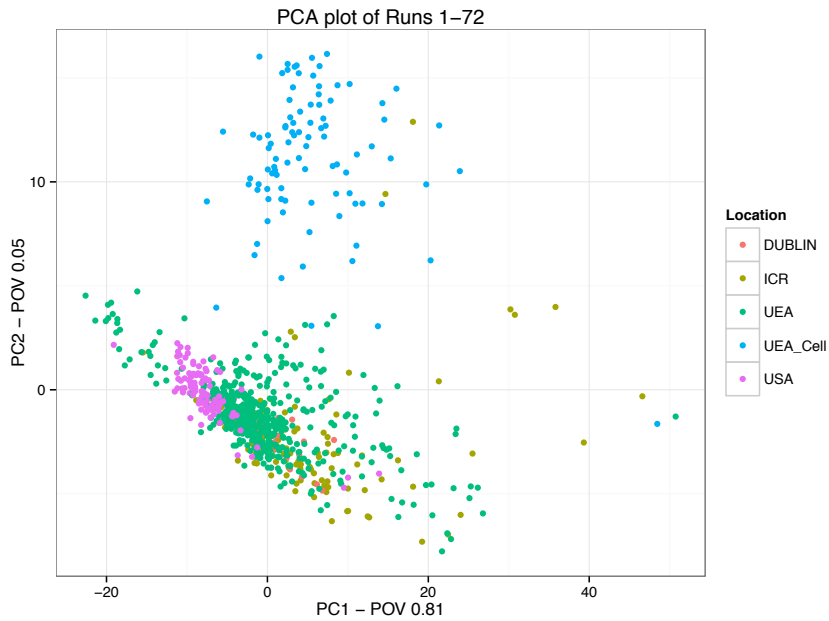


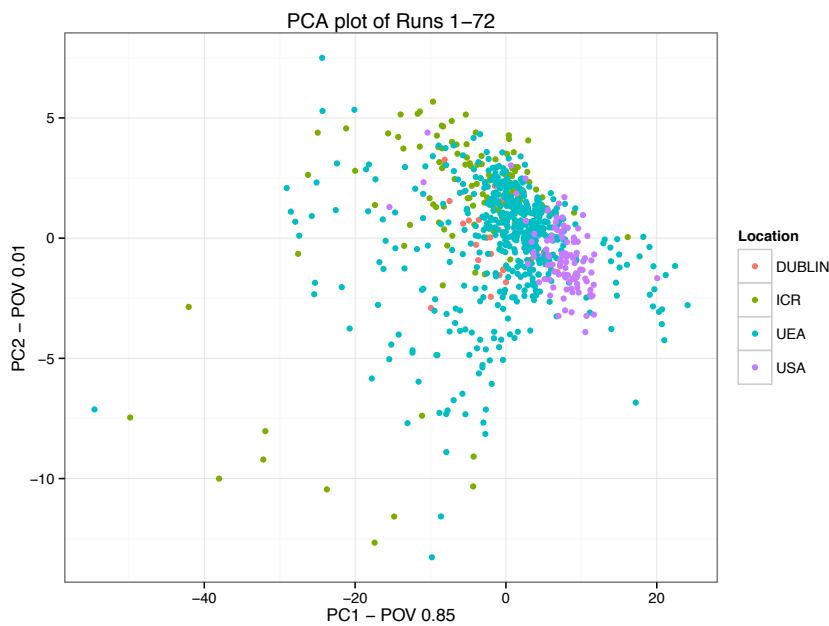
Figure 5.2 Median Vs. IQR of samples on the second NanoString study. Six samples were identified with low medians and/or IQRs, which could be problematic to further analyses.

### 5.3.3.2 NanoString Cartridges

NanoStringNorm showed significant differences between the mean and standard deviation of the normalised data between some cartridges; indicating there might be batch effects. Cartridge dependent variations were therefore examined using boxplots (Figure 5.6) and there was significant association between mean expression per sample and cartridge (Kruskal-Wallis rank sum test:  $p < 2.2 \times 10^{-16}$ ,  $\chi = 329.25$ ). As samples from the same collection centres were loaded consecutively, there was no surprise that there was a significant association between centre and cartridge also (Chi-square test;  $p$ -value  $< 2.2 \times 10^{-16}$ ,  $\chi = 2036.5$ ). As location was also significantly associated with median expression of samples, it was not a leap to believe this issue with cartridge discrepancies was due to location.



**Figure 5.3** DNA extracted from EVs was collected from four different centres (Dublin, ICR, UEA, and the USA). DNA extracted from the cell pellet was only collected at UEA (UEA\_Cell). PCA plot clearly identifies cell sediment derived samples as a separate cluster from EV derived samples.



**Figure 5.4** PCA plot of only EV derived DNA shows evidence of collection-centre of origin based clustering.

**Table 5.5** Expression values from different collection-centres of origin compared by Mann Whitney U tests show that all centres are significantly different.

	<i>USA</i>	<i>ICR</i>	<i>DUBLIN</i>
<i>UEA</i>	$< 2.2 \times 10^{-16}$	$< 2.2 \times 10^{-16}$	$2.311 \times 10^{-07}$
<i>USA</i>	-	$< 2.2 \times 10^{-16}$	$< 2.2 \times 10^{-16}$

ICR                    -                    -                     $1.704 \times 10^{-11}$

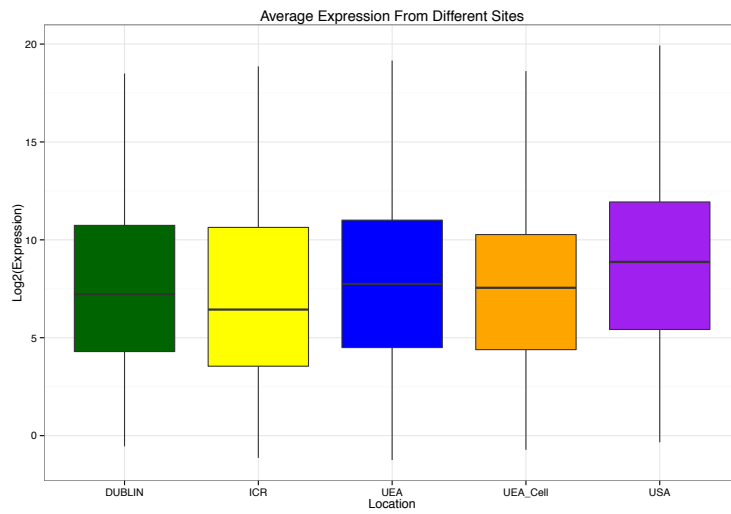


Figure 5.5 Average Log<sub>2</sub> expression across centres shows similar expression levels.

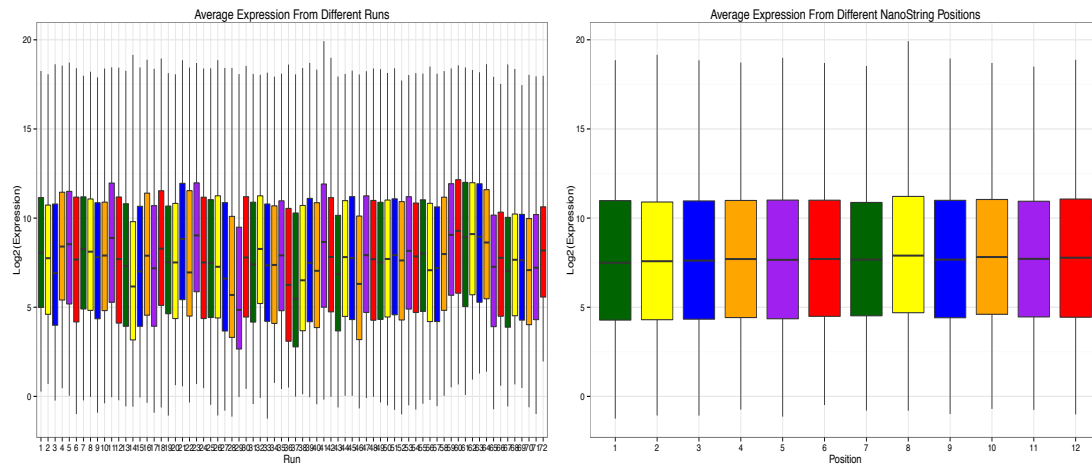
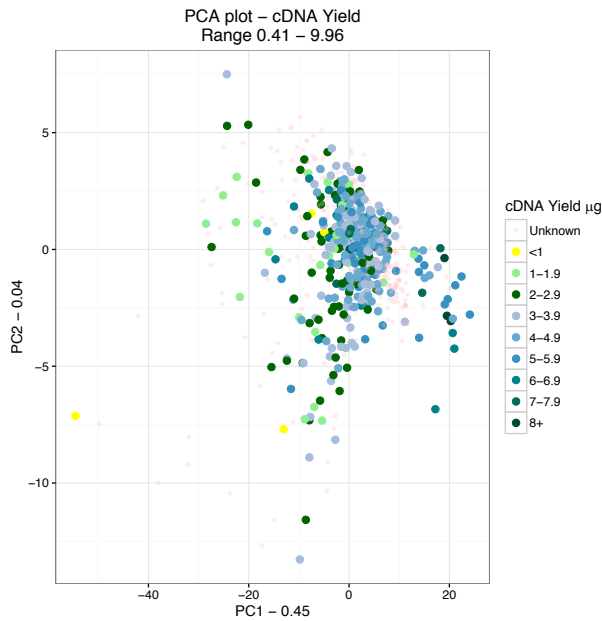


Figure 5.6 Boxplots showing average expression across cartridge and position on cartridge are similar and are showing no batch effects.





**Figure 5.7** PCA plot of EV derived samples, showing a lack of clustering by cDNA yield.

### 5.3.3.3 RNA Amplification to cDNA

As 100 ng of RNA or cDNA is required for NanoString analysis, and the amounts of EV RNA harvestable from urine were limiting in a large proportion of samples, 15-20ng RNA from each sample was amplified using a Nugen Ovation WTA2 cDNA amplification kit. The amount of cDNA obtained after amplification (in µg) was investigated for clustering affects using PCA (Figure 5.7). cDNA yields were split into groups; <1 µg, 1-1.9µg, 2-2.9µg, 3-3.9µg, 4-4.9µg, 5-5.9µg, 6-6.9µg, 7-7.9µg and >8µg. Mild clustering affects were observed, and a significant correlation was found between cDNA yield and median  $\log_2$  expression per sample ( $p < 2.2 \times 10^{-16}$ ,  $r = 0.44$ , Pearson's correlation). The distribution of clinical categories within each amplification yield group was not statistically significant; ( $\chi = 125.3$ ,  $p > 0.05$ ,  $\chi^2$  test (section 2.4.4)).

### 5.3.4 ComBat – Removing collection-centre based significance

Batch effects caused by location of sample origin (centre) were accounted for by using the ComBat function of the *sva* package. PCA was then used to visualise clustering in the post-ComBat data (Figure 5.9). There was no significant difference between median  $\log_2$  expression across location (Kruskal-Wallis rank sum test:  $p = 0.6488$ ,  $\chi = 1.647$ ).

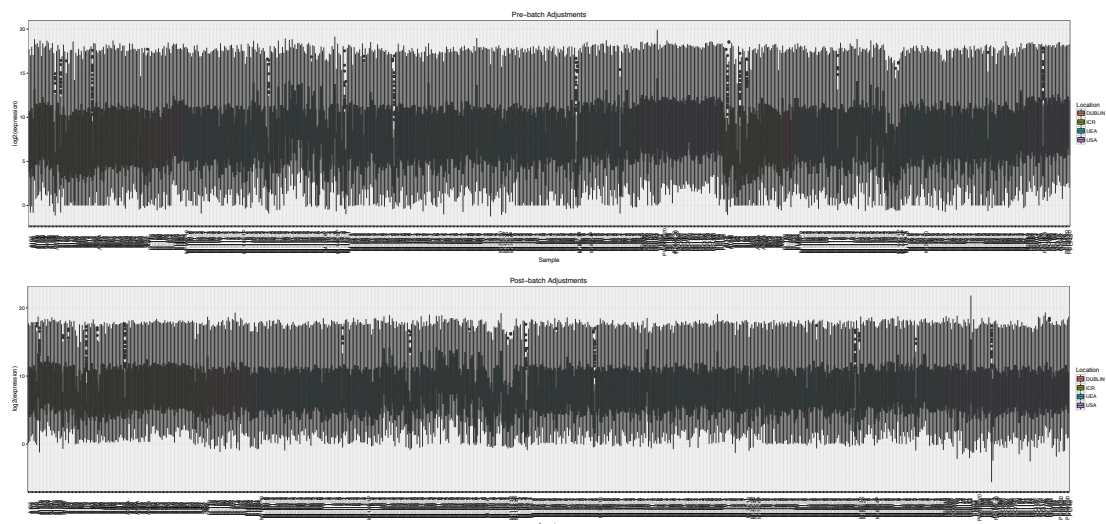


Figure 5.8 Boxplots show the log<sub>2</sub> expression across each sample, coloured by location before and after the application of ComBat.

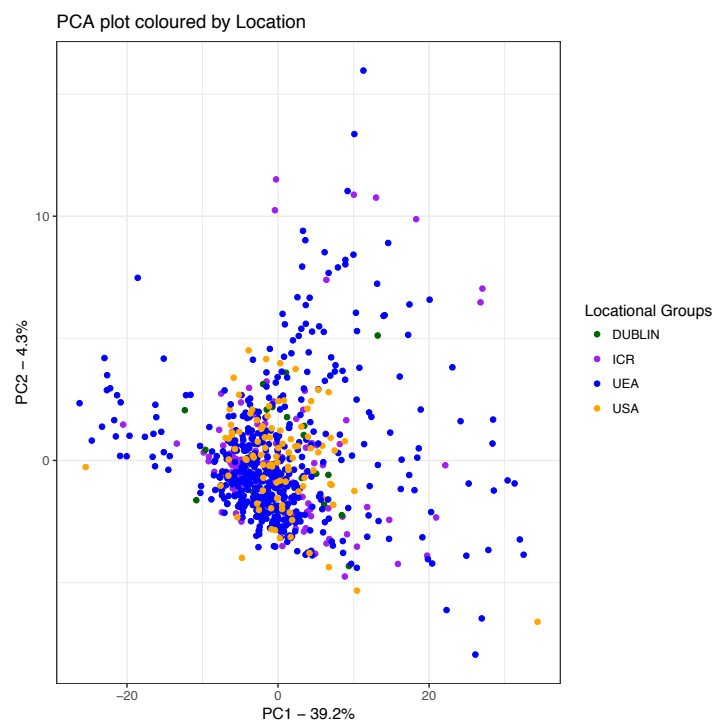


Figure 5.9 PCA plots of post-ComBat data, shows no clustering by location of origin.

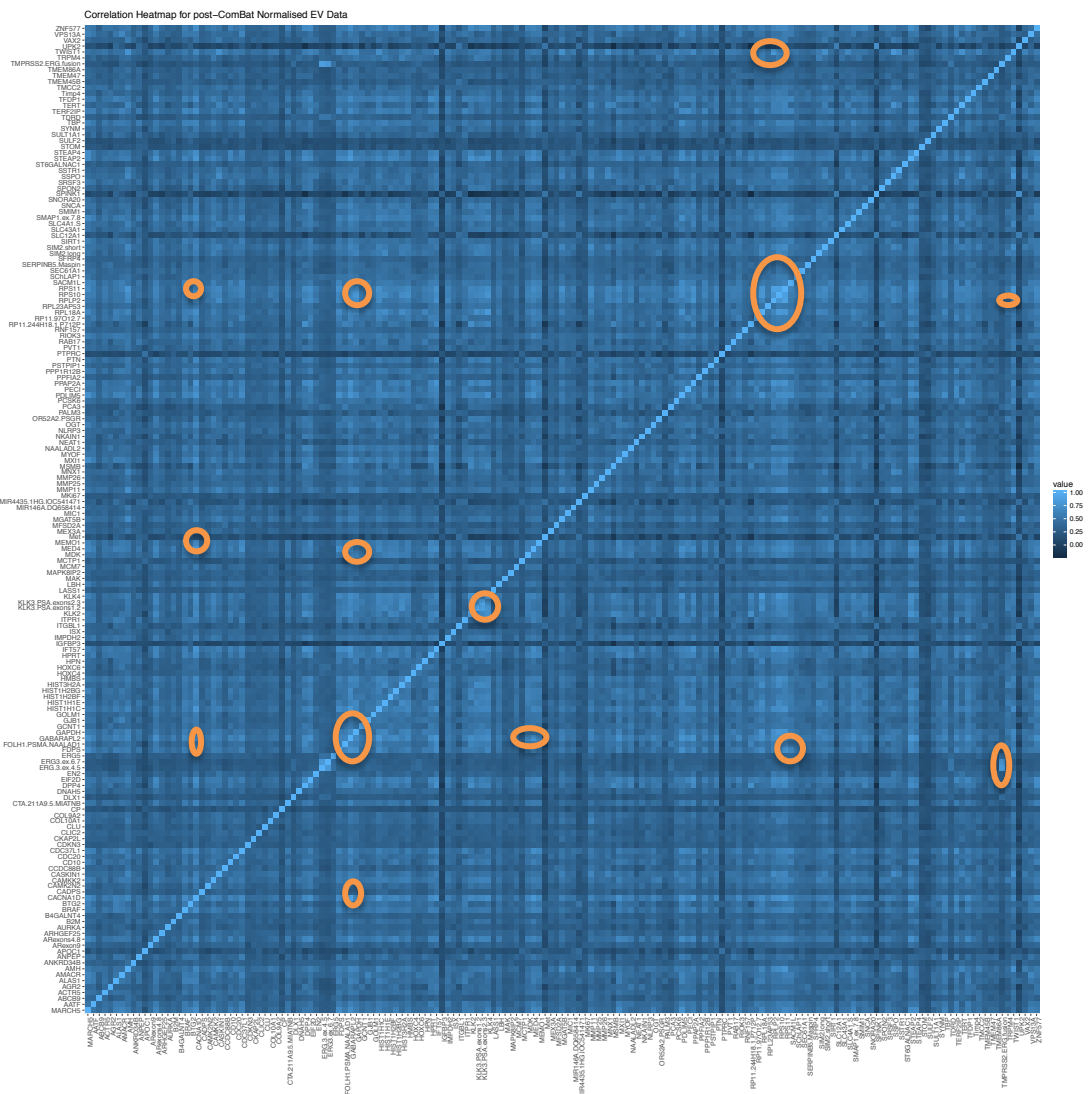
### 5.3.5 Correlating Gene Probes

Pearson’s correlation was used to identify correlating probes (Figure 5.10). There were a number of probes that correlated with  $R > 0.8$ . The correlations were: *CACNAID* with *GABARAPL2* ( $R = 0.965$ ). *ERG3'* exons 4-5 with *TMPRSS2:ERG* ( $R = 0.843$ ). *GABARAPL2* with *CACNAID*, *MED4* and *RPS11* ( $R = 0.965$ ,  $R = 0.805$ , and  $R =$

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

0.804, respectively). *RPLP2* with *RPS11* and *TWIST1* ( $R = 0.859$  and  $R = 0.814$ , respectively). *RPS10* with *RPS11* ( $R = 0.857$ ). *RPS11* with *GABARAPL2*, *RPLP2* and *RPS10* ( $R = 0.804$ ,  $R = 0.859$ , and  $R = 0.857$ , respectively). *TWIST1* with *RPLP2* ( $R = 0.814$ ). Whilst *KLK3* exons 1-2 and *KLK3* exons 2-3 correlated with each other ( $R = 0.814$ ). Whilst *KLK3* exons 1-2 and *KLK3* exons 2-3 correlated with each other ( $R = 0.839$  and  $R = 0.839$ , respectively).

These data correlations were encouraging as many of them fitted with published expression data, for example, expression of *TMPRSS2:ERG* and *ERG3'*, and the two *KLK3* probes. *RPL11* is known to be co-expressed with *RPL10*.



**Figure 5.10** Heatmap showing correlation between NanoString Probes in post-ComBat data.  $R$ -values between 0 (darker) and 1 (lighter). Correlations with  $R > 0.8$  have been highlighted.

### 5.3.6 Comparison of NanoString2 with NanoString1

Comparing the forty-nine common probes across the one hundred and thirty one common samples between NanoString1 and NanoString2 yielded three probes with a Pearson's correlation  $R < 0.6$ : *Timp4* ( $R = 0.14$ ), *TMPRSS2:ERG* ( $R = 0.18$ ), and *TERT* ( $R = 0.38$ ). Twenty-one of the probes showed high correlation, with  $R > 0.9$  (Table 5.6).

**Table 5.6 Pearson's Correlation between the 49 common probes and 131 common samples between NanoString1 and NanoString2.**

<i>Probe</i>	<i>R</i>	<i>Probe</i>	<i>R</i>	<i>Probe</i>	<i>R</i>
<i>HOXC6</i>	0.98	<i>CLU</i>	0.92	<i>HPN</i>	0.83
<i>ERG3' exons 6-7</i>	0.97	<i>KLK3 exons 2-3</i>	0.92	<i>GAPDH</i>	0.83
<i>SPINK1</i>	0.97	<i>KLK3 exons 1-2</i>	0.92	<i>HOXC4</i>	0.82
<i>SULT1A1</i>	0.97	<i>CAMKK2</i>	0.91	<i>AURKA</i>	0.82
<i>KLK2</i>	0.96	<i>STEAP4</i>	0.90	<i>BRAF</i>	0.81
<i>AR exons 4-8</i>	0.96	<i>ANPEP</i>	0.90	<i>PCA3</i>	0.80
<i>KLK4</i>	0.95	<i>AGR2</i>	0.90	<i>PPAP2A</i>	0.78
<i>AR exon 9</i>	0.95	<i>B2M</i>	0.89	<i>IMPDH2</i>	0.78
<i>UPK2</i>	0.95	<i>PECI</i>	0.89	<i>OGT</i>	0.77
<i>FOLH1</i>	0.95	<i>PTPRC</i>	0.89	<i>CDC20</i>	0.71
<i>ALAS1</i>	0.94	<i>DLX1</i>	0.89	<i>MKi67</i>	0.67
<i>AMACR</i>	0.94	<i>MDK</i>	0.89	<i>ERG5'</i>	0.63
<i>TDRD</i>	0.93	<i>MMP26</i>	0.87	<i>TERT</i>	0.37
<i>SLC12A1</i>	0.93	<i>NAALADL2</i>	0.87	<i>TMPRSS2:ERG</i>	0.18
<i>SERPINB5</i>	0.93	<i>TBP</i>	0.86	<i>Timp4</i>	0.14
<i>GOLM1</i>	0.93	<i>CDKN3</i>	0.85		
<i>STEAP2</i>	0.93	<i>HPRT</i>	0.83		

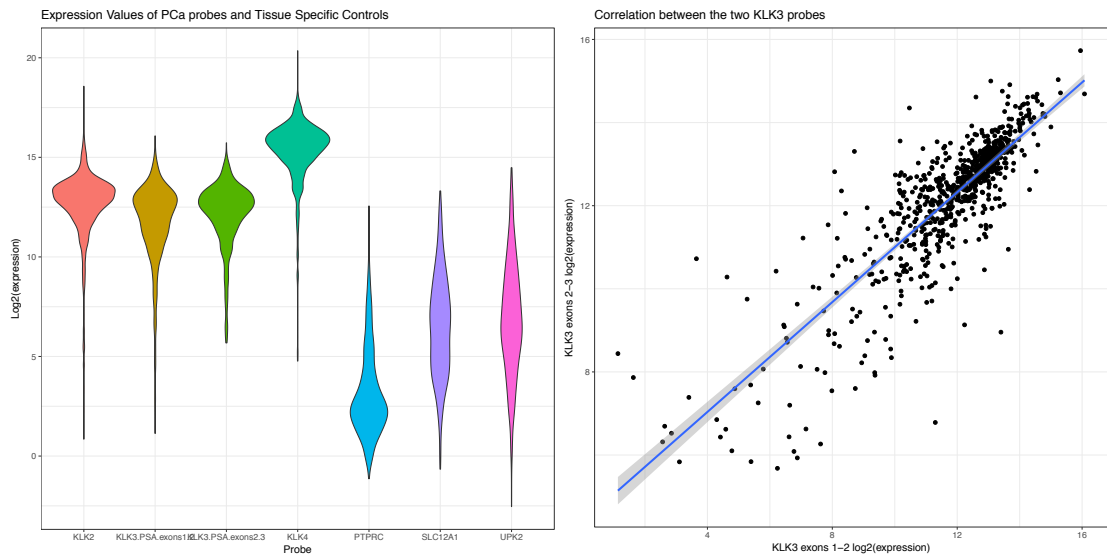
## 5.4 Identification of Prostate and Cancer Specific Transcripts and

### DRE relevance

#### 5.4.1 Kallikrein identification

NanoString median signals for the *KLK2*, *KLK3* exons 1-2, *KLK3* exons 2-3 and *KLK4* probes were again at significantly higher levels than those for the control tissue probes for blood, kidney and bladder (*PTPRC*, *SLC12A1* and *UPK2* respectively) (Mann Whitney U test:  $p < 2.2 \times 10^{-16}$  in each case, Figure 5.11). This was seen previously in NanoString1 (section 3.4.1) and shows that some of the material collected did originate

from the prostate. Once again, similar expression levels and a correlation, is observed between the two *KLK3* probes (Pearson's correlation:  $R = 0.84$ ,  $p < 2.2 \times 10^{-16}$ ).



**Figure 5.11** *KLK2*, *KLK3* and *KLK4* expression is higher than the tissue specific controls for blood, kidney and bladder. The two *KLK3* probes are highly correlated (Pearson's correlation:  $R = 0.84$ ,  $p < 2.2 \times 10^{-16}$ ).

#### 5.4.2 *TMPRSS2:ERG* Identification

Similar results can be seen in regards to the *TMPRSS2:ERG* fusion gene, the *ERG3'* probes and *ERG5'*, as in NanoString1 (section 3.4.2). *TMPRSS2:ERG* fusions, *ERG 3'* and *ERG 5'* expression are linked to PCa, and are therefore expected to be seen more prevalently in samples obtained from men with known PCa compared to those with no clinical evidence of PCa (CBN samples) (Mann Whitney U test between respective probe's expression values and local cancer (low-, intermediate- and high-risk cancer)/CBN groupings. (*TMPRSS2:ERG*:  $p < 2.2 \times 10^{-16}$ , *ERG 3'* exons 4-5:  $p < 2.2 \times 10^{-16}$ ; *ERG 3'* exons 6-7:  $p < 2.2 \times 10^{-16}$ ; and *ERG 5'*:  $p = 1.572 \times 10^{-08}$ ). The density plots for *TMPRSS2:ERG* and the *ERG3'* probes (Figure 5.12) have two peaks which would be compatible with an on/off pattern for that probe suggesting that approximately 50% of the samples from men with cancer have detectable *TMPRSS2:ERG* fusions (which is in agreement with the literature available (section 1.4.6) and the results from NanoString1, (section 3.4.2)).

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

A larger proportion of the CBN and raised PSA negative Bx (S) samples do not have high expression of *TMPRSS2:ERG*, compared to the cancer samples. The cancer samples across all clinical categories and abnormal (including HG:PIN, prostatitis and atypia samples) have fewer samples with lower *TMPRSS2:ERG* expression. The CBN samples also show lower numbers with high *TMPRSS2:ERG* expression, however there are a few (as expected).

The *ERG5'* probe, which is not part of the *TMPRSS2:ERG* fusion transcript, is not significantly different between clinical risk categories. This is also seen in NanoString1 (section 3.4.2). These results suggest that the second set of NanoString data is detecting transcripts accurately and that a proportion of the genetic material identified is coming from PCa or HG-PIN, again.

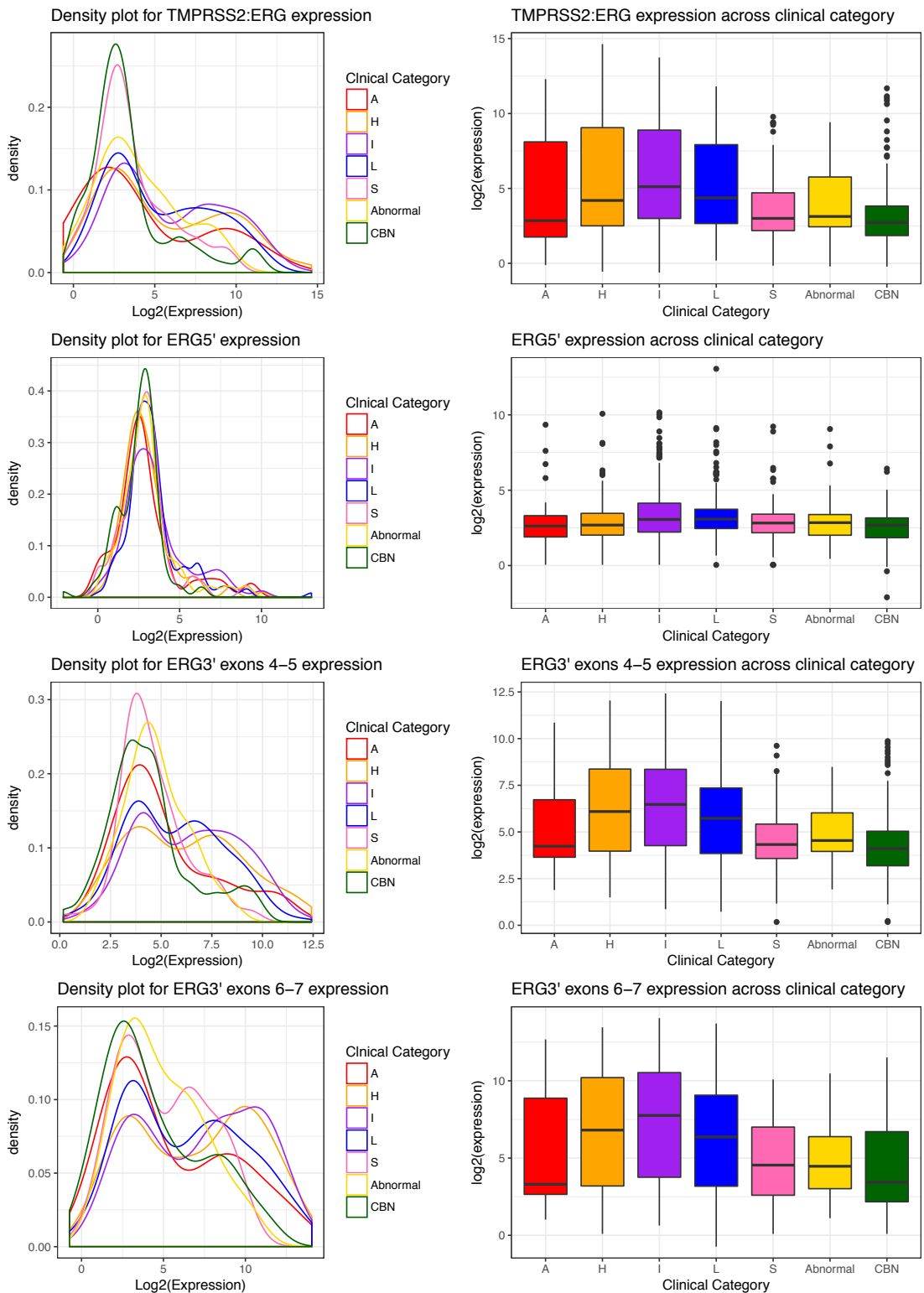


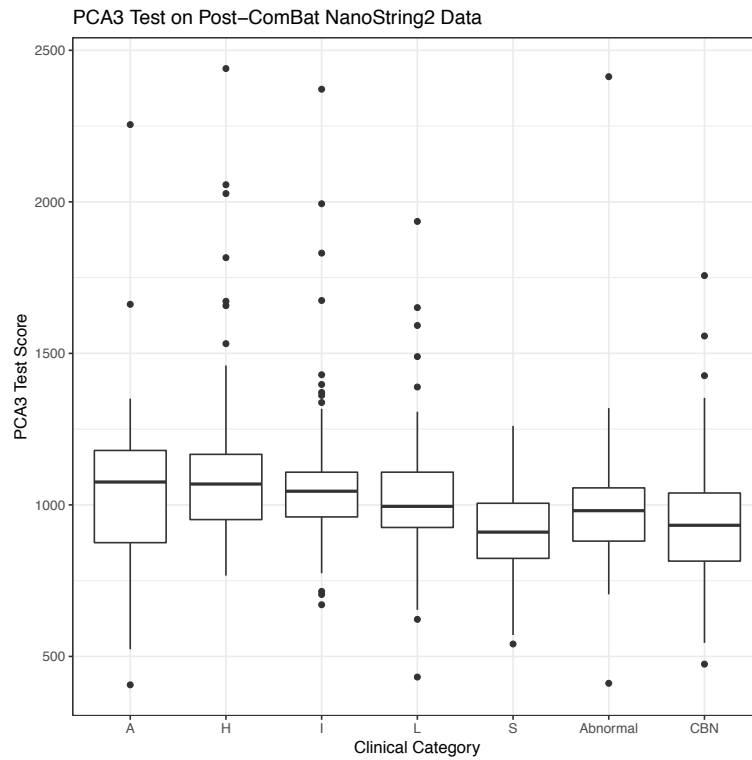
Figure 5.12 Density plots and Boxplots showing the expression changes of *TPMRSS2:ERG*, two *ERG* 3' probes, and *ERG* 5' across clinical categories.

### 5.4.3 PCA3 Test

As in the NanoString1 (section 3.4.3) data, the PCA3 test was significantly different between PCa (Advanced, high-risk, intermediate-risk and low-risk) samples and CBN

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

samples (Kruskal-Wallis rank sum test:  $p = 6.2 \times 10^{-09}$ ,  $\chi^2 = 33.76$  and Mann Whitney U test:  $p < 2.2 \times 10^{-16}$ , Figure 5.13). There are some significant differences across clinical categories also ( $p < 0.05$ ; Mann-Whitney U test; Table 5.7).



**Figure 5.13 PCA3 Test on post-ComBat NanoString2 data (PCA3 transcript expression/average KLK3 transcript expression \* 1000)**



**Table 5.7 Mann Whitney U test of PCA3 Test scores between the different clinical categories.**

<i>p-value</i>	<i>Advanced</i>	<i>High-risk</i>	<i>Intermediate-Risk</i>	<i>Low-Risk</i>	<i>High PSA negative Bx</i>	<i>Abnormal</i>	<i>CBN</i>
<i>Advanced</i>		<b>0.657</b>	<b>0.756</b>	<b>0.255</b>	<b>0.003</b> (Up in A)	<b>0.095</b>	<b>0.021</b> (Up in A)
<i>High-Risk</i>	<b>0.657</b>		<b>0.126</b>	<b>0.004</b> (Up in H)	<b>3.14x10<sup>-10</sup></b> (Up in H)	<b>5.5x10<sup>-04</sup></b> (Up in H)	<b>1.3x10<sup>-07</sup></b> (Up in H)
<i>Intermediate-Risk</i>	<b>0.756</b>	<b>0.126</b>		<b>0.024</b> (Up in I)	<b>2.7x10<sup>-12</sup></b> (Up in I)	<b>0.001</b> (Up in I)	<b>1.2x10<sup>-08</sup></b> (Up in I)
<i>Low-Risk</i>	<b>0.255</b>	<b>0.004</b> (Up in H)	<b>0.024</b> (Up in I)		<b>3.4x10<sup>-07</sup></b> (Up in L)	<b>0.101</b>	<b>1.0x10<sup>-04</sup></b> (Up in L)
<i>High PSA negative Bx</i>	<b>0.003</b> (Up in A)	<b>3.14x10<sup>-10</sup></b> (Up in H)	<b>2.7x10<sup>-12</sup></b> (Up in I)	<b>3.4x10<sup>-07</sup></b> (Up in L)		<b>0.029</b> (Up in Abnormal)	<b>0.408</b>
<i>Abnormal</i>	<b>0.095</b> (Up in A)	<b>5.5x10<sup>-04</sup></b> (Up in H)	<b>0.001</b> (Up in I)	<b>0.101</b>	<b>0.029</b> (Up in Abnormal)		<b>0.189</b>
<i>CBN</i>	<b>0.021</b> (Up in A)	<b>1.3x10<sup>-07</sup></b> (Up in H)	<b>1.2x10<sup>-08</sup></b> (Up in I)	<b>1.0x10<sup>-04</sup></b>	<b>0.408</b>	<b>0.189</b>	

## 5.5 Clustering

### 5.5.1 Principal Component Analysis

PCA (section 2.5.1) shows no significant clustering by clinical category (Kruskal-Wallis rank sum test:  $p = 0.2064$ ,  $\chi = 8.5$ ).

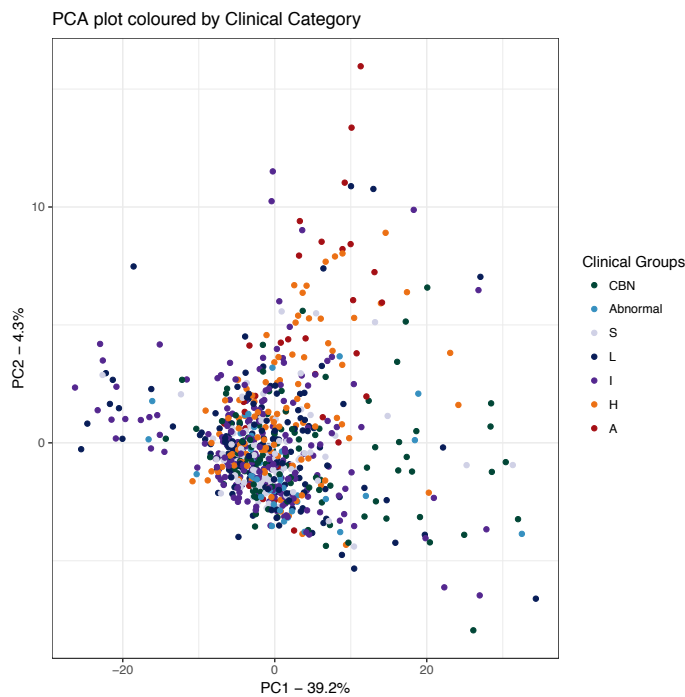


Figure 5.14 PCA plot of post-ComBat data, shows no clustering by clinical category.

### 5.5.2 Latent Process Decomposition (LPD)

LPD (section 2.5.5) was applied to the dataset for three hundred and forty-six of the training samples. There were predicted to be five clusters in the data, with a sigma parameter of -1. LPD analysis was then performed 100 times using these parameters. A significant association was found between LPD group and clinical risk group (Chi-square:  $p = 7.46 \times 10^{-14}$ ,  $\chi = 115$ , Figure 5.15) but not the sample origin (Chi-square:  $p = 0.095$ ,  $\chi = 18.7$ , Figure 5.17, Table 5.8, Figure 5.18). This suggests that this data set is picking up on underlying processes in the NanoString2 data that effects clinical risk. Figure 5.16 shows the clinical breakdown of each LPD group. There appeared to

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

be an over-representation of CBN samples in LPD 1 but this was not significant (Chi-square test CBN vs. low-, intermediate- and high-risk cancer:  $p$ -value = 0.09,  $X^2 = 2.8$ ). LPD2 had an over representation of localised cancer (low-risk and intermediate-risk) Chi-square test:  $p$ -value = 0.037,  $X^2 = 4.3$ . Whilst LPD3 showed a significant over-representation of more progressed cancer (high-risk/advanced cancer) Chi-square test:  $p$ -value =  $1.671 \times 10^{-07}$ ,  $X^2 = 31.2$ . There was no significant over-representation of cancer (advanced, high-, intermediate- and low-risk) or CBN samples in either LPD4 or LPD5. All cancer vs. CBN, more progressed cancer vs. localised cancer vs. CBN were both tested.

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

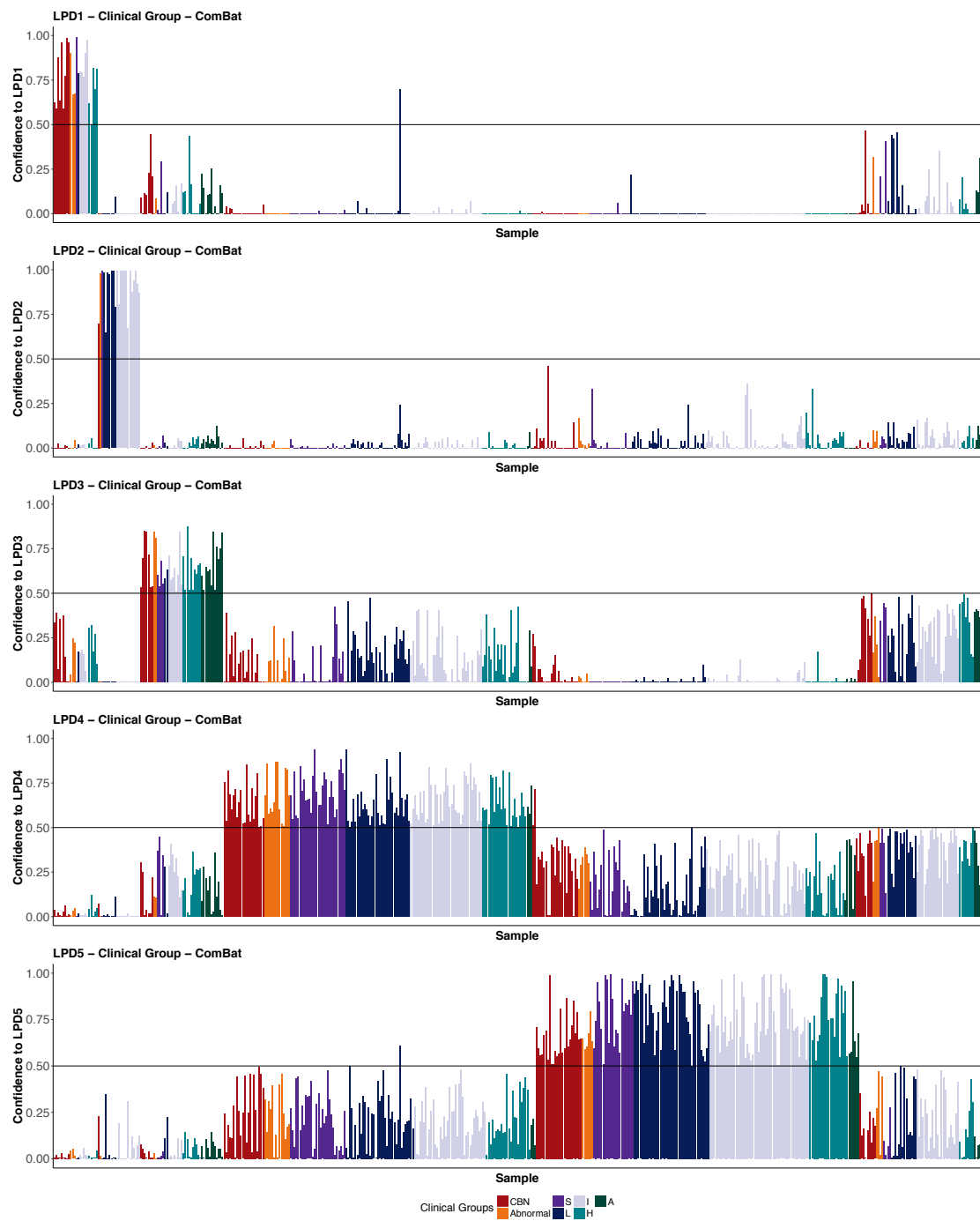


Figure 5.15 LPD of post-ComBat data separated into five processes and coloured by clinical category.

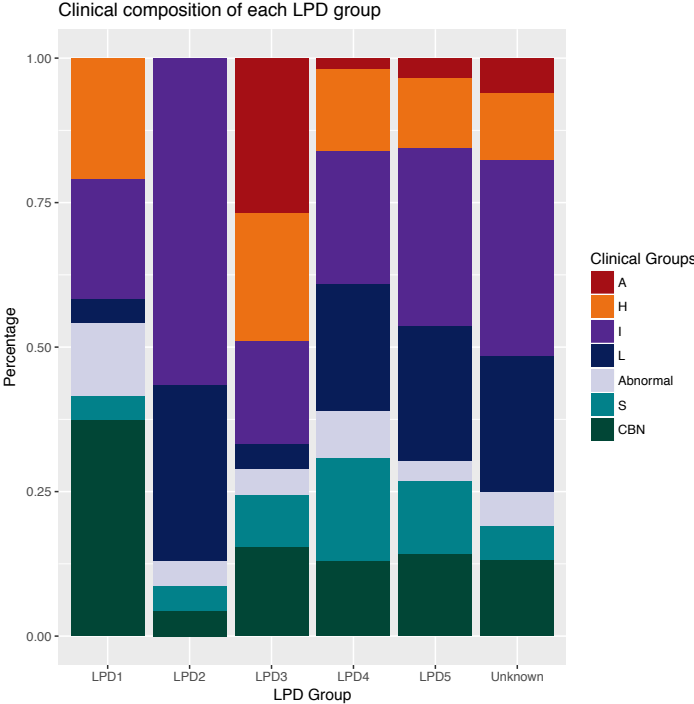
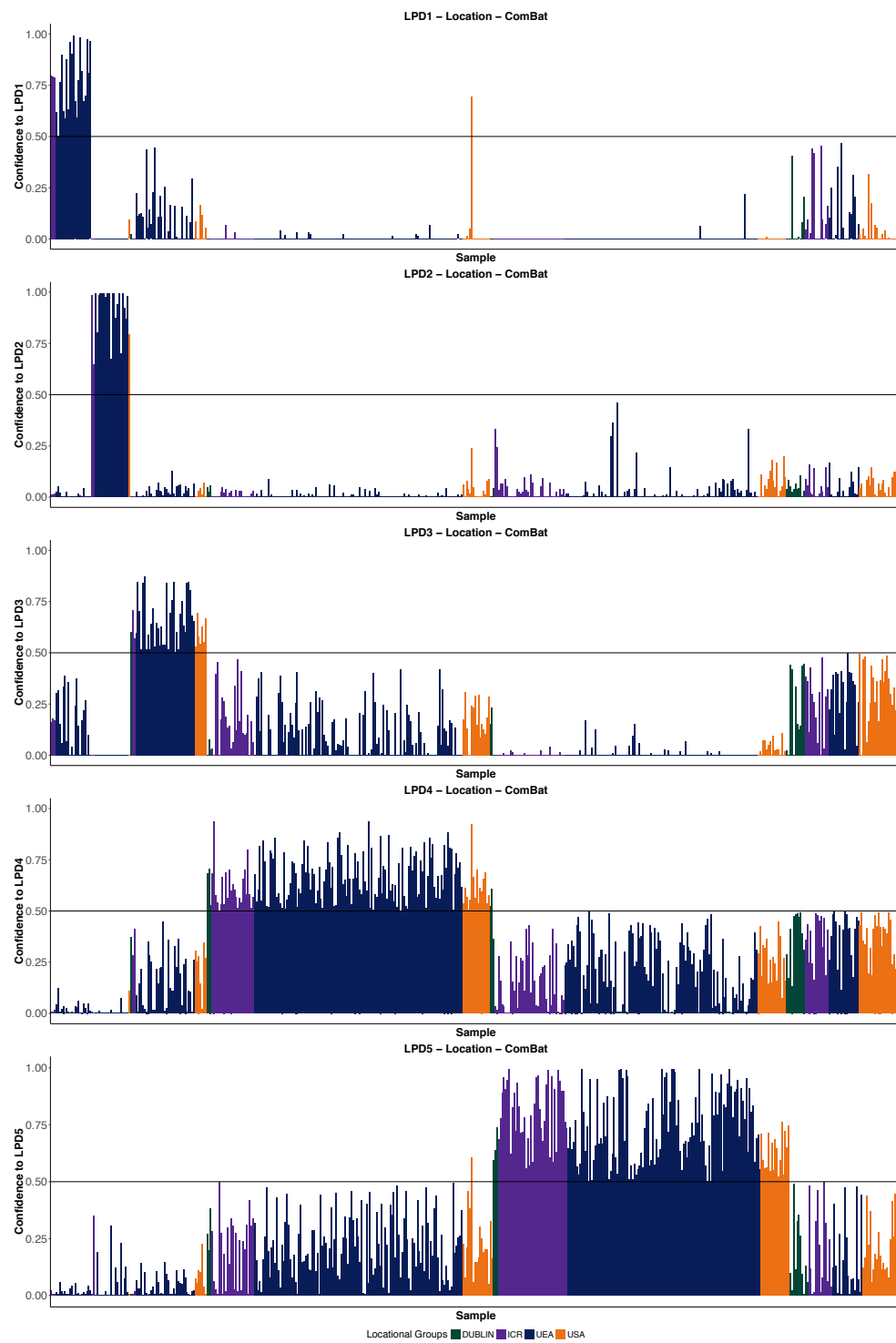


Figure 5.16 Clinical breakdown of each LPD group. Chi-square test:  $p$ -value =  $7.46 \times 10^{-14}$ ,  $X^2 = 115$  (ignoring samples from unknown LPD group).

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

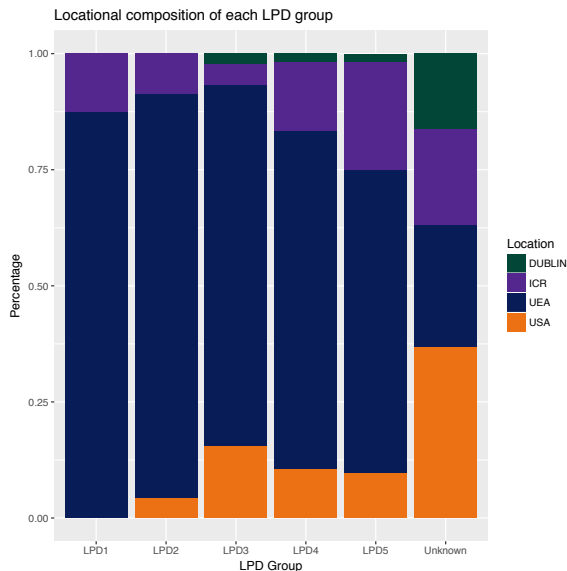


**Figure 5.17** LPD of post-ComBat data separated into five processes and coloured by location of origin.

**Table 5.8 Location of origin breakdown of LPD groups.**

	<i>LPD1</i>	<i>LPD2</i>	<i>LPD3</i>	<i>LPD4</i>	<i>LPD5</i>
<i>DUBLIN</i>	0	0	1	3	3
<i>ICR</i>	3	2	2	25	41
<i>UEA</i>	21	20	35	123	114
<i>USA</i>	0	1	7	18	17
<i>Total</i>	24	23	45	169	175

There were 167, 166, 131, 61, & 153 transcripts that were significantly differentially expressed in LPD processes 1-5 respectively vs. the rest ( $p < 0.05$  after multiple testing correction, Mann-Whitney U test: section 2.4.1). Looking at the top 10 most significant associated transcripts shows a decrease in expression in LPD groups 1, 3 and 4 and an increase in expression in LPD groups 2 and 5 (Table 5.9).



**Figure 5.18 Location of origin breakdown of each LPD group. Chi-square test:  $p$ -value = 0.095,  $X^2 = 18.7$  (ignoring unknown LPD group samples).**

**Table 5.9 Top ten significantly associated transcripts involved in the separation of samples into LPD groups. The p-value shown is adjusted using Benjamin Hochberg multiple testing correction.**

LPD Group	LPD1	p-value	Log <sub>2</sub> (FC)	LPD2	p-value	Log <sub>2</sub> (FC)	LPD3	p-value	Log <sub>2</sub> (FC)
<b># Sig Genes</b>	<b>167</b>			<b>166</b>			<b>131</b>		
<b>Top 10:</b>	<i>CAMKK2</i>	$4.80 \times 10^{-14}$	-1.17	<i>IFT57</i>	$9.50 \times 10^{-14}$	0.18	<i>KLK2</i>	$1.97 \times 10^{-12}$	-0.20
	<i>CACNA1D</i>	$1.10 \times 10^{-13}$	-0.49	<i>OGT</i>	$1.26 \times 10^{-13}$	0.27	<i>DPP4</i>	$2.20 \times 10^{-12}$	-0.23
	<i>GABARAPL</i>	$1.10 \times 10^{-13}$	-0.32	<i>GABARAPL</i>	$1.31 \times 10^{-13}$	0.17	<i>CASKIN1</i>	$1.29 \times 10^{-10}$	-0.21
	2			2					
	<i>RPS11</i>	$3.13 \times 10^{-13}$	-0.13	<i>DPP4</i>	$1.56 \times 10^{-13}$	0.19	<i>MSMB</i>	$1.34 \times 10^{-10}$	-0.08
	<i>RPL23AP53</i>	$3.74 \times 10^{-13}$	-1.35	<i>IMPDH2</i>	$1.65 \times 10^{-13}$	0.26	<i>CACNA1D</i>	$1.55 \times 10^{-10}$	-0.20
	<i>PPAP2A</i>	$3.94 \times 10^{-13}$	-0.33	<i>HPRT</i>	$1.68 \times 10^{-13}$	0.30	<i>GABARAPL</i>	$1.71 \times 10^{-10}$	-0.14
							2		
	<i>CTA.211A9.5/MIATNB</i>	$4.44 \times 10^{-13}$	-2.43	<i>EIF2D</i>	$1.69 \times 10^{-13}$	0.25	<i>TERT</i>	$2.02 \times 10^{-10}$	-0.24
	<i>STEAP2</i>	$5.07 \times 10^{-13}$	-0.60	<i>MXI1</i>	$2.05 \times 10^{-13}$	0.22	<i>ZNF577</i>	$2.69 \times 10^{-10}$	-0.26
	<i>IFT57</i>	$8.73 \times 10^{-13}$	-0.33	<i>PECI</i>	$2.09 \times 10^{-13}$	0.25	<i>SSPO</i>	$3.12 \times 10^{-10}$	-0.20
	<i>MIC1</i>	$8.77 \times 10^{-13}$	-1.20	<i>RP11.97012</i>	$2.10 \times 10^{-13}$	0.28	<i>CAMK2N2</i>	$3.32 \times 10^{-10}$	-0.52
				.7					
LPD Group	LPD4	p-value	Log <sub>2</sub> (FC)	LPD5	p-value	Log <sub>2</sub> (FC)			
<b># Sig Genes</b>	<b>61</b>			<b>153</b>					
<b>Top 10:</b>	<i>VPS13A</i>	$3.38 \times 10^{-06}$	-0.11	<i>GABARA</i>	$2.26 \times 10^{-22}$	0.07			
				<i>PL2</i>					
	<i>TERF2IP</i>	$3.79 \times 10^{-06}$	-0.05	<i>CACNA1D</i>	$2.71 \times 10^{-21}$	0.09			
				<i>D</i>					
	<i>ABCB9</i>	$1.47 \times 10^{-05}$	-0.21	<i>STEAP2</i>	$3.26 \times 10^{-17}$	0.09			
	<i>MARCH5</i>	$1.64 \times 10^{-05}$	-0.08	<i>KLK2</i>	$4.09 \times 10^{-17}$	0.07			



CHAPTER 5: NANOSTRING DATA ANALYSIS 2

<i>MMP25</i>	$1.89 \times 10^{-05}$	-0.25	<i>MED4</i>	$2.31 \times 10^{-16}$	0.09
<i>TMEM45B</i>	$1.92 \times 10^{-05}$	-0.14	<i>CASKIN1</i>	$1.66 \times 10^{-15}$	0.13
<i>RPLP2</i>	$1.93 \times 10^{-05}$	-0.03	<i>DPP4</i>	$7.40 \times 10^{-15}$	0.07
<i>PECI</i>	$2.41 \times 10^{-05}$	-0.06	<i>IFT57</i>	$8.66 \times 10^{-15}$	0.07
<i>CASKIN1</i>	$2.64 \times 10^{-05}$	-0.10	<i>RPS11</i>	$8.75 \times 10^{-14}$	0.03
<i>MEMO1</i>	$3.23 \times 10^{-05}$	-0.08	<i>MARCH5</i>	$9.67 \times 10^{-14}$	0.09

## 5.6 Further processing techniques

After positive control normalisation and  $\log_2$  transformation, four further processing techniques were used. These included adjusting the data to focus on the prostate derived proportion by: using *KLK3* as per the PCA3 test (section 2.1.1), using *KLK2* in a similar way, and using a *KLK2* ratio (section 2.1.1). In addition, *RPLP2* and *GAPDH* were identified as novel housekeeper genes and used to normalise to the amount of material. *RPLP2* and *GAPDH* did not have any significant association with clinical category ( $p < 0.05$ ; Tukey test (section 2.4.7)) and had a strong correlation ( $r = 2.2 \times 10^{-16}$ , Pearson's correlation, section 2.4.3). Each of these methods were used to create a data set and subsequently to build clinical prediction models (Table 5.10). The *KLK2* and *KLK3* adjusted data also included the removal of CBN with high *TMPRSS2:ERG*. As CBN samples were from patients with no clinical evidence of cancer rather than strictly benign, it was expected that there would be some cancer present in some of the samples. Removal of high *TMPRSS2:ERG* CBN samples, was a step towards correcting for this.

Samples with low *KLK2* and *KLK3* values were also removed. These are prostate-expression specific control transcripts. Eliminating these data, removed samples where the majority of the RNA was not originating from the prostate.

**Table 5.10** The different normalisations of the data that the predictive models were built using (separately).

<i>Data</i>	<i>Description</i>
<i>KLK2 ratio</i>	<i>The ratio of KLK2 was used to normalise the data</i>
<i>KLK2 adjusted</i>	<i>Low KLK2 removed and high TMPRSS2:ERG removed. Median and IQR used to adjust data</i>
<i>KLK3 adjusted</i>	<i>Low KLK3 removed and high TMPRSS2:ERG removed. Median and IQR used to adjust data</i>
<i>Housekeeper normalised – GAPDH and RPLP2</i>	<i>KLK2 ratio data, further normalised via GAPDH and RPLP2</i>

### **5.7 Clinical Prediction models**

The data were stratified into test and training sets in the ratio 1:2 (Table 5.1) weighted according to sample origin and clinical risk category. Models were built to predict four different response variables i.e. clinical questions (Table 5.11) using each of the four different processed datasets (Section 5.6) using the training samples.

For models predicting a binary variable, logistic regression (section 2.6.1) and Mann Whitney U (section 2.4.1) tests were used to identify transcripts that individually could predict the two groups ( $p < 0.05$ ). For models predicting an ordinal variable, univariate proportional odds models (polr) were used to identify significant transcripts ( $p < 0.05$ ). Multiple testing correction using Benjamin Hochberg was applied.

For each clinical question, final models were built using LASSO using three input criteria:

1. All 167 probes
2. Probes that were identified as significant in univariate analyses ( $p < 0.05$ ; no multiple testing correction)
3. Probes that were identified as significant in univariate analyses when multiple testing correction was applied (Benjamin Hochberg corrected  $p < 0.05$ )

Models were then applied to the test datasets, where the specificity, sensitivity and PPV of each model was determined (Table 5.13, Table 5.15, Table 5.17, Table 5.21).

**Table 5.11 Clinical predictive models built using the training set and tested using the test set.**

<i>Model</i>	<i>Samples</i>	<i>Model type</i>
<i>CB vs. Cancer</i>	<i>Clinically benign samples Vs low-, intermediate-, and high-risk cancer samples grouped together</i>	<i>Binary</i>
<i>CB vs. High risk cancer</i>	<i>Clinically benign Vs. high-risk cancer (extreme ends of no evidence of cancer and and those with higher grade)</i>	<i>Binary</i>
<i>CB, low-, intermediate-, and high-risk trend</i>	<i>Each sample category is a separate group and ordered</i>	<i>Ordinal</i>
<i>CB, cancer, metastatic cancer trend</i>	<i>Clinically benign samples, with low-, intermediate-, and high-risk cancer samples grouped together, and metastatic cancer samples in groups ordered by severity</i>	<i>Ordinal</i>

### 5.7.1 Models predicting presence of cancer CB and cancer (L, I, H) samples

Expression of 80, 63, 49, 55 probes had a significant association with whether a sample had no evidence for cancer (CB) or not (L, I, H) in the four processed datasets (*KLK2* ratio, *KLK2* adjusted, *KLK3* adjusted, HK normalised, respectively) (Supplementary Table 4). The top probe was *ERG3*' exons 4-5 ( $p = 1.54 \times 10^{-09}$ ,  $\log_2FC = 1.58$ ), *PCA3* ( $p = 4.5 \times 10^{-07}$ ,  $\log_2FC = 0.19$ ), *PCA3* ( $p = 1.61 \times 10^{-06}$ ,  $\log_2FC = 0.14$ ), and *ERG3*' exons 4-5 ( $p = 4.5 \times 10^{-09}$ ,  $\log_2FC = 0.699$ ), respectively.

Multivariate models were built to predict whether a patient had cancer (L, I, H samples) or had no evidence for cancer (CB) (Table 5.12, Table 5.13). The ROC curves and probes involved in each model can be found in the supplementary figures (Supplementary Figure 2, Supplementary Figure 3, Supplementary Figure 4,

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

Supplementary Table 5, Supplementary Table 6, Supplementary Table 7 and Supplementary Table 8, respectively).

In this comparison there were large differences in the number of samples in each of the two categories, with CB having approximately only a quarter of the sample size of cancer. Therefore, random sampling was used to select a similar number of cancer samples to CB samples, and the model predictive process was run iteratively 1,000 times. The model with the mean AUC was selected to be applied to the test dataset. Again, the AUC, Sensitivity, Specificity and PPV and the selected probes were recorded for each model on the training set (Table 5.14) and the test set (Table 5.15) and the curves and probes involved in each model can be observed in the supplementary figures (Supplementary Figure 5, Supplementary Figure 6, Supplementary Figure 7, Supplementary Figure 8, Supplementary Table 9, Supplementary Table 10, Supplementary Table 11, and Supplementary Table 12, respectively).

The models were generally good predictors of whether cancer was present or not (median AUC = 0.8045, IQR = 0.06). In general, AUC in the test data was better in the *KLK2* ratio and the *GAPDH* and *RPLP2* normalised data (all had AUC > 0.8) compared to the *KLK2* and *KLK3* adjusted data (mostly AUC > 0.7). There was not much difference observed between those with the randomly selected cancer samples (median AUC = 0.847, IQR = 0.11), and those with all of the cancer samples (median AUC = 0.846, IQR = 0.098).

The accuracy measures remained very high in the test sets (median AUC = 0.915, IQR = 0.05, but were slightly lower than the training data set (median AUC = 0.8045, IQR = 0.06), showing the models in general were robust and useful.

The model with the best AUC in the training data, was when using all of the probes from the *RPLP2* and *GAPDH* normalised data (Training AUC = 0.925, Test AUC = 0.851) in detail as an example. 18 transcripts were selected by Lasso and went into these models; *TMPRSS2:ERG*, *ERG3'* exons 4-5, *APOC1*, *ISX*, *SLC12A1*, *HOXC6*, *MCTP1*, *TDRD*, *PDLIM5*, *CD10*, *GABARAPL2*, *PTN*, *AR* exon 9, *PPP1R12B*, *CP*,

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

*MXII*, and *KLK4*. The training model had 85% sensitivity, 73% specificity and 94% PPV (Figure 5.19).

CHAPTER 5: NANOSTRING DATA ANALYSIS 2

**Table 5.12 Training model outcomes comparing CB with Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.**

	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>AUC</i>	<b>0.949</b>	<b>0.886</b>	<b>0.891</b>	<b>0.91</b>	<b>0.929</b>	<b>0.849</b>	<b>0.966</b>	<b>0.935</b>	<b>0.824</b>	<b>0.925</b>	<b>0.902</b>	<b>0.859</b>
<i>Sensitivity</i>	<b>89%</b>	<b>77%</b>	<b>71%</b>	<b>95%</b>	<b>93%</b>	<b>72%</b>	<b>95%</b>	<b>89%</b>	<b>68%</b>	<b>88%</b>	<b>81%</b>	<b>75%</b>
<i>Specificity</i>	<b>89%</b>	<b>87%</b>	<b>92%</b>	<b>71%</b>	<b>81%</b>	<b>87%</b>	<b>89%</b>	<b>89%</b>	<b>86%</b>	<b>85%</b>	<b>87%</b>	<b>90%</b>
<i>PPV</i>	<b>97%</b>	<b>96%</b>	<b>97%</b>	<b>92%</b>	<b>94%</b>	<b>95%</b>	<b>97%</b>	<b>96%</b>	<b>94%</b>	<b>97%</b>	<b>97%</b>	<b>98%</b>
<i>Threshold</i>	<b>0.6899804</b>	<b>0.7713957</b>	<b>0.8137426</b>	<b>0.6489155</b>	<b>0.6545943</b>	<b>0.7788927</b>	<b>0.6751793</b>	<b>0.7230226</b>	<b>0.7683476</b>	<b>0.7735114</b>	<b>0.8235587</b>	<b>0.8317314</b>
<i>Number of Probes</i>	<b>21</b>	<b>4</b>	<b>8</b>	<b>26</b>	<b>31</b>	<b>6</b>	<b>50</b>	<b>29</b>	<b>4</b>	<b>18</b>	<b>10</b>	<b>6</b>

CHAPTER 5: NANOSTRING DATA ANALYSIS 2

Table 5.13 Test model outcomes comparing CB with Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.

	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>AUC</i>	<b>0.846</b>	<b>0.819</b>	<b>0.816</b>	<b>0.772</b>	<b>0.776</b>	<b>0.775</b>	<b>0.745</b>	<b>0.762</b>	<b>0.718</b>	<b>0.851</b>	<b>0.838</b>	<b>0.816</b>
<i>Sensitivity</i>	<b>89%</b>	<b>89%</b>	<b>68%</b>	<b>59%</b>	<b>69%</b>	<b>62%</b>	<b>72%</b>	<b>74%</b>	<b>68%</b>	<b>85%</b>	<b>83%</b>	<b>60%</b>
<i>Specificity</i>	<b>67%</b>	<b>63%</b>	<b>83%</b>	<b>91%</b>	<b>82%</b>	<b>82%</b>	<b>77%</b>	<b>77%</b>	<b>68%</b>	<b>73%</b>	<b>73%</b>	<b>93%</b>
<i>PPV</i>	<b>91%</b>	<b>90%</b>	<b>92%</b>	<b>96%</b>	<b>93%</b>	<b>92%</b>	<b>91%</b>	<b>90%</b>	<b>87%</b>	<b>94%</b>	<b>93%</b>	<b>97%</b>
<i>Threshold</i>	<b>0.63388</b>	<b>0.632465</b>	<b>0.832557</b>	<b>0.79834</b>	<b>0.787580</b>	<b>0.817720</b>	<b>0.76251</b>	<b>0.712675</b>	<b>0.768974</b>	<b>0.74827</b>	<b>0.759985</b>	<b>0.878513</b>
<i>Number of Probes</i>	<b>98</b>	<b>5</b>	<b>8</b>	<b>29</b>	<b>3</b>	<b>5</b>	<b>93</b>	<b>5</b>	<b>4</b>	<b>62</b>	<b>1</b>	<b>9</b>
	<b>21</b>	<b>4</b>	<b>8</b>	<b>26</b>	<b>31</b>	<b>6</b>	<b>50</b>	<b>29</b>	<b>4</b>	<b>18</b>	<b>10</b>	<b>6</b>



CHAPTER 5: NANOSTRING DATA ANALYSIS 2

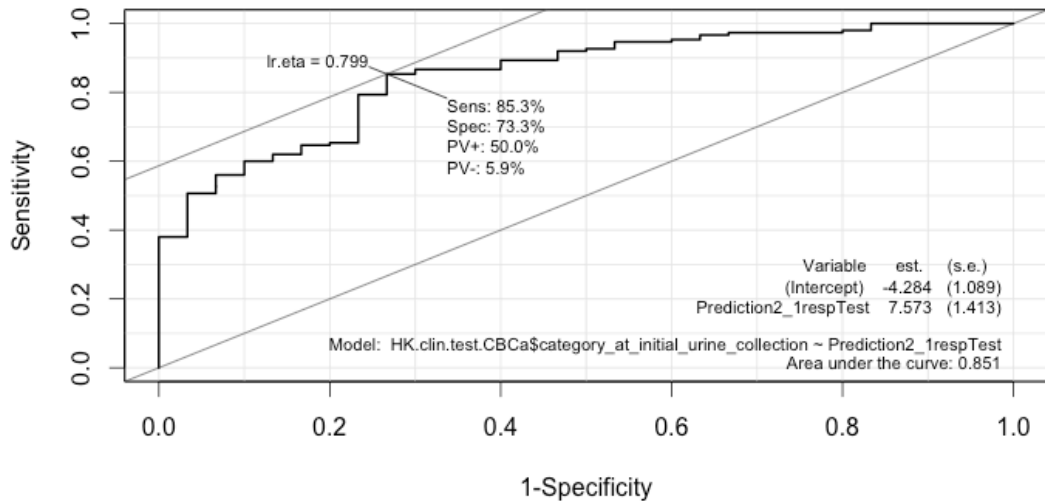
**Table 5.14 Training model outcomes comparing CB with randomly selected Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.**

	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>AUC</i>	<b>0.957</b>	<b>0.916</b>	<b>0.876</b>	<b>0.991</b>	<b>0.924</b>	<b>0.893</b>	<b>0.915</b>	<b>0.943</b>	<b>0.851</b>	<b>0.936</b>	<b>0.915</b>	<b>0.87</b>
<i>Sensitivity</i>	<b>87%</b>	<b>94%</b>	<b>73%</b>	<b>98%</b>	<b>94%</b>	<b>83%</b>	<b>86%</b>	<b>94%</b>	<b>81%</b>	<b>87%</b>	<b>85%</b>	<b>75%</b>
<i>Specificity</i>	<b>94%</b>	<b>71%</b>	<b>87%</b>	<b>94%</b>	<b>79%</b>	<b>90%</b>	<b>87%</b>	<b>83%</b>	<b>79%</b>	<b>90%</b>	<b>85%</b>	<b>86%</b>
<i>PPV</i>	<b>92%</b>	<b>92%</b>	<b>84%</b>	<b>92%</b>	<b>82%</b>	<b>88%</b>	<b>86%</b>	<b>84%</b>	<b>94%</b>	<b>88%</b>	<b>89%</b>	<b>85%</b>
<i>Threshold</i>	<b>0.4473</b>	<b>0.362903</b>	<b>0.487065</b>	<b>0.40367</b>	<b>0.404678</b>	<b>0.462369</b>	<b>0.44877</b>	<b>0.379279</b>	<b>0.403381</b>	<b>0.43971</b>	<b>0.469126</b>	<b>0.517129</b>
<i>Number of Probes</i>	<b>12</b>	<b>8</b>	<b>7</b>	<b>53</b>	<b>8</b>	<b>8</b>	<b>93</b>	<b>4</b>	<b>59</b>	<b>3</b>	<b>9</b>	<b>5</b>
	<b>17</b>	<b>9</b>	<b>5</b>	<b>35</b>	<b>19</b>	<b>6</b>	<b>16</b>	<b>20</b>	<b>4</b>	<b>8</b>	<b>7</b>	<b>5</b>

CHAPTER 5: NANOSTRING DATA ANALYSIS 2

Table 5.15 Test model outcomes comparing CB with randomly selected Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.

	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>AUC</i>	<b>0.843</b>	<b>0.803</b>	<b>0.813</b>	<b>0.806</b>	<b>0.751</b>	<b>0.768</b>	<b>0.72</b>	<b>0.713</b>	<b>0.695</b>	<b>0.821</b>	<b>0.828</b>	<b>0.808</b>
<i>Sensitivity</i>	<b>78.00%</b>	<b>71%</b>	<b>69%</b>	<b>67%</b>	<b>57%</b>	<b>64%</b>	<b>85%</b>	<b>67%</b>	<b>74%</b>	<b>65%</b>	<b>82%</b>	<b>63%</b>
<i>Specificity</i>	<b>80.00%</b>	<b>80%</b>	<b>83%</b>	<b>91%</b>	<b>88%</b>	<b>88%</b>	<b>56%</b>	<b>71%</b>	<b>59%</b>	<b>90%</b>	<b>70%</b>	<b>87%</b>
<i>PPV</i>	<b>92.00%</b>	<b>92%</b>	<b>93%</b>	<b>96%</b>	<b>92%</b>	<b>93%</b>	<b>85%</b>	<b>87%</b>	<b>85%</b>	<b>96%</b>	<b>93%</b>	<b>96</b>
<i>Threshold</i>	<b>0.41796</b>	<b>0.459550</b>	<b>0.523588</b>	<b>0.47894</b>	<b>0.591645</b>	<b>0.585591</b>	<b>0.39739</b>	<b>0.490522</b>	<b>0.372621</b>	<b>0.5083</b>	<b>0.401851</b>	<b>0.594466</b>
<i>Number of Probes</i>	<b>17</b>	<b>9</b>	<b>5</b>	<b>35</b>	<b>19</b>	<b>6</b>	<b>16</b>	<b>20</b>	<b>4</b>	<b>8</b>	<b>7</b>	<b>5</b>



**Figure 5.19** ROC curve of top performing model for the prediction of CB vs. Cancer (Low-, Intermediate- and High-risk).

### 5.7.2 Models to distinguish the extreme categories i.e. CB and high-risk cancer samples

Expression of 98, 43, 39, 39 probes had a significant association with whether a sample was high-risk (H) or there was no evidence for cancer (CBN) in the four processed datasets (*KLK2* ratio, *KLK2* adjusted, *KLK3* adjusted, HK normalised, respectively) (Supplementary Table 13). The top probe was *ERG3'* exons 4-5 ( $p = 6.995 \times 10^{-07}$ , logFC = 1.87), *HPN* ( $p = 3.767 \times 10^{-06}$ , logFC = 0.24), *HPN* ( $p = 1.317 \times 10^{-05}$ , logFC = 0.19), and *ERG3'* exons 4-5 ( $p = 1.42 \times 10^{-06}$ , logFC = 0.79), respectively.

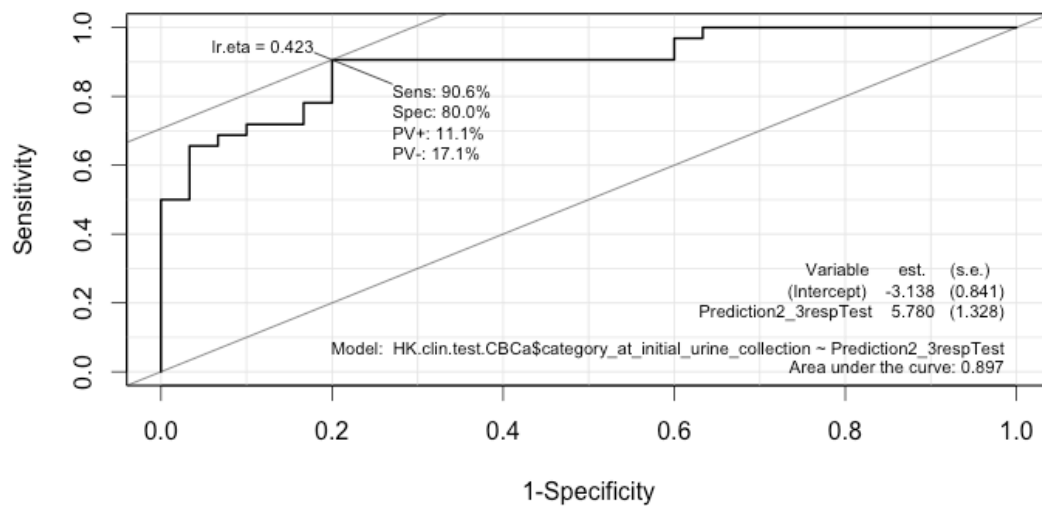
Binomial models were built to predict whether a patient was at high risk of cancer (H) or had no evidence for cancer (CB) (Table 5.16, Table 5.17, see Supplementary Table 14, Supplementary Table 15, Supplementary Table 16 and Supplementary Table 17).

The models were decent predictors (test model median AUC = 0.957, IQR = 0.036, training model median AUC = 0.831, IQR = 0.07). In general, the metrics of the models didn't seem to differ much between the different normalisations (slightly lower AUCs in the *KLK3* adjusted data), or the input probe subset. Models with AUC of up to 0.9 were seen in the training sets, and models with AUC of up to 0.8 were seen when applying the models to the test data. Sensitivities in the 90% and PPVs in the 80%

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

were observed on the test data, suggesting these models were capable to distinguishing well between the CB and high-risk cancer samples.

The model built using the adjusted significant probe lists from the *GAPDH* and *RPLP2* normalised data gave a high AUC of 0.897 in the training data (AUC = 0.924 in the test data). This model had high sensitivity (91%), 80% specificity and 83% PPV (ROC - Figure 5.20). The transcripts used to build this model were *PCA3*, *APOC1*, *HPN*, *ERG3*' exons 4-5 and *TMPRSS2:ERG*.



**Figure 5.20** ROC curve of the training set for the *GAPDH* and *RPLP2* normalised model built using the 5 significant probes post multiple testing correction.

CHAPTER 5: NANOSTRING DATA ANALYSIS 2

**Table 5.16 Training model outcomes comparing CB with high-risk Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.**

	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>AUC</i>	<b>0.991</b>	<b>0.97</b>	<b>0.94</b>	<b>0.952</b>	<b>0.955</b>	<b>0.866</b>	<b>0.962</b>	<b>0.959</b>	<b>0.85</b>	<b>0.976</b>	<b>0.992</b>	<b>0.924</b>
<i>Sensitivity</i>	<b>100%</b>	<b>84%</b>	<b>86%</b>	<b>86%</b>	<b>86%</b>	<b>71%</b>	<b>91%</b>	<b>84%</b>	<b>74%</b>	<b>94%</b>	<b>97%</b>	<b>96%</b>
<i>Specificity</i>	<b>92%</b>	<b>98%</b>	<b>90%</b>	<b>94%</b>	<b>94%</b>	<b>91%</b>	<b>90%</b>	<b>97%</b>	<b>90%</b>	<b>92%</b>	<b>96%</b>	<b>77%</b>
<i>PPV</i>	<b>93%</b>	<b>97%</b>	<b>92%</b>	<b>92%</b>	<b>94%</b>	<b>89%</b>	<b>91%</b>	<b>95%</b>	<b>90%</b>	<b>94%</b>	<b>97%</b>	<b>84%</b>
<i>Threshold</i>	<b>0.40659</b>	<b>0.578702</b>	<b>0.541775</b>	<b>0.50800</b>	<b>0.503871</b>	<b>0.554859</b>	<b>0.44880</b>	<b>0.538367</b>	<b>0.532871</b>	<b>0.47027</b>	<b>0.509973</b>	<b>0.403401</b>
<i>Number of Probes</i>	<b>26</b>	<b>16</b>	<b>9</b>	<b>19</b>	<b>17</b>	<b>5</b>	<b>21</b>	<b>19</b>	<b>3</b>	<b>13</b>	<b>21</b>	<b>5</b>

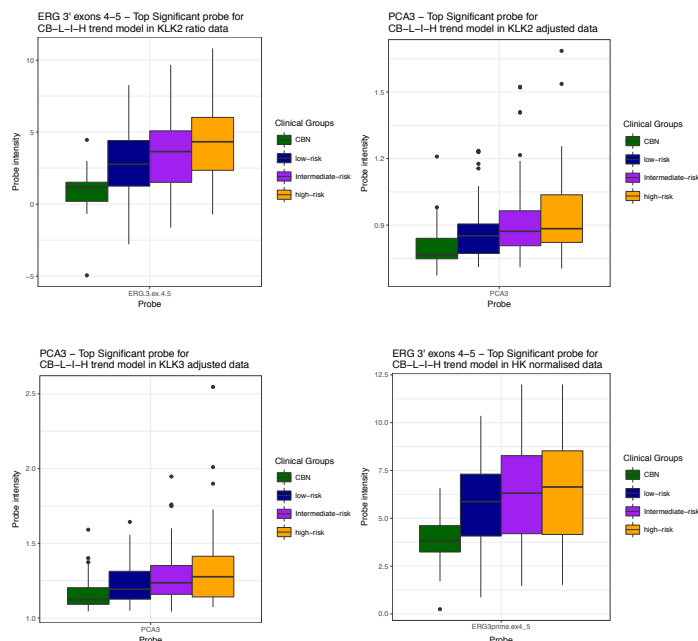
CHAPTER 5: NANOSTRING DATA ANALYSIS 2

**Table 5.17 Test model outcomes comparing CB with high –risk Cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.**

	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>AUC</i>	<b>0.851</b>	<b>0.859</b>	<b>0.832</b>	<b>0.822</b>	<b>0.829</b>	<b>0.738</b>	<b>0.789</b>	<b>0.796</b>	<b>0.738</b>	<b>0.897</b>	<b>0.883</b>	<b>0.897</b>
<i>Sensitivity</i>	<b>88%</b>	<b>97%</b>	<b>94%</b>	<b>91%</b>	<b>91%</b>	<b>91%</b>	<b>97%</b>	<b>97%</b>	<b>91%</b>	<b>88%</b>	<b>84%</b>	<b>91%</b>
<i>Specificity</i>	<b>77%</b>	<b>63%</b>	<b>60%</b>	<b>65%</b>	<b>77%</b>	<b>59%</b>	<b>65%</b>	<b>65%</b>	<b>59%</b>	<b>83%</b>	<b>83%</b>	<b>80%</b>
<i>PPV</i>	<b>80%</b>	<b>73%</b>	<b>70%</b>	<b>71%</b>	<b>76%</b>	<b>71%</b>	<b>70%</b>	<b>72%</b>	<b>67%</b>	<b>82%</b>	<b>84%</b>	<b>83%</b>
<i>Threshold</i>	<b>0.38343</b> <b>88</b>	<b>0.264001</b> <b>1</b>	<b>0.229727</b> <b>1</b>	<b>0.35623</b> <b>61</b>	<b>0.458083</b> <b>8</b>	<b>0.402334</b> <b>6</b>	<b>0.33008</b> <b>49</b>	<b>0.286475</b> <b>5</b>	<b>0.402334</b> <b>6</b>	<b>0.52862</b> <b>75</b>	<b>0.435489</b> <b>3</b>	<b>0.488914</b> <b>3</b>
<i>Number of Probes</i>	<b>26</b>	<b>16</b>	<b>9</b>	<b>19</b>	<b>17</b>	<b>5</b>	<b>21</b>	<b>19</b>	<b>3</b>	<b>13</b>	<b>21</b>	<b>5</b>

### 5.7.3 Models to predict risk categories using trends in expression

Expression of 114, 45, 50, 53 probes had a significant association with increasing risk category (CB->L->I->H) in the four processed datasets (*KLK2* ratio, *KLK2* adjusted, *KLK3* adjusted, HK normalised, respectively) (Supplementary Table 18). The top probe was *ERG3*' exons 4-5 ( $p = 1.86 \times 10^{-13}$ ), *PCA3* ( $p = 1.45 \times 10^{-08}$ ), *PCA3* ( $p = 1.52 \times 10^{-07}$ ), and *ERG3*' exons 4-5 ( $p = 1.44 \times 10^{-08}$ ) respectively (Figure 5.21).



**Figure 5.21 Top Significant Probe for CB, low-risk, intermediate-risk and high-risk cancer trend in all four data normalisations.**

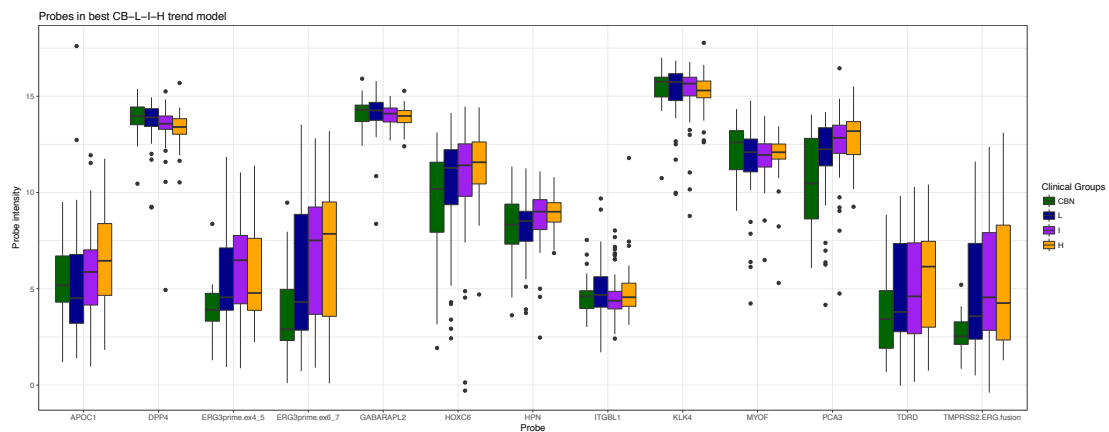
Multivariate proportional odds models were built to predict CB samples, the low-, intermediate- and high-risk cancer samples (section 2.6.1) (Table 5.18, Table 5.19). The probes involved in each model can be observed (Supplementary Table 19, Supplementary Table 20, Supplementary Table 21, and Supplementary Table 22).

The metrics of the models for the *KLK2* ratio and *KLK2* adjusted data were very similar (median = 0.67015, IQR = 0.06 and median AUC 0.6689, IQR = 0.08). Slightly lower AUCs were observed in the *KLK3* adjusted data (median AUC = 0.669, IQR = 0.1), and slightly higher AUCs were observed in the *GAPDH* and *RPLP2* normalised data (median AUC = 0.73385, IQR = 0.05). The average model metrics for the test data were poorer than with previous clinical questions (median AUC = 0.65005, IQR = 0.05). The

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

sensitivity of all of the models were fairly low (median= 29%, IQR=0.45), whilst specificity fairly high (median= 91% IQR=0.23). This suggested that separating between the different risk categories of cancer can be difficult.

The model built using the *GAPDH* and *RPLP2* normalised data and only the probes still significant post multiple testing correction has the highest AUC = 0.7088. The probes used to build this model were *APOC1*, *DPP4*, *ERG 3' exons 4-5*, *ERG 3' exons 6-7*, *GABARAPL2*, *HOXC6*, *HPN*, *ITGBL1*, *KLK4*, *MYOF*, *PCA3*, *TDRD*, and *TMPRSS2:ERG* (Figure 5.22). The Sensitivities of this model ranged from 9%-79% and the specificities ranged from 46%-95%.



**Figure 5.22** Boxplot showing the expression level of each transcript featured in the CB-L-I-H model built using the multiple tested correction significant probes from the *GAPDH* and *RPLP2* normalised data. This model showed the best test data AUC (0.7008).



Table 5.18 Training model outcomes comparing CB, low-, intermediate- and high- risk cancer samples for the four different normalisations of data.

Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction

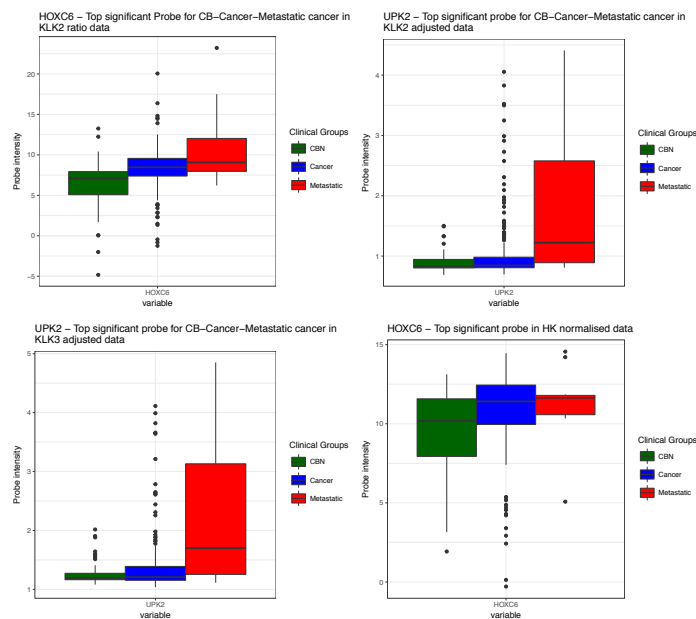
	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>Accuracy</i>	<b>0.5112</b>	<b>0.4581</b>	<b>0.4413</b>	<b>0.5412</b>	<b>0.522</b>	<b>0.4643</b>	<b>0.6749</b>	<b>0.576</b>	<b>0.5124</b>	<b>0.5</b>	<b>0.4944</b>	<b>0.5112</b>
<i>AUC</i>	<b>0.7663</b>	<b>0.6757</b>	<b>0.6196</b>	<b>0.7802</b>	<b>0.7469</b>	<b>0.6856</b>	<b>0.8146</b>	<b>0.7606</b>	<b>0.6929</b>	<b>0.7587</b>	<b>0.7728</b>	<b>0.7608</b>
<i>Sensitivity:</i>												
<i>CB</i>	<b>52%</b>	<b>38%</b>	<b>27%</b>	<b>45%</b>	<b>48%</b>	<b>36%</b>	<b>84%</b>	<b>74%</b>	<b>60%</b>	<b>50%</b>	<b>56%</b>	<b>54%</b>
<i>L</i>	<b>19%</b>	<b>8%</b>	<b>4%</b>	<b>26%</b>	<b>24%</b>	<b>12%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>15%</b>	<b>19%</b>	<b>25%</b>
<i>I</i>	<b>84%</b>	<b>88%</b>	<b>94%</b>	<b>84%</b>	<b>81%</b>	<b>83%</b>	<b>89%</b>	<b>80%</b>	<b>75%</b>	<b>84%</b>	<b>74%</b>	<b>76%</b>
<i>H</i>	<b>25%</b>	<b>10%</b>	<b>7%</b>	<b>37%</b>	<b>32%</b>	<b>25%</b>	<b>59%</b>	<b>41%</b>	<b>36%</b>	<b>25%</b>	<b>34%</b>	<b>31%</b>
<i>Specificity:</i>												
<i>CB</i>	<b>97%</b>	<b>96%</b>	<b>96%</b>	<b>97%</b>	<b>96%</b>	<b>94%</b>	<b>95%</b>	<b>91%</b>	<b>85%</b>	<b>97%</b>	<b>95%</b>	<b>96%</b>
<i>L</i>	<b>87%</b>	<b>91%</b>	<b>96%</b>	<b>87%</b>	<b>89%</b>	<b>92%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>86%</b>	<b>84%</b>	<b>83%</b>
<i>I</i>	<b>42%</b>	<b>26%</b>	<b>18%</b>	<b>47%</b>	<b>45%</b>	<b>36%</b>	<b>59%</b>	<b>53%</b>	<b>48%</b>	<b>40%</b>	<b>47%</b>	<b>50%</b>
<i>H</i>	<b>98%</b>	<b>98%</b>	<b>100%</b>	<b>97%</b>	<b>97%</b>	<b>96%</b>	<b>96%</b>	<b>92%</b>	<b>93%</b>	<b>98%</b>	<b>97%</b>	<b>97%</b>
<i>PPV:</i>												
<i>CB</i>	<b>73%</b>	<b>65%</b>	<b>54%</b>	<b>74%</b>	<b>68%</b>	<b>54%</b>	<b>81%</b>	<b>69%</b>	<b>53%</b>	<b>74%</b>	<b>67%</b>	<b>68%</b>
<i>L</i>	<b>32%</b>	<b>23%</b>	<b>25%</b>	<b>44%</b>	<b>42%</b>	<b>32%</b>	<i>NA</i>	<i>NA</i>	<i>NA</i>	<b>27%</b>	<b>28%</b>	<b>34%</b>
<i>I</i>	<b>49%</b>	<b>44%</b>	<b>44%</b>	<b>51%</b>	<b>49%</b>	<b>46%</b>	<b>58%</b>	<b>52%</b>	<b>48%</b>	<b>49%</b>	<b>48%</b>	<b>51%</b>
<i>H</i>	<b>78%</b>	<b>58%</b>	<b>83%</b>	<b>72%</b>	<b>72%</b>	<b>58%</b>	<b>82%</b>	<b>63%</b>	<b>63%</b>	<b>78%</b>	<b>73%</b>	<b>69%</b>
<i>Number of Probes</i>	<b>36</b>	<b>13</b>	<b>5</b>	<b>12</b>	<b>37</b>	<b>14</b>	<b>78</b>	<b>39</b>	<b>12</b>	<b>37</b>	<b>34</b>	<b>13</b>

**Table 5.19 Test model outcomes comparing CB, low-, intermediate- and high- risk cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction.**

	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>Accuracy</i>	<b>0.4611</b>	<b>0.45</b>	<b>0.4278</b>	<b>0.4444</b>	<b>0.4222</b>	<b>0.3944</b>	<b>0.3944</b>	<b>0.222</b>	<b>0.4056</b>	<b>0.4716</b>	<b>0.4659</b>	<b>0.4773</b>
<i>AUC</i>	<b>0.6894</b>	<b>0.6646</b>	<b>0.6115</b>	<b>0.6479</b>	<b>0.6522</b>	<b>0.6273</b>	<b>0.6372</b>	<b>0.4993</b>	<b>0.6468</b>	<b>0.6791</b>	<b>0.709</b>	<b>0.7088</b>
<i>Sensitivity:</i>												
<i>CB</i>	<b>37%</b>	<b>37%</b>	<b>33%</b>	<b>30%</b>	<b>43%</b>	<b>27%</b>	<b>47%</b>	<b>90%</b>	<b>53%</b>	<b>35%</b>	<b>35%</b>	<b>42%</b>
<i>L</i>	<b>28%</b>	<b>15%</b>	<b>9%</b>	<b>28%</b>	<b>15%</b>	<b>7%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>28%</b>	<b>35%</b>	<b>28%</b>
<i>I</i>	<b>79%</b>	<b>85%</b>	<b>88%</b>	<b>75%</b>	<b>76%</b>	<b>78%</b>	<b>68%</b>	<b>18%</b>	<b>72%</b>	<b>82%</b>	<b>76%</b>	<b>79%</b>
<i>H</i>	<b>6%</b>	<b>6%</b>	<b>0%</b>	<b>13%</b>	<b>31%</b>	<b>13%</b>	<b>25%</b>	<b>0%</b>	<b>16%</b>	<b>6%</b>	<b>6%</b>	<b>9%</b>
<i>Specificity:</i>												
<i>CB</i>	<b>97%</b>	<b>94%</b>	<b>95%</b>	<b>97%</b>	<b>91%</b>	<b>89%</b>	<b>85%</b>	<b>15%</b>	<b>81%</b>	<b>96%</b>	<b>94%</b>	<b>93%</b>
<i>L</i>	<b>84%</b>	<b>90%</b>	<b>93%</b>	<b>86%</b>	<b>90%</b>	<b>92%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>85%</b>	<b>83%</b>	<b>86%</b>
<i>I</i>	<b>42%</b>	<b>33%</b>	<b>23%</b>	<b>39%</b>	<b>39%</b>	<b>31%</b>	<b>42%</b>	<b>90%</b>	<b>40%</b>	<b>39%</b>	<b>46%</b>	<b>46%</b>
<i>H</i>	<b>95%</b>	<b>97%</b>	<b>97%</b>	<b>93%</b>	<b>93%</b>	<b>95%</b>	<b>84%</b>	<b>99%</b>	<b>91%</b>	<b>97%</b>	<b>95%</b>	<b>95%</b>
<i>PPV:</i>												
<i>CB</i>	<b>69%</b>	<b>55%</b>	<b>59%</b>	<b>64%</b>	<b>50%</b>	<b>33%</b>	<b>38%</b>	<b>18%</b>	<b>36%</b>	<b>60%</b>	<b>50%</b>	<b>50%</b>
<i>L</i>	<b>37%</b>	<b>35%</b>	<b>31%</b>	<b>41%</b>	<b>33%</b>	<b>21%</b>	<b>NA</b>	<b>NA</b>	<b>NA</b>	<b>41%</b>	<b>42%</b>	<b>42%</b>
<i>I</i>	<b>48%</b>	<b>46%</b>	<b>43%</b>	<b>45%</b>	<b>45%</b>	<b>43%</b>	<b>44%</b>	<b>54%</b>	<b>44%</b>	<b>48%</b>	<b>50%</b>	<b>50%</b>
<i>H</i>	<b>22%</b>	<b>29%</b>	<b>0%</b>	<b>29%</b>	<b>8%</b>	<b>33%</b>	<b>26%</b>	<b>0%</b>	<b>28%</b>	<b>29%</b>	<b>22%</b>	<b>30%</b>
<i>Number of Probes</i>	<b>36</b>	<b>13</b>	<b>5</b>	<b>12</b>	<b>37</b>	<b>14</b>	<b>78</b>	<b>39</b>	<b>12</b>	<b>37</b>	<b>34</b>	<b>13</b>

### 5.7.4 Models to predict patient type using trends in expression

Expression of 152, 57, 56, 45 probes had a significant association with increasing severity of disease type i.e. no evidence for cancer (CB), organ confined cancer (L, I, & H) and metastatic disease (A) in the four processed datasets (*KLK2* ratio, *KLK2* adjusted, *KLK3* adjusted, HK normalised respectively) (Supplementary Table 23). The top probe was *HOXC6* ( $p = 5.19 \times 10^{-10}$ ), *UPK2* ( $p = 2.91 \times 10^{-08}$ ), *UPK2* ( $p = 2.4 \times 10^{-08}$ ), and *HOXC6* ( $p = 3.39 \times 10^{-06}$ ) respectively.



**Figure 5.23 Top Significant Probe for CB, Cancer, Metastatic trend in all four data normalisations.**

Multivariate proportional odds models were built to predict clinical categories (section 2.6.1), no evidence for cancer (CB), organ confined cancer (L, I, & H) and metastatic disease (Table 5.20, Table 5.21). The probes involved in each model can be observed (Supplementary Table 24, Supplementary Table 25, Supplementary Table 26 and Supplementary Table 27).

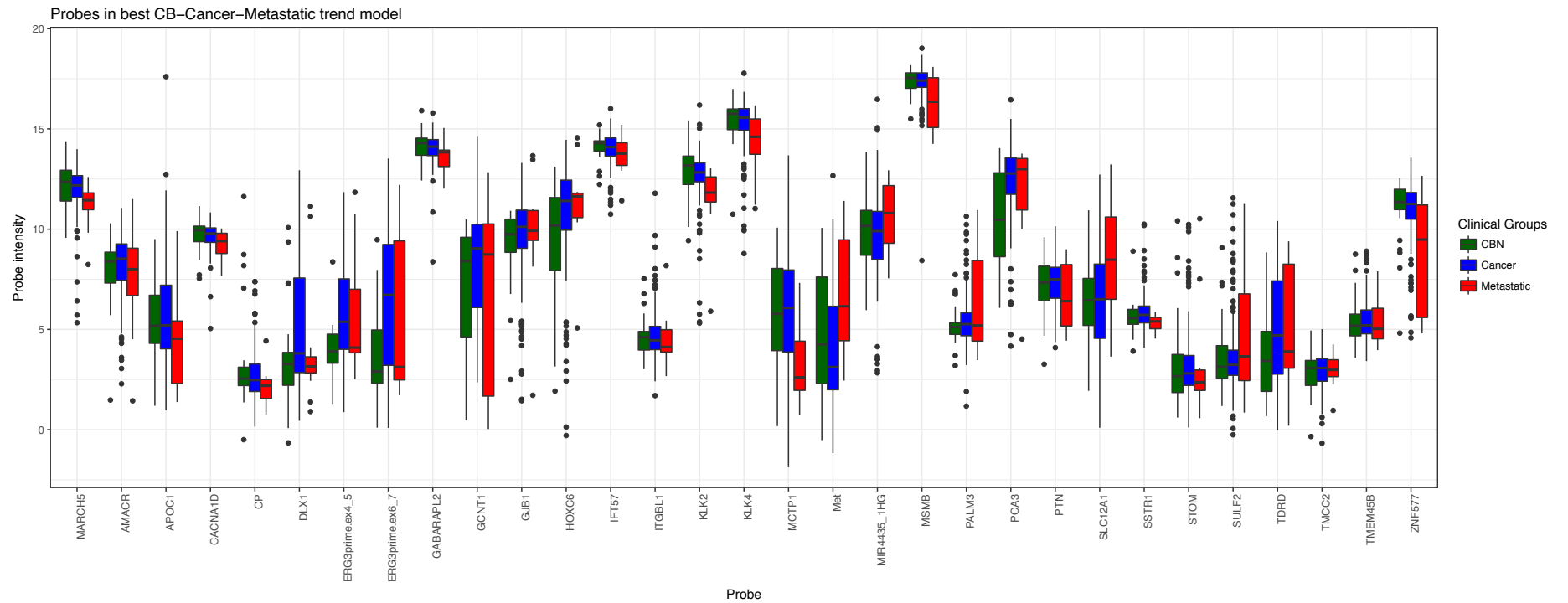
Low AUCs were observed across all inputs and data sets (median AUC = 0.57365, IQR = 0.08). The *GAPDH* and *RPLP2* normalised data showed slightly higher AUCs (median AUC = 0.6388, IQR = 0.997). The sensitivity of the sample categories in all of the models were fairly low (median = 18%, IQR = 87%). Whilst the specificity is fairly

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

high but not uniformly across the models (median = 98%, IQR = 75%). Inclusion of the advanced samples could be a reason for this poor model quality. Advanced tumours tend to be firm to the touch and it is thought that upon compression tend to release fewer cells into the urine (section 1.3.4.2). This is further supported by the lower levels of prostate specific transcripts observed in advanced samples (section 3.4) and *UPK2* (the bladder specific marker) is one of the most significant differential probes comparing these samples.

Again, the model with the best AUC (0.6469) is from the *GAPDH* and *RPLP2* (HK) normalised data. The model was built using the significant probes (*MARCH5*, *AMACR*, *APOC1*, *CACNA1D*, *CP*, *DLX1*, *ERG* 3' exons 4-5, *ERG* 3' exons 6-7, *GABARAPL2*, *GCNT1*, *GJB1*, *HOXC6*, *IFT57*, *ITGBL1*, *KLK2*, *KLK4*, *MCTP1*, *Met*, *MIR4435\_1HG*, *MSMB*, *PALM3*, *PCA3*, *PTN*, *SLC12A1*, *SSTR1*, *STOM*, *SULF2*, *TDRD*, *TMCC1*, *TMEM45B*, *ZNF577*). The model's sensitivity ranged from 17% - 93% and it's specificity ranged from 26%- 98%.

CHAPTER 5: NANOSTRING DATA ANALYSIS 2



**Figure 5.24** Boxplot showing the expression level of each transcript featured in the CB-Cancer-Metastatic cancer model built using the significant probes from the *GAPDH* and *RPLP2* normalised data. This model showed the best test data AUC (0.6469).

CHAPTER 5: NANOSTRING DATA ANALYSIS 2

**Table 5.20 Training model outcomes comparing CB, Cancer (low-, intermediate- and high- risk) and metastatic (A) cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction**

	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>Accuracy</i>	<b>0.8136</b>	<b>0.811</b>	<b>0.811</b>	<b>0.8114</b>	<b>0.8372</b>	<b>0.7855</b>	<b>0.781</b>	<b>0.8072</b>	<b>0.7353</b>	<b>0.8819</b>	<b>0.8504</b>	<b>0.8005</b>
<i>AUC</i>	<b>0.5554</b>	<b>0.5495</b>	<b>0.5495</b>	<b>0.5878</b>	<b>0.6566</b>	<b>0.5201</b>	<b>0.6267</b>	<b>0.7126</b>	<b>0.5913</b>	<b>0.7375</b>	<b>0.6685</b>	<b>0.541</b>
<i>Sensitivity:</i>												
<i>CB</i>	<b>4%</b>	<b>2%</b>	<b>2%</b>	<b>5%</b>	<b>24%</b>	<b>2%</b>	<b>23%</b>	<b>35%</b>	<b>6%</b>	<b>44%</b>	<b>37%</b>	<b>4%</b>
<i>Cancer</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>99%</b>	<b>99%</b>	<b>100%</b>	<b>98%</b>	<b>98%</b>	<b>100%</b>	<b>98%</b>	<b>98%</b>
<i>Metastatic</i>	<b>13%</b>	<b>13%</b>	<b>13%</b>	<b>22%</b>	<b>26%</b>	<b>4%</b>	<b>17%</b>	<b>35%</b>	<b>22%</b>	<b>35%</b>	<b>17%</b>	<b>9%</b>
<i>Specificity:</i>												
<i>CB</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>99%</b>	<b>100%</b>	<b>100%</b>	<b>99%</b>	<b>99%</b>
<i>Cancer</i>	<b>7%</b>	<b>5%</b>	<b>5%</b>	<b>10%</b>	<b>25%</b>	<b>2%</b>	<b>21%</b>	<b>35%</b>	<b>11%</b>	<b>41%</b>	<b>31%</b>	<b>5%</b>
<i>Metastatic</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>99%</b>	<b>100%</b>	<b>99%</b>	<b>99%</b>	<b>100%</b>	<b>99%</b>	<b>100%</b>
<i>PPV:</i>												
<i>CB</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>93%</b>	<b>100%</b>	<b>100%</b>	<b>92%</b>	<b>80%</b>	<b>96%</b>	<b>86%</b>	<b>33%</b>
<i>Cancer</i>	<b>81%</b>	<b>81%</b>	<b>81%</b>	<b>81%</b>	<b>83%</b>	<b>79%</b>	<b>77%</b>	<b>80%</b>	<b>74%</b>	<b>87%</b>	<b>85%</b>	<b>81%</b>
<i>Metastatic</i>	<b>75%</b>	<b>75%</b>	<b>75%</b>	<b>100%</b>	<b>86%</b>	<b>20%</b>	<b>100%</b>	<b>80%</b>	<b>56%</b>	<b>100%</b>	<b>67%</b>	<b>67%</b>
<i>Number of Probes</i>	<b>11</b>	<b>7</b>	<b>8</b>	<b>39</b>	<b>39</b>	<b>11</b>	<b>35</b>	<b>39</b>	<b>9</b>	<b>69</b>	<b>31</b>	<b>9</b>

CHAPTER 5: NANOSTRING DATA ANALYSIS 2

**Table 5.21 Test model outcomes comparing CB, Cancer (low-, intermediate- and high- risk) and metastatic (A) cancer samples for the four different normalisations of data. Three input probe sets were used: all probes, those significant via GLM testing and those significant post - multiple testing correction**

	<i>KLK2 ratio</i>			<i>KLK2 Adjusted</i>			<i>KLK3 Adjusted</i>			<i>GAPDH and RPLP2 normalised</i>		
	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>	<i>All probes</i>	<i>Significant probes</i>	<i>Adjusted Significant Probes</i>
<i>Accuracy</i>	<b>0.7865</b>	<b>0.7812</b>	<b>0.7812</b>	<b>0.7917</b>	<b>0.776</b>	<b>0.776</b>	<b>0.8021</b>	<b>0.3542</b>	<b>0.7656</b>	<b>0.7819</b>	<b>0.7926</b>	<b>0.7819</b>
<i>AUC</i>	<b>0.5111</b>	<b>0.5</b>	<b>0.5333</b>	<b>0.5595</b>	<b>0.5799</b>	<b>0.5657</b>	<b>0.5911</b>	<b>0.5778</b>	<b>0.5695</b>	<b>0.6307</b>	<b>0.6469</b>	<b>0.5</b>
<i>Sensitivity:</i>												
<i>CB</i>	<b>3%</b>	<b>0%</b>	<b>0%</b>	<b>10%</b>	<b>17%</b>	<b>0%</b>	<b>20%</b>	<b>90%</b>	<b>13%</b>	<b>27%</b>	<b>31%</b>	<b>0%</b>
<i>Cancer</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>99%</b>	<b>95%</b>	<b>99%</b>	<b>98%</b>	<b>28%</b>	<b>95%</b>	<b>91%</b>	<b>93%</b>	<b>98%</b>
<i>Metastatic</i>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>8%</b>	<b>8%</b>	<b>17%</b>	<b>8%</b>	<b>19%</b>	<b>8%</b>	<b>25%</b>	<b>17%</b>	<b>0%</b>
<i>Specificity:</i>												
<i>CB</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>99%</b>	<b>100%</b>	<b>99%</b>	<b>28%</b>	<b>96%</b>	<b>93%</b>	<b>96%</b>	<b>98%</b>
<i>Cancer</i>	<b>24%</b>	<b>0%</b>	<b>0%</b>	<b>10%</b>	<b>14%</b>	<b>2%</b>	<b>17%</b>	<b>83%</b>	<b>12%</b>	<b>29%</b>	<b>26%</b>	<b>0%</b>
<i>Metastatic</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>99%</b>	<b>97%</b>	<b>99%</b>	<b>99%</b>	<b>100%</b>	<b>99%</b>	<b>98%</b>	<b>98%</b>	<b>100%</b>
<i>PPV:</i>												
<i>CB</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>71%</b>	<b>NA</b>	<b>86%</b>	<b>19%</b>	<b>40%</b>	<b>39%</b>	<b>53%</b>	<b>0%</b>
<i>Cancer</i>	<b>79%</b>	<b>78%</b>	<b>78%</b>	<b>80%</b>	<b>80%</b>	<b>78%</b>	<b>81%</b>	<b>85%</b>	<b>79%</b>	<b>84%</b>	<b>83%</b>	<b>79%</b>
<i>Metastatic</i>	<b>NA</b>	<b>NA</b>	<b>NA</b>	<b>33%</b>	<b>17%</b>	<b>3%</b>	<b>33%</b>	<b>NA</b>	<b>33%</b>	<b>50%</b>	<b>33%</b>	<b>NA</b>
<i>Number of Probes</i>	<b>11</b>	<b>7</b>	<b>8</b>	<b>39</b>	<b>39</b>	<b>11</b>	<b>35</b>	<b>39</b>	<b>9</b>	<b>69</b>	<b>31</b>	<b>9</b>

### 5.7.5 Conclusions

Use of the housekeeping probes *GAPDH* and *RPLP2* provided normalised data that produced good prediction models (this data provided the best AUC for prediction models for all four clinical questions (Table 5.11)). Data was otherwise treated similarly to NanoString1 (chapter 3). Identification of *GAPDH* and *RPLP2* as housekeepers to normalise urinary EV RNA derived NanoString data increased the robustness of my prediction models.

All models were built using a training set that included samples from all four centres, and particularly the binomial tests were robust (high AUCs). The models therefore, can predict cancer from samples with no evidence of cancer (CB) regardless of sample origin.

Optimal models built from the expression of 167 markers for risk stratification and detection of cancer were found using the *GAPDH* and *RPLP2* normalised data, however, input lists varied from all probes, significant probes (identified by polr) and adjusted significant probes (Benjamin Hochberg multiple testing correction).

The Prostate Cancer Prevention Trial risk calculator (PCPTrc) and the Prostate Cancer Prevention Trial high-grade risk calculator (PCPThg) are logistic regression models, which incorporate PSA level, PSA velocity, DRE result, previous biopsy results, age at biopsy, race and family history of PCa<sup>212</sup>. These models have been combined with urinary (whole cell) *TMPRSS2:ERG* and urine *PCA3* levels to improve model AUC: PCPTrc alone had an AUC of 0.639, whilst inclusion of urinary *TMPRSS2:ERG* and *PCA3* improved the AUC to 0.762. Urinary *TMPRSS2:ERG* and *PCA3* also improved the predictive power of serum PSA (AUC = 0.651 increased to AUC = 0.772)<sup>213</sup>. Similarly, *PCA3*, which is used in the *PCA3* test, which was the first commercially available urinary test for PCa, is capable of predicting cancer from non-cancer samples (AUC = 0.98)<sup>214</sup>. The models achieved similar AUCs, when predicting cancer (L, I and H) from samples with no evidence of cancer (CB): AUC = 0.851 for the best model and



## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

median AUC = 0.8045. I found that in EV harvested material the *ERG3'* exons 4-5 and *PCA3* data were the most highly differentiating between cancer and samples with no evidence of cancer. The probes used in the top model were *TMPRSS2:ERG*, *ERG3'* exons 4-5, *APOC1*, *ISX*, *SLC12A1*, *HOXC6*, *MCTP1*, *TDRD*, *PDLIM5*, *CD10*, *GABARAPL2*, *PTN*, *AR* exon 9, *PPP1R12B*, *CP*, *MXII*, and *KLK4*.

The high-grade predictor also benefitted from the addition of urinary *TMPRSS2:ERG* and *PCA3* data (AUC = 0.707 increased to AUC = 0.779)<sup>213</sup>. A second high-grade predictor was produced by Van Neste *et al.*, which used whole urine mRNA levels of *HOXC4*, *HOXC6*, *TDRD1*, *DLX1* and *PCA3* (with *KLK3* as a reference) alongside clinical factors (including PSA density, previous biopsies, PSA, age and family history)<sup>215</sup>. This model reached an overall AUC of 0.9 in their validation set. The high-risk (H) Vs. no evidence of cancer (CB) models also achieved high AUCs (top AUC = 0.897, median AUC = 0.831). The model also used *PCA3* and *TMPRSS2:ERG* levels, along with *APOC1*, *HPN* and *ERG3'* exons 4-5 (from EV harvested RNA). The top most significant probes when comparing high-risk cancer with samples with no evidence of cancer was *HPN* and *ERG 3'* exons 4-5.

The ExoDx Prostate IntelliScore urine exosome assay uses *ERG* and *PCA3* data normalised using *SPDEF* combined with clinical factors (including PSA level, age, race and family history) to predict between Gleason 6 and Gleason 7 PCa with AUC = 0.73<sup>216</sup>. The models to predict between different risk categories (CB->L->I->H) had similar AUCs (median = 0.67015, highest AUC = 0.709). The model was built using *APOC1*, *DPP4*, *ERG 3'* exons 4-5, *ERG 3'* exons 6-7, *GABARAPL2*, *HOXC6*, *HPN*, *ITGBL1*, *KLK4*, *MYOF*, *PCA3*, *TDRD*, and *TMPRSS2:ERG*. It is not surprising that *PCA3* and *ERG3'* exons 4-5 were also the most highly significant in all four data normalisations.

I have shown that EV derived material from multiple centres can be quantified by NanoString to produce models that can predict cancer presence and aggression without biopsy. However, much greater numbers and model refinements (such as RF etc.)

## CHAPTER 5: NANOSTRING DATA ANALYSIS 2

would be needed to strengthen the models into a test that could be used in the clinic. A multivariate regression with the combination of RNA signatures and clinical factors should also be investigated.

# 6

## **Expression Profile of the Cell Sediment Urine Fraction**

### **6.1 Summary**

In this chapter, I compared the transcriptome profiles of two urine fractions from prostate cancer patients and controls, and examined whether the transcriptomes from cell sediment were better than EV transcriptomes for PCa diagnosis. I found that the cell sediments have a very different transcriptome profile to the EV fractions, which is similar to what was found in renal cancer<sup>217</sup>. Transcripts found by microarray analysis to be significantly more abundant in the EV fraction compared to the cell sediment were more commonly expressed in prostatic tissue and also had more known associations with prostate cancer. This suggested that the majority of RNA within the extracellular vesicle fraction comes from prostatic tissue, both normal and cancerous. These analyses support the hypothesis that EVs are a better fraction to study for biomarkers in prostate cancer.

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

Analysis of cell sediment NanoString data to identify transcripts that could be used diagnostically to identify D'Amico clinically categories show common transcripts being selected by both Lasso and Random forest analyses when i) different input subsets of transcripts were used, and ii) when data was normalised with different control genes. Different probes for *ERG* 3' sequences were identified in cell (*ERG* probe targeting exons 6-7) and EV models (*ERG* probe targeting exons 4-5). *HOXC6* and *TDRD* were found in both cell and EV models. Interestingly, *PCA3* and *TMPRSS2:ERG* were found in EV models and not cell models. This supports other work that the majority of PCa RNA content in whole urine is originating from EVs and not whole cells.

### **6.2 Introduction**

NanoString technology was applied to cell and extracellular vesicle (EV) fractions of urine from prostate cancer patients to form the NanoString 2 data set. Urine samples were divided into two fractions by centrifugation: i) cell sediment and ii) supernatant containing extracellular vesicles (section 2.1.2). In this chapter, analysis of the cell fraction will be completed. The investigation of the EV fraction can be found in Chapter 5.

#### **6.2.1 The Research Gap**

Since the production of the PCA3 test<sup>214</sup>, urine has been investigated for PCa biomarkers. Whole urine and cell sediment are commonly used and many models have been developed or built upon to include urinary expression of transcripts as biomarkers<sup>213,215</sup>. However, little work has been done on the EV fraction. The EV fraction has been identified to be a useful source of biomarkers in renal cancer<sup>217</sup> and PCa associated transcripts have been quantified from PCa urine EVs<sup>218</sup>. No comparisons between transcript expression levels in EV fractions and cell fractions

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

have been published. It is therefore, unknown which may be the better source of PCa biomarkers.

### **6.2.2 The Aims**

In the first part of this chapter, I will examine the NanoString data for differences in the expression profiles of the cell sediment between different clinical categories and try to construct models to predict clinical categories: These comparisons include comparing i) CB (no evidence of cancer) samples with D'Amico cancer risk groups: Low, Intermediate and High, and ii) CB vs high-risk cancer samples. Two trends will be investigated, CB, Low-, Intermediate-, High-risk cancers, ordered as such and CB, cancer and metastatic samples; ordered as such. In each comparison I have used two methods of analysis (logistic regression analysis and Mann Whitney U test), and will compare and contrast the selected gene transcripts from each. These investigations have already been presented for the extracellular vesicle fraction (chapter 5).

In the second part of this chapter, I will compare and contrast the matched EV and cell sediment fraction data from microarray and NanoString analyses. Other studies have observed that the transcriptomes of urinary extracellular vesicles and whole urine are different in renal cancer<sup>27</sup>. I will identify transcripts that are significantly differentially expressed between the cell sediment and EVs in both NanoString data and microarray data.

### **6.2.3 The Data**

The cell and EV fractions were analysed in 95 samples from a range of clinical categories (Table 6.1) based on the D'Amico classification using 167 NanoString probes. Three of these samples were taken pre-DRE, and as shown previously, these samples are not fully comparable with those obtained post-DRE and were not used in this chapter. These data were normalised with the spiked in positive controls as per the

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

NanoString manual (section 2.3.1) and  $\log_2$  transformed (section 2.3.3) to produce the baseline normalised data. Investigations into use of housekeeper transcripts were completed and data was normalised (section 2.3.2) using *RPLP2* and *TWIST1*. These probes showed no association to clinical categories ( $p > 0.05$  ANOVA-Tukey test) and were heavily correlated ( $p < 2.2 \times 10^{-16}$ ,  $r = 0.83$ ). Then similarly to the EV fraction *KLK2* ratio normalisation (section 2.1.1) was also performed. PCA plots (not shown) were used to visualise the *RPLP2* and *TWIST1* normalised (“HK normalised data”) and the *KLK2* ratio data. There were two outlier samples; M\_86\_1 (an Intermediate – risk sample) and M\_147\_1 (a CB sample), which were removed from the HK normalised data. M\_147\_1 was also removed from the *KLK2* ratio normalised data, as forty-six of the one hundred and sixty-seven values for M\_147\_1 were zero. The values for M\_88\_5 looked normal in the *KLK2* ratio data. Four clinical questions (Table 6.2) were investigated in the data and prediction models were produced accordingly. Due to limited numbers of samples, the data was not divided into test and training data.

**Table 6.1 Clinical breakdown of cell sediment fraction samples subjected to NanoString (within the second NanoString set). Twelve samples were CB (no evidence of cancer). Thirty raised PSA samples were negative for PCa on biopsy, but other abnormalities were found such as, HGPIN, prostatitis and atypia. Forty-six had localised cancer on TRUS biopsy of which four were D’Amico graded as Low risk, twenty-eight Intermediate risk and fourteen High-risk. Four samples had shown signs of metastasis.**

	<i>CB</i>	<i>Abnormal</i>	<i>L</i>	<i>I</i>	<i>H</i>	<i>A</i>	<i>Total</i>
<i>Number of Samples</i>	<b>12</b>	<b>30</b>	<b>4</b>	<b>28</b>	<b>14</b>	<b>4</b>	<b>92</b>
<i>Percentage</i>	<b>13%</b>	<b>33%</b>	<b>4%</b>	<b>30%</b>	<b>15%</b>	<b>4%</b>	<b>100%</b>
<i>Median Age</i>	<b>65</b>	<b>66</b>	<b>63.5</b>	<b>71</b>	<b>67.5</b>	<b>82</b>	<b>68</b>
<i>Median PSA</i>	<b>0.9</b>	<b>7.9</b>	<b>6.4</b>	<b>7.8</b>	<b>16.8</b>	<b>377</b>	<b>8.1</b>

**Table 6.2 Clinical predictive models built using the cell dataset.**

<i>Model</i>	<i>Samples</i>	<i>Model type</i>
<i>CB vs. Cancer</i>	<b><i>Clinically benign samples Vs low-, intermediate-, and high-risk cancer samples grouped together</i></b>	<b><i>Binary</i></b>
<i>CB vs. High risk cancer</i>	<b><i>Clinically benign Vs. high-risk cancer (extreme ends of no evidence of cancer and</i></b>	<b><i>Binary</i></b>

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

	<i>those with higher grade)</i>	
<i>CB, low-, intermediate-, and high-risk trend</i>	<i>Each sample category is a separate group and ordered</i>	<i>Ordinal</i>
<i>CB, cancer, metastatic cancer trend</i>	<i>Clinically benign samples, with low-, intermediate-, and high-risk cancer samples grouped together, and metastatic cancer samples in groups ordered by severity</i>	<i>Ordinal</i>

**6.3 Models predicting presence of cancer CB and cancer (L, I, H) samples using cell sediment data**

***6.3.1.1 Differentially expressed transcripts***

Expression of 85, 28, and 24 transcripts had a significant association (via logistic regression section 2.6.1) with whether a sample had no evidence for cancer (CB) or not (L, I, H) in the three processed datasets (the baseline data, *KLK2* ratio, and HK normalised, respectively) (Supplementary Table 31). Only *MCTP1* remained significant post multiple testing correction (adjusted  $p = 0.04$ ) in the baseline data and none remained significant in the *KLK2* ratio and HK normalised data. The top significant probe in these datasets was *ERG* 3' exons 6-7 ( $p = 0.001$ ) and *NAALADL2* ( $p = 3.33 \times 10^{-05}$ ), respectively.

Expression of 94, 33, and 56 transcripts had a significant association (via Mann Whitney U (MWU) testing, section 2.4.1) with whether a sample had no evidence for cancer (CB) or not (L, I, H) in the three processed datasets (the baseline data, *KLK2* ratio, and HK normalised, respectively) (Supplementary Table 32). The top significant probes identified by MWU were *SULF2* ( $p = 9.18 \times 10^{-06}$ ), *PCA3* ( $p = 3.72 \times 10^{-05}$ ) and *SPINK1* ( $p = 3.72 \times 10^{-05}$ ), respectively.

Between the two tests 79, 20 and 26 transcripts were common between the two methods suggesting a good level of robustness. The top ten transcripts with the biggest log2 fold

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

change for the baseline data, *KLK2* ratio, and HK normalised data are shown (Table 6.3, Table 6.4, Table 6.5, respectively).



CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

Table 6.3 Top ten transcripts with biggest log2 fold change in the baseline normalised data.

Transcript	MWU		glm		Log <sub>2</sub> (FC)
	p-value	Adjusted p-value	p-value	Adjusted p-value	
HOXC6	0.0002	0.024	0.0014	0.2049	1.64
ERG3' exons 6-7	2.84x10 <sup>-07</sup>	4.74x10 <sup>-05</sup>	0.0008	0.128	1.38
TMPRSS2:ERG	4.52x10 <sup>-05</sup>	0.0069	0.0013	0.1979	1.31
SLC43A1	0.0003	0.0406	0.0019	0.2745	1.17
CLIC2	2.66x10 <sup>-05</sup>	0.0042	0.001	0.1645	1.05
B4GALNT4	3.38x10 <sup>-05</sup>	0.0053	0.0012	0.1807	1.04
CADPS	1.37x10 <sup>-05</sup>	0.0022	0.0004	0.0682	1.04
CKAP2L	0.0116	1	0.0033	0.4318	1.01
HPN	7.04x10 <sup>-05</sup>	0.0103	0.0006	0.1041	0.97
LASS1	0.0002	0.022	0.0011	0.1703	0.97

Table 6.4 Top ten transcripts with biggest log2 fold change in the KLK2 ratio data.

Transcript	MWU		glm		Log <sub>2</sub> (FC)
	p-value	Adjusted p-value	p-value	Adjusted p-value	
HOXC6	6.80x10 <sup>-05</sup>	0.01	0.004	0.63	0.21
ERG3' exons 6-7	7.80x10 <sup>-05</sup>	0.01	0.001	0.24	0.18
TDRD	0.0004	0.06	0.004	0.72	0.18
SLC43A1	0.002	0.32			0.17
CADPS	0.004	0.67	0.01	1	0.16
ERG5'	0.01	0.99			0.15
B4GALNT4	0.01	0.87			0.14
SLC12A1	0.003	0.54	0.03	1	0.13
TMCC2	0.05	0.99	0.05	1	0.13
TMPRSS2:ERG	0.001	0.17	0.01	1	0.13

Table 6.5 Top ten transcripts with biggest log2 fold change in the HK normalised data.

Transcript	MWU		glm		Log <sub>2</sub> (FC)
	p-value	Adjusted p-value	p-value	Adjusted p-value	
HOXC6	0.0002	0.0374	0.0019	0.3087	1.5
ERG3' exons 6-7	0.0006	0.1045	0.0228	0.9861	1.1
TMPRSS2:ERG	0.0036	0.5527	0.0069	0.9861	1.1
CP	0.0146	0.9924	0.0109	0.9861	-1
TDRD	0.001	0.153	0.0105	0.9861	0.9
NAALADL2	3.33x10 <sup>-05</sup>	0.0056	0.0012	0.2012	-0.8
SLC43A1	0.0005	0.0895	0.0168	0.9861	0.8
ST6GALNAC1	0.0008	0.1311	0.0238	0.9861	-0.8
SPINK1	7.80x10 <sup>-05</sup>	0.0129			-0.7
UPK2	0.0007	0.1128	0.0026	0.4313	-0.7

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

### ***6.3.1.2 Models and gene selection***

A number of different transcript subsets were input to Lasso for probe shrinkage and selection, these included i) all of the transcripts ( $n = 167$ ). Transcripts identified as having significantly different expression between cancer and CB using ii) Mann Whitney U ( $n = 94$ ,  $n = 33$  and  $n = 56$ ) and iii) logistic regression ( $n = 85$ ,  $n = 28$  and  $n = 24$ ), separately, and iv) transcripts common to both those identified by Mann Whitney U and logistic regression ( $n = 79$ ,  $n = 20$  and  $n = 26$ ) for each of the three normalisations (the baseline data, *KLK2* ratio, and HK normalised), respectively. The AUC, sensitivity and specificity of each model on the same data was collected (Table 6.6) and transcript lists (Table 6.7, Table 6.8, Table 6.9) and boxplots of the Lasso selected probes were produced (Supplementary Figure 13, Supplementary Figure 14 and Supplementary Figure 15).

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.6 AUC, Sensitivity and Specificity of models to predict CB vs. Cancer (L, I, H) in different data normalisations of cell NanoString data.**

	<i>HK normalised</i>				<i>KLK2 ratio</i>				<i>Baseline</i>			
	<i>All Transcripts</i>	<i>MWU</i>	<i>glm</i>	<i>Both</i>	<i>All Transcripts</i>	<i>MWU</i>	<i>glm</i>	<i>Both</i>	<i>All Transcripts</i>	<i>MWU</i>	<i>glm</i>	<i>Both</i>
<i>AUC</i>	<b>0.989</b>	<b>0.998</b>	<b>0.998</b>	<b>0.998</b>	<b>0.996</b>	<b>0.998</b>	<b>0.993</b>	<b>0.995</b>	<b>0.998</b>	<b>1</b>	<b>1</b>	<b>1</b>
<i>Sensitivity</i>	<b>100%</b>	<b>98%</b>	<b>98%</b>	<b>98%</b>	<b>98%</b>	<b>98%</b>	<b>94%</b>	<b>96%</b>	<b>98%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>
<i>Specificity</i>	<b>92%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>
<i>Number of Probes</i>	<b>8*</b>	<b>8*</b>	<b>8*</b>	<b>8*</b>	<b>9</b>	<b>7</b>	<b>7**</b>	<b>7**</b>	<b>13</b>	<b>17</b>	<b>14***</b>	<b>14***</b>

\*, \*\* and \*\*\* selected probes are identical in model.

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.7** Beta values of individual transcripts within models suggested by Lasso using different input transcripts for the baseline normalised data.

<i>Transcript</i>	<i>Beta – All Transcripts</i>	<i>Beta - MWU</i>	<i>Beta - glm</i>	<i>Beta - both</i>
<i>ACTR5</i>		<i>0.572</i>	<i>0.206</i>	<i>0.226</i>
<i>APOC1</i>	<i>0.045</i>	<i>0.022</i>	<i>0.054</i>	<i>0.052</i>
<i>ARHGEF25</i>		<i>-0.176</i>	<i>0.396</i>	
<i>CADPS</i>	<i>0.273</i>	<i>0.481</i>		<i>0.399</i>
<i>CAMKK2</i>	<i>0.055</i>			
<i>ERG 3' exons 6-7</i>	<i>0.082</i>	<i>0.203</i>	<i>0.135</i>	<i>0.137</i>
<i>EN2</i>		<i>0.32</i>	<i>0.146</i>	<i>0.164</i>
<i>HIST1H2BG</i>	<i>0.006</i>		<i>0.015</i>	<i>0.013</i>
<i>HOXC6</i>	<i>0.096</i>	<i>0.138</i>	<i>0.114</i>	<i>0.116</i>
<i>IGFBP3</i>		<i>-0.148</i>		
<i>LASS1</i>	<i>0.115</i>	<i>0.314</i>	<i>0.26</i>	<i>0.263</i>
<i>MCTP1</i>	<i>0.159</i>			
<i>MMP25</i>	<i>0.042</i>	<i>0.3470</i>	<i>0.219</i>	<i>0.224</i>
<i>MMP26</i>	<i>-0.124</i>	<i>-0.137</i>		
<i>NAALADL2</i>		<i>-0.515</i>	<i>-0.356</i>	<i>-0.371</i>
<i>PCA3</i>	<i>0.019</i>	<i>0.084</i>	<i>0.076</i>	<i>0.078</i>
<i>RIOK3</i>	<i>0.095</i>	<i>0.0290</i>	<i>0.012</i>	<i>0.003</i>
<i>SPINK1</i>	<i>-0.05</i>	<i>-0.0220</i>	<i>-0.056</i>	<i>-0.058</i>
<i>SLC12A1</i>		<i>0.1020</i>		
<i>TDRD</i>		<i>0.1260</i>	<i>0.041</i>	<i>0.044</i>

**Table 6.8** Beta values of individual transcripts within models suggested by Lasso using different input transcripts for the *KLK2* ratio data.

<i>Transcript</i>	<i>Beta – All transcripts</i>	<i>Beta - MWU</i>	<i>Beta - glm</i>	<i>Beta - Both</i>
<i>CADPS</i>	<i>0.075</i>	<i>0.1795</i>	<i>0.1052</i>	<i>0.1691</i>
<i>CKAP2L</i>	<i>0.1992</i>	<i>0.2766</i>	<i>0.1894</i>	<i>0.2467</i>
<i>EN2</i>	<i>0.0828</i>			
<i>ERG 3' exons 6-7</i>	<i>0.7197</i>	<i>0.7699</i>	<i>0.7411</i>	<i>0.8396</i>
<i>HOXC6</i>	<i>0.3855</i>	<i>0.6533</i>	<i>0.3772</i>	<i>0.5718</i>
<i>MFSD2A</i>	<i>0.0104</i>			
<i>NAALADL2</i>	<i>-1.456</i>	<i>-2.084</i>	<i>-1.3254</i>	<i>-1.994</i>
<i>SFRP4</i>	<i>0.1328</i>	<i>0.228</i>		
<i>SIM2 long</i>			<i>0.0251</i>	<i>0.2771</i>
<i>TDRD</i>	<i>0.157</i>	<i>0.3875</i>	<i>0.1587</i>	<i>0.258</i>

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.9 Beta values of individual transcripts within models suggested by Lasso using different input transcripts for the HK normalised data.**

<i>Transcript</i>	<i>Beta – All transcripts</i>	<i>Beta - MWU</i>	<i>Beta - glm</i>	<i>Beta - Both</i>
<i>CADPS</i>	<b>0.1096</b>	<b>0.1433</b>	<b>0.1971</b>	<b>0.2285</b>
<i>CLIC2</i>	<b>0.1062</b>	<b>0.1143</b>	<b>0.1241</b>	<b>0.128</b>
<i>ERG 3' exons 6-7</i>	<b>0.0716</b>	<b>0.0819</b>	<b>0.0967</b>	<b>0.1051</b>
<i>HOXC6</i>	<b>0.1962</b>	<b>0.2044</b>	<b>0.2182</b>	<b>0.2272</b>
<i>NAALADL2</i>	<b>-0.287</b>	<b>-0.315</b>	<b>-0.353</b>	<b>-0.3719</b>
<i>SIM2 long</i>	<b>0.0411</b>	<b>0.0536</b>	<b>0.0602</b>	<b>0.0575</b>
<i>TDRD</i>	<b>0.0502</b>	<b>0.064</b>	<b>0.088</b>	<b>0.1031</b>
<i>UPK2</i>	<b>-0.081</b>	<b>-0.073</b>	<b>-0.061</b>	<b>-0.055</b>

Random forest was also applied to i) all transcripts, ii) significant transcripts identified by MWU and iii) significant transcripts identified by glm for the three different normalisations (the baseline data, *KLK2* ratio, and HK normalised), respectively (Supplementary Table 34, Supplementary Table 35 and Supplementary Table 36). The random forest model with the least error (Table 6.11) was built using the glm identified significant probes from the *KLK2* ratio data (the mean square of residuals = 0.088). *ERG 3' exon 6-7* was in the top 5 transcripts in 8/9 random forests, whilst *APOC1* was in the top 5 transcripts in 6/9 random forests. *HOXC6*, *CADPS*, *RIOK3*, *TMPRSS2:ERG*, *SLC12A1* and *SPINK1* occur in 3/9 random forests (Table 6.10).

**Table 6.10 Frequency of transcripts in top 5 for random forests.**

<i>Transcript</i>	<i>Frequency in top 5 random forest important transcripts</i>	<i>Data</i>
<i>APOC1</i>	<b>6</b>	<b>Baseline, <i>KLK2</i> and HK</b>
<i>CADPS</i>	<b>3</b>	<b>HK</b>
<i>CCDC88B</i>	<b>2</b>	<b>HK</b>
<i>ERG 3' exons 6-7</i>	<b>8</b>	<b>Baseline, <i>KLK2</i> and HK</b>
<i>NEAT1</i>	<b>2</b>	<b>Baseline</b>
<i>RIOK3</i>	<b>3</b>	<b>Baseline</b>
<i>TMPRSS2:ERG</i>	<b>3</b>	<b>Baseline</b>
<i>SLC12A1</i>	<b>3</b>	<b><i>KLK2</i></b>
<i>SPINK1</i>	<b>3</b>	<b>Baseline + HK</b>
<i>HOXC6</i>	<b>4</b>	<b><i>KLK2</i> + HK</b>
<i>PCA3</i>	<b>2</b>	<b><i>KLK2</i></b>

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.11 Mean square of residuals for random forest models for predicting CB vs. cancer (L, I and H) samples using different input probes across three different normalisations.**

	<i>HK</i>			<i>KLK2</i>			<i>Baseline</i>		
<i>Input:</i>	<i>All</i>	<i>glm</i>	<i>MWU</i>	<i>All</i>	<i>glm</i>	<i>MWU</i>	<i>All</i>	<i>glm</i>	<i>MWU</i>
<i>Mean square of residuals</i>	<b>0.11</b>	<b>0.11</b>	<b>0.138</b>	<b>0.115</b>	<b>0.08</b>	<b>0.104</b>	<b>0.106</b>	<b>0.099</b>	<b>0.103</b>
	4	8			8				

### 6.3.2 CB vs High risk cancer patients

#### 6.3.2.1 Differentially expressed transcripts

The 12 samples with no evidence of cancer (CB) were compared to the 14 high-risk cancer samples (H) using glm and MWU tests. 51, 12, and 20 transcripts had a significant association (via logistic regression, section 2.6.1) with whether a sample had no evidence for cancer (CB) or was high-risk cancer (H) in the three processed datasets (the baseline data, *KLK2* ratio, and HK normalised, respectively, Table 6.12, Table 6.13, Table 6.14, Supplementary Table 37). None remained significant post multiple testing correction in the baseline, *KLK2* ratio or the HK normalised data. The top significant probe in these datasets was *NEATI* ( $p = 0.004$ ), *ERG 3' exons 6-7* ( $p = 0.008$ ), and *HOXC6* ( $p = 0.005$ ), respectively.

Expression of 65, 25, and 35 transcripts had a significant association (via Mann Whitney U (MWU) testing, section 2.4.1) with whether a sample had no evidence for cancer (CB) or if the samples were high-risk cancer (H) in the three processed datasets (the baseline data, *KLK2* ratio, and HK normalised, respectively, Table 6.12, Table 6.13, Table 6.14, Supplementary Table 38). Post multiple testing correction, the expression of 10 (*ERG 3' exons 6-7*, *BAGALNT4*, *RIOK3*, *CADPS*, *MCTP1*, *HOXC6*, *NEATI*, *CLIC2*, *APOC1* and *SIM2* long), 1 (*HOXC6*) and 0, remained significant. The top significant probes identified by MWU were *ERG 3' exons 6-7* ( $p = 6.21 \times 10^{-06}$ ), *HOXC6* ( $p = 4.28 \times 10^{-05}$  and adjusted  $p = 0.007$ ) and *HOXC6* ( $p = 0.0005$ ), respectively.

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

Table 6.12 Top 10 transcripts with biggest log<sub>2</sub> fold change between CB and HR-cancer in the baseline data.

<i>Transcript</i>	<i>glm</i>		<i>MWU</i>		<i>Log<sub>2</sub>(FC)</i>
	<i>p - value</i>	<i>adjusted p - value</i>	<i>p - value</i>	<i>adjusted p - value</i>	
<i>HOXC6</i>	0.0002	0.0299	0.004	0.6711	2
<i>ERG3' exons 6-7</i>	6.21x10 <sup>-06</sup>	0.001	0.0371	0.9942	1.6
<i>TDRD</i>	0.0011	0.1558	0.0333	0.9942	1.5
<i>TMPRSS2:ERG</i>	0.0004	0.0668	0.0386	0.9942	1.3
<i>B4GALNT4</i>	2.88x10 <sup>-05</sup>	0.0048	0.0409	0.9942	1.2
<i>SLC43A1</i>	0.002	0.2897	0.0117	0.9942	1.2
<i>CADPS</i>	6.70x10 <sup>-05</sup>	0.011	0.02	0.9942	1.1
<i>CLIC2</i>	0.0002	0.0386	0.0087	0.9942	1
<i>HPN</i>	0.0008	0.1258	0.0092	0.9942	0.9
<i>LASS1</i>	0.0011	0.1558	0.0103	0.9942	0.9

Table 6.13 Top 10 transcripts with biggest log<sub>2</sub> fold change between CB and HR-cancer in the *KLK2* ratio data.

<i>Transcript</i>	<i>glm</i>		<i>MWU</i>		<i>Log<sub>2</sub>(FC)</i>
	<i>p - value</i>	<i>adjusted p - value</i>	<i>p - value</i>	<i>adjusted p - value</i>	
<i>TMPRSS2:ERG</i>	0.004	0.68	0.028	1.000	0.25
<i>ERG 3' exons 6-7</i>	0.000	0.07	0.008	1.000	0.25
<i>HOXC6</i>	4.28E-05	0.01			0.25
<i>TDRD</i>	0.001	0.09	0.017	1.000	0.24
<i>SLC43A1</i>	0.002	0.27	0.022	1.000	0.21
<i>CADPS</i>	0.007	1			0.18
<i>B4GALNT4</i>	0.002	0.33	0.035	1.000	0.17
<i>ERG 5'</i>	0.027	1			0.16
<i>SLC12A1</i>	0.013	1			0.15
<i>ERG 3' exons 4-5</i>	0.046	1	0.050	1.000	0.14

Table 6.14 Top 10 transcripts with biggest log<sub>2</sub> fold change between CB and HR-cancer in the HK normalised data.

<i>Transcript</i>	<i>glm</i>		<i>MWU</i>		<i>Log<sub>2</sub>(FC)</i>
	<i>p - value</i>	<i>adjusted p - value</i>	<i>p - value</i>	<i>adjusted p - value</i>	
<i>HOXC6</i>	0.0005	0.0882	0.0059	0.9765	1.6
<i>ERG3' exons 6-7</i>	0.0013	0.2186	0.0266	0.9765	1.4
<i>TDRD</i>	0.0031	0.4948	0.0272	0.9765	1.1
<i>TMPRSS2:ERG</i>	0.0094	1	0.033	0.9765	1.1
<i>ST6GALNAC1</i>	0.0037	0.5969	0.0168	0.9765	-1
<i>SLC43A1</i>	0.0013	0.2186	0.0197	0.9765	0.9
<i>B4GALNT4</i>	0.0202	1			0.8
<i>HPN</i>	0.0077	1	0.0314	0.9765	0.8
<i>CADPS</i>	0.0145	1	0.0326	0.9765	0.7
<i>CCDC88B</i>	0.031	1	0.0482	0.9765	0.7

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

### **6.3.2.2 Models and gene selection**

A number of different transcript subsets were input to Lasso for probe shrinkage and selection, these included i) all of the transcripts ( $n = 167$ ). Transcripts identified as having significantly different expression between cancer and CB using ii) Mann Whitney U ( $n = 65$ ,  $n = 25$  and  $n = 35$ ) and iii) logistic regression ( $n = 51$ ,  $n = 12$ , and  $n = 20$ ), separately, and iv) transcripts common to both those identified by Mann Whitney U and logistic regression ( $n = 49$ ,  $n = 12$  and  $n = 20$ ) for each of the three normalisations (the baseline data, *KLK2* ratio, and HK normalised), respectively. The AUC, sensitivity and specificity of each model on the same training data was collected (Table 6.15) and transcript lists (Table 6.16, Table 6.17, Table 6.18) and boxplots of the Lasso selected probes were produced (section 2.6.1, Supplementary Figure 16, Supplementary Figure 17 and Supplementary Figure 18).



CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.15 AUC, Sensitivity and Specificity of models to predict CB vs. high-risk cancer (H) in different data normalisations of cell NanoString data.**

	<i>HK normalised</i>				<i>KLK2 ratio</i>				<i>Baseline</i>			
	<i>All Transcripts</i>	<i>MWU</i>	<i>glm</i>	<i>Both</i>	<i>All Transcripts</i>	<i>MWU</i>	<i>glm</i>	<i>Both</i>	<i>All Transcripts</i>	<i>MWU</i>	<i>glm</i>	<i>Both</i>
<i>AUC</i>	<i>1</i>	<i>1</i>	<i>1</i>	<i>1</i>	<i>0.952</i>	<i>0.905</i>	<i>0.958</i>	<i>0.905</i>	<i>1</i>	<i>1</i>	<i>1</i>	<i>1</i>
<i>Sensitivity</i>	<i>100%</i>	<i>100%</i>	<i>100%</i>	<i>100%</i>	<i>93%</i>	<i>86%</i>	<i>93%</i>	<i>86%</i>	<i>100%</i>	<i>100%</i>	<i>100%</i>	<i>100%</i>
<i>Specificity</i>	<i>100%</i>	<i>100%</i>	<i>100%</i>	<i>100%</i>	<i>92%</i>	<i>83%</i>	<i>92%</i>	<i>83%</i>	<i>100%</i>	<i>100%</i>	<i>100%</i>	<i>100%</i>
<i>Number of Probes</i>	<i>6*</i>	<i>6*</i>	<i>7</i>	<i>6*</i>	<i>2**</i>	<i>2***</i>	<i>2**</i>	<i>2***</i>	<i>9</i>	<i>9</i>	<i>10</i>	<i>4</i>

*\*, \*\* and \*\*\** have identical probes selected for the model.

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.16 Beta values of individual transcripts within HR cancer and CB models suggested by Lasso using different input transcripts for the baseline normalised data.**

<i>Transcript</i>	<i>Beta – All Transcripts</i>	<i>Beta - MWU</i>	<i>Beta - glm</i>	<i>Beta - both</i>
<i>AATF</i>	<b>0.094</b>	<b>0.696</b>	<b>0.982</b>	
<i>CADPS</i>	<b>0.773</b>	<b>1.255</b>	<b>1.798</b>	<b>0.337</b>
<i>CAMKK2</i>		<b>0.055</b>		
<i>CCDC88B</i>		<b>0.037</b>		
<i>CDKN3</i>			<b>-0.101</b>	
<i>CKAP2L</i>	<b>0.042</b>		<b>0.358</b>	
<i>ERG 3' exons 6-7</i>	<b>0.135</b>	<b>0.219</b>	<b>0.193</b>	<b>0.096</b>
<i>HOXC6</i>	<b>0.197</b>	<b>0.218</b>	<b>0.168</b>	<b>0.115</b>
<i>IGFBP3</i>	<b>-0.051</b>	<b>-0.288</b>		
<i>LASS1</i>	<b>0.187</b>	<b>0.186</b>	<b>0.623</b>	
<i>MCTP1</i>	<b>0.003</b>			<b>0.029</b>
<i>MMP25</i>		<b>0.092</b>	<b>0.197</b>	
<i>NAALADL2</i>	<b>-0.121</b>			
<i>SIM2 long</i>			<b>0.337</b>	
<i>TDRD</i>		<b>0.084</b>		

**Table 6.17 Beta values of individual transcripts within HR cancer and CB models suggested by Lasso using different input transcripts for the *KLK2* ratio data.**

<i>Transcript</i>	<i>Beta – All Transcripts</i>	<i>Beta - MWU</i>	<i>Beta - glm</i>	<i>Beta - both</i>
<i>ERG 3' exons 6-7</i>	<b>0.3394</b>	<b>0.5927</b>	<b>0.391</b>	<b>0.391</b>
<i>HOXC6</i>	<b>0.0287</b>	<b>0.1029</b>		
<i>SIM2 long</i>			<b>0.0349</b>	<b>0.0349</b>

**Table 6.18 Beta values of individual transcripts within HR cancer and CB models suggested by Lasso using different input transcripts for the HK normalised data.**

<i>Transcript</i>	<i>Beta – All transcripts</i>	<i>Beta - MWU</i>	<i>Beta – glm</i>	<i>Beta - both</i>
<i>CADPS</i>	<b>0.3134</b>	<b>0.8661</b>	<b>0.4861</b>	<b>1.2907</b>
<i>ERG 3' exons 6-7</i>	<b>0.0636</b>	<b>0.0684</b>	<b>0.0784</b>	<b>0.0329</b>
<i>GJB1</i>	<b>-0.0503</b>	<b>-0.009</b>	<b>-0.0518</b>	
<i>HOXC6</i>	<b>0.1835</b>	<b>0.2967</b>	<b>0.2332</b>	<b>0.3504</b>
<i>NAALADL2</i>	<b>-0.1273</b>	<b>-0.3281</b>	<b>-0.1872</b>	<b>-0.4819</b>
<i>SIM2 long</i>		<b>0.1225</b>		<b>0.3251</b>
<i>SPINK1</i>	<b>-0.0754</b>	<b>-0.0949</b>	<b>-0.0812</b>	<b>-0.1063</b>

Random forest was also applied to i) all transcripts, ii) significant transcripts identified by MWU and iii) significant transcripts identified by glm for the three different normalisations (the baseline data, *KLK2* ratio, and HK normalised), respectively (Supplementary Table 40, Supplementary Table 41 and Supplementary Table 42).

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

Using the glm identified significant probes in the HK normalised data gives the models with the smallest error (mean square of residuals: 0.117), although all models are very similar (Table 6.20).

*HOXC6* was in the top 5 transcripts in 8/9 random forests, whilst *CADPS* was in the top 5 transcripts in 7/9 random forests. *ERG3*' exons 6-7 and *SPINK1* occur in 4/9, and *ST6GALNAC1* and *TDRD* occur in 3/9 random forests (Table 6.19).

**Table 6.19** Frequency of transcripts in top 5 for random forests (CB vs high-risk cancer models).

<i>Transcript</i>	<i>Frequency in top 5 random forest important transcripts</i>	<i>Data</i>
<i>CADPS</i>	7	<i>Baseline + KLK2 + HK</i>
<i>CCDC88B</i>	2	<i>Baseline</i>
<i>ERG3' exons 6-7</i>	4	<i>Baseline + KLK2 + HK</i>
<i>HOXC6</i>	8	<i>Baseline + KLK2 + HK</i>
<i>SIM2 long</i>	2	<i>KLK2</i>
<i>SLC43A1</i>	2	<i>KLK2</i>
<i>SPINK1</i>	4	<i>Baseline + HK</i>
<i>ST6GALNAC1</i>	3	<i>HK</i>
<i>TDRD</i>	3	<i>KLK2</i>
<i>VAX2</i>	2	<i>HK</i>

**Table 6.20** Mean Square of residuals error for each random forest model produced using different input probes in three different normalisations.

	<i>HK</i>			<i>KLK2</i>			<i>Baseline</i>		
<i>Input:</i>	<i>All</i>	<i>glm</i>	<i>MWU</i>	<i>All</i>	<i>glm</i>	<i>MWU</i>	<i>All</i>	<i>glm</i>	<i>MWU</i>
<i>Mean square of residuals</i>	0.18	0.117	0.138	0.149	0.18	0.14	0.146	0.138	0.145

### 6.3.3 Trend CBN-L-I-H

#### 6.3.3.1 Significant transcripts

Trend (increase or decrease) in expression across the 12 CB samples, 4 low-risk, 28 intermediate risk and 14 high-risk samples was investigated. Two methods of ordered multinomial regression were used: i) proportional odds logistic regression (polr) and ii)

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

logistic regression setting clinical group to an ordered integer. Using polr there were 70, 20 and 15 transcripts that significantly modelled the trend in the three processed datasets (the baseline data, *KLK2* ratio, and HK normalised, respectively), ( $p < 0.05$ , Supplementary Table 43). Of these only 7 (*B4GALNT4*, *HOXC6*, *ERG3*’ exons 6-7, *APOC1*, *TMPRSS2:ERG*, *NEAT1*, and *MCTP1*), 2 (*HOXC6* and *ERG 3*’ exons 6-7) and 1 (*HOXC6*) remained significant post multiple testing correction. The top significant probes identified by polr were *APOC1*, *HOXC6* and *ERG 3*’ exons 6-7 jointly ( $p = 0.0001$ ), *HOXC6* ( $p = 1.36 \times 10^{-05}$ ) and *HOXC6* ( $p = 4.54 \times 10^{-6}$ ), respectively.

Using logistic regression there were 87, 36 and 19 transcripts that modelled trend with statistical significance in the three processed datasets (the baseline data, *KLK2* ratio, and HK normalised, respectively), ( $p < 0.05$ , Supplementary Table 44). Of these 19, 4 (*HOXC6*, *ERG3*’ exons 6-7, *TMPRSS2:ERG* and *TDRD*), and 1 (*HOXC6*) remained significant post multiple testing correction, respectively. The top significant probes identified by polr were *APOC1* ( $p = 2.90 \times 10^{-06}$ , adjusted p-value = 0.0005), *ERG 3*’ exons 6-7 and *HOXC6* jointly ( $p = 0.0002$ ) and *HOXC6* ( $p = 6.37 \times 10^{-05}$ ), respectively (Table 6.21, Table 6.22, Table 6.23).

Polr identifies fewer transcripts than glm but all but one transcript identified by polr were also identified by logistic regression in each case, showing robustness in their identification. Similar probes were identified as most significant by the two methods also.

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.21 Top 15 significant transcripts identified by polr to have trend across CB - L - I - H clinical categories in the baseline normalised cell data.**

Transcript	glm <i>p</i> -value	glm adjusted <i>p</i> -	polr <i>p</i> -value	polr adjusted <i>p</i> -
<i>APOC1</i>	2.90x10 <sup>-06</sup>	0.0005	0.0001	0.0246
<i>ERG3'</i> exons 6-7	5.87x10 <sup>-06</sup>	0.001	0.0001	0.0191
<i>HOXC6</i>	1.85x10 <sup>-05</sup>	0.003	0.0001	0.0191
<i>TMPRSS2:ERG</i>	6.00x10 <sup>-05</sup>	0.0095	0.0002	0.0385
<i>MCTP1</i>	4.10x10 <sup>-06</sup>	0.0007	0.0003	0.0481
<i>NEAT1</i>	0.0002	0.0265	0.0003	0.047
<i>RIOK3</i>	6.64x10 <sup>-06</sup>	0.0011	0.0003	0.0515
<i>ISX</i>	3.91x10 <sup>-05</sup>	0.0063	0.0004	0.066
<i>HPN</i>	6.51x10 <sup>-05</sup>	0.0102	0.0007	0.1079
<i>GCNT1</i>	0.0002	0.025	0.0008	0.1318
<i>SULF2</i>	2.23x10 <sup>-05</sup>	0.0036	0.0008	0.1288
<i>CAMKK2</i>	5.44x10 <sup>-05</sup>	0.0087	0.0012	0.1813
<i>MMP25</i>	0.0002	0.0277	0.0014	0.2178
<i>CADPS</i>	0.0001	0.019	0.0017	0.261
<i>LASS1</i>	0.0002	0.0331	0.0019	0.2854

**Table 6.22 Top 15 significant transcripts identified by polr to have trend across CB - L - I - H clinical categories in the *KLK2* ratio cell data.**

Transcript	glm <i>p</i> -value	glm adjusted <i>p</i> -	polr <i>p</i> -value	polr adjusted <i>p</i> -
<i>ERG3'</i> exons 6-7	2.18 x10 <sup>-05</sup>	0.0036	0.0002	0.0283
<i>HOXC6</i>	1.36x10 <sup>-05</sup>	0.0023	0.0002	0.0406
<i>TMPRSS2:ERG</i>	6.86 x10 <sup>-05</sup>	0.0112	0.0007	0.1136
<i>TDRD</i>	0.0002	0.0263	0.0011	0.1757
<i>SIM2</i> long	0.0031	0.5021	0.0056	0.9028
<i>HPN</i>	0.0027	0.4419	0.0081	0.994
<i>GCNT1</i>	0.0066	0.998	0.0104	0.994
<i>CADPS</i>	0.0052	0.8154	0.0112	0.994
<i>TMEM86A</i>	0.0079	0.998	0.0184	0.994
<i>CKAP2L</i>	0.0046	0.731	0.0187	0.994
<i>LASS1</i>	0.005	0.7872	0.0209	0.994
<i>ERG3'</i> exons 4-5	0.0275	0.998	0.0243	0.994
<i>FOLH1</i>	0.0124	0.998	0.0348	0.994
<i>ISX</i>	0.0022	0.3607	0.0367	0.994
<i>ANKRD34B</i>	0.0112	0.998	0.0374	0.994

**Table 6.23 Top 15 significant transcripts identified by polr to have trend across CB - L - I - H clinical categories in the HK normalised cell data.**

Transcript	glm <i>p</i> -value	glm adjusted <i>p</i> -	polr <i>p</i> -value	polr adjusted <i>p</i> -
<i>HOXC6</i>	4.54 x10 <sup>-6</sup>	0.0008	6.37x10 <sup>-05</sup>	0.0106
<i>TDRD</i>	0.0012	0.2024	0.0034	0.564
<i>SIM2</i> long	0.0032	0.5147	0.0043	0.7056
<i>SLC43A1</i>	0.0011	0.1895	0.006	0.978
<i>UPK2</i>	0.0028	0.4609	0.0077	0.9994
<i>ERG 3'</i> exons 6-7	0.0043	0.6877	0.0098	0.9994
<i>NAALADL2</i>	0.0018	0.2913	0.0098	0.9994
<i>TMPRSS2:ERG</i> fusion	0.004	0.6414	0.0127	0.9994
<i>ST6GALNAC1</i>	0.0049	0.7755	0.0179	0.9994
<i>FOLH1</i>	0.0174	0.9941	0.0191	0.9994
<i>MEX3A</i>	0.0243	0.9941	0.0337	0.9994
<i>TMEM86A</i>	0.0107	0.9941	0.0337	0.9994
<i>SERPINB5</i>	0.0162	0.9941	0.0425	0.9994
<i>PALM3</i>	0.027	0.9941	0.0461	0.9994
<i>EN2</i>			0.0463	0.9994

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

### 6.3.3.2 *Models and gene selection*

A number of different transcript subsets were input to Lasso for probe shrinkage and selection, these included i) all of the transcripts ( $n = 167$ ). Transcripts identified as having significant decrease or increase in expression across CB->L->I->H clinical categories using ii) polr ( $n = 70$ ,  $n = 20$  and  $n = 15$ ), and iii) logistic regression ( $n = 87$ ,  $n = 36$ , and  $n = 19$ ), separately, for each of the three normalisations (the baseline data, *KLK2* ratio, and HK normalised), respectively (Supplementary Table 45). In addition, the transcripts common to both those identified by polr and glm for the HK normalised data only was also submitted to Lasso ( $n = 14$ ), these were the only significant transcript lists where polr did not contain all of the glm identified probes. *APOC1* was the only probe selected by Lasso in all three transcript inputs for the baseline normalised data. *HOXC6* was the only probe selected by Lasso in the *KLK2* ratio data. *HOXC6*, *NAALADL2* and *UPK2* were common probes selected by Lasso in the HK normalised data.

The AUC, sensitivity and specificity of each model on the same training data was collected (Table 6.27) and transcript lists (Table 6.24, Table 6.25, Table 6.26) and boxplots of the Lasso selected probes were produced (Supplementary Figure 19, Supplementary Figure 20 and Supplementary Figure 21).

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.24 Optimal multinomial models for predicting clinical category (CB, low-risk, intermediate-risk, and high-risk cancer) with different subsets of input transcripts (from preliminary ordered glm and polr tests) in the baseline normalised cell data.**

<i>Transcript</i>	<i>All transcripts (n = 167) - Beta</i>	<i>glm (n = 87) - Beta</i>	<i>polr (n = 70) - Beta</i>
<i>AATF</i>	<b>0.1</b>	<b>0.115</b>	<b>0.109</b>
<i>APOC1</i>		<b>0.056</b>	
<i>B4GALNT4</i>	<b>0.004</b>	<b>0.121</b>	
<i>CADPS</i>	<b>0.068</b>		<b>0.105</b>
<i>CAMKK2</i>		<b>0.062</b>	
<i>CCDC88B</i>		<b>0.036</b>	
<i>EN2</i>			<b>0.024</b>
<i>ERG 3' exons 6-7</i>	<b>0.154</b>	<b>0.139</b>	<b>0.144</b>
<i>HOXC6</i>	<b>0.16</b>		<b>0.121</b>
<i>KLK3 exons 2-3</i>	<b>-0.022</b>		
<i>LASS1</i>	<b>0.034</b>		<b>0.017</b>
<i>MCTP1</i>	<b>0.037</b>		<b>0.095</b>
<i>MMP25</i>	<b>0.007</b>		<b>0.088</b>
<i>NAALADL2</i>	<b>-0.104</b>	<b>-0.034</b>	
<i>RIOK3</i>	<b>0.016</b>		<b>0.095</b>
<i>SPINK1</i>			<b>-0.1</b>
<i>SULF2</i>	<b>0.088</b>		<b>0.046</b>
<i>VAX1</i>	<b>-0.131</b>		
<i>Cp1</i>	<b>1.198</b>	<b>0.702</b>	<b>1.283</b>
<i>Cp2</i>	<b>2.018</b>	<b>1.786</b>	<b>2.132</b>
<i>Cp3</i>	<b>-0.414</b>	<b>-0.368</b>	<b>-0.436</b>

**Table 6.25 Optimal multinomial models for predicting clinical category (CB, low-risk, intermediate-risk, and high-risk cancer) with different subsets of input transcripts (from preliminary ordered glm and polr tests) in *KLK2* ratio cell data**

<i>Transcript</i>	<i>All transcripts (n = 94) - Beta</i>	<i>glm (n = 36) - Beta</i>	<i>polr (n = 20) - Beta</i>
<i>CADPS</i>	<b>0.027</b>	<b>0.075</b>	<b>0.026</b>
<i>CKAP2L</i>		<b>0.074</b>	
<i>ERG3' exons 4-5</i>		<b>-0.121</b>	
<i>ERG3' exons 6-7</i>	<b>0.809</b>	<b>0.784</b>	<b>0.591</b>
<i>HOXC4</i>	<b>-0.134</b>		
<i>HOXC6</i>	<b>0.264</b>	<b>0.475</b>	<b>0.34</b>
<i>ITGBL1</i>	<b>-0.118</b>		
<i>NAALADL2</i>	<b>-0.479</b>		
<i>PALM3</i>	<b>-0.017</b>		
<i>RIOK3</i>		<b>-0.487</b>	
<i>TDRD</i>		<b>0.031</b>	
<i>TMPRSS2:ERG</i>	<b>0.008</b>	<b>0.12</b>	<b>0.131</b>
<i>Cp1</i>	<b>0.998</b>	<b>1.094</b>	<b>0.837</b>
<i>Cp2</i>	<b>1.938</b>	<b>2.065</b>	<b>1.937</b>
<i>Cp3</i>	<b>-0.401</b>	<b>-0.427</b>	<b>-0.400</b>

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.26 Optimal multinomial models for predicting clinical category (CB, low-risk, intermediate-risk, and high-risk cancer) with different subsets of input transcripts (from preliminary ordered glm and polr tests) in HK normalised cell data.**

<i>Transcript</i>	<i>All transcripts (n = 167) - Beta</i>	<i>glm (n = 19) - Beta</i>	<i>polr (n = 15) - Beta</i>	<i>glm and polr (n = 14) - Beta</i>
<i>CADPS</i>		<b>0.035</b>		
<i>CLIC2</i>		<b>0.043</b>		
<i>ERG 3' exons 6-7</i>	<b>0.088</b>	<b>0.19</b>	<b>0.179</b>	<b>0.211</b>
<i>GJB1</i>	<b>-0.057</b>	<b>-0.233</b>		
<i>HOXC6</i>	<b>0.151</b>	<b>0.266</b>	<b>0.205</b>	<b>0.234</b>
<i>NAALADL2</i>	<b>-0.094</b>	<b>-0.183</b>	<b>-0.235</b>	<b>-0.285</b>
<i>PALM3</i>		<b>-0.028</b>	<b>-0.045</b>	<b>-0.074</b>
<i>SLC43A1</i>		<b>0.015</b>	<b>0.033</b>	<b>0.058</b>
<i>TDRD</i>		<b>0.045</b>		
<i>TMEM86A</i>				<b>0.023</b>
<i>UPK2</i>	<b>-0.014</b>		<b>-0.003</b>	
<i>Cp1</i>	<b>0.65</b>	<b>1.622</b>	<b>1.276</b>	<b>1.563</b>
<i>Cp2</i>	<b>1.76</b>	<b>2.321</b>	<b>2.089</b>	<b>2.26</b>
<i>Cp3</i>	<b>-0.361</b>	<b>-0.466</b>	<b>-0.425</b>	<b>-0.456</b>



CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

Table 6.27 AUC, Sensitivity and Specificity of models to predict trend across clinical categories: CB > L- > I- > H-risk cancer in different data normalisations of cell NanoString data.

<i>Data type:</i>	<i>Baseline</i>			<i>KLK2</i>			<i>HK</i>			
<i>Model Input:</i>	<i>All transcripts</i>	<i>Glm</i>	<i>Polr</i>	<i>All transcripts</i>	<i>Glm</i>	<i>Polr</i>	<i>All transcripts</i>	<i>Glm</i>	<i>Polr</i>	<i>glm + polr</i>
<i>Accuracy</i>	<b>0.7069</b>	<b>0.6552</b>	<b>0.6897</b>	<b>0.7069</b>	<b>0.6379</b>	<b>0.6207</b>	<b>0.6552</b>	<b>0.6897</b>	<b>0.6897</b>	<b>0.6897</b>
<i>AUC</i>	<b>0.7604</b>	<b>0.7242</b>	<b>0.7669</b>	<b>0.7504</b>	<b>0.7252</b>	<b>0.702</b>	<b>0.6944</b>	<b>0.7609</b>	<b>0.7609</b>	<b>0.7609</b>
<i>Sensitivity: CB</i>	<b>83%</b>	<b>67%</b>	<b>92%</b>	<b>75%</b>	<b>67%</b>	<b>58%</b>	<b>67%</b>	<b>83%</b>	<b>83%</b>	<b>83%</b>
<i>L</i>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>	<b>0%</b>
<i>I</i>	<b>96%</b>	<b>93%</b>	<b>93%</b>	<b>96%</b>	<b>89%</b>	<b>89%</b>	<b>100%</b>	<b>93%</b>	<b>93%</b>	<b>93%</b>
<i>H</i>	<b>29%</b>	<b>26%</b>	<b>21%</b>	<b>36%</b>	<b>29%</b>	<b>29%</b>	<b>14%</b>	<b>29%</b>	<b>29%</b>	<b>29%</b>
<i>Specificity: CB</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>98%</b>	<b>100%</b>	<b>98%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>
<i>L</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>
<i>I</i>	<b>47%</b>	<b>40%</b>	<b>47%</b>	<b>47%</b>	<b>40%</b>	<b>37%</b>	<b>33%</b>	<b>47%</b>	<b>47%</b>	<b>47%</b>
<i>H</i>	<b>98%</b>	<b>95%</b>	<b>95%</b>	<b>100%</b>	<b>93%</b>	<b>95%</b>	<b>100%</b>	<b>95%</b>	<b>95%</b>	<b>95%</b>
<i>PPV: CB</i>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>90%</b>	<b>100%</b>	<b>88%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>	<b>100%</b>
<i>L</i>	<b>NA</b>	<b>NA</b>	<b>NA</b>	<b>NA</b>	<b>NA</b>	<b>NA</b>	<b>NA</b>	<b>NA</b>	<b>NA</b>	<b>NA</b>
<i>I</i>	<b>63%</b>	<b>59%</b>	<b>62%</b>	<b>63%</b>	<b>58%</b>	<b>57%</b>	<b>58%</b>	<b>62%</b>	<b>62%</b>	<b>62%</b>
<i>H</i>	<b>80%</b>	<b>67%</b>	<b>60%</b>	<b>100%</b>	<b>57%</b>	<b>67%</b>	<b>100%</b>	<b>67%</b>	<b>67%</b>	<b>67%</b>
<i>Number of Probes</i>	<b>13</b>	<b>7</b>	<b>11</b>	<b>8</b>	<b>8</b>	<b>4</b>	<b>5</b>	<b>9</b>	<b>6</b>	<b>6</b>

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

Random forest was also applied to i) all transcripts, ii) significant transcripts identified by polr and iii) significant transcripts identified by glm for the three different normalisations (the baseline data, *KLK2* ratio, and HK normalised), respectively (Supplementary Table 46, Supplementary Table 47 and Supplementary Table 48). Using all probes in the baseline normalised data gives the models with the smallest error (OOB error: 27.6%, Table 6.28). *ERG3'* exons 6-7 was present in 8/9 random forest models, whilst *TMPRSS2:ERG* and *HOXC6* were present in 6/9 RF models (Table 6.29).

**Table 6.28 OOB error rates for random forest models built to predict trend over clinical categories: CB > L > I > H**

<i>Input:</i>	<i>Baseline</i>			<i>KLK2</i>			<i>HK</i>			<i>Glm + polr</i>
	<i>All</i>	<i>glm</i>	<i>polr</i>	<i>All</i>	<i>glm</i>	<i>polr</i>	<i>All</i>	<i>glm</i>	<i>polr</i>	
<i>OOB error</i>	27.6 %	50%	50%	46.6 %	48.3 %	51.7 %	44.8 %	43.1 %	41.4 %	44.8 %

**Table 6.29 Frequency of transcripts in top 5 for random forests (CB > L > I > H trend models).**

<i>Transcript</i>	<i>Frequency in top 5 random forest important transcripts</i>	<i>Data</i>
<i>ERG3' exons 6-7</i>	8	<i>Baseline, KLK2 and HK</i>
<i>TMPRSS2:ERG</i>	6	<i>Baseline, KLK2 and HK</i>
<i>HOXC6</i>	6	<i>KLK2 and HK</i>
<i>PCA3</i>	3	<i>KLK2 and HK</i>
<i>PALM3</i>	3	<i>HK</i>
<i>RIOK3</i>	2	<i>Baseline</i>
<i>NEAT1</i>	2	<i>Baseline</i>
<i>CADPS</i>	2	<i>Baseline</i>
<i>APOC1</i>	2	<i>Baseline</i>
<i>FOLH1</i>	2	<i>KLK2</i>
<i>NAALADL2</i>	2	<i>HK</i>
<i>UPK2</i>	2	<i>HK</i>

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**6.4 Summary of Predictive Models**

In the *KLK2* ratio data half of the models had better AUCs on the training set from the cell data and half from the EV data (Table 6.30). However, in the HK data, the AUCs were higher in the cell data. I am limited in the number of samples for the cell fraction and so the models have not been applied to a test data set, this means that the models could be over fitting the data.

Table 6.30 Comparison of AUCs from models using cell and EV data.

	Cell	EV
<b>KLK2 ratio data</b>		
CB vs Cancer (L, I, H) All transcripts	0.996	0.949
CB vs Cancer (L, I, H) Significant transcripts	0.998	0.886
CB vs HR Cancer (H) All transcripts	0.952	0.991
CB vs HR Cancer (H) Significant transcripts	0.958	0.97
CB > L > I > H All transcripts	0.7504	0.7663
CB > L > I > H Significant transcripts	0.702	0.6757
<b>HK normalised data</b>		
CB vs Cancer (L, I, H) All transcripts	0.989	0.925
CB vs Cancer (L, I, H) Significant transcripts	0.998	0.902
CB vs HR Cancer (H) All transcripts	1	0.976
CB vs HR Cancer (H) Significant transcripts	1	0.992
CB > L > I > H All transcripts	0.7609	0.7587
CB > L > I > H Significant transcripts	0.7609	0.7728

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.31** Transcripts identified by all selection models for the different clinical category tests across the different normalisations on the cell NanoString data.

<i>Normalisation</i>	<i>Clinically Benign vs. Cancer</i>	<i>Clinically Benign vs. High risk cancer</i>	<i>Trend Clinically Benign, low-risk, intermediate-risk and high-risk</i>
<i>Baseline</i>	<i>ACTR5 APOC1 ARHGEF25 CADPS CAMKK2 ERG3' exon 6-7 EN2 HIST1H2BG HOXC6 IGFBP3 LASS1 MCTP1 MMP25 MMP26 NAALADL2 PCA3 RIOK3 SPINK1 SLC12A1 TDRD</i>	<i>AATF CADPS CAMKK2 CCDC88B CDKN3 CKAP2L ERG3' exon 6-7 HOXC6 ITGFBP3 LASS1 MCTP1 MMP25 NAALADL2 SIM2 long TDRD</i>	<i>AATF APOC1 B4GALNT4 CADPS CAMKK2 CCDC88B EN2 ERG3' exon 6-7 HOXC6 KLK3 exons 2-3 LASS1 MCTP1 MMP25 NAALADL2 RIOK3 SPINK1 SULF2 VAX1</i>
<i>KLK2 ratio</i>	<i>CADPS CKAP2L EN2 ERG3' exons 6-7 HOXC6 MFSD2A NAALADL2 SFRP4 SIM2 long TDRD</i>	<i>ERG3' exons 6-7 HOXC6 SIM2 long</i>	<i>CADPS CKAP2L ERG3' exons 4-5 ERG3' exons 6-7 HOXC4 HOXC6 ITGBL1 NAALADL2 PALM3 RIOK3 TDRD TMPRSS2:ERG</i>
<i>RPLP2 and TWIST1 normalised</i>	<i>CADPS CLIC2 ERG3' exons 6-7 HOXC6 NAALADL2 SIM2 long TDRD UPK2</i>	<i>CADPS ERG3' exon 6-7 GJB1 HOXC6 NAALADL2 SIM2 long SPINK1</i>	<i>CADPS CLIC2 ERG3' exons 6-7 GJB1 HOXC6 NAALADL2 PALM3 SLC43A1 TDRD TMEM86A UPK2</i>

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

**Table 6.32** Transcripts selected for models in EV data.

<i>Normalisation</i>	<i>Clinically Benign vs. Cancer</i>	<i>Clinically Benign vs. High risk cancer</i>	<i>Trend Clinically Benign, low-risk, intermediate-risk and high-risk</i>
<i>KLK2 ratio</i>	<i>AMACR</i>	<i>ACTR5</i>	<i>AMACR</i>
	<i>APOC1</i>	<i>ALAS1</i>	<i>ANKRD34B</i>
	<i>AR exon 9</i>	<i>AMACR</i>	<i>APOC1</i>
	<i>CP</i>	<i>ANKRD34B</i>	<i>AR exon 9</i>
	<i>DLX1</i>	<i>APOC1</i>	<i>AR exons 4-8</i>
	<i>ERG3' exon 4-5</i>	<i>AR exon 9</i>	<i>BTG2</i>
	<i>GJB1</i>	<i>AR exons 4-8</i>	<i>CD10</i>
	<i>HOXC6</i>	<i>AURKA</i>	<i>CP</i>
	<i>IGFBP3</i>	<i>BTG2</i>	<i>DLX1</i>
	<i>ISX</i>	<i>CD10</i>	<i>DPP4</i>
	<i>KLK4</i>	<i>CKAP2L</i>	<i>ERG 3' exons 4-5</i>
	<i>MXII</i>	<i>CP</i>	<i>ERG 3' exons 6-7</i>
	<i>NEAT1</i>	<i>DLX1</i>	<i>GABARAPL2</i>
	<i>PCA3</i>	<i>DPP4</i>	<i>HIST1H1E</i>
	<i>PPP1R12B</i>	<i>ERG 3' exons 4-5</i>	<i>HOXC6</i>
	<i>RNF157</i>	<i>4-5</i>	<i>HPN</i>
	<i>ST6GALNAC</i>	<i>HOXC6</i>	<i>IGFBP3</i>
	<i>SULT1A1</i>	<i>HPN</i>	<i>ISX</i>
	<i>TDRD</i>	<i>IGFBP3</i>	<i>ITGBL1</i>
	<i>TMEM47</i>	<i>ISX</i>	<i>KLK4</i>
	<i>TMPRSS2:ERG</i>	<i>KLK4</i>	<i>MED4</i>
		<i>MAK</i>	<i>MEMO1</i>
		<i>MED4</i>	<i>MXII</i>
		<i>MMP25</i>	<i>MYOF</i>
		<i>NEAT1</i>	<i>NEAT1</i>
		<i>PCA3</i>	<i>PCA3</i>
		<i>PDLIM5</i>	<i>PPP1R12B</i>
		<i>PPFIA2</i>	<i>PSGR</i>
		<i>PSTPIP1</i>	<i>PSTPIP1</i>
		<i>PTPRC</i>	<i>SLC12A1</i>
		<i>RPL18A</i>	<i>SRSF3</i>
		<i>SRSF3</i>	<i>SULT1A1</i>
		<i>STEAP4</i>	<i>TDRD</i>
		<i>TMEM47</i>	<i>Timp4</i>
	<i>TMPRSS2:ERG</i>	<i>TMEM47</i>	
		<i>TMPRSS2:ERG</i>	
		<i>ZNF577</i>	
<i>RPLP2 and GAPDH normalised</i>	<i>APOC1</i>	<i>AMACR</i>	<i>ACT5R</i>
	<i>AR exon 9</i>	<i>ANKRD34B</i>	<i>AMH</i>
	<i>CD10</i>	<i>APOC1</i>	<i>ANKRD34B</i>
	<i>CP</i>	<i>AR exon 9</i>	<i>APOC1</i>
	<i>ERG3' exons 4-5</i>	<i>AR exon 4-8</i>	<i>AR exon 9</i>
	<i>GABARAPL2</i>	<i>CD10</i>	<i>AR exon 4-8</i>
	<i>HOXC6</i>	<i>DLX1</i>	<i>CD10</i>
	<i>HPN</i>	<i>DPP4</i>	<i>CP</i>
	<i>ISX</i>	<i>ERG3' exons 4-5</i>	<i>DPP4</i>
	<i>KLK4</i>	<i>5</i>	<i>ERG3' exons 4-5</i>
	<i>MCTP1</i>	<i>GABARAPL2</i>	<i>ERG3' exons 6-7</i>
		<i>HOXC6</i>	<i>FDPS</i>

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

<i>PCA3</i>	<i>HPN</i>	<i>GABARAPL2</i>
<i>PDLIM5</i>	<i>KLK4</i>	<i>GCNT1</i>
<i>PPP1R12B</i>	<i>MYOF</i>	<i>GJB1</i>
<i>PTN</i>	<i>NEAT1</i>	<i>HIST1H1E</i>
<i>SLC12A1</i>	<i>PCA3</i>	<i>HIST1H2BF</i>
<i>SULT1A1</i>	<i>PDLIM5</i>	<i>HOXC6</i>
<i>TDRD</i>	<i>SLC12A1</i>	<i>HPN</i>
<i>TMPRSS2:ERG</i>	<i>SRSF3</i>	<i>IGFBP3</i>
	<i>STOM</i>	<i>ISX</i>
	<i>SULT1A1</i>	<i>ITGBL1</i>
	<i>TMPRSS2:ERG</i>	<i>KLK4</i>
		<i>MED4</i>
		<i>MEMO1</i>
		<i>MIATNB</i>
		<i>MSMB</i>
		<i>MXI1</i>
		<i>MYOF</i>
		<i>NEAT1</i>
		<i>PCA3</i>
		<i>PPP1R12B</i>
		<i>RPS10</i>
		<i>SLC12A1</i>
		<i>SPINK1</i>
		<i>SRSF3</i>
		<i>SULT1A1</i>
		<i>TDRD</i>
		<i>Timp4</i>
		<i>TMPRSS2:ERG</i>
		<i>TRPM4</i>
		<i>UPK2</i>
		<i>ZNF577</i>

Comparing the transcripts selected for models in the cell *KLK2* ratio data (Table 6.31) and the EV *KLK2* ratio data (Table 6.32), only 5 transcripts were selected for both sets of models (*CKAP2L*, *HOXC6*, *TDRD*, *ITGBL1* and *TMPRSS2:ERG*). The same comparison for the HK normalised data yielded a different 5 transcripts in common (*ERG* 3' exons 6-7, *HOXC6*, *TDRD*, *GJB1* and *UPK2*). This shows that different probes are selected as important for predictive models between the different fractions of urine (cell vs EV).

**6.5 Comparison of the urine expression profiles of Extracellular vesicle and Cell fractions in Prostate Cancer**

**6.5.1 Microarray comparison of the global expression profile of Extracellular vesicle and Cell fractions.**

I examined Affymetrix microarray expression data from the cell sediment and EV fraction of urine collected from prostate cancer patients from either the NNUH or the Royal Marsden Hospital NHS foundation trust ( $n = 3$ ). Genes that were significantly differentially expressed between the two fractions were determined by Dr. Daniel Brewer using the Limma package and the method proposed by Mootha *et al.*, 2003, to give a value for variance of expression and if it significantly differs between fractions<sup>219</sup>. 98 genes were found to be up-regulated in the extracellular vesicles and 116 up-regulated in the cell sediment fraction.

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

Table 6.33 A list of the top 20 microarray detected transcripts out of 98 that were found to be significantly more abundant in extracellular vesicles compared with sediment from the same urine.

<i>Rank</i>	<i>Gene</i>	<i>Log<sub>2</sub>(FC)</i>	<i>p-value</i>	<i>Tissue Expression</i>	<i>Known Cancer Associations</i>
1	<i>TMSB15A</i>	5.19	0.048	<i>Prostate</i>	<i>Prostate</i> <sup>220</sup> , <i>Other</i> <sup>221,222,223</sup>
2	<i>PRKG2</i>	5.02	0.035	<i>Prostate, Other</i>	<i>N</i>
3	<i>TCEA3</i>	4.87	0.048	<i>Other</i>	<i>N</i>
4	<i>PRAC</i>	4.85	0.035	<i>Prostate, Other</i>	<i>Prostate</i> <sup>224</sup> , <i>Other</i> <sup>225</sup>
5	<i>KLK4</i>	4.82	0.041	<i>Prostate</i>	<i>Prostate</i> <sup>226,227</sup>
6	<i>FOLH1,</i> <i>FOLH1B</i>	4.59	0.048	<i>Prostate, Other</i>	<i>Prostate</i> <sup>228</sup>
7	<i>EPHX2</i>	4.58	0.041	<i>Expressed in all</i>	<i>Prostate</i> <sup>229</sup> , <i>Other</i> <sup>230,231,232</sup>
8	<i>GMPR</i>	4.57	0.042	<i>Expressed in all (higher expression in Prostate, Other)</i>	<i>Prostate</i> <sup>233</sup>
9	<i>RANBP3L</i>	4.5	0.046	<i>Prostate, Kidney, Other</i>	<i>Multiple</i> <sup>234,235</sup>
10	<i>MPPED2</i>	4.39	0.047	<i>Prostate, Other</i>	<i>Other</i> <sup>236</sup>
11	<i>CKB</i>	4.17	0.035	<i>Expressed in all (highest expression in Prostate)</i>	<i>Prostate</i> <sup>237</sup> , <i>Other</i> <sup>238</sup>
12	<i>MLPH</i>	4.09	0.045	<i>Prostate, Other</i>	<i>Prostate</i> <sup>239</sup>
13	<i>NFIA</i>	4.06	0.048	<i>Expressed in all</i>	<i>Prostate</i> <sup>240</sup> , <i>Other</i> <sup>241</sup>
14	<i>GLYATL1</i>	4.00	0.049	<i>Prostate, Kidney, Other</i>	<i>Other</i> <sup>242</sup>
15	<i>NFIB</i>	3.98	0.048	<i>Mixed</i>	<i>Prostate</i> <sup>243</sup> , <i>Other</i> <sup>244</sup>
16	<i>CCDC88C</i>	3.98	0.043	<i>Expressed in all</i>	<i>Other</i> <sup>245</sup>
17	<i>HOXB13</i>	3.97	0.042	<i>Prostate</i>	<i>Prostate</i> <sup>246</sup> , <i>Other</i> <sup>246</sup>
18	<i>PART1</i>	3.95	0.035	<i>Prostate*</i>	<i>Prostate</i> <sup>247</sup>
19	<i>AZGP1</i>	3.85	0.043	<i>Prostate, Other</i>	<i>Prostate</i> <sup>248</sup> , <i>Other</i> <sup>249,250</sup>
20	<i>TCEAL2</i>	3.73	0.048	<i>Tissue Enhanced (glands, reproductive including prostate and cerebral cortex)</i>	<i>Other</i> <sup>251</sup>



## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

I researched the top twenty up-regulated transcripts in the two fractions to examine the link between the genes, prostate tissue and cancer (Table 6.33, Table 6.34). Information about normal tissue expression was usually acquired from ‘protein atlas’<sup>252</sup> but when this was not available, data was instead acquired from ‘Genecards’<sup>253</sup>. Known cancer associations were determined using a literature search using the gene ID and the words ‘cancer’ or ‘prostate cancer’. 80% of the top 20 genes up-regulated in extracellular vesicles were associated with prostate tissue, compared with 25% from the cell fraction. 65% of the top 20 genes up-regulated in extracellular vesicles were linked with prostate cancer and 65% cancer generally. The equivalent figures for the cell fraction were 30% and 65%. This is a strong indication that the extracellular vesicles contain RNA from prostate cancer cells and it is a better source of biomarkers than the cell fraction.

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

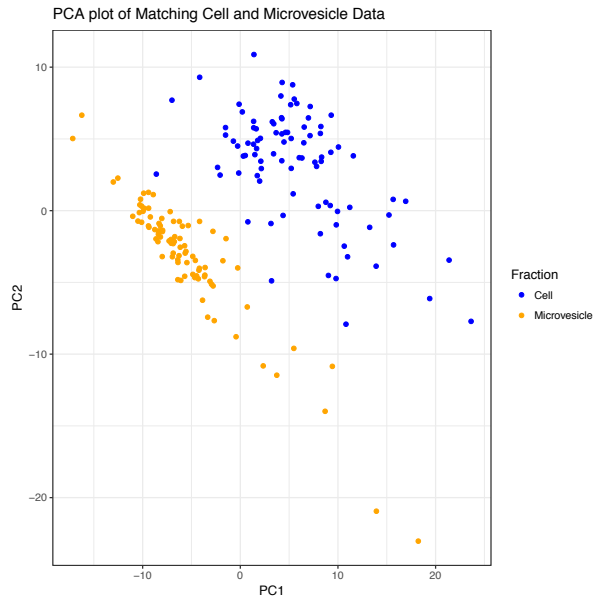
**Table 6.34** A list of the top 20 microarray detected gene-transcripts out of 116 that were found to be significantly more abundant in the cell sediment compared with extracellular vesicles from the same urine.

<i>Rank</i>	<i>Gene</i>	<i>Log<sub>2</sub>(FC)</i>	<i>p-value</i>	<i>Tissue Expression</i>	<i>Known Cancer Associations</i>
1	<i>SCARNA9</i>	7.86	0.035		None
2	<i>SNORD58A, SNORD58B</i>	7.77	0.043		None
3	<i>ALOX5AP</i>	7.16	0.042	Other	Prostate <sup>254</sup> , Other <sup>255</sup>
4	<i>LYZ</i>	7.12	0.035	Other	Other <sup>256</sup>
5	<i>FCER1G</i>	7.00	0.035	Prostate, Other	Other <sup>257</sup>
6	<i>FCGR2A</i>	6.76	0.048	Prostate, Other	None
7	<i>CYBB</i>	6.72	0.045	Other	Prostate <sup>258</sup> , Other <sup>259</sup>
8	<i>TNFRSF1B</i>	6.71	0.045	Other	Other <sup>260</sup>
9	<i>SCARNA9</i>	6.45	0.044		None
10	<i>SRGN</i>	6.43	0.045	Other	Other <sup>261</sup>
11	<i>IL8</i>	6.12	0.035	Other	Prostate <sup>262</sup>
12	<i>EVI2B</i>	6.03	0.043	Other	Other <sup>263</sup>
13	<i>TREM1</i>	5.97	0.044	Other	Other <sup>264</sup>
14	<i>MIR21</i>	5.67	0.049	Not found	Prostate <sup>265</sup> , Other <sup>266, 267</sup>
15	<i>SCARNA7</i>	5.66	0.043	Not found	None
16	<i>HNRNPK</i>	5.62	0.050	Prostate, Other	Prostate <sup>268</sup> , Other <sup>269</sup>
17	<i>GNS</i>	5.58	0.035	Prostate, Other	None
18	<i>CBX3</i>	5.32	0.045	Prostate, Other	Other <sup>270, 271</sup>
19	<i>CTSS</i>	5.32	0.045	Other	Prostate <sup>272</sup> , Other <sup>273, 274</sup>
20	<i>ERO1L</i>	5.23	0.041	Other	Other <sup>275</sup>

### 6.5.2 NanoString comparison of the global expression profile of Extracellular vesicle and Cell fractions.

#### 6.5.2.1 Visualisation of expression differences between fractions

NanoString data (167 probes) from both extracellular vesicle and cell fractions were available for 92 patients. In this section NanoString internal positive control normalised data was used. A PCA plot (section 2.5.1) was produced to visualise the variance of the cell sediment expression against extracellular vesicles expression (Figure 6.1). The expression profiles for the fractions cluster together, indicating that fraction has a bigger influence on the expression profile than the patient.



**Figure 6.1** PCA plot of the expression levels for samples taken from the cell sediment and the extracellular vesicle fraction of urine.

### 6.5.2.2 Differentially expressed transcripts

Expression of 142/167 transcripts were significantly different between extracellular vesicle and cell fractions (adjusted  $p < 0.05$ , paired Mann Whitney U test). 100 were up-regulated in the extracellular vesicle fractions and 42 in the cell sediment fractions (Table 6.35, Table 6.36). *HOXC6* is a known PCa biomarker that can be identified in patient urine<sup>276</sup>, it is therefore very interesting that it is found in abundance in EVs over whole urine. *PTPRC* is a positive regulator of T-cell coactivation and is found in immune cells<sup>277</sup>.

CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

Table 6.35 NanoString top twenty transcripts that were up-regulated in extracellular vesicle fractions compared to cell sediment fractions.

<b>Transcript</b>	<b><i>p</i>-value</b>	<b>Adjusted <i>p</i>-value</b>	<b>Log Fold change</b>
<i>HOXC6</i>	1.12x10 <sup>-10</sup>	1.07x10 <sup>-08</sup>	0.81
<i>SERPINB5</i>	9.78x10 <sup>-14</sup>	1.23x10 <sup>-11</sup>	0.79
<i>OR52A2</i>	3.57x10 <sup>-13</sup>	4.22x10 <sup>-11</sup>	0.77
<i>PTN</i>	1.17x10 <sup>-15</sup>	1.73x10 <sup>-13</sup>	0.76
<i>SChLAP1</i>	8.67x10 <sup>-10</sup>	7.89x10 <sup>-08</sup>	0.67
<i>P712P</i>	6.29x10 <sup>-15</sup>	8.69x10 <sup>-13</sup>	0.67
<i>PPFIA2</i>	2.07x10 <sup>-12</sup>	2.32x10 <sup>-10</sup>	0.65
<i>SIM2 long</i>	1.68x10 <sup>-11</sup>	1.74x10 <sup>-09</sup>	0.65
<i>ERG3' exons 4-5</i>	2.31x10 <sup>-05</sup>	0.0013	0.64
<i>SMIM1</i>	1.69x10 <sup>-13</sup>	2.09x10 <sup>-11</sup>	0.62
<i>TMEM47</i>	0.001	0.0434	0.61
<i>CLU</i>	2.99x10 <sup>-06</sup>	0.0002	0.61
<i>Timp4</i>	6.74x10 <sup>-11</sup>	6.67x10 <sup>-09</sup>	0.61
<i>ARHGEF25</i>	6.56x10 <sup>-10</sup>	6.10x10 <sup>-08</sup>	0.58
<i>RNF157</i>	3.58x10 <sup>-07</sup>	2.47x10 <sup>-05</sup>	0.58
<i>PCA3</i>	7.66x10 <sup>-14</sup>	9.73x10 <sup>-12</sup>	0.58
<i>NKAIN1</i>	1.07x10 <sup>-13</sup>	1.34x10 <sup>-11</sup>	0.57
<i>DNAH5</i>	5.02x10 <sup>-09</sup>	4.22x10 <sup>-07</sup>	0.57
<i>KLK2</i>	8.19x10 <sup>-16</sup>	1.24x10 <sup>-13</sup>	0.55
<i>SYNM</i>	4.87x10 <sup>-08</sup>	3.55x10 <sup>-06</sup>	0.54

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

Table 6.36 NanoString top twenty transcripts that were up-regulated in cell sediment fractions compared to extracellular vesicle fractions.

Transcript	<i>p</i> -value	Adjusted <i>p</i> -value	Log Fold change
<i>PTPRC</i>	1.77x10 <sup>-16</sup>	2.94x10 <sup>-14</sup>	-1.97
<i>STOM</i>	2.93x10 <sup>-16</sup>	4.66x10 <sup>-14</sup>	-1.73
<i>SULF2</i>	2.48x10 <sup>-16</sup>	4.07x10 <sup>-14</sup>	-1.69
<i>MFSD2A</i>	3.24x10 <sup>-16</sup>	5.12x10 <sup>-14</sup>	-1.66
<i>NLRP3</i>	6.72x10 <sup>-16</sup>	1.03x10 <sup>-13</sup>	-1.64
<i>PSTPIP1</i>	3.96x10 <sup>-16</sup>	6.17x10 <sup>-14</sup>	-1.44
<i>MMP25</i>	2.17x10 <sup>-16</sup>	3.58x10 <sup>-14</sup>	-1.43
<i>CLIC2</i>	1.52x10 <sup>-15</sup>	2.20x10 <sup>-13</sup>	-1.35
<i>CCDC88B</i>	2.93x10 <sup>-16</sup>	4.66x10 <sup>-14</sup>	-1.27
<i>TMEM86A</i>	9.51x10 <sup>-15</sup>	1.27x10 <sup>-12</sup>	-1.19
<i>MKi67</i>	3.85x10 <sup>-09</sup>	3.27x10 <sup>-07</sup>	-1.18
<i>MAK</i>	1.39x10 <sup>-14</sup>	1.85x10 <sup>-12</sup>	-1.15
<i>MCTP1</i>	2.83x10 <sup>-16</sup>	4.59x10 <sup>-14</sup>	-1.09
<i>APOC1</i>	6.72x10 <sup>-16</sup>	1.03x10 <sup>-13</sup>	-1.07
<i>CP</i>	4.49x10 <sup>-11</sup>	4.54x10 <sup>-09</sup>	-0.99
<i>MIR146A</i>	1.74x10 <sup>-15</sup>	2.47x10 <sup>-13</sup>	-0.96
<i>NEAT1</i>	1.77x10 <sup>-16</sup>	2.94x10 <sup>-14</sup>	-0.88
<i>Met</i>	9.62x10 <sup>-12</sup>	1.01x10 <sup>-09</sup>	-0.88
<i>MIC1</i>	4.15x10 <sup>-13</sup>	4.85 x10 <sup>-11</sup>	-0.67
<i>COL10A1</i>	2.09x10 <sup>-11</sup>	2.15x10 <sup>-09</sup>	-0.59

### 6.6 Discussion

I found that the AUCs of the cell sediment models were marginally higher in the baseline normalised data for CB vs cancer models, CB vs high-risk cancer models and CB > L > I > H trend models (Table 6.30). However, these AUCs need to be taken with caution as the models have not been tested in a validation dataset and so overfitting may be occurring. There was a low number of samples used to build the cell predictive models and they all came from the same centre, so it is possible that the cell models are not as robust as one would desire. Comparing the transcripts identified via glm and those identified by Mann Whitney U, there were a large percentage of transcripts in common, suggesting a level of robustness when using these methods.

For the cell sediment, the transcripts identified as significantly different between the “No Evidence for Cancer” samples and the cancer samples differed depending on the normalisation (Supplementary Table 49). However, *CADPS* and *ERG3*’ exons 6-7, *HOXC6*, *NAALADL2* and *TDRD* are present in all analyses. This shows a robustness of these

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

transcripts and indicates a level of importance when using cell sediment from urine samples. All are up-regulated in the cancer samples, with the exception of *NAALADL2*. *ERG3'* exons 6-7, *HOXC6* and *SIM2* long were the only probes that consistently distinguished high-risk cancer from CB samples. All three probes were up-regulated in the high-risk cancer samples. Looking at the trend across clinically benign, low-risk, intermediate-risk and high-risk cancer, *CADPS*, *ERG3'* exons 6-7, *HOXC6* and *NAALADL2* probes were again the common transcripts across all of the different normalisations. *CADPS* increases as risk increases, with lowest expression in CB and highest in high-risk cancer. *ERG 3'* exons 6-7 and *HOXC6* increase in low risk cancer but then have a decreased expression in intermediate and high-risk cancer with lowest expression found in CB. *NAALADL2* expression decreases in trend with advancement of cancer.

*CADPS* is a cytosolic and peripheral membrane protein required for vesicle docking and priming steps that precede vesicle exocytosis<sup>278</sup>. Down-regulation of *CADPS* has been associated with poor outcome in pancreatic ductal adenocarcinoma<sup>279</sup> and a genome wide molecular characterisation of central nervous system primitive neuroectodermal tumour and pineoblastoma found that the *CADPS* locus (3p14.2) was lost in 27.6% of cases and was also associated with poor prognosis<sup>280</sup>. Searching for “*CADPS* prostate cancer” yields no results during a literature search. *ERG 3'* is a proto-oncogene known to be associated with PCa, it is also involved in the *TMPRSS2:ERG* fusion but has been shown to be increased in PCa via alternate mechanisms to the fusion also<sup>281</sup>. The *TMPRSS2:ERG* fusion is identified in ~50% of PCa samples but has not been identified as a key biomarker for PCa prediction in cell data in this study. *HOXC6* is known to be associated with PCa, there is a urine based test that utilises the identification of *HOXC6* mRNA called the SelectMDx<sup>276</sup>. Therefore, our findings support other work showing its association with prostate cancer and its identification in PCa patient urine. *NAALADL2* is known to be overexpressed in PCa tissue compared to benign tissue using IHC. Expression of *NAALADL2* has been shown to impact on a number of pro-oncogenic pathways such as cell migration, invasion and colony-forming potential. Leading to the belief that *NAALADL2* is a useful biomarker for diagnosis

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

and prognosis<sup>282</sup>. In PCa cell urine, a lower expression has been associated with PCa, here. *TDRD* is a cancer/testis (CT) antigen that has previously been associated with liver cancer<sup>283</sup> and breast cancer<sup>284</sup> but not PCa. *SIM2* long has been found to be over-expressed in PCa tissue when compared with CB tissue<sup>285</sup> and its up regulation has also been associated with biochemical recurrence post-radical prostatectomy<sup>286</sup>. However, it has not previously been identified as a PCa urine biomarker.

In the EV models, *ERG* 3' exon 4-5 was more highly selected as a biomarker over *ERG* 3' 6-7 like in the cell models. *HOXC6*, *TDRD* also appear in all of the EV models. However, *TMPRSS2:ERG*, *PCA3* also appear in EV models and not cell models. It has previously been observed that most of the RNA content in whole urine is actually coming from EVs and not from cells – this is shown by a comparison of the RNA yields from cells and EVs from the same urine samples (data not shown). Further to this, NanoString analysis was only performed on 95 cell RNA fractions out of the 756 Samples because amounts of cell RNA were on the whole so limiting that expression analysis was not deemed viable. Consistently higher EV RNA yields explains how *TMPRSS2:ERG* and *PCA3* are highly detectable in whole urine and EV fractions of urine. *CADPS* is not selected in EV models, which makes sense as *CADPS* is an EV making gene. However, many more transcripts are commonly selected such as *APOC1*, *KLK4* and *HPN*. Showing EVs are a good source of urinary biomarkers for PCa.

Comparing the expression of transcripts in the cell fraction to the EV fraction via microarray has shown that a high proportion of prostate, and PCa associated transcripts are more abundant in the EV fraction. It also showed that *PTPRC*, which is a blood immune associated transcript is more abundant in the cell fraction. This us to believe that many of the cells in the cell fraction of the PCa patient urine are actually immune cells and not prostate/PCa cells, which is in support of other literature<sup>281</sup>.

Our findings support previous research that the genetic content of cell sediment and that of extracellular vesicles differs. Expression levels of many transcripts that were both expressed in the prostate tissue and known to be prostate cancer associated were found in

## CHAPTER 6: EXPRESSION PROFILE OF THE CELL SEDIMENT URINE FRACTION

increased levels in the extracellular vesicle compared to the cell sediment. This highlights that the extracellular vesicle fraction is indeed of great interest to investigate further for PCa biomarkers.



# 7

## Discussion

### **6.7 Summary**

Prostate Cancer (PCa) is a major clinical problem worldwide with considerable variability in clinical outcome of patients. PCa diagnostics and prognostics currently lack specific and sensitive clinical biomarkers and treatment is not well individualised. The *PCA3* test, amongst others, highlights the utility of urine in PCa diagnostics and prognostics<sup>214</sup>. The extracellular vesicle (EV) fraction contains exosomes and is obtainable from urine. Exosome levels are known to be increased during malignancy and those produced by tumours contain nucleic material from malignant cells<sup>104</sup>. EVs from tumour cells have roles involved in tumourigenesis, metastasis, and response to therapy by triggering signalling cascades and transferring mRNA, miRNA and proteins between cancer cells and the tumour microenvironment<sup>105</sup>. Our aim was to interrogate PCa patient's urine samples, mostly the EV fraction to identify novel biomarkers or sets of biomarkers to aid in PCa management. This study was completed as part of the Movember GAP1 global PCa biomarker initiative, which involved multiple collaborators and samples collected from four different centres worldwide, for the identification of urinary biomarkers for the risk-stratification of PCa.

#### **7.1.1. Chapter 3: NanoString Data Analysis 1: The Pilot Study**

## CHAPTER 7: DISCUSSION

In a pilot study, NanoString technology was able to detect PCa specific markers in 196 samples, such as *TMPRSS2:ERG*, which was detected in 58% of all PCa samples and in 19% of samples from men with no clinical evidence of PCa (CB). Latent Process Decomposition unsupervised analysis clustered the EV expression data into four groups, which was associated with clinical risk categories ( $p < 0.05$ ). Transcripts were identified that were differentially expressed and models were built that could distinguish between PCa and samples that showed no evidence of PCa (CB) with an AUC of 0.937, high-risk PCa and samples showing no evidence of PCa (CB) with an AUC of 0.852 and metastatic PCa (A) and samples showing no evidence of PCa (CB) with an AUC of 0.983. These findings highlight that the transcript data collected from urinary EVs in PCa patients comes, at least in part, from the prostate and holds clinically relevant structure.

### **6.8 Chapter 4: NanoString2 Analysis: The Movember GAP1 Project**

Following on from the pilot study, further samples ( $n = 756$ ) obtained from four centres worldwide were sent to NanoString for the quantification of 167 transcripts. The aims were to primarily identify optimal models capable of predicting PCa and to risk-stratify PCa without the need for biopsy. Models were built to answer four important clinical questions:

- 1) Determine which samples were from PCa and which were from samples with no evidence of Ca (AUC = 0.851).
- 2) Determine which samples were from high-risk PCa only and which were from samples with no evidence of cancer, (AUC = 0.897).
- 3) Determine if there was a trend in expression that corresponds to a trend in risk category (CB>L>I>H), (AUC = 0.709).
- 4) Determine if there was a trend in expression that corresponds to a trend in patient type (CB>Ca>Metastatic cancer), (AUC = 0.6469).

## CHAPTER 7: DISCUSSION

The data was stratified into training and test sets in the ratio 2:1, models were built with the training set and validated using the test set. I used four different normalisations of the data, which included using *KLK2* ratio, *KLK2* adjusting, *KLK3* adjusting, and *GAPDH & RPLP2* normalisation. Models built using the *GAPDH & RPLP2* normalised data generally had higher AUCs. These models are improvements on existing tests and have the potential to be developed in to clinical tests.

### 6.8.1 Chapter 5: Response to treatment

Many cancers have benefitted from treatment stratification due to expression of certain genes, however with the exception of the DESNT poor prognosis expression group, this has not yet been done for PCa. With hormone therapy (HT) it is known that patients will inevitably progress to castration resistant prostate cancer (CRPC). How long each patient will last on HT varies widely from months to years. Samples from the advanced patients in the NanoString pilot study ( $n = 32$ ) were used to identify a significant predictor of early progression in patients on HT: A signature of seven transcripts was identified that could optimally predict progression of patients on hormone therapy (cox-regression model;  $p = 2.3 \times 10^{-05}$ ; HR = 0.04288). The transcripts in the predictor were *AGR2*, *DLX1*, *KLK2*, *NAALADL2*, *AR* exons 4-8, *PPAP2A* and *AMACR*. This model was an independent predictor of progression when established clinical variables initial PSA, age, Gleason score and initial bone scan result were taken into account (cox-regression model;  $p = 0.003$ ; HR = 0.03). When the data was adjusted to *KLK2* levels, similar to *KLK3* adjustment used in the PCA3 test, an optimal model of three transcripts (*CAMKK2*, *PSGR* and *UPK*) was identified (cox-regression model;  $p = 0.007$ , HR = 1.0028). This model does not remain significant predictor when adding clinical factors (cox-regression model;  $p = 0.14$ ; HR= 1.009). However when both of these models were applied to the second NanoString dataset but they were not validated. Despite this, I have shown the potential of using urine extracellular vesicles from prostate cancer patients with NanoString measurements of expression to

predict patient response to treatments. A larger cohort with longer follow up would be required to further develop these models in to something usable in the clinic.

### **6.8.2 Chapter 6: Analysis of Cell Fraction and comparison with EV fraction**

The transcriptome profiles of cell sediment and EV fractions were compared from PCa patients and controls (taken from patients with no evidence of cancer (CB)). Data from microarray of samples collected from NNUH, Norwich and Royal Marsden Hospital, London was used for this comparison. 98 genes were found to be significantly ( $p < 0.05$ ) up-regulated in the extracellular vesicles and 116 up-regulated in the cell sediment fraction. 92 samples from the NanoString 2 experiment were also EV and cell sediment matched and were also used to compare transcriptome profiles. 100 genes were found to be significantly ( $p < 0.05$ ) up-regulated in the extracellular vesicles and 42 genes were up-regulated in the cell sediment fractions. The top twenty of each set of these genes were investigated for known prostate expression and PCa association. The EV fraction contained higher levels of prostate expressed and PCa associated transcripts. This is a strong indication that the EVs contain RNA from prostate cancer cells and it is a better source of biomarkers than the cell fraction.

The NanoString data from cell sediment was used to produce models able to predict PCa (low, intermediate and high-risk) from CB samples, high-risk PCa from CB samples and trend in expression across clinical category. These models had similar AUCs in the training set to the EV fractions but we were unable to validate them at this stage. The power of these cell fraction models is also reduced due to a much lower sample size.

### **6.9 Discussion**

There is an urgent clinical need for biomarkers to determine which patients have PCa, which patients have disease that will progress rapidly, and to individualise treatment to optimise response. Lung cancer and breast cancer are already benefitting from individualised treatment based on expression levels<sup>10,14</sup>. For PCa, stratification models have been produced that include a number of clinical factors, these include D'Amico, and nomograms or points systems such as CAPRA<sup>287</sup>, the Prostate Health Index<sup>288</sup>, the European Randomised Study of Screening for Prostate Cancer (ERSPC) Risk Calculator<sup>289</sup>, the Prostate Cancer Prevention Trial Risk Calculator (PCPT-RC)<sup>290</sup>. However, apart from D'Amico, none of these risk calculators are in general use in the clinic, and effectiveness varies with the cohort<sup>290</sup>. The production of the PCA3 test<sup>214</sup> has led to an increase in studies investigating urine as a source of PCa biomarkers for clinical tests that may prevent unnecessary biopsies. Further research has merged clinical data with urine expression information such as MiPS<sup>213</sup>, which built on the PCA3 urine test to include other urine expression data (*TMPRSS2:ERG*) and PSA. There is also a model for predicting high-grade PCa using *HOXC6* and *DLX1* urinary expression levels along with clinical factors such as prostate volume<sup>215</sup>, the ExoDx Prostate (IntelliScore "EPI")<sup>216</sup> which can be used in conjunction with clinical data, and the PCRT-RC which is designed to incorporate future biomarker information as it becomes available. *TMPRSS2:ERG* fusions are only found in ~50% of PCa tumours and PCA3 is not expressed in all PCa tumours also. A panel of more transcripts may improve the diagnostic and prognostic abilities of these tests.

EVs have been investigated as a source of urinary biomarkers in renal cancer studies and it was found that the RNA profile was better preserved in urinary

## CHAPTER 7: DISCUSSION

microvesicles compared with whole cells<sup>217</sup>. It has been suggested that this is because the EV membrane may protect RNA from degradation in urine<sup>217</sup>.

Identification of PCa biomarkers in EVs was subsequently observed<sup>218</sup>.

Biomarker discovery in urine of prostate cancer patients has so far focused on just a few gene targets. In this study we have taken a more holistic, but not transcriptome wide approach increasing the number of probes significantly. I have shown that NanoString is a viable technology to measure 100s of probes in urine efficiently and a viable solution for biomarker discovery and potential implementation in the clinic. I have produced potentially important combinations of biomarkers to predict prostate cancer, aggressive prostate cancer, and response to hormone therapy treatment. These gave AUCs up to 0.897, which is an improvement on published tests in the literature. The translational appeal of NanoString analysis can be seen in the ProSigna PAM50 test for aggressive breast cancer<sup>291</sup>, which uses NanoString technology and is commercially available.

In our study I have shown that the transcriptome profile of whole cells and EVs differs and that EVs are a potential better source of PCa biomarkers as they contain more prostate derived transcripts as well as more PCa associated and cancer associated transcripts. This indicates that using EVs in biomarker discovery in urine will improve results, but it is likely that whole urine could be used in a final test. Biologically it is likely that EVs can find their way into the urine more easily than the bulky cell counterparts. We also observed that the whole urine includes many white blood cells. Recent research has shown that WBC can be utilised as prognostic markers in BCa showing capability for predicting distant metastasis preoperatively over a 65-month timeline. Increased platelet indices and decreased neutrophil numbers were associated with a poorer prognosis<sup>292</sup>.

I have identified urinary EV models from NanoString data capable of predicting PCa, and PCa risk categories with AUCs similar to previously published urine models. Which include

## CHAPTER 7: DISCUSSION

both known PCa associated transcripts from whole urine and novel transcripts that may be EV specific. The AUCs of cell models and EV models are very similar and thus it may be that a combinatory model could be better to predict PCa and its prognosis.

The need for cancer specific biomarkers for assessing response to hormonal treatments in metastatic PCa has been acknowledged<sup>293</sup>, yet very little work appears to have been completed in this area. I identified two signatures capable of distinguishing early relapse to HT in two different data normalisations. However these signatures were not validated in a second dataset.

A urine test would aid clinicians and patients for the management of PCa in a few areas. Firstly, there is a decision of whether a biopsy needs to be undertaken. Usually, this is based on serum PSA level and DRE findings. A urine test could help limit the amount of unnecessary biopsies conducted. Secondly, it is known that biopsies generally under grade the PCa, and higher Gleason scores are identified on whole prostates from radical prostatectomies. Therefore, a urine test may help to identify which patients can safely go on to active surveillance.

A third area where a urine test could aid in the clinic is alongside MRIs. MRIs have shown great potential in the diagnosis of PCa but does suffer from a high false positive rate (~50%)<sup>294</sup>. Introduction of a urine test alongside MRI could help to reduce the false positive rate especially for PIRADS  $\leq 4$ <sup>295</sup> (Figure 0.1).

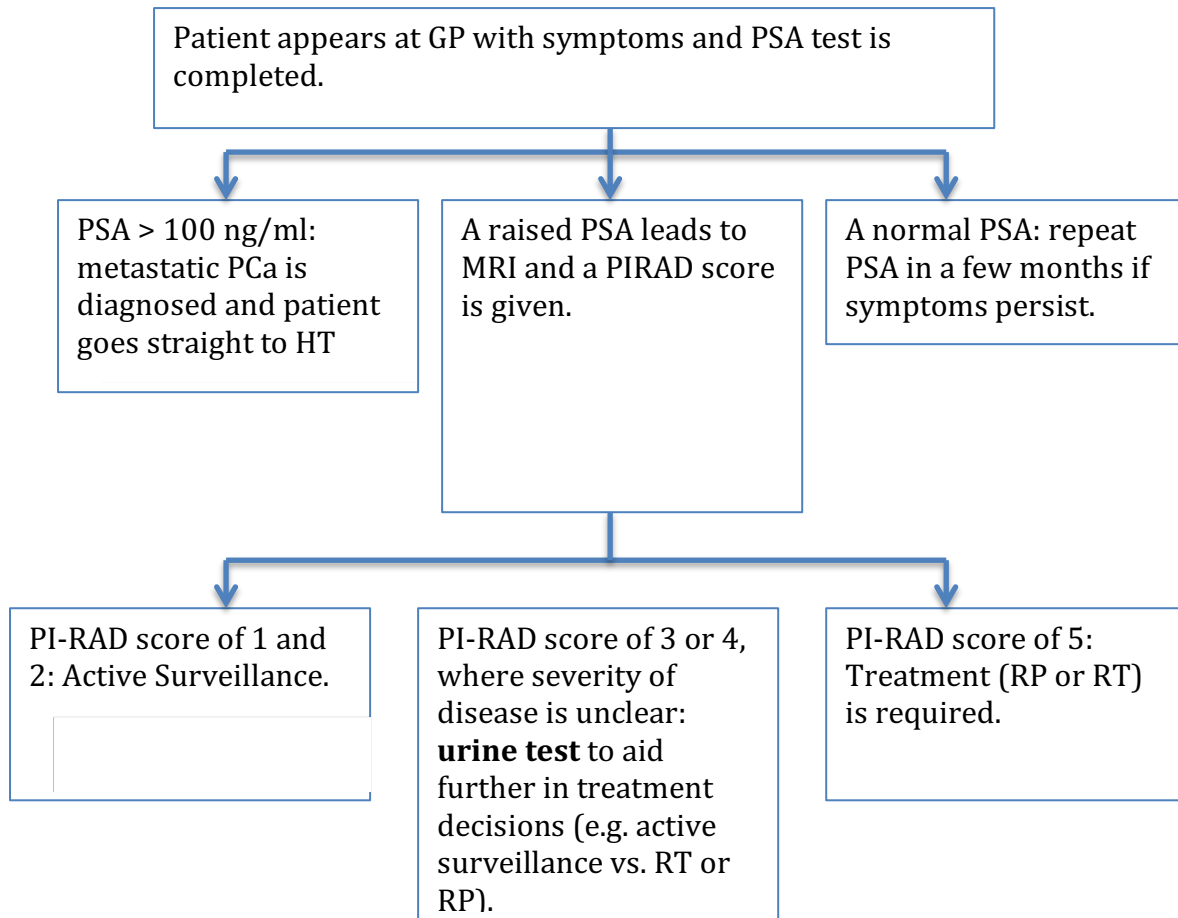


Figure 0.1 A flow diagram showing where the urine test would be best utilised in current diagnostic procedures.

### 6.10 Future Work

Further work needs to be performed on the models produced for determining between PCa and CB samples, high-risk PCa and CB samples as well as risk category (CB>L>I>H) and patient type (CB> cancer > metastatic cancer) trends. In particular, an immediate next step should be that they are incorporated with clinical factors to identify if they outperform clinical factors alone, as well as if they have a better prediction when including the clinical factors (including but not limited to prostate volume, age, family history, previous biopsy results and serum PSA). The next major step would be to validate these models in an independent large-scale trial such as PROMIS or PROTECT. If this was successful, then the



## CHAPTER 7: DISCUSSION

models would be evaluated in a large multi-centre prospective study. This is necessary to obtain FDA approval and translation in to a test used in the clinic.

Additional work to optimise the methodology used to collect urine, to standardise it, (simplify and make it more robust). Currently, samples have to be processed within 2 hours of collection, the introduction of urine preservatives could also streamline procedures. The models would need to be tested to see if they worked in whole urine and without DRE, which would make the collection and processing methodology a lot simpler. Comparing alternative methods for the quantification of transcripts from urinary EVs may also help to improve the reliability and clinical use of the models.

Further work needs to be completed to identify a robust and validated signature for the prediction of early relapse to CRPC. This is a vital area that needs improvement for the clinical management of PCa. A larger cohort with longer follow up is required. I would also like to look at the data from patients on active surveillance in the NanoString 2 data set. There is considerable potential to develop a predictor of time to treatment in these patients. Another response to treatment that should be investigated is biochemical recurrence (BCR) after radical prostatectomy or radiation therapy. Unfortunately, our follow up was not long enough to have sufficient numbers of patients that suffered from BCR to be able to perform any of these experiments at this time. I would also like to examine whether models that were developed for the prediction of aggressiveness could be applied to predict response to treatment. For example, could the optimal model for predicting risk category also be used to predict time to treatment for patients on active surveillance. In this whole project I have been reliant on the 167 gene probes used in the NanoString assay. It is not clear whether these are the optimal probes to use,

## CHAPTER 7: DISCUSSION

although it is apparent that they are at least sufficient for some clinical questions. I would like to perform a similar scale project but using a global transcriptome approach using microarrays, or for an exon splice variant analysis then RNAseq would be ideal, though highly analysis intensive. This would allow us to identify the very best probes to use in a clinical test to answer the important clinical questions in prostate cancer.

Due to my work in this thesis re-funding has been awarded for the development of a clinically implementable Prostate Urine Risk test. This has resulted in two further PhD posts one for lab work and one for bioinformatics.

### **6.11 Conclusions**

O1: To determine whether RNA expression from urine extracellular vesicles in prostate cancer patients are a viable target for the development of biomarkers through the use of NanoString technology.

I have shown that urine extracellular vesicles from prostate cancer patients contain information from tumours and are a viable area to investigate for non-invasive biomarkers. I have shown that NanoString technology is sensitive and specific enough to use as a semi-high throughput approach for discovery and potentially for clinical use.

O2: To determine an optimal combination of probes to predict cancer presence and aggression in prostate cancer patients.

I have determined a number of models that work extremely well in predicting both cancer presence and the aggressiveness of disease. These have the potential, with further work, to have an impact in the clinic. Models to accurately stratify patients' disease into D'Amico risk groupings were less satisfactory and may require alternative probes or other techniques.

## CHAPTER 7: DISCUSSION

O3: To determine whether an optimal combination of probes can predict response to hormone therapy treatment.

I have shown that there may be some information in urine extracellular vesicles to predict patient response to treatments. I have developed some potential tests, but for confidence in these a much bigger data set with longer follow up would be required.

O4: To evaluate the differences between urine fractions (extracellular vesicles and cell sediment) and determine whether cell sediment can be used to predict cancer presence and aggression in prostate cancer patients.

I have shown that there are considerable differences between the extracellular vesicles fraction and the cell sediment fraction of urine collected from prostate cancer patients. There is a strong indication that the EVs contain more RNA from prostate cancer cells and it is a better source of biomarkers than the cell fraction. Despite this, I was able to produce some models that were reasonable good at detecting the presence and aggressiveness of prostate cancer.

In this thesis, I have shown that by interrogating the EV fraction of PCa patient's urine samples using NanoString technology that novel biomarkers or sets of biomarkers can be identified to aid in PCa management in a non-invasive test.

# 8

## References

1. Leongamornlert, D. *et al.* Frequent germline deleterious mutations in DNA repair genes in familial prostate cancer cases are associated with advanced disease. **110**, 1663–1672 (2014).
2. Akin, O. *et al.* Transition zone prostate cancers: features, detection, localization, and staging at endorectal MR imaging. *Radiology* **239**, 784–792 (2006).
3. Olson, S., Robinson, S. & Giffin, R. *Accelerating the Development of Biomarkers for Drug Safety: Workshop Summary. Sciences-New York* (2009). doi:20464768
4. World Health Organisation. WHO International Programme on Chemical Safety Biomarkers and Risk Assessment: Concepts and Principles. (1993).
5. Strimbu, K. & Tavel, J. A. What are biomarkers? *Curr. Opin. HIV AIDS* **5**, 463–466 (2010).
6. Fuzery, A. K., Levin, J., Chan, M. M. & Chan, D. W. Translation of proteomic biomarkers into FDA approved cancer diagnostics: issues and challenges. *Clin. Proteomics* **10**, 13 (2013).
7. Antoniou, A. *et al.* Average risks of breast and ovarian cancer associated with BRCA1 or BRCA2 mutations detected in case Series unselected for family history: a combined analysis of 22 studies. *Am. J. Hum. Genet.* **72**, 1117–1130 (2003).
8. Markert, E. K., Mizuno, H., Vazquez, A. & Levine, A. J. Molecular classi

- fication of prostate cancer using curated expression signatures. (2011). doi:10.1073/pnas.1117029108/-/DCSupplemental.www.pnas.org/cgi/doi/10.1073/pnas.1117029108
9. Heinemann, V., Stintzing, S., Kirchner, T., Boeck, S. & Jung, A. Clinical relevance of EGFR- and KRAS-status in colorectal cancer patients treated with monoclonal antibodies directed against the EGFR. *Cancer Treat. Rev.* **35**, 262–271 (2009).
  10. Rosell, R., Bivona, T. G. & Karachaliou, N. Genetics and biomarkers in personalisation of lung cancer treatment. *Lancet* **382**, 720–31 (2013).
  11. Misale, S. *et al.* Emergence of KRAS mutations and acquired resistance to anti EGFR therapy in colorectal cancer. *Nature* **486**, 532–536 (2014).
  12. Diamandis, E. P. Cancer Biomarkers: Can We Turn Recent Failures into Success? *JNCI J. Natl. Cancer Inst.* **102**, 1462–1467 (2010).
  13. Sawyers, C. L. The cancer biomarker problem. *Nature* **452**, 548–552 (2008).
  14. Eaton, A. A. *et al.* Estimating the OncotypeDX score : validation of an inexpensive estimation tool. *Breast Cancer Res. Treat.* (2016). doi:10.1007/s10549-016-4069-4
  15. Marrone, M., Potosky, A. L., Penson, D. & Freedman, A. N. A 22 Gene-expression Assay, Decipher® (GenomeDx Biosciences) to Predict Five-year Risk of Metastatic Prostate Cancer in Men Treated with Radical Prostatectomy. *PLoS Curr.* **7**, ecurrents.eogt.761b81608129ed61b0b48d42c04f92a4 (2015).
  16. Cuzick, J. *et al.* Prognostic value of a cell cycle progression signature for prostate cancer death in a conservatively managed needle biopsy cohort. 1095–1099 (2012). doi:10.1038/bjc.2012.39
  17. Cancer Research UK. Molecular diagnostic provision in England for targeted cancer medicines (solid tumours) in the NHS. (2015).
  18. Olzscha, H., New, M. & La Thangue, N. B. Personalised Cancer Medicine: Fulfilling the Promise. *Encycl. Life Sci.* 1–10 (2013). doi:10.1002/9780470015902.a0025180
  19. Mitri, Z., Constantine, T. & O'Regan, R. The HER2 Receptor in Breast Cancer: Pathophysiology, Clinical Use, and New Advances in Therapy. *Chemother. Res. Pract.* **2012**, 1–7 (2012).
  20. Crown, J., O'Shaughnessy, J. & Gullo, G. Emerging targeted therapies in triple-negative breast cancer. *Ann. Oncol.* **23**, vi56-vi65 (2012).
  21. Giaccone, G. & Rodriguez, J. A. EGFR inhibitors: what have we learned from the treatment of lung cancer? *Nat. Clin. Pract. Oncol.* **2**, 554–61 (2005).
  22. Dziadziuszko, R. & Jassem, J. Epidermal growth factor receptor (EGFR) inhibitors and derived treatments. *Ann. Oncol.* **23 Suppl 1**, x193-6 (2012).
  23. Velonas, V. M., Woo, H. H., dos Remedios, C. G. & Assinder, S. J. Current status of biomarkers for prostate cancer. *Int. J. Mol. Sci.* **14**, 11034–11060 (2013).
  24. Ferlay, J. *et al.* Cancer incidence and mortality patterns in Europe: Estimates for 40 countries in 2012. *Eur. J. Cancer* **49**, 1374–1403

## 9: APPENDICES

- (2013).
25. Lloyd, T. *et al.* Lifetime risk of being diagnosed with, or dying from, prostate cancer by major ethnic group in England 2008–2010. *BMC Med.* **13**, 171 (2015).
  26. UK, C. R. Cancer Research UK. *Prostate cancer statistics* (2014). Available at: <http://www.cancerresearchuk.org/health-professional/cancer-statistics/statistics-by-cancer-type/prostate-cancer>. (Accessed: 1st January 2015)
  27. Oh WK, Hurwitz M, D. A. *Biology of Prostate Cancer. Hollan-Frei Cancer Medicine.* (BC Decker, 2003).
  28. Cancer Research UK. *Prostate Cancer: Types and grades* (2016). Available at: <http://www.cancerresearchuk.org/about-cancer/prostate-cancer/types-grades>. (Accessed: 1st January 2016)
  29. UK, C. R. Cancer Research UK. *Prostate cancer statistics* (2014).
  30. Vellekoop, A. & Loeb, S. More Aggressive Prostate Cancer in Elderly Men. *Rev. Urol.* **15**, 202–204 (2013).
  31. Gann, P. H. Risk factors for prostate cancer. *Rev. Urol.* **4 Suppl 5**, S3–S10 (2002).
  32. Madu, C. O. & Lu, Y. Novel diagnostic biomarkers for prostate cancer. *J. Cancer* **1**, 150–177 (2010).
  33. Ito, K. Prostate cancer in Asian men. *Nat Rev Urol* **11**, 197–212 (2014).
  34. Ito, K. Prostate cancer in Asian men. *Nat Rev Urol* **11**, 197–212 (2014).
  35. Eeles, R. *et al.* The genetic epidemiology of prostate cancer and its clinical implications. *Nat Rev Urol* **11**, 18–31 (2014).
  36. Al Olama, A. A. *et al.* A meta-analysis of 87,040 individuals identifies 23 new susceptibility loci for prostate cancer. *Nat. Genet.* **46**, 1103–9 (2014).
  37. Al Olama, A. A. *et al.* A meta-analysis of 87,040 individuals identifies 23 new susceptibility loci for prostate cancer. *Nat. Genet.* **46**, 1103–9 (2014).
  38. Catalona, W. J. *et al.* Comparison of digital rectal examination and serum prostate specific antigen in the early detection of prostate cancer: results of a multicenter clinical trial of 6,630 men. *J. Urol.* **151**, 1283–1290 (1994).
  39. Humphrey, P. A. Gleason grading and prognostic factors in carcinoma of the prostate. *Mod. Pathol.* **17**, 292–306 (2004).
  40. National Collaborating Centre for Cancer. Prostate Cancer : diagnosis and treatment. Clinical guideline. *Natl. Inst. Heal. Care Excell.* 1–480 (2014).
  41. (Uk), N. C. C. for C. Prostate Cancer: Diagnosis and Treatment. *Natl. Inst. Heal. Clin. Excell. Guid.* 2007–2009 (2008).
  42. Paller, C. J. & Antonarakis, E. S. Management of biochemically recurrent prostate cancer after local therapy: evolving standards of care and new directions. *Clin. Adv. Hematol. Oncol.* **11**, 14–23 (2013).
  43. Paller, C. J. & Antonarakis, E. S. Management of biochemically recurrent prostate cancer after local therapy: evolving standards of care and new directions. *Clin. Adv. Hematol. Oncol.* **11**, 14–23 (2013).
  44. Wan, X. *et al.* UHRF1 overexpression is involved in cell proliferation

- and biochemical recurrence in prostate cancer after radical prostatectomy. *J. Exp. Clin. Cancer Res.* **35**, 34 (2016).
45. Raatikainen, S., Aaltonen, S., Karja, V. & Soini, Y. Increased Peroxiredoxin 6 Expression Predicts Biochemical Recurrence in Prostate Cancer Patients After Radical Prostatectomy. *Anticancer Res.* **35**, 6465–6470 (2015).
  46. Qu, X. *et al.* Identification of Combinatorial Genomic Abnormalities Associated with Prostate Cancer Early Recurrence. *J. Mol. Diagn.* **18**, 215–224 (2016).
  47. Ma, D. *et al.* Association of molecular biomarkers expression with biochemical recurrence in prostate cancer through tissue microarray immunostaining. *Oncol. Lett.* **10**, 2185–2191 (2015).
  48. Meng, Y., Li, H., Xu, P. & Wang, J. Do tumor volume , percent tumor volume predict biochemical recurrence after radical prostatectomy ? A meta-analysis. **8**, 22319–22327 (2015).
  49. Lorente, J. A., Morote, J., Raventos, C., Encabo, G. & Valenzuela, H. Clinical efficacy of bone alkaline phosphatase and prostate specific antigen in the diagnosis of bone metastasis in prostate cancer. *J. Urol.* **155**, 1348–1351 (1996).
  50. (Uk), N. C. C. for C. Prostate Cancer: Diagnosis and Treatment. *Natl. Inst. Heal. Clin. Excell. Guid.* 2007–2009 (2008).
  51. Perlmutter, M. A. & Lepor, H. Androgen deprivation therapy in the treatment of advanced prostate cancer. *Rev. Urol.* **9 Suppl 1**, S3-8 (2007).
  52. Mostaghel, E. a. Abiraterone in the treatment of metastatic castration-resistant prostate cancer. *Cancer Manag. Res.* **6**, 39–51 (2014).
  53. Schoenborn, J. R., Nelson, P. & Fang, M. Genomic profiling defines subtypes of prostate cancer with the potential for therapeutic stratification. *Clin. Cancer Res.* **19**, 4058–66 (2013).
  54. Mostaghel, E. a. Abiraterone in the treatment of metastatic castration-resistant prostate cancer. *Cancer Manag. Res.* **6**, 39–51 (2014).
  55. Miyamoto, H., Messing, E. M. & Chang, C. Androgen deprivation therapy for prostate cancer: current status and future prospects. *Prostate* **61**, 332–353 (2004).
  56. Chen, C. D. *et al.* Molecular determinants of resistance to antiandrogen therapy. *Nat. Med.* **10**, 33–39 (2004).
  57. Nice. Single Technology Appraisal ( STA ) Cabazitaxel for the second-line treatment of metastatic hormone refractory prostate cancer. 1–164 (2011).
  58. Mikolajczyk, S. D. *et al.* A Precursor Form of Prostate-specific Antigen Is More Highly Elevated in Prostate Cancer Compared with Benign Transition Zone Prostate Tissue A Precursor Form of Prostate-specific Antigen Is More Highly Elevated in Prostate Cancer Compared with Benign Tra. 756–759 (2000).
  59. Thompson, Ian M, Lucia, M. S. *et al.* new england journal. 2239–2246 (2004).
  60. DeAntoni, E. P. *et al.* Age- and race-specific reference ranges for prostate-specific antigen from a large community-based study. *Urology* **48**, 234–239 (1996).

## 9: APPENDICES

61. Romero Otero, J., Garcia Gomez, B., Campos Juanatey, F. & Touijer, K. a. Prostate cancer biomarkers: An update. *Urol. Oncol. Semin. Orig. Investig.* **32**, 252–260 (2014).
62. National Cancer Institute (NCI). Prostate-Specific Antigen (PSA) Test. *Prostate-Specific Antigen (PSA) Test* (2012). Available at: <https://www.cancer.gov/types/prostate/psa-fact-sheet>. (Accessed: 1st January 2015)
63. Tinzl, M., Marberger, M., Horvath, S. & Chypre, C. DD3PCA3 RNA analysis in urine--a new perspective for detecting prostate cancer. *Eur. Urol.* **46**, 182–6; discussion 187 (2004).
64. Catalona, W. J. *et al.* Use of the percentage of free prostate-specific antigen to enhance differentiation of prostate cancer from benign prostatic disease: a prospective multicenter clinical trial. *JAMA* **279**, 1542–1547 (1998).
65. Mikolajczyk, S. D., Marks, L. S., Partin, A. W. & Rittenhouse, H. G. Free prostate-specific antigen in serum is becoming more complex. *Urology* **59**, 797–802 (2002).
66. Mikolajczyk, S. D. *et al.* A Truncated Precursor Form of Prostate-specific Antigen Is a More Specific Serum Marker of Prostate Cancer A Truncated Precursor Form of Prostate-specific Antigen Is a More Specific Serum Marker of Prostate Cancer. 6958–6963 (2001).
67. Guazzoni, G. *et al.* Prostate-specific antigen (PSA) isoform p2PSA significantly improves the prediction of prostate cancer at initial extended prostate biopsies in patients with total PSA between 2.0 and 10 ng/ml: results of a prospective study in a clinical setting. *Eur. Urol.* **60**, 214–22 (2011).
68. Stephan, C. *et al.* A [-2]proPSA-based artificial neural network significantly improves differentiation between prostate cancer and benign prostatic diseases. *Prostate* **69**, 198–207 (2009).
69. National Cancer Institute (NCI). Prostate-Specific Antigen (PSA) Test. *Prostate-Specific Antigen (PSA) Test* (2012).
70. Makarov, D. V *et al.* Management for Prostate Cancer. **15**, 7316–7321 (2010).
71. Wang, Y., Liu, X.-J. & Yao, X.-D. Function of PCA3 in prostate tissue and clinical research progress on developing a PCA3 score. *Chin. J. Cancer Res.* **26**, 493–500 (2014).
72. Sokoll, L. J. *et al.* A multicenter evaluation of the PCA3 molecular urine test: pre-analytical effects, analytical performance, and diagnostic accuracy. *Clin. Chim. Acta.* **389**, 1–6 (2008).
73. Hessels, D. *et al.* Detection of TMPRSS2-ERG fusion transcripts and prostate cancer antigen 3 in urinary sediments may improve diagnosis of prostate cancer. *Clin. Cancer Res.* **13**, 5103–8 (2007).
74. van Gils, M. P. M. Q. *et al.* The time-resolved fluorescence-based PCA3 test on urinary sediments after digital rectal examination; a Dutch multicenter validation of the diagnostic performance. *Clin. Cancer Res.* **13**, 939–43 (2007).
75. Schmidt, U. *et al.* Quantitative Multi-Gene Expression Profiling of Primary Prostate Cancer. **1534**, (2006).
76. Gittelman, M. C. *et al.* PCA3 molecular urine test as a predictor of



- repeat prostate biopsy outcome in men with previous negative biopsies: a prospective multicenter clinical study. *J. Urol.* **190**, 64–69 (2013).
77. Manuscript, A., Variants, A. R. S. & Tissues, M. P. Expression in Normal and Malignant Prostate Tissues. **77**, 1–12 (2012).
  78. Sambasivarao, S. V. NIH Public Access. **18**, 1199–1216 (2013).
  79. Ozgur, T., Atik, E., Hakverdi, S. & Yaldiz, M. The expressions of {AMACR} and {iNOS} in prostate adenocarcinomas. *Pakistan J. Med. Sci.* **29**, 610–613 (2013).
  80. Lloyd, M. D., Darley, D. J., Wierzbicki, A. S. & Threadgill, M. D. Alpha-methylacyl-CoA racemase--an 'obscure' metabolic enzyme takes centre stage. *FEBS J.* **275**, 1089–102 (2008).
  81. Rubin, M. a *et al.* Decreased alpha-methylacyl CoA racemase expression in localized prostate cancer is associated with an increased rate of biochemical recurrence and cancer-specific death. *Cancer Epidemiol. Biomarkers Prev.* **14**, 1424–1432 (2005).
  82. Ware, K. E., Garcia-Blanco, M. a., Armstrong, A. J. & Dehm, S. M. Biologic and clinical significance of androgen receptor variants in castration resistant prostate cancer. *Endocr. Relat. Cancer* **21**, 87–103 (2014).
  83. Geng, C. *et al.* Prostate cancer-associated mutations in speckle-type POZ protein (SPOP) regulate steroid receptor coactivator 3 protein turnover. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 6997–7002 (2013).
  84. An, J., Wang, C., Deng, Y., Yu, L. & Huang, H. Destruction of full-length androgen receptor by wild-type SPOP, but not prostate-cancer-associated mutants. *Cell Rep.* **6**, 657–669 (2014).
  85. Schoenborn, J. R., Nelson, P. & Fang, M. Genomic profiling defines subtypes of prostate cancer with the potential for therapeutic stratification. *Clin. Cancer Res.* **19**, 4058–66 (2013).
  86. Barbieri, C. E. *et al.* Exome sequencing identifies recurrent SPOP, FOXA1 and MED12 mutations in prostate cancer. *Nat. Genet.* **44**, 685–9 (2012).
  87. Barbieri, C. E. *et al.* Exome sequencing identifies recurrent SPOP, FOXA1 and MED12 mutations in prostate cancer. *Nat. Genet.* **44**, 685–9 (2012).
  88. Grasso, C. S. *et al.* The mutational landscape of lethal castration-resistant prostate cancer. *Nature* **487**, 239–43 (2012).
  89. Grasso, C. S. *et al.* The mutational landscape of lethal castration-resistant prostate cancer. *Nature* **487**, 239–43 (2012).
  90. Hossain, D. & Bostwick, D. G. Significance of the TMPRSS2:ERG gene fusion in prostate cancer. *BJU Int.* **111**, 834–5 (2013).
  91. Clark, J. P. & Cooper, C. S. ETS gene fusions in prostate cancer. *Nat. Rev. Urol.* **6**, 429–39 (2009).
  92. Clark, J. P. & Cooper, C. S. ETS gene fusions in prostate cancer. *Nat. Rev. Urol.* **6**, 429–39 (2009).
  93. Li, L. *et al.* Targeting Poly(ADP-Ribose) Polymerase and the c-Myb-Regulated DNA Damage Response Pathway in Castration-Resistant Prostate Cancer. *Sci. Signal.* **7**, ra47 (2014).
  94. Bisen, A. & Claxton, D. F. Tyrosine kinase targeted treatment of

## 9: APPENDICES

- chronic myelogenous leukemia and other myeloproliferative neoplasms. *Adv. Exp. Med. Biol.* **779**, 179–196 (2013).
95. Tomlins, S. a. *et al.* Role of the TMPRSS2-ERG Gene Fusion in Prostate Cancer. *Neoplasia* **10**, 177-IN9 (2008).
  96. Tomlins, S. a. *et al.* Role of the TMPRSS2-ERG Gene Fusion in Prostate Cancer. *Neoplasia* **10**, 177-IN9 (2008).
  97. Shao, L. *et al.* Highly specific targeting of the TMPRSS2/ERG fusion gene using liposomal nanovectors. *Clin. Cancer Res.* **18**, 6648–57 (2012).
  98. Goh, C. L. *et al.* Genetic variants associated with predisposition to prostate cancer and potential clinical implications. *J. Intern. Med.* **271**, 353–65 (2012).
  99. Goh, C. L. *et al.* Genetic variants associated with predisposition to prostate cancer and potential clinical implications. *J. Intern. Med.* **271**, 353–65 (2012).
  100. Levy-Lahad, E. & Friedman, E. Cancer risks among BRCA1 and BRCA2 mutation carriers. *Br. J. Cancer* **96**, 11–15 (2007).
  101. Drake, Richard, Vogl, Wayne & Mitchell, A. *Gray's Anatomy for Students.* (Elsevier Inc., 2015).
  102. Boutros, P. C. *et al.* Spatial genomic heterogeneity within localized , multifocal prostate cancer. (2015). doi:10.1038/ng.3315
  103. Bátiz, L. F. *et al.* Exosomes as Novel Regulators of Adult Neurogenic Niches . *Frontiers in Cellular Neuroscience* **9**, 501 (2016).
  104. Nilsson, J. *et al.* Prostate cancer-derived urine exosomes : a novel approach to biomarkers for prostate cancer. **100**, 1603–1607 (2009).
  105. Lin, J. *et al.* Exosomes : Novel Biomarkers for Clinical Diagnosis. *Sci. World J.* **2015**, (2015).
  106. Hanahan, D. & Weinberg, R. A. Hallmarks of cancer: The next generation. *Cell* **144**, 646–674 (2011).
  107. Azmi, A. S., Bao, B. & Sarkar, F. H. Exosomes in cancer development, metastasis, and drug resistance: A comprehensive review. *Cancer Metastasis Rev.* **32**, 623–642 (2013).
  108. Hoshino, A. *et al.* Tumour exosome integrins determine organotropic metastasis. *Nature* **527**, 329–335 (2015).
  109. Costa-silva, B. *et al.* Pancreatic cancer exosomes initiate pre-metastatic niche formation in the liver. **17**, (2015).
  110. Skog, J. *et al.* Glioblastoma microvesicles transport RNA and protein that promote tumor growth and provide diagnostic biomarkers. *Nat. Cell Biol.* **10**, 1470–1476 (2008).
  111. Nilsson, J. *et al.* Prostate cancer-derived urine exosomes : a novel approach to biomarkers for prostate cancer. **100**, 1603–1607 (2009).
  112. Dijkstra, S. *et al.* Prostate cancer biomarker profiles in urinary sediments and exosomes. *J. Urol.* **191**, 1132–1138 (2014).
  113. Hu, L. *et al.* Fluorescence in situ hybridization (FISH): an increasingly demanded tool for biomarker research and personalized medicine. *Biomark. Res.* **2**, 3 (2014).
  114. Geiss, G. K. *et al.* Direct multiplexed measurement of gene expression with color-coded probe pairs. *Nat. Biotechnol.* **26**, 317–25 (2008).
  115. Geiss, G. K. *et al.* Direct multiplexed measurement of gene expression

- with color-coded probe pairs. *Nat. Biotechnol.* **26**, 317–25 (2008).
116. NanoString Technologies®. nCounter® Gene Expression CodeSets.
  117. Sanger, F., Nicklen, S. & Coulson, a R. DNA sequencing with chain-terminating inhibitors. *Proc. Natl. Acad. Sci. U. S. A.* **74**, 5463–7 (1977).
  118. National Human Genome Research Insitute (NIH). The Cost of Sequencing a Human Genome. (2016). Available at: <https://www.genome.gov/sequencingcosts/>. (Accessed: 1st January 2016)
  119. National Human Genome Research Insitute (NIH). The Cost of Sequencing a Human Genome. (2016).
  120. Heather, J. M. & Chain, B. The sequence of sequencers: The history of sequencing DNA. *Genomics* **107**, 1–8 (2016).
  121. Heather, J. M. & Chain, B. The sequence of sequencers: The history of sequencing DNA. *Genomics* **107**, 1–8 (2016).
  122. Holt, R. A. & Jones, S. J. M. The new paradigm of flow cell sequencing. 839–846 (2008). doi:10.1101/gr.073262.107.cell
  123. Fullwood, M. J., Wei, C., Liu, E. T. & Ruan, Y. Next-generation DNA sequencing of paired-end tags ( PET ) for transcriptome and genome analyses. 521–532 (2009). doi:10.1101/gr.074906.107.Freely
  124. Metzker, M. L. Sequencing technologies - the next generation. *Nat. Rev. Genet.* **11**, 31–46 (2010).
  125. Voelkerding, K. V, Dames, S. a & Durtschi, J. D. Next-generation sequencing: from basic research to diagnostics. *Clin. Chem.* **55**, 641–58 (2009).
  126. Voelkerding, K. V, Dames, S. a & Durtschi, J. D. Next-generation sequencing: from basic research to diagnostics. *Clin. Chem.* **55**, 641–58 (2009).
  127. Lee, H. *et al.* Third-generation sequencing and the future of genomics. *bioRxiv* 048603 (2016). doi:10.1101/048603
  128. Lee, H. *et al.* Third-generation sequencing and the future of genomics. *bioRxiv* 048603 (2016). doi:10.1101/048603
  129. Rabbani, B., Tekin, M. & Mahdih, N. The promise of whole-exome sequencing in medical genetics. **59**, 5–15 (2013).
  130. Warr, A. *et al.* Exome Sequencing: Current and Future Perspectives. *G3&#58; Genes/Genomes/Genetics* **5**, 1543–1550 (2015).
  131. Warr, A. *et al.* Exome Sequencing: Current and Future Perspectives. *G3&#58; Genes/Genomes/Genetics* **5**, 1543–1550 (2015).
  132. Wilhelm, B. T. & Landry, J. R. RNA-Seq-quantitative measurement of expression through massively parallel RNA-sequencing. *Methods* **48**, 249–257 (2009).
  133. Zhao, S., Bittner, A., Ngo, K. & Liu, X. Comparison of RNA-Seq and Microarray in Transcriptome Profiling of Activated T Cells. **9**, (2014).
  134. Ozsolak, F. & Milos, P. M. RNA sequencing: advances, challenges and opportunities. *Nat. Rev. Genet.* **12**, 87–98 (2011).
  135. Liu, S. & Trapnell, C. Single-cell transcriptome sequencing : recent advances and remaining challenges [ version 1 ; referees : 2 approved ] Referee Status : **5**, (2016).
  136. Liu, E. T., Pott, S. & Huss, M. Q & A : ChIP-seq technologies and the study of gene regulation. 4–9 (2010).

## 9: APPENDICES

137. Mamanova, L. *et al.* Target-enrichment strategies for next- generation sequencing. **7**, (2010).
138. Clark-langone, K. M. *et al.* assay. **18**, 1–18 (2007).
139. Knuutila, S. Biomarker analysis in human neoplasias : superior next-generation sequencing on frozen bone marrow cells and on formalin-fixed , paraffin- embedded tumor tissues. *BMC Proc.* **7**, K18 (2013).
140. Saini, S. PSA and beyond: alternative prostate cancer biomarkers. *Cell. Oncol. (Dordr).* **39**, 97–106 (2016).
141. Murphy, D. Gene expression studies using microarrays: principles, problems, and prospects. *Adv. Physiol. Educ.* **26**, 256–270 (2002).
142. Vermeeren, V. & Michiels, L. *Evolution Towards the Implementation of Point-Of-Care Biosensors.* (INTECH Open Access Publisher, 2011).
143. Kretschmer, A. & Tilki, D. Biomarkers in prostate cancer – Current clinical utility and future perspectives. *Crit. Rev. Oncol. Hematol.* **120**, 180–193 (2017).
144. Dalela, D. *et al.* Genomic Classifier Augments the Role of Pathological Features in Identifying Optimal Candidates for Adjuvant Radiation Therapy in Patients With Prostate Cancer: Development and Internal Validation of a Multivariable Prognostic Model. *J. Clin. Oncol.* **35**, 1982–1990 (2017).
145. Burgess, A., Shah, K., Hough, O. & Hynynen, K. HHS Public Access. **15**, 477–491 (2016).
146. Klein, E. A. *et al.* Decipher Genomic Classifier Measured on Prostate Biopsy Predicts Metastasis Risk. *Urology* **90**, 148–152 (2016).
147. Graves, P. R. & Haystead, T. A. J. Molecular Biologist ' s Guide to Proteomics. **66**, 39–63 (2002).
148. Semmes, O. J. Defining the Role of Mass Spectrometry in Cancer Diagnostics Defining the Role of Mass Spectrometry in Cancer Diagnostics. 1555–1557 (2004).
149. Kakimoto, Y., Tsuruyama, T., Yamamoto, T., Furuta, M. & Kotani, H. Novel In Situ Pretreatment Method for Significantly Enhancing the Signal In MALDI-TOF MS of Formalin- Fixed Paraffin-Embedded Tissue Sections. **7**, 1–7 (2012).
150. Comparison, P. Validation Processes of Protein Biomarkers in Serum — A Cross Platform Comparison. 12710–12728 (2012). doi:10.3390/s120912710
151. Manuscript, A. targeted mass spectrometry. **10**, 28–34 (2014).
152. Bayani, J. & Squire, J. a. Application and interpretation of FISH in biomarker studies. *Cancer Lett.* **249**, 97–109 (2007).
153. Bayani, J. & Squire, J. a. Application and interpretation of FISH in biomarker studies. *Cancer Lett.* **249**, 97–109 (2007).
154. Hu, L. *et al.* Fluorescence in situ hybridization (FISH): an increasingly demanded tool for biomarker research and personalized medicine. *Biomark. Res.* **2**, 3 (2014).
155. Dunstan, R. W., Wharton, K. a, Quigley, C. & Lowe, A. The use of immunohistochemistry for biomarker assessment--can it compete with other technologies? *Toxicol. Pathol.* **39**, 988–1002 (2011).
156. Dunstan, R. W., Wharton, K. a, Quigley, C. & Lowe, A. The use of immunohistochemistry for biomarker assessment--can it compete

- with other technologies? *Toxicol. Pathol.* **39**, 988–1002 (2011).
157. Ferrier, C. M. *et al.* Comparison of immunohistochemistry with immunoassay (ELISA) for the detection of components of the plasminogen activation system in human tumour tissue. *Br. J. Cancer* **79**, 1534–1541 (1999).
  158. Sant, K. E., Nahar, M. S. & Dolinoy, D. C. DNA methylation screening and analysis. *Methods Mol. Biol.* **889**, 385–406 (2012).
  159. Shen, L. & Waterland, R. A. Methods of DNA methylation analysis. *Curr. Opin. Clin. Nutr. Metab. Care* **10**, 576–581 (2007).
  160. Jin, B., Li, Y. & Robertson, K. D. DNA Methylation : Superior or Subordinate in the Epigenetic Hierarchy ? 607–617 (2011). doi:10.1177/1947601910393957
  161. Hoque, M. O. & Surgery, N. HHS Public Access. **9**, 243–257 (2015).
  162. Kurdyukov, S. & Bullock, M. DNA Methylation Analysis: Choosing the Right Method. *Biology (Basel)*. **5**, 3 (2016).
  163. Habeeb, N. M. A. W. *et al.* - ntegrated analysis of epigenomic and genomic changes by DNA methylation dependent mechanisms provides potential novel biomarkers for prostate cancer. **5**,
  164. Yang, M. & Park, J. Y. DNA methylation in promoter region as biomarkers in prostate cancer. *Methods Mol. Biol.* **863**, 67–109 (2012).
  165. Murphy, T. M., Perry, A. S. & Lawler, M. The emergence of DNA methylation as a key modulator of aberrant cell death in prostate cancer. (2008). doi:10.1677/ERC-07-0208
  166. Goering, W., Kloth, M. & Schulz, W. A. DNA methylation changes in prostate cancer. *Methods Mol. Biol.* **863**, 47–66 (2012).
  167. Kurdyukov, S. & Bullock, M. DNA Methylation Analysis: Choosing the Right Method. *Biology (Basel)*. **5**, 3 (2016).
  168. Hastie, T., Tibshirani, R. & Friedman, J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. (Springer US, 2011).
  169. Hahne, F., Huber, W., Gentleman, R. & Falcon, S. *Bioconductor Case Studies*. (2008).
  170. Knuutila, S. Biomarker analysis in human neoplasias: superior next-generation sequencing on frozen bone marrow cells and on formalin-fixed, paraffin-embedded tumor tissues. *BMC Proc.* **7**, K18–K18 (2013).
  171. Protocol, Q. Ovation Pico WTA System ®. *Ovation 4–7* (2010).
  172. Waggott, D. *et al.* NanoStringNorm: an extensible R package for the pre-processing of NanoString mRNA and miRNA data. *Bioinformatics* **28**, 1546–8 (2012).
  173. Nickles, A. D., Sandmann, T., Ziman, R. & Bourgon, R. Package ‘ NanoStringQCPro ’. (2018).
  174. Waggott, D. *et al.* NanoStringNorm: an extensible R package for the pre-processing of NanoString mRNA and miRNA data. *Bioinformatics* **28**, 1546–8 (2012).
  175. Jung, S.-H. & Sohn, I. Statistical Issues in the Design and Analysis of nCounter Projects. *Cancer Inform.* **13**, 35–43 (2014).
  176. Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data using empirical Bayes methods.

- Biostatistics* **8**, 118–127 (2007).
177. Hart, A. Mann-Whitney test is not just a test of medians: differences in spread can be important. *Bmj* **323**, 391–393 (2001).
  178. Mukaka, M. M. Statistics corner: A guide to appropriate use of correlation coefficient in medical research. *Malawi Med. J.* **24**, 69–71 (2012).
  179. Howell, D. C. Chi-Square Test - Analysis of Contingency Tables. *Test* 1–4 (2000). doi:10.1007/978-3-642-04898-2\_174
  180. Bland, J. M. & Altman, D. G. The logrank test. *BMJ* **328**, 1073 (2004).
  181. Royston, P. Remark AS R94: A Remark on Algorithm AS 181: The W-test for Normality. *J. R. Stat. Soc. Ser. C (Applied Stat.)* **44**, 547–551 (1995).
  182. Yokoyama, M., Nishi, Y., Yoshii, J., Okubo, K. & Matsubara, K. Identification and cloning of neuroblastoma-specific and nerve tissue-specific genes through compiled expression profiles. *DNA Res.* **3**, 311–320 (1996).
  183. Carey, M. V. Package 'ROC'. (2018).
  184. Bendix, A. & Carstensen, M. B. Package 'Epi'. (2018).
  185. Jeong, D. H., Ziemkiewicz, C., Ribarsky, W. & Chang, R. Understanding Principal Component Analysis Using a Visual Analytics Tool. *Proc. UKC 2009, Math. Fundam. Appl. 2009* 1–10 (2009). doi:10.1.1.157.1469
  186. Legendre, P. & Legendre, L. Numerical Ecology, Volume 24. (*Developments Environ. Model.* **24**, 870 (1988).
  187. Suzuki, R. & Shimodaira, H. Pvcust: an R package for assessing the uncertainty in hierarchical clustering. *Bioinformatics* **22**, 1540–1542 (2006).
  188. Hartigan, J. A. & Wong, M. A. Algorithm AS 136: A K-Means Clustering Algorithm. *J. R. Stat. Soc. Ser. C (Applied Stat.)* **28**, 100–108 (1979).
  189. Singh, K., Malik, D. & Sharma, N. Evolving limitations in K-means algorithm in data mining and their removal. *IJCEM Int. J. Comput. Eng. Manag. ISSN* **12**, 2230–7893 (2011).
  190. Singh, K., Malik, D. & Sharma, N. Evolving limitations in K-means algorithm in data mining and their removal. *IJCEM Int. J. Comput. Eng. Manag. ISSN* **12**, 2230–7893 (2011).
  191. Charrad, M., Ghazzali, N., Boiteau, V. & Niknafs, A. **NbClust** : An R Package for Determining the Relevant Number of Clusters in a Data Set. *J. Stat. Softw.* **61**, (2014).
  192. Rogers, S., Girolami, M., Campbell, C. & Breitling, R. The latent process decomposition of cDNA microarray data sets. *IEEE/ACM Trans. Comput. Biol. Bioinforma.* **2**, 143–156 (2005).
  193. Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear Models via Coordinate Descent. *J. Stat. Softw.* **33**, 1–22 (2010).
  194. Tibshirani, R. Regression Selection and Shrinkage via the Lasso. *Journal of the Royal Statistical Society B* **58**, 267–288 (1994).
  195. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).
  196. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).
  197. Zhang, Z. Variable selection with stepwise and best subset

- approaches. *Ann. Transl. Med.* **4**, 136 (2016).
198. Zhang, Z. Variable selection with stepwise and best subset approaches. *Ann. Transl. Med.* **4**, 136 (2016).
  199. Huang, D. W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat. Protoc.* **4**, 44–57 (2009).
  200. Crawley, M. J. *The R book*. (Second edition. Chichester, West Sussex, United Kingdom : Wiley, 2013., 2013).
  201. Lumley, T. & S-, R. Package ‘ survival ’. (2018).
  202. Schloerke, B. Package ‘ GGally ’ R topics documented : (2018).
  203. Ahmed, H. U. *et al.* Diagnostic accuracy of multi-parametric MRI and TRUS biopsy in prostate cancer (PROMIS): a paired validating confirmatory study. *Lancet* **389**, 815–822 (2018).
  204. Malkov, V. *a et al.* Multiplexed measurements of gene signatures in different analytes using the Nanostring nCounter Assay System. *BMC Res. Notes* **2**, 80 (2009).
  205. David, A. *et al.* Unusual alternative splicing within the human kallikrein genes KLK2 and KLK3 gives rise to novel prostate-specific proteins. (2002).
  206. Gleason, D. F. & Mellinger, G. T. Prediction of prognosis for prostatic adenocarcinoma by combined histological grading and clinical staging. *J. Urol.* **111**, 58–64 (1974).
  207. Gordetsky, J. & Epstein, J. Grading of prostatic adenocarcinoma: current state and prognostic implications. *Diagn. Pathol.* **11**, 25 (2016).
  208. Sorlie, T. *et al.* Repeated observation of breast tumor subtypes in independent gene expression data sets. *Proc. Natl. Acad. Sci. U. S. A.* **100**, 8418–23 (2003).
  209. Rosell, R. *et al.* Erlotinib versus standard chemotherapy as first-line treatment for European patients with advanced EGFR mutation-positive non-small-cell lung cancer (EURTAC): a multicentre, open-label, randomised phase 3 trial. *Lancet Oncol.* **13**, 239–246 (2012).
  210. Bethune, G., Bethune, D., Ridgway, N. & Xu, Z. Epidermal growth factor receptor (EGFR) in lung cancer: an overview and update. *J. Thorac. Dis.* **2**, 48–51 (2010).
  211. Ein-Dor, L., Kela, I., Getz, G., Givol, D. & Domany, E. Outcome signature genes in breast cancer: is there a unique set? *Bioinformatics* **21**, 171–178 (2005).
  212. Thompson, I. M. *et al.* Assessing prostate cancer risk: results from the Prostate Cancer Prevention Trial. *J. Natl. Cancer Inst.* **98**, 529–534 (2006).
  213. Tomlins, S. A. *et al.* Urine TMPRSS2:ERG Plus PCA3 for Individualized Prostate Cancer Risk Assessment. *Eur. Urol.* **70**, 45–53 (2016).
  214. Hessels, D. *et al.* DD3(PCA3)-based molecular urine analysis for the diagnosis of prostate cancer. *Eur. Urol.* **44**, 6–8 (2003).
  215. Van Neste, L. *et al.* Detection of High-grade Prostate Cancer Using a Urinary Molecular Biomarker-Based Risk Score. *Eur. Urol.* **70**, 740–748 (2016).
  216. McKiernan, J. *et al.* A Novel Urine Exosome Gene Expression Assay to

- Predict High-grade Prostate Cancer at Initial Biopsy. *JAMA Oncol.* **2**, 882–889 (2016).
217. Miranda, K. C. *et al.* Nucleic acids within urinary exosomes/microvesicles are potential biomarkers for renal disease. *Kidney Int.* **78**, 191–199 (2010).
  218. Dijkstra, S. *et al.* Prostate cancer biomarker profiles in urinary sediments and exosomes. *J. Urol.* **191**, 1132–1138 (2014).
  219. Mootha, V. K. *et al.* PGC-1[alpha]-responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes. *Nat Genet* **34**, 267–273 (2003).
  220. Dhaese, S. *et al.* Functional and profiling studies prove that prostate cancer upregulated neuroblastoma thymosin beta is the true human homologue of rat thymosin beta15. *FEBS Lett.* **581**, 4809–4815 (2007).
  221. Yokoyama, M., Nishi, Y., Yoshii, J., Okubo, K. & Matsubara, K. Identification and cloning of neuroblastoma-specific and nerve tissue-specific genes through compiled expression profiles. *DNA Res.* **3**, 311–320 (1996).
  222. Theunissen, W. *et al.* Thymosin Beta 4 and Thymosin Beta 10 Expression in Hepatocellular Carcinoma. *Eur. J. Histochem.* **58**, 2242 (2014).
  223. Darb-Esfahani, S. *et al.* Thymosin beta 15A (TMSB15A) is a predictor of chemotherapy response in triple-negative breast cancer. *Br. J. Cancer* **107**, 1892–1900 (2012).
  224. Lenka, G., Weng, W.-H., Chuang, C.-K., Ng, K.-F. & Pang, S.-T. Aberrant expression of the PRAC gene in prostate cancer. *Int. J. Oncol.* **43**, 1960–1966 (2013).
  225. Koestler, D. C. *et al.* Distinct patterns of DNA methylation in conventional adenomas involving the right and left colon. *Mod. Pathol. an Off. J. United States Can. Acad. Pathol. Inc* **27**, 145–155 (2014).
  226. Lose, F. *et al.* Genetic association of the KLK4 locus with risk of prostate cancer. *PLoS One* **7**, e44520 (2012).
  227. Vaananen, R.-M. *et al.* Association of transcript levels of 10 established or candidate-biomarker gene targets with cancerous versus non-cancerous prostate tissue from radical prostatectomy specimens. *Clin. Biochem.* **46**, 670–674 (2013).
  228. Bostwick, D. G., Pacelli, A., Blute, M., Roche, P. & Murphy, G. P. Prostate specific membrane antigen expression in prostatic intraepithelial neoplasia and adenocarcinoma: a study of 184 cases. *Cancer* **82**, 2256–2261 (1998).
  229. Vainio, P. *et al.* Arachidonic acid pathway members PLA2G7, HPGD, EPHX2, and CYP4F8 identified as putative novel therapeutic targets in prostate cancer. *Am. J. Pathol.* **178**, 525–536 (2011).
  230. Lorenzo, Y. *et al.* Differential genetic and functional markers of second neoplasias in Hodgkin's disease patients. *Clin. Cancer Res.* **15**, 4823–4828 (2009).
  231. Figueroa, J. D. *et al.* Bladder cancer risk and genetic variation in AKR1C3 and other metabolizing genes. *Carcinogenesis* **29**, 1955–1962 (2008).



## 9: APPENDICES

232. Brenner, D. R. *et al.* Hierarchical modeling identifies novel lung cancer susceptibility variants in inflammation pathways among 10,140 cases and 11,012 controls. *Hum. Genet.* **132**, 579–589 (2013).
233. Tomatis, S. *et al.* Late rectal bleeding after 3D-CRT for prostate cancer: development of a neural-network-based predictive model. *Phys. Med. Biol.* **57**, 1399–1412 (2012).
234. Hendriksen, J. *et al.* RanBP3 enhances nuclear export of active (beta)-catenin independently of CRM1. *J. Cell Biol.* **171**, 785–797 (2005).
235. Polakis, P. Wnt signaling and cancer. *Genes Dev.* **14**, 1837–1851 (2000).
236. Liguori, L. *et al.* The metallophosphodiesterase Mpped2 impairs tumorigenesis in neuroblastoma. *Cell Cycle* **11**, 569–581 (2012).
237. Chen, W.-Z., Pang, B., Yang, B., Zhou, J.-G. & Sun, Y.-H. Differential proteome analysis of conditioned medium of BPH-1 and LNCaP cells. *Chin. Med. J. (Engl)*. **124**, 3806–3809 (2011).
238. Loo, J. M. *et al.* Extracellular metabolic energetics can promote cancer progression. *Cell* **160**, 393–406 (2015).
239. Schumacher, F. R. *et al.* Genome-wide association study identifies new prostate cancer susceptibility loci. *Hum. Mol. Genet.* **20**, 3867–3875 (2011).
240. Grabowska, M. M. *et al.* NFI Transcription Factors Interact with FOXA1 to Regulate Prostate-Specific Gene Expression. *Mol. Endocrinol.* **28**, 949–964 (2014).
241. Lee, J. S. *et al.* A novel tumor-promoting role for nuclear factor IA in glioblastomas is mediated through negative regulation of p53, p21, and PAI1. *Neuro. Oncol.* **16**, 191–203 (2014).
242. Matsuo, M. *et al.* Designation of enzyme activity of glycine-N-acyltransferase family genes and depression of glycine-N-acyltransferase in human hepatocellular carcinoma. *Biochem. Biophys. Res. Commun.* **420**, 901–906 (2012).
243. Grabowska, M. M. *et al.* NFI Transcription Factors Interact with FOXA1 to Regulate Prostate-Specific Gene Expression. *Mol. Endocrinol.* **28**, 949–964 (2014).
244. Mirabello, L. *et al.* A Genome-Wide Scan Identifies Variants in NFIB Associated with Metastasis in Patients with Osteosarcoma. *Cancer Discov.* **5**, 920–931 (2015).
245. Gosenca, D. *et al.* Identification and functional characterization of imatinib-sensitive DTD1-PDGFRB and CCDC88C-PDGFRB fusion genes in eosinophilia-associated myeloid/lymphoid neoplasms. *Genes. Chromosomes Cancer* **53**, 411–421 (2014).
246. Bhatlekar, S., Fields, J. Z. & Boman, B. M. HOX genes and their role in the development of human cancers. *J. Mol. Med. (Berl)*. **92**, 811–823 (2014).
247. Henderson, D. J. P. *et al.* The cAMP phosphodiesterase-4D7 (PDE4D7) is downregulated in androgen-independent prostate cancer cells and mediates proliferation by compartmentalising cAMP at the plasma membrane of VCaP prostate cancer cells. *Br. J. Cancer* **110**, 1278–1287 (2014).
248. Severi, G. *et al.* A three-protein biomarker panel assessed in

- diagnostic tissue predicts death from prostate cancer for men with localized disease. *Cancer Med.* **3**, 1266–1274 (2014).
249. Ji, D. *et al.* Prognostic role of serum AZGP1, PEDF and PRDX2 in colorectal cancer patients. *Carcinogenesis* **34**, 1265–1272 (2013).
  250. Huang, C. *et al.* Decreased expression of AZGP1 is associated with poor prognosis in primary gastric cancer. *PLoS One* **8**, e69155 (2013).
  251. Kim, Y.-S., Do Hwan, J., Bae, S., Bae, D.-H. & Ahn Shick, W. Identification of differentially expressed genes using an annealing control primer system in stage III serous ovarian carcinoma. *BMC Cancer* **10**, 576 (2010).
  252. Protein Atlas. No Title. Available at: <https://www.proteinatlas.org/>.
  253. GeneCards.
  254. Amirian, E. S., Ittmann, M. M. & Scheurer, M. E. Associations between arachidonic acid metabolism gene polymorphisms and prostate cancer risk. *Prostate* **71**, 1382–1389 (2011).
  255. White, K. L. *et al.* Ovarian cancer risk associated with inherited inflammation-related variants. *Cancer Res.* **72**, 1064–1069 (2012).
  256. Sanchez-Espiridion, B. *et al.* Immunohistochemical markers for tumor associated macrophages and survival in advanced classical Hodgkin's lymphoma. *Haematologica* **97**, 1080–1084 (2012).
  257. Huang, H., Hara, A., Homma, T., Yonekawa, Y. & Ohgaki, H. Altered expression of immune defense genes in pilocytic astrocytomas. *J. Neuropathol. Exp. Neurol.* **64**, 891–901 (2005).
  258. Lu, J. P. *et al.* Androgens induce oxidative stress and radiation resistance in prostate cancer cells through NADPH oxidase. *Prostate Cancer Prostatic Dis.* **13**, 39–46 (2010).
  259. Kikuchi, H., Hikage, M., Miyashita, H. & Fukumoto, M. NADPH oxidase subunit, gp91(phox) homologue, preferentially expressed in human colon epithelial cells. *Gene* **254**, 237–243 (2000).
  260. Speeckaert, M. M., Speeckaert, R., Laute, M., Vanholder, R. & Delanghe, J. R. Tumor necrosis factor receptors: biology and therapeutic potential in kidney diseases. *Am. J. Nephrol.* **36**, 261–270 (2012).
  261. He, L. *et al.* Serglycin (SRGN) overexpression predicts poor prognosis in hepatocellular carcinoma patients. *Med. Oncol.* **30**, 707 (2013).
  262. Baron, V. T., Pio, R., Jia, Z. & Mercola, D. Early Growth Response 3 regulates genes of inflammation and directly activates IL6 and IL8 expression in prostate cancer. *Br. J. Cancer* **112**, 755–764 (2015).
  263. Huang, M.-Y. *et al.* EVI2B, ATP2A2, S100B, TM4SF3, and OLFM4 as potential prognostic markers for postoperative Taiwanese colorectal cancer patients. *DNA Cell Biol.* **31**, 625–635 (2012).
  264. Kim, S.-K. *et al.* A nineteen gene-based risk score classifier predicts prognosis of colorectal cancer patients. *Mol. Oncol.* **8**, 1653–1666 (2014).
  265. Leite, K. R. M. *et al.* Controlling RECK miR21 Promotes Tumor Cell Invasion and Is Related to Biochemical Recurrence in Prostate Cancer. *J. Cancer* **6**, 292–301 (2015).
  266. Mao, B., Xiao, H., Zhang, Z., Wang, D. & Wang, G. MicroRNA21 regulates the expression of BTG2 in HepG2 liver cancer cells. *Mol. Med. Rep.* **12**, 4917–4924 (2015).

## 9: APPENDICES

267. Maachani, U. B., Tandle, A., Shankavaram, U., Kramp, T. & Camphausen, K. Modulation of miR-21 signaling by MPS1 in human glioblastoma. *Oncotarget* (2015).
268. Ciarlo, M. *et al.* Regulation of neuroendocrine differentiation by AKT/hnRNP/AR/beta-catenin signaling in prostate cancer cells. *Int. J. Cancer* **131**, 582–590 (2012).
269. Chauhan, S. S. *et al.* Prediction of recurrence-free survival using a protein expression-based risk classifier for head and neck cancer. *Oncogenesis* **4**, e147 (2015).
270. Han, S.-S. *et al.* RNA sequencing identifies novel markers of non-small cell lung cancer. *Lung Cancer* **84**, 229–235 (2014).
271. Saini, V. *et al.* Identification of CBX3 and ABCA5 as putative biomarkers for tumor stem cells in osteosarcoma. *PLoS One* **7**, e41401 (2012).
272. Fernandez, P. L. *et al.* Expression of cathepsins B and S in the progression of prostate carcinoma. *Int. J. Cancer* **95**, 51–55 (2001).
273. Liu, P., Jiang, W., Ren, H., Zhang, H. & Hao, J. Exploring the Molecular Mechanism and Biomarkers of Liver Cancer Based on Gene Expression Microarray. *Pathol. Oncol. Res.* **21**, 1077–1083 (2015).
274. Tsai, J.-Y. *et al.* Effects of novel human cathepsin S inhibitors on cell migration in human cancer cells. *J. Enzyme Inhib. Med. Chem.* **29**, 538–546 (2014).
275. Kutomi, G. *et al.* Human endoplasmic reticulum oxidoreductin 1-alpha is a novel predictor for poor prognosis of breast cancer. *Cancer Sci.* **104**, 1091–1096 (2013).
276. Fujita, K. & Nonomura, N. Urinary biomarkers of prostate cancer. *Int. J. Urol.* **25**, 770–779 (2018).
277. NCBI: Gene. (2018). Available at: <https://www.ncbi.nlm.nih.gov/gene/5788>.
278. Petrie, M. *et al.* The Vesicle Priming Factor CAPS Functions as a Homodimer via C2 Domain Interactions to Promote Regulated Vesicle Exocytosis. *J. Biol. Chem.* **291**, 21257–21270 (2016).
279. Haider, S. *et al.* A multi-gene signature predicts outcome in patients with pancreatic ductal adenocarcinoma. *Genome Med.* **6**, 105 (2014).
280. Miller, S. *et al.* Genome-wide molecular characterization of central nervous system primitive neuroectodermal tumor and pineoblastoma. *Neuro. Oncol.* **13**, 866–879 (2011).
281. Pal, R. P. *et al.* Immunocytochemical detection of ERG expression in exfoliated urinary cells identifies with high specificity patients with prostate cancer. *BJU Int.* **117**, 686–696 (2016).
282. Whitaker, H. C. *et al.* N-acetyl-L-aspartyl-L-glutamate peptidase-like 2 is overexpressed in cancer and promotes a pro-migratory and pro-metastatic phenotype. *Oncogene* **33**, 5274–5287 (2014).
283. Yoon, H. *et al.* Tudor domain-containing protein 4 as a potential cancer/testis antigen in liver cancer. *Tohoku J. Exp. Med.* **224**, 41–46 (2011).
284. Jiang, Y., Liu, L., Shan, W. & Yang, Z.-Q. An integrated genomic analysis of Tudor domain-containing proteins identifies PHD finger protein 20-like 1 (PHF20L1) as a candidate oncogene in breast cancer. *Mol.*

## 9: APPENDICES

- Oncol.* **10**, 292–302 (2016).
285. Halvorsen, O. J. *et al.* Increased expression of SIM2-s protein is a novel marker of aggressive prostate cancer. *Clin. Cancer Res.* **13**, 892–897 (2007).
  286. Long, Q. *et al.* Protein-coding and microRNA biomarkers of recurrence of prostate cancer following radical prostatectomy. *Am. J. Pathol.* **179**, 46–54 (2011).
  287. Cooperberg, M. R. *et al.* Multiinstitutional validation of the UCSF cancer of the prostate risk assessment for prediction of recurrence after radical prostatectomy. *Cancer* **107**, 2384–2391 (2006).
  288. Roobol, M. J. *et al.* Prediction of prostate cancer risk: the role of prostate volume and digital rectal examination in the ERSPC risk calculators. *Eur. Urol.* **61**, 577–583 (2012).
  289. Foley, R. W. *et al.* European Randomised Study of Screening for Prostate Cancer (ERSPC) risk calculators significantly outperform the Prostate Cancer Prevention Trial (PCPT) 2.0 in the prediction of prostate cancer: a multi-institutional study. *BJU Int.* **118**, 706–713 (2016).
  290. De Nunzio, C. *et al.* External validation of Chun, PCPT, ERSPC, Kawakami, and Karakiewicz nomograms in the prediction of prostate cancer: A single center cohort-study. *Urol. Oncol.* **36**, 364.e1-364.e7 (2018).
  291. Wallden, B. *et al.* Development and verification of the PAM50-based Prosigna breast cancer gene signature assay. *BMC Med. Genomics* **8**, 54 (2015).
  292. Mantas, D., Kostakis, I. D., Machairas, N. & Markopoulos, C. White blood cell and platelet indices as prognostic markers in patients with invasive ductal breast carcinoma. 1610–1614 (2016). doi:10.3892/ol.2016.4760
  293. Schalken, J., Dijkstra, S., Baskin-Bey, E. & Van Oort, I. Potential utility of cancer-specific biomarkers for assessing response to hormonal treatments in metastatic prostate cancer. *Ther. Adv. Urol.* **6**, 245–252 (2014).
  294. Ahmed, H. U. *et al.* Diagnostic accuracy of multi-parametric MRI and TRUS biopsy in prostate cancer (PROMIS): a paired validating confirmatory study. *Lancet* **389**, 815–822 (2017).
  295. Purysko, A. S., Rosenkrantz, A. B., Barentsz, J. O., Weinreb, J. C. & Macura, K. J. PI-RADS Version 2: A Pictorial Update. *Radiographics* **36**, 1354–1372 (2016).

# 9

## **Appendices**

9: APPENDICES

Supplementary Table 1 Probe list for NanoString2 (n = 167).

<i>Transcript</i>	<i>Accession</i>	<i>Capture Probe</i>	<i>Reporter Probe</i>	<i>Biomarker Type</i>	<i>Source</i>
<i>AATF</i>	<i>NM_012138.3:1175</i>	<i>TCATCATCTTCACTAGAAATCTCCTCA CTTCCCGCATTGGGCTTTGTCCC</i>	<i>CTCTTTGCAGGGACCCTTCTTCGTTGCT GCTTCTTCTTCTACCAGC</i>	<i>test</i>	<i>Cooper NGS</i>
<i>ABCB9</i>	<i>NM_001243013.1:48 8</i>	<i>GGGCCCCAGCGCACTGTTCTTGCCAC ACCAATGGTGG</i>	<i>ACGAAGAGGCACACGAGGGTGATGACC AGCCACGAGGCCCGCAGCCGCCG</i>	<i>test</i>	<i>Cooper NGS</i>
<i>ACTR5</i>	<i>NM_024855.3:1840</i>	<i>CAAGGCATGGCGTGCAGGGCAGTCTC TCTGGAGGG</i>	<i>GGCAGGTACATCTAGCACAATCACAGT CCTGTCACACTGCCAACGTGGCC</i>	<i>test</i>	<i>Cooper NGS</i>
<i>AGR2</i>	<i>NM_006408.2:1365</i>	<i>TGCCTCATCAACACGTCACCACCCTTT GCTCTTCTCCAATTAGTCACAT</i>	<i>TGCCACAGCCTTTCACGTTTCTAAACC CTAGTAACCTCTGATCTCCATC</i>	<i>test</i>	<i>Mills</i>
<i>ALAS1</i>	<i>NM_000688.4:1615</i>	<i>AGTGTTCAGAAATGATGTCCATTTT GGCATGACTCCATCCCGATCCCC</i>	<i>GAGAACTCGTGCTGGCGATGTACCCTC CAACACAACCAAAGGCTTTGCCA</i>	<i>housekeeper</i>	<i>Cooper</i>
<i>AMACR</i>	<i>NM_014324.4:2145</i>	<i>TGGAATCTACCCCTCCTCACATGCCT TTAGGAAGTTGAGTCCAGGGAAG</i>	<i>CAACATCCATTCTACTCCCTCTACTC TGATGGCACCCGGATTAGATTG</i>	<i>PCa positive control</i>	<i>Cooper</i>
<i>AMH</i>	<i>NM_000479.3:1626</i>	<i>TTGGCCTGGTAGGTCTCGGGGATGAG TACGGAGCG</i>	<i>CGGACTGAGGCCAGCCGCACACGCCCT GGCAATTG</i>	<i>test</i>	<i>Sanda</i>
<i>ANKRD34 B</i>	<i>NM_001004441.2:14 60</i>	<i>TTTATAGGATAGTTCTCCTCTGGTGT AATATCCTGGAGCTCCTCTTGCA</i>	<i>ATGCTTTGGTGCCTAGTGATGAACCGC TTGGAAAGTGCCAGCCCATTGGT</i>	<i>test</i>	<i>Sanda</i>
<i>ANPEP</i>	<i>NM_001150.1:2670</i>	<i>GTAATGCTGATGATGGTAGAGGTGGC GTCCTGCTTCCGGATTAAAGTC</i>	<i>AGTTGCTCTGGACAAAGTCCCAGACCA GACCTTGCCCAATGACGTTGTTG</i>	<i>test</i>	<i>Mills</i>
<i>APOC1</i>	<i>NM_001645.3:32</i>	<i>CGGAGGGGCACTCTGAATCCTTGCTG GAGGGCTTGGTTGGGAGGTC</i>	<i>CAGAACCACCACCAGGACCGGGAGCGA CAGGAAGAGCCTCATGGCGAGGC</i>	<i>test</i>	<i>Sanda</i>
<i>ARexon9</i>	<i>NM_000044.2:3401</i>	<i>GACTTGTGCATGCGGTAATGAAA ACCAGATCAGGGGCGAAGTAGAG</i>	<i>CAAACCTTTGAGAGAGGTGCCTCATTC GGACACACTGGCTGTACATCCGG</i>	<i>test</i>	<i>Cooper</i>
<i>ARexons4- 8</i>	<i>ENST00000514029.1 :3171</i>	<i>TTTGAAGAGAGGGGTTGGCTGGCTTCT TCTCCTGGAGAAGCAGAAATCTG</i>	<i>CAGTAAGGCTAGATGTAAGAGGGAAAG TCGGACTGTAGTCTCTCAGTGTG</i>	<i>test</i>	<i>Cooper</i>
<i>ARHGEF2 5</i>	<i>NM_001111270.2:11 02</i>	<i>CAGCGCTTGGGCACAAAGCACATGAC CTCCACAGCTTG</i>	<i>CTCAAATCCCCGCAATCTCCCCAGCGT CATCATATCGTTG</i>	<i>test</i>	<i>Cooper</i>
<i>AURKA</i>	<i>NM_003600.2:405</i>	<i>AAGGAAATTGCTGAGTCACGAGAACAC GTTTTGGACCTCCAAGTGGAGCT</i>	<i>ACACAAGACCCGCTGAGCCTGGCCACT ATTTACAGGTAATGGATTCTGAC</i>	<i>test</i>	<i>Cooper</i>
<i>B2M</i>	<i>NM_004048.2:25</i>	<i>CACGGAGCGAGACATCTCGGCCCGAA TGCTGTCAGCTT</i>	<i>CAGGCCAGAAAGAGAGAGTAGCGCGA GCACAGCTAAGGC</i>	<i>housekeeper</i>	<i>Cooper</i>
<i>B4GALNT</i>	<i>NM_178537.4:492</i>	<i>TCCCTCGCCGGGTGGATGAAACCAAAA</i>	<i>CAGAACTCCGAGTTGTCGTCTGAGGCC</i>	<i>test</i>	<i>Sanda</i>

9: APPENDICES

4		<i>ATACGGAGTCCATAGTTCTTCCA</i>	<i>ACAGAAAACCTGGACGTCTCCG</i>		
<i>BRAF</i>	<i>NM_004333.3:565</i>	<i>AGTGTCTTTCTTTAGACTGTCTCGGACT GTAACCTCCACACCTTGCAGGTAC</i>	<i>CCTGAATTCTGTAAACAGCACAGCACT CTGGGATTAGACCTCTCATCATC</i>	<i>test</i>	<i>Cooper</i>
<i>BTG2</i>	<i>NM_006763.2:1700</i>	<i>CAAGGAATACATGCAAGGCTGACTAGC CAGCCATCATCCCAAGGAGAG</i>	<i>ACAAGAATACCAAGTAGTCTTGCAGAA CATGGGGCACTCTCCATTTCAGC</i>	<i>test</i>	<i>Sanda</i>
<i>CACNA1D</i>	<i>NM_000720.3:6044</i>	<i>GTACTTCTGGGCTTACTTGAATCTAG GCCGGCAACTGCCATGATCTGTT</i>	<i>GTTGCTGGAGGGGTGGCCACGACCG GGTCGAGTGACTCGGTGA</i>	<i>test</i>	<i>Sanda</i>
<i>CADPS</i>	<i>NM_183394.2:1870</i>	<i>TTGAGGCTTATCCATTCCGACAGCAAG TTTGATTTGAGATCTTGGTCCG</i>	<i>TTCCAGACATTCTTACCGATGGCCATA AATACCCAGAATGCTTCATGTT</i>	<i>test</i>	<i>Sanda</i>
<i>CAMK2N2</i>	<i>NM_033259.2:908</i>	<i>AAATACAAATGTGCTGAGGAAGTCCCT TAGAAAGAGGCTGAGGCTGGGGT</i>	<i>GGGAGGGCAGGAACCATGAGCAGAGC CAGTAAACAAAGAGTCGGATATAA</i>	<i>test</i>	<i>Sanda</i>
<i>CAMKK2</i>	<i>NM_006549.3:1710</i>	<i>GGTGGATGATCTTCTGGTAGTGTAAAGT ACTCGATGCCTTTGATCAGATCC</i>	<i>CTTGATGTGCCATCTTCTCCGACCAG GAGGTTGGAAGGTTTGTATGTCAC</i>	<i>test</i>	<i>Cooper</i>
<i>CASKIN1</i>	<i>NM_020764.3:1664</i>	<i>ACCTTGTAGTACTGGGCCAGGCCGATC ATGGACAG</i>	<i>AGGTGATGTCGGTGATGAAATCAATGT TCTCGTAGCCATTGTCCACCAAC</i>	<i>test</i>	<i>Sanda</i>
<i>CCDC88B</i>	<i>NM_032251.5:400</i>	<i>TCCACCGCTTCTTCTGAGAGAGGGTCA AATCCCAATGTCTG</i>	<i>TGACGCTCCCAACAGTAGCCGAAGAAC GCCTTCCAGCTGC</i>	<i>test</i>	<i>Cooper - NGS</i>
<i>CDC10</i>	<i>NM_000902.2:5059</i>	<i>TAGGGCTGGAACAAGGACTCTTTTCTC TGGACAGCTTGCACCTACAATCC</i>	<i>CCAAAGGAATATTGCAAATACCCAAGG TCACCCTGTCAGGAGTGGCAGAA</i>	<i>test</i>	<i>Whitaker</i>
<i>CDC20</i>	<i>NM_001255.2:430</i>	<i>CCTCTACATCAAAACCGTTCAGGTTCA AAGCCCAGGCTTTCTGATGTTCC</i>	<i>ACCCTCTGGCGCATTTTGTGGTTTTCCA CTGAGCCGAAGGATCTTGGCTT</i>	<i>test</i>	<i>Mills</i>
<i>CDC37L1</i>	<i>NM_017913.2:1146</i>	<i>TCATCTTCTTTATGTACCACCGAGTTTA AGCTGCAGAGAGCTGTACTGAT</i>	<i>GGCCTCAGCAGTCTTAACCAAATTATA CAGTGTCCATCATTTTGGGTTCA</i>	<i>test</i>	<i>Sanda</i>
<i>CDKN3</i>	<i>NM_005192.3:510</i>	<i>AGACAAGATCTCCCAAGTCTCCATAG CAGTGTATTAAGGTTTTTCGGTA</i>	<i>CTCTGGTGATATTGTGTGACAGAGGTA TAGTAGGAGACAAGCAGCTACA</i>	<i>test</i>	<i>Mills</i>
<i>CKAP2L</i>	<i>NM_152515.3:1120</i>	<i>TGAGGTATACAAACTTGGCTGGACTTC TGATCTTGCTTGATGTTTGGATG</i>	<i>AATTAGGCCTCTGGCTTATGGCTTTTGA CTTTTGCAGTACACATGATGTC</i>	<i>test</i>	<i>Cooper</i>
<i>CLIC2</i>	<i>NM_001289.4:50</i>	<i>CCAGTCTCTTCTCTCAAGAGGTGTGAC GCAGAAAATTCTAGATGCTTAAG</i>	<i>TGCTTTAAGAAGACCGTCTAGCTTGTA GTGGACTGAGTCAGACCTGGAG</i>	<i>test</i>	<i>Cooper</i>
<i>CLU</i>	<i>NM_203339.1:2460</i>	<i>GCCTGTGGTCCAGGGAAAGGTATGAA GATCATATAAACCGGCGGTGGACA</i>	<i>AGCGTAGGGTACTGCAGCCCAGCTATG GTTCAGACTAAAAGCCGAGAAAC</i>	<i>test</i>	<i>Cooper</i>
<i>COL10A1</i>	<i>NM_000493.3:135</i>	<i>CCTGTGGGCATTTGGTATCGTTCAGCG TAAAACACTCCATGAACCAAGTT</i>	<i>TGTAGGGAATGAAGAACTGTGTCTTGG TGTTGGGTAGTGGCCTTTTATG</i>	<i>test</i>	<i>Sanda</i>
<i>COL9A2</i>	<i>NM_001852.3:795</i>	<i>CGATAGCGCCACCATGCCTTTATATC</i>	<i>CCTAGGACCTTCTCACCCGGTGGCC</i>	<i>test</i>	<i>Sanda</i>

9: APPENDICES

		<i>CATGAGGGCCCGTCTCTCCCTTG</i>	<i>AGTGGCAC</i>		
CP	NM_000096.3:1110	<i>CTTGCCCGTGAAAGAAAGCTGCGTGCA</i> <i>CATCAACTTCATTACCCATACCA</i>	<i>AGCAGGAAAGAGGGTTGATTGTGTCAAT</i> <i>ACGGTAGTTCTTGTTAGTCAGTG</i>	test	Cooper – NGS
CTA- 211A9.5/M LATNB	CTA_211A95.1:407	<i>CTGGAGGTATCCAAGAGTCTGCCGAG</i> <i>GGACTTCAAGTATTCAGGAAGGGG</i>	<i>GAAGAGCCCAAACCTGCCTGGCTTCAA</i> <i>AACAGGTGGTGAGCTCCCCATTG</i>	test	Cooper – NGS
DLX1	NM_001038493.1:13 35	<i>CAGCCTCAGGCGAAGTCCATTTCTCAA</i> <i>TAAATAAAACCCCTCCCTCCAA</i>	<i>CGTTTGAACAGTGCGTTCCTTGCGCCC</i> <i>AGCAGAACCCTGAATTGGCAA</i>	test	Schalke n
DNAH5	NM_001369.2:12374	<i>GGCGGAACGCATCATGTACAAGCTCA</i> <i>GTTTCTATGATTATGTCCATCAGC</i>	<i>CTGAAGGAGTGTAAATGGGAAACTGCTT</i> <i>ATGAGCCTCGGTGGTCATCCAGA</i>	test	Sanda
DPP4	NM_001935.3:2700	<i>AAATCCACTCCAACATCGACCAGGGCT</i> <i>TTGGAGATCTGAGCTGACTGCTG</i>	<i>CTGCTAGCTATTCCATGGTCTTCATCAG</i> <i>TATACCACATTGCCCTGG</i>	test	
EIF2D	NM_006893.2:1600	<i>GCTCTTGTCGGGAAGGGTCACTTGAT</i> <i>AGGCAGGCTGTAATTTTTCCAAA</i>	<i>TTGTGCTAGGGTGATGTCAATTGGACA</i> <i>GATTCTCCCTTTCTTCACAATGG</i>	test	Sanda
EN2	NM_001427.3:2576	<i>AAGGTAGCCACATGTTTCAGAAGTGTG</i> <i>GACTCAAACACGCCTGGTGTGTG</i>	<i>CTTCTTCCTTCTTCTAGATCCTGGAGG</i> <i>ATTCTGAGTTCTTTTGAAAGAC</i>	test	Pandha
ERG 3' ex 4-5	NM_001243428.1:17 7	<i>CCATCTTTTTTCTCTGTGAGTCATTTGT</i> <i>CTTGCTTTTGGTCAACACGGCT</i>	<i>CCATCTACCAGCTGTTCAGAACCTGAC</i> <i>GGCTTTAGTTGCCCTTGGTTCTG</i>	test	Cooper
ERG3' ex 6-7	NM_004449.4:477	<i>TGAGCCATTACCTGGCTAGGGTTACA</i> <i>TTCCATTTTGATGGTGACCCTGG</i>	<i>CCACCATCTTCCCGCCTTTGGCCACACT</i> <i>GCATTCATCAGGAGAGTTCTT</i>	test	Cooper
ERG5'	NM_182918.3:697	<i>ACATCATCTGAAGTCAAATGTGGAAGA</i> <i>GGAGTCTCTCTGAGGTAGTGGAG</i>	<i>CTGTGTTTCTAGCATGCATTAACCGTG</i> <i>GAGAGTTTTGTAAGGCTTTATCA</i>	PCa positive control	Cooper
FDPS	NM_001135822.1:40 4	<i>CATCCTGTTTCCTTGGCTCCACCAGCT</i> <i>CCCGGAATGCTACTAC</i>	<i>CCAGCCCACAGTCCAGGCCCGCTGGAG</i> <i>ACTATCAG</i>	test	Sanda
FOLH1	NM_004476.1:695	<i>TGAAAGGTGGTACAATATCCGAAACAT</i> <i>TTTCATATCCTGGAGGAGGTGGT</i>	<i>GTTAACATACACTAGATCGCCCTCTGG</i> <i>CATTCCTTGAGGAGAGAAAGCAC</i>	test	Mills
GABARAP L2	NM_007285.6:340	<i>GGGACTGTCTTATCCACAAACAGGAAG</i> <i>ATCGCCTTTTCAGAAGGAAGCTG</i>	<i>CTTCATCTTTTTCTTCTCGTAAAGCTG</i> <i>TCCCATAGTTAGGCTGGACTGT</i>	test	Clark
GAPDH	NM_002046.3:972	<i>AAGTGGTCGTTGAGGGCAATGCCAGC</i> <i>CCCAGCGTCAAAG</i>	<i>CCCTGTTGCTGTAGCCAAATTCGTTGTC</i> <i>ATACCAGGAAATGAGCTTGACA</i>	housekeeper	Cooper
GCNT1	NM_001097633.1:39 4	<i>TTTCAAACAATAATCAGGGATTTCTT</i> <i>TGTGAAGGGCAGTCTTCTATGCT</i>	<i>GTATTTGGTGGGATAAGAAAAAAGTCT</i> <i>CCTTCGCAGCAACGTCTCAGCA</i>	Test	Sanda
GJB1	NM_000166.5:190	<i>TGAAGATGAAGATGACCGAGAGCCAT</i> <i>ACTCGGCCAATGGCAGTAGAATGC</i>	<i>TTTCTCATCACCCACACACTCTCTGCA</i> <i>GCCACCACCAGCCATGATTC</i>	test	Sanda



9: APPENDICES

<i>GOLM1</i>	NM_016548.3:508	GGATGAGCCTCTCACCTGTGGTGATGT TATTCACCAAAACCGC	TAATTCCTCTGCAGGGTCTTTAACTGGT CTTGCACTC	test	Cooper
<i>HIST1H1C</i>	NM_005319.3:401	CTTGGCTGCCCAACTGGCTTCTTAGG TTTGGTTCCGCCCGCCTTTTAA	TTCGGAGTTGCGCCGCCAGCCGCCTTC TTGGGCTT	test	Sanda
<i>HIST1H1E</i>	NM_005321.2:172	GCGCTCCTTGGAGGCGGCAACAGCTTT AGTAATGAGCTCGG	CTGCCAGCGCTTCTTGAGAGCGGCCA AAGATACGCCGCT	test	Sanda
<i>HIST1H2B</i> <i>F</i>	NM_003522.3:313	CTTGGTGACGGCCTTGGTGCCCTCTGA CACGGCGTG	AGCCTTTGGGATTGGGTATGAAGACGT TAGAATTACTTAGAGCTGGTGTA	test	Cooper – NGS
<i>HIST1H2B</i> <i>G</i>	NM_003518.3:318	TATACTTGGTGACAGCCTTGGTACCTT CGGACACTGCGTGCTTGG	AAGAGCCTTTGAGTTTTAAAGCACCTA AGCACACATTTACTTGGAGCTTG	test	Sanda
<i>HIST3H2A</i>	NM_033445.2:114	CGGAGCAACCGGTGCACGCGGCCAC GGGAACTG	CGCCGGCGCCACGCGCTCCGAATAGT TGCCCTTG	test	Sanda
<i>HMBS</i>	NM_000190.3:1020	GCTGGGCAGGGACATGGATGGTAGCC TGCATGGTC	AGTGATGCCTACCAACTGTGGGTCATC CTCAGGGCCATCTTCAT	test	Clark
<i>HOXC4</i>	NM_014620.4:1058	TGAATTTTTTTCATCCATGGGTAGACT ATGGGTTGCTTGCTGGCGGCG	CGCTTGGGTTCCCTCCGTTATAATTG GGGTTACCGTGCTAACG	test	Schalke n
<i>HOXC6</i>	NM_153693.3:570	GGTCGAGAAATGCCTCACTGGATCATA GGCGGTGGAATTGAGGGCGACGT	GAATAAAAGGAGTCGAGTAGATCCGG TTCTGGGCAACGGCCGCTCCATA	test	Schalke n
<i>HPN</i>	NM_182983.1:1870	CCGAGAGATGCTGTCTCACACAAAA GGGACCACCGCTG	CCAATCACAATGCCACACAGCCGCCA ACGTGGCGT	test	Cooper
<i>HPRT</i>	NM_000194.1:240	TGAGCACACAGAGGGCTACAATGTGAT GGCCTCCCATCTCCTTCATCACA	CAGTGCTTTGATGTAATCCAGCAGGTC AGCAAAGAATTTATAGCCCCCT	housekeeper	Cooper
<i>IFT57</i>	NM_018010.2:790	AATCGTGACTTTCAAGTTGCGGTAGTAC ACGTTCCACTTCTAGGCTCCATT	TGCTGGTGCAATTTGGTCAACATGGATT CTCCAATCCTTATTGTCAGTCT	test	Sanda
<i>IGFBP3</i>	NM_000598.4:1255	CGGGCGCATGAAGTCTGGGTGCTGTG CTCGAGTCTCTGAATATTTTGATA	TGGTCGGCCGCTTCGACCAACATGTGG TGAGCATTCCA	test	Sanda
<i>IMPDH2</i>	NM_000884.2:545	TCTTTGAGAAAATCAATGTCCCTGGAG GAGATGATGCCACCAAGCGGCT	TCCCTCTTTGTCATTATCTCTCCAAGA AACAGTCATGTTCTCC	test	Mills
<i>ISX</i>	NM_001008494.1:31 40	ATCTGGCATTTTAAGATGGCAAAGCA CTTTGTCATCCTGTGGGCTGTTG	TGCTAGAGACCTGGTGTGATATCCAC ATTCATAGGCTCTGAGTG	test	Sanda
<i>ITGBL1</i>	NM_004791.2:1317	AGACCACACCATCGAGGTCTTCACAGC GGCGATCATCACTCAAGTC	TCCTCTCTCAAACACAGCGACCACA GGAACATGTGCCGTGGCCTCCAC	test	Sanda
<i>ITPR1</i>	NM_001099952.1:67 75	GACAATCTCTATCTGCGCCGTGTGCTT GGCATAAACTCCAGGGC	CATATGCTGGGCACGGGAAAGACTATC TGTTCCATTGTTCCGGTCTAATCT	test	Sanda

9: APPENDICES

KLK2	NM_005551.3:1820	CTTGGACACTAAGGATCAGGTGAGCTT CCTCAGTTGGAATTACTTTGTAC	GTC AATTATTCAAGTACTCCATACTCGT CCTACAGACCCCCAGTAAAAAC	test	Cooper
KLK3/PSA (exons1-2)	NM_001030048.1:16	TGAGGAAGACAACCGGGACCCACATG GTGACACAGCTCTCCGGGTG	AATCCGAGACAGGATGAGGGGTGCAGC ACCAATCCACGTCACGGACAGGG	Prostate control	Doll
KLK3/PSA (exons2-3)	NM_001648.2:209	ATCACGCTTTTGTTCCTGATGCAGTGG GCAGCTGTG	CCTGTGTCTTCAGGATGAAACAGGCTG TGCCGACCCAGCAAG	Prostate control	Cooper
KLK4	NM_004917.3:410	CCCAGCCAGAAACGAGGCAAGAGTTC CCCGCGGTAG	CAGCACGGTAGGCATTCTGCCGTTCCG CAGCAGAC	test	Cooper
LASS1	NM_198207.2:1918	GCATCTCGCACCTCCCGTTCCAAAAA CGTCACGGAGCTCTGAG	CTGCCTGGCTACAGCCCCGGATGTGTT AAATGTCT	test	Sanda
LBH	NM_030915.3:2340	GAGAGTATGGATGAACCACTCTCTGCA GCCAAAACAGAACGAAGCGGGGA	ACAGGAATTGAAAAGGCAAGACCCCG TCCACAAGGGGAGGCGAGGGAAAT	test	Sanda
MAK	NM_005906.3:1395	TATCTCCAGACTTGAAGATAGTCTGAC CCCAACGCCTCTACCACCTTTA	CTTCTTGGAATGGGAGGCTCCGAAATC ATAGTCTCCA ACTCTTCCAGC	test	Sanda
MAPK8IP 2	NM_012324.2:1885	CTCTCGCTCCTCGCCGTTGACCAGACA GGAGAAAAGGCCAAAAGGACTCG	CCGCGGGATGAACCTGAACACAGCCCG GTGAGTCTG	test	Sanda
MARCH5	NM_017824.4:2136	TGTGCTGAAACTAGACTGTCAACTCTG TAAGAGCTTGGACCAAGTCTGTC	AAACAAAGAGCTCAAGGCCTCACCTTG GTTTATTCACTGCTGGTTTTCTA	test	Sanda
MCM7	NM_182776.1:1325	TGTGTTCTCTCCTTCTACCAGCACCGT GATACTACGAGGGATATT	CAAGAAAATACCAGTGACGCTGACGTG GTCTCCAGGCTGGGCAATCCT	test	Perry
MCTP1	NM_024717.4:1005	AACTCCAATTGTGTCAGATCCAGAAAG GCTGAGCCATAAAGTCATCCTG	GATAATGAGGATCTTTCAGAGTAAGGG TCACATCTGTGGGCCTGTTT	test	Cooper – NGS
MDK	NM_001012334.1:71 1	CGAGCAGACAGAAGGCACTGGTGGGT CACATCTCGGGC	GGGGCTGGGGAGTGAGAGGGACAAGG CAGGGCATGATTGATTAAAGCTAA	test	Cooper
MED4	NM_001270629.1:32 4	TCTTGCTTTTTCTATTGACTTGAGTTTC TCCTTCGCTTGGTAAACAGCTG	CTGATCCTATGTGCATACTTAATTATTT CTTCAGAGGAGATAGCACTTT	test	Sanda
MEMO1	NM_001137602.1:11 92	GAATGTGCAGGTGGCATCCCTGAGGA TTCAGAGCT	TATCGTGGTAAAGGCTAGGCTGGGACC CCGGACAGAGTATGA	test	Sanda
Met	NM_001127500.1:19 25	AAATTTATTATTCCTCCGAAATCCAAA GTCCCAGCCACATATGGTCAGCC	GTCAAGGTGCAGCTCTCATTTC AAGG AGAACTCTAGTTTTCTTTAAATC	test	Cooper
MEX3A	NM_001093725.1:20 90	GATCTATGCAACTTCTGATAGGACTCC AACTCCCTTACACTGCTGGAAAC	CCTTTCAGCCACAGAAACGATTGACAT GCTTCTCTCCCAACCCCTAGAA	test	Sanda
MFSD2A	NM_032793.4:592	AAGAGGCAATAGAAAAGCAGGTACCA ATAGGTCTGGCCGTGTGGGAAGTC	ACATGGTGAGAGCCGAGTAGGGAACAT GGAAACACGTGACCATTGTTTCA	test	Cooper – NGS

9: APPENDICES

MGAT5B	NM_144677.2:3392	GGTTGGAACAAGCAGGAGAGAGAAAC AATTC AACAGGGTCTGGGTGGTC	CAGGTCATGCCAGGATGGGTTTTGGGA GAAGCCCAGAGTGAAAAAG	test	Sanda
MIC1	NM_004864.2:180	CCTGGTTAGCAGGTCCTCGTAGCGTTT CCGCAACTC	GTGTTCCAATCTTCCCAGCTCTGGTTG GCCCCGAG	test	Whitaker
MIR146A/ DQ658414	ENST00000517927.1 :1642	CGGTTGAGATTCACCAAGGTTCTGGT TCTGGAATGAGTCACTGGCTAAG	TTCTGGATTTTCTCCATCAGTCTAGGAC TGAAGACACCGATCTCTGGTGT	test	Cooper – NGS
MIR4435- IHG/IOC5 41471	ENST00000409569b. 1:45	AAAGCAGCGACCATCCAGTCATTTATT TCCCTCCATTCCCAATGATGTAC	CAGGCACGGGCTCAGGCACCGCTTGTC TGGAATGTCAATTTGAAACTTAA	test	Cooper – NGS
MKi67	NM_002417.2:4020	CTGATGGCATTAGATTCTGCACGCTA AGAGTTCTCCCTCTACATCTG	GTCTTTCTCTTACCTACTGATGGTTTA GGCGTGTGCATGGCTTTGCCTG	test	Cooper
MMP11	NM_005940.3:702	TCAGTGGGTAGCGAAAGGTGTAGAAG GCGGACATCAGGGCCTTGG	ATATAGGTGTTGAACGCCCTGCAGTC ATCTGGGCTGAGAC	test	Sanda
MMP25	NM_022468.4:2955	CATTTAGATCCTAAAAGTGTGGGGAGT GGGGACAGGGTGAACGAGGTGCC	CCCAGTGATTCTGATGTGGGATAGTCT AGAAGAATAGTTCCAGAGGCAAT	test	Cooper – NGS
MMP26	NM_021801.3:515	CAGGATTTCCAGAATTTGGTAAAAAGG CATGGCCTAAGATACCACCTGGC	TCCAGTGTCTGAAGCTGACCAGTGTC ATTCTTGTCAAAATGGACAATC	test	Cooper
MNX1	NM_005515.3:1680	TTTCTGAAGAGCAGGTGAGGCGCCCT TGCTTAAAAGGGAAGCGCCAGG	TTAAAAGAACCAGAGTTCAAGTTTCAG CCCCCTGGGTCTCCCTCTCGCTG	test	Sanda
MSMB	NM_002443.2:295	TTTTTGGGTCTTCTTCTCCACCACGA TATACTTGCACTCCTCCTTCTTG	GTGCCTACTAGAAGCACATTAGATTAT CCATTCCTGACAGAACAGGTCT	test	Whitaker
MXI1	NM_001008541.1:61 5	GAAGTGAATGAAAGTTTGACACTGGCA CTGGAGTAAACCCTCGTCACTCCC	TGGCCCAGTGAATATTTGCCCTGCAC TGTTATGTCATGCTGGGTTCTAT	test	Sanda
MYOF	NM_013451.3:5805	ATGATCGTGTGACGCAAGTCAAGTTCT AGGAAACCCAAGTAGTCATCCAG	TGAGGTCCGGAATCATGTCCAATCTGC ATTTCTCTGGTGATTTGCAGGA	test	Cooper – NGS
NAALADL 2	NM_207015.2:250	ATTCTCAGCACCGTCTAGCTGGAATTG GTCAAAACCAGACTCCTCTAGTT	TGAATGGAATCAAGATTGAGGTCTATA GTCTCTGAATGCCCTAGGTTCTG	test	Mills
NEAT1	NR_028272.1:1850	TTTCTCACACACAGATTTAGGAATGAC CAACTTGTACCTCCAGCGTTT	TTCTCCTAGTAATCTGCAATGCAATCAC AATGCCCAAACCTAGACCTGCCA	test	Sanda
NKAIN1	NM_024522.2:1620	CACTGTGTTCAAGGCCCACTTCCACCA AAAATCTAGCTGTGTGGCCTCAA	GAAGTCAAGAGAGCAGACTGGGTTTT ACAGTCAGAACTGCAGAAAGTA	test	Sanda
NLRP3	NM_001079821.2:41 5	CTGGCATATCACAGTGGGATTCGAAAC ACGTGCATTATCTGAACCCACT	CTCGAAAGGTACTCCAGTAAACCCATC CACTCCTCTTCAATGCTGTCTTC	test	Cooper – NGS
OGT	NM_181672.1:1080	CTTGAGAGCATTGGCTAGGTTGCAGT	ACGGAGAGCTGTATTATAACAATCTTCT	test	Cooper

9: APPENDICES

		<i>AAGCATCAGGGAAATGTGGTTGT</i>	<i>GCTTCAGCAACACTGCCCTTCT</i>		
PSGR	NM_030774.2:360	<i>GAGCGTGCAGGCTGCGTTCGGTCCCTTA CGATGAAGACCACGATGCAGTTT</i>	<i>GGATAAGGCCAGGTCAATGGCTGCAAG CATGCAGAGAAAGAGGTACATCG</i>	test	Doll
PALM3	NM_001145028.1:23 4	<i>AGCTGGGACTGGAGTGTGAACAAACT GTCTTCCAGGTTCCG</i>	<i>GCTGGGCACCTGTGGAAGCACTTTGCA ACAGTTGC</i>	test	Cooper – NGS
PCA3	NR_015342.1:362	<i>TAAGGAACACATCAATTCATTTTCTAA TGTCCTTCCCTCACAAAGCGGGAC</i>	<i>TCCCGTTCAAATAAATATCCACAACAG GATCTGTTTTCTGCCCATCCTT</i>	test	Cooper
PCSK6	NM_138320.1:1112	<i>ACATCGCCGTCCAGCATGCGGATGCCT CCTATTTTGGCATTGTACGCTAT</i>	<i>CGATGTAGTTGGGTCTGATGCCAGCG ACTTTCCTCGACCACATCTGTG</i>	test	Sanda
PDLIM5	NR_046186.1:120	<i>CTCAAAGTCCAATGACAGAAAATGAAA TATGCTCGGGTCCGGCGCGGCGC</i>	<i>GGCCAACCAGTGACACACTGTAGTTGC TCATGGTTCTAATGG</i>	test	Sanda
PECI	NM_006117.2:940	<i>GAAAACCTTCAGTAACAAGTCTTGAGC ACATGCCTCTCCCGCTGTAACT</i>	<i>CAAATGCCTTCAGCCTGGTCCAGACTT CTTCTGAAAAGTGCTATCAGG</i>	test	Mills
PPAP2A	NM_176895.1:1215	<i>GTGATTGCTCGGATAGTGATTCCCAGT TGTTGGTGTTCATGCAGAGTTG</i>	<i>TTAGAAAACAGGCCAGCTTCACCTGGG CACCTGCTGCCTTCAAGGCTG</i>	test	Mills
PPFIA2	NM_003625.2:3670	<i>CACTTTCATCCAGTCGCCTTTCAGTTC CCAGGGCCAAGAGGTTATTGTAT</i>	<i>AGGAGGAACTGCCTTCTCCAGGTTGA TCCACGTCTGAAGTTCTTGTCAT</i>	test	Sanda
PPP1R12 B	NM_001167857.1:13 05	<i>TGCTCTGTGATACTACTTTGCTTTCA GAGTTGGAATGATTGACAAAGGC</i>	<i>CTAGCAGAAGAGGCAGAGAAGGTATTT TGAGCTGGTGCTGGTATC</i>	test	Cooper – NGS
PSTPIP1	XM_006720737.1:35 2	<i>TCAAAGGAGGCCCTCAGGGAGTTGAT CTCCGTCTG</i>	<i>AGCTGCCACATTCTCCATTTGCTGCTT CAAGGAG</i>	test	Cooper – NGS
PTN	NM_002825.5:418	<i>TTTCTTCCCTGCTTCAGCAGTATCCAC AGCTGCCAGTATGAAAATGAATG</i>	<i>CCATTCTCCACAGTCAGACTTCTTCACT TTTTTTCTGGTTTCTC</i>	test	Sanda
PTPRC	NM_080923.2:154	<i>CAAGAGTTTAAGCCACAAATACATGGT CATATCTGGAAGTCAGCCGTGTC</i>	<i>CTTGGCCCTGTCACAAATACTTCTGTGT CCAGAAAAGGCAAAGCCAAATGC</i>	Blood control	Cooper
PVT1	NR_003367.2:0	<i>AAAATACTTGAACGAAGCTCCATGCAG CTGACAGGCACAGCCATCTTGAG</i>	<i>AGCGTTATTCCCAGACCACTGAAGAT CACTGTAAATCCATCAGGCTCAG</i>	test	Sanda
RAB17	NR_033308.1:1310	<i>ACAGCACTTCTCTGGGAGCCATGTGAC GCCAGATCTTCTCTGGCAGTTC</i>	<i>GGAACAGGCACAGGCATCGGGGAATCA GATGGTATCAGTGGGGATAGGGC</i>	test	Sanda
RIOK3	NM_003831.3:1920	<i>CTGGA AAAACTGCGAGACATTCCTGCA GTCCCGGAACAAGA ACTCCAGGC</i>	<i>ACAGCATTGAAGAGTTCTCGTTCATA AGGGCTTCTTGACTCCTCCTT</i>	test	Clark
RNF157	NM_052916.2:618	<i>ACTAGAGGGTAAACTTCTCGGTCTAAA TCAAAGCCAAGTCTCTTCGGC</i>	<i>CATGGCAATGGCCAAAATACTCGTCTC CTTCATCCACCACGGCATGTACC</i>	test	Sanda
RP11-	ENST00000561140.1	<i>TTGCCAGTCGCTGGTTTTTCATCCAGAG</i>	<i>CAGCAATATATCCTGTTTCATCTTCTCA</i>	test	Cooper

9: APPENDICES

97O12.7	:110	CACGAAGCTCGTGGTCTGAATAC	TCATGAAGGTCAGCTTTCTTCT		– NGS
RPL18A	NM_000980.3:177	GAGATACAAAGTACCAGAAAGCGGGAC TTGGCGACGACATGATTAGGCGCA	CTGCCACAGTAGACAATCTCCCCTGA AGACTTCTTCATCTTCTTTAACT	test	Cooper – NGS
RPL23AP5 3	NR_003572.2:3226	AAATCCGAAAGGATCTCATCCATTAG GACCCTTGTCTCCTTTTCTGTTG	CATTTATGGCTGTCAACCCGCCAGTTCT CAGGAGTTTGTATAAAAGCCT	test	Sanda
RPLP2	NM_001004.3:186	CTGATAACCTTGTGAGCCGGTCGTCG TCCGCCTCGATAC	TGCCAATACCTGGGCAATGACGTCTT CAATGTTTTTCCATTCAGCTCA	housekeeper	Whitaker
RPS10	NM_001014.3:219	GAAATGTCTCCAGGCAAACTGTTCTT CACGTAGCCTCGGGACTTGAGAG	TGAAGGTAATCACGGAGATACTGGATA CCCTCATTGGTAAGGTACCAGTA	test	Cooper – NGS
RPS11	NM_001015.3:105	CAGCAGGACCCTCTTCTTGTGTTTAAA GATGGTCGGCTGCTTTTGGTAGG	AGACCGATGTTCTTGTAGTACCGCGGG AGCTTCTCCTTGCCAGTTTCTCC	test	Cooper – NGS
SACMIL	NM_014016.3:685	AGAAAGTTCTCTTAGAAGATGACCATT CCATACAAACCGCTGATCTGCCC	ATAAAGCCATGTAACACTGGAAGGGCA AACCGATGAACCTCTGGCTGTGC	test	Sanda
SChLAP1	NR_104320.1:359	CCAGGTACATGGTAAAAGTGCCTTATA CAGGTTGAATAAAAATCACTGCC	ACTTGTGTCCCAGCATCTAGATTGCT GAAAAAGATGTAGATGTTGCTT	test	Sanda
SEC61A1	NM_013336.3:2245	CTCTAAGCCCAACCAGAAAGAGTCAGCT AGAAGAGCCAATAGGTGCACAGA	GAGCTGATGACCCAAGTGGACTAAACA CGGAGCTAGCAGAAAACAGGCAGA	test	?
SERPINB5	NM_002639.4:90	CGGGCCTGGAGTCACAGTTATCCTGGA AAATGCGTGGAAAAGGAACAGGC	GAACAGATCAACGGCAAAAGCCGAATT TGCTAGTTGCAGGGCATCCATTG	test	Cooper
SFRP4	NM_003014.2:1060	CAGCCTCTCTCCACTGTATGGATCT TTACTAAGCTGATCTCTCCATT	CCCGGCTGTTTTCTTCTTGTCTGAACT GTTCTCCGCTGTTCTG	test	Sanda
SIM2.long	NM_005069.3:2099	TTAATGTAGGTCGTGCGCATTGCCGG GCTCGGTGGCGCCGAGCC	ATCCGCAAGTCGGCGGCGGGTCCAAT TCAAACAGCTGTCTCTGCATAAA	test	Sanda
SIM2.short	NM_009586.3:2220	CTGCCACCCACCGCCATGGCTGCTTCG GCTCCCGG	GAAGCAGAAAGAGGGCAAGTTTGCCCA AAGCGTGAGGGTTCTGTCTCCAT	test	Sanda
SIRT1	NM_012238.4:1595	GGTGTGGGTGGCAACTCTGACAAATAA GCCAATTCTTTTTGTGTTCTGTTG	CTGGTGGTGAAGTTCTTTCTGGTGAAC TTGAGTCTTCTGAAACATGAAGA	test	Sanda
SLC12A1	NM_000338.2:3380	CCATATACAACAAATCCGATATGGATC CCTTCTTGCCACGGGAAGGCTC	TCTAACTAGTAAGACAGGTGGGAGGTT CTTTGTGAGGATTCCAACCAAG	Kidney control	Cooper, Mills
SLC43A1	NM_003627.5:925	TTGACTTCTCAGGGGCAGGAAAGGCT TCGATGGGCCAGTTGAGGGTGCA	CTTGTGGTCCAGGGCCAGCCACTCAG CTTGATCTTCTTCGTGTA	test	Sanda
SLC4A1 S	NM_000342.3:2770	CATCATCAGCATCCAGACACTGAAGCT CCACGTTCTGAAGATGAGCGG	CACTTCGTCGTATTCATCCCGACCTTC TCCTCATCAAAGGTTGCCTTG	test	Clark
SMAP1, ex	NM_021940.3:1075	GAGTACTTTGCTGTTGAATGGTCTCTG	TGGTCTTGTGAGGTAATGGTATATT	test	Cooper

9: APPENDICES

7-8		<i>TGCCATACAGAGATAAGATGGAG</i>	<i>TGTGGGTCCCATAAATACACCAG</i>		
<i>SMIM1</i>	<i>ENST00000444870.1:353</i>	<i>TTCATGGCGATGCCAGCTTGCCCGTG CACAGCCTCTGGGAGAT</i>	<i>GGTAGCCCAGGATGAAGATGATCCAGA AGAGGGCCACGCCGCCAGCACC</i>	<i>test</i>	<i>Cooper – NGS</i>
<i>SNCA</i>	<i>NM_007308.2:568</i>	<i>ACTGGGAGCAAAGATATTTCTTAGGCT TCAGGTTTCGTAGTCTTGATACCC</i>	<i>GGAAGTCTGACACTTGTACAGGATGGAA CATCTGTCAGCAGATCTCAAGAA</i>	<i>test</i>	<i>Clark</i>
<i>SNORA20</i>	<i>NR_002960.1:2</i>	<i>CGTATAACTGCTCGTATCACTGTGAGA CTACAAGCAGCAAATAAATGGGA</i>	<i>ATGGTACTTCATCTCAATTTACAGTGG CCCAATGTTATTTTATCCCATG</i>	<i>test</i>	<i>Sanda</i>
<i>SPINK1</i>	<i>NM_003122.2:65</i>	<i>AAGTTCTGCGTCCAGAGGTCAGTTGAA AACTGCACCGCACTTACCACGTC</i>	<i>CAACAGGGCCAAGGCACTGAGAAAGAAA GATGCCTGTTACCTTCATGGCTG</i>	<i>test</i>	<i>Cooper</i>
<i>SPON2</i>	<i>NM_012445.1:1680</i>	<i>CATTTATTCACCTTCTCAAGTGGCCCC GCTTGGATGCGCCCTCG</i>	<i>AACGCAGAGAGATCCATAACATGGAAA CACTGACGCTTCCGAAACCGCCC</i>	<i>test</i>	<i>Whitaker</i>
<i>SRSF3</i>	<i>NM_003017.4:2640</i>	<i>TAAAGTAACTGCCAACTGGGACTGTAT GTCACCTAAGTCAGGATAACTCC</i>	<i>CCATGTTCTAAAGTTTCTAAGAGTCTTG AGGTTATGCTAGGGCTCCTGGT</i>	<i>test</i>	<i>Sanda</i>
<i>SSPO</i>	<i>NM_198455.2:7270</i>	<i>CCACAAGGCAGGGAGAGAAGGGAGCC ACATAAGTAGATTCCTGGCG</i>	<i>ATGGTAGGCATCATGAAGGGCACAGTG CTCGCTGC</i>	<i>test</i>	<i>Sanda</i>
<i>SSTR1</i>	<i>NM_001049.2:2575</i>	<i>TCCGACCCCGCAATCTTATAAAAACTC CTCATTCGGCTTGTCTCAGCTC</i>	<i>GGTCTTTGAAAACGCGCAGTAGGAGGG TGATTCCTATTACGCGCCCACAC</i>	<i>test</i>	<i>Sanda</i>
<i>ST6GALN AC1</i>	<i>ENST00000592042.1:1036</i>	<i>TTTTTCCTCAAATCCACCGAGGCTC AGATTTGAAGTTGGCGGCCTTCA</i>	<i>TTCACAGAGTCAGGGCAAGTCGTCTGA AGGCCTCTATTTCGAAGCTGTA</i>	<i>test</i>	<i>Cooper – NGS</i>
<i>STEAP2</i>	<i>NM_152999.2:845</i>	<i>ATATATAAACCTGCCGGCTGGCATCCT TAGGTCCTAACTGAAGTGCCCAA</i>	<i>CTGGCGGGCAAGTTCAATAACCTGTTG TCGCGCTTGAATATTGTTGCTGC</i>	<i>test</i>	<i>Mills</i>
<i>STEAP4</i>	<i>NM_024636.2:3555</i>	<i>ATCAAAGATAAGTTGAAGGAGCGTGTG TTCTGTGTACCTTTGCAACCAGT</i>	<i>CCATGACTCTACTCAATGTCGTCCAAC TTTTGTATCCTTGCTTGGGTTT</i>	<i>test</i>	<i>Mills</i>
<i>STOM</i>	<i>NM_004099.5:120</i>	<i>GAGTCGGGGAGCCGCTGGGCTTCGGA GTCCCGTGT</i>	<i>CCAAAATCCATCCGCAAGGTCCAAGGC CCTTACTGGGGCTGTCCTTGAAG</i>	<i>test</i>	<i>Clark</i>
<i>SULF2</i>	<i>NM_001161841.1:1206</i>	<i>ATGAGGTCTGTGAGGTAATCCTTGGAG TAGTCGGAGC</i>	<i>GTACATCTTCTTGACGTGCGGAAGAA GCTCACGCTGTCATTGGTG</i>	<i>test</i>	<i>Cooper – NGS</i>
<i>SULT1A1</i>	<i>NM_177534.2:1393</i>	<i>CCCTCAATTCATATTTTATTCTTGAGCC GCTTGGTCAGGTTTGATTGCA</i>	<i>TCAGCCTCCAAATTGCTGGGATTACAG ACATGACCTACCGTCCCGG</i>	<i>test</i>	<i>Mills</i>
<i>SYNM</i>	<i>NM_015286.4:2460</i>	<i>AATGTGACATCGCTTCTCCATAACCT TCCTCCTCCTTAACCAACCCCA</i>	<i>TCGTGTTCTCCTGAGGCTGCTTGGTCC TTCGATGCTGATTAAGTACTGAG</i>	<i>test</i>	<i>Sanda</i>
<i>TBP</i>	<i>NM_001172085.1:587</i>	<i>GCACGAAGTGCAATGGTCTTTAGGTCA AGTTTACAACCAAGATTCAGTGT</i>	<i>TCCTCATGATTACCGCAGCAAACCGCT TGGGATTATATTCGGCGTTTCGG</i>	<i>housekeeper</i>	<i>Cooper</i>
<i>TDRD</i>	<i>NM_198795.1:2615</i>	<i>TGTTTCTAGACTGTATATCTGCTAACT</i>	<i>CCCAGCAACACACATCTGGAATCTTGT</i>	<i>test</i>	<i>Schalke</i>

9: APPENDICES

		<i>GGCACCGTATCCCTGAAAGGGA</i>	<i>TATGGCTTCTTCAGACCAATGTT</i>		<i>n</i>
<i>TERF2IP</i>	<i>NM_018975.3:1100</i>	<i>GCCTGTGTAAGTGTGATAGATCCAAG TAAACTTCTCCATTAAGTCCG</i>	<i>ACGCTAAGAAAGGCGGAAGTAGCCTCCA GCTCACCACATTTTTTAGGAAG</i>	<i>test</i>	<i>Clark</i>
<i>TERT</i>	<i>NM_198253.1:2570</i>	<i>CGCAAGACCCCAAAGAGTTGCGACG CATGTTCTCCAGCCTTGAAGCC</i>	<i>TCTGGAGGCTGTTACCTGCAAATCCA GAAACAGGCTGTGACACTTCAGC</i>	<i>test</i>	<i>Cooper</i>
<i>TFDP1</i>	<i>NM_007111.4:551</i>	<i>TTCCTCTGCACCTTCTCGCAGACCTTC ATGGAGAAATGCCGTAGGCCCTT</i>	<i>TGAACTCCGCAACCAGCTCGTCTGCCA CTTCGTTGTAGGAAGTGGTCCCT</i>	<i>test</i>	<i>Clark</i>
<i>Timp4</i>	<i>NM_003256.2:1000</i>	<i>TCTGCAGGGAAGGAGAAGTGGCTTGA TCTTCAGGACTCTTGAAGGGATGT</i>	<i>GGCACTTCTTATTAGCTGGCAGCAAGA GGTCAGGTGGTAATGGCCAAAGC</i>	<i>test</i>	<i>Cooper</i>
<i>TMCC2</i>	<i>NM_014858.3:1312</i>	<i>ACGTTGCTGCCGTCGGCCAGCAGCAG AGCAGTGTCGGTG</i>	<i>CCCCGATGCCTTCGGCCTCCTCAGCCA GGAGGTAC</i>	<i>test</i>	<i>Clark</i>
<i>TMEM45B</i>	<i>NM_138788.3:469</i>	<i>GCATACAGCAGGAGTGAGTGGATGTG CTGGTCCAGCGGAGGCCGG</i>	<i>GGTCCCGGAAGATCACCTCTAGGGAGA TACTAACACACCCTCCGAACAGA</i>	<i>test</i>	<i>Sanda</i>
<i>TMEM47</i>	<i>NM_031442.3:1215</i>	<i>AGCAAATAACCAACAGCCAATGTAGTC ATTGGGTAGGATAAGCAGGCCGT</i>	<i>CCCATTAGATGCTGAAGGGCAGTTCAT TTTCAAGGGCTCACTCA</i>	<i>test</i>	<i>Cooper – NGS</i>
<i>TMEM86A</i>	<i>NM_153347.1:2320</i>	<i>AATGAATCAGCCAATCTAATCCCATTG CTCCCAGCTGTTCAACTAAGCCC</i>	<i>GCTCCTGGAGCAGAGTGATGTATTATT CTGCCAGGGCTTTACAACATAATG</i>	<i>test</i>	<i>Cooper – NGS</i>
<i>TMPRSS2: ERG fusion</i>	<i>Fusion_0120.1:0</i>	<i>CTGCCGCGCTCCAGGCGGCGCTCCCC GCCCTCGC</i>	<i>TAGGCACACTCAAACAACGACTGGTCC TCACTCACAACCTGATAAGGCTTC</i>	<i>PCa positive control</i>	<i>Schalke n, Cooper</i>
<i>TRPM4</i>	<i>NM_001195227.1:2800</i>	<i>CTTCCAGTAGAGATCGCTGTTGCCCTG TACTTTGCCGAATGTGTAAGTGA</i>	<i>GCCAGCGCGGGCCGAGAGTGGAAATCC CGGATGAGGCGGTAACGCTGCGC</i>	<i>test</i>	<i>Sanda</i>
<i>TWIST1</i>	<i>NM_000474.3:393</i>	<i>CTCGGCGGCTGCTGCCGGTCTGGCTCT TCCTCGCTG</i>	<i>TGCTGCTGCGCCGCTTGGTCCCCCGC GCTTGCCG</i>	<i>test</i>	<i>Sanda</i>
<i>UPK2</i>	<i>NM_006760.3:332</i>	<i>ACGAGGTTTGTACCTGGTATGCACTG AGCCGAGTGAAGT</i>	<i>TCCCCTTCTTCACTAGGTAGGAAATGTA GAATTTGGTTCTGGC</i>	<i>Bladder control</i>	<i>Cooper</i>
<i>VAX2</i>	<i>NM_012476.2:871</i>	<i>TCACAGGGTGGGAGTCTTAAGTGTTAG CTTCTTGCAG</i>	<i>ACAGGAGACTGGGAAGGTGCTGTGCTC GGGACTCAGT</i>	<i>test</i>	<i>Sanda</i>
<i>VPS13A</i>	<i>NM_033305.2:8260</i>	<i>TAAAGGGCTTTGGTGCTGAATCCATGG TGACCGACTTTGGAGGTTAACA</i>	<i>ACGTGATATCTGGGAATGTCCTGCAGA TCTCATGACAATACTGACATCTG</i>	<i>Test</i>	<i>?</i>
<i>ZNF577</i>	<i>NM_032679.2:268</i>	<i>TCTCTCTCTGTCTATTCTGGGCCTTCC CAGAAGTGGTGGTCAG</i>	<i>GCCTTGCCCATTTCTGTTCAACTCTTAGG GGCTAGCAACTCTAGTATGTT</i>	<i>Test</i>	<i>Sanda</i>

## 9: APPENDICES

Supplementary Table 2 Samples flagged by quality checks on Nanostring2 data set

<i>Samples flagged by Quality Checks</i>			
<i>Samples detected by</i>	<i>M_83_7</i>	<i>M_26_6</i>	<i>pc145</i>
<i>NanoStringNorm where</i>	<i>M_84_2</i>	<i>M_27_1</i>	<i>pc1008_0</i>
<i>normalization parameters</i>	<i>M_84_5</i>	<i>M_31_1</i>	<i>pc017</i>
<i>extended beyond 100% from</i>	<i>M_84_6</i>	<i>M_31_3</i>	<i>a293</i>
<i>the mean</i>	<i>M_85_1</i>	<i>M_73_1</i>	<i>a316</i>
	<i>M_85_2</i>	<i>M_76_6</i>	<i>a303</i>
	<i>M_86_1</i>	<i>M_77_1</i>	<i>a1316</i>
	<i>M_86_2</i>	<i>M_77_2</i>	<i>a1319</i>
	<i>M_86_3</i>	<i>M_78_3</i>	<i>a1329</i>
	<i>M_133_7</i>	<i>M_78_5</i>	<i>a138a</i>
	<i>M_142_7</i>	<i>M_78_6</i>	<i>C113_1</i>
	<i>M_120_5</i>	<i>M_78_7</i>	<i>C118_4</i>
	<i>M_122_2</i>	<i>M_78_9</i>	<i>C110_1</i>
	<i>M_127_6</i>	<i>M_79_2</i>	<i>C112_4</i>
	<i>M_129_3</i>	<i>M_79_4</i>	<i>C107_2</i>
	<i>M_131_4</i>	<i>M_81_1</i>	<i>C111_1</i>
	<i>M_75_3</i>	<i>M_81_2</i>	<i>C109_4</i>
	<i>M_42_7</i>	<i>M_81_4</i>	<i>C107_1</i>
	<i>M_80_3</i>	<i>M_81_5</i>	<i>C118_3</i>
	<i>M_73_7</i>	<i>M_68_8</i>	<i>C106_8</i>
	<i>M_129_5</i>	<i>M_92_5</i>	<i>C116_5</i>
	<i>M_131_8</i>	<i>M_54_7</i>	<i>C116_2</i>
	<i>M_132_2</i>	<i>M_58_5</i>	
	<i>M_132_5</i>	<i>M_67_5</i>	
<i>Samples detected by</i>	<i>M_91-6</i>	<i>pc135</i>	
<i>NanoStringQCPro which</i>	<i>M_97-3</i>	<i>pc137</i>	
<i>were found to have</i>	<i>M_97-4</i>	<i>pc139</i>	
<i>overlapping barcodes</i>	<i>M_142-7</i>	<i>pc145</i>	
	<i>M_120-5</i>	<i>pc146</i>	
	<i>M_122-2</i>	<i>pc1008</i>	
	<i>M_127-6</i>	<i>pc1013</i>	
	<i>M_129-3</i>	<i>pc1029</i>	
	<i>M_131-4</i>	<i>pc1043</i>	
	<i>pc105-5</i>	<i>pc118</i>	
	<i>pc130-2</i>	<i>pc119</i>	
	<i>pc140-20</i>	<i>pc121</i>	
	<i>pc140-5</i>	<i>a1316</i>	
	<i>pc140-2</i>	<i>a1331</i>	



9: APPENDICES

Supplementary Table 3 Differentially expressed probes for each LPD group determined in the Nanostring2 data set. A) LPD groups 1-3. B) LPD groups 4-5.

A

LPD Group 1			LPD Group 2			LPD Group 3		
	Adjusted $p$ -value	Log2(FC)		Adjusted $p$ -value	Log2(FC)		Adjusted $p$ -value	Log2(FC)
<i>CAMKK2</i>	4.80X10-14	-1.17	<i>IFT57</i>	9.50X10-14	0.18	<i>KLK2</i>	1.97X10-12	-0.20
<i>CACNA1D</i>	1.10X10-13	-0.49	<i>OGT</i>	1.26X10-13	0.27	<i>DPP4</i>	2.20X10-12	-0.23
<i>GABARAPL2</i>	1.10X10-13	-0.32	<i>GABARAPL2</i>	1.31X10-13	0.17	<i>CASKIN1</i>	1.29X10-10	-0.21
<i>RPS11</i>	3.13X10-13	-0.13	<i>DPP4</i>	1.56X10-13	0.19	<i>MSMB</i>	1.34X10-10	-0.08
<i>RPL23AP53</i>	3.74X10-13	-1.35	<i>IMPDH2</i>	1.65X10-13	0.26	<i>CACNA1D</i>	1.55X10-10	-0.20
<i>PPAP2A</i>	3.94X10-13	-0.33	<i>HPRT</i>	1.68X10-13	0.30	<i>GABARAPL2</i>	1.71X10-10	-0.14
<i>CTA.211A9.5.MIATN</i>	4.44X10-13	-2.43	<i>EIF2D</i>	1.69X10-13	0.25	<i>TERT</i>	2.02X10-10	-0.24
<b>B</b>								
<i>STEAP2</i>	5.07X10-13	-0.60	<i>MXI1</i>	2.05X10-13	0.22	<i>ZNF577</i>	2.69X10-10	-0.26
<i>IFT57</i>	8.73X10-13	-0.33	<i>PECI</i>	2.09X10-13	0.25	<i>SSPO</i>	3.12X10-10	-0.20
<i>MIC1</i>	8.77X10-13	-1.20	<i>RP11.97012.7</i>	2.10X10-13	0.28	<i>CAMK2N2</i>	3.32X10-10	-0.52
<i>CASKIN1</i>	1.05X10-12	-0.41	<i>CACNA1D</i>	2.14X10-13	0.22	<i>IFT57</i>	5.68X10-10	-0.15
<i>HMBS</i>	1.25X10-12	-0.84	<i>FDPS</i>	3.25X10-13	0.19	<i>FOLH1</i>	5.86X10-10	-0.32
<i>MED4</i>	1.25X10-12	-0.83	<i>MYOF</i>	4.00X10-13	0.26	<i>MNX1</i>	6.38X10-10	-0.22
<i>RPLP2</i>	1.32X10-12	-0.16	<i>BRAF</i>	4.31X10-13	0.30	<i>MXI1</i>	1.56X10-09	-0.16
<i>HIST1H1C</i>	1.72X10-12	-0.31	<i>ALAS1</i>	4.49X10-13	0.21	<i>RP11.244H18.1.P71</i>	3.02X10-09	-0.28
<b>2P</b>								
<i>PCSK6</i>	1.97X10-12	-0.47	<i>MARCH5</i>	4.73X10-13	0.24	<i>STEAP2</i>	4.33X10-09	-0.24
<i>MMP11</i>	1.98X10-12	-0.66	<i>KLK2</i>	4.87X10-13	0.20	<i>TWIST1</i>	4.66X10-09	-0.21
<i>Timp4</i>	2.11X10-12	-1.43	<i>RIOK3</i>	6.19X10-13	0.35	<i>RPLP2</i>	8.02X10-09	-0.08
<i>SIM2.short</i>	2.17X10-12	-0.76	<i>ZNF577</i>	6.44X10-13	0.27	<i>HPRT</i>	8.93X10-09	-0.19
<i>SLC43A1</i>	2.17X10-12	-2.98	<i>HIST1H1C</i>	6.62X10-13	0.17	<i>MED4</i>	9.74X10-09	-0.21
<i>SSTR1</i>	2.28X10-12	-0.45	<i>MED4</i>	6.69X10-13	0.29	<i>RPS11</i>	2.70X10-08	-0.05
<i>SYNM</i>	2.48X10-12	-1.69	<i>TWIST1</i>	6.69X10-13	0.32	<i>PCSK6</i>	3.33X10-08	-0.27
<i>RPS10</i>	2.52X10-12	-0.18	<i>TFDP1</i>	8.58X10-13	0.28	<i>MMP26</i>	4.93X10-08	-0.90
<i>MEX3A</i>	2.59X10-12	-1.06	<i>STEAP2</i>	9.68X10-13	0.27	<i>PSGR</i>	7.51X10-08	-0.33
<i>HIST1H2BG</i>	2.78X10-12	-1.35	<i>GAPDH</i>	1.04X10-12	0.18	<i>NKAIN1</i>	1.44X10-07	-0.31
<i>VAX2</i>	3.63X10-12	-0.55	<i>LBH</i>	1.23X10-12	0.46	<i>ARexons4.8</i>	1.46X10-07	-0.19
<i>HOXC4</i>	3.91X10-12	-1.20	<i>SSPO</i>	1.33X10-12	0.38	<i>SSTR1</i>	1.75X10-07	-0.18

9: APPENDICES

<i>RPL18A</i>	4.06X10 <sup>-12</sup>	-0.33	<i>TERF2IP</i>	1.33X10 <sup>-12</sup>	0.18	<i>EN2</i>	1.76X10 <sup>-07</sup>	-0.31
<i>PALM3</i>	4.82X10 <sup>-12</sup>	-0.82	<i>HIST1H2BF</i>	1.50X10 <sup>-12</sup>	0.22	<i>PPAP2A</i>	1.97X10 <sup>-07</sup>	-0.17
<i>FDPS</i>	5.40X10 <sup>-12</sup>	-0.35	<i>MEMO1</i>	1.67X10 <sup>-12</sup>	0.28	<i>CDC20</i>	2.70X10 <sup>-07</sup>	-0.34
<i>TWIST1</i>	5.54X10 <sup>-12</sup>	-0.51	<i>NAALADL2</i>	1.67X10 <sup>-12</sup>	0.30	<i>MGAT5B</i>	3.52X10 <sup>-07</sup>	-0.25
<i>EN2</i>	7.01X10 <sup>-12</sup>	-0.56	<i>RPL18A</i>	1.67X10 <sup>-12</sup>	0.15	<i>TFDP1</i>	3.65X10 <sup>-07</sup>	-0.19
<i>MXI1</i>	1.01X10 <sup>-11</sup>	-0.36	<i>CASKIN1</i>	1.86X10 <sup>-12</sup>	0.33	<i>SLC4A1.S</i>	7.77X10 <sup>-07</sup>	-0.63
<i>STEAP4</i>	1.11X10 <sup>-11</sup>	-0.64	<i>ITPR1</i>	1.86X10 <sup>-12</sup>	0.28	<i>KLK3 exons 2-3</i>	7.99X10 <sup>-07</sup>	-0.25
<i>ISX</i>	1.17X10 <sup>-11</sup>	-1.67	<i>RPLP2</i>	2.27X10 <sup>-12</sup>	0.10	<i>EIF2D</i>	8.19X10 <sup>-07</sup>	-0.19
<i>KLK3 exons 2-3</i>	1.26X10 <sup>-11</sup>	-0.69	<i>HMBS</i>	2.34X10 <sup>-12</sup>	0.34	<i>KLK4</i>	1.21X10 <sup>-06</sup>	-0.14
<i>TMEM86A</i>	1.36X10 <sup>-11</sup>	-1.51	<i>PPAP2A</i>	2.34X10 <sup>-12</sup>	0.18	<i>ARHGEF25</i>	1.36X10 <sup>-06</sup>	-0.60
<i>KLK4</i>	1.49X10 <sup>-11</sup>	-0.36	<i>SLC4A1.S</i>	2.46X10 <sup>-12</sup>	0.79	<i>COL9A2</i>	1.47X10 <sup>-06</sup>	-0.79
<i>MXN1</i>	1.54X10 <sup>-11</sup>	-0.43	<i>MMP11</i>	2.71X10 <sup>-12</sup>	0.44	<i>HIST1H2BF</i>	2.43X10 <sup>-06</sup>	-0.21
<i>AMH</i>	1.82X10 <sup>-11</sup>	-0.59	<i>SIM2.short</i>	2.81X10 <sup>-12</sup>	0.45	<i>SNCA</i>	4.94X10 <sup>-06</sup>	-0.28
<i>AR exons 4-8</i>	2.58X10 <sup>-11</sup>	-0.67	<i>Ar exons 4-8</i>	2.85X10 <sup>-12</sup>	0.27	<i>MIR146A.DQ65841 4</i>	5.41X10 <sup>-06</sup>	-0.26
<i>SMIM1</i>	2.64X10 <sup>-11</sup>	-1.17	<i>RP11.244H18.1.P7 12P</i>	3.07X10 <sup>-12</sup>	0.25	<i>COL10A1</i>	5.98X10 <sup>-06</sup>	-0.30
<i>RIOK3</i>	3.77X10 <sup>-11</sup>	-0.83	<i>MXN1</i>	3.52X10 <sup>-12</sup>	0.42	<i>CADPS</i>	6.01X10 <sup>-06</sup>	-0.59
<i>BRAF</i>	3.82X10 <sup>-11</sup>	-0.46	<i>HIST1H2BG</i>	3.66X10 <sup>-12</sup>	0.25	<i>VAX2</i>	6.99X10 <sup>-06</sup>	-0.21
<i>SChLAP1</i>	4.04X10 <sup>-11</sup>	-1.22	<i>STEAP4</i>	3.67X10 <sup>-12</sup>	0.33	<i>HIST1H1C</i>	7.08X10 <sup>-06</sup>	-0.11
<i>DLX1</i>	4.75X10 <sup>-11</sup>	-1.94	<i>AMH</i>	4.12X10 <sup>-12</sup>	0.41	<i>PTN</i>	7.92X10 <sup>-06</sup>	-0.39
<i>PVT1</i>	4.81X10 <sup>-11</sup>	-0.68	<i>SMAP1 exons 7-8</i>	4.43X10 <sup>-12</sup>	0.36	<i>PSTPIP1</i>	8.67X10 <sup>-06</sup>	-0.57
<i>IMPDH2</i>	4.92X10 <sup>-11</sup>	-0.59	<i>FOLH1</i>	5.02X10 <sup>-12</sup>	0.27	<i>MARCH5</i>	9.65X10 <sup>-06</sup>	-0.13
<i>MGAT5B</i>	5.09X10 <sup>-11</sup>	-0.68	<i>RPS10</i>	5.90X10 <sup>-12</sup>	0.13	<i>SIM2.short</i>	9.74X10 <sup>-06</sup>	-0.26
<i>CD10</i>	5.79X10 <sup>-11</sup>	-0.54	<i>GCNT1</i>	7.08X10 <sup>-12</sup>	0.39	<i>AMH</i>	1.07X10 <sup>-05</sup>	-0.21
<i>HPRT</i>	6.13X10 <sup>-11</sup>	-0.66	<i>BTG2</i>	7.11X10 <sup>-12</sup>	0.31	<i>MMP11</i>	1.57X10 <sup>-05</sup>	-0.20
<i>ITGBL1</i>	6.13X10 <sup>-11</sup>	-0.86	<i>CCDC88B</i>	7.11X10 <sup>-12</sup>	0.78	<i>RPS10</i>	1.62X10 <sup>-05</sup>	-0.07
<i>EIF2D</i>	7.83X10 <sup>-11</sup>	-0.46	<i>CDC37L1</i>	8.55X10 <sup>-12</sup>	0.39	<i>STEAP4</i>	1.62X10 <sup>-05</sup>	-0.13
<i>DPP4</i>	8.53X10 <sup>-11</sup>	-0.44	<i>TMCC2</i>	8.55X10 <sup>-12</sup>	0.65	<i>IMPDH2</i>	1.67X10 <sup>-05</sup>	-0.16
<i>TFDP1</i>	9.69X10 <sup>-11</sup>	-0.88	<i>DNAH5</i>	9.37X10 <sup>-12</sup>	0.67	<i>SPINK1</i>	2.08X10 <sup>-05</sup>	0.29
<i>TERF2IP</i>	1.01X10 <sup>-10</sup>	-0.31	<i>NLRP3</i>	9.81X10 <sup>-12</sup>	0.80	<i>ISX</i>	3.99X10 <sup>-05</sup>	-0.55
<i>ZNF577</i>	1.22X10 <sup>-10</sup>	-0.84	<i>RPS11</i>	1.05X10 <sup>-11</sup>	0.08	<i>ACTR5</i>	4.01X10 <sup>-05</sup>	-0.45
<i>ARHGEF25</i>	1.39X10 <sup>-10</sup>	-1.39	<i>PDLIM5</i>	1.13X10 <sup>-11</sup>	0.24	<i>MMP25</i>	5.23X10 <sup>-05</sup>	-0.58
<i>SIM2.long</i>	1.42X10 <sup>-10</sup>	-1.29	<i>SSTR1</i>	1.72X10 <sup>-11</sup>	0.42	<i>CAMKK2</i>	5.77X10 <sup>-05</sup>	-0.39
<i>RP11.97O12.7</i>	1.48X10 <sup>-10</sup>	-0.97	<i>SRSF3</i>	2.05X10 <sup>-11</sup>	0.51	<i>MYOF</i>	6.98X10 <sup>-05</sup>	-0.15
<i>SRSF3</i>	1.51X10 <sup>-10</sup>	-1.16	<i>ABCB9</i>	2.15X10 <sup>-11</sup>	0.77	<i>NAALADL2</i>	7.35X10 <sup>-05</sup>	-0.22
<i>SFRP4</i>	1.55X10 <sup>-10</sup>	-0.89	<i>PPP1R12B</i>	2.45X10 <sup>-11</sup>	0.33	<i>TERF2IP</i>	8.59X10 <sup>-05</sup>	-0.09

9: APPENDICES

<i>PTN</i>	1.58X10-10	-0.71	<i>B2M</i>	2.51X10-11	0.27	<i>CD10</i>	8.75X10-05	-0.15
<i>COL10A1</i>	1.95X10-10	-0.59	<i>MDK</i>	3.09X10-11	0.27	<i>SChLAP1</i>	9.41X10-05	-0.49
<i>GOLM1</i>	1.95X10-10	-2.39	<i>MGAT5B</i>	3.56X10-11	0.33	<i>MEMO1</i>	0.000124295	-0.15
<i>PECI</i>	2.39X10-10	-0.38	<i>CD10</i>	4.43X10-11	0.25	<i>GCNT1</i>	0.000127837	-0.35
<i>SSPO</i>	2.47X10-10	-0.38	<i>KLK4</i>	4.53X10-11	0.16	<i>KLK3 exons 1-2</i>	0.000138743	-0.22
<i>MSMB</i>	2.62X10-10	-0.15	<i>VPS13A</i>	6.07X10-11	0.37	<i>MAPK8IP2</i>	0.00014341	-0.57
<i>SNCA</i>	2.88X10-10	-0.88	<i>HIST1H1E</i>	6.14X10-11	0.32	<i>DNAH5</i>	0.00014901	-0.39
<i>TBP</i>	2.91X10-10	-2.06	<i>SEC61A1</i>	6.14X10-11	0.61	<i>ST6GALNAC1</i>	0.000152737	-1.26
<i>PPP1R12B</i>	3.16X10-10	-0.95	<i>TERT</i>	6.15X10-11	0.31	<i>CTA.211A9.5.MIAT NB</i>	0.000157375	-0.65
<i>ANKRD34B</i>	4.26X10-10	-1.57	<i>PSGR</i>	8.66X10-11	0.29	<i>SLC43A1</i>	0.00017283	-0.98
<i>TMPRSS2:ERG</i>	4.88X10-10	-1.95	<i>SMIM1</i>	9.84X10-11	0.26	<i>NLRP3</i>	0.000178965	-0.46
<i>SACM1L</i>	5.03X10-10	-0.62	<i>COL9A2</i>	1.20X10-10	0.91	<i>SRSF3</i>	0.000264469	-0.36
<i>CADPS</i>	5.12X10-10	-1.63	<i>AMACR</i>	1.23X10-10	0.32	<i>ANKRD34B</i>	0.000277869	-0.61
<i>KLK2</i>	5.12X10-10	-0.46	<i>CKAP2L</i>	1.41X10-10	0.82	<i>FDPS</i>	0.00030588	-0.11
<i>COL9A2</i>	5.17X10-10	-1.60	<i>PSTPIP1</i>	1.52X10-10	0.78	<i>SIM2.long</i>	0.000332184	-0.43
<i>TMEM47</i>	5.25X10-10	-2.48	<i>PTN</i>	1.52X10-10	0.36	<i>UPK2</i>	0.000344195	0.53
<i>MARCH5</i>	5.48X10-10	-0.58	<i>VAX2</i>	1.52X10-10	0.43	<i>PDLIM5</i>	0.000389202	-0.12
<i>PPFIA2</i>	5.70X10-10	-1.30	<i>MMP25</i>	1.67X10-10	0.78	<i>ERG5</i>	0.000408807	-0.47
<i>LASS1</i>	6.33X10-10	-0.69	<i>SACM1L</i>	1.91X10-10	0.35	<i>OGT</i>	0.000408807	-0.13
<i>PDLIM5</i>	8.66X10-10	-0.44	<i>CLIC2</i>	1.93X10-10	0.83	<i>PVT1</i>	0.000436352	-0.33
<i>FOLH1</i>	8.98X10-10	-0.62	<i>GOLM1</i>	1.99X10-10	0.43	<i>TMEM47</i>	0.000458492	-1.08
<i>SNORA20</i>	8.98X10-10	-2.62	<i>PCSK6</i>	2.05X10-10	0.38	<i>RPL18A</i>	0.000499086	-0.07
<i>GCNT1</i>	9.05X10-10	-1.97	<i>CDC20</i>	2.19X10-10	0.73	<i>PCA3</i>	0.00052407	-0.19
<i>CLIC2</i>	9.17X10-10	-1.42	<i>KLK3 exons 2-3</i>	2.19X10-10	0.19	<i>SNORA20</i>	0.000536406	-1.10
<i>GAPDH</i>	9.19X10-10	-0.23	<i>ACTR5</i>	3.15X10-10	0.38	<i>LBH</i>	0.000583164	-0.34
<i>HOXC6</i>	1.14X10-09	-1.38	<i>SNCA</i>	3.25X10-10	0.23	<i>HIST1H2BG</i>	0.000703794	-0.15
<i>RP11.244H18.1.P712 P</i>	1.14X10-09	-0.52	<i>HPN</i>	3.31X10-10	0.38	<i>MEX3A</i>	0.000862201	-0.30
<i>SMAP1 exons 7-8</i>	1.44X10-09	-1.35	<i>MEX3A</i>	3.33X10-10	0.68	<i>PECI</i>	0.000862201	-0.12
<i>TMEM45B</i>	1.51X10-09	-0.53	<i>RPL23AP53</i>	4.01X10-10	0.53	<i>Timp4</i>	0.000976934	-0.68
<i>MDK</i>	1.56X10-09	-1.02	<i>TMEM45B</i>	4.21X10-10	0.44	<i>TMEM86A</i>	0.001450492	-0.56
<i>TERT</i>	2.17X10-09	-0.36	<i>CAMKK2</i>	4.60X10-10	0.33	<i>CKAP2L</i>	0.001943582	-0.42
<i>CDC20</i>	2.41X10-09	-0.75	<i>SYNM</i>	4.93X10-10	0.54	<i>RNF157</i>	0.002251894	-0.52
<i>MMP26</i>	2.47X10-09	-1.45	<i>SIRT1</i>	5.07X10-10	0.54	<i>RIOK3</i>	0.002581105	-0.15
<i>AGR2</i>	2.87X10-09	-1.52	<i>MFS2A</i>	5.26X10-10	0.78	<i>AGR2</i>	0.002794445	-0.32
<i>OR52A2.PSGR</i>	2.97X10-09	-1.37	<i>COL10A1</i>	6.30X10-10	0.44	<i>CCDC88B</i>	0.003237426	-0.54

9: APPENDICES

<i>MYOF</i>	3.29X10-09	-0.36	<i>MCM7</i>	6.67X10-10	0.47	<i>GAPDH</i>	0.003237426	-0.07
<i>SLC4A1.S</i>	3.42X10-09	-2.05	<i>SIM2.long</i>	6.92X10-10	0.50	<i>SMAP1 exons 7-8</i>	0.003237426	-0.26
<i>KLK3 exons 1-2</i>	3.66X10-09	-0.93	<i>NKAIN1</i>	7.55X10-10	0.25	<i>SACM1L</i>	0.003611413	-0.15
<i>AATF</i>	3.68X10-09	-1.94	<i>MIR146A.DQ6584</i> <i>14</i>	8.31X10-10	0.44	<i>LASS1</i>	0.003758809	-0.32
<i>NLRP3</i>	3.68X10-09	-1.21	<i>KLK3 exons 1-2</i>	9.53X10-10	0.21	<i>MCM7</i>	0.003892978	-0.36
<i>SPON2</i>	3.93X10-09	-0.79	<i>EN2</i>	1.00X10-09	0.34	<i>GOLM1</i>	0.004104272	-0.29
<i>ACTR5</i>	4.55X10-09	-0.87	<i>GJB1</i>	1.00X10-09	0.26	<i>MIC1</i>	0.004950076	-0.36
<i>HPN</i>	5.48X10-09	-0.99	<i>TMEM86A</i>	1.11X10-09	1.06	<i>SFRP4</i>	0.005187855	-0.43
<i>MAK</i>	5.71X10-09	-2.15	<i>ANKRD34B</i>	1.13X10-09	0.91	<i>SMIM1</i>	0.005244137	-0.16
<i>B4GALNT4</i>	5.82X10-09	-1.13	<i>TBP</i>	1.34X10-09	0.42	<i>MAK</i>	0.005275769	-0.75
<i>CDKN3</i>	5.82X10-09	-0.94	<i>AGR2</i>	1.39X10-09	0.45	<i>IGFBP3</i>	0.005353681	0.36
<i>GJB1</i>	5.82X10-09	-1.53	<i>MMP26</i>	1.39X10-09	1.02	<i>SYNM</i>	0.005479266	-0.45
<i>TRPM4</i>	5.82X10-09	-1.59	<i>PVT1</i>	1.83X10-09	0.45	<i>TMCC2</i>	0.006177628	-0.31
<i>ITPR1</i>	5.88X10-09	-0.64	<i>ANPEP</i>	1.91X10-09	0.38	<i>DLX1</i>	0.006513656	-0.51
<i>CKAP2L</i>	6.54X10-09	-0.97	<i>Timp4</i>	2.60X10-09	0.40	<i>SERPINB5</i>	0.007855423	-0.32
<i>NAALADL2</i>	9.10X10-09	-0.81	<i>RAB17</i>	2.66X10-09	0.45	<i>TRPM4</i>	0.009529009	-0.41
<i>AMACR</i>	9.36X10-09	-1.31	<i>AATF</i>	2.69X10-09	0.46	<i>RP11.97012.7</i>	0.010060703	-0.14
<i>HIST3H2A</i>	1.00X10-08	-1.05	<i>CTA.211A9.5.MIAT</i> <i>NB</i>	2.91X10-09	0.36	<i>CLU</i>	0.010087038	-0.45
<i>MIR146A.DQ658414</i>	1.04X10-08	-0.74	<i>SERPINB5</i>	3.68X10-09	0.47	<i>SLC12A1</i>	0.010326711	0.41
<i>SEC61A1</i>	1.15X10-08	-1.99	<i>NEAT1</i>	3.68X10-09	0.41	<i>MDK</i>	0.010838865	-0.16
<i>ERG3 exons 4-5</i>	1.50X10-08	-1.34	<i>B4GALNT4</i>	4.93X10-09	0.87	<i>CDKN3</i>	0.011131753	-0.39
<i>MCM7</i>	1.50X10-08	-1.48	<i>CADPS</i>	7.75X10-09	1.22	<i>BRAF</i>	0.011920795	-0.11
<i>HIST1H2BF</i>	1.69X10-08	-0.53	<i>SchLAP1</i>	8.52X10-09	0.63	<i>HMBS</i>	0.011920795	-0.22
<i>SIRT1</i>	1.85X10-08	-1.41	<i>MIC1</i>	8.65X10-09	0.58	<i>MIR4435.1HG.IOC5</i> <i>41471</i>	0.013229338	0.15
<i>ERG3 exons 6-7</i>	1.93X10-08	-1.96	<i>PALM3</i>	8.65X10-09	0.46	<i>RPL23AP53</i>	0.013779284	-0.23
<i>CDC37L1</i>	2.09X10-08	-0.71	<i>HIST3H2A</i>	9.40X10-09	0.28	<i>AATF</i>	0.023593676	-0.23
<i>ARexon9</i>	2.38X10-08	-1.50	<i>ARHGEF25</i>	1.09X10-08	0.73	<i>B4GALNT4</i>	0.023593676	-0.28
<i>HIST1H1E</i>	2.50X10-08	-0.58	<i>MSMB</i>	1.31X10-08	0.06	<i>HIST3H2A</i>	0.023593676	-0.22
<i>STOM</i>	2.86X10-08	-2.23	<i>ISX</i>	2.15X10-08	0.95	<i>MKi67</i>	0.023593676	-0.61
<i>SERPINB5</i>	3.66X10-08	-1.39	<i>TRPM4</i>	2.75X10-08	0.54	<i>B2M</i>	0.024443101	-0.10
<i>LBH</i>	3.98X10-08	-1.22	<i>CLU</i>	2.93X10-08	0.84	<i>AR exon 9</i>	0.027370444	-0.57
<i>TMCC2</i>	4.17X10-08	-1.84	<i>HOXC6</i>	2.93X10-08	0.32	<i>CLIC2</i>	0.029401344	-0.28
<i>MMP25</i>	4.69X10-08	-1.22	<i>RNF157</i>	2.93X10-08	0.51	<i>SPON2</i>	0.029401344	-0.10
<i>ERG5</i>	5.45X10-08	-1.67	<i>ST6GALNAC1</i>	2.93X10-08	0.65	<i>ABC9</i>	0.047180253	-0.18

9: APPENDICES

<b>CAMK2N2</b>	7.63X10-08	-0.79	<b>AR exon 9</b>	3.18X10-08	1.02
<b>VPS13A</b>	8.47X10-08	-0.88	<b>HOXC4</b>	4.42X10-08	0.63
<b>MFS2A</b>	8.52X10-08	-1.44	<b>ERG5</b>	5.96X10-08	0.79
<b>ST6GALNAC1</b>	9.99X10-08	-2.46	<b>SLC43A1</b>	7.19X10-08	0.70
<b>AURKA</b>	1.35X10-07	-2.28	<b>MAPK8IP2</b>	9.88X10-08	0.93
<b>CLU</b>	2.46X10-07	-0.87	<b>PCA3</b>	9.92X10-08	0.25
<b>MEMO1</b>	3.02X10-07	-0.58	<b>CAMK2N2</b>	1.00X10-07	0.43
<b>PCA3</b>	3.02X10-07	-0.31	<b>MKI67</b>	1.33X10-07	1.01
<b>ALAS1</b>	3.42X10-07	-0.26	<b>CP</b>	2.10X10-07	0.75
<b>CCDC88B</b>	3.75X10-07	-1.09	<b>ITGBL1</b>	2.10X10-07	0.45
<b>OGT</b>	3.75X10-07	-0.38	<b>SNORA20</b>	2.84X10-07	0.56
<b>SULF2</b>	3.75X10-07	-1.39	<b>PPFIA2</b>	2.90X10-07	0.60
<b>Met</b>	3.77X10-07	NA	<b>TDRD</b>	3.25X10-07	1.05
<b>TDRD</b>	4.18X10-07	-1.95	<b>STOM</b>	3.73X10-07	0.94
<b>B2M</b>	9.21X10-07	-0.38	<b>CDKN3</b>	4.26X10-07	0.68
<b>RNF157</b>	9.31X10-07	-1.40	<b>SULF2</b>	6.30X10-07	0.66
<b>RAB17</b>	9.77X10-07	-0.67	<b>AURKA</b>	6.74X10-07	0.49
<b>SULT1A1</b>	2.23X10-06	-1.48	<b>SFRP4</b>	6.74X10-07	0.73
<b>MAPK8IP2</b>	3.59X10-06	-0.96	<b>MIR4435.1HG.IOC 541471</b>	2.25X10-06	0.26
<b>DNAH5</b>	3.84X10-06	-1.08	<b>LASS1</b>	2.34X10-06	0.64
<b>MCTP1</b>	4.54X10-06	-1.97	<b>SLC12A1</b>	1.18X10-05	0.56
<b>MKI67</b>	5.94X10-06	-3.13	<b>SPON2</b>	1.18X10-05	0.18
<b>BTG2</b>	6.71X10-06	-1.07	<b>TMEM47</b>	1.26X10-05	0.75
<b>NKAIN1</b>	6.71X10-06	-1.27	<b>ERG3 exons 4-5</b>	2.02X10-05	0.73
<b>CP</b>	6.85X10-06	-1.89	<b>ERG3 exons 6-7</b>	2.02X10-05	0.92
<b>ANPEP</b>	7.42X10-06	-1.24	<b>MAK</b>	2.57X10-05	1.04
<b>PSTPIP1</b>	7.42X10-06	-0.98	<b>DLX1</b>	2.70X10-05	1.20
<b>NEAT1</b>	2.28X10-05	-1.31	<b>SULT1A1</b>	6.39X10-05	0.23
<b>ABC9</b>	0.000117259	-0.68	<b>MCTP1</b>	6.98X10-05	0.48
<b>APOC1</b>	0.000171548	-1.63	<b>Met</b>	9.24X10-05	0.91
<b>SLC12A1</b>	0.000171548	-1.10	<b>PTPRC</b>	9.24X10-05	0.90
<b>MIR4435.1HG.IOC54 1471</b>	0.000399737	-1.01	<b>TMPRSS2:ERG</b>	9.24X10-05	1.19
<b>UPK2</b>	0.001159053	-1.38	<b>APOC1</b>	0.001191952	0.41
<b>PTPRC</b>	0.001376488	-1.44	<b>SPINK1</b>	0.004238366	0.17
<b>IGFBP3</b>	0.009734383	-1.56	<b>IGFBP3</b>	0.005829551	0.41

## 9: APPENDICES

---

<i>SPINK1</i>	0.039101303	-0.22
---------------	-------------	-------

---

## 9: APPENDICES

B

LPD Group 4			LPD Group 5		
	Adjusted <i>p</i> -value	Log2(FC)		Adjusted <i>p</i> -value	Log2(FC)
<i>VPS13A</i>	3.38X10-06	-0.11	<i>GABARAPL2</i>	2.26X10-22	0.07
<i>TERF2IP</i>	3.79X10-06	-0.05	<i>CACNA1D</i>	2.71X10-21	0.09
<i>ABCB9</i>	1.47X10-05	-0.21	<i>STEAP2</i>	3.26X10-17	0.09
<i>X05.Mar</i>	1.64X10-05	-0.08	<i>KLK2</i>	4.09X10-17	0.07
<i>MMP25</i>	1.89X10-05	-0.25	<i>MED4</i>	2.31X10-16	0.09
<i>TMEM45B</i>	1.92X10-05	-0.14	<i>CASKIN1</i>	1.66X10-15	0.13
<i>RPLP2</i>	1.93X10-05	-0.03	<i>DPP4</i>	7.40X10-15	0.07
<i>PECI</i>	2.41X10-05	-0.06	<i>IFT57</i>	8.66X10-15	0.07
<i>CASKIN1</i>	2.64X10-05	-0.10	<i>RPS11</i>	8.75X10-14	0.03
<i>MEMO1</i>	3.23X10-05	-0.08	<i>MARCH5</i>	9.67X10-14	0.09
<i>AMACR</i>	4.96X10-05	-0.10	<i>MMP25</i>	1.56X10-13	0.33
<i>SLC4A1.S</i>	5.73X10-05	-0.24	<i>STEAP4</i>	2.42X10-13	0.09
<i>GABARAPL2</i>	6.12X10-05	-0.04	<i>TWIST1</i>	2.71X10-13	0.12
<i>TWIST1</i>	6.76X10-05	-0.08	<i>MMP26</i>	4.57X10-13	0.47
<i>FDPS</i>	8.19X10-05	-0.04	<i>SYNM</i>	5.46X10-13	0.31
<i>CACNA1D</i>	0.00010192	-0.06	<i>TERF2IP</i>	6.23X10-13	0.05
<i>CP</i>	0.000320555	-0.30	<i>FDPS</i>	6.52X10-13	0.05
<i>RPS11</i>	0.000453771	-0.02	<i>PCSK6</i>	8.32X10-13	0.13
<i>TMEM86A</i>	0.000537454	-0.27	<i>SSTR1</i>	1.15X10-12	0.12
<i>PPP1R12B</i>	0.00056124	-0.08	<i>MNX1</i>	1.79X10-12	0.15
<i>TDRD</i>	0.000623636	-0.71	<i>HPRT</i>	2.47X10-12	0.10
<i>TMCC2</i>	0.001057497	-0.18	<i>FOLH1</i>	2.60X10-12	0.11
<i>BTG2</i>	0.001079391	-0.07	<i>CDC20</i>	2.84X10-12	0.21
<i>BRAF</i>	0.001466681	-0.07	<i>SLC4A1.S</i>	3.09X10-12	0.32
<i>ITPR1</i>	0.00149178	-0.07	<i>COL10A1</i>	9.32X10-12	0.20
<i>MFSD2A</i>	0.001912749	-0.24	<i>TDRD</i>	1.16X10-11	0.89
<i>EN2</i>	0.001944902	-0.10	<i>TERT</i>	1.77X10-11	0.14
<i>SLC12A1</i>	0.001954045	-0.20	<i>EN2</i>	3.39X10-11	0.14
<i>CDC37L1</i>	0.002355197	-0.09	<i>ZNF577</i>	4.49X10-11	0.08
<i>PSTPIP1</i>	0.002736615	-0.18	<i>SSPO</i>	4.81X10-11	0.12
<i>COL10A1</i>	0.002804912	-0.14	<i>VAX2</i>	5.56X10-11	0.16
<i>ITGBL1</i>	0.00290803	-0.16	<i>MGAT5B</i>	8.97X10-11	0.13
<i>ALAS1</i>	0.003470784	-0.04	<i>RPL23AP53</i>	1.42X10-10	0.29
<i>DPP4</i>	0.003974	-0.04	<i>CAMK2N2</i>	1.56X10-10	0.26
<i>CCDC88B</i>	0.004126726	-0.22	<i>ERG5</i>	2.20X10-10	0.36
<i>SPINK1</i>	0.004287134	-0.15	<i>MXI1</i>	2.44X10-10	0.06
<i>HOXC4</i>	0.004817751	-0.20	<i>HIST1H2BG</i>	3.78X10-10	0.11
<i>IGFBP3</i>	0.006132992	-0.25	<i>PPAP2A</i>	4.44X10-10	0.06
<i>UPK2</i>	0.006309542	-0.30	<i>TMEM86A</i>	6.42X10-10	0.36
<i>RP11.97012.7</i>	0.00684069	-0.05	<i>MEMO1</i>	9.10X10-10	0.09
<i>GJB1</i>	0.007619494	-0.09	<i>RP11.244H18.1.P712P</i>	9.28X10-10	0.09
<i>SSTR1</i>	0.007727274	-0.08	<i>ARHGEF25</i>	1.20X10-09	0.39
<i>EIF2D</i>	0.008645702	-0.06	<i>RPLP2</i>	1.25X10-09	0.03
<i>MED4</i>	0.009155071	-0.05	<i>SFRP4</i>	1.34X10-09	0.34
<i>OGT</i>	0.010622928	-0.06	<i>HIST1H1C</i>	1.50X10-09	0.05
<i>MIR4435.1HG.IO</i>	0.011952383	-0.13	<i>COL9A2</i>	1.82X10-09	0.38
<i>C541471</i>					
<i>AMH</i>	0.012837851	-0.08	<i>PECI</i>	1.82X10-09	0.06
<i>MGAT5B</i>	0.012837851	-0.09	<i>BRAF</i>	2.21X10-09	0.09
<i>RIOK3</i>	0.013178355	-0.07	<i>CAMKK2</i>	3.30X10-09	0.16
<i>MXI1</i>	0.01326158	-0.05	<i>SIM2.short</i>	3.59X10-09	0.15
<i>PPAP2A</i>	0.01326158	-0.04	<i>SchLAP1</i>	3.59X10-09	0.36
<i>STEAP4</i>	0.01326158	-0.06	<i>RIOK3</i>	3.88X10-09	0.11
<i>PTPRC</i>	0.017081934	-0.45	<i>AMH</i>	4.07X10-09	0.13
<i>VAX2</i>	0.017081934	-0.09	<i>LBH</i>	4.60X10-09	0.20
<i>SSPO</i>	0.021717641	-0.06	<i>SACM1L</i>	4.74X10-09	0.11
<i>SACM1L</i>	0.022028919	-0.08	<i>PDLIM5</i>	8.54X10-09	0.08

## 9: APPENDICES

<b>SULF2</b>	0.02375717	-0.22	<b>ERG3 exons 4-5</b>	9.30X10-09	0.51
<b>TMPRSS2:ERG</b>	0.023967171	-0.54	<b>LASS1</b>	9.35X10-09	0.30
<b>CKAP2L</b>	0.029163389	-0.14	<b>GCNT1</b>	1.20X10-08	0.15
<b>KLK2</b>	0.035832883	-0.04	<b>MIR146A.DQ65 8414</b>	1.48X10-08	0.16
<b>HIST1H2BG</b>	0.041078353	-0.07	<b>MMP11</b>	1.90X10-08	0.12
			<b>MEX3A</b>	1.90X10-08	0.19
			<b>ANKRD34B</b>	2.35X10-08	0.27
			<b>EIF2D</b>	2.38X10-08	0.07
			<b>OGT</b>	2.54X10-08	0.07
			<b>PSTPIP1</b>	2.79X10-08	0.27
			<b>TFDP1</b>	3.91X10-08	0.07
			<b>DLX1</b>	4.64X10-08	0.65
			<b>B4GALNT4</b>	5.18X10-08	0.31
			<b>TMPRSS2:ERG</b>	5.95X10-08	0.93
			<b>MSMB</b>	7.42X10-08	0.03
			<b>CADPS</b>	1.01X10-07	0.37
			<b>CKAP2L</b>	1.23X10-07	0.24
			<b>SNORA20</b>	1.42X10-07	0.35
			<b>MAK</b>	1.67X10-07	0.72
			<b>ERG3 exons 6-7</b>	1.95X10-07	0.66
			<b>ITPR1</b>	3.32X10-07	0.09
			<b>RP11.97012.7</b>	4.08X10-07	0.06
			<b>AMACR</b>	5.11X10-07	0.12
			<b>ISX</b>	5.11X10-07	0.36
			<b>PVT1</b>	5.11X10-07	0.16
			<b>HOXC4</b>	6.68X10-07	0.23
			<b>AR exons 4-8</b>	6.96X10-07	0.07
			<b>MKI67</b>	7.96X10-07	0.44
			<b>TMEM47</b>	7.96X10-07	0.42
			<b>HMBS</b>	8.12X10-07	0.14
			<b>IMPDH2</b>	8.35X10-07	0.05
			<b>NAALADL2</b>	8.35X10-07	0.07
			<b>AR exon 9</b>	9.27X10-07	0.58
			<b>TMCC2</b>	9.91X10-07	0.22
			<b>HIST1H2BF</b>	1.56X10-06	0.07
			<b>SMAP1 exons 7- 8</b>	1.62X10-06	0.12
			<b>CLIC2</b>	1.66X10-06	0.34
			<b>TMEM45B</b>	2.10X10-06	0.12
			<b>KLK4</b>	2.14X10-06	0.06
			<b>ABCB9</b>	2.50X10-06	0.25
			<b>GJB1</b>	2.50X10-06	0.10
			<b>MYOF</b>	2.52X10-06	0.06
			<b>DNAH5</b>	2.53X10-06	0.28
			<b>CDC37L1</b>	3.63X10-06	0.09
			<b>PTN</b>	6.39X10-06	0.11
			<b>PPP1R12B</b>	6.95X10-06	0.08
			<b>NKAIN1</b>	8.02X10-06	0.10
			<b>PPFIA2</b>	1.11X10-05	0.28
			<b>SRSF3</b>	1.31X10-05	0.14
			<b>CP</b>	1.41X10-05	0.29
			<b>TBP</b>	1.64X10-05	0.10
			<b>RNF157</b>	1.79X10-05	0.22
			<b>CCDC88B</b>	1.99X10-05	0.29
			<b>NLRP3</b>	2.01X10-05	0.21
			<b>ACTR5</b>	2.28X10-05	0.13
			<b>GOLM1</b>	2.46X10-05	0.14
			<b>VPS13A</b>	2.47X10-05	0.08
			<b>CD10</b>	2.99X10-05	0.06
			<b>MAPK8IP2</b>	2.99X10-05	0.29
			<b>PCA3</b>	2.99X10-05	0.11
			<b>CTA.211A9.5.MI</b>	4.66X10-05	0.17



9: APPENDICES

<b>ATNB</b>			
<b>SLC43A1</b>	4.66X10-05		0.27
<b>RPL18A</b>	5.11X10-05		0.04
<b>SNCA</b>	5.11X10-05		0.08
<b>MIC1</b>	8.46X10-05		0.19
<b>Timp4</b>	9.62X10-05		0.15
<b>MFSD2A</b>	0.000117774		0.28
<b>RPS10</b>	0.000141176		0.04
<b>CLU</b>	0.000315682		0.32
<b>MCM7</b>	0.000317762		0.16
<b>SMIM1</b>	0.000387115		0.09
<b>SEC61A1</b>	0.000429125		0.20
<b>CDKN3</b>	0.000442848		0.19
<b>PALM3</b>	0.000442848		0.13
<b>SIRT1</b>	0.000556554		0.13
<b>TRPM4</b>	0.000941039		0.18
<b>AATF</b>	0.00110258		0.09
<b>ALAS1</b>	0.001662622		0.03
<b>ST6GALNAC1</b>	0.002493816		0.35
<b>ITGBL1</b>	0.003929639		0.12
<b>KLK3 exons 2-3</b>	0.004011089		0.04
<b>AGR2</b>	0.00404698		0.11
<b>SERPIN5</b>	0.00406767		0.13
<b>MDK</b>	0.004201206		0.05
<b>SULF2</b>	0.005040354		0.13
<b>GAPDH</b>	0.010197511		0.02
<b>KLK3 exons 1-2</b>	0.010197511		0.05
<b>MCTP1</b>	0.010197511		0.21
<b>PSGR</b>	0.010197511		0.07
<b>HOXC6</b>	0.015699677		0.06
<b>STOM</b>	0.015699677		0.24
<b>AURKA</b>	0.016070673		0.14
<b>PTPRC</b>	0.019655708		0.35
<b>SIM2.long</b>	0.021194067		0.17
<b>NEAT1</b>	0.030247935		0.09
<b>HIST3H2A</b>	0.030856327		0.07
<b>SPON2</b>	0.043657569		0.04
<b>B2M</b>	0.050443121		0.04
<b>RAB17</b>	0.050443121		0.10
<b>BTG2</b>	0.05175315		0.03
<b>SPINK1</b>	0.08344823		-0.09
<b>ANPEP</b>	0.131696047		0.07
<b>HIST1H1E</b>	0.221627267		0.03
<b>HPN</b>	0.261795514		0.05
<b>SULT1A1</b>	0.298725188		0.03

**6.12 Binomial Testing between CB and Ca****Supplementary Table 4 Glm binomial tests – significant probes between CB and Ca (L I H)**

<i>KLK2 Ratio data</i>				<i>KLK2 adjusted data</i>			
<i>Transcript</i>	<i>p-value</i>	<i>Log<sub>2</sub>(FC)</i>	<i>Adjusted p-value</i>	<i>Transcript</i>	<i>p-value</i>	<i>Log<sub>2</sub>(FC)</i>	<i>Adjusted p-value</i>
<i>ERG3' exons 4-5</i>	<i>1.54x10-09</i>	<i>1.582</i>	<i>2.55x10-07</i>	<i>PCA3</i>	<i>4.49 x10-07</i>	<i>0.192</i>	<i>7.46 x10-05</i>
<i>TMPRSS2:ERG</i>	<i>8.73x10-09</i>	<i>NA</i>	<i>1.44x10-06</i>	<i>HPN</i>	<i>4.82 x10-06</i>	<i>0.180</i>	<i>0.001</i>
<i>PCA3</i>	<i>1.10x10-08</i>	<i>0.321</i>	<i>1.81x10-06</i>	<i>SIM2.short</i>	<i>6.21 x10-05</i>	<i>0.124</i>	<i>0.010</i>
<i>ERG3' exons 6-7</i>	<i>2.44x10-08</i>	<i>2.808</i>	<i>3.97x10-06</i>	<i>AMACR</i>	<i>6.40 x10-05</i>	<i>0.124</i>	<i>0.010</i>
<i>HOXC6</i>	<i>8.04x10-07</i>	<i>0.295</i>	<i>0.0001</i>	<i>ERG3' exons 4-5</i>	<i>0.0001</i>	<i>0.103</i>	<i>0.018</i>
<i>TDRD</i>	<i>2.14x10-06</i>	<i>3.683</i>	<i>0.0003</i>	<i>SMIM1</i>	<i>0.0003</i>	<i>0.142</i>	<i>0.048</i>
<i>DLX1</i>	<i>2.76x10-05</i>	<i>4.219</i>	<i>0.004</i>	<i>ERG3' exons 6-7</i>	<i>0.0003</i>	<i>0.101</i>	<i>0.056</i>
<i>ERG5</i>	<i>0.0002</i>	<i>NA</i>	<i>0.025</i>	<i>HOXC6</i>	<i>0.0004</i>	<i>0.130</i>	<i>0.058</i>
<i>ISX</i>	<i>0.0002</i>	<i>2.227</i>	<i>0.028</i>	<i>GJB1</i>	<i>0.0004</i>	<i>0.129</i>	<i>0.061</i>
<i>HOXC4</i>	<i>0.0002</i>	<i>0.900</i>	<i>0.031</i>	<i>TMPRSS2:ERG</i>	<i>0.0004</i>	<i>0.098</i>	<i>0.061</i>
<i>TRPM4</i>	<i>0.0002</i>	<i>0.652</i>	<i>0.032</i>	<i>CAMKK2</i>	<i>0.001</i>	<i>0.098</i>	<i>0.079</i>
<i>PPFIA2</i>	<i>0.0002</i>	<i>0.613</i>	<i>0.032</i>	<i>GAPDH</i>	<i>0.001</i>	<i>0.119</i>	<i>0.116</i>
<i>HPN</i>	<i>0.0003</i>	<i>0.270</i>	<i>0.046</i>	<i>MMP11</i>	<i>0.001</i>	<i>0.083</i>	<i>0.132</i>
<i>GJB1</i>	<i>0.0003</i>	<i>0.234</i>	<i>0.050</i>	<i>TRPM4</i>	<i>0.001</i>	<i>0.103</i>	<i>0.143</i>
<i>APOC1</i>	<i>0.001</i>	<i>1.001</i>	<i>0.093</i>	<i>AMH</i>	<i>0.001</i>	<i>0.112</i>	<i>0.164</i>
<i>AMACR</i>	<i>0.001</i>	<i>0.261</i>	<i>0.099</i>	<i>SIM2.long</i>	<i>0.001</i>	<i>0.121</i>	<i>0.216</i>
<i>DNAH5</i>	<i>0.001</i>	<i>0.485</i>	<i>0.105</i>	<i>RAB17</i>	<i>0.002</i>	<i>0.164</i>	<i>0.286</i>
<i>MCTP1</i>	<i>0.001</i>	<i>1.025</i>	<i>0.116</i>	<i>IMPDH2</i>	<i>0.002</i>	<i>0.101</i>	<i>0.291</i>
<i>SIM2.long</i>	<i>0.001</i>	<i>0.132</i>	<i>0.121</i>	<i>DNAH5</i>	<i>0.002</i>	<i>0.086</i>	<i>0.330</i>
<i>ANKRD34B</i>	<i>0.001</i>	<i>5.060</i>	<i>0.219</i>	<i>TDRD</i>	<i>0.003</i>	<i>0.061</i>	<i>0.390</i>

## 9: APPENDICES

<i>SLC12A1</i>	<b>0.002</b>	<b>0.759</b>	<b>0.280</b>	<i>RIOK3</i>	<b>0.003</b>	<b>0.065</b>	<b>0.445</b>
<i>MEX3A</i>	<b>0.002</b>	<b>0.782</b>	<b>0.286</b>	<i>RP11.97O12.7</i>	<b>0.004</b>	<b>0.104</b>	<b>0.578</b>
<i>PVT1</i>	<b>0.002</b>	<b>0.361</b>	<b>0.328</b>	<i>ISX</i>	<b>0.004</b>	<b>0.075</b>	<b>0.596</b>
<i>CDKN3</i>	<b>0.002</b>	<b>0.275</b>	<b>0.331</b>	<i>TWIST1</i>	<b>0.005</b>	<b>0.055</b>	<b>0.656</b>
<i>RP11.97O12.7</i>	<b>0.002</b>	<b>0.045</b>	<b>0.345</b>	<i>CLU</i>	<b>0.005</b>	<b>0.042</b>	<b>0.690</b>
<i>SSTR1</i>	<b>0.003</b>	<b>0.249</b>	<b>0.411</b>	<i>DLX1</i>	<b>0.007</b>	<b>0.067</b>	<b>0.974</b>
<i>NAALADL2</i>	<b>0.003</b>	<b>0.077</b>	<b>0.417</b>	<i>ANKRD34B</i>	<b>0.007</b>	<b>0.082</b>	<b>0.994</b>
<i>CAMKK2</i>	<b>0.003</b>	<b>0.141</b>	<b>0.429</b>	<i>RNF157</i>	<b>0.007</b>	<b>0.067</b>	<b>0.994</b>
<i>SMIMI</i>	<b>0.003</b>	<b>0.143</b>	<b>0.434</b>	<i>KLK4</i>	<b>0.008</b>	<b>-0.058</b>	<b>0.994</b>
<i>RAB17</i>	<b>0.004</b>	<b>0.207</b>	<b>0.542</b>	<i>ERG5</i>	<b>0.009</b>	<b>0.089</b>	<b>0.994</b>
<i>NEAT1</i>	<b>0.004</b>	<b>0.116</b>	<b>0.554</b>	<i>MYOF</i>	<b>0.009</b>	<b>-0.091</b>	<b>0.994</b>
<i>RIOK3</i>	<b>0.004</b>	<b>0.081</b>	<b>0.576</b>	<i>EN2</i>	<b>0.010</b>	<b>0.069</b>	<b>0.994</b>
<i>SIM2.short</i>	<b>0.005</b>	<b>0.375</b>	<b>0.625</b>	<i>SULT1A1</i>	<b>0.012</b>	<b>0.090</b>	<b>0.994</b>
<i>ST6GALNAC1</i>	<b>0.005</b>	<b>0.389</b>	<b>0.648</b>	<i>CASKINI</i>	<b>0.013</b>	<b>0.056</b>	<b>0.994</b>
<i>GOLM1</i>	<b>0.005</b>	<b>0.193</b>	<b>0.662</b>	<i>PVT1</i>	<b>0.013</b>	<b>0.121</b>	<b>0.994</b>
<i>RPL23AP53</i>	<b>0.005</b>	<b>0.348</b>	<b>0.696</b>	<i>APOC1</i>	<b>0.016</b>	<b>0.095</b>	<b>0.994</b>
<i>SULT1A1</i>	<b>0.005</b>	<b>0.123</b>	<b>0.700</b>	<i>RPS11</i>	<b>0.016</b>	<b>-0.026</b>	<b>0.994</b>
<i>MIC1</i>	<b>0.006</b>	<b>0.395</b>	<b>0.781</b>	<i>MNX1</i>	<b>0.016</b>	<b>0.051</b>	<b>0.994</b>
<i>IMPDH2</i>	<b>0.006</b>	<b>0.074</b>	<b>0.831</b>	<i>GABARAPL2</i>	<b>0.016</b>	<b>-0.078</b>	<b>0.994</b>
<i>RNF157</i>	<b>0.007</b>	<b>0.439</b>	<b>0.885</b>	<i>SLC12A1</i>	<b>0.018</b>	<b>0.080</b>	<b>0.994</b>
<i>SYNM</i>	<b>0.008</b>	<b>0.250</b>	<b>0.946</b>	<i>PSGR</i>	<b>0.018</b>	<b>0.067</b>	<b>0.994</b>
<i>COL9A2</i>	<b>0.008</b>	<b>-2.292</b>	<b>0.967</b>	<i>ITGBL1</i>	<b>0.022</b>	<b>0.068</b>	<b>0.994</b>
<i>AMH</i>	<b>0.008</b>	<b>0.246</b>	<b>0.981</b>	<i>SSPO</i>	<b>0.022</b>	<b>0.077</b>	<b>0.994</b>
<i>CLU</i>	<b>0.008</b>	<b>0.406</b>	<b>0.981</b>	<i>MIC1</i>	<b>0.024</b>	<b>0.094</b>	<b>0.994</b>
<i>MMP11</i>	<b>0.009</b>	<b>0.295</b>	<b>0.992</b>	<i>HMBS</i>	<b>0.024</b>	<b>0.067</b>	<b>0.994</b>
<i>MKi67</i>	<b>0.010</b>	<b>-1.896</b>	<b>0.992</b>	<i>IGFBP3</i>	<b>0.024</b>	<b>-0.016</b>	<b>0.994</b>
<i>MMP26</i>	<b>0.010</b>	<b>0.438</b>	<b>0.992</b>	<i>RPLP2</i>	<b>0.025</b>	<b>-0.069</b>	<b>0.994</b>
<i>SULF2</i>	<b>0.010</b>	<b>1.538</b>	<b>0.992</b>	<i>MFSD2A</i>	<b>0.026</b>	<b>0.078</b>	<b>0.994</b>
<i>MCM7</i>	<b>0.010</b>	<b>0.290</b>	<b>0.992</b>	<i>SYNM</i>	<b>0.026</b>	<b>0.066</b>	<b>0.994</b>
<i>MIR146A.DQ65</i>	<b>0.010</b>	<b>0.324</b>	<b>0.992</b>	<i>NEAT1</i>	<b>0.027</b>	<b>0.061</b>	<b>0.994</b>

## 9: APPENDICES

<i>8414</i>							
<i>EN2</i>	<i>0.011</i>	<i>0.239</i>	<i>0.992</i>	<i>CD10</i>	<i>0.027</i>	<i>-0.046</i>	<i>0.994</i>
<i>TMCC2</i>	<i>0.011</i>	<i>4.237</i>	<i>0.992</i>	<i>CDKN3</i>	<i>0.028</i>	<i>0.045</i>	<i>0.994</i>
<i>ITGBL1</i>	<i>0.011</i>	<i>0.388</i>	<i>0.992</i>	<i>SSTR1</i>	<i>0.029</i>	<i>0.061</i>	<i>0.994</i>
<i>PECI</i>	<i>0.014</i>	<i>0.018</i>	<i>0.992</i>	<i>TMCC2</i>	<i>0.031</i>	<i>0.020</i>	<i>0.994</i>
<i>MIR146A.DQ65</i>							
<i>MMP25</i>	<i>0.015</i>	<i>0.889</i>	<i>0.992</i>	<i>8414</i>	<i>0.031</i>	<i>0.076</i>	<i>0.994</i>
<i>LASS1</i>	<i>0.015</i>	<i>0.199</i>	<i>0.992</i>	<i>ST6GALNAC1</i>	<i>0.032</i>	<i>0.030</i>	<i>0.994</i>
<i>CASKIN1</i>	<i>0.016</i>	<i>0.107</i>	<i>0.992</i>	<i>MMP26</i>	<i>0.033</i>	<i>0.047</i>	<i>0.994</i>
<i>PALM3</i>	<i>0.016</i>	<i>0.123</i>	<i>0.992</i>	<i>HIST1H1E</i>	<i>0.034</i>	<i>0.071</i>	<i>0.994</i>
<i>HPRT</i>	<i>0.017</i>	<i>0.060</i>	<i>0.992</i>	<i>TBP</i>	<i>0.036</i>	<i>0.063</i>	<i>0.994</i>
<i>TMEM45B</i>	<i>0.018</i>	<i>0.272</i>	<i>0.992</i>	<i>MKi67</i>	<i>0.038</i>	<i>0.070</i>	<i>0.994</i>
<i>TMEM86A</i>	<i>0.018</i>	<i>0.898</i>	<i>0.992</i>	<i>STOM</i>	<i>0.041</i>	<i>0.054</i>	<i>0.994</i>
<i>MIR4435.1HG.1</i>							
<i>OC541471</i>	<i>0.019</i>	<i>0.102</i>	<i>0.992</i>	<i>CADPS</i>	<i>0.048</i>	<i>0.032</i>	<i>0.994</i>
<i>SChLAP1</i>	<i>0.019</i>	<i>0.452</i>	<i>0.992</i>	<i>PTN</i>	<i>0.049</i>	<i>-0.048</i>	<i>0.994</i>
<i>STOM</i>	<i>0.021</i>	<i>NA</i>	<i>0.992</i>				
<i>SFRP4</i>	<i>0.022</i>	<i>0.456</i>	<i>0.992</i>				
<i>FOLH1</i>	<i>0.024</i>	<i>0.077</i>	<i>0.992</i>				
<i>MNX1</i>	<i>0.025</i>	<i>0.127</i>	<i>0.992</i>				
<i>TWIST1</i>	<i>0.026</i>	<i>0.103</i>	<i>0.992</i>				
<i>CLIC2</i>	<i>0.027</i>	<i>NA</i>	<i>0.992</i>				
<i>VAX2</i>	<i>0.034</i>	<i>0.170</i>	<i>0.992</i>				
<i>PCSK6</i>	<i>0.036</i>	<i>0.210</i>	<i>0.992</i>				
<i>ACTR5</i>	<i>0.036</i>	<i>0.153</i>	<i>0.992</i>				
<i>CAMK2N2</i>	<i>0.042</i>	<i>0.163</i>	<i>0.992</i>				
<i>ABCB9</i>	<i>0.042</i>	<i>NA</i>	<i>0.992</i>				
<i>EIF2D</i>	<i>0.042</i>	<i>0.054</i>	<i>0.992</i>				
<i>HMBS</i>	<i>0.043</i>	<i>0.107</i>	<i>0.992</i>				
<i>B4GALNT4</i>	<i>0.046</i>	<i>NA</i>	<i>0.992</i>				

9: APPENDICES

<i>Met</i>	<b>0.046</b>	<b>1.819</b>	<b>0.992</b>
<i>HIST3H2A</i>	<b>0.047</b>	<b>0.065</b>	<b>0.992</b>
<i>COL10A1</i>	<b>0.048</b>	<b>0.191</b>	<b>0.992</b>

<i>KLK3 Adjusted data</i>				<i>HK normalised data</i>			
<i>Transcript</i>	<i>p-value</i>	<i>Log<sub>2</sub>(FC)</i>	<i>Adjusted p-value</i>	<i>Transcript</i>	<i>p-value</i>	<i>Log<sub>2</sub>(FC)</i>	<i>Adjusted p-value</i>
<i>PCA3</i>	<b>1.61x10-06</b>	<b>0.14</b>	<b>0.0003</b>	<i>ERG3' exons 4-5</i>	<b>4.58x10-09</b>	<b>0.699</b>	<b>7.64x10-07</b>
<i>HPN</i>	<b>3.27x10-05</b>	<b>0.13</b>	<b>0.01</b>	<i>PCA3</i>	<b>1.40x10-08</b>	<b>0.191</b>	<b>2.32x10-06</b>
<i>SIM2.short</i>	<b>0.0002</b>	<b>0.091</b>	<b>0.029</b>	<i>TMPRSS2:ERG</i>	<b>4.02x10-08</b>	<b>1.006</b>	<b>6.63x10-06</b>
<i>ERG3' exons 4-5</i>	<b>0.0002</b>	<b>0.080</b>	<b>0.031</b>	<i>ERG3' exons 6-7</i>	<b>4.79x10-07</b>	<b>1.130</b>	<b>7.86x10-05</b>
<i>HOXC6</i>	<b>0.001</b>	<b>0.113</b>	<b>0.084</b>	<i>HOXC6</i>	<b>3.71x10-06</b>	<b>0.178</b>	<b>0.001</b>
<i>ERG3' exons 6-7</i>	<b>0.001</b>	<b>0.062</b>	<b>0.117</b>	<i>TDRD</i>	<b>2.70x10-05</b>	<b>0.848</b>	<b>0.004</b>
<i>AMACR</i>	<b>0.001</b>	<b>0.086</b>	<b>0.118</b>	<i>HPN</i>	<b>0.0002</b>	<b>0.123</b>	<b>0.028</b>
<i>TMPRSS2:ERG</i>	<b>0.001</b>	<b>0.062</b>	<b>0.119</b>	<i>HOXC4</i>	<b>0.0003</b>	<b>0.200</b>	<b>0.046</b>
<i>SMIMI</i>	<b>0.001</b>	<b>0.124</b>	<b>0.217</b>	<i>DLX1</i>	<b>0.0004</b>	<b>0.424</b>	<b>0.057</b>
<i>KLK4</i>	<b>0.001</b>	<b>-0.099</b>	<b>0.219</b>	<i>APOC1</i>	<b>0.0004</b>	<b>0.390</b>	<b>0.057</b>
<i>GJB1</i>	<b>0.002</b>	<b>0.103</b>	<b>0.324</b>	<i>ERG5'</i>	<b>0.001</b>	<b>0.175</b>	<b>0.157</b>
<i>TRPM4</i>	<b>0.004</b>	<b>0.079</b>	<b>0.572</b>	<i>GJB1</i>	<b>0.001</b>	<b>0.129</b>	<b>0.182</b>
<i>IMPDH2</i>	<b>0.004</b>	<b>0.102</b>	<b>0.627</b>	<i>MCTP1</i>	<b>0.001</b>	<b>0.333</b>	<b>0.183</b>
<i>MYOF</i>	<b>0.005</b>	<b>-0.098</b>	<b>0.732</b>	<i>ISX</i>	<b>0.002</b>	<b>0.190</b>	<b>0.247</b>
<i>RAB17</i>	<b>0.005</b>	<b>0.106</b>	<b>0.738</b>	<i>SSTR1</i>	<b>0.002</b>	<b>0.035</b>	<b>0.252</b>
<i>SIM2.long</i>	<b>0.005</b>	<b>0.105</b>	<b>0.754</b>	<i>PPFIA2</i>	<b>0.002</b>	<b>0.312</b>	<b>0.255</b>
<i>AMH</i>	<b>0.005</b>	<b>0.075</b>	<b>0.787</b>	<i>TRPM4</i>	<b>0.002</b>	<b>0.294</b>	<b>0.326</b>
<i>PTN</i>	<b>0.007</b>	<b>-0.071</b>	<b>0.998</b>	<i>RAB17</i>	<b>0.002</b>	<b>0.111</b>	<b>0.366</b>
<i>CAMKK2</i>	<b>0.008</b>	<b>0.052</b>	<b>0.998</b>	<i>SIM2.long</i>	<b>0.003</b>	<b>0.079</b>	<b>0.432</b>
<i>ISX</i>	<b>0.008</b>	<b>0.062</b>	<b>0.998</b>	<i>SLC12A1</i>	<b>0.004</b>	<b>0.233</b>	<b>0.560</b>
<i>DLX1</i>	<b>0.009</b>	<b>0.042</b>	<b>0.998</b>	<i>SIM2.short</i>	<b>0.004</b>	<b>0.064</b>	<b>0.566</b>
<i>MMP11</i>	<b>0.009</b>	<b>0.070</b>	<b>0.998</b>	<i>AMACR</i>	<b>0.004</b>	<b>0.139</b>	<b>0.591</b>
<i>TDRD</i>	<b>0.010</b>	<b>0.042</b>	<b>0.998</b>	<i>MMP11</i>	<b>0.005</b>	<b>0.066</b>	<b>0.665</b>

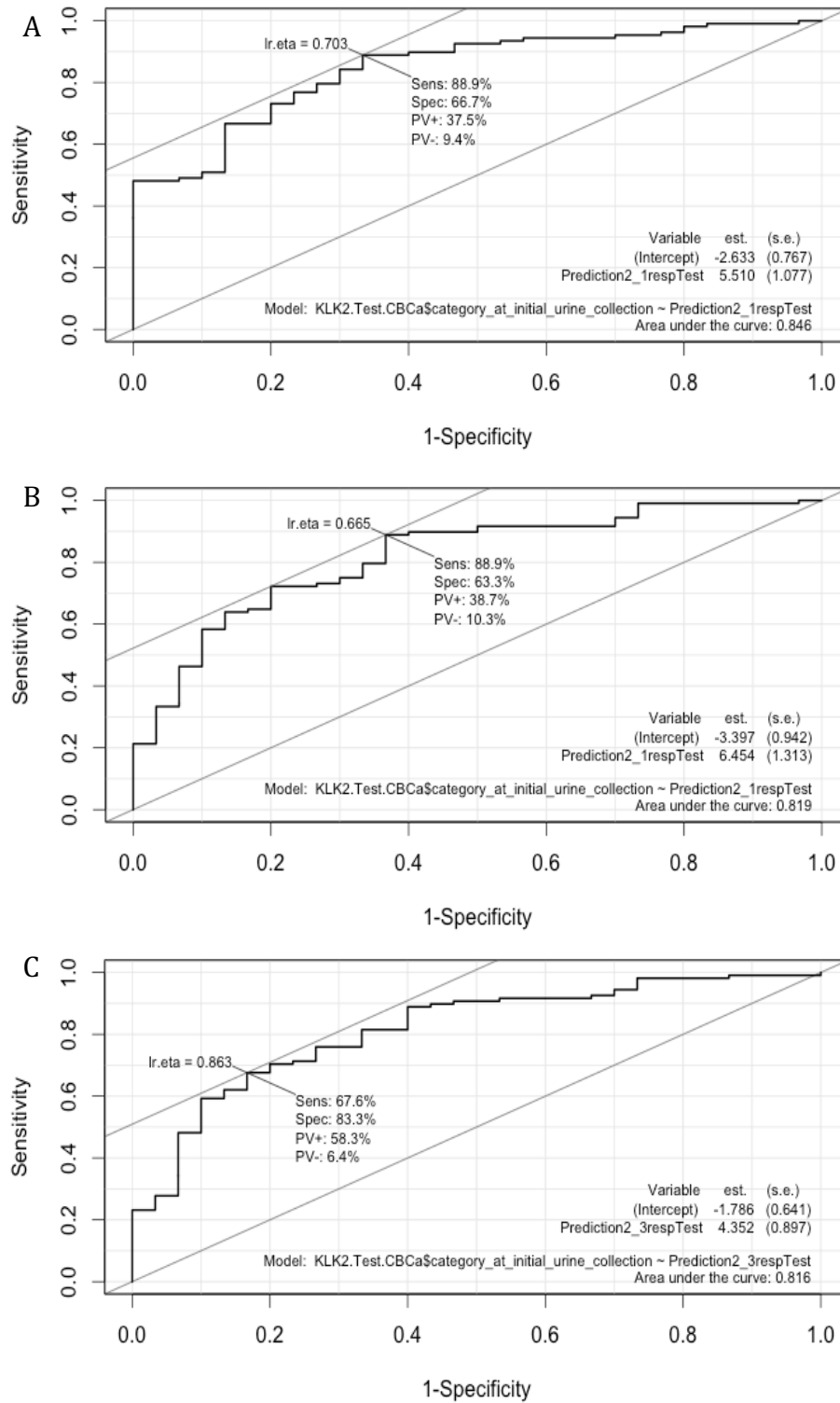
9: APPENDICES

<i>ERG5</i>	<i>0.012</i>	<i>0.055</i>	<i>0.998</i>	<i>ANKRD34B</i>	<i>0.005</i>	<i>0.099</i>	<i>0.724</i>
<i>GAPDH</i>	<i>0.012</i>	<i>0.078</i>	<i>0.998</i>	<i>DNAH5</i>	<i>0.006</i>	<i>0.266</i>	<i>0.795</i>
<i>SULT1A1</i>	<i>0.013</i>	<i>0.063</i>	<i>0.998</i>	<i>AMH</i>	<i>0.007</i>	<i>0.046</i>	<i>0.994</i>
<i>IGFBP3</i>	<i>0.014</i>	<i>-0.033</i>	<i>0.998</i>	<i>RP11_97O12.7</i>	<i>0.008</i>	<i>0.050</i>	<i>0.994</i>
<i>RIOK3</i>	<i>0.015</i>	<i>0.063</i>	<i>0.998</i>	<i>MEX3A</i>	<i>0.009</i>	<i>0.184</i>	<i>0.994</i>
<i>TWIST1</i>	<i>0.016</i>	<i>0.032</i>	<i>0.998</i>	<i>PVT1</i>	<i>0.011</i>	<i>0.079</i>	<i>0.994</i>
<i>RP11.97O12.7</i>	<i>0.016</i>	<i>0.066</i>	<i>0.998</i>	<i>SMIM1</i>	<i>0.011</i>	<i>0.082</i>	<i>0.994</i>
<i>ANKRD34B</i>	<i>0.016</i>	<i>0.069</i>	<i>0.998</i>	<i>EN2</i>	<i>0.011</i>	<i>0.059</i>	<i>0.994</i>
<i>DNAH5</i>	<i>0.017</i>	<i>0.052</i>	<i>0.998</i>	<i>CASKIN1</i>	<i>0.013</i>	<i>0.035</i>	<i>0.994</i>
<i>CD10</i>	<i>0.017</i>	<i>-0.040</i>	<i>0.998</i>	<i>KLK4</i>	<i>0.014</i>	<i>-0.031</i>	<i>0.994</i>
<i>MARCH5</i>	<i>0.018</i>	<i>-0.077</i>	<i>0.998</i>	<i>ITGBL1</i>	<i>0.015</i>	<i>0.092</i>	<i>0.994</i>
<i>GABARAPL2</i>	<i>0.019</i>	<i>-0.073</i>	<i>0.998</i>	<i>NEAT1</i>	<i>0.015</i>	<i>0.112</i>	<i>0.994</i>
<i>APOC1</i>	<i>0.019</i>	<i>0.051</i>	<i>0.998</i>	<i>SULT1A1</i>	<i>0.017</i>	<i>0.068</i>	<i>0.994</i>
<i>SLC12A1</i>	<i>0.022</i>	<i>0.049</i>	<i>0.998</i>	<i>CDKN3</i>	<i>0.019</i>	<i>0.080</i>	<i>0.994</i>
<i>SSTR1</i>	<i>0.022</i>	<i>0.025</i>	<i>0.998</i>	<i>RIOK3</i>	<i>0.022</i>	<i>0.023</i>	<i>0.994</i>
<i>CLU</i>	<i>0.025</i>	<i>0.033</i>	<i>0.998</i>	<i>MIR146A</i>	<i>0.023</i>	<i>0.079</i>	<i>0.994</i>
<i>ITGBL1</i>	<i>0.025</i>	<i>0.063</i>	<i>0.998</i>	<i>TMEM45B</i>	<i>0.023</i>	<i>0.038</i>	<i>0.994</i>
<i>EN2</i>	<i>0.026</i>	<i>0.049</i>	<i>0.998</i>	<i>NAALADL2</i>	<i>0.025</i>	<i>0.061</i>	<i>0.994</i>
<i>RPS11</i>	<i>0.026</i>	<i>-0.083</i>	<i>0.998</i>	<i>TWIST1</i>	<i>0.029</i>	<i>0.003</i>	<i>0.994</i>
<i>RNF157</i>	<i>0.026</i>	<i>0.043</i>	<i>0.998</i>	<i>RPL23AP53</i>	<i>0.030</i>	<i>0.181</i>	<i>0.994</i>
<i>MNX1</i>	<i>0.026</i>	<i>0.018</i>	<i>0.998</i>	<i>PALM3</i>	<i>0.033</i>	<i>0.061</i>	<i>0.994</i>
<i>PVT1</i>	<i>0.035</i>	<i>0.054</i>	<i>0.998</i>	<i>SULF2</i>	<i>0.035</i>	<i>0.055</i>	<i>0.994</i>
<i>MIC1</i>	<i>0.043</i>	<i>0.057</i>	<i>0.998</i>	<i>COL9A2</i>	<i>0.035</i>	<i>0.140</i>	<i>0.994</i>
<i>CASKIN1</i>	<i>0.044</i>	<i>0.038</i>	<i>0.998</i>	<i>RNF157</i>	<i>0.035</i>	<i>0.180</i>	<i>0.994</i>
<i>MIR146A.DQ658414</i>	<i>0.044</i>	<i>0.072</i>	<i>0.998</i>	<i>CLU</i>	<i>0.037</i>	<i>0.078</i>	<i>0.994</i>
<i>STOM</i>	<i>0.046</i>	<i>0.019</i>	<i>0.998</i>	<i>MIR4435_1HG</i>	<i>0.039</i>	<i>0.086</i>	<i>0.994</i>
				<i>MMP25</i>	<i>0.040</i>	<i>0.030</i>	<i>0.994</i>
				<i>MIC1</i>	<i>0.040</i>	<i>0.102</i>	<i>0.994</i>
				<i>RPS11</i>	<i>0.040</i>	<i>-0.007</i>	<i>0.994</i>
				<i>IMPDH2</i>	<i>0.041</i>	<i>0.055</i>	<i>0.994</i>

9: APPENDICES

	<i>MKi67</i>	<i>0.042</i>	<i>0.371</i>	<i>0.994</i>
	<i>TMCC2</i>	<i>0.043</i>	<i>0.029</i>	<i>0.994</i>

9: APPENDICES



Supplementary Figure 1 KLK2 Ratio Data ROC curves for test data using models detecting between CB and Ca (L I H) for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.

Supplementary Table 5 Lasso output for models detecting between CB and Ca (L I H) using KLK2 ratio data.

*All Transcripts*

*Significant Transcripts*

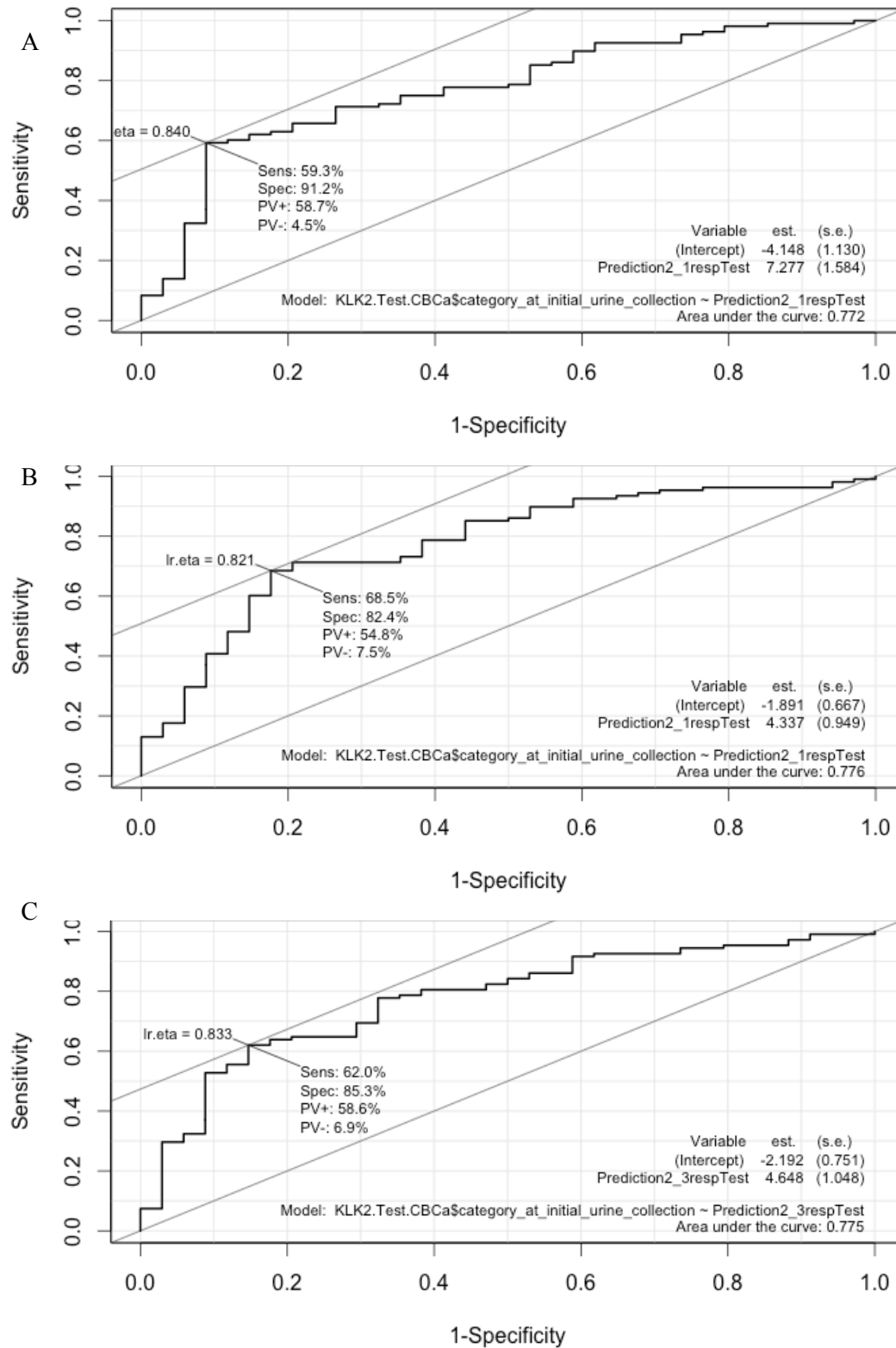
*Multiple testing corrected*



## 9: APPENDICES

<i>Transcripts</i>					
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>ERG3' exons 4-5</i>	<b>0.28</b>	<i>ERG3' exons 4-5</i>	<b>0.21</b>	<i>ERG3' exons 4-5</i>	<b>0.25</b>
<i>TMPRSS2:ERG</i>	<b>0.22</b>	<i>TMPRSS2:ERG</i>	<b>0.20</b>	<i>TMPRSS2:ERG</i>	<b>0.24</b>
<i>PCA3</i>	<b>0.20</b>	<i>PCA3</i>	<b>0.17</b>	<i>PCA3</i>	<b>0.20</b>
<i>HOXC6</i>	<b>0.08</b>	<i>HOXC6</i>	<b>0.07</b>	<i>HOXC6</i>	<b>0.08</b>
<i>ISX</i>	<b>0.08</b>			<i>ISX</i>	<b>0.03</b>
<i>APOC1</i>	<b>0.06</b>			<i>GJB1</i>	<b>0.02</b>
<i>GJB1</i>	<b>0.06</b>			<i>DLX1</i>	<b>0.01</b>
<i>AMACR</i>	<b>0.05</b>			<i>TDRD</i>	<b>0.01</b>
<i>NEAT1</i>	<b>0.03</b>				
<i>DLX1</i>	<b>0.02</b>				
<i>TDRD</i>	<b>0.02</b>				
<i>TMEM47</i>	<b>0.01</b>				
<i>SULT1A1</i>	<b>0.01</b>				
<i>RNF157</i>	<b>0.01</b>				
<i>ST6GALNAC1</i>	<b>0.00</b>				
<i>IGFBP3</i>	<b>-0.01</b>				
<i>ARexon9</i>	<b>-0.06</b>				
<i>PPP1R12B</i>	<b>-0.08</b>				
<i>CP</i>	<b>-0.11</b>				
<i>MXII</i>	<b>-0.16</b>				
<i>KLK4</i>	<b>-0.24</b>				

9: APPENDICES



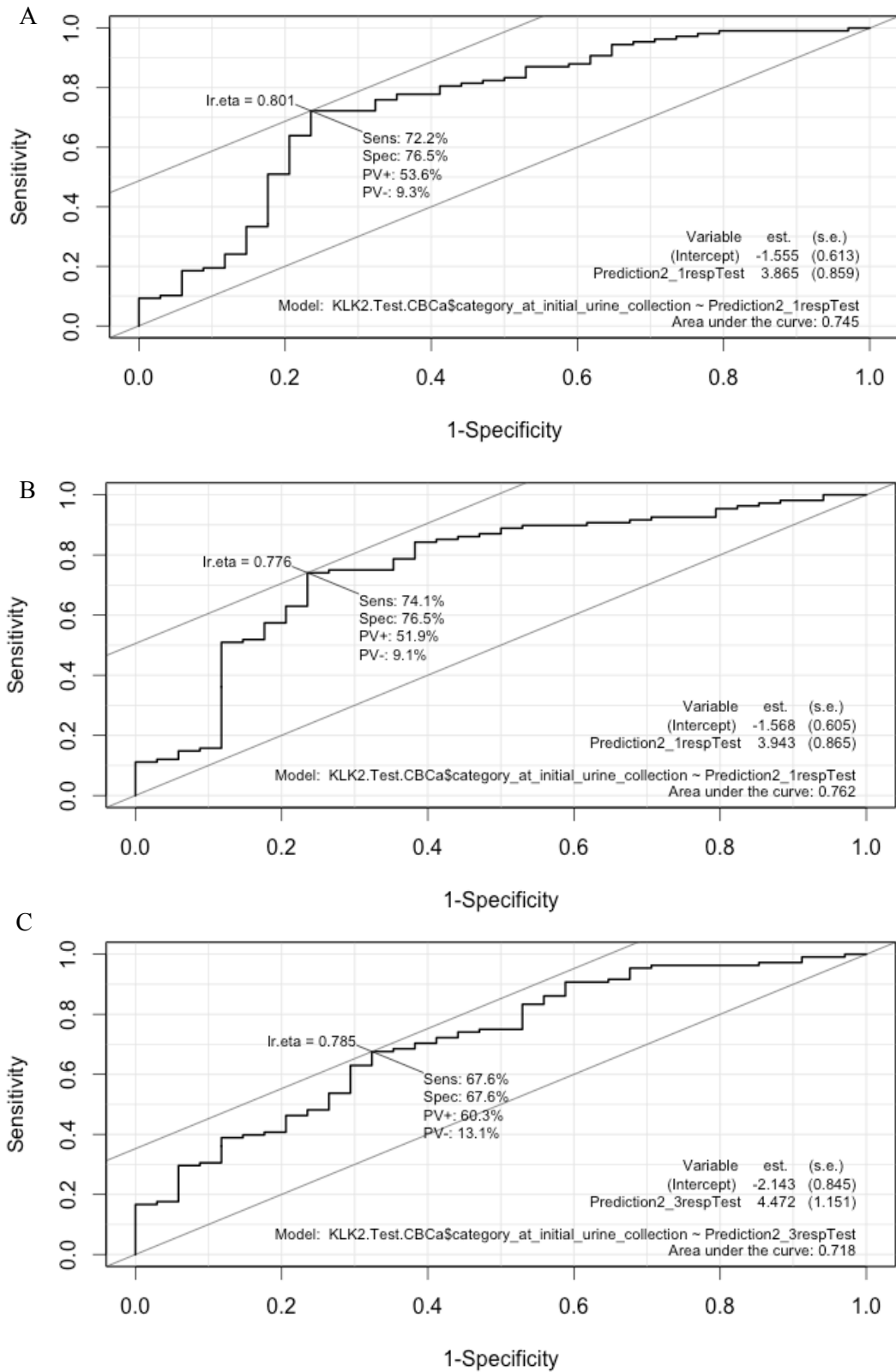
Supplementary Figure 2 KLK2 Adjusted Data ROC curves for test data using models detecting between CB and Ca (L I H) for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.

9: APPENDICES

Supplementary Table 6 Lasso output for models detecting between CB and Ca (L I H) using KLK2 adjusted data.

<i>All Transcript</i>		<i>Significant Transcripts</i>		<i>Multiple Testing correction Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>PCA3</i>	<b>2.89</b>	<i>PCA3</i>	<b>2.85</b>	<i>PCA3</i>	<b>4.59</b>
<i>AMACR</i>	<b>1.73</b>	<i>AMACR</i>	<b>2.68</b>	<i>SIM2.short</i>	<b>2.88</b>
<i>SIM2.short</i>	<b>1.07</b>	<i>ERG3' exons 4-5</i>	<b>1.77</b>	<i>ERG3' exons 4-5</i>	<b>2.50</b>
<i>SMIM1</i>	<b>0.89</b>	<i>SMIM1</i>	<b>1.43</b>	<i>AMACR</i>	<b>2.44</b>
<i>AMH</i>	<b>0.89</b>	<i>RP11.97O12.7</i>	<b>1.38</b>	<i>HPN</i>	<b>2.02</b>
<i>ERG3' exons 4-5</i>	<b>0.70</b>	<i>SIM2.short</i>	<b>1.25</b>	<i>SMIM1</i>	<b>1.89</b>
<i>CLU</i>	<b>0.49</b>	<i>CAMKK2</i>	<b>1.03</b>		
<i>HPN</i>	<b>0.46</b>	<i>AMH</i>	<b>0.93</b>		
<i>CAMKK2</i>	<b>0.39</b>	<i>CLU</i>	<b>0.90</b>		
<i>GAPDH</i>	<b>0.29</b>	<i>RNF157</i>	<b>0.68</b>		
<i>RP11.97O12.7</i>	<b>0.24</b>	<i>DNAH5</i>	<b>0.63</b>		
<i>DNAH5</i>	<b>0.20</b>	<i>RIOK3</i>	<b>0.52</b>		
<i>RNF157</i>	<b>0.18</b>	<i>NEAT1</i>	<b>0.43</b>		
<i>APOC1</i>	<b>0.17</b>	<i>APOC1</i>	<b>0.42</b>		
<i>ERG3' exons 6-7</i>	<b>0.17</b>	<i>DLX1</i>	<b>0.38</b>		
<i>RIOK3</i>	<b>0.05</b>	<i>TBP</i>	<b>0.33</b>		
<i>MMP25</i>	<b>0.02</b>	<i>MMP11</i>	<b>0.23</b>		
<i>CP</i>	<b>-0.01</b>	<i>SYNM</i>	<b>0.23</b>		
<i>CD10</i>	<b>-0.03</b>	<i>CADPS</i>	<b>0.18</b>		
<i>AR exon 9</i>	<b>-0.08</b>	<i>SLC12A1</i>	<b>0.15</b>		
<i>PTN</i>	<b>-0.42</b>	<i>MIC1</i>	<b>0.11</b>		
<i>IGFBP3</i>	<b>-0.45</b>	<i>HPN</i>	<b>0.10</b>		
<i>MYOF</i>	<b>-0.68</b>	<i>STOM</i>	<b>0.08</b>		
<i>GABARAPL2</i>	<b>-0.75</b>	<i>MKi67</i>	<b>0.05</b>		
<i>KLK4</i>	<b>-1.09</b>	<i>RPS11</i>	<b>-0.13</b>		
<i>MARCH5</i>	<b>-1.11</b>	<i>CD10</i>	<b>-0.98</b>		
		<i>IGFBP3</i>	<b>-1.17</b>		
		<i>PTN</i>	<b>-1.25</b>		
		<i>MYOF</i>	<b>-1.66</b>		
		<i>KLK4</i>	<b>-1.87</b>		
		<i>GABARAPL2</i>	<b>-2.90</b>		

9: APPENDICES



**Supplementary Figure 3** KLK3 Adjusted Data ROC curves for test data using models detecting between CB and Ca (L I H) for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.

## 9: APPENDICES

## Supplementary Table 7 Lasso output for models detecting between CB and Ca (L I H)

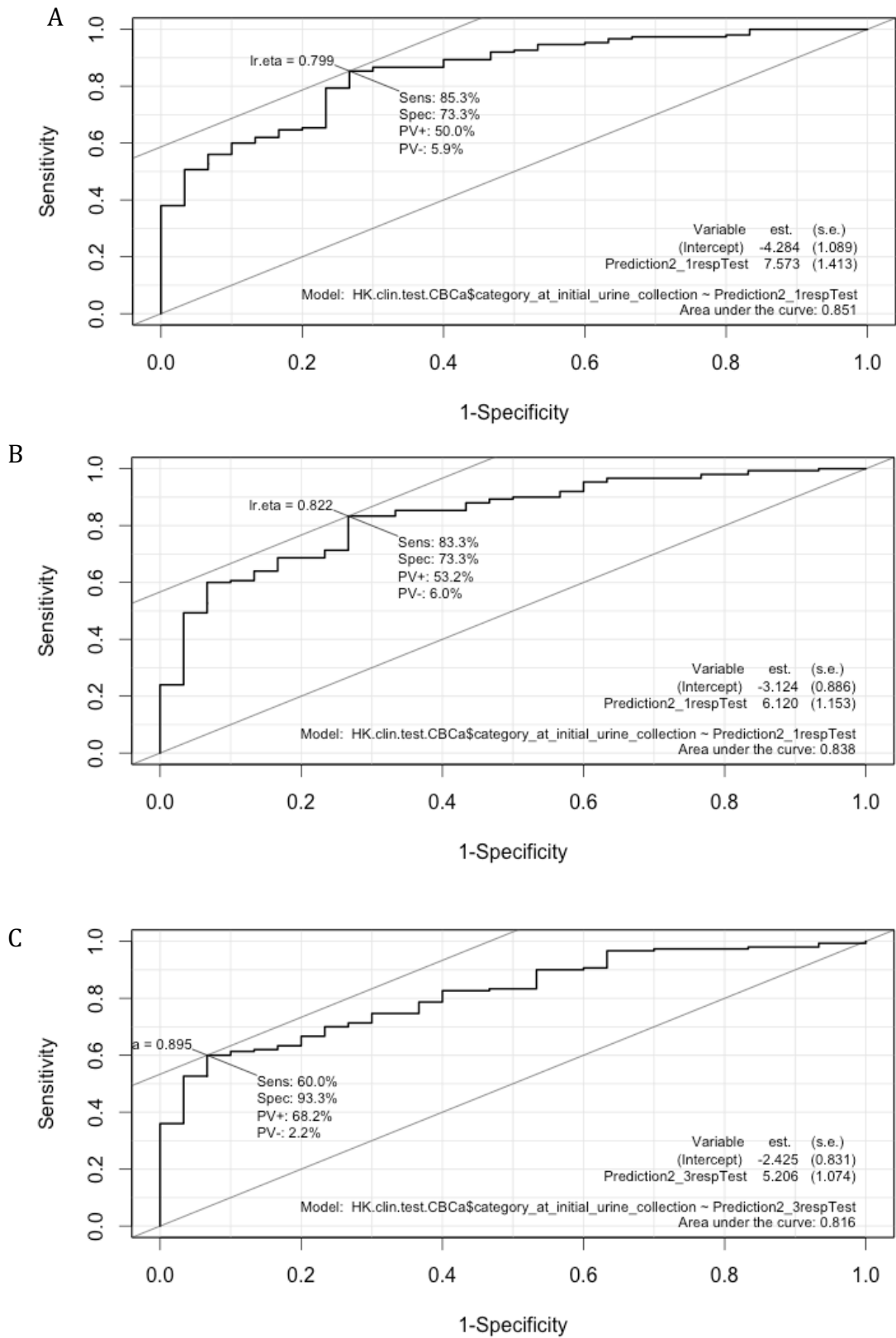
using KLK3 adjusted data

<i>All Transcript</i>		<i>Significant Transcripts</i>		<i>Multiple Testing correction Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>PCA3</i>	<b>2.94</b>	<i>PCA3</i>	<b>2.88</b>	<i>PCA3</i>	<b>4.40</b>
<i>ERG3' exons 4-5</i>	<b>1.74</b>	<i>AMACR</i>	<b>2.30</b>	<i>SIM2.short</i>	<b>2.87</b>
<i>AMACR</i>	<b>1.57</b>	<i>ERG3' exons 4-5</i>	<b>2.01</b>	<i>HPN</i>	<b>2.15</b>
<i>SIM2.short</i>	<b>1.55</b>	<i>SMIM1</i>	<b>1.25</b>	<i>ERG3' exons 4-5</i>	<b>1.97</b>
<i>SMIM1</i>	<b>1.24</b>	<i>SIM2.short</i>	<b>1.19</b>		
<i>AMH</i>	<b>1.23</b>	<i>AMH</i>	<b>0.82</b>		
<i>APOC1</i>	<b>0.63</b>	<i>CLU</i>	<b>0.71</b>		
<i>NEAT1</i>	<b>0.60</b>	<i>RIOK3</i>	<b>0.70</b>		
<i>MMP25</i>	<b>0.59</b>	<i>CAMKK2</i>	<b>0.69</b>		
<i>TBP</i>	<b>0.55</b>	<i>APOC1</i>	<b>0.64</b>		
<i>SERPINB5</i>	<b>0.52</b>	<i>HPN</i>	<b>0.62</b>		
<i>HPN</i>	<b>0.47</b>	<i>RNF157</i>	<b>0.52</b>		
<i>CLU</i>	<b>0.34</b>	<i>DLX1</i>	<b>0.36</b>		
<i>DLX1</i>	<b>0.31</b>	<i>MMP11</i>	<b>0.30</b>		
<i>RNF157</i>	<b>0.30</b>	<i>SLC12A1</i>	<b>0.22</b>		
<i>CAMKK2</i>	<b>0.29</b>	<i>SULT1A1</i>	<b>0.22</b>		
<i>PPAP2A</i>	<b>0.26</b>	<i>ISX</i>	<b>0.09</b>		
<i>MMP11</i>	<b>0.25</b>	<i>DNAH5</i>	<b>0.09</b>		
<i>SLC12A1</i>	<b>0.19</b>	<i>EN2</i>	<b>0.07</b>		
<i>STOM</i>	<b>0.17</b>	<i>STOM</i>	<b>0.06</b>		
<i>CADPS</i>	<b>0.15</b>	<i>ANKRD34B</i>	<b>0.02</b>		
<i>RIOK3</i>	<b>0.15</b>	<i>CD10</i>	<b>-0.47</b>		
<i>EN2</i>	<b>0.13</b>	<i>RPS11</i>	<b>-0.53</b>		
<i>ISX</i>	<b>0.13</b>	<i>IGFBP3</i>	<b>-1.08</b>		
<i>COL10A1</i>	<b>0.12</b>	<i>KLK4</i>	<b>-1.17</b>		
<i>ST6GALNAC1</i>	<b>0.12</b>	<i>MYOF</i>	<b>-1.46</b>		
<i>MNX1</i>	<b>0.11</b>	<i>PTN</i>	<b>-1.57</b>		
<i>DNAH5</i>	<b>0.11</b>	<i>GABARAPL2</i>	<b>-2.02</b>		
<i>SULT1A1</i>	<b>0.08</b>	<i>MARCH5</i>	<b>-2.04</b>		
<i>HOXC6</i>	<b>0.07</b>				
<i>GJB1</i>	<b>0.04</b>				
<i>ERG5</i>	<b>0.03</b>				
<i>RP11.244H18.1.P712P</i>	<b>-0.14</b>				
<i>SPON2</i>	<b>-0.16</b>				
<i>CLIC2</i>	<b>-0.20</b>				
<i>PPP1R12B</i>	<b>-0.21</b>				
<i>CD10</i>	<b>-0.23</b>				
<i>CP</i>	<b>-0.27</b>				
<i>AR exon 9</i>	<b>-0.30</b>				
<i>MXII</i>	<b>-0.32</b>				
<i>CDC20</i>	<b>-0.39</b>				
<i>CKAP2L</i>	<b>-0.43</b>				

## 9: APPENDICES

<i>Timp4</i>	<b>-0.45</b>
<i>RPS11</i>	<b>-0.61</b>
<i>IGFBP3</i>	<b>-0.90</b>
<i>MYOF</i>	<b>-1.18</b>
<i>PTN</i>	<b>-1.29</b>
<i>KLK4</i>	<b>-1.36</b>
<i>MARCH5</i>	<b>-1.51</b>
<i>GABARAPL2</i>	<b>-1.54</b>

9: APPENDICES



**Supplementary Figure 4 HK normalised data ROC curves for test data using models detecting between CB and Ca (L I H) for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.**

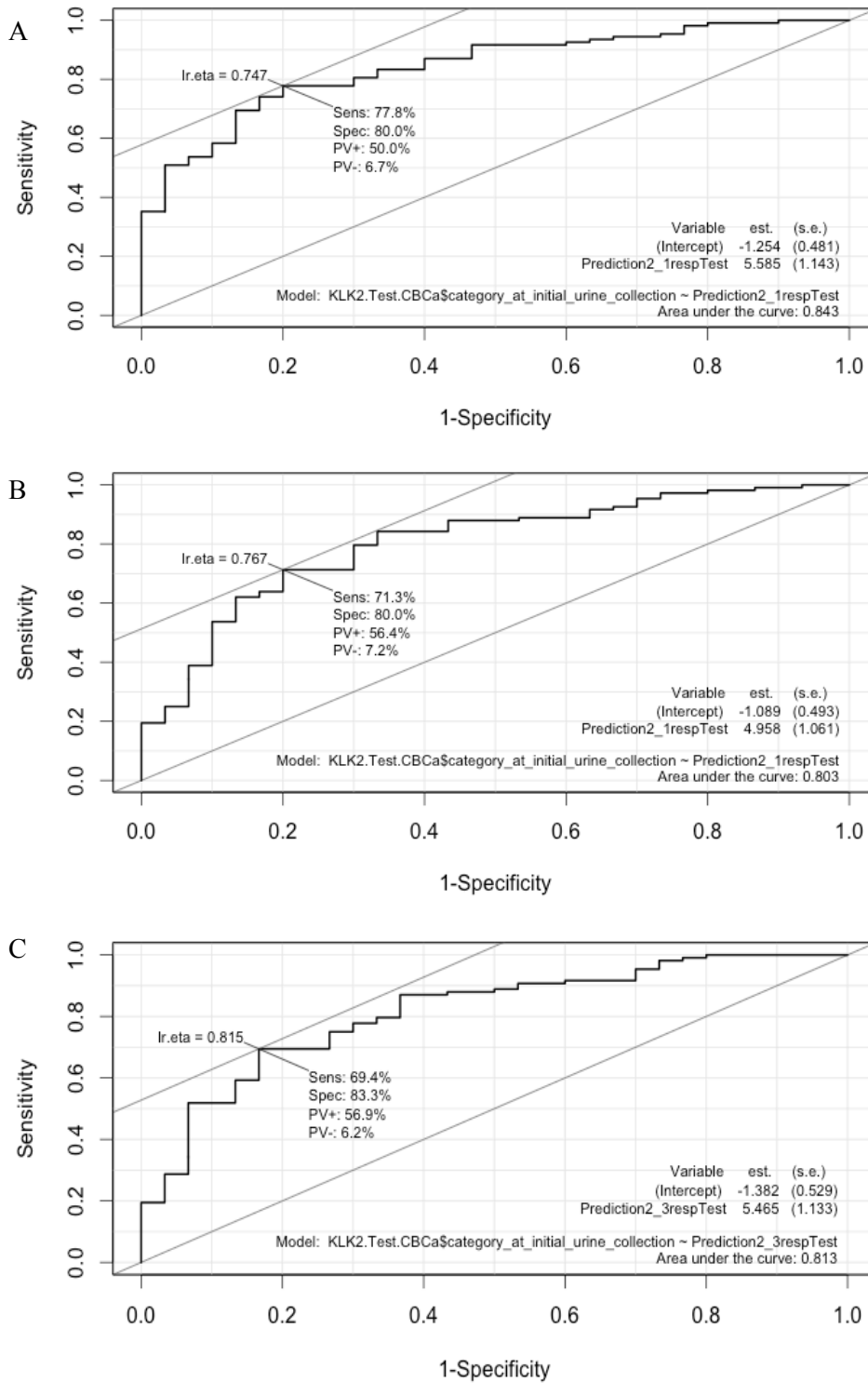
9: APPENDICES

**Supplementary Table 8** Lasso output for models detecting between CB and Ca (L I H) using HK normalised data.

<i>All Transcript</i>		<i>Significant Transcripts</i>		<i>Multiple Testing correction Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>PCA3</i>	<b>0.29</b>	<i>PCA3</i>	<b>0.35</b>	<i>PCA3</i>	<b>0.29</b>
<i>TMPRSS2:ERG</i>	<b>0.18</b>	<i>ERG3' exons 4-5</i>	<b>0.19</b>	<i>TMPRSS2:ERG</i>	<b>0.23</b>
<i>ERG3' exons 4-5</i>	<b>0.17</b>	<i>TMPRSS2:ERG</i>	<b>0.18</b>	<i>ERG3' exons 4-5</i>	<b>0.19</b>
<i>APOC1</i>	<b>0.11</b>	<i>APOC1</i>	<b>0.13</b>	<i>HPN</i>	<b>0.04</b>
<i>ISX</i>	<b>0.04</b>	<i>SLC12A1</i>	<b>0.05</b>	<i>HOXC6</i>	<b>0.03</b>
<i>SLC12A1</i>	<b>0.04</b>	<i>ISX</i>	<b>0.04</b>		
<i>HOXC6</i>	<b>0.04</b>	<i>MCTP1</i>	<b>0.03</b>		
<i>MCTP1</i>	<b>0.03</b>	<i>HOXC6</i>	<b>0.02</b>		
<i>TDRD</i>	<b>0.00</b>	<i>SULT1A1</i>	<b>0.00</b>		
<i>PDLIM5</i>	<b>-0.01</b>	<i>KLK4</i>	<b>-0.41</b>		
<i>CD10</i>	<b>-0.02</b>				
<i>GABARAPL2</i>	<b>-0.02</b>				
<i>PTN</i>	<b>-0.02</b>				
<i>AR exon 9</i>	<b>-0.04</b>				
<i>PPP1R12B</i>	<b>-0.04</b>				
<i>CP</i>	<b>-0.08</b>				
<i>MXII</i>	<b>-0.15</b>				
<i>KLK4</i>	<b>-0.20</b>				



**6.13 Binomial Testing between CB and Ca (Random Sampling)**



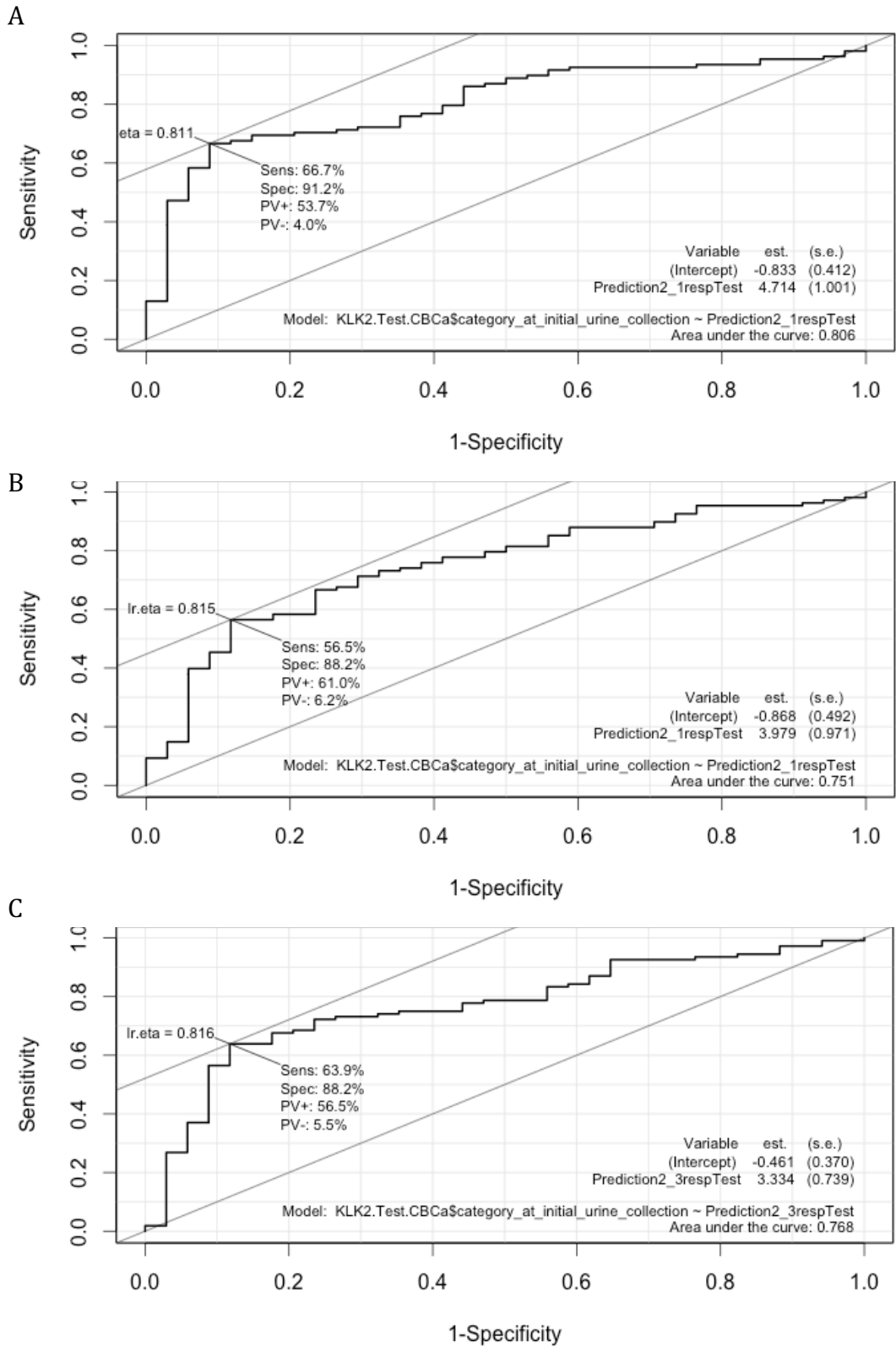
**Supplementary Figure 5 KLK2 ratio data ROC curves for test data using models (random sampling) detecting between CB and Ca (L I H) for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.**

9: APPENDICES

Supplementary Table 9 Lasso output for models (random sampling detecting between CB and Ca (L I H) using KLK2 ratio data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>ERG3' exons 4-5</i>	<b>0.51</b>	<i>PCA3</i>	<b>0.21</b>	<i>ERG3' exons 4-5</i>	<b>0.25</b>
<i>PCA3</i>	<b>0.14</b>	<i>ERG3' exons 4-5</i>	<b>0.20</b>	<i>PCA3</i>	<b>0.24</b>
<i>TMPRSS2:ERG</i>	<b>0.14</b>	<i>TMPRSS2:ERG</i>	<b>0.15</b>	<i>TMPRSS2:ERG</i>	<b>0.17</b>
<i>SLC12A1</i>	<b>0.06</b>	<i>AMACR</i>	<b>0.08</b>	<i>HOXC6</i>	<b>0.02</b>
<i>ERG5</i>	<b>0.05</b>	<i>GJB1</i>	<b>0.06</b>	<i>GJB1</i>	<b>0.01</b>
<i>GJB1</i>	<b>0.04</b>	<i>NEAT1</i>	<b>0.03</b>		
<i>HOXC6</i>	<b>0.04</b>	<i>TDRD</i>	<b>0.03</b>		
<i>TDRD</i>	<b>0.01</b>	<i>DLX1</i>	<b>0.02</b>		
<i>LASS1</i>	<b>0.00</b>	<i>TRPM4</i>	<b>0.01</b>		
<i>HIST1H2BF</i>	<b>-0.01</b>				
<i>CP</i>	<b>-0.02</b>				
<i>CKAP2L</i>	<b>-0.03</b>				
<i>DPP4</i>	<b>-0.04</b>				
<i>PTN</i>	<b>-0.07</b>				
<i>ZNF577</i>	<b>-0.08</b>				
<i>MYOF</i>	<b>-0.10</b>				
<i>GABARAPL2</i>	<b>-0.31</b>				

9: APPENDICES



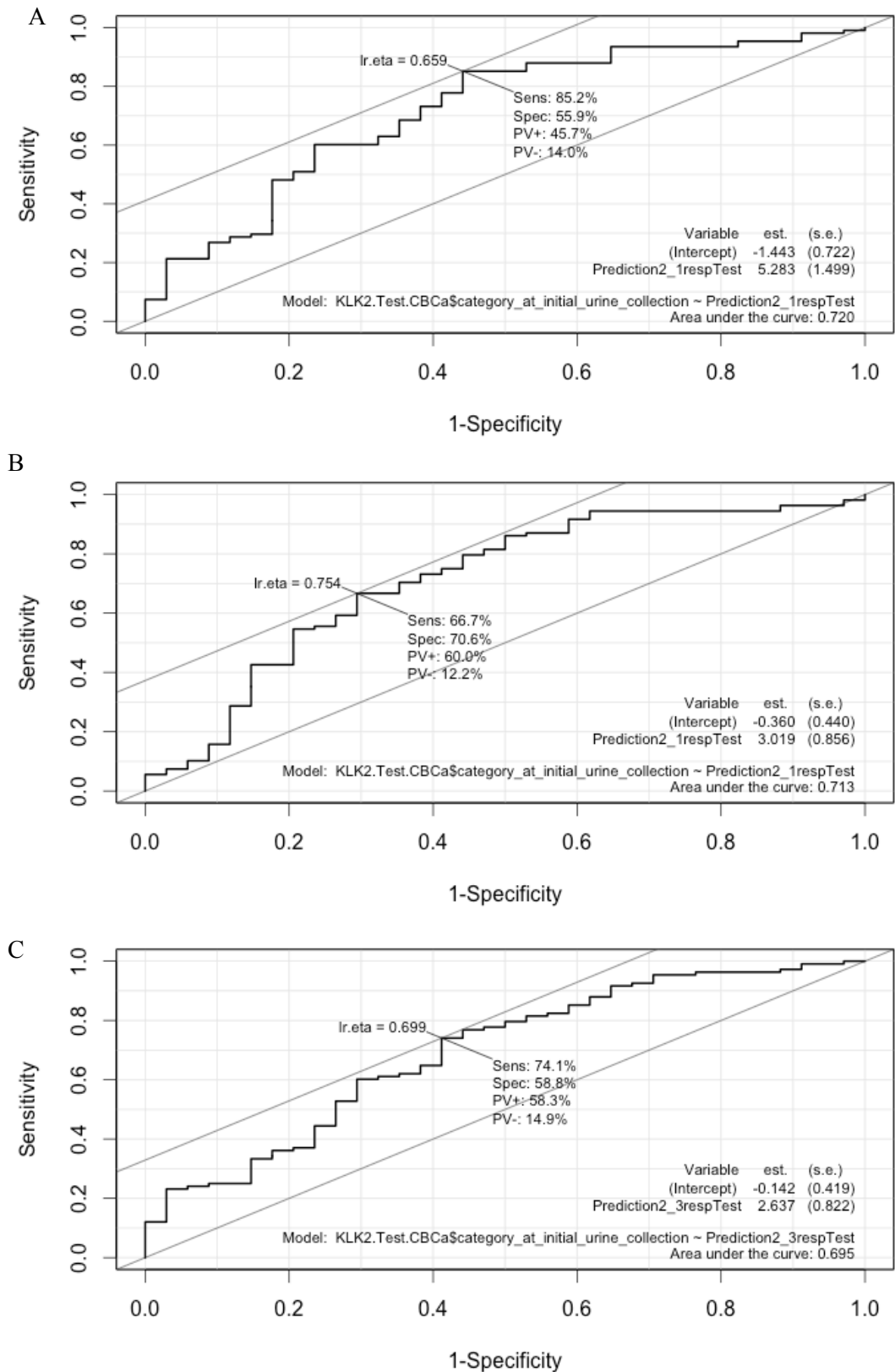
**Supplementary Figure 6** KLK2 Adjusted Data ROC curves for test data using models (random sampling) detecting between CB and Ca (L I H) for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.

9: APPENDICES

Supplementary Table 10 Lasso output for models (random sampling) detecting between CB and Ca (L I H) using KLK2 adjusted data.

<i>All Transcript</i>		<i>Significant Transcripts</i>		<i>Multiple Testing correction Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>SIM2.short</i>	<b>4.95</b>	<i>AMACR</i>	<b>2.32</b>	<i>SMIMI</i>	<b>6.89</b>
<i>SMIMI</i>	<b>2.78</b>	<i>SMIMI</i>	<b>2.20</b>	<i>PCA3</i>	<b>6.59</b>
<i>ERG3' exons 6-7</i>	<b>1.95</b>	<i>MMP11</i>	<b>1.69</b>	<i>SIM2.short</i>	<b>3.42</b>
<i>AMH</i>	<b>1.76</b>	<i>SIM2.short</i>	<b>1.54</b>	<i>AMACR</i>	<b>2.95</b>
<i>HPN</i>	<b>1.26</b>	<i>TMPRSS2:ERG</i>	<b>1.33</b>	<i>HPN</i>	<b>1.88</b>
				<i>ERG3'</i>	
<i>PCA3</i>	<b>1.20</b>	<i>HPN</i>	<b>1.31</b>	<i>exons 4-5</i>	<b>1.18</b>
<i>NEAT1</i>	<b>1.06</b>	<i>ISX</i>	<b>0.93</b>		
<i>PCSK6</i>	<b>1.02</b>	<i>CLU</i>	<b>0.89</b>		
<i>DNAH5</i>	<b>0.68</b>	<i>DLX1</i>	<b>0.46</b>		
<i>TMPRSS2:ERG</i>	<b>0.66</b>	<i>APOC1</i>	<b>0.36</b>		
<i>SEC61A1</i>	<b>0.53</b>	<i>GJB1</i>	<b>0.30</b>		
<i>HIST1H2BF</i>	<b>0.47</b>	<i>CASKIN1</i>	<b>0.22</b>		
<i>CADPS</i>	<b>0.46</b>	<i>MIR146A.DQ658414</i>	<b>0.15</b>		
<i>APOC1</i>	<b>0.45</b>	<i>HOXC6</i>	<b>0.08</b>		
<i>TBP</i>	<b>0.37</b>	<i>PTN</i>	<b>-0.38</b>		
<i>ERG5</i>	<b>0.34</b>	<i>IGFBP3</i>	<b>-0.44</b>		
<i>CAMKK2</i>	<b>0.32</b>	<i>GABARAPL2</i>	<b>-0.69</b>		
<i>CAMK2N2</i>	<b>0.22</b>	<i>MYOF</i>	<b>-1.43</b>		
<i>TMCC2</i>	<b>0.19</b>	<i>KLK4</i>	<b>-1.86</b>		
<i>SERPINB5</i>	<b>0.12</b>				
<i>EN2</i>	<b>0.12</b>				
<i>ERG3' exons 4-5</i>	<b>0.03</b>				
<i>SChLAP1</i>	<b>0.00</b>				
<i>PTN</i>	<b>-0.01</b>				
<i>PPP1R12B</i>	<b>-0.15</b>				
<i>SIRT1</i>	<b>-0.30</b>				
<i>PTPRC</i>	<b>-0.36</b>				
<i>IGFBP3</i>	<b>-0.46</b>				
<i>CD10</i>	<b>-0.61</b>				
<i>SNCA</i>	<b>-0.68</b>				
<i>MEMO1</i>	<b>-0.75</b>				
<i>RPLP2</i>	<b>-1.46</b>				
<i>MYOF</i>	<b>-1.88</b>				
<i>SACMIL</i>	<b>-2.85</b>				
<i>KLK4</i>	<b>-4.13</b>				

9: APPENDICES



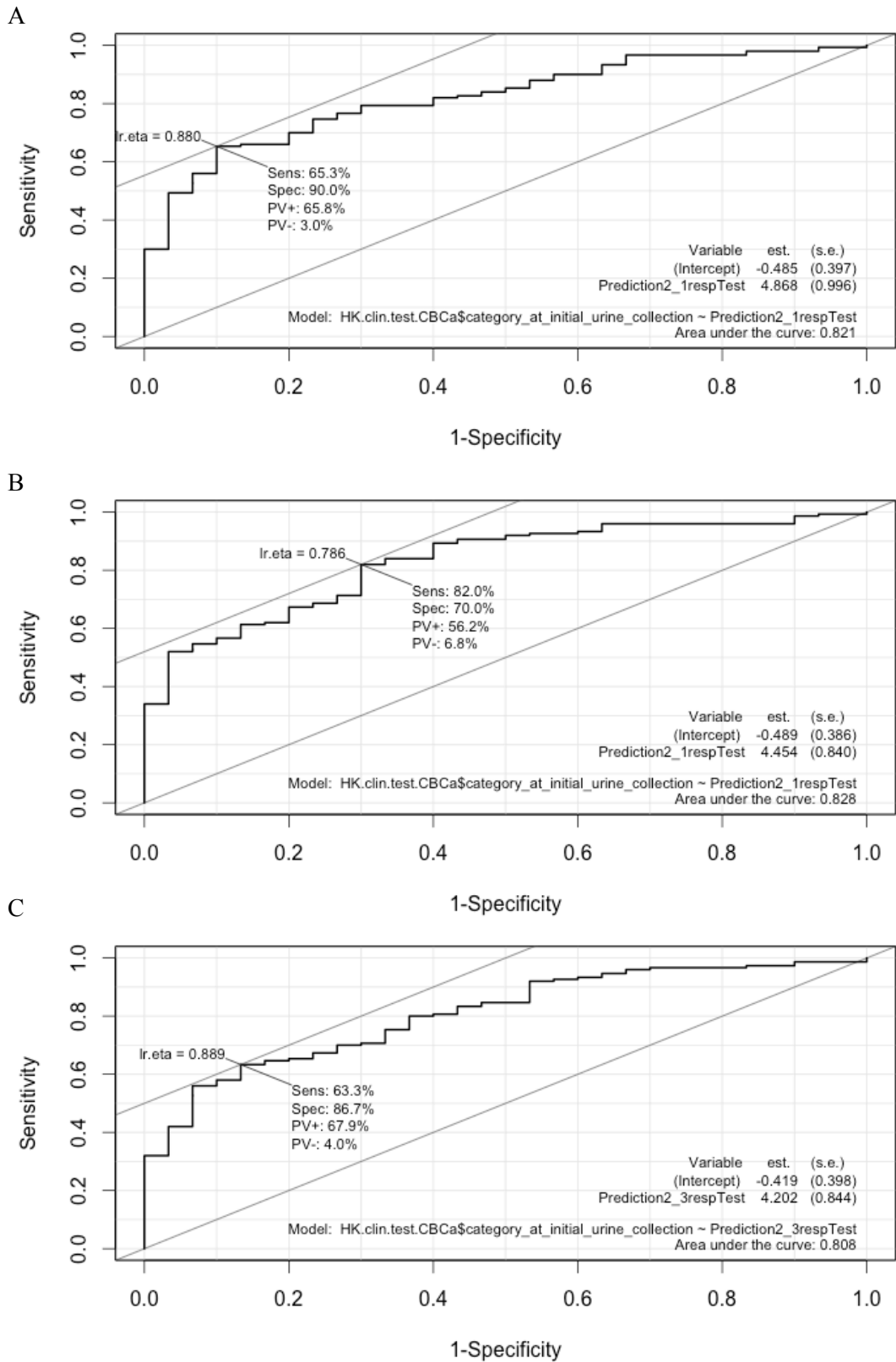
Supplementary Figure 7 KLK3 Adjusted Data ROC curves for test data using models (random sampling) detecting between CB and Ca (L I H) for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.

9: APPENDICES

Supplementary Table 11 Lasso output for models (random sampling) detecting between CB and Ca (L I H) using KLK3 adjusted data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>PCA3</i>	<i>1.57</i>	<i>SMIM1</i>	<i>2.26</i>	<i>HPN</i>	<i>3.49</i>
<i>SMIM1</i>	<i>1.27</i>	<i>SIM2.short</i>	<i>1.84</i>	<i>SIM2.short</i>	<i>3.49</i>
<i>SIM2.short</i>	<i>1.21</i>	<i>SULT1A1</i>	<i>1.48</i>	<i>PCA3</i>	<i>3.33</i>
<i>HPN</i>	<i>0.76</i>	<i>ERG3' exons 4-5</i>	<i>1.42</i>	<i>ERG3' exons 4-5</i>	<i>2.55</i>
<i>ERG3' exons 4-5</i>	<i>0.57</i>	<i>GAPDH</i>	<i>1.37</i>		
<i>GAPDH</i>	<i>0.43</i>	<i>PCA3</i>	<i>1.34</i>		
<i>AMH</i>	<i>0.41</i>	<i>AMACR</i>	<i>0.85</i>		
<i>HOXC6</i>	<i>0.29</i>	<i>HPN</i>	<i>0.69</i>		
<i>CLU</i>	<i>0.16</i>	<i>MMP25</i>	<i>0.64</i>		
<i>ISX</i>	<i>0.15</i>	<i>ERG3' exons 6-7</i>	<i>0.52</i>		
<i>MMP25</i>	<i>0.09</i>	<i>CLU</i>	<i>0.48</i>		
<i>APOC1</i>	<i>0.05</i>	<i>GJB1</i>	<i>0.45</i>		
<i>TMPRSS2:ERG</i>	<i>0.04</i>	<i>ANKRD34B</i>	<i>0.27</i>		
<i>MYOF</i>	<i>-0.08</i>	<i>STOM</i>	<i>0.20</i>		
<i>GABARAPL2</i>	<i>-0.36</i>	<i>RAB17</i>	<i>0.06</i>		
<i>KLK4</i>	<i>-0.45</i>	<i>IGFBP3</i>	<i>-0.33</i>		
		<i>RPS11</i>	<i>-0.48</i>		
		<i>PTN</i>	<i>-0.78</i>		
		<i>MYOF</i>	<i>-1.08</i>		
		<i>GABARAPL2</i>	<i>-2.60</i>		

9: APPENDICES



Supplementary Figure 8 GAPDH and RPLP2 Normalised Data ROC curves for test data using models (random sampling) detecting between CB and Ca (L I H) for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.

9: APPENDICES

Supplementary Table 12 Lasso output for models (random sampling) detecting between CB and Ca (L I H) using HK normalised data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>PCA3</i>	<i>0.33</i>	<i>PCA3</i>	<i>0.63</i>	<i>PCA3</i>	<i>0.39</i>
<i>ERG3' exons 4-5</i>	<i>0.30</i>	<i>TMPRSS2:ERG</i>	<i>0.27</i>	<i>ERG3' exons 4-5</i>	<i>0.22</i>
<i>TMPRSS2:ERG</i>	<i>0.24</i>	<i>SMIM1</i>	<i>0.18</i>	<i>TMPRSS2:ERG</i>	<i>0.17</i>
<i>TDRD</i>	<i>0.04</i>	<i>TDRD</i>	<i>0.02</i>	<i>HOXC6</i>	<i>0.12</i>
<i>CLU</i>	<i>0.01</i>	<i>HOXC6</i>	<i>0.02</i>	<i>TDRD</i>	<i>0.06</i>
<i>MMP25</i>	<i>0.00</i>	<i>ERG5'</i>	<i>0.01</i>		
<i>ALAS1</i>	<i>-0.01</i>	<i>KLK4</i>	<i>-0.21</i>		
<i>PDLIM5</i>	<i>-0.09</i>				



**6.14 Binomial Testing between CB and High-risk Ca**

Supplementary Table 13 Glm test significant probes between CB and High-risk Ca

KLK2 ratio data				KLK2 adjusted data			
Transcript	p-value	Log <sub>2</sub> (FC)	Adjusted p-value	Transcript	p-value	Log <sub>2</sub> (FC)	Adjusted p-value
<i>ERG3' exons 4-5</i>	7.00E-07	1.868	0.0001	<i>HPN</i>	3.77x10-06	0.241	0.001
<i>ERG3' exons 6-7</i>	7.15E-07	2.966	0.0001	<i>PCA3</i>	4.78x10-06	0.222	0.001
<i>PCA3</i>	9.17E-07	0.376	0.0002	<i>GJB1</i>	0.0001	0.159	0.018
<i>APOC1</i>	7.71E-06	1.472	0.001	<i>AMACR</i>	0.0001	0.130	0.021
<i>HPN</i>	8.27E-06	0.352	0.001	<i>KLK4</i>	0.0003	-0.127	0.044
<i>TMPRSS2:ERG</i>	9.85E-06	NA	0.002	<i>ERG3' exons 4-5</i>	0.0004	0.109	0.063
<i>HOXC6</i>	2.14E-05	0.301	0.003	<i>ERG3' exons 6-7</i>	0.001	0.112	0.098
<i>TDRD</i>	2.54E-05	3.689	0.004	<i>TMPRSS2:ERG</i>	0.001	0.114	0.208
<i>DLX1</i>	4.09E-05	4.487	0.006	<i>HOXC6</i>	0.001	0.129	0.221
<i>AMACR</i>	7.47E-05	0.341	0.012	<i>RAB17</i>	0.002	0.227	0.272
<i>GJB1</i>	9.73E-05	0.320	0.015	<i>APOC1</i>	0.002	0.209	0.372
<i>ANKRD34B</i>	0.0002	5.892	0.025	<i>DLX1</i>	0.002	0.074	0.372
<i>TRPM4</i>	0.0002	0.730	0.029	<i>SPINK1</i>	0.003	0.163	0.445
<i>MCTP1</i>	0.0003	1.149	0.041	<i>MYOF</i>	0.003	-0.155	0.463
<i>PPFIA2</i>	0.0003	0.807	0.041	<i>SULT1A1</i>	0.003	0.126	0.509
<i>ITGBL1</i>	0.0003	0.799	0.042	<i>DPP4</i>	0.004	-0.105	0.552
<i>HOXC4</i>	0.0004	0.938	0.063	<i>ITGBL1</i>	0.004	0.087	0.611
<i>SLC12A1</i>	0.0004	1.022	0.064	<i>AR exons 4-8</i>	0.004	-0.121	0.637
<i>ISX</i>	0.001	2.371	0.077	<i>TRPM4</i>	0.004	0.080	0.640
<i>RAB17</i>	0.001	0.303	0.097	<i>CD10</i>	0.005	-0.092	0.771
<i>VPS13A</i>	0.001	0.110	0.118	<i>GABARAPL2</i>	0.006	-0.137	0.863
<i>NEAT1</i>	0.001	0.186	0.131	<i>RP11.244H18.1.P712P</i>	0.006	-0.100	0.890
<i>STOM</i>	0.001	NA	0.134	<i>TDRD</i>	0.007	0.064	0.983
<i>PVT1</i>	0.001	0.383	0.208	<i>UPK2</i>	0.007	0.108	0.996
<i>SSTR1</i>	0.001	0.369	0.209	<i>SLC12A1</i>	0.007	0.148	0.999

## 9: APPENDICES

<i>Met</i>	0.002	2.490	0.266	<i>MIR4435.1HG.IOC541471</i>	0.008	0.087	0.999
<i>SIM2.short</i>	0.002	0.507	0.274	<i>GAPDH</i>	0.011	0.111	0.999
<i>CDKN3</i>	0.002	0.380	0.275	<i>RP11.97O12.7</i>	0.011	0.116	0.999
<i>ERG5</i>	0.002	NA	0.276	<i>STOM</i>	0.011	0.116	0.999
<i>SPINK1</i>	0.002	0.265	0.294	<i>SMIM1</i>	0.012	0.111	0.999
<i>SULT1A1</i>	0.002	0.225	0.323	<i>ANKRD34B</i>	0.012	0.134	0.999
<i>TMEM45B</i>	0.002	0.410	0.325	<i>NEAT1</i>	0.012	0.094	0.999
<i>UPK2</i>	0.002	0.676	0.334	<i>SIM2.short</i>	0.019	0.077	0.999
<i>AMH</i>	0.003	0.404	0.348	<i>MCTP1</i>	0.024	0.094	0.999
<i>MIR146A.DQ658414</i>	0.003	0.549	0.352	<i>MED4</i>	0.029	-0.054	0.999
<i>SULF2</i>	0.003	2.000	0.352	<i>DNAH5</i>	0.029	0.048	0.999
<i>RP11.97O12.7</i>	0.003	0.086	0.355	<i>ISX</i>	0.029	0.081	0.999
<i>MMP11</i>	0.003	0.451	0.382	<i>PPFIA2</i>	0.031	0.072	0.999
<i>TMCC2</i>	0.003	4.761	0.436	<i>Met</i>	0.035	0.098	0.999
<i>PALM3</i>	0.004	0.269	0.474	<i>SNCA</i>	0.038	-0.053	0.999
<i>MIR4435.1HG.IOC541471</i>	0.005	0.199	0.605	<i>VPS13A</i>	0.041	0.040	0.999
<i>MIC1</i>	0.005	0.445	0.642	<i>PTN</i>	0.044	-0.082	0.999
<i>LASS1</i>	0.005	0.473	0.665	<i>PVT1</i>	0.049	0.113	0.999
<i>RIOK3</i>	0.005	0.104	0.676				
<i>MEX3A</i>	0.006	0.801	0.694				
<i>RPL23AP53</i>	0.006	0.417	0.758				
<i>CASKIN1</i>	0.007	0.178	0.872				
<i>TWIST1</i>	0.008	0.192	0.893				
<i>IMPDH2</i>	0.008	0.099	0.894				
<i>SIM2.long</i>	0.009	0.133	0.966				
<i>PECI</i>	0.009	0.063	0.966				
<i>GAPDH</i>	0.009	0.067	0.966				
<i>DNAH5</i>	0.009	0.409	0.966				
<i>EN2</i>	0.009	0.358	0.966				
<i>MKI67</i>	0.010	-1.667	0.966				
<i>NAALADL2</i>	0.010	0.081	0.966				
<i>SMIM1</i>	0.010	0.137	0.966				
<i>MMP26</i>	0.011	0.411	0.966				
<i>MXN1</i>	0.011	0.241	0.966				

## 9: APPENDICES

<i>MMP25</i>	0.012	0.932	0.966
<i>HIST1H1C</i>	0.012	0.055	0.966
<i>SChLAP1</i>	0.013	0.525	0.966
<i>MGAT5B</i>	0.013	0.258	0.966
<i>PCSK6</i>	0.014	0.219	0.966
<i>CLIC2</i>	0.014	NA	0.966
<i>MCM7</i>	0.015	0.273	0.966
<i>MFSD2A</i>	0.016	-2.398	0.966
<i>TERT</i>	0.017	0.153	0.966
<i>HPRT</i>	0.017	0.073	0.966
<i>SSPO</i>	0.017	0.221	0.966
<i>HIST3H2A</i>	0.020	0.091	0.966
<i>ITPR1</i>	0.022	0.069	0.966
<i>B4GALNT4</i>	0.022	NA	0.966
<i>SLC4A1.S</i>	0.022	NA	0.966
<i>RPLP2</i>	0.023	0.053	0.966
<i>SACM1L</i>	0.025	0.058	0.966
<i>SYNM</i>	0.025	0.214	0.966
<i>VAX2</i>	0.026	0.270	0.966
<i>TMEM86A</i>	0.026	0.795	0.966
<i>RPS11</i>	0.027	0.035	0.966
<i>ABCB9</i>	0.028	NA	0.966
<i>CLU</i>	0.030	0.248	0.966
<i>CCDC88B</i>	0.030	-5.601	0.966
<i>HIST1H2BG</i>	0.032	0.124	0.966
<i>FOLH1</i>	0.032	0.063	0.966
<i>COL9A2</i>	0.034	-2.208	0.966
<i>BRAF</i>	0.035	0.072	0.966
<i>RPL18A</i>	0.035	0.045	0.966
<i>CAMKK2</i>	0.036	0.085	0.966
<i>AURKA</i>	0.036	0.428	0.966
<i>ARHGEF25</i>	0.036	0.278	0.966
<i>ALAS1</i>	0.037	0.021	0.966
<i>SFRP4</i>	0.039	0.502	0.966

9: APPENDICES

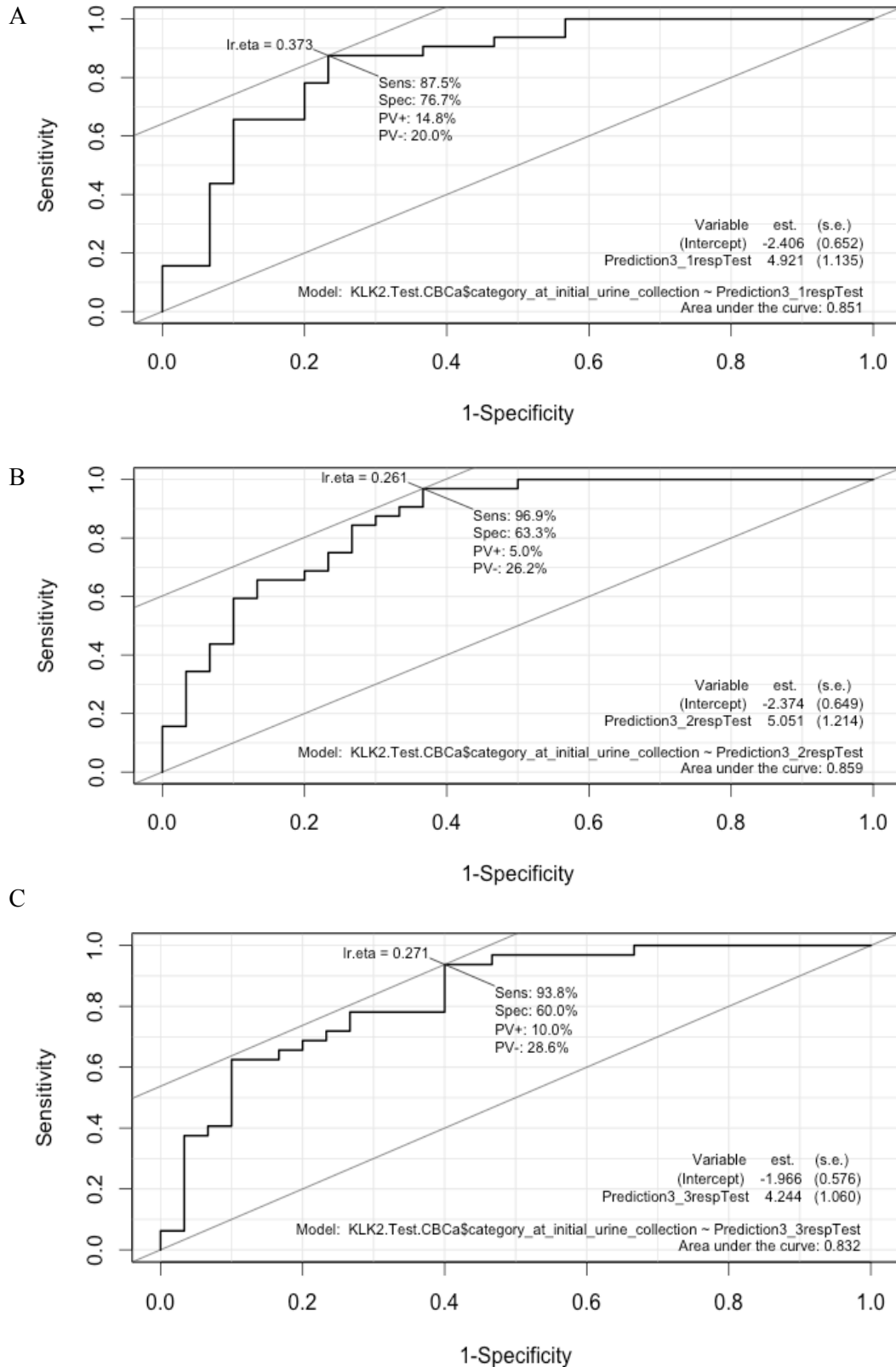
<i>TERF2IP</i>	0.041	0.031	0.966
<i>PTPRC</i>	0.046	NA	0.966
<i>COL10A1</i>	0.047	0.274	0.966
<i>ACTR5</i>	0.049	0.194	0.966
<i>PSTPIP1</i>	0.050	NA	0.966

KLK3 adjusted data				GAPDH and RPLP2 Normalised data			
Transcript	p-value	Log <sub>2</sub> (FC)	Adjusted p-value	Transcript	p-value	Log <sub>2</sub> (FC)	Adjusted p-value
<i>HPN</i>	1.32E-05	0.190	0.002	<i>ERG3' exons 4-5</i>	1.43E-06	0.793	0.000
<i>PCA3</i>	1.45E-05	0.184	0.002	<i>PCA3</i>	6.29E-06	0.196	0.001
<i>KLK4</i>	0.0002	-0.159	0.034	<i>TMPRSS2:ERG</i>	1.05E-05	0.953	0.002
<i>ERG3' exons 4-5</i>	0.0004	0.125	0.066	<i>ERG3' exons 6-7</i>	1.40E-05	1.265	0.002
<i>GJB1</i>	0.0005	0.142	0.075	<i>APOC1</i>	2.12E-05	0.505	0.003
<i>AMACR</i>	0.001	0.121	0.105	<i>HPN</i>	8.13E-05	0.149	0.013
<i>ERG3' exons 6-7</i>	0.001	0.088	0.106	<i>KLK4</i>	0.0005	-0.069	0.074
<i>MYOF</i>	0.001	-0.128	0.126	<i>HOXC6</i>	0.001	0.170	0.125
<i>TMPRSS2:ERG</i>	0.001	0.084	0.156	<i>TDRD</i>	0.002	0.800	0.245
<i>ARexons4.8</i>	0.001	-0.072	0.163	<i>SLC12A1</i>	0.002	0.371	0.273
<i>HOXC6</i>	0.002	0.120	0.273	<i>DLX1</i>	0.002	0.391	0.279
<i>RP11.244H18.1.P712P</i>	0.002	-0.096	0.300	<i>ITGBL1</i>	0.002	0.144	0.383
<i>DPP4</i>	0.002	-0.107	0.329	<i>MYOF</i>	0.005	-0.055	0.758
<i>APOC1</i>	0.002	0.162	0.347	<i>DPP4</i>	0.005	-0.039	0.762
<i>DLX1</i>	0.003	0.057	0.525	<i>SPINK1</i>	0.005	0.126	0.808
<i>SULT1A1</i>	0.004	0.110	0.532	<i>GABARAPL2</i>	0.005	-0.050	0.821
<i>SPINK1</i>	0.004	0.131	0.565	<i>RAB17</i>	0.005	0.125	0.826
<i>ITGBL1</i>	0.005	0.108	0.676	<i>CD10</i>	0.006	-0.074	0.967
<i>RAB17</i>	0.005	0.142	0.676	<i>HOXC4</i>	0.008	0.160	0.995
<i>CD10</i>	0.006	-0.095	0.856	<i>AR exons 4-8</i>	0.008	-0.059	0.995
<i>KLK2</i>	0.007	-0.080	0.952	<i>NEAT1</i>	0.010	0.126	0.995
<i>GABARAPL2</i>	0.007	-0.143	0.974	<i>UPK2</i>	0.010	0.306	0.995
<i>SLC12A1</i>	0.007	0.092	0.985	<i>PPFIA2</i>	0.011	0.325	0.995
<i>UPK2</i>	0.008	0.098	0.998	<i>GJB1</i>	0.012	0.127	0.995
<i>STOM</i>	0.012	0.089	0.998	<i>SRSF3</i>	0.014	-0.178	0.995

## 9: APPENDICES

<i>PTN</i>	0.014	-0.104	0.998	<i>MCTP1</i>	0.014	0.370	0.995
<i>MIR4435.1HG.IOC541471</i>	0.014	0.122	0.998	<i>Met</i>	0.016	0.924	0.995
<i>MED4</i>	0.015	-0.061	0.998	<i>KLK2</i>	0.016	-0.043	0.995
<i>TDRD</i>	0.016	0.040	0.998	<i>AMACR</i>	0.016	0.140	0.995
<i>SNCA</i>	0.020	-0.058	0.998	<i>ANKRD34B</i>	0.018	0.083	0.995
<i>TRPM4</i>	0.022	0.053	0.998	<i>STOM</i>	0.023	0.180	0.995
<i>NEAT1</i>	0.027	0.059	0.998	<i>AR.ex9</i>	0.024	-0.448	0.995
<i>MARCH5</i>	0.030	-0.077	0.998	<i>MXI1</i>	0.025	-0.043	0.995
<i>ANKRD34B</i>	0.032	0.082	0.998	<i>P712P</i>	0.026	-0.054	0.995
<i>MEMO1</i>	0.032	-0.085	0.998	<i>STEAP2</i>	0.028	-0.032	0.995
<i>SMIM1</i>	0.035	0.103	0.998	<i>SULT1A1</i>	0.029	0.078	0.995
<i>SIM2.short</i>	0.039	0.028	0.998	<i>PDLIM5</i>	0.030	-0.042	0.995
<i>RP11.97O12.7</i>	0.040	0.063	0.998	<i>PTN</i>	0.031	-0.101	0.995
<i>SRSF3</i>	0.040	-0.094	0.998	<i>TRPM4</i>	0.032	0.219	0.995

9: APPENDICES



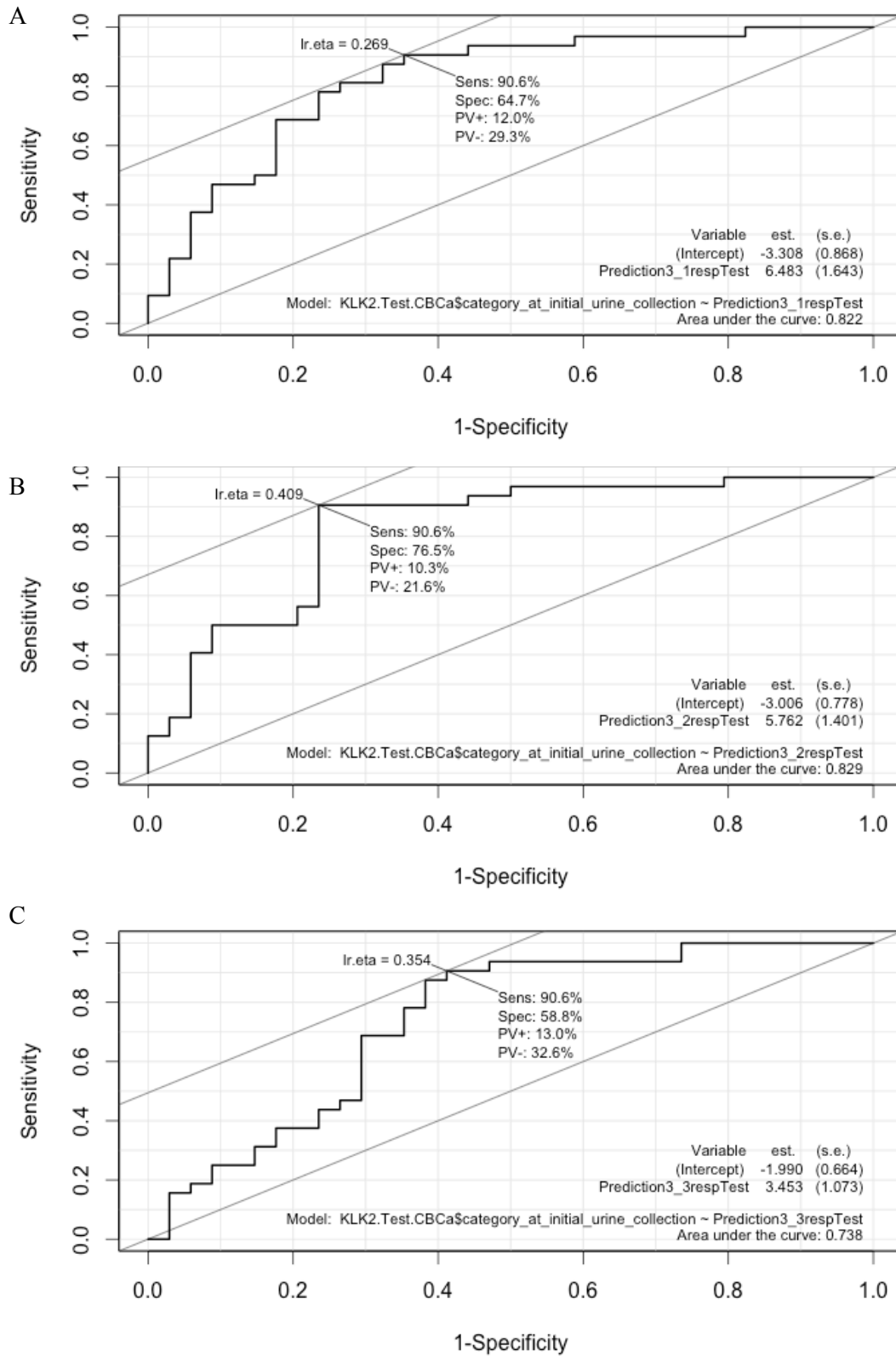
Supplementary Figure 9 KLK2 Ratio Data ROC curves for test data using models detecting between CB and high risk Ca for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.

## 9: APPENDICES

Supplementary Table 14 Lasso output for models detecting between CB and high risk Ca using KLK2 ratio data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>ERG3' exons 4-5</i>	<i>0.55</i>	<i>ERG3' exons 4-5</i>	<i>0.49</i>	<i>ERG3' exons 4-5</i>	<i>0.44</i>
<i>PCA3</i>	<i>0.22</i>	<i>APOC1</i>	<i>0.23</i>	<i>PCA3</i>	<i>0.20</i>
<i>ANKRD34B</i>	<i>0.17</i>	<i>PCA3</i>	<i>0.19</i>	<i>APOC1</i>	<i>0.19</i>
<i>APOC1</i>	<i>0.16</i>	<i>AMACR</i>	<i>0.15</i>	<i>AMACR</i>	<i>0.16</i>
<i>AMACR</i>	<i>0.14</i>	<i>HOXC6</i>	<i>0.10</i>	<i>TMPRSS2:ERG</i>	<i>0.14</i>
<i>HOXC6</i>	<i>0.09</i>	<i>TMPRSS2:ERG</i>	<i>0.10</i>	<i>HOXC6</i>	<i>0.12</i>
<i>TMPRSS2:ERG</i>	<i>0.08</i>	<i>ANKRD34B</i>	<i>0.07</i>	<i>ANKRD34B</i>	<i>0.05</i>
<i>TMEM47</i>	<i>0.07</i>	<i>HPN</i>	<i>0.06</i>	<i>DLX1</i>	<i>0.03</i>
<i>MMP25</i>	<i>0.05</i>	<i>NEAT1</i>	<i>0.04</i>	<i>PPF1A2</i>	<i>-0.01</i>
<i>DLX1</i>	<i>0.03</i>	<i>DLX1</i>	<i>0.03</i>		
<i>NEAT1</i>	<i>0.03</i>	<i>AURKA</i>	<i>-0.02</i>		
<i>ISX</i>	<i>0.01</i>	<i>PTPRC</i>	<i>-0.03</i>		
<i>MAK</i>	<i>0.00</i>	<i>ALAS1</i>	<i>-0.05</i>		
<i>MED4</i>	<i>-0.02</i>	<i>PSTPIP1</i>	<i>-0.06</i>		
<i>CP</i>	<i>-0.02</i>	<i>ACTR5</i>	<i>-0.14</i>		
<i>CKAP2L</i>	<i>-0.02</i>	<i>RPL18A</i>	<i>-0.22</i>		
<i>IGFBP3</i>	<i>-0.02</i>				
<i>AR exon 9</i>	<i>-0.03</i>				
<i>SRSF3</i>	<i>-0.04</i>				
<i>PDLIM5</i>	<i>-0.07</i>				
<i>BTG2</i>	<i>-0.07</i>				
<i>STEAP4</i>	<i>-0.08</i>				
<i>CD10</i>	<i>-0.14</i>				
<i>AR exons 4-8</i>	<i>-0.17</i>				
<i>KLK4</i>	<i>-0.27</i>				
<i>DPP4</i>	<i>-0.29</i>				

9: APPENDICES



Supplementary Figure 10 KLK2 Adjusted Data ROC curves for test data using models detecting between CB and high risk Ca for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.

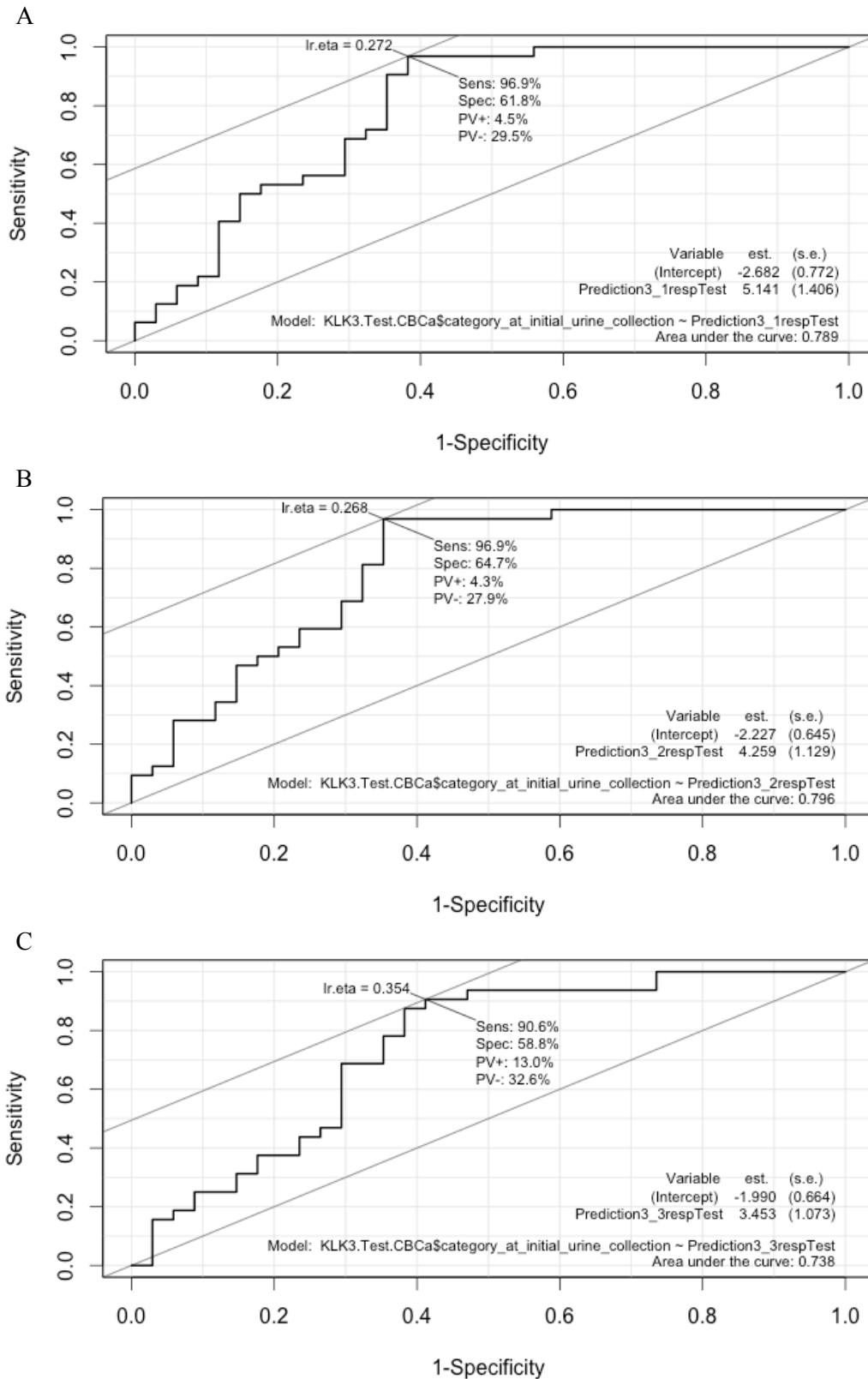


9: APPENDICES

Supplementary Table 15 Lasso output for models detecting between CB and high risk Ca using KLK2 adjusted data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>PCA3</i>	<b>3.08</b>	<i>PCA3</i>	<b>3.34</b>	<i>PCA3</i>	<b>5.34</b>
<i>HPN</i>	<b>2.79</b>	<i>HPN</i>	<b>2.40</b>	<i>HPN</i>	<b>5.02</b>
<i>AMACR</i>	<b>1.09</b>	<i>AMACR</i>	<b>1.93</b>	<i>GJB1</i>	<b>2.32</b>
<i>ERG3' exons 6-7</i>	<b>0.83</b>	<i>SIM2.short</i>	<b>1.45</b>	<i>AMACR</i>	<b>1.85</b>
<i>SIM2.short</i>	<b>0.72</b>	<i>DNAH5</i>	<b>1.03</b>	<i>KLK4</i>	<b>-2.59</b>
<i>RAB17</i>	<b>0.37</b>	<i>ERG3' exons 6-7</i>	<b>0.94</b>		
<i>APOC1</i>	<b>0.34</b>	<i>RAB17</i>	<b>0.50</b>		
<i>MMP25</i>	<b>0.34</b>	<i>ANKRD34B</i>	<b>0.45</b>		
<i>ANKRD34B</i>	<b>0.27</b>	<i>APOC1</i>	<b>0.44</b>		
<i>DLX1</i>	<b>0.25</b>	<i>DLX1</i>	<b>0.43</b>		
<i>CLU</i>	<b>0.23</b>	<i>SLC12A1</i>	<b>0.38</b>		
<i>DNAH5</i>	<b>0.16</b>	<i>STOM</i>	<b>0.16</b>		
<i>SLC12A1</i>	<b>0.14</b>	<i>ERG3' exons 4-5</i>	<b>0.09</b>		
<i>ERG3' exons 4-5</i>	<b>0.03</b>	<i>KLK4</i>	<b>-0.46</b>		
<i>STOM</i>	<b>0.03</b>	<i>MYOF</i>	<b>-0.71</b>		
<i>MYOF</i>	<b>-0.29</b>	<i>DPP4</i>	<b>-1.37</b>		
<i>KLK4</i>	<b>-0.40</b>	<i>CD10</i>	<b>-2.13</b>		
<i>DPP4</i>	<b>-0.98</b>				
<i>CD10</i>	<b>-1.43</b>				

9: APPENDICES



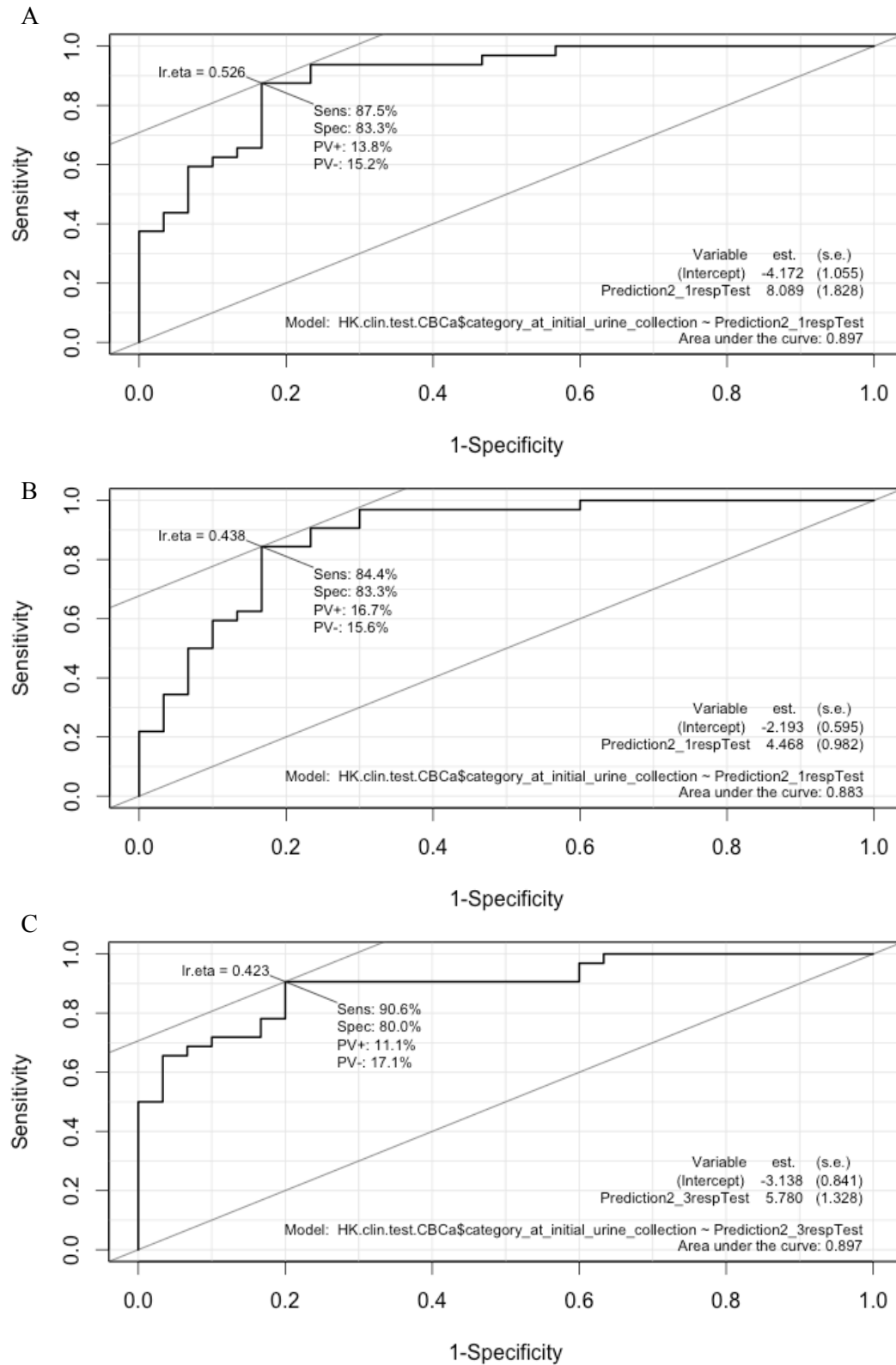
Supplementary Figure 11 KLK3 Adjusted Data ROC curves for test data using models detecting between CB and high risk Ca for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.

9: APPENDICES

Supplementary Table 16 Lasso output for models detecting between CB and high risk Ca using KLK3 adjusted data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>PCA3</i>	<b>2.56</b>	<i>PCA3</i>	<b>2.96</b>	<i>HPN</i>	<b>5.06</b>
<i>HPN</i>	<b>2.37</b>	<i>HPN</i>	<b>2.13</b>	<i>PCA3</i>	<b>4.95</b>
<i>ERG3' exons 4-5</i>	<b>0.77</b>	<i>SIM2.short</i>	<b>1.46</b>	<i>KLK4</i>	<b>-3.20</b>
<i>SIM2.short</i>	<b>0.64</b>	<i>ERG3' exons 4-5</i>	<b>1.17</b>		
<i>APOC1</i>	<b>0.44</b>	<i>AMACR</i>	<b>1.10</b>		
<i>MMP25</i>	<b>0.41</b>	<i>APOC1</i>	<b>0.64</b>		
<i>AMACR</i>	<b>0.30</b>	<i>ANKRD34B</i>	<b>0.57</b>		
<i>ANKRD34B</i>	<b>0.29</b>	<i>SLC12A1</i>	<b>0.43</b>		
<i>SLC12A1</i>	<b>0.18</b>	<i>SULT1A1</i>	<b>0.41</b>		
<i>ERG3' exons 6-7</i>	<b>0.17</b>	<i>DLX1</i>	<b>0.30</b>		
<i>DLX1</i>	<b>0.13</b>	<i>RAB17</i>	<b>0.27</b>		
<i>RAB17</i>	<b>0.11</b>	<i>STOM</i>	<b>0.17</b>		
<i>SULT1A1</i>	<b>0.07</b>	<i>ERG3' exons 6-7</i>	<b>0.16</b>		
<i>STOM</i>	<b>0.06</b>	<i>PTN</i>	<b>-0.14</b>		
<i>PTN</i>	<b>-0.01</b>	<i>RP11.244H18.1.P712P</i>	<b>-0.37</b>		
<i>KLK4</i>	<b>-0.18</b>	<i>MARCH5</i>	<b>-0.72</b>		
<i>MARCH5</i>	<b>-0.23</b>	<i>MYOF</i>	<b>-1.00</b>		
<i>RP11.244H18.1.P712P</i>	<b>-0.39</b>	<i>DPP4</i>	<b>-1.25</b>		
<i>MYOF</i>	<b>-0.65</b>	<i>CD10</i>	<b>-1.68</b>		
<i>DPP4</i>	<b>-0.82</b>				
<i>CD10</i>	<b>-1.17</b>				

9: APPENDICES



**Supplementary Figure 12 HK Normalised Data ROC curves for test data using models detecting between CB and high risk Ca for models using the following inputs A) all probes, B) significant probes, C) adjusted significant probes.**

9: APPENDICES

Supplementary Table 17 Lasso output for models detecting between CB and high risk Ca using HK normalised data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>ERG3' exons 4-5</i>	<b>0.35</b>	<i>ERG3' exons 4-5</i>	<b>0.67</b>	<i>PCA3</i>	<b>0.50</b>
<i>PCA3</i>	<b>0.27</b>	<i>ANKRD34B</i>	<b>0.34</b>	<i>APOC1</i>	<b>0.39</b>
<i>APOC1</i>	<b>0.18</b>	<i>PCA3</i>	<b>0.33</b>	<i>HPN</i>	<b>0.37</b>
<i>HPN</i>	<b>0.13</b>	<i>APOC1</i>	<b>0.24</b>	<i>ERG3' exons 4-5</i>	<b>0.23</b>
<i>SLC12A1</i>	<b>0.03</b>	<i>AMACR</i>	<b>0.19</b>	<i>TMPRSS2:ERG</i>	<b>0.20</b>
<i>TMPRSS2:ERG</i>	<b>0.02</b>	<i>HPN</i>	<b>0.12</b>		
<i>ANKRD34B</i>	<b>0.02</b>	<i>SULT1A1</i>	<b>0.08</b>		
<i>HOXC6</i>	<b>0.01</b>	<i>NEATI</i>	<b>0.08</b>		
<i>AR exons 4-8</i>	<b>-0.05</b>	<i>TMPRSS2:ERG</i>	<b>0.07</b>		
<i>GABARAPL2</i>	<b>-0.06</b>	<i>DLX1</i>	<b>0.04</b>		
<i>CD10</i>	<b>-0.14</b>	<i>HOXC6</i>	<b>0.03</b>		
<i>KLK4</i>	<b>-0.15</b>	<i>STOM</i>	<b>0.02</b>		
<i>DPP4</i>	<b>-0.18</b>	<i>SLC12A1</i>	<b>0.01</b>		
		<i>AR exon 9</i>	<b>-0.03</b>		
		<i>MYOF</i>	<b>-0.06</b>		
		<i>SRSF3</i>	<b>-0.07</b>		
		<i>AR exons 4-8</i>	<b>-0.20</b>		
		<i>CD10</i>	<b>-0.26</b>		
		<i>GABARAPL2</i>	<b>-0.28</b>		
		<i>DPP4</i>	<b>-0.36</b>		
		<i>PDLIM5</i>	<b>-0.56</b>		

9: APPENDICES

**6.15 Multinomial CBLIH Trend**

**Supplementary Table 18 Glm test significant probes for CB, L, I, H trend**

<i>KLK2</i> ratio data			<i>KLK2</i> adjusted data		
<b>Transcript</b>	<b><i>p</i>-value</b>	<b>Adjusted <i>p</i>-value</b>	<b>Transcript</b>	<b><i>p</i>-value</b>	<b>Adjusted <i>p</i>-value</b>
<i>ERG3' exons 4-5</i>	1.86x10 <sup>-13</sup>	3.09x10 <sup>-11</sup>	<i>PCA3</i>	1.45x10 <sup>-08</sup>	2.41x10 <sup>-06</sup>
<i>PCA3</i>	8.90x10 <sup>-13</sup>	1.47x10 <sup>-10</sup>	<i>ERG3' exons 4-5</i>	1.18x10 <sup>-07</sup>	1.94x10 <sup>-05</sup>
<i>TMPRSS2:ERG</i>	1.88x10 <sup>-11</sup>	3.09x10 <sup>-09</sup>	<i>ERG3' exons 6-7</i>	3.73x10 <sup>-07</sup>	6.11x10 <sup>-05</sup>
<i>ERG3' exons 6-7</i>	6.66x10 <sup>-10</sup>	1.09x10 <sup>-07</sup>	<i>SPINK1</i>	1.03x10 <sup>-06</sup>	0.0002
<i>HOXC6</i>	7.88x10 <sup>-09</sup>	1.28x10 <sup>-06</sup>	<i>HOXC6</i>	3.85x10 <sup>-06</sup>	0.0006
<i>HPN</i>	4.19x10 <sup>-08</sup>	6.75x10 <sup>-06</sup>	<i>HPN</i>	7.34x10 <sup>-06</sup>	0.0012
<i>APOC1</i>	6.38x10 <sup>-08</sup>	1.02x10 <sup>-05</sup>	<i>TMPRSS2:ERG</i>	7.35x10 <sup>-06</sup>	0.0012
<i>TDRD</i>	1.63x10 <sup>-07</sup>	2.59x10 <sup>-05</sup>	<i>KLK4</i>	3.65x10 <sup>-05</sup>	0.0058
<i>ANKRD34B</i>	1.47x10 <sup>-06</sup>	0.0002	<i>SLC12A1</i>	4.71x10 <sup>-05</sup>	0.0074
<i>ITGBL1</i>	3.19x10 <sup>-06</sup>	0.001	<i>UPK2</i>	8.47x10 <sup>-05</sup>	0.0133
<i>SLC12A1</i>	5.82x10 <sup>-06</sup>	0.001	<i>TDRD</i>	0.0002	0.0242
<i>DLX1</i>	7.26x10 <sup>-06</sup>	0.001	<i>ITGBL1</i>	0.0002	0.0281
<i>RAB17</i>	9.26x10 <sup>-06</sup>	0.001	<i>RP11.244H18.1.P712P</i>	0.0002	0.0320
<i>HOXC4</i>	1.07x10 <sup>-05</sup>	0.002	<i>GABARAPL2</i>	0.0002	0.0363
<i>GJB1</i>	1.49x10 <sup>-05</sup>	0.002	<i>GJB1</i>	0.0002	0.0366
<i>PPFIA2</i>	1.84x10 <sup>-05</sup>	0.003	<i>AMACR</i>	0.0005	0.0695
<i>SPINK1</i>	2.62x10 <sup>-05</sup>	0.004	<i>MYOF</i>	0.0005	0.0775
<i>AMACR</i>	3.58x10 <sup>-05</sup>	0.005	<i>APOC1</i>	0.0005	0.0816
<i>AMH</i>	4.60x10 <sup>-05</sup>	0.007	<i>MED4</i>	0.0010	0.1423
<i>TRPM4</i>	5.74x10 <sup>-05</sup>	0.008	<i>SULT1A1</i>	0.0011	0.1624
<i>NEAT1</i>	6.14x10 <sup>-05</sup>	0.009	<i>RAB17</i>	0.0020	0.2865
<i>SIM2.short</i>	6.84x10 <sup>-05</sup>	0.010	<i>ANKRD34B</i>	0.0025	0.3583
<i>SSTR1</i>	7.31x10 <sup>-05</sup>	0.011	<i>SNCA</i>	0.0035	0.4989
<i>UPK2</i>	7.83x10 <sup>-05</sup>	0.011	<i>MMP26</i>	0.0047	0.6782
<i>SULT1A1</i>	8.22x10 <sup>-05</sup>	0.012	<i>PTN</i>	0.0055	0.7804
<i>MEX3A</i>	9.10x10 <sup>-05</sup>	0.013	<i>DLX1</i>	0.0055	0.7804

## 9: APPENDICES

<i>MIR146A.DQ658414</i>	0.0001	0.015	<i>IFT57</i>	0.0058	0.8138
<i>TMEM45B</i>	0.0001	0.015	<i>SIM2.short</i>	0.0061	0.8487
<i>ISX</i>	0.0001	0.017	<i>DPP4</i>	0.0073	0.9916
<i>MIC1</i>	0.0001	0.019	<i>STOM</i>	0.0080	0.9916
<i>TWIST1</i>	0.0002	0.021	<i>GAPDH</i>	0.0105	0.9916
<i>Met</i>	0.0002	0.021	<i>VPS13A</i>	0.0135	0.9916
<i>MMP11</i>	0.0002	0.023	<i>MIR146A.DQ658414</i>	0.0168	0.9916
<i>CDKN3</i>	0.0002	0.023	<i>PPAP2A</i>	0.0182	0.9916
<i>RP11.97012.7</i>	0.0002	0.023	<i>ZNF577</i>	0.0185	0.9916
<i>STOM</i>	0.0003	0.035	<i>SMIM1</i>	0.0233	0.9916
<i>PALM3</i>	0.0003	0.043	<i>PPFIA2</i>	0.0249	0.9916
<i>LASS1</i>	0.0003	0.043	<i>Met</i>	0.0251	0.9916
<i>SSPO</i>	0.0003	0.044	<i>MIC1</i>	0.0268	0.9916
<i>MMP26</i>	0.000	0.049	<i>EIF2D</i>	0.0316	0.9916
<i>VPS13A</i>	0.000	0.049	<i>CD10</i>	0.0336	0.9916
<i>PECI</i>	0.000	0.050	<i>STEAP2</i>	0.0432	0.9916
<i>PCSK6</i>	0.000	0.054	<i>MIR4435.1HG.IOC541471</i>	0.0437	0.9916
<i>GAPDH</i>	0.000	0.056	<i>ITPR1</i>	0.0445	0.9916
<i>PVT1</i>	0.000	0.056	<i>MXI1</i>	0.0487	0.9916
<i>TERT</i>	0.000	0.060			
<i>CASKIN1</i>	0.001	0.061			
<i>TMCC2</i>	0.001	0.064			
<i>RPLP2</i>	0.001	0.080			
<i>MXN1</i>	0.001	0.102			
<i>SIM2.long</i>	0.001	0.106			
<i>RPS11</i>	0.001	0.106			
<i>SULF2</i>	0.001	0.130			
<i>HIST1H1C</i>	0.001	0.133			
<i>EN2</i>	0.001	0.133			
<i>DNAH5</i>	0.001	0.166			
<i>MMP25</i>	0.002	0.198			
<i>MFSD2A</i>	0.002	0.212			
<i>MIR4435.1HG.IOC541471</i>	0.002	0.226			
<i>SMIM1</i>	0.002	0.239			

## 9: APPENDICES

<i>MGAT5B</i>	0.003	0.266
<i>RIOK3</i>	0.003	0.267
<i>MCTP1</i>	0.003	0.315
<i>RPS10</i>	0.003	0.329
<i>VAX2</i>	0.003	0.337
<i>TMEM86A</i>	0.003	0.340
<i>ERG5</i>	0.004	0.358
<i>IMPDH2</i>	0.004	0.368
<i>COL10A1</i>	0.004	0.400
<i>ABCB9</i>	0.004	0.424
<i>B4GALNT4</i>	0.005	0.471
<i>Mki67</i>	0.005	0.472
<i>CLIC2</i>	0.006	0.526
<i>SChLAP1</i>	0.007	0.671
<i>CCDC88B</i>	0.009	0.807
<i>PTPRC</i>	0.009	0.809
<i>CAMKK2</i>	0.009	0.838
<i>NAALADL2</i>	0.009	0.844
<i>HIST3H2A</i>	0.010	0.873
<i>HPRT</i>	0.010	0.897
<i>TERF2IP</i>	0.011	0.949
<i>ITPR1</i>	0.014	0.994
<i>SLC4A1.S</i>	0.014	0.994
<i>COL9A2</i>	0.014	0.994
<i>MCM7</i>	0.015	0.994
<i>CKAP2L</i>	0.017	0.994
<i>RPL18A</i>	0.017	0.994
<i>BRAF</i>	0.017	0.994
<i>MAPK8IP2</i>	0.017	0.994
<i>SFRP4</i>	0.018	0.994
<i>FDPS</i>	0.018	0.994
<i>SACM1L</i>	0.019	0.994
<i>MSMB</i>	0.020	0.994
<i>HMBS</i>	0.020	0.994



## 9: APPENDICES

SPON2	0.021	0.994
ANPEP	0.021	0.994
CACNA1D	0.022	0.994
SYNM	0.023	0.994
ALAS1	0.026	0.994
RNF157	0.027	0.994
HIST1H1E	0.027	0.994
ARHGEF25	0.028	0.994
RPL23AP53	0.028	0.994
AURKA	0.031	0.994
PSTPIP1	0.032	0.994
FOLH1	0.032	0.994
GOLM1	0.033	0.994
EIF2D	0.035	0.994
IFT57	0.039	0.994
SLC43A1	0.039	0.994
CDC20	0.039	0.994
CAMK2N2	0.047	0.994
GABARAPL2	0.049	0.994
CDC37L1	0.050	0.994

KLK3 adjusted data			GAPDH and RPLP2 normalised data		
Transcript	<i>p</i> -value	Adjusted <i>p</i> -value	Transcript	<i>p</i> -value	Adjusted <i>p</i> -value
PCA3	1.52x10-07	2.52x10-05	ERG3' exons 4-5	1.44x10-08	2.41x10-06
SPINK1	5.80x10-06	0.001	TMPRSS2:ERG	1.18x10-07	1.96x10-05
ERG3' exons 4-5	6.32x10-06	0.001	PCA3	2.06x10-07	3.39x10-05
ERG3' exons 6-7	7.48x10-06	0.001	ERG3' exons 6-7	2.28x10-06	0.0004
KLK4	8.86x10-06	0.001	APOC1	9.64x10-06	0.002
SLC12A1	4.36x10-05	0.007	HOXC6	1.34x10-05	0.002
HOXC6	4.72x10-05	0.008	HPN	2.01x10-05	0.003
UPK2	5.48x10-05	0.009	DPP4	9.43x10-05	0.015
HPN	7.32x10-05	0.012	GABARAPL2	0.0001	0.017
TMPRSS2:ERG	0.0001	0.017	ITGBL1	0.0001	0.017
SULT1A1	0.0002	0.036	MYOF	0.0001	0.018

9: APPENDICES

<i>APOC1</i>	0.0004	0.056	<i>KLK2</i>	0.0004	0.065
<i>GJB1</i>	0.0004	0.064	<i>SLC12A1</i>	0.0004	0.069
<i>MYOF</i>	0.0005	0.074	<i>TDRD</i>	0.0004	0.069
<i>CD10</i>	0.001	0.081	<i>SRSF3</i>	0.001	0.088
<i>ITGBL1</i>	0.001	0.138	<i>SPINK1</i>	0.001	0.097
<i>RP11.244H18.1.P712P</i>	0.001	0.138	<i>P712P</i>	0.001	0.098
<i>DLX1</i>	0.001	0.174	<i>KLK4</i>	0.001	0.110
<i>RAB17</i>	0.001	0.179	<i>RAB17</i>	0.001	0.156
<i>GABARAPL2</i>	0.001	0.189	<i>AR exons 4-8</i>	0.001	0.191
<i>STOM</i>	0.001	0.215	<i>IFT57</i>	0.001	0.191
<i>TDRD</i>	0.002	0.266	<i>CD10</i>	0.002	0.276
<i>PTN</i>	0.002	0.293	<i>PTN</i>	0.003	0.375
<i>AMACR</i>	0.002	0.327	<i>DLX1</i>	0.003	0.418
<i>MED4</i>	0.003	0.378	<i>ANKRD34B</i>	0.003	0.419
<i>SNCA</i>	0.004	0.622	<i>ZNF577</i>	0.003	0.434
<i>NEAT1</i>	0.005	0.743	<i>UPK2</i>	0.003	0.443
<i>ANKRD34B</i>	0.005	0.743	<i>MXI1</i>	0.004	0.606
<i>MIR4435.1HG.IOC541471</i>	0.006	0.783	<i>HOXC4</i>	0.006	0.767
<i>KLK2</i>	0.007	0.928	<i>SNCA</i>	0.006	0.769
<i>Met</i>	0.012	0.980	<i>STEAP2</i>	0.006	0.800
<i>AURKA</i>	0.014	0.980	<i>MEMO1</i>	0.006	0.834
<i>SIM2.short</i>	0.014	0.980	<i>CACNA1D</i>	0.006	0.834
<i>MIC1</i>	0.019	0.980	<i>STEAP4</i>	0.007	0.954
<i>PPFIA2</i>	0.020	0.980	<i>PPAP2A</i>	0.008	0.997
<i>MEMO1</i>	0.021	0.980	<i>Met</i>	0.011	0.997
<i>ZNF577</i>	0.023	0.980	<i>MED4</i>	0.011	0.997
<i>CACNA1D</i>	0.025	0.980	<i>MIATNB</i>	0.011	0.997
<i>AR exon 9</i>	0.026	0.980	<i>GJB1</i>	0.013	0.997
<i>PDLIM5</i>	0.027	0.980	<i>AR exon 9</i>	0.014	0.997
<i>RP11.97O12.7</i>	0.033	0.980	<i>SULT1A1</i>	0.018	0.997
<i>IFT57</i>	0.033	0.980	<i>PPFIA2</i>	0.018	0.997
<i>MMP26</i>	0.033	0.980	<i>FDPS</i>	0.019	0.997
<i>MARCH5</i>	0.034	0.980	<i>MARCH5</i>	0.019	0.997
<i>RPS10</i>	0.036	0.980	<i>MSMB</i>	0.020	0.997

## 9: APPENDICES

<i>AR exons 4-8</i>	0.036	0.980	<i>KLK3 exons 2-3</i>	0.024	0.997
<i>ITPR1</i>	0.039	0.980	<i>SNORA20</i>	0.026	0.997
<i>SMIM1</i>	0.043	0.980	<i>NEAT1</i>	0.027	0.997
<i>SNORA20</i>	0.044	0.980	<i>RPS10</i>	0.028	0.997
<i>VPS13A</i>	0.046	0.980	<i>SERPINB5</i>	0.035	0.997
			<i>TRPM4</i>	0.038	0.997
			<i>NLRP3</i>	0.040	0.997
			<i>HIST1H2BF</i>	0.049	0.997

9: APPENDICES

Supplementary Table 19 Lasso output for models detecting CB, L, I, H trend using KLK2 ratio data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>PCA3</i>	<b>0.21</b>	<i>PCA3</i>	<b>0.18</b>	<i>PCA3</i>	<b>0.14</b>
<i>ERG3' exons 4-5</i>	<b>0.11</b>	<i>ERG3' exons 4-5</i>	<b>0.12</b>	<i>ERG3' exons 4-5</i>	<b>0.11</b>
<i>APOC1</i>	<b>0.11</b>	<i>APOC1</i>	<b>0.08</b>	<i>APOC1</i>	<b>0.05</b>
<i>ANKRD34B</i>	<b>0.07</b>	<i>TMPRSS2:ERG</i>	<b>0.04</b>	<i>TMPRSS2:ERG</i>	<b>0.04</b>
<i>NEAT1</i>	<b>0.05</b>	<i>SLC12A1</i>	<b>0.03</b>	<i>HOXC6</i>	<b>0.01</b>
<i>HOXC6</i>	<b>0.05</b>	<i>HOXC6</i>	<b>0.02</b>	<i>cp1</i>	<b>2.05</b>
<i>HPN</i>	<b>0.05</b>	<i>NEAT1</i>	<b>0.02</b>	<i>cp2</i>	<b>1.26</b>
<i>TMPRSS2:ERG</i>	<b>0.04</b>	<i>HPN</i>	<b>0.01</b>	<i>cp3</i>	<b>-0.36</b>
<i>ITGBL1</i>	<b>0.03</b>	<i>ANKRD34B</i>	<b>0.01</b>		
<i>SLC12A1</i>	<b>0.03</b>	<i>DLX1</i>	<b>0.00</b>		
<i>SULT1A1</i>	<b>0.03</b>	<i>PSTPIP1</i>	<b>0.00</b>		
<i>ISX</i>	<b>0.03</b>	<i>HIST1H1E</i>	<b>-0.05</b>		
<i>DLX1</i>	<b>0.02</b>	<i>GABARAPL2</i>	<b>-0.16</b>		
<i>ERG3' exons 6-7</i>	<b>0.01</b>	<i>cp1</i>	<b>2.19</b>		
<i>TMEM47</i>	<b>0.01</b>	<i>cp2</i>	<b>1.32</b>		
<i>TDRD</i>	<b>0.01</b>	<i>cp3</i>	<b>-0.38</b>		
<i>AMACR</i>	<b>0.01</b>				
<i>HIST1H1E</i>	<b>-0.01</b>				
<i>IGFBP3</i>	<b>-0.01</b>				
<i>PSGR</i>	<b>-0.01</b>				
<i>BTG2</i>	<b>-0.01</b>				
<i>MED4</i>	<b>-0.02</b>				
<i>AR exons 4-8</i>	<b>-0.02</b>				
<i>PPP1R12B</i>	<b>-0.02</b>				
<i>AR exon 9</i>	<b>-0.02</b>				
<i>Timp4</i>	<b>-0.03</b>				
<i>DPP4</i>	<b>-0.03</b>				
<i>CP</i>	<b>-0.04</b>				
<i>MYOF</i>	<b>-0.04</b>				
<i>GCNT1</i>	<b>-0.04</b>				
<i>MEMO1</i>	<b>-0.05</b>				
<i>SRSF3</i>	<b>-0.06</b>				
<i>ZNF577</i>	<b>-0.06</b>				
<i>CD10</i>	<b>-0.06</b>				
<i>MXII</i>	<b>-0.10</b>				
<i>KLK4</i>	<b>-0.14</b>				
<i>cp1</i>	<b>2.45</b>				
<i>cp2</i>	<b>1.43</b>				
<i>cp3</i>	<b>-0.43</b>				

## 9: APPENDICES

Supplementary Table 20 Lasso output for models detecting CB, L, I, H trend using KLK2 adjusted data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Bet a</i>	<i>Transcript</i>	<i>Bet a</i>	<i>Transcript</i>	<i>Beta</i>
<i>AMACR</i>	<b>0.14</b>	<i>AMACR</i>	<b>0.46</b>	<i>ERG3' exons 4-5</i>	<b>0.75</b>
<i>ERG3' exons 4-5</i>	<b>0.74</b>	<i>ANKRD34B</i>	<b>0.33</b>	<i>GABARAPL2</i>	<b>1.05</b>
<i>GJB1</i>	<b>0.60</b>	<i>APOC1</i>	<b>0.53</b>	<i>GJB1</i>	<b>0.63</b>
<i>HOXC6</i>	<b>0.36</b>	<i>CD10</i>	<b>1.04</b>	<i>HOXC6</i>	<b>0.38</b>
<i>HPN</i>	<b>0.74</b>	<i>DLX1</i>	<b>0.13</b>	<i>HPN</i>	<b>0.70</b>
<i>ITGBL1</i>	<b>0.12</b>	<i>DPP4</i>	<b>0.14</b>	<i>ITGBL1</i>	<b>0.18</b>
<i>KLK4</i>	<b>1.22</b>	<i>ERG3' exons 4-5</i>	<b>0.80</b>	<i>KLK4</i>	<b>1.10</b>
<i>PCA3</i>	<b>2.38</b>	<i>GABARAPL2</i>	<b>0.88</b>	<i>PCA3</i>	<b>2.27</b>
<i>SLC12A1</i>	<b>0.25</b>	<i>GAPDH</i>	<b>0.07</b>	<i>RP11.244H18.1.P712P</i>	<b>1.34</b>
<i>SPINK1</i>	<b>0.44</b>	<i>GJB1</i>	<b>0.08</b>	<i>SLC12A1</i>	<b>0.26</b>
<i>TMPRSS2:ERG</i>	<b>0.36</b>	<i>HOXC6</i>	<b>0.22</b>	<i>SPINK1</i>	<b>0.38</b>
<i>UPK2</i>	<b>0.24</b>	<i>IFT57</i>	<b>0.95</b>	<i>TDRD</i>	<b>0.14</b>
<i>cp1</i>	<b>2.06</b>	<i>ITPR1</i>	<b>0.13</b>	<i>TMPRSS2:ERG</i>	<b>0.38</b>
<i>cp2</i>	<b>1.35</b>	<i>KLK4</i>	<b>0.60</b>	<i>UPK2</i>	<b>0.19</b>
<i>cp3</i>	<b>0.42</b>	<i>MED4</i>	<b>0.85</b>	<i>cp1</i>	<b>2.21</b>
		<i>Met</i>	<b>0.11</b>	<i>cp2</i>	<b>1.41</b>
		<i>MIC1</i>	<b>0.28</b>	<i>cp3</i>	<b>0.43</b>
		<i>MIR146A.DQ658414</i>	<b>0.24</b>		
		<i>MMP26</i>	<b>0.50</b>		
		<i>MX11</i>	<b>1.15</b>		
		<i>MYOF</i>	<b>1.45</b>		
		<i>PCA3</i>	<b>2.69</b>		
		<i>PPAP2A</i>	<b>0.08</b>		
		<i>PPFIA2</i>	<b>0.65</b>		
		<i>PTN</i>	<b>0.79</b>		
		<i>RP11.244H18.1.P712P</i>	<b>0.72</b>		
		<i>SIM2.short</i>	<b>1.12</b>		
		<i>SLC12A1</i>	<b>0.14</b>		
		<i>SMIM1</i>	<b>0.40</b>		

9: APPENDICES

	-
<i>SNCA</i>	<i>0.77</i>
<i>SPINK1</i>	<i>0.47</i>
<i>STEAP2</i>	<i>0.82</i>
<i>STOM</i>	<i>0.11</i>
<i>SULT1A1</i>	<i>0.85</i>
<i>TMPRSS2:ERG</i>	<i>0.13</i>
<i>UPK2</i>	<i>0.25</i>
	-
<i>ZNF577</i>	<i>0.49</i>
<i>cp1</i>	<i>2.47</i>
<i>cp2</i>	<i>1.52</i>
	-
<i>cp3</i>	<i>0.48</i>

## 9: APPENDICES

Supplementary Table 21 Lasso output for models detecting CB, L, I, H trend using KLK3 adjusted data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>ACTR5</i>	<b>-0.40</b>	<i>MARCH5</i>	<b>-0.88</b>	<i>APOC1</i>	<b>0.78</b>
<i>AMH</i>	<b>1.60</b>	<i>AMACR</i>	<b>0.27</b>	<i>ERG3'</i> <i>exons 4-5</i>	<b>0.56</b>
<i>ANKRD34B</i>	<b>0.20</b>	<i>ANKRD34B</i>	<b>0.21</b>	<i>ERG3'</i> <i>exons 6-7</i>	<b>0.07</b>
<i>APOC1</i>	<b>0.96</b>	<i>APOC1</i>	<b>0.60</b>	<i>GJB1</i>	<b>0.51</b>
<i>AR exon 9</i>	<b>-0.05</b>	<i>AR exon 9</i>	<b>-0.49</b>	<i>HOXC6</i>	<b>0.36</b>
<i>AURKA</i>	<b>0.07</b>	<i>AURKA</i>	<b>0.70</b>	<i>HPN</i>	<b>0.52</b>
<i>B2M</i>	<b>-0.23</b>	<i>CACNA1D</i>	<b>0.42</b>	<i>KLK4</i>	<b>-1.25</b>
<i>BRAF</i>	<b>0.48</b>	<i>CD10</i>	<b>-0.75</b>	<i>PCA3</i>	<b>2.21</b>
<i>BTG2</i>	<b>-1.11</b>	<i>DLX1</i>	<b>0.27</b>	<i>SLC12A1</i>	<b>0.19</b>
<i>CASKINI</i>	<b>-0.01</b>	<i>ERG3'</i> <i>4-5</i>	<b>0.64</b>	<i>SULT1A1</i>	<b>0.65</b>
<i>CCDC88B</i>	<b>-0.27</b>	<i>GABARAPL2</i>	<b>-1.14</b>	<i>TMPRSS2:E</i> <i>RG</i>	<b>0.25</b>
<i>CD10</i>	<b>-0.95</b>	<i>GJB1</i>	<b>0.28</b>	<i>UPK2</i>	<b>0.50</b>
<i>CDC20</i>	<b>-0.44</b>	<i>HOXC6</i>	<b>0.45</b>	<i>cp1</i>	<b>1.14</b>
<i>CKAP2L</i>	<b>-0.39</b>	<i>ITGBL1</i>	<b>-0.31</b>	<i>cp2</i>	<b>1.35</b>
<i>CLIC2</i>	<b>-0.52</b>	<i>ITPR1</i>	<b>0.68</b>	<i>cp3</i>	<b>-0.36</b>
<i>CLU</i>	<b>0.11</b>	<i>KLK2</i>	<b>0.42</b>		
<i>CP</i>	<b>-0.46</b>	<i>KLK4</i>	<b>-0.53</b>		
<i>CTA.211A9.5.MIATNB</i>	<b>-0.41</b>	<i>MED4</i>	<b>-0.86</b>		
<i>DLX1</i>	<b>0.33</b>	<i>MEMO1</i>	<b>-0.90</b>		
<i>DNAH5</i>	<b>0.23</b>	<i>MIC1</i>	<b>0.40</b>		
<i>ERG3'</i> <i>exons 4-5</i>	<b>0.54</b>	<i>MIR4435.1H</i> <i>G.IOC541471</i>	<b>-0.35</b>		
<i>ERG3'</i> <i>exons 6-7</i>	<b>0.36</b>	<i>MMP26</i>	<b>0.24</b>		
<i>GABARAPL2</i>	<b>-1.42</b>	<i>MYOF</i>	<b>-1.84</b>		
<i>GOLM1</i>	<b>-0.02</b>	<i>NEAT1</i>	<b>0.80</b>		
<i>HIST3H2A</i>	<b>-0.34</b>	<i>PCA3</i>	<b>3.73</b>		
<i>HOXC4</i>	<b>-0.97</b>	<i>PPFIA2</i>	<b>-0.97</b>		
<i>HOXC6</i>	<b>0.76</b>	<i>PTN</i>	<b>-1.10</b>		
<i>HPRT</i>	<b>0.81</b>	<i>RAB17</i>	<b>-0.89</b>		
<i>IGFBP3</i>	<b>-0.87</b>	<i>RP11.244H18</i> <i>.1.P712P</i>	<b>-0.34</b>		
<i>IMPDH2</i>	<b>0.03</b>	<i>RP11.97O12.</i> <i>7</i>	<b>0.48</b>		
<i>ITPR1</i>	<b>0.17</b>	<i>SIM2.short</i>	<b>1.35</b>		
<i>KLK4</i>	<b>-1.18</b>	<i>SLC12A1</i>	<b>0.39</b>		
<i>LASS1</i>	<b>-0.16</b>	<i>SMIM1</i>	<b>0.47</b>		
<i>LBH</i>	<b>0.51</b>	<i>SNCA</i>	<b>-0.20</b>		
<i>MCM7</i>	<b>0.27</b>	<i>STOM</i>	<b>0.31</b>		
<i>MDK</i>	<b>0.01</b>	<i>SULT1A1</i>	<b>1.05</b>		
<i>MED4</i>	<b>-0.38</b>	<i>TMPRSS2:E</i> <i>RG</i>	<b>-0.02</b>		

9: APPENDICES

<i>MEMO1</i>	<b>-0.97</b>	<i>UPK2</i>	<b>0.94</b>
<i>MGAT5B</i>	<b>-0.24</b>	<i>ZNF577</i>	<b>-0.14</b>
<i>MIC1</i>	<b>0.64</b>	<i>cp1</i>	<b>1.77</b>
<i>MIR146A.DQ658414</i>	<b>0.20</b>	<i>cp2</i>	<b>1.63</b>
<i>MIR4435.1HG.IOC5414</i> <i>71</i>	<b>0.02</b>	<i>cp3</i>	<b>-0.46</b>
<i>MMP11</i>	<b>0.61</b>		
<i>MMP25</i>	<b>0.83</b>		
<i>MNX1</i>	<b>0.42</b>		
<i>MX11</i>	<b>-0.18</b>		
<i>MYOF</i>	<b>-1.42</b>		
<i>NAALADL2</i>	<b>-0.32</b>		
<i>NEAT1</i>	<b>0.98</b>		
<i>PSGR</i>	<b>-0.15</b>		
<i>PALM3</i>	<b>0.21</b>		
<i>PCA3</i>	<b>3.16</b>		
<i>PPAP2A</i>	<b>0.89</b>		
<i>PSTPIP1</i>	<b>-0.62</b>		
<i>PTN</i>	<b>-0.71</b>		
<i>PVT1</i>	<b>0.17</b>		
<i>RPL23AP53</i>	<b>0.06</b>		
<i>RPS11</i>	<b>-0.02</b>		
<i>SACMIL</i>	<b>-0.41</b>		
<i>SERPINB5</i>	<b>0.28</b>		
<i>SIM2.short</i>	<b>0.90</b>		
<i>SIRT1</i>	<b>-0.68</b>		
<i>SLC12A1</i>	<b>0.23</b>		
<i>SMIM1</i>	<b>0.20</b>		
<i>SNCA</i>	<b>-0.22</b>		
<i>SPINK1</i>	<b>0.49</b>		
<i>SPON2</i>	<b>-0.43</b>		
<i>SRSF3</i>	<b>-0.01</b>		
<i>ST6GALNAC1</i>	<b>0.25</b>		
<i>STOM</i>	<b>0.36</b>		
<i>SULT1A1</i>	<b>0.61</b>		
<i>SYNM</i>	<b>0.19</b>		
<i>TDRD</i>	<b>0.03</b>		
<i>Timp4</i>	<b>-0.91</b>		
<i>TWIST1</i>	<b>0.65</b>		
<i>UPK2</i>	<b>0.75</b>		
<i>VAX2</i>	<b>-0.27</b>		
<i>ZNF577</i>	<b>-0.02</b>		
<i>cp1</i>	<b>2.04</b>		
<i>cp2</i>	<b>1.75</b>		
<i>cp3</i>	<b>-0.50</b>		

Supplementary Table 22 Lasso output for models detecting CB, L, I, H trend using HK normalised data.

*All Transcripts*

*Significant Transcripts*

*Multiple testing*



## 9: APPENDICES

				<i>corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
<i>ACTR5</i>	<b>0.00</b>	<i>ANKRD34B</i>	<b>0.11</b>	<i>APOC1</i>	<b>0.14</b>
<i>AMH</i>	<b>0.04</b>	<i>APOC1</i>	<b>0.12</b>	<i>DPP4</i>	<b>-0.25</b>
<i>ANKRD34B</i>	<b>0.07</b>	<i>AR exon 9</i>	<b>-0.04</b>	<i>ERG3' exons 4-5</i>	<b>0.11</b>
<i>APOC1</i>	<b>0.11</b>	<i>AR exons 4-8</i>	<b>-0.04</b>	<i>ERG3' exons 6-7</i>	<b>0.01</b>
<i>AR exon 9</i>	<b>0.02</b>	<i>CD10</i>	<b>-0.08</b>	<i>GABARAPL2</i>	<b>-0.26</b>
<i>AR exons 4-8</i>	<b>0.02</b>	<i>MIATNB</i>	<b>0.00</b>	<i>HOXC6</i>	<b>0.08</b>
<i>CD10</i>	<b>0.05</b>	<i>DLX1</i>	<b>0.03</b>	<i>HPN</i>	<b>0.10</b>
<i>CP</i>	<b>0.05</b>	<i>DPP4</i>	<b>-0.07</b>	<i>ITGBL1</i>	<b>0.10</b>
<i>DLX1</i>	<b>0.01</b>	<i>ERG3' exons 4-5</i>	<b>0.10</b>	<i>KLK4</i>	<b>-0.30</b>
<i>DPP4</i>	<b>0.06</b>	<i>ERG3' exons 6-7</i>	<b>0.01</b>	<i>MYOF</i>	<b>-0.17</b>
<i>ERG3' exons 4-5</i>	<b>0.10</b>	<i>FDPS</i>	<b>-0.03</b>	<i>PCA3</i>	<b>0.23</b>
<i>ERG3' exons 6-7</i>	<b>0.01</b>	<i>GABARAPL2</i>	<b>-0.03</b>	<i>TDRD</i>	<b>0.04</b>
<i>GABARAPL2</i>	<b>0.08</b>	<i>GJB1</i>	<b>0.02</b>	<i>TMPRSS2:ERG</i>	<b>0.04</b>
<i>GCNT1</i>	<b>0.03</b>	<i>HOXC6</i>	<b>0.06</b>	<i>cp1</i>	<b>2.63</b>
<i>HIST1H1E</i>	<b>0.03</b>	<i>HPN</i>	<b>0.08</b>	<i>cp2</i>	<b>1.52</b>
<i>HIST1H2BF</i>	<b>0.00</b>	<i>ITGBL1</i>	<b>0.07</b>	<i>cp3</i>	<b>-0.46</b>
<i>HOXC6</i>	<b>0.05</b>	<i>KLK4</i>	<b>-0.21</b>		
<i>HPN</i>	<b>0.05</b>	<i>MED4</i>	<b>-0.13</b>		
<i>IGFBP3</i>	<b>0.01</b>	<i>MEMO1</i>	<b>-0.09</b>		
<i>ISX</i>	<b>0.02</b>	<i>MSMB</i>	<b>0.10</b>		
<i>ITGBL1</i>	<b>0.07</b>	<i>MXI1</i>	<b>-0.24</b>		
<i>KLK4</i>	<b>0.16</b>	<i>MYOF</i>	<b>-0.06</b>		
<i>MED4</i>	<b>0.02</b>	<i>NEAT1</i>	<b>0.06</b>		
<i>MEMO1</i>	<b>0.05</b>	<i>PCA3</i>	<b>0.21</b>		
<i>MXI1</i>	<b>0.13</b>	<i>RPS10</i>	<b>-0.02</b>		
<i>MYOF</i>	<b>0.05</b>	<i>SLC12A1</i>	<b>0.04</b>		
<i>NEAT1</i>	<b>0.04</b>	<i>SPINK1</i>	<b>0.00</b>		
<i>PCA3</i>	<b>0.19</b>	<i>SRSF3</i>	<b>-0.09</b>		
<i>PPP1R12B</i>	<b>0.02</b>	<i>SULT1A1</i>	<b>0.06</b>		
<i>SLC12A1</i>	<b>0.03</b>	<i>TDRD</i>	<b>0.03</b>		
<i>SRSF3</i>	<b>0.05</b>	<i>TMPRSS2:ERG</i>	<b>0.03</b>		
<i>SULT1A1</i>	<b>0.03</b>	<i>TRPM4</i>	<b>-0.02</b>		

9: APPENDICES

<i>TDRD</i>	<b>0.01</b>	<i>UPK2</i>	<b>0.02</b>
	-		
<i>Timp4</i>	<b>0.02</b>	<i>ZNF577</i>	<b>-0.07</b>
<i>TMPRSS2:ERG</i>	<b>0.04</b>	<i>cp1</i>	<b>2.67</b>
<i>UPK2</i>	<b>0.00</b>	<i>cp2</i>	<b>1.55</b>
	-		
<i>ZNF577</i>	<b>0.04</b>	<i>cp3</i>	<b>-0.47</b>
<i>cp1</i>	<b>2.42</b>		
<i>cp2</i>	<b>1.41</b>		
	-		
<i>cp3</i>	<b>0.42</b>		

9: APPENDICES

**6.16 Multinomial CBCaM Trend**

**Supplementary Table 23 Glm test significant probes for CB, Ca, Mets trend**

<i>KLK2</i> ratio data			<i>KLK2</i> adjusted data		
<b>Transcript</b>	<b><i>p</i>-value</b>	<b>Adjusted <i>p</i>-value</b>	<b>Transcript</b>	<b><i>p</i>-value</b>	<b>Adjusted <i>p</i>-value</b>
<i>HOXC6</i>	5.20X10-10	8.62X10-08	<i>UPK2</i>	2.91X10-08	4.83X10-06
<i>ERG3' exons 4-5</i>	5.70X10-10	9.41X10-08	<i>SPINK1</i>	2.04X10-07	3.36X10-05
<i>PCA3</i>	4.58X10-09	7.51X10-07	<i>SLC12A1</i>	1.61X10-06	0.0003
<i>TMPRSS2:ERG</i>	1.27X10-08	2.07X10-06	<i>HOXC6</i>	5.86X10-05	0.0096
<i>APOC1</i>	1.42X10-08	2.30X10-06	<i>HPN</i>	7.26X10-05	0.0118
<i>TDRD</i>	2.07X10-08	3.34X10-06	<i>MFSD2A</i>	7.70X10-05	0.0124
<i>SLC12A1</i>	2.82X10-08	4.51X10-06	<i>GAPDH</i>	0.0001	0.0196
<i>HPN</i>	3.41X10-08	5.42X10-06	<i>RAB17</i>	0.0002	0.0285
<i>HOXC4</i>	2.04X10-07	3.23X10-05	<i>KLK4</i>	0.0002	0.0313
<i>RAB17</i>	2.07X10-07	3.26X10-05	<i>PCA3</i>	0.0005	0.0758
<i>GJB1</i>	2.14X10-07	3.34X10-05	<i>GJB1</i>	0.0005	0.0826
<i>ERG3' exons 6-7</i>	2.32X10-07	3.59X10-05	<i>MIR4435.1HG.IOC541471</i>	0.0007	0.1054
<i>AMACR</i>	4.06X10-07	6.25X10-05	<i>GABARAPL2</i>	0.0007	0.1085
<i>SPINK1</i>	4.92X10-07	7.53X10-05	<i>TMEM45B</i>	0.0009	0.1315
<i>SSTR1</i>	5.08X10-07	7.71X10-05	<i>APOC1</i>	0.0009	0.1376
<i>UPK2</i>	5.69X10-07	8.60X10-05	<i>AURKA</i>	0.0012	0.1751
<i>TMCC2</i>	6.63X10-07	9.94X10-05	<i>ANPEP</i>	0.0012	0.1760
<i>TMEM45B</i>	6.75X10-07	0.0001	<i>SULT1A1</i>	0.0017	0.2479
<i>PPFIA2</i>	7.28X10-07	0.0001	<i>RP11.244H18.1.P712P</i>	0.0020	0.2933
<i>DLX1</i>	8.97X10-07	0.0001	<i>ERG3' exons 4-5</i>	0.0023	0.3385
<i>PALM3</i>	9.33X10-07	0.0001	<i>PALM3</i>	0.0023	0.3421
<i>SULT1A1</i>	1.03X10-06	0.0001	<i>TMPRSS2:ERG</i>	0.0024	0.3487
<i>RP11.97012.7</i>	1.60X10-06	0.0002	<i>TDRD</i>	0.0028	0.4042
<i>SULF2</i>	2.21X10-06	0.0003	<i>ERG3' exons 6-7</i>	0.0029	0.4165
<i>AMH</i>	2.21X10-06	0.0003	<i>TBP</i>	0.0031	0.4376
<i>EN2</i>	2.70X10-06	0.0004	<i>HMBS</i>	0.0035	0.4878

9: APPENDICES

<i>CASKIN1</i>	2.79X10-06	0.0004	<i>ITGBL1</i>	0.0041	0.5687
<i>MIR4435.1HG.IOC541471</i>	2.91X10-06	0.0004	<i>AMACR</i>	0.0042	0.5824
<i>HIST1H1C</i>	3.23X10-06	0.0004	<i>TMCC2</i>	0.0044	0.6065
<i>MEX3A</i>	3.24X10-06	0.0004	<i>MYOF</i>	0.0058	0.7895
<i>PECI</i>	3.81X10-06	0.0005	<i>RNF157</i>	0.0073	0.9921
<i>SIM2.short</i>	3.92X10-06	0.0005	<i>PTPRC</i>	0.0079	0.9921
<i>ISX</i>	4.01X10-06	0.0005	<i>SIM2.short</i>	0.0081	0.9921
<i>TMEM86A</i>	4.87X10-06	0.0006	<i>SULF2</i>	0.0088	0.9921
<i>ERG5</i>	5.11X10-06	0.0007	<i>EN2</i>	0.0094	0.9921
<i>TWIST1</i>	5.24X10-06	0.0007	<i>PTN</i>	0.0095	0.9921
<i>ITGBL1</i>	6.47X10-06	0.0008	<i>ALAS1</i>	0.0104	0.9921
<i>MGAT5B</i>	6.56X10-06	0.0008	<i>TMEM86A</i>	0.0115	0.9921
<i>MMP11</i>	7.01X10-06	0.0009	<i>RP11.97O12.7</i>	0.0117	0.9921
<i>HMBS</i>	7.33X10-06	0.0009	<i>PPFIA2</i>	0.0122	0.9921
<i>MCTP1</i>	8.97X10-06	0.0011	<i>DPP4</i>	0.0125	0.9921
<i>GAPDH</i>	1.06X10-05	0.0013	<i>STOM</i>	0.0132	0.9921
<i>STOM</i>	1.09X10-05	0.0014	<i>Met</i>	0.0139	0.9921
<i>HIST3H2A</i>	1.26X10-05	0.0016	<i>ZNF577</i>	0.0153	0.9921
<i>RPL23AP53</i>	1.29X10-05	0.0016	<i>ERG5</i>	0.0251	0.9921
<i>MFSD2A</i>	1.49X10-05	0.0018	<i>ITPR1</i>	0.0280	0.9921
<i>TERT</i>	1.72X10-05	0.0021	<i>MARCH5</i>	0.0299	0.9921
<i>Met</i>	2.06X10-05	0.0025	<i>HIST1H1E</i>	0.0342	0.9921
<i>B4GALNT4</i>	2.08X10-05	0.0025	<i>SMIM1</i>	0.0380	0.9921
<i>NLRP3</i>	2.10X10-05	0.0025	<i>DLX1</i>	0.0384	0.9921
<i>PVT1</i>	2.14X10-05	0.0025	<i>RPL23AP53</i>	0.0433	0.9921
<i>MIR146A.DQ658414</i>	2.34X10-05	0.0027	<i>CASKIN1</i>	0.0457	0.9921
<i>CCDC88B</i>	2.70X10-05	0.0031	<i>SEC61A1</i>	0.0474	0.9921
<i>PPAP2A</i>	2.71X10-05	0.0031	<i>AMH</i>	0.0478	0.9921
<i>ITPR1</i>	3.03X10-05	0.0034	<i>IFT57</i>	0.0481	0.9921
<i>ABCB9</i>	3.22X10-05	0.0035	<i>CLU</i>	0.0497	0.9921
<i>ANPEP</i>	3.25X10-05	0.0035	<i>VPS13A</i>	0.0499	0.9921
<i>VPS13A</i>	3.25X10-05	0.0035			
<i>MMP25</i>	3.30X10-05	0.0036			
<i>PSTPIP1</i>	3.39X10-05	0.0036			

## 9: APPENDICES

<i>AURKA</i>	3.44X10-05	0.0036
<i>VAX2</i>	3.72X10-05	0.0039
<i>TRPM4</i>	3.81X10-05	0.0039
<i>PTPRC</i>	3.82X10-05	0.0039
<i>RIOK3</i>	3.85X10-05	0.0039
<i>OGT</i>	4.06X10-05	0.0041
<i>MXN1</i>	4.10X10-05	0.0041
<i>SLC4A1.5</i>	4.52X10-05	0.0045
<i>HPRT</i>	4.84X10-05	0.0047
<i>TBP</i>	4.91X10-05	0.0048
<i>HIST1H1E</i>	5.41X10-05	0.0052
<i>NAALADL2</i>	5.44X10-05	0.0052
<i>SIM2.long</i>	5.69X10-05	0.0053
<i>CLIC2</i>	6.01X10-05	0.0056
<i>DNAH5</i>	6.67X10-05	0.0061
<i>SMIM1</i>	7.71X10-05	0.0070
<i>PCSK6</i>	8.13X10-05	0.0073
<i>MKI67</i>	8.24X10-05	0.0073
<i>COL9A2</i>	8.76X10-05	0.0077
<i>BRAF</i>	8.85X10-05	0.0077
<i>COL10A1</i>	9.42X10-05	0.0081
<i>TERF2IP</i>	9.55X10-05	0.0081
<i>SSPO</i>	0.0001	0.0091
<i>RPLP2</i>	0.0001	0.0097
<i>SFRP4</i>	0.0001	0.0097
<i>MAPK8IP2</i>	0.0001	0.0097
<i>CDC37L1</i>	0.0001	0.0097
<i>RNF157</i>	0.0001	0.0101
<i>ACTR5</i>	0.0001	0.0103
<i>RPS11</i>	0.0001	0.0110
<i>RPS10</i>	0.0001	0.0110
<i>SYNM</i>	0.0001	0.0111
<i>CDKN3</i>	0.0002	0.0120
<i>AATF</i>	0.0002	0.0127

## 9: APPENDICES

<i>EIF2D</i>	0.0002	0.0132
<i>ALAS1</i>	0.0002	0.0139
<i>IMPDH2</i>	0.0002	0.0145
<i>FDPS</i>	0.0002	0.0157
<i>SACM1L</i>	0.0002	0.0165
<i>TFDP1</i>	0.0003	0.0170
<i>MCM7</i>	0.0003	0.0202
<i>NEAT1</i>	0.0003	0.0208
<i>CAMKK2</i>	0.0004	0.0232
<i>IFT57</i>	0.0004	0.0242
<i>MEMO1</i>	0.0005	0.0299
<i>ANKRD34B</i>	0.0005	0.0309
<i>RPL18A</i>	0.0005	0.0310
<i>PPP1R12B</i>	0.0006	0.0364
<i>CACNA1D</i>	0.0006	0.0370
<i>MIC1</i>	0.0008	0.0436
<i>SPON2</i>	0.0008	0.0441
<i>CLU</i>	0.0008	0.0441
<i>BTG2</i>	0.0010	0.0526
<i>GABARAPL2</i>	0.0014	0.0709
<i>SMAP1 exons 7-8</i>	0.0014	0.0709
<i>STEAP4</i>	0.0014	0.0709
<i>CKAP2L</i>	0.0015	0.0727
<i>KLK3 exons 2-3</i>	0.0015	0.0727
<i>ARHGEF25</i>	0.0017	0.0796
<i>LASS1</i>	0.0017	0.0796
<i>STEAP2</i>	0.0019	0.0887
<i>B2M</i>	0.0021	0.0964
<i>MMP26</i>	0.0023	0.0996
<i>MXI1</i>	0.0024	0.1046
<i>SChLAP1</i>	0.0025	0.1056
<i>CDC20</i>	0.0026	0.1085
<i>CAMK2N2</i>	0.0028	0.1104
<i>SIRT1</i>	0.0032	0.1231

## 9: APPENDICES

<i>GOLM1</i>	0.0044	0.1676
<i>LBH</i>	0.0062	0.2206
<i>SEC61A1</i>	0.0062	0.2206
<i>AR exons 4-8</i>	0.0064	0.2206
<i>MAK</i>	0.0066	0.2206
<i>PDLIM5</i>	0.0067	0.2206
<i>SRSF3</i>	0.0074	0.2378
<i>SNCA</i>	0.0087	0.2628
<i>FOLH1</i>	0.0089	0.2628
<i>CADPS</i>	0.0091	0.2628
<i>CD10</i>	0.0103	0.2879
<i>MDK</i>	0.0133	0.3598
<i>KLK3 exons 1-2</i>	0.0145	0.3767
<i>MED4</i>	0.0152	0.3792
<i>HIST1H2BG</i>	0.0193	0.4640
<i>PTN</i>	0.0245	0.5637
<i>IGFBP3</i>	0.0268	0.5789
<i>DPP4</i>	0.0276	0.5789
<i>SERPINB5</i>	0.0302	0.6032
<i>ST6GALNAC1</i>	0.0343	0.6512
<i>MARCH5</i>	0.0372	0.6697
<i>HIST1H2BF</i>	0.0405	0.6880
<i>MSMB</i>	0.0471	0.7190
<i>SLC43A1</i>	0.0479	0.7190

KLK3 adjusted data			<i>GAPDH</i> and <i>RPLP2</i> normalised data		
Transcript	<i>p</i> -value	Adjusted <i>p</i> -value	Transcript	<i>p</i> -value	Adjusted <i>p</i> -value
<i>UPK2</i>	2.39x10-08	3.97x10-06	<i>HOXC6</i>	3.39X10-06	0.0006
<i>SPINK1</i>	1.87x10-06	0.0003	<i>SLC12A1</i>	3.93X10-06	0.0007
<i>SLC12A1</i>	2.58x10-06	0.0004	<i>APOC1</i>	7.43X10-06	0.0012
<i>RAB17</i>	4.04x10-06	0.0007	<i>ERG3' exons 4-5</i>	2.17X10-05	0.0036
<i>MIR4435.1HG.IOC541471</i>	3.58x10-05	0.0058	<i>SPINK1</i>	2.71X10-05	0.0044
<i>HPN</i>	5.22x10-05	0.0084	<i>KLK2</i>	3.80X10-05	0.0062
<i>KLK4</i>	7.61x10-05	0.0122	<i>TMPRSS2:ERG</i>	5.96X10-05	0.0096

## 9: APPENDICES

<i>HOXC6</i>	0.0001	0.0168	<i>HPN</i>	6.42X10-05	0.0103
<i>GABARAPL2</i>	0.0003	0.0404	<i>UPK2</i>	6.50X10-05	0.0103
<i>MFS2A</i>	0.0005	0.0832	<i>RAB17</i>	6.91X10-05	0.0109
<i>SULT1A1</i>	0.0007	0.1103	<i>TDRD</i>	0.0003	0.0531
<i>GJB1</i>	0.0007	0.1155	<i>KLK4</i>	0.0005	0.0831
<i>APOC1</i>	0.0010	0.1464	<i>ERG3' exons 6-7</i>	0.0005	0.0843
<i>PCA3</i>	0.0012	0.1782	<i>GABARAPL2</i>	0.0008	0.1191
<i>TMEM45B</i>	0.0012	0.1888	<i>HOXC4</i>	0.0010	0.1498
<i>RP11.244H18.1.P712P</i>	0.0017	0.2624	<i>TMEM45B</i>	0.0010	0.1563
<i>MYOF</i>	0.0018	0.2720	<i>P712P</i>	0.0012	0.1870
<i>SULF2</i>	0.0018	0.2720	<i>SULT1A1</i>	0.0016	0.2338
<i>MARCH5</i>	0.0019	0.2796	<i>GJB1</i>	0.0021	0.3090
<i>GAPDH</i>	0.0019	0.2826	<i>DLX1</i>	0.0023	0.3475
<i>ERG3' exons 4-5</i>	0.0023	0.3348	<i>PCA3</i>	0.0024	0.3598
<i>TMPRSS2:ERG</i>	0.0023	0.3376	<i>MSMB</i>	0.0027	0.3995
<i>PTN</i>	0.0025	0.3557	<i>MIR4435_1HG</i>	0.0028	0.4081
<i>AURKA</i>	0.0036	0.5179	<i>ITGBL1</i>	0.0033	0.4719
<i>ERG3' exons 6-7</i>	0.0041	0.5766	<i>DPP4</i>	0.0037	0.5275
<i>TDRD</i>	0.0054	0.7671	<i>SULF2</i>	0.0048	0.6756
<i>PALM3</i>	0.0056	0.7794	<i>ZNF577</i>	0.0051	0.7173
<i>ITGBL1</i>	0.0065	0.9059	<i>PPFIA2</i>	0.0052	0.7324
<i>CD10</i>	0.0071	0.9744	<i>PALM3</i>	0.0061	0.8496
<i>TBP</i>	0.0074	0.9897	<i>Met</i>	0.0064	0.8839
<i>KLK2</i>	0.0093	0.9897	<i>PTN</i>	0.0096	0.9961
<i>ZNF577</i>	0.0106	0.9897	<i>MCTP1</i>	0.0105	0.9961
<i>RNF157</i>	0.0107	0.9897	<i>CACNA1D</i>	0.0152	0.9961
<i>PTPRC</i>	0.0107	0.9897	<i>AMACR</i>	0.0174	0.9961
<i>ANPEP</i>	0.0136	0.9897	<i>CP</i>	0.0181	0.9961
<i>NEAT1</i>	0.0137	0.9897	<i>SSTR1</i>	0.0215	0.9961
<i>Met</i>	0.0137	0.9897	<i>GCNT1</i>	0.0221	0.9961
<i>STOM</i>	0.0148	0.9897	<i>PTPRC</i>	0.0310	0.9961
<i>BTG2</i>	0.0150	0.9897	<i>STOM</i>	0.0317	0.9961
<i>AMACR</i>	0.0159	0.9897	<i>IFT57</i>	0.0322	0.9961
<i>MCTP1</i>	0.0175	0.9897	<i>HIST1H2BF</i>	0.0333	0.9961



## 9: APPENDICES

<i>CACNA1D</i>	0.0178	0.9897	<i>RP11_97O12.7</i>	0.0360	0.9961
<i>ALAS1</i>	0.0201	0.9897	<i>STEAP2</i>	0.0398	0.9961
<i>ERG5</i>	0.0223	0.9897	<i>TMCC2</i>	0.0409	0.9961
<i>EN2</i>	0.0225	0.9897	<i>MARCH5</i>	0.0431	0.9961
<i>PPFIA2</i>	0.0270	0.9897			
<i>SIM2.short</i>	0.0274	0.9897			
<i>DLX1</i>	0.0287	0.9897			
<i>ITPR1</i>	0.0289	0.9897			
<i>TMEM86A</i>	0.0296	0.9897			
<i>TMCC2</i>	0.0325	0.9897			
<i>RP11.97O12.7</i>	0.0336	0.9897			
<i>HMBS</i>	0.0343	0.9897			
<i>MSMB</i>	0.0355	0.9897			
<i>IFT57</i>	0.0375	0.9897			
<i>HOXC4</i>	0.0390	0.9897			

9: APPENDICES

Supplementary Table 24 Lasso output for models detecting CB,Ca, Mets trend using KLK2 ratio data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
		<i>ERG3'</i>		<i>ERG3'</i>	
<i>PCA3</i>	<i>0.12</i>	<i>exons 4-5</i>	<i>0.11</i>	<i>exons 4-5</i>	<i>0.11</i>
<i>ERG3' exons 4-5</i>	<i>0.11</i>	<i>PCA3</i>	<i>0.09</i>	<i>PCA3</i>	<i>0.10</i>
<i>APOC1</i>	<i>0.09</i>	<i>HOXC6</i>	<i>0.08</i>	<i>HOXC6</i>	<i>0.08</i>
<i>HOXC6</i>	<i>0.08</i>	<i>APOC1</i>	<i>0.06</i>	<i>APOC1</i>	<i>0.06</i>
<i>SLC12A1</i>	<i>0.05</i>	<i>DLX1</i>	<i>0.03</i>	<i>DLX1</i>	<i>0.04</i>
<i>DLX1</i>	<i>0.04</i>	<i>SLC12A1</i>	<i>0.02</i>	<i>SLC12A1</i>	<i>0.03</i>
<i>TDRD</i>	<i>0.03</i>	<i>TDRD</i>	<i>0.01</i>	<i>TDRD</i>	<i>0.01</i>
				<i>ERG3'</i>	
<i>TMPRSS2:ERG</i>	<i>0.00</i>	<i>cp1</i>	<i>5.00</i>	<i>exons 6-7</i>	<i>0.01</i>
			<i>4.85x10<sup>-14</sup></i>		
<i>ERG3' exons 6-7</i>	<i>0.00</i>	<i>cp2</i>		<i>cp1</i>	<i>5.09</i>
					<i>4.69x10<sup>-14</sup></i>
<i>ZNF577</i>	<i>-0.03</i>			<i>cp2</i>	
<i>GCNT1</i>	<i>-0.04</i>				
<i>CP</i>	<i>-0.09</i>				
<i>cp1</i>	<i>5.09</i>				
	<i>3.99x10<sup>-14</sup></i>				
<i>cp2</i>					

Supplementary Table 25 Lasso output for models detecting CB,Ca, Mets trend using KLK2 adjusted data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
					-
<i>MARCH5</i>	<i>-2.11</i>	<i>MARCH5</i>	<i>-4.03</i>	<i>GABARAPL2</i>	<i>3.18</i>
<i>AMACR</i>	<i>0.09</i>	<i>AMACR</i>	<i>1.17</i>	<i>GAPDH</i>	<i>0.82</i>
<i>APOC1</i>	<i>0.16</i>	<i>ANPEP</i>	<i>0.06</i>	<i>GJB1</i>	<i>1.07</i>
<i>CASKIN1</i>	<i>0.46</i>	<i>APOC1</i>	<i>0.28</i>	<i>HOXC6</i>	<i>1.52</i>
<i>CDC20</i>	<i>-0.15</i>	<i>AURKA</i>	<i>0.05</i>	<i>HPN</i>	<i>1.12</i>
					-
<i>EN2</i>	<i>0.15</i>	<i>DLX1</i>	<i>0.02</i>	<i>KLK4</i>	<i>1.13</i>
<i>ERG5</i>	<i>0.08</i>	<i>EN2</i>	<i>0.33</i>	<i>MFSD2A</i>	<i>0.97</i>
		<i>ERG3' exons 6-7</i>		<i>MIR4435.1HG</i>	
<i>GABARAPL2</i>	<i>-2.25</i>		<i>0.08</i>	<i>.IOC541471</i>	<i>0.79</i>
<i>GJB1</i>	<i>0.56</i>	<i>ERG5</i>	<i>0.23</i>	<i>RAB17</i>	<i>0.23</i>
					-
<i>HIST1H1C</i>	<i>1.00</i>	<i>GABARAPL2</i>	<i>-2.92</i>	<i>SPINK1</i>	<i>0.50</i>
<i>HMBS</i>	<i>0.94</i>	<i>GJB1</i>	<i>0.67</i>	<i>UPK2</i>	<i>0.75</i>
<i>HOXC6</i>	<i>0.92</i>	<i>HMBS</i>	<i>1.97</i>	<i>cp1</i>	<i>5.36</i>
					-
					<i>4.06</i>
<i>HPN</i>	<i>0.32</i>	<i>HOXC6</i>	<i>1.02</i>	<i>cp2</i>	<i>x10<sup>-14</sup></i>

9: APPENDICES

14			
<i>IGFBP3</i>	-0.12	<i>IFT57</i>	1.08
<i>KLK4</i>	-0.68	<i>ITGBL1</i>	-0.46
<i>MFSD2A</i>	0.69	<i>KLK4</i>	-0.93
<i>MIR4435.1HG.1OC5</i>			
<i>41471</i>	0.28	<i>Met</i>	-0.30
<i>MYOF</i>	-1.04	<i>MFSD2A</i>	0.96
<i>MIR4435.1HG.1</i>			
<i>NLRP3</i>	0.05	<i>OC541471</i>	0.35
<i>PALM3</i>	0.11	<i>MYOF</i>	-1.57
<i>PCA3</i>	0.81	<i>PALM3</i>	0.19
<i>PPAP2A</i>	0.21	<i>PCA3</i>	1.49
<i>PTN</i>	-0.25	<i>PPFIA2</i>	-0.65
<i>PTPRC</i>	0.12	<i>PTN</i>	-0.72
<i>RNF157</i>	0.38	<i>PTPRC</i>	0.21
<i>RP11.244H18.1.P71</i>			
<i>2P</i>	-0.23	<i>RNF157</i>	0.83
<i>RP11.244H18.1</i>			
<i>RPL23AP53</i>	0.49	<i>.P712P</i>	-0.88
<i>SFRP4</i>	0.13	<i>RP11.97O12.7</i>	0.24
<i>SIM2.short</i>	1.04	<i>SIM2.short</i>	1.77
<i>SLC12A1</i>	0.07	<i>SLC12A1</i>	0.35
<i>TBP</i>	0.22	<i>STOM</i>	0.11
<i>TDRD</i>	0.10	<i>TBP</i>	1.04
<i>Timp4</i>	-0.66	<i>TDRD</i>	0.17
<i>TMCC2</i>	0.41	<i>TMCC2</i>	0.56
<i>TMEM45B</i>	0.20	<i>TMEM45B</i>	0.36
<i>TMEM86A</i>	0.13	<i>TMEM86A</i>	0.45
<i>TMPRSS2:ER</i>			
<i>TMPRSS2:ERG</i>	0.08	<i>G</i>	0.15
<i>UPK2</i>	0.56	<i>UPK2</i>	0.73
<i>ZNF577</i>	-0.32	<i>ZNF577</i>	-1.06
<i>cp1</i>	5.23	<i>cp1</i>	5.95
<i>cp2</i>	$-3.58 \times 10^{-14}$	<i>cp2</i>	$-3.55 \times 10^{-14}$

Supplementary Table 26 Lasso output for models detecting CB,Ca, Mets trend using KLK3 adjusted data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>
	-2.05		-3.75		-
<i>MARCH5</i>		<i>MARCH5</i>		<i>GABARAPL2</i>	1.96
<i>AMH</i>	0.42	<i>AMACR</i>	0.30	<i>GJB1</i>	0.66
<i>APOC1</i>	0.24	<i>APOC1</i>	0.45	<i>HOXC6</i>	1.57
<i>CP</i>	-0.20	<i>DLX1</i>	0.05	<i>HPN</i>	0.90
					-
<i>EN2</i>	0.09	<i>EN2</i>	0.32	<i>KLK4</i>	1.24
<i>ERG3' exons 4-5</i>	0.09	<i>ERG3' exons 4-5</i>	0.37	<i>MIR4435.1HG</i>	
				<i>.IOC541471</i>	0.18

## 9: APPENDICES

<i>ERG3' exons 6-7</i>	<b>0.00</b>	<i>ERG5</i>	<b>0.22</b>	<i>RAB17</i>	<b>0.66</b>
<i>ERG5</i>	<b>0.05</b>	<i>GABARAPL2</i>	<b>-1.42</b>	<i>SPINK1</i>	<b>0.68</b>
<i>GABARAPL2</i>	<b>-1.46</b>	<i>GJB1</i>	<b>0.48</b>	<i>UPK2</i>	<b>1.15</b>
<i>GJB1</i>	<b>0.18</b>	<i>HMBS</i>	<b>1.05</b>	<i>cp1</i>	<b>4.65</b>
					<b>6.23</b>
<i>HIST1H1C</i>	<b>0.40</b>	<i>HOXC4</i>	<b>-0.40</b>	<i>cp2</i>	<b><math>7.49 \times 10^{14}</math></b>
<i>HMBS</i>	<b>0.41</b>	<i>HOXC6</i>	<b>1.00</b>		
<i>HOXC6</i>	<b>0.79</b>	<i>IFT57</i>	<b>0.30</b>		
<i>HPN</i>	<b>0.10</b>	<i>ITGBL1</i>	<b>-0.59</b>		
<i>IGFBP3</i>	<b>-0.35</b>	<i>ITPR1</i>	<b>0.29</b>		
<i>KLK4</i>	<b>-0.46</b>	<i>KLK2</i>	<b>0.06</b>		
<i>MFSD2A</i>	<b>0.33</b>	<i>KLK4</i>	<b>-0.42</b>		
<i>MIR4435.1HG.1OC5</i>					
<i>41471</i>	<b>0.02</b>	<i>Met</i>	<b>-0.12</b>		
<i>MMP25</i>	<b>0.15</b>	<i>MFSD2A</i>	<b>0.51</b>		
		<i>MIR4435.1HG.1</i>			
<i>MXII</i>	<b>-0.05</b>	<i>OC541471</i>	<b>0.21</b>		
<i>MYOF</i>	<b>-0.41</b>	<i>MYOF</i>	<b>-0.63</b>		
<i>PALM3</i>	<b>0.06</b>	<i>NEAT1</i>	<b>0.22</b>		
<i>PCA3</i>	<b>1.17</b>	<i>PALM3</i>	<b>0.09</b>		
<i>PSTPIP1</i>	<b>0.05</b>	<i>PCA3</i>	<b>2.10</b>		
<i>PTN</i>	<b>-0.26</b>	<i>PPF1A2</i>	<b>-0.46</b>		
<i>RNF157</i>	<b>0.27</b>	<i>PTN</i>	<b>-0.75</b>		
<i>RP11.244H18.1.P71</i>					
<i>2P</i>	<b>-0.04</b>	<i>PTPRC</i>	<b>0.03</b>		
<i>RPL23AP53</i>	<b>0.71</b>	<i>RNF157</i>	<b>0.59</b>		
		<i>RP11.244H18.1</i>			
<i>SIM2.short</i>	<b>0.87</b>	<i>.P712P</i>	<b>-0.45</b>		
<i>SLC12A1</i>	<b>0.23</b>	<i>RP11.97O12.7</i>	<b>0.37</b>		
<i>TBP</i>	<b>0.65</b>	<i>SIM2.short</i>	<b>1.56</b>		
<i>TDRD</i>	<b>0.05</b>	<i>SLC12A1</i>	<b>0.49</b>		
<i>Timp4</i>	<b>-0.29</b>	<i>STOM</i>	<b>0.02</b>		
<i>TMPRSS2-ERG</i>	<b>0.02</b>	<i>TBP</i>	<b>1.70</b>		
<i>UPK2</i>	<b>0.84</b>	<i>TMCC2</i>	<b>0.23</b>		
<i>cp1</i>	<b>4.52</b>	<i>TMEM45B</i>	<b>0.24</b>		
	<b><math>7.49 \times 10^{14}</math></b>	<i>TMEM86A</i>	<b>0.32</b>		
<i>cp2</i>		<i>UPK2</i>	<b>0.91</b>		
		<i>ZNF577</i>	<b>-0.55</b>		
		<i>cp1</i>	<b>5.18</b>		
		<i>cp2</i>	<b><math>9.67 \times 10^{14}</math></b>		

Supplementary Table 27 Lasso output for models detecting CB,Ca, Mets trend using HK normalised data.

<i>All Transcripts</i>		<i>Significant Transcripts</i>		<i>Multiple testing corrected Transcripts</i>	
<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>	<i>Transcript</i>	<i>Beta</i>

9: APPENDICES

<i>MARCH5</i>	-0.02	<i>MARCH5</i>	-0.15	<i>APOC1</i>	0.12
			0.07	<i>ERG3' exons</i>	
<i>ACTR5</i>	-0.02	<i>AMACR</i>	0.07	4-5	0.17
<i>AMACR</i>	0.11	<i>APOC1</i>	0.20	<i>HOXC6</i>	0.12
<i>AMH</i>	0.15	<i>CACNA1D</i>	-0.34	<i>HPN</i>	0.05
				-	
<i>ANKRD34B</i>	-0.05	<i>CP</i>	-0.33	<i>KLK2</i>	0.29
<i>APOC1</i>	0.22	<i>DLX1</i>	0.09	<i>SLC12A1</i>	0.08
		<i>ERG3' exons 4-</i>			
<i>AR exon 9</i>	-0.07	5	0.06	<i>SPINK1</i>	0.04
		<i>ERG3' exons 6-</i>		<i>TMPRSS2:ER</i>	
<i>AURKA</i>	0.02	7	0.04	G	0.03
<i>CACNA1D</i>	-0.30	<i>GABARAPL2</i>	-0.41	<i>UPK2</i>	0.00
<i>CASKINI</i>	0.18	<i>GCNT1</i>	-0.13	<i>cp1</i>	5.27
					3.65
<i>CD10</i>	-0.01	<i>GJB1</i>	0.07	<i>cp2</i>	$\times 10^{14}$
<i>CDC20</i>	-0.04	<i>HOXC6</i>	0.13		
<i>CP</i>	-0.30	<i>IFT57</i>	0.10		
<i>DLX1</i>	0.09	<i>ITGBL1</i>	-0.10		
<i>EN2</i>	0.03	<i>KLK2</i>	-0.11		
<i>ERG3' exons 4-5</i>	0.02	<i>KLK4</i>	-0.24		
<i>ERG3' exons 6-7</i>	0.06	<i>MCTP1</i>	0.05		
<i>GABARAPL2</i>	-0.87	<i>Met</i>	-0.01		
<i>GCNT1</i>	-0.13	<i>MIR4435 1HG</i>	0.09		
<i>GJB1</i>	0.08	<i>MSMB</i>	0.07		
<i>GOLM1</i>	-0.02	<i>PALM3</i>	0.13		
<i>HIST1H1C</i>	0.41	<i>PCA3</i>	0.21		
<i>HMBS</i>	0.24	<i>PTN</i>	-0.01		
<i>HOXC6</i>	0.17	<i>SLC12A1</i>	0.05		
<i>IGFBP3</i>	-0.01	<i>SSTR1</i>	0.11		
<i>ISX</i>	0.09	<i>STOM</i>	0.03		
<i>ITGBL1</i>	-0.05	<i>SULF2</i>	0.02		
<i>KLK2</i>	-0.01	<i>TDRD</i>	0.03		
<i>KLK4</i>	-0.23	<i>TMCC2</i>	0.28		
<i>LASS1</i>	-0.11	<i>TMEM45B</i>	0.30		
<i>MCTP1</i>	0.01	<i>ZNF577</i>	-0.17		
<i>MED4</i>	-0.03	<i>cp1</i>	6.15		
<i>MEMO1</i>	-0.08	<i>cp2</i>	$3.32 \times 10^{-14}$		
<i>MEX3A</i>	0.03				
<i>MGAT5B</i>	0.15				
<i>MIC1</i>	-0.03				
<i>MIR146A</i>	-0.10				
<i>MIR4435 1HG</i>	0.09				
<i>MMP25</i>	0.11				
<i>MMP26</i>	-0.06				
<i>MXI1</i>	-0.07				
<i>NEAT1</i>	0.04				
<i>NLRP3</i>	0.12				
<i>OR52A2</i>	-0.08				
<i>PALM3</i>	0.09				

9: APPENDICES

<i>PCA3</i>	<b>0.15</b>
<i>PDLIM5</i>	<b>-0.17</b>
<i>PPAP2A</i>	<b>0.42</b>
<i>PPFIA2</i>	<b>0.03</b>
<i>PPP1R12B</i>	<b>-0.19</b>
<i>RAB17</i>	<b>0.00</b>
<i>RNF157</i>	<b>0.01</b>
<i>RPL23AP53</i>	<b>0.12</b>
<i>SChLAP1</i>	<b>0.00</b>
<i>SEC61A1</i>	<b>-0.02</b>
<i>SFRP4</i>	<b>0.07</b>
<i>SIRT1</i>	<b>-0.02</b>
<i>SLC12A1</i>	<b>0.04</b>
<i>SLC43A1</i>	<b>-0.12</b>
<i>SSPO</i>	<b>-0.30</b>
<i>SSTR1</i>	<b>0.07</b>
<i>SULT1A1</i>	<b>0.01</b>
<i>TDRD</i>	<b>0.03</b>
<i>Timp4</i>	<b>-0.13</b>
<i>TMCC2</i>	<b>0.20</b>
<i>TMEM45B</i>	<b>0.29</b>
<i>TMEM86A</i>	<b>0.07</b>
<i>TMPRSS2:ERG</i>	<b>0.03</b>
<i>ZNF577</i>	<b>-0.11</b>
<i>cp1</i>	<b>6.38</b>
<i>cp2</i>	<b>2.17x10<sup>14</sup></b>

**6.17 Looking for Housekeepers****Supplementary Table 28 Top twenty transcripts with the lowest variance in cell sediment urine fraction data**

<i>Transcript</i>	<i>Variance</i>
<i>MNX1</i>	<i>0.94</i>
<i>TWIST1</i>	<i>0.95</i>
<i>SSPO</i>	<i>1.03</i>
<i>SLC4A1 S</i>	<i>1.04</i>
<i>COL9A2</i>	<i>1.09</i>
<i>TERT</i>	<i>1.12</i>
<i>SSTR1</i>	<i>1.14</i>
<i>ABCB9</i>	<i>1.15</i>
<i>CASKIN1</i>	<i>1.27</i>
<i>MMP11</i>	<i>1.28</i>
<i>TMCC2</i>	<i>1.33</i>
<i>HIST1H1C</i>	<i>1.42</i>
<i>AMH</i>	<i>1.42</i>
<i>ISX</i>	<i>1.50</i>
<i>RPS11</i>	<i>1.50</i>
<i>AATF</i>	<i>1.51</i>
<i>HIST1H1E</i>	<i>1.53</i>
<i>VAX2</i>	<i>1.55</i>
<i>ARHGEF25</i>	<i>1.66</i>
<i>FDPS</i>	<i>1.66</i>

**Supplementary Table 29 Top twenty transcripts with the lowest IQR in cell sediment urine fraction data**

<i>Transcript</i>	<i>IQR</i>
<i>SSPO</i>	<i>0.95</i>
<i>RPS11</i>	<i>1.01</i>
<i>SLC4A1 S</i>	<i>1.03</i>
<i>TWIST1</i>	<i>1.04</i>
<i>ABCB9</i>	<i>1.05</i>
<i>HIST1H1E</i>	<i>1.05</i>
<i>B2M</i>	<i>1.08</i>
<i>VAX2</i>	<i>1.13</i>
<i>CASKIN1</i>	<i>1.13</i>
<i>RIOK3</i>	<i>1.15</i>
<i>CADPS</i>	<i>1.16</i>
<i>RP11_97012.7</i>	<i>1.18</i>
<i>COL9A2</i>	<i>1.18</i>
<i>MMP26</i>	<i>1.20</i>
<i>MGAT5B</i>	<i>1.20</i>
<i>ISX</i>	<i>1.21</i>
<i>TFDP1</i>	<i>1.21</i>

9: APPENDICES

<i>MNX1</i>	1.21
<i>SSTR1</i>	1.22
<i>RPL18A</i>	1.23

Supplementary Table 30 Comparing the expression between all clinical categories using Tukey test, looking for potential house keeping transcripts.

<i>Transcript</i>	<i>Significan t</i>	<i>Transcript</i>	<i>Significan t</i>	<i>Transcript</i>	<i>Significan t</i>
<i>MARCH5</i>	0	<i>PTN</i>	0	<i>OR52A2</i>	2
<i>ABCB9</i>	0	<i>PVT1</i>	0	<i>PCA3</i>	2
<i>ACTR5</i>	0	<i>RAB17</i>	0	<i>PDLIM5</i>	2
<i>AMACR</i>	0	<i>RNF157</i>	0	<i>PSTPIP1</i>	2
<i>AMH</i>	0	<i>RP11_97012.7</i>	0	<i>SIM2 long</i>	2
<i>AR exon 9</i>	0	<i>RPL18A</i>	0	<i>SIM2 short</i>	2
<i>AR exons 4-8</i>	0	<i>RPL23AP53</i>	0	<i>SLC12A1</i>	2
<i>ARHGEF25</i>	0	<i>RPLP2</i>	0	<i>SNORA20</i>	2
<i>BRAF</i>	0	<i>RPS10</i>	0	<i>ST6GALNAC1</i>	2
<i>CAMK2N2</i>	0	<i>SACMIL</i>	0	<i>TMCC2</i>	2
<i>CASKIN1</i>	0	<i>SChLAP1</i>	0	<i>TMEM86A</i>	2
<i>CDC20</i>	0	<i>SLC4A1 S</i>	0	<i>AGR2</i>	3
<i>CDC37L1</i>	0	<i>SMAP1 exons7-8</i>	0	<i>BTG2</i>	3
<i>CDKN3</i>	0	<i>SMIM1</i>	0	<i>FOLH1</i>	3
<i>CLU</i>	0	<i>SPINK1</i>	0	<i>GABARAPL2</i>	3
<i>COL10A1</i>	0	<i>SPON2</i>	0	<i>MAPK8IP2</i>	3
<i>CP</i>	0	<i>STEAP2</i>	0	<i>SLC43A1</i>	3
<i>MIATNB</i>	0	<i>STEAP4</i>	0	<i>SNCA</i>	3
<i>DLX1</i>	0	<i>SYNM</i>	0	<i>TDRD</i>	3
<i>ERG3' exons 4-5</i>	0	<i>TERT</i>	0	<i>ANPEP</i>	4
<i>ERG5'</i>	0	<i>TFDP1</i>	0	<i>B2M</i>	4
<i>FDPS</i>	0	<i>Timp4</i>	0	<i>CLIC2</i>	4
<i>GOLM1</i>	0	<i>TMEM47</i>	0	<i>EIF2D</i>	4
<i>HIST1H1C</i>	0	<i>TRPM4</i>	0	<i>GAPDH</i>	4
<i>HIST1H1E</i>	0	<i>TWIST1</i>	0	<i>LASS1</i>	4
<i>HIST1H2B F</i>	0	<i>VAX2</i>	0	<i>MIR146A</i>	4
<i>HIST3H2A</i>	0	<i>VPS13A</i>	0	<i>MIR4435_1HG</i>	4
<i>HMBS</i>	0	<i>ZNF577</i>	0	<i>NLRP3</i>	4
<i>HOXC4</i>	0	<i>ALAS1</i>	1	<i>SRSF3</i>	4
<i>IFT57</i>	0	<i>ANKRD34B</i>	1	<i>SSPO</i>	4
<i>IGFBP3</i>	0	<i>AURKA</i>	1	<i>TERF2IP</i>	4
<i>IMPDH2</i>	0	<i>CKAP2L</i>	1	<i>TMPRSS2:ER G fusion</i>	4
<i>ITGBL1</i>	0	<i>COL9A2</i>	1	<i>AATF</i>	5
<i>KLK2</i>	0	<i>DNAH5</i>	1	<i>B4GALNT4</i>	5
<i>KLK3 exons 2-3</i>	0	<i>GJB1</i>	1	<i>CADPS</i>	5



9: APPENDICES

<i>KLK4</i>	0	<i>KLK3 exons 1-2</i>	1	<i>CAMKK2</i>	5
<i>LBH</i>	0	<i>MED4</i>	1	<i>CCDC88B</i>	5
<i>MAK</i>	0	<i>MEMO1</i>	1	<i>HPN</i>	5
<i>MCM7</i>	0	<i>MMP26</i>	1	<i>ISX</i>	5
<i>MDK</i>	0	<i>MNX1</i>	1	<i>ITPR1</i>	5
<i>Met</i>	0	<i>NAALADL2</i>	1	<i>MFSD2A</i>	5
<i>MEX3A</i>	0	<i>P712P</i>	1	<i>MMP25</i>	5
<i>MGAT5B</i>	0	<i>RPS11</i>	1	<i>SEC61A1</i>	5
<i>MIC1</i>	0	<i>SFRP4</i>	1	<i>HPRT</i>	6
<i>MKi67</i>	0	<i>SSTR1</i>	1	<i>RIOK3</i>	6
<i>MMP11</i>	0	<i>SULT1A1</i>	1	<i>SERPINB5</i>	6
<i>MSMB</i>	0	<i>TBP</i>	1	<i>SIRT1</i>	6
<i>MYOF</i>	0	<i>TMEM45B</i>	1	<i>ERG3' exons 6-7</i>	7
<i>NKAIN1</i>	0	<i>UPK2</i>	1	<i>HOXC6</i>	7
<i>OGT</i>	0	<i>CACNA1D</i>	2	<i>STOM</i>	9
<i>PALM3</i>	0	<i>CD10</i>	2	<i>APOC1</i>	11
<i>PCSK6</i>	0	<i>DPP4</i>	2	<i>MCTP1</i>	11
<i>PECI</i>	0	<i>EN2</i>	2	<i>NEAT1</i>	11
<i>PPAP2A</i>	0	<i>GCNT1</i>	2	<i>PTPRC</i>	12
<i>PPFIA2</i>	0	<i>HIST1H2BG</i>	2	<i>SULF2</i>	12
<i>PPP1R12B</i>	0	<i>MXII</i>	2		

**6.18 Cancer Vs CB**

Supplementary Table 31 Transcripts that have significant differential expression between CB and cancer samples (L, I, H) in the baseline normalised NanoString data.

<i>Transcript</i>	<i>MWU</i>		<i>glm</i>		<i>Log<sub>2</sub>(FC)</i>
	<i>p-value</i>	<i>Adjusted p-value</i>	<i>p-value</i>	<i>Adjusted p-value</i>	
<i>HOXC6</i>	0.0002	0.024	0.0014	0.2049	1.64
<i>ERG3' exons 6-7</i>	2.84x10 <sup>-07</sup>	4.74x10 <sup>-05</sup>	0.0008	0.128	1.38
<i>TMPRSS2:ERG</i>	4.52x10 <sup>-05</sup>	0.0069	0.0013	0.1979	1.31
<i>SLC43A1</i>	0.0003	0.0406	0.0019	0.2745	1.17
<i>CLIC2</i>	2.66x10 <sup>-05</sup>	0.0042	0.001	0.1645	1.05
<i>B4GALNT4</i>	3.38x10 <sup>-05</sup>	0.0053	0.0012	0.1807	1.04
<i>CADPS</i>	1.37x10 <sup>-05</sup>	0.0022	0.0004	0.0682	1.04
<i>CKAP2L</i>	0.0116	1	0.0033	0.4318	1.01
<i>HPN</i>	7.04x10 <sup>-05</sup>	0.0103	0.0006	0.1041	0.97
<i>LASS1</i>	0.0002	0.022	0.0011	0.1703	0.97
<i>TDRD</i>	0.0002	0.022	0.0047	0.5935	0.97
<i>SFRP4</i>	0.0004	0.0478	0.0031	0.4076	0.87
<i>OR52A2</i>	0.0343	1	0.0284	0.9777	0.85
<i>ANKRD34B</i>	0.0013	0.1706	0.0093	0.9777	0.83
<i>MAPK8IP2</i>	0.0063	0.6758	0.0067	0.7955	0.8
<i>PCA3</i>	0.0004	0.0565	0.0019	0.2745	0.8
<i>CDKN3</i>	0.024	1	0.011	0.9777	0.76

## 9: APPENDICES

<i>ERG5'</i>	<b>0.0016</b>	<b>0.2005</b>	<b>0.009</b>	<b>0.9777</b>	<b>0.68</b>
<i>MFSD2A</i>	<b>1.32x10<sup>-05</sup></b>	<b>0.0021</b>	<b>0.001</b>	<b>0.1645</b>	<b>0.66</b>
<i>MMP25</i>	<b>7.80x10<sup>-05</sup></b>	<b>0.0113</b>	<b>0.0008</b>	<b>0.1278</b>	<b>0.66</b>
<i>APOC1</i>	<b>1.85x10<sup>-06</sup></b>	<b>0.0003</b>	<b>0.0004</b>	<b>0.0586</b>	<b>0.65</b>
<i>TMCC2</i>	<b>0.0126</b>	<b>1</b>	<b>0.0075</b>	<b>0.8811</b>	<b>0.65</b>
<i>NKAIN1</i>	<b>0.0379</b>	<b>1</b>	<b>0.0429</b>	<b>0.9777</b>	<b>0.62</b>
<i>SIM2 long</i>	<b>3.72x10<sup>-05</sup></b>	<b>0.0057</b>	<b>0.0031</b>	<b>0.4076</b>	<b>0.62</b>
<i>MCTP1</i>	<b>3.97x10<sup>-07</sup></b>	<b>6.59x10<sup>-05</sup></b>	<b>0.0002</b>	<b>0.0406</b>	<b>0.61</b>
<i>ISX</i>	<b>0.0007</b>	<b>0.086</b>	<b>0.0024</b>	<b>0.3365</b>	<b>0.6</b>
<i>X?</i>	<b>0.0032</b>	<b>0.3694</b>			<b>0.59</b>
<i>MMP26</i>	<b>0.0448</b>	<b>1</b>			<b>0.56</b>
<i>AMH</i>	<b>0.0032</b>	<b>0.3694</b>	<b>0.0335</b>	<b>0.9777</b>	<b>0.55</b>
<i>SLC12A1</i>	<b>0.0014</b>	<b>0.1798</b>	<b>0.0064</b>	<b>0.7769</b>	<b>0.55</b>
<i>SULF2</i>	<b>9.18x10<sup>-06</sup></b>	<b>0.0015</b>	<b>0.0011</b>	<b>0.1754</b>	<b>0.55</b>
<i>CCDC88B</i>	<b>6.34x10<sup>-05</sup></b>	<b>0.0094</b>	<b>0.0012</b>	<b>0.1785</b>	<b>0.54</b>
<i>NLRP3</i>	<b>0.0024</b>	<b>0.29</b>	<b>0.002</b>	<b>0.2817</b>	<b>0.54</b>
<i>UPK2</i>	<b>0.0071</b>	<b>0.7532</b>	<b>0.0147</b>	<b>0.9777</b>	<b>-0.54</b>
<i>TMEM86A</i>	<b>0.0001</b>	<b>0.0202</b>	<b>0.0019</b>	<b>0.2742</b>	<b>0.53</b>
<i>CAMKK2</i>	<b>2.13x10<sup>-06</sup></b>	<b>0.0003</b>	<b>0.0005</b>	<b>0.0859</b>	<b>0.51</b>
<i>FOLH1</i>	<b>0.013</b>	<b>1</b>	<b>0.0121</b>	<b>0.9777</b>	<b>0.49</b>
<i>ANPEP</i>	<b>0.0011</b>	<b>0.1472</b>	<b>0.003</b>	<b>0.4076</b>	<b>0.46</b>
<i>SRSF3</i>	<b>0.0002</b>	<b>0.0311</b>	<b>0.0038</b>	<b>0.4879</b>	<b>0.45</b>
<i>MIR146A</i>	<b>0.0019</b>	<b>0.235</b>	<b>0.0023</b>	<b>0.3187</b>	<b>0.44</b>
<i>GCNT1</i>	<b>0.003</b>	<b>0.353</b>	<b>0.0035</b>	<b>0.4645</b>	<b>0.43</b>
<i>SIRT1</i>	<b>5.71x10<sup>-05</sup></b>	<b>0.0085</b>	<b>0.0012</b>	<b>0.1785</b>	<b>0.41</b>
<i>SERPINB5</i>	<b>0.031</b>	<b>1</b>			<b>-0.4</b>
<i>NAALADL2</i>	<b>0.0059</b>	<b>0.6513</b>	<b>0.0308</b>	<b>0.9777</b>	<b>-0.38</b>
<i>SNORA20</i>	<b>0.0081</b>	<b>0.8455</b>	<b>0.0206</b>	<b>0.9777</b>	<b>0.37</b>
<i>CDC20</i>	<b>0.0063</b>	<b>0.6758</b>	<b>0.0117</b>	<b>0.9777</b>	<b>0.35</b>
<i>TMEM45B</i>	<b>0.0253</b>	<b>1</b>			<b>-0.35</b>
<i>AATF</i>	<b>5.14x10<sup>-05</sup></b>	<b>0.0077</b>	<b>0.0014</b>	<b>0.2012</b>	<b>0.34</b>
<i>IGFBP3</i>	<b>0.0067</b>	<b>0.7136</b>			<b>-0.34</b>
<i>AURKA</i>	<b>0.0016</b>	<b>0.1909</b>	<b>0.0056</b>	<b>0.6828</b>	<b>0.33</b>
<i>CD10</i>	<b>0.0014</b>	<b>0.1798</b>	<b>0.0053</b>	<b>0.6649</b>	<b>0.33</b>
<i>PTPRC</i>	<b>4.62x10<sup>-05</sup></b>	<b>0.007</b>	<b>0.0014</b>	<b>0.2012</b>	<b>0.33</b>
<i>SSTR1</i>	<b>0.0164</b>	<b>1</b>	<b>0.0177</b>	<b>0.9777</b>	<b>0.33</b>
<i>SEC61A1</i>	<b>0.0002</b>	<b>0.0285</b>	<b>0.0055</b>	<b>0.682</b>	<b>0.32</b>
<i>SIM2 short</i>	<b>0.0193</b>	<b>1</b>	<b>0.0164</b>	<b>0.9777</b>	<b>0.32</b>
<i>SNCA</i>	<b>0.0016</b>	<b>0.1909</b>	<b>0.007</b>	<b>0.8231</b>	<b>0.32</b>
<i>MMP11</i>	<b>0.013</b>	<b>1</b>	<b>0.0222</b>	<b>0.9777</b>	<b>0.31</b>
<i>SPINK1</i>	<b>0.0032</b>	<b>0.3694</b>	<b>0.0133</b>	<b>0.9777</b>	<b>-0.31</b>
<i>HPRT</i>	<b>0.0005</b>	<b>0.0666</b>	<b>0.0067</b>	<b>0.7955</b>	<b>0.29</b>
<i>PSTPIP1</i>	<b>0.046</b>	<b>1</b>	<b>0.0226</b>	<b>0.9777</b>	<b>0.28</b>
<i>HOXC4</i>	<b>0.0356</b>	<b>1</b>			<b>0.26</b>
<i>RIOK3</i>	<b>1.03x10<sup>-06</sup></b>	<b>0.0002</b>	<b>0.0013</b>	<b>0.1979</b>	<b>0.26</b>
	<b>0.0469</b>	<b>1</b>			<b>0.25</b>
	<b>0.0266</b>	<b>1</b>			<b>0.24</b>
<i>EN2</i>	<b>0.0192</b>	<b>1</b>	<b>0.0282</b>	<b>0.9777</b>	<b>0.24</b>
<i>MEX3A</i>	<b>0.0224</b>	<b>1</b>	<b>0.0252</b>	<b>0.9777</b>	<b>0.24</b>
<i>CACNA1D</i>	<b>0.0014</b>	<b>0.1798</b>	<b>0.0054</b>	<b>0.681</b>	<b>0.23</b>
<i>MIR4435 IHG</i>	<b>3.72x10<sup>-05</sup></b>	<b>0.0057</b>	<b>0.0026</b>	<b>0.3533</b>	<b>0.23</b>

9: APPENDICES

<i>MXI1</i>	<b>0.0001</b>	<b>0.0168</b>	<b>0.0026</b>	<b>0.3537</b>	<b>0.23</b>
<i>DPP4</i>	<b>0.0048</b>	<b>0.54</b>	<b>0.0108</b>	<b>0.9777</b>	<b>0.22</b>
<i>CASKIN1</i>	<b>0.0343</b>	<b>1</b>	<b>0.0315</b>	<b>0.9777</b>	<b>0.21</b>
<i>HIST3H2A</i>	<b>0.0193</b>	<b>1</b>			<b>0.21</b>
<i>ITPR1</i>	<b>4.15x10<sup>-05</sup></b>	<b>0.0063</b>	<b>0.0019</b>	<b>0.2745</b>	<b>0.21</b>
<i>NEAT1</i>	<b>1.89x10<sup>-05</sup></b>	<b>0.003</b>	<b>0.0009</b>	<b>0.1377</b>	<b>0.21</b>
<i>STOM</i>	<b>0.0116</b>	<b>1</b>	<b>0.0219</b>	<b>0.9777</b>	<b>0.21</b>
<i>PDLIM5</i>	<b>0.0007</b>	<b>0.0916</b>	<b>0.0046</b>	<b>0.5856</b>	<b>0.2</b>
<i>BTG2</i>	<b>0.0037</b>	<b>0.4198</b>	<b>0.0173</b>	<b>0.9777</b>	<b>0.19</b>
<i>GABARAPL2</i>	<b>0.0003</b>	<b>0.044</b>	<b>0.0025</b>	<b>0.3502</b>	<b>0.19</b>
<i>HIST1H2BG</i>	<b>0.0048</b>	<b>0.54</b>	<b>0.0064</b>	<b>0.7753</b>	<b>0.19</b>
<i>MAK</i>	<b>0.0183</b>	<b>1</b>			<b>0.19</b>
<i>EIF2D</i>	<b>0.0026</b>	<b>0.309</b>	<b>0.0125</b>	<b>0.9777</b>	<b>0.18</b>
<i>MGAT5B</i>			<b>0.0446</b>	<b>0.9777</b>	<b>0.18</b>
<i>TBP</i>	<b>0.0063</b>	<b>0.6758</b>	<b>0.0358</b>	<b>0.9777</b>	<b>0.18</b>
<i>TWIST1</i>	<b>0.0138</b>	<b>1</b>			<b>0.18</b>
<i>MED4</i>	<b>0.0146</b>	<b>1</b>	<b>0.0361</b>	<b>0.9777</b>	<b>0.17</b>
<i>TERF2IP</i>	<b>7.80x10<sup>-05</sup></b>	<b>0.0113</b>	<b>0.0019</b>	<b>0.2745</b>	<b>0.17</b>
<i>GAPDH</i>	<b>2.98x10<sup>-05</sup></b>	<b>0.0047</b>	<b>0.0007</b>	<b>0.109</b>	<b>0.16</b>
<i>ACTR5</i>	<b>0.0164</b>	<b>1</b>	<b>0.0123</b>	<b>0.9777</b>	<b>0.15</b>
<i>B2M</i>	<b>0.0002</b>	<b>0.0285</b>	<b>0.0044</b>	<b>0.5721</b>	<b>0.14</b>
	<b>0.0438</b>	<b>1</b>			<b>0.11</b>
<i>SACMIL</i>	<b>0.0266</b>	<b>1</b>	<b>0.0347</b>	<b>0.9777</b>	<b>0.11</b>
<i>RPL18A</i>	<b>0.0123</b>	<b>1</b>			<b>0.1</b>
	<b>0.0138</b>	<b>1</b>			<b>0.09</b>
<i>STEAP4</i>			<b>0.0474</b>	<b>0.9777</b>	<b>0.09</b>
<i>HIST1H2BF</i>			<b>0.0424</b>	<b>0.9777</b>	<b>0.08</b>
<i>MEMO1</i>	<b>0.0361</b>	<b>1</b>	<b>0.0357</b>	<b>0.9777</b>	<b>0.08</b>
<i>RPS11</i>	<b>0.0173</b>	<b>1</b>			<b>0.08</b>
<i>HMBS</i>			<b>0.0391</b>	<b>0.9777</b>	<b>0.06</b>
<i>SLC4A1 S</i>			<b>0.0223</b>	<b>0.9777</b>	<b>0.03</b>
<i>SMAP1 exons 7-8</i>			<b>0.0458</b>	<b>0.9777</b>	<b>0.03</b>

Supplementary Table 32 Transcripts that have significant differential expression between CB and cancer samples (L, I, H) in the *KLK2* ratio NanoString data.

<i>Transcript</i>	<i>MWU</i>		<i>glm</i>		<i>Log<sub>2</sub>(FC)</i>
	<i>p-value</i>	<i>Adjusted p-value</i>	<i>p-value</i>	<i>Adjusted p-value</i>	
<i>HOXC6</i>	<b>6.80x10<sup>-05</sup></b>	<b>0.01</b>	<b>0.004</b>	<b>0.63</b>	<b>0.21</b>
<i>ERG3' exons 6-7</i>	<b>7.80x10<sup>-05</sup></b>	<b>0.01</b>	<b>0.001</b>	<b>0.24</b>	<b>0.18</b>
<i>TDRD</i>	<b>0.0004</b>	<b>0.06</b>	<b>0.004</b>	<b>0.72</b>	<b>0.18</b>
<i>SLC43A1</i>	<b>0.002</b>	<b>0.32</b>			<b>0.17</b>
<i>CADPS</i>	<b>0.004</b>	<b>0.67</b>	<b>0.01</b>	<b>1</b>	<b>0.16</b>
<i>ERG5'</i>	<b>0.01</b>	<b>0.99</b>			<b>0.15</b>
<i>B4GALNT4</i>	<b>0.01</b>	<b>0.87</b>			<b>0.14</b>
<i>SLC12A1</i>	<b>0.003</b>	<b>0.54</b>	<b>0.03</b>	<b>1</b>	<b>0.13</b>
<i>TMCC2</i>	<b>0.05</b>	<b>0.99</b>	<b>0.05</b>	<b>1</b>	<b>0.13</b>

9: APPENDICES

<i>TMPRSS2-ERG</i>	<b>0.001</b>	<b>0.17</b>	<b>0.01</b>	<b>1</b>	<b>0.13</b>
<i>CKAP2L</i>	<b>0.02</b>	<b>0.99</b>	<b>0.02</b>	<b>1</b>	<b>0.12</b>
<i>MFSD2A</i>	<b>0.01</b>	<b>0.99</b>	<b>0.02</b>	<b>1</b>	<b>0.12</b>
<i>CLIC2</i>	<b>0.004</b>	<b>0.58</b>	<b>0.01</b>	<b>1</b>	<b>0.11</b>
<i>LASS1</i>	<b>0.01</b>	<b>0.89</b>	<b>0.01</b>	<b>1</b>	<b>0.11</b>
<i>MMP25</i>	<b>0.01</b>	<b>0.99</b>	<b>0.02</b>	<b>1</b>	<b>0.11</b>
<i>PCA3</i>	<b>3.72x10<sup>-05</sup></b>	<b>0.01</b>	<b>0.003</b>	<b>0.41</b>	<b>0.1</b>
<i>ANKRD34B</i>	<b>0.04</b>	<b>0.99</b>			<b>0.09</b>
<i>HPN</i>	<b>0.001</b>	<b>0.23</b>	<b>0.01</b>	<b>1</b>	<b>0.09</b>
<i>TMEM86A</i>	<b>0.01</b>	<b>0.99</b>	<b>0.03</b>	<b>1</b>	<b>0.09</b>
<i>NAALADL2</i>	<b>0.03</b>	<b>0.99</b>	<b>0.04</b>	<b>1</b>	<b>-0.09</b>
<i>UPK2</i>	<b>0.03</b>	<b>0.99</b>			<b>-0.09</b>
<i>APOC1</i>	<b>0.004</b>	<b>0.65</b>			<b>0.08</b>
<i>CCDC88B</i>	<b>0.01</b>	<b>0.99</b>	<b>0.05</b>	<b>1</b>	<b>0.08</b>
<i>ST6GALNAC1</i>	<b>0.03</b>	<b>0.99</b>			<b>-0.08</b>
<i>ISX</i>	<b>0.01</b>	<b>0.99</b>			<b>0.07</b>
<i>MCTP1</i>	<b>0.01</b>	<b>0.99</b>	<b>0.03</b>	<b>1</b>	<b>0.07</b>
<i>MIR146A</i>	<b>0.02</b>	<b>0.99</b>			<b>0.07</b>
<i>NLRP3</i>			<b>0.03</b>	<b>1</b>	<b>0.07</b>
<i>SULF2</i>	<b>0.01</b>	<b>0.99</b>			<b>0.07</b>
<i>SERPINB5</i>	<b>0.03</b>	<b>0.99</b>			<b>-0.06</b>
<i>SFRP4</i>	<b>0.04</b>	<b>0.99</b>			<b>0.06</b>
<i>FOLH1</i>	<b>0.03</b>	<b>0.99</b>			<b>0.05</b>
<i>OR52A2</i>			<b>0.03</b>	<b>1</b>	<b>0.05</b>
<i>SIM2.long</i>	<b>0.003</b>	<b>0.51</b>	<b>0.01</b>	<b>1</b>	<b>0.05</b>
<i>CAMKK2</i>	<b>0.004</b>	<b>0.62</b>	<b>0.04</b>	<b>1</b>	<b>0.04</b>
<i>GCNT1</i>			<b>0.02</b>	<b>1</b>	<b>0.03</b>
<i>HIST1H2BG</i>			<b>0.04</b>	<b>1</b>	<b>0.03</b>

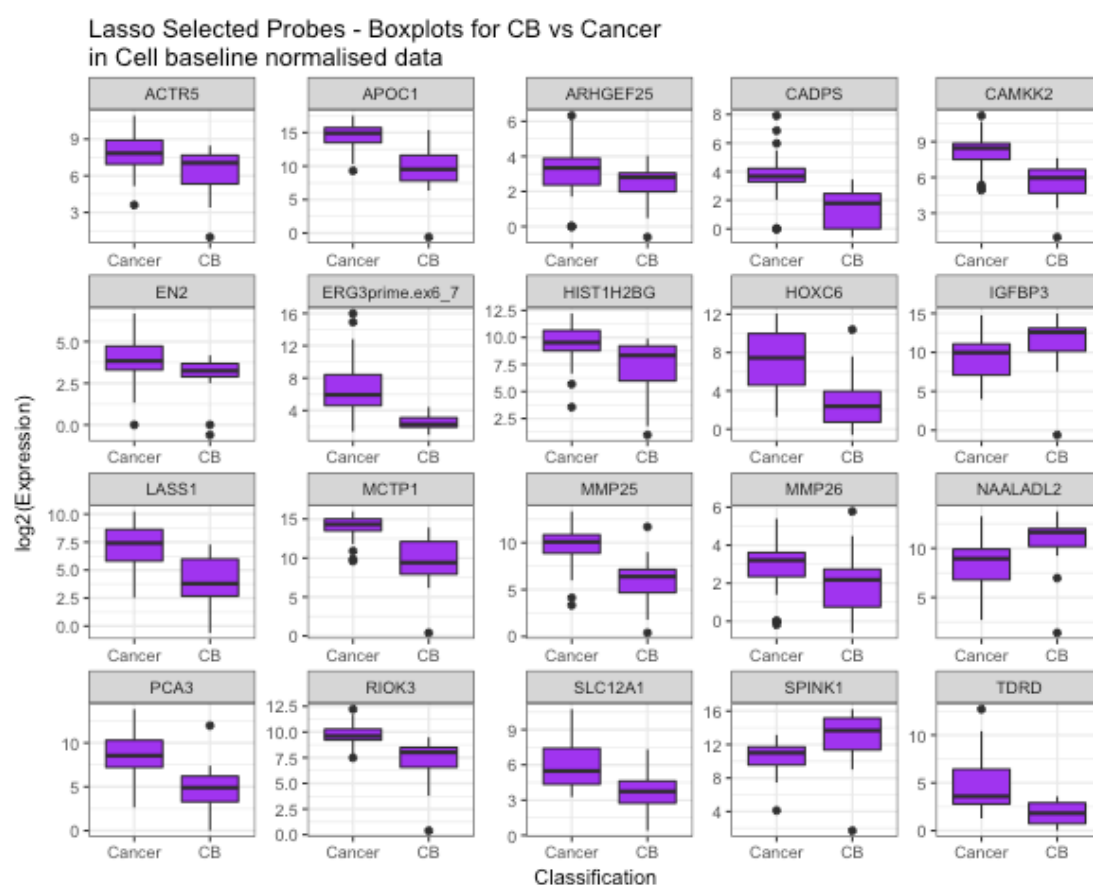
## 9: APPENDICES

Supplementary Table 33 Transcripts that have significant differential expression between CB and cancer samples (L, I, H) in the HK normalised NanoString data.

<i>Transcript</i>	<i>MWU</i>		<i>glm</i>		<i>Log<sub>2</sub>(FC)</i>
	<i>p-value</i>	<i>Adjusted p-value</i>	<i>p-value</i>	<i>Adjusted p-value</i>	
<i>HOXC6</i>	0.0002	0.0374	0.0019	0.3087	1.5
<i>ERG3' exons 6-7</i>	0.0006	0.1045	0.0228	0.9861	1.1
<i>TMPRSS2:ERG</i>	0.0036	0.5527	0.0069	0.9861	1.1
<i>CP</i>	0.0146	0.9924	0.0109	0.9861	-1
<i>TDRD</i>	0.001	0.153	0.0105	0.9861	0.9
<i>NAALADL2</i>	3.33x10 <sup>-05</sup>	0.0056	0.0012	0.2012	-0.8
<i>SLC43A1</i>	0.0005	0.0895	0.0168	0.9861	0.8
<i>ST6GALNAC1</i>	0.0008	0.1311	0.0238	0.9861	-0.8
<i>SPINK1</i>	7.80x10 <sup>-05</sup>	0.0129			-0.7
<i>UPK2</i>	0.0007	0.1128	0.0026	0.4313	-0.7
<i>CADPS</i>	0.0083	0.9924	0.0076	0.9861	0.7
<i>HPN</i>	0.0022	0.3485	0.0072	0.9861	0.7
<i>MFS2A</i>	0.0123	0.9924	0.0082	0.9861	0.7
<i>DNAH5</i>	0.0116	0.9924			-0.7
<i>IGFBP3</i>	0.0086	0.9924			-0.7
<i>SERPIN5</i>	0.0003	0.0489	0.0205	0.9861	-0.6
<i>B4GALNT4</i>	0.0273	0.9924			0.6
<i>CLIC2</i>	0.0055	0.8138	0.0097	0.9861	0.6
<i>LASS1</i>	0.0055	0.8138	0.0182	0.9861	0.6
<i>PCA3</i>	0.0006	0.1045	0.005	0.8153	0.6
<i>ITGBL1</i>	0.0227	0.9924			-0.6
<i>CCDC88B</i>	0.024	0.9924	0.0349	0.9861	0.5
<i>ERG5'</i>	0.0164	0.9924			0.5
<i>ISX</i>	0.0169	0.9924	0.0124	0.9861	0.5
<i>MMP25</i>	0.0071	0.9924			0.5
<i>AGR2</i>	0.0227	0.9924	0.0348	0.9861	-0.5
<i>GJB1</i>	0.0024	0.3722	0.0136	0.9861	-0.5
<i>MCTP1</i>	0.0164	0.9924	0.0349	0.9861	0.4
<i>PPAP2A</i>	0.0008	0.1311			-0.4
<i>PPP1R12B</i>	0.0138	0.9924			-0.4
<i>TMEM86A</i>	0.0037	0.5561	0.0227	0.9861	-0.4
<i>APOC1</i>	0.003	0.4577	0.0172	0.9861	0.3
<i>CAMKK2</i>	0.024	0.9924			0.3
<i>GCNT1</i>	0.0482	0.9924	0.0425	0.9861	0.3
<i>SIM2 long</i>	0.0109	0.9924	0.0125	0.9861	0.3
<i>SLC12A1</i>	0.0398	0.9924			0.3
<i>SULF2</i>			0.0391	0.9861	0.3
<i>MDK</i>	0.0022	0.3485			-0.3
<i>MNX1</i>	0.0193	0.9924			-0.3
<i>OGT</i>	0.0253	0.9924			-0.3
<i>PALM3</i>	0.0343	0.9924			-0.3
<i>RAB17</i>	0.0052	0.7725			-0.3
<i>RPS10</i>	0.0155	0.9924			-0.3
<i>STEAP2</i>	0.0076	0.9924			-0.3
<i>MIR146A</i>			0.0347	0.9861	0.2

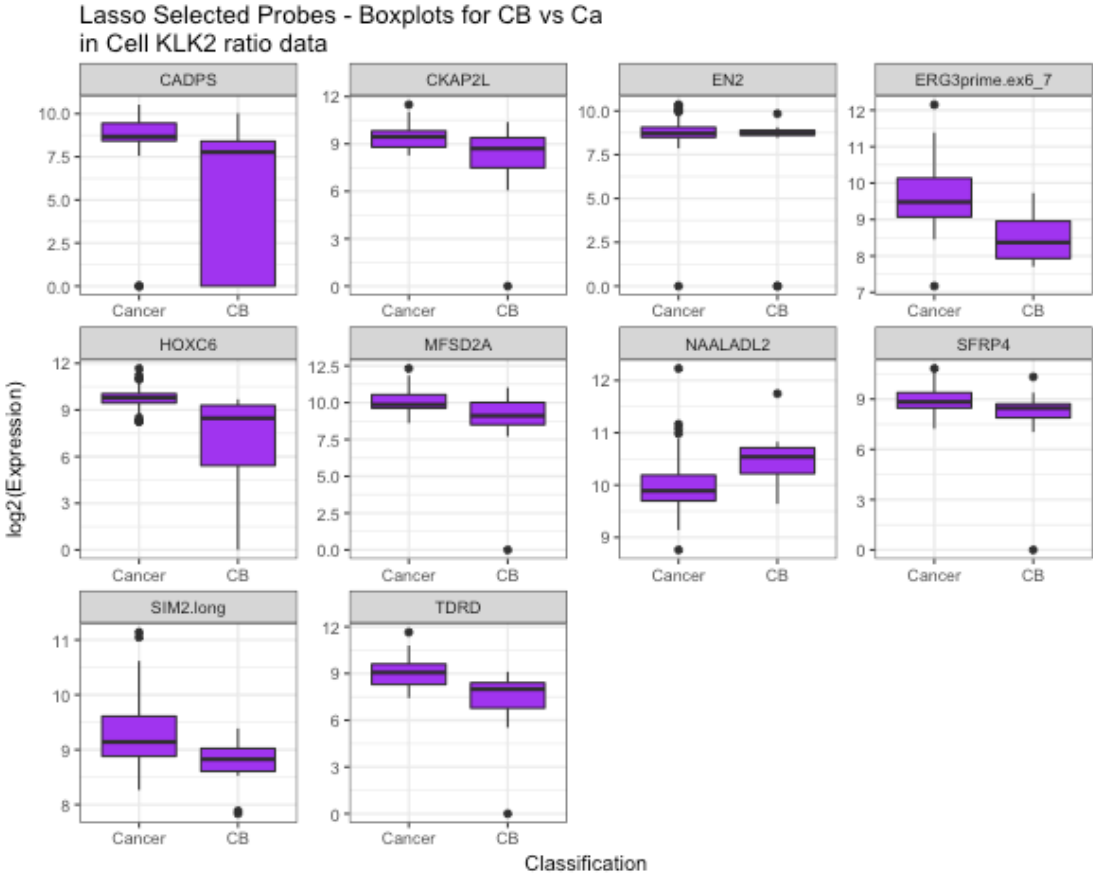
9: APPENDICES

<i>RIOK3</i>	<b>0.0379</b>	<b>0.9924</b>	<b>0.0486</b>	<b>0.9861</b>	<b>0.2</b>
<i>TMEM86A</i>	<b>0.0155</b>	<b>0.9924</b>			<b>0.2</b>
<i>HIST1H1C</i>	<b>0.0379</b>	<b>0.9924</b>			<b>-0.2</b>
<i>IFT57</i>	<b>0.0081</b>	<b>0.9924</b>			<b>-0.2</b>
<i>IMPDH2</i>	<b>0.0361</b>	<b>0.9924</b>			<b>-0.2</b>
<i>MSMB</i>	<b>0.0091</b>	<b>0.9924</b>			<b>-0.2</b>
<i>MYOF</i>	<b>0.0081</b>	<b>0.9924</b>			<b>-0.2</b>
<i>PTN</i>	<b>0.0266</b>	<b>0.9924</b>			<b>-0.2</b>
<i>RPL18A</i>	<b>0.0253</b>	<b>0.9924</b>			<b>-0.2</b>
<i>RPLP2</i>	<b>0.0138</b>	<b>0.9924</b>			<b>-0.2</b>
<i>RPS11</i>	<b>0.0253</b>	<b>0.9924</b>			<b>-0.2</b>
<i>ZNF577</i>	<b>0.0193</b>	<b>0.9924</b>			<b>-0.2</b>
<i>HIST1H1E</i>	<b>0.0438</b>	<b>0.9924</b>			<b>-0.1</b>



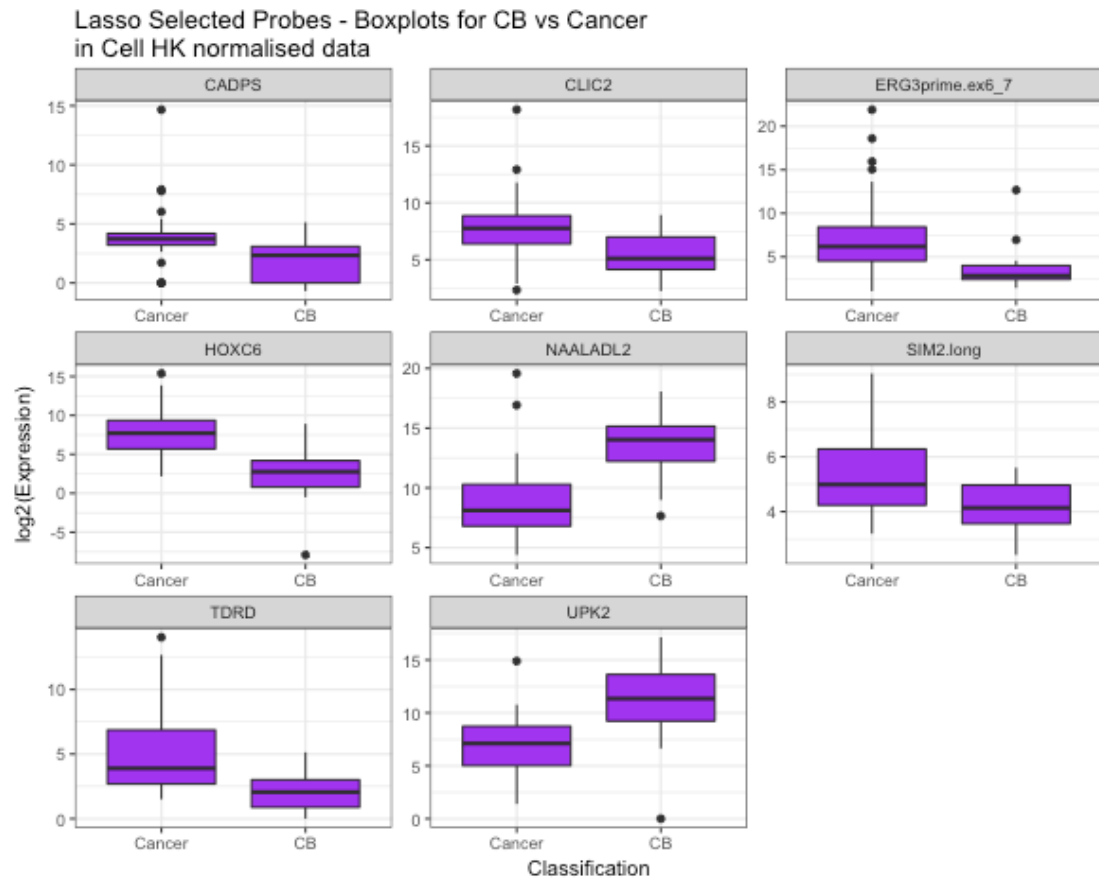
Supplementary Figure 13 Boxplots of all of the Lasso selected probes involved in CB vs. cancer (L, I, and H) models from the baseline normalised data.

9: APPENDICES



Supplementary Figure 14 Boxplots of all of the Lasso selected probes involved in CB vs. cancer (L, I, and H) models from the *KLK2* ratio data.

## 9: APPENDICES



Supplementary Figure 15 Boxplots of all of the Lasso selected probes involved in CB vs. cancer (L, I, and H) models from the HK normalised data.



9: APPENDICES

Supplementary Table 34 Random Forest results for Ca vs CBN baseline normalisation

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 85)</i>			<i>Transcripts identified by Mann Whitney U (n = 94)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>TMPRSS2:ERG</i>	<i>0.54</i>	<i>167</i>	<i>APOC1</i>	<i>0.57</i>	<i>85</i>	<i>TMPRSS2:ERG</i>	<i>0.78</i>	<i>94</i>
<i>APOC1</i>	<i>0.53</i>	<i>166</i>	<i>SPINK1</i>	<i>0.47</i>	<i>84</i>	<i>APOC1</i>	<i>0.59</i>	<i>93</i>
<i>ERG3' exons 6-7</i>	<i>0.53</i>	<i>165</i>	<i>TMPRSS2:ERG</i>	<i>0.47</i>	<i>83</i>	<i>NEAT1</i>	<i>0.49</i>	<i>92</i>
<i>NEAT1</i>	<i>0.45</i>	<i>164</i>	<i>RIOK3</i>	<i>0.46</i>	<i>82</i>	<i>RIOK3</i>	<i>0.47</i>	<i>91</i>
<i>RIOK3</i>	<i>0.41</i>	<i>163</i>	<i>ERG3' exons 6-7</i>	<i>0.44</i>	<i>81</i>	<i>ERG3' exons 6-7</i>	<i>0.39</i>	<i>90</i>
<i>SIM2 long</i>	<i>0.38</i>	<i>162</i>	<i>NEAT1</i>	<i>0.44</i>	<i>80</i>	<i>MCTP1</i>	<i>0.36</i>	<i>89</i>
<i>MCTP1</i>	<i>0.35</i>	<i>161</i>	<i>SIM2 long</i>	<i>0.39</i>	<i>79</i>	<i>MFSD2A</i>	<i>0.35</i>	<i>88</i>
<i>SPINK1</i>	<i>0.32</i>	<i>160</i>	<i>MCTP1</i>	<i>0.38</i>	<i>78</i>	<i>CADPS</i>	<i>0.34</i>	<i>87</i>
<i>CCDC88B</i>	<i>0.28</i>	<i>159</i>	<i>CADPS</i>	<i>0.36</i>	<i>77</i>	<i>SIM2 long</i>	<i>0.28</i>	<i>86</i>
<i>CADPS</i>	<i>0.27</i>	<i>158</i>	<i>MFSD2A</i>	<i>0.32</i>	<i>76</i>	<i>SPINK1</i>	<i>0.25</i>	<i>85</i>
<i>MXII</i>	<i>0.22</i>	<i>157</i>	<i>CCDC88B</i>	<i>0.32</i>	<i>75</i>	<i>CCDC88B</i>	<i>0.24</i>	<i>84</i>
<i>MFSD2A</i>	<i>0.21</i>	<i>156</i>	<i>CKAP2L</i>	<i>0.24</i>	<i>74</i>	<i>CKAP2L</i>	<i>0.23</i>	<i>83</i>
<i>MMP25</i>	<i>0.21</i>	<i>155</i>	<i>CAMKK2</i>	<i>0.21</i>	<i>73</i>	<i>GAPDH</i>	<i>0.22</i>	<i>82</i>
<i>CKAP2L</i>	<i>0.19</i>	<i>154</i>	<i>SLC43A1</i>	<i>0.20</i>	<i>72</i>	<i>MXII</i>	<i>0.21</i>	<i>81</i>
<i>HOXC6</i>	<i>0.16</i>	<i>153</i>	<i>MIR4435 IHG</i>	<i>0.19</i>	<i>71</i>	<i>CAMKK2</i>	<i>0.20</i>	<i>80</i>
<i>SULF2</i>	<i>0.16</i>	<i>152</i>	<i>MMP25</i>	<i>0.19</i>	<i>70</i>	<i>SLC43A1</i>	<i>0.17</i>	<i>79</i>
<i>MIR4435 IHG</i>	<i>0.15</i>	<i>151</i>	<i>SULF2</i>	<i>0.17</i>	<i>69</i>	<i>MMP25</i>	<i>0.16</i>	<i>78</i>
<i>SIRT1</i>	<i>0.15</i>	<i>150</i>	<i>HOXC6</i>	<i>0.17</i>	<i>68</i>	<i>AURKA</i>	<i>0.14</i>	<i>77</i>
<i>CAMKK2</i>	<i>0.14</i>	<i>149</i>	<i>GAPDH</i>	<i>0.16</i>	<i>67</i>	<i>HOXC6</i>	<i>0.14</i>	<i>76</i>
<i>AURKA</i>	<i>0.13</i>	<i>148</i>	<i>UPK2</i>	<i>0.15</i>	<i>66</i>	<i>HPRT</i>	<i>0.13</i>	<i>75</i>
<i>GAPDH</i>	<i>0.13</i>	<i>147</i>	<i>ISX</i>	<i>0.14</i>	<i>65</i>	<i>SIRT1</i>	<i>0.13</i>	<i>74</i>
<i>B4GALNT4</i>	<i>0.11</i>	<i>146</i>	<i>AATF</i>	<i>0.13</i>	<i>64</i>	<i>MIR4435 IHG</i>	<i>0.12</i>	<i>73</i>
<i>TDRD</i>	<i>0.11</i>	<i>145</i>	<i>TDRD</i>	<i>0.13</i>	<i>63</i>	<i>HPN</i>	<i>0.12</i>	<i>72</i>
<i>SLC43A1</i>	<i>0.10</i>	<i>144</i>	<i>MXII</i>	<i>0.12</i>	<i>62</i>	<i>TDRD</i>	<i>0.11</i>	<i>71</i>
<i>UPK2</i>	<i>0.10</i>	<i>143</i>	<i>AURKA</i>	<i>0.11</i>	<i>61</i>	<i>IGFBP3</i>	<i>0.10</i>	<i>70</i>

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 85)</i>			<i>Transcripts identified by Mann Whitney U (n = 94)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>HPN</i>	<b>0.09</b>	<b>142</b>	<i>HPRT</i>	<b>0.11</b>	<b>60</b>	<i>SULF2</i>	<b>0.10</b>	<b>69</b>
<i>IGFBP3</i>	<b>0.09</b>	<b>141</b>	<i>SLC4A1 S</i>	<b>0.10</b>	<b>59</b>	<i>UPK2</i>	<b>0.09</b>	<b>68</b>
<i>SNCA</i>	<b>0.09</b>	<b>140</b>	<i>HPN</i>	<b>0.09</b>	<b>58</b>	<i>MMP26</i>	<b>0.09</b>	<b>67</b>
<i>TMEM45B</i>	<b>0.07</b>	<b>139</b>	<i>PCA3</i>	<b>0.09</b>	<b>57</b>	<i>AATF</i>	<b>0.09</b>	<b>66</b>
<i>AATF</i>	<b>0.07</b>	<b>138</b>	<i>MAPK8IP2</i>	<b>0.09</b>	<b>56</b>	<i>MAPK8IP2</i>	<b>0.09</b>	<b>65</b>
<i>LASS1</i>	<b>0.07</b>	<b>137</b>	<i>GCNT1</i>	<b>0.09</b>	<b>55</b>	<i>SNCA</i>	<b>0.09</b>	<b>64</b>
<i>GCNT1</i>	<b>0.06</b>	<b>136</b>	<i>STOM</i>	<b>0.09</b>	<b>54</b>	<i>LASS1</i>	<b>0.09</b>	<b>63</b>
<i>HPRT</i>	<b>0.06</b>	<b>135</b>	<i>B4GALNT4</i>	<b>0.08</b>	<b>53</b>	<i>CD10</i>	<b>0.08</b>	<b>62</b>
<i>NAALADL2</i>	<b>0.06</b>	<b>134</b>	<i>SLC12A1</i>	<b>0.08</b>	<b>52</b>	<i>GCNT1</i>	<b>0.08</b>	<b>61</b>
<i>ISX</i>	<b>0.06</b>	<b>133</b>	<i>PTPRC</i>	<b>0.07</b>	<b>51</b>	<i>TMCC2</i>	<b>0.07</b>	<b>60</b>
<i>AMH</i>	<b>0.05</b>	<b>132</b>	<i>DPP4</i>	<b>0.07</b>	<b>50</b>	<i>SFRP4</i>	<b>0.07</b>	<b>59</b>
<i>SLC12A1</i>	<b>0.05</b>	<b>131</b>	<i>CD10</i>	<b>0.06</b>	<b>49</b>	<i>ITPR1</i>	<b>0.06</b>	<b>58</b>
<i>SLC4A1 S</i>	<b>0.05</b>	<b>130</b>	<i>EN2</i>	<b>0.06</b>	<b>48</b>	<i>EN2</i>	<b>0.06</b>	<b>57</b>
<i>CACNAID</i>	<b>0.05</b>	<b>129</b>	<i>SNCA</i>	<b>0.06</b>	<b>47</b>	<i>B4GALNT4</i>	<b>0.06</b>	<b>56</b>
<i>RPL23AP53</i>	<b>0.05</b>	<b>128</b>	<i>PDLIM5</i>	<b>0.05</b>	<b>46</b>	<i>ERG5'</i>	<b>0.06</b>	<b>55</b>
<i>CDC37L1</i>	<b>0.05</b>	<b>127</b>	<i>TMCC2</i>	<b>0.05</b>	<b>45</b>	<i>SLC12A1</i>	<b>0.06</b>	<b>54</b>
<i>PCA3</i>	<b>0.05</b>	<b>126</b>	<i>SIRT1</i>	<b>0.05</b>	<b>44</b>	<i>ISX</i>	<b>0.06</b>	<b>53</b>
<i>ACTR5</i>	<b>0.05</b>	<b>125</b>	<i>MGAT5B</i>	<b>0.05</b>	<b>43</b>	<i>PCA3</i>	<b>0.05</b>	<b>52</b>
<i>PTPRC</i>	<b>0.04</b>	<b>124</b>	<i>SNORA20</i>	<b>0.05</b>	<b>42</b>	<i>PDLIM5</i>	<b>0.05</b>	<b>51</b>
<i>MMP26</i>	<b>0.04</b>	<b>123</b>	<i>TMEM86A</i>	<b>0.05</b>	<b>41</b>	<i>STOM</i>	<b>0.05</b>	<b>50</b>
<i>RNF157</i>	<b>0.04</b>	<b>122</b>	<i>LASS1</i>	<b>0.04</b>	<b>40</b>	<i>ACTR5</i>	<b>0.05</b>	<b>49</b>
<i>MAPK8IP2</i>	<b>0.04</b>	<b>121</b>	<i>HIST1H2BG</i>	<b>0.04</b>	<b>39</b>	<i>DPP4</i>	<b>0.05</b>	<b>48</b>
<i>STOM</i>	<b>0.04</b>	<b>120</b>	<i>SRSF3</i>	<b>0.04</b>	<b>38</b>	<i>ERG3' exons 4-5</i>	<b>0.05</b>	<b>47</b>
<i>CDC20</i>	<b>0.04</b>	<b>119</b>	<i>NAALADL2</i>	<b>0.04</b>	<b>37</b>	<i>TMEM45B</i>	<b>0.04</b>	<b>46</b>
<i>EN2</i>	<b>0.04</b>	<b>118</b>	<i>AMH</i>	<b>0.04</b>	<b>36</b>	<i>AMH</i>	<b>0.04</b>	<b>45</b>
<i>SRSF3</i>	<b>0.04</b>	<b>117</b>	<i>STEAP4</i>	<b>0.04</b>	<b>35</b>	<i>NAALADL2</i>	<b>0.04</b>	<b>44</b>
<i>ERG5'</i>	<b>0.04</b>	<b>116</b>	<i>CACNAID</i>	<b>0.04</b>	<b>34</b>	<i>HIST3H2A</i>	<b>0.04</b>	<b>43</b>

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 85)</i>			<i>Transcripts identified by Mann Whitney U (n = 94)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>SFRP4</i>	<b>0.04</b>	<b>115</b>	<i>ACTR5</i>	<b>0.04</b>	<b>33</b>	<i>TMEM86A</i>	<b>0.04</b>	<b>42</b>
<i>MYOF</i>	<b>0.04</b>	<b>114</b>	<i>ANKRD34B</i>	<b>0.03</b>	<b>32</b>	<i>FOLH1</i>	<b>0.04</b>	<b>41</b>
<i>CLIC2</i>	<b>0.04</b>	<b>113</b>	<i>SFRP4</i>	<b>0.03</b>	<b>31</b>	<i>HIST1H2BG</i>	<b>0.04</b>	<b>40</b>
<i>HIST1H2BG</i>	<b>0.04</b>	<b>112</b>	<i>CDC20</i>	<b>0.02</b>	<b>30</b>	<i>FDPS</i>	<b>0.03</b>	<b>39</b>
<i>MAK</i>	<b>0.04</b>	<b>111</b>	<i>SEC61A1</i>	<b>0.02</b>	<b>29</b>	<i>CLIC2</i>	<b>0.03</b>	<b>38</b>
<i>Timp4</i>	<b>0.03</b>	<b>110</b>	<i>CLIC2</i>	<b>0.02</b>	<b>28</b>	<i>MIR146A</i>	<b>0.03</b>	<b>37</b>
<i>TMEM86A</i>	<b>0.03</b>	<b>109</b>	<i>HIST1H2BF</i>	<b>0.02</b>	<b>27</b>	<i>SRSF3</i>	<b>0.03</b>	<b>36</b>
<i>PPP1R12B</i>	<b>0.03</b>	<b>108</b>	<i>FOLH1</i>	<b>0.02</b>	<b>26</b>	<i>PTPRC</i>	<b>0.03</b>	<b>35</b>
<i>STEAP4</i>	<b>0.03</b>	<b>107</b>	<i>ANPEP</i>	<b>0.02</b>	<b>25</b>	<i>MAK</i>	<b>0.03</b>	<b>34</b>
<i>DPP4</i>	<b>0.03</b>	<b>106</b>	<i>ERG5'</i>	<b>0.02</b>	<b>24</b>	<i>SEC61A1</i>	<b>0.03</b>	<b>33</b>
<i>CD10</i>	<b>0.03</b>	<b>105</b>	<i>MIR146A</i>	<b>0.02</b>	<b>23</b>	<i>TWIST1</i>	<b>0.03</b>	<b>32</b>
<i>SULT1A1</i>	<b>0.03</b>	<b>104</b>	<i>TERF2IP</i>	<b>0.02</b>	<b>22</b>	<i>SERPINB5</i>	<b>0.02</b>	<b>31</b>
<i>PDLIM5</i>	<b>0.03</b>	<b>103</b>	<i>MED4</i>	<b>0.02</b>	<b>21</b>	<i>NLRP3</i>	<b>0.02</b>	<b>30</b>
<i>P712P</i>	<b>0.03</b>	<b>102</b>	<i>ITPR1</i>	<b>0.01</b>	<b>20</b>	<i>CDC20</i>	<b>0.02</b>	<b>29</b>
<i>MSMB</i>	<b>0.03</b>	<b>101</b>	<i>BTG2</i>	<b>0.01</b>	<b>19</b>	<i>RPS11</i>	<b>0.02</b>	<b>28</b>
<i>ERG3' exons 4-5</i>	<b>0.03</b>	<b>100</b>	<i>NKAIN1</i>	<b>0.01</b>	<b>18</b>	<i>CACNA1D</i>	<b>0.02</b>	<b>27</b>
<i>AGR2</i>	<b>0.02</b>	<b>99</b>	<i>MEMO1</i>	<b>0.01</b>	<b>17</b>	<i>SACMIL</i>	<b>0.02</b>	<b>26</b>
<i>PECI</i>	<b>0.02</b>	<b>98</b>	<i>CASKIN1</i>	<b>0.01</b>	<b>16</b>	<i>RPL18A</i>	<b>0.02</b>	<b>25</b>
<i>MNX1</i>	<b>0.02</b>	<b>97</b>	<i>SMAP1 exons 7-8</i>	<b>0.01</b>	<b>15</b>	<i>ANKRD34B</i>	<b>0.01</b>	<b>24</b>
<i>PPAP2A</i>	<b>0.02</b>	<b>96</b>	<i>TBP</i>	<b>0.01</b>	<b>14</b>	<i>TERF2IP</i>	<b>0.01</b>	<b>23</b>
<i>PPFIA2</i>	<b>0.02</b>	<b>95</b>	<i>SIM2 short</i>	<b>0.01</b>	<b>13</b>	<i>GABARAPL2</i>	<b>0.01</b>	<b>22</b>
<i>PALM3</i>	<b>0.02</b>	<b>94</b>	<i>MEX3A</i>	<b>0.01</b>	<b>12</b>	<i>SNORA20</i>	<b>0.01</b>	<b>21</b>
<i>ITPR1</i>	<b>0.02</b>	<b>93</b>	<i>CDKN3</i>	<b>0.01</b>	<b>11</b>	<i>MEX3A</i>	<b>0.01</b>	<b>20</b>
<i>RPS11</i>	<b>0.02</b>	<b>92</b>	<i>SACMIL</i>	<b>0.01</b>	<b>10</b>	<i>HOXC4</i>	<b>0.01</b>	<b>19</b>
<i>VAX2</i>	<b>0.02</b>	<b>91</b>	<i>MMP11</i>	<b>0.01</b>	<b>9</b>	<i>ALAS1</i>	<b>0.01</b>	<b>18</b>
<i>EIF2D</i>	<b>0.02</b>	<b>90</b>	<i>OR52A2</i>	<b>0.00</b>	<b>8</b>	<i>CAMK2N2</i>	<b>0.01</b>	<b>17</b>

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 85)</i>			<i>Transcripts identified by Mann Whitney U (n = 94)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>FOLH1</i>	<b>0.02</b>	<b>89</b>	<i>GABARAPL2</i>	<b>0.00</b>	<b>7</b>	<i>MED4</i>	<b>0.01</b>	<b>16</b>
<i>PVT1</i>	<b>0.02</b>	<b>88</b>	<i>EIF2D</i>	<b>0.00</b>	<b>6</b>	<i>NKAIN1</i>	<b>0.01</b>	<b>15</b>
<i>AR.ex9</i>	<b>0.02</b>	<b>87</b>	<i>PSTPIP1</i>	<b>0.00</b>	<b>5</b>	<i>MMP11</i>	<b>0.01</b>	<b>14</b>
<i>ANKRD34B</i>	<b>0.02</b>	<b>86</b>	<i>SSTR1</i>	<b>0.00</b>	<b>4</b>	<i>ANPEP</i>	<b>0.01</b>	<b>13</b>
<i>MKi67</i>	<b>0.02</b>	<b>85</b>	<i>NLRP3</i>	<b>0.00</b>	<b>3</b>	<i>ARHGEF25</i>	<b>0.01</b>	<b>12</b>
<i>MGAT5B</i>	<b>0.02</b>	<b>84</b>	<i>HMBS</i>	<b>0.00</b>	<b>2</b>	<i>CASKIN1</i>	<b>0.01</b>	<b>11</b>
<i>SNORA20</i>	<b>0.02</b>	<b>83</b>	<i>B2M</i>	<b>0.00</b>	<b>1</b>	<i>B2M</i>	<b>0.01</b>	<b>10</b>
<i>IMPDH2</i>	<b>0.01</b>	<b>82</b>				<i>OR52A2</i>	<b>0.00</b>	<b>9</b>
<i>MED4</i>	<b>0.01</b>	<b>81</b>				<i>BTG2</i>	<b>0.00</b>	<b>8</b>
<i>GJB1</i>	<b>0.01</b>	<b>80</b>				<i>SSTR1</i>	<b>0.00</b>	<b>7</b>
<i>HIST3H2A</i>	<b>0.01</b>	<b>79</b>				<i>SIM2 short</i>	<b>0.00</b>	<b>6</b>
<i>CAMK2N2</i>	<b>0.01</b>	<b>78</b>				<i>EIF2D</i>	<b>0.00</b>	<b>5</b>
<i>OGT</i>	<b>0.01</b>	<b>77</b>				<i>MEMO1</i>	<b>0.00</b>	<b>4</b>
<i>HIST1H2BF</i>	<b>0.01</b>	<b>76</b>				<i>CDKN3</i>	<b>0.00</b>	<b>3</b>
<i>DLX1</i>	<b>0.01</b>	<b>75</b>				<i>TBP</i>	<b>0.00</b>	<b>2</b>
<i>MCM7</i>	<b>0.01</b>	<b>74</b>				<i>PSTPIP1</i>	<b>0.00</b>	<b>1</b>
<i>SEC61A1</i>	<b>0.01</b>	<b>73</b>						
<i>PSTPIP1</i>	<b>0.01</b>	<b>72</b>						
<i>ARHGEF25</i>	<b>0.01</b>	<b>71</b>						
<i>IFT57</i>	<b>0.01</b>	<b>70</b>						
<i>GOLM1</i>	<b>0.01</b>	<b>69</b>						
<i>TMCC2</i>	<b>0.01</b>	<b>68</b>						
<i>SERPINB5</i>	<b>0.01</b>	<b>67</b>						
<i>TERF2IP</i>	<b>0.01</b>	<b>66</b>						
<i>SPON2</i>	<b>0.01</b>	<b>65</b>						
<i>SSPO</i>	<b>0.01</b>	<b>64</b>						
<i>TMEM47</i>	<b>0.01</b>	<b>63</b>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 85)</i>			<i>Transcripts identified by Mann Whitney U (n = 94)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>GABARAPL2</i>	<i>0.01</i>	<i>62</i>						
<i>COL9A2</i>	<i>0.01</i>	<i>61</i>						
<i>RPS10</i>	<i>0.01</i>	<i>60</i>						
<i>SIM2 short</i>	<i>0.01</i>	<i>59</i>						
<i>MIR146A</i>	<i>0.01</i>	<i>58</i>						
<i>MEX3A</i>	<i>0.01</i>	<i>57</i>						
<i>ALAS1</i>	<i>0.01</i>	<i>56</i>						
<i>AMACR</i>	<i>0.01</i>	<i>55</i>						
<i>ITGBL1</i>	<i>0.01</i>	<i>54</i>						
<i>FDPS</i>	<i>0.01</i>	<i>53</i>						
<i>TWIST1</i>	<i>0.01</i>	<i>52</i>						
<i>HMBS</i>	<i>0.01</i>	<i>51</i>						
<i>KLK3 exons 1-2</i>	<i>0.01</i>	<i>50</i>						
<i>KLK4</i>	<i>0.01</i>	<i>49</i>						
<i>TFDP1</i>	<i>0.01</i>	<i>48</i>						
<i>VPS13A</i>	<i>0.01</i>	<i>47</i>						
<i>MEMO1</i>	<i>0.01</i>	<i>46</i>						
<i>ANPEP</i>	<i>0.01</i>	<i>45</i>						
<i>RAB17</i>	<i>0.01</i>	<i>44</i>						
<i>TRPM4</i>	<i>0.01</i>	<i>43</i>						
<i>HIST1H1C</i>	<i>0.01</i>	<i>42</i>						
<i>TBP</i>	<i>0.01</i>	<i>41</i>						
<i>RPL18A</i>	<i>0.01</i>	<i>40</i>						
<i>KLK2</i>	<i>0.01</i>	<i>39</i>						
<i>NKAIN1</i>	<i>0.01</i>	<i>38</i>						
<i>ZNF577</i>	<i>0.01</i>	<i>37</i>						
<i>BTG2</i>	<i>0.01</i>	<i>36</i>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 85)</i>			<i>Transcripts identified by Mann Whitney U (n = 94)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>SChLAP1</i>	<i>0.01</i>	<i>35</i>						
<i>PCSK6</i>	<i>0.00</i>	<i>34</i>						
<i>CLU</i>	<i>0.00</i>	<i>33</i>						
<i>RPLP2</i>	<i>0.00</i>	<i>32</i>						
<i>ST6GALNAC1</i>	<i>0.00</i>	<i>31</i>						
<i>OR52A2</i>	<i>0.00</i>	<i>30</i>						
<i>SMIMI</i>	<i>0.00</i>	<i>29</i>						
<i>CDKN3</i>	<i>0.00</i>	<i>28</i>						
<i>MIC1</i>	<i>0.00</i>	<i>27</i>						
<i>ABCB9</i>	<i>0.00</i>	<i>26</i>						
<i>AR.ex4_8</i>	<i>0.00</i>	<i>25</i>						
<i>HIST1H1E</i>	<i>0.00</i>	<i>24</i>						
<i>DNAH5</i>	<i>0.00</i>	<i>23</i>						
<i>SMAP1 exons 7-8</i>	<i>0.00</i>	<i>22</i>						
<i>SYNM</i>	<i>0.00</i>	<i>21</i>						
<i>TERT</i>	<i>0.00</i>	<i>20</i>						
<i>PTN</i>	<i>0.00</i>	<i>19</i>						
<i>NLRP3</i>	<i>0.00</i>	<i>18</i>						
<i>CASKIN1</i>	<i>0.00</i>	<i>17</i>						
<i>BRAF</i>	<i>0.00</i>	<i>16</i>						
<i>Met</i>	<i>0.00</i>	<i>15</i>						
<i>MIATNB</i>	<i>0.00</i>	<i>14</i>						
<i>COL10A1</i>	<i>0.00</i>	<i>13</i>						
<i>HOXC4</i>	<i>0.00</i>	<i>12</i>						
<i>MDK</i>	<i>0.00</i>	<i>11</i>						
<i>SSTR1</i>	<i>0.00</i>	<i>10</i>						
<i>LBH</i>	<i>0.00</i>	<i>9</i>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 85)</i>			<i>Transcripts identified by Mann Whitney U (n = 94)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>RP11_97012.7</i>	<b>0.00</b>	<b>8</b>						
<i>STEAP2</i>	<b>0.00</b>	<b>7</b>						
<i>KLK3 exons 2-3</i>	<b>0.00</b>	<b>5.5</b>						
<i>SACMIL</i>	<b>0.00</b>	<b>5.5</b>						
<i>MARCH5</i>	<b>0.00</b>	<b>4</b>						
<i>CP</i>	<b>0.00</b>	<b>3</b>						
<i>B2M</i>	<b>0.00</b>	<b>2</b>						
<i>MMP11</i>	<b>0.00</b>	<b>1</b>						

Supplementary Table 35 Random Forest results for comparing cancer samples with clinically benign samples in *KLK2* factorised cell data.

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 24)</i>			<i>Transcripts identified by Mann Whitney U (n = 33)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ERG3' exons 6-7</i>	<b>0.85</b>	<b>166</b>	<i>SLC12A1</i>	<b>0.98</b>	<b>24</b>	<i>HOXC6</i>	<b>1.20</b>	<b>33</b>
<i>SLC12A1</i>	<b>0.80</b>	<b>165</b>	<i>ERG3' exons 6-7</i>	<b>0.92</b>	<b>23</b>	<i>SLC12A1</i>	<b>0.76</b>	<b>32</b>
<i>HOXC6</i>	<b>0.69</b>	<b>164</b>	<i>HOXC6</i>	<b>0.92</b>	<b>22</b>	<i>ERG3' exons 6-7</i>	<b>0.74</b>	<b>31</b>
<i>APOC1</i>	<b>0.41</b>	<b>163</b>	<i>PCA3</i>	<b>0.63</b>	<b>21</b>	<i>PCA3</i>	<b>0.51</b>	<b>30</b>
<i>CKAP2L</i>	<b>0.38</b>	<b>162</b>	<i>HIST1H2BG</i>	<b>0.59</b>	<b>20</b>	<i>APOC1</i>	<b>0.50</b>	<b>29</b>
<i>HIST1H2BG</i>	<b>0.36</b>	<b>161</b>	<i>CADPS</i>	<b>0.46</b>	<b>19</b>	<i>CKAP2L</i>	<b>0.42</b>	<b>28</b>
			<i>TMPRSS2:ERG fusion</i>			<i>TMPRSS2:ERG fusion</i>		
<i>CADPS</i>	<b>0.27</b>	<b>160</b>	<i>fusion</i>	<b>0.44</b>	<b>18</b>	<i>fusion</i>	<b>0.39</b>	<b>27</b>
<i>LASS1</i>	<b>0.25</b>	<b>159</b>	<i>CKAP2L</i>	<b>0.43</b>	<b>17</b>	<i>CADPS</i>	<b>0.34</b>	<b>26</b>
<i>SLC43A1</i>	<b>0.24</b>	<b>158</b>	<i>NAALADL2</i>	<b>0.37</b>	<b>16</b>	<i>HPN</i>	<b>0.34</b>	<b>25</b>
<i>NAALADL2</i>	<b>0.23</b>	<b>157</b>	<i>SIM2 long</i>	<b>0.36</b>	<b>15</b>	<i>NAALADL2</i>	<b>0.31</b>	<b>24</b>
<i>PCA3</i>	<b>0.23</b>	<b>156</b>	<i>TDRD</i>	<b>0.35</b>	<b>14</b>	<i>TMEM86A</i>	<b>0.30</b>	<b>23</b>

## 9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 24)</i>			<i>Transcripts identified by Mann Whitney U (n = 33)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>HPN</i>	<b>0.23</b>	<b>155</b>	<i>HPN</i>	<b>0.35</b>	<b>13</b>	<i>UPK2</i>	<b>0.28</b>	<b>22</b>
<i>SIM2 long</i>	<b>0.21</b>	<b>154</b>	<i>GCNT1</i>	<b>0.34</b>	<b>12</b>	<i>TDRD</i>	<b>0.26</b>	<b>21</b>
<i>TMPRSS2:ERG fusion</i>	<b>0.19</b>	<b>153</b>	<i>TMEM86A</i>	<b>0.24</b>	<b>11</b>	<i>SIM2 long</i>	<b>0.25</b>	<b>20</b>
<i>TMEM86A</i>	<b>0.19</b>	<b>152</b>	<i>LASS1</i>	<b>0.24</b>	<b>10</b>	<i>SLC43A1</i>	<b>0.24</b>	<b>19</b>
<i>ANKRD34B</i>	<b>0.17</b>	<b>151</b>	<i>TMCC2</i>	<b>0.23</b>	<b>9</b>	<i>ST6GALNAC1</i>	<b>0.20</b>	<b>18</b>
<i>AMACR</i>	<b>0.17</b>	<b>150</b>	<i>CLIC2</i>	<b>0.21</b>	<b>8</b>	<i>LASS1</i>	<b>0.20</b>	<b>17</b>
<i>TDRD</i>	<b>0.17</b>	<b>149</b>	<i>MMP25</i>	<b>0.20</b>	<b>7</b>	<i>TMCC2</i>	<b>0.18</b>	<b>16</b>
<i>GCNT1</i>	<b>0.14</b>	<b>148</b>	<i>MFSD2A</i>	<b>0.16</b>	<b>6</b>	<i>ERG5'</i>	<b>0.18</b>	<b>15</b>
<i>MFSD2A</i>	<b>0.12</b>	<b>147</b>	<i>MCTP1</i>	<b>0.14</b>	<b>5</b>	<i>SERPINB5</i>	<b>0.18</b>	<b>14</b>
<i>MCTP1</i>	<b>0.10</b>	<b>146</b>	<i>OR52A2</i>	<b>0.14</b>	<b>4</b>	<i>CLIC2</i>	<b>0.17</b>	<b>13</b>
<i>CAMKK2</i>	<b>0.10</b>	<b>145</b>	<i>CAMKK2</i>	<b>0.11</b>	<b>3</b>	<i>SFRP4</i>	<b>0.17</b>	<b>12</b>
<i>CLIC2</i>	<b>0.09</b>	<b>144</b>	<i>CCDC88B</i>	<b>0.09</b>	<b>2</b>	<i>B4GALNT4</i>	<b>0.17</b>	<b>11</b>
<i>TMCC2</i>	<b>0.09</b>	<b>143</b>	<i>NLRP3</i>	<b>0.05</b>	<b>1</b>	<i>ANKRD34B</i>	<b>0.13</b>	<b>10</b>
<i>B4GALNT4</i>	<b>0.09</b>	<b>142</b>				<i>CAMKK2</i>	<b>0.10</b>	<b>9</b>
<i>Timp4</i>	<b>0.09</b>	<b>141</b>				<i>MCTP1</i>	<b>0.09</b>	<b>8</b>
<i>UPK2</i>	<b>0.09</b>	<b>140</b>				<i>MMP25</i>	<b>0.09</b>	<b>7</b>
<i>ERG5'</i>	<b>0.08</b>	<b>139</b>				<i>ISX</i>	<b>0.08</b>	<b>6</b>
<i>DLX1</i>	<b>0.08</b>	<b>138</b>				<i>FOLH1</i>	<b>0.08</b>	<b>5</b>
<i>MMP25</i>	<b>0.08</b>	<b>137</b>				<i>MFSD2A</i>	<b>0.07</b>	<b>4</b>
<i>RNF157</i>	<b>0.08</b>	<b>136</b>				<i>CCDC88B</i>	<b>0.05</b>	<b>3</b>
<i>AURKA</i>	<b>0.08</b>	<b>135</b>				<i>SULF2</i>	<b>0.03</b>	<b>2</b>
<i>TERT</i>	<b>0.08</b>	<b>134</b>				<i>MIR146A</i>	<b>0.03</b>	<b>1</b>
<i>SFRP4</i>	<b>0.07</b>	<b>133</b>						
<i>CP</i>	<b>0.06</b>	<b>132</b>						
<i>NKAIN1</i>	<b>0.06</b>	<b>131</b>						
<i>CCDC88B</i>	<b>0.05</b>	<b>130</b>						
<i>OR52A2</i>	<b>0.05</b>	<b>129</b>						



9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 24)</i>			<i>Transcripts identified by Mann Whitney U (n = 33)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>AR exons 4-8</i>	<b>0.05</b>	<b>128</b>						
<i>STOM</i>	<b>0.04</b>	<b>127</b>						
<i>ABCB9</i>	<b>0.04</b>	<b>126</b>						
<i>ERG3' exons 4-5</i>	<b>0.04</b>	<b>125</b>						
<i>SERPINB5</i>	<b>0.04</b>	<b>124</b>						
<i>SULF2</i>	<b>0.04</b>	<b>123</b>						
<i>MAPK8IP2</i>	<b>0.03</b>	<b>122</b>						
<i>AGR2</i>	<b>0.03</b>	<b>121</b>						
<i>ISX</i>	<b>0.03</b>	<b>120</b>						
<i>STEAP2</i>	<b>0.03</b>	<b>119</b>						
<i>CDKN3</i>	<b>0.03</b>	<b>118</b>						
<i>FOLH1</i>	<b>0.03</b>	<b>117</b>						
<i>MMP11</i>	<b>0.03</b>	<b>116</b>						
<i>TMEM45B</i>	<b>0.03</b>	<b>115</b>						
<i>SPINK1</i>	<b>0.03</b>	<b>114</b>						
<i>ITGBL1</i>	<b>0.03</b>	<b>113</b>						
<i>PPAP2A</i>	<b>0.02</b>	<b>112</b>						
<i>MEX3A</i>	<b>0.02</b>	<b>111</b>						
<i>IGFBP3</i>	<b>0.02</b>	<b>110</b>						
<i>PVT1</i>	<b>0.02</b>	<b>109</b>						
<i>P712P</i>	<b>0.02</b>	<b>108</b>						
<i>PPFIA2</i>	<b>0.02</b>	<b>107</b>						
<i>TRPM4</i>	<b>0.02</b>	<b>106</b>						
<i>MSMB</i>	<b>0.02</b>	<b>105</b>						
<i>SLC4A1.S</i>	<b>0.02</b>	<b>104</b>						
<i>PPP1R12B</i>	<b>0.02</b>	<b>103</b>						
<i>AMH</i>	<b>0.02</b>	<b>102</b>						

## 9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 24)</i>			<i>Transcripts identified by Mann Whitney U (n = 33)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ST6GALNAC1</i>	<b>0.02</b>	<b>101</b>						
<i>DPP4</i>	<b>0.02</b>	<b>100</b>						
<i>SNORA20</i>	<b>0.02</b>	<b>99</b>						
<i>TMEM47</i>	<b>0.02</b>	<b>98</b>						
<i>VAX2</i>	<b>0.02</b>	<b>97</b>						
<i>HMBS</i>	<b>0.02</b>	<b>96</b>						
<i>VPS13A</i>	<b>0.01</b>	<b>95</b>						
<i>RPL23AP53</i>	<b>0.01</b>	<b>94</b>						
<i>EN2</i>	<b>0.01</b>	<b>93</b>						
<i>MKi67</i>	<b>0.01</b>	<b>92</b>						
<i>KLK4</i>	<b>0.01</b>	<b>91</b>						
<i>PALM3</i>	<b>0.01</b>	<b>90</b>						
<i>ALAS1</i>	<b>0.01</b>	<b>89</b>						
<i>RPL18A</i>	<b>0.01</b>	<b>88</b>						
<i>SEC61A1</i>	<b>0.01</b>	<b>87</b>						
<i>PTN</i>	<b>0.01</b>	<b>86</b>						
<i>MNX1</i>	<b>0.01</b>	<b>85</b>						
<i>TWIST1</i>	<b>0.01</b>	<b>84</b>						
<i>MGAT5B</i>	<b>0.01</b>	<b>83</b>						
<i>RPS11</i>	<b>0.01</b>	<b>82</b>						
<i>ZNF577</i>	<b>0.01</b>	<b>81</b>						
<i>PSTPIP1</i>	<b>0.01</b>	<b>80</b>						
<i>RIOK3</i>	<b>0.01</b>	<b>79</b>						
<i>KLK3 exons 2-3</i>	<b>0.01</b>	<b>78</b>						
<i>COL10A1</i>	<b>0.01</b>	<b>77</b>						
<i>OGT</i>	<b>0.01</b>	<b>76</b>						
<i>CASKIN1</i>	<b>0.01</b>	<b>75</b>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 24)</i>			<i>Transcripts identified by Mann Whitney U (n = 33)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>RPS10</i>	<b>0.01</b>	<b>74</b>						
<i>NLRP3</i>	<b>0.01</b>	<b>73</b>						
<i>CLU</i>	<b>0.01</b>	<b>72</b>						
<i>HIST1H1C</i>	<b>0.01</b>	<b>71</b>						
<i>SMIM1</i>	<b>0.01</b>	<b>70</b>						
<i>GJB1</i>	<b>0.01</b>	<b>69</b>						
<i>MIATNB</i>	<b>0.01</b>	<b>68</b>						
<i>CD10</i>	<b>0.01</b>	<b>67</b>						
<i>PDLIM5</i>	<b>0.01</b>	<b>66</b>						
<i>TBP</i>	<b>0.009</b>	<b>65</b>						
<i>MMP26</i>	<b>0.009</b>	<b>64</b>						
<i>CACNA1D</i>	<b>0.009</b>	<b>63</b>						
<i>SPON2</i>	<b>0.009</b>	<b>62</b>						
<i>MCM7</i>	<b>0.009</b>	<b>61</b>						
<i>MEMO1</i>	<b>0.009</b>	<b>60</b>						
<i>ACTR5</i>	<b>0.008</b>	<b>59</b>						
<i>RP11_97O12.7</i>	<b>0.008</b>	<b>58</b>						
<i>ITPR1</i>	<b>0.008</b>	<b>57</b>						
<i>TERF2IP</i>	<b>0.008</b>	<b>56</b>						
<i>STEAP4</i>	<b>0.008</b>	<b>55</b>						
<i>MAK</i>	<b>0.008</b>	<b>54</b>						
<i>SULT1A1</i>	<b>0.007</b>	<b>53</b>						
<i>NEAT1</i>	<b>0.007</b>	<b>52</b>						
<i>MYOF</i>	<b>0.006</b>	<b>51</b>						
<i>MIC1</i>	<b>0.006</b>	<b>50</b>						
<i>KLK3 exons 1-2</i>	<b>0.006</b>	<b>49</b>						
<i>HOXC4</i>	<b>0.005</b>	<b>48</b>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 24)</i>			<i>Transcripts identified by Mann Whitney U (n = 33)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>SRSF3</i>	<i>0.005</i>	<i>47</i>						
<i>GAPDH</i>	<i>0.005</i>	<i>46</i>						
<i>MDK</i>	<i>0.005</i>	<i>45</i>						
<i>SACMIL</i>	<i>0.005</i>	<i>44</i>						
<i>HIST1H1E</i>	<i>0.005</i>	<i>43</i>						
<i>GABARAPL2</i>	<i>0.005</i>	<i>42</i>						
<i>MIR4435 IHG</i>	<i>0.005</i>	<i>41</i>						
<i>FDPS</i>	<i>0.005</i>	<i>40</i>						
<i>COL9A2</i>	<i>0.004</i>	<i>39</i>						
<i>DNAH5</i>	<i>0.004</i>	<i>38</i>						
<i>LBH</i>	<i>0.004</i>	<i>37</i>						
<i>RAB17</i>	<i>0.003</i>	<i>36</i>						
<i>SChLAPI</i>	<i>0.003</i>	<i>35</i>						
<i>BRAF</i>	<i>0.003</i>	<i>34</i>						
<i>TFDP1</i>	<i>0.003</i>	<i>33</i>						
<i>IFT57</i>	<i>0.003</i>	<i>32</i>						
<i>RPLP2</i>	<i>0.003</i>	<i>31</i>						
<i>HIST3H2A</i>	<i>0.003</i>	<i>30</i>						
<i>SIM2 short</i>	<i>0.002</i>	<i>29</i>						
<i>ANPEP</i>	<i>0.002</i>	<i>28</i>						
<i>AATF</i>	<i>0.002</i>	<i>27</i>						
<i>BTG2</i>	<i>0.002</i>	<i>26</i>						
<i>MXI1</i>	<i>0.002</i>	<i>25</i>						
<i>MED4</i>	<i>0.002</i>	<i>24</i>						
<i>IMPDH2</i>	<i>0.002</i>	<i>23</i>						
<i>SSTR1</i>	<i>0.002</i>	<i>22</i>						
<i>MIR146A</i>	<i>0.002</i>	<i>21</i>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 24)</i>			<i>Transcripts identified by Mann Whitney U (n = 33)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>Mar-05</i>	<b>0.002</b>	<b>20</b>						
<i>SIRT1</i>	<b>0.002</b>	<b>19</b>						
<i>AR exon 9</i>	<b>0.002</b>	<b>18</b>						
<i>PECI</i>	<b>0.002</b>	<b>17</b>						
<i>SYNM</i>	<b>1.15x10<sup>-17</sup></b>	<b>16</b>						
<i>PTPRC</i>	<b>1.07x10<sup>-17</sup></b>	<b>15</b>						
<i>GOLM1</i>	<b>8.99x10<sup>-18</sup></b>	<b>14</b>						
<i>ARHGEF25</i>	<b>7.55x10<sup>-18</sup></b>	<b>13</b>						
<i>CDC37L1</i>	<b>7.11x10<sup>-18</sup></b>	<b>12</b>						
<i>CDC20</i>	<b>5.77x10<sup>-18</sup></b>	<b>11</b>						
<i>SSPO</i>	<b>4.00x10<sup>-18</sup></b>	<b>10</b>						
<i>SMAP1 exons 7-8</i>	<b>3.55x10<sup>-18</sup></b>	<b>9</b>						
<i>EIF2D</i>	<b>2.66x10<sup>-18</sup></b>	<b>8</b>						
<i>SNCA</i>	<b>1.78x10<sup>-18</sup></b>	<b>7</b>						
<i>B2M</i>	<b>4.44x10<sup>-19</sup></b>	<b>6</b>						
<i>CAMK2N2</i>	<b>0</b>	<b>3</b>						
<i>HIST1H2BF</i>	<b>0</b>	<b>3</b>						
<i>HPRT</i>	<b>0</b>	<b>3</b>						
<i>Met</i>	<b>0</b>	<b>3</b>						
<i>PCSK6</i>	<b>0</b>	<b>3</b>						

Supplementary Table 36 Random Forest results for CB vs Cancer in the *RPLP2* and *TWIST1* normalised data.

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>HOXC6</i>	<b>0.76</b>	<b>167</b>	<i>ERG3' exons 6-7</i>	<b>1.05</b>	<b>87</b>	<i>ERG3' exons 6-7</i>	<b>0.76</b>	<b>65</b>
<i>SPINK1</i>	<b>0.73</b>	<b>166</b>	<i>APOC1</i>	<b>0.81</b>	<b>86</b>	<i>CCDC88B</i>	<b>0.53</b>	<b>64</b>

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>NAALADL2</i>	<i>0.67</i>	<i>165</i>	<i>SPINK1</i>	<i>0.65</i>	<i>85</i>	<i>CADPS</i>	<i>0.51</i>	<i>63</i>
<i>UPK2</i>	<i>0.61</i>	<i>164</i>	<i>CCDC88B</i>	<i>0.57</i>	<i>84</i>	<i>B4GALNT4</i>	<i>0.29</i>	<i>62</i>
<i>CADPS</i>	<i>0.43</i>	<i>163</i>	<i>CADPS</i>	<i>0.48</i>	<i>83</i>	<i>HOXC6</i>	<i>0.27</i>	<i>61</i>
						<i>TMPRSS2:ERG</i>		
<i>ERG3' exons 6-7</i>	<i>0.34</i>	<i>162</i>	<i>CKAP2L</i>	<i>0.43</i>	<i>82</i>	<i>fusion</i>	<i>0.24</i>	<i>60</i>
<i>HPN</i>	<i>0.32</i>	<i>161</i>	<i>GAPDH</i>	<i>0.37</i>	<i>81</i>	<i>RIOK3</i>	<i>0.19</i>	<i>59</i>
<i>ISX</i>	<i>0.29</i>	<i>160</i>	<i>CAMKK2</i>	<i>0.36</i>	<i>80</i>	<i>SIM2 long</i>	<i>0.19</i>	<i>58</i>
<i>CP</i>	<i>0.25</i>	<i>159</i>	<i>AURKA</i>	<i>0.22</i>	<i>79</i>	<i>MIR4435 1HG</i>	<i>0.18</i>	<i>57</i>
<i>TMPRSS2:ERG fusion</i>	<i>0.24</i>	<i>158</i>	<i>HPN</i>	<i>0.22</i>	<i>78</i>	<i>NEAT1</i>	<i>0.18</i>	<i>56</i>
<i>PCA3</i>	<i>0.20</i>	<i>157</i>	<i>UPK2</i>	<i>0.21</i>	<i>77</i>	<i>AATF</i>	<i>0.15</i>	<i>55</i>
<i>TDRD</i>	<i>0.19</i>	<i>156</i>	<i>AATF</i>	<i>0.21</i>	<i>76</i>	<i>SIRT1</i>	<i>0.13</i>	<i>54</i>
<i>B4GALNT4</i>	<i>0.18</i>	<i>155</i>	<i>B4GALNT4</i>	<i>0.20</i>	<i>75</i>	<i>APOC1</i>	<i>0.12</i>	<i>53</i>
<i>CKAP2L</i>	<i>0.16</i>	<i>154</i>	<i>IGFBP3</i>	<i>0.18</i>	<i>74</i>	<i>HPRT</i>	<i>0.11</i>	<i>52</i>
<i>ST6GALNAC1</i>	<i>0.15</i>	<i>153</i>	<i>ISX</i>	<i>0.18</i>	<i>73</i>	<i>MMP25</i>	<i>0.11</i>	<i>51</i>
<i>SFRP4</i>	<i>0.14</i>	<i>152</i>	<i>TDRD</i>	<i>0.17</i>	<i>72</i>	<i>TDRD</i>	<i>0.10</i>	<i>50</i>
<i>GCNT1</i>	<i>0.14</i>	<i>151</i>	<i>PCA3</i>	<i>0.17</i>	<i>71</i>	<i>MCTP1</i>	<i>0.09</i>	<i>49</i>
<i>Timp4</i>	<i>0.13</i>	<i>150</i>	<i>CD10</i>	<i>0.16</i>	<i>70</i>	<i>TMEM86A</i>	<i>0.09</i>	<i>48</i>
<i>APOC1</i>	<i>0.12</i>	<i>149</i>	<i>CLIC2</i>	<i>0.13</i>	<i>69</i>	<i>CLIC2</i>	<i>0.09</i>	<i>47</i>
<i>SLC43A1</i>	<i>0.12</i>	<i>148</i>	<i>NAALADL2</i>	<i>0.11</i>	<i>68</i>	<i>SFRP4</i>	<i>0.09</i>	<i>46</i>
<i>CLIC2</i>	<i>0.11</i>	<i>147</i>	<i>ERG3' exons 4-5</i>	<i>0.10</i>	<i>67</i>	<i>ERG5'</i>	<i>0.09</i>	<i>45</i>
<i>TMCC2</i>	<i>0.11</i>	<i>146</i>	<i>SLC4A1 S</i>	<i>0.10</i>	<i>66</i>	<i>MEX3A</i>	<i>0.09</i>	<i>44</i>
<i>EN2</i>	<i>0.09</i>	<i>145</i>	<i>CDC37L1</i>	<i>0.09</i>	<i>65</i>	<i>SLC43A1</i>	<i>0.08</i>	<i>43</i>
<i>AR exon 9</i>	<i>0.08</i>	<i>144</i>	<i>CACNA1D</i>	<i>0.09</i>	<i>64</i>	<i>MFS2A</i>	<i>0.08</i>	<i>42</i>
<i>RNF157</i>	<i>0.08</i>	<i>143</i>	<i>HIST1H2BG</i>	<i>0.09</i>	<i>63</i>	<i>SEC61A1</i>	<i>0.07</i>	<i>41</i>
<i>ANKRD34B</i>	<i>0.08</i>	<i>142</i>	<i>SNORA20</i>	<i>0.08</i>	<i>62</i>	<i>MAK</i>	<i>0.07</i>	<i>40</i>
<i>CLU</i>	<i>0.08</i>	<i>141</i>	<i>ACTR5</i>	<i>0.08</i>	<i>61</i>	<i>HPN</i>	<i>0.07</i>	<i>39</i>
<i>MMP25</i>	<i>0.07</i>	<i>140</i>	<i>TERF2IP</i>	<i>0.07</i>	<i>60</i>	<i>SULF2</i>	<i>0.07</i>	<i>38</i>
<i>SIM2 long</i>	<i>0.07</i>	<i>139</i>	<i>TMCC2</i>	<i>0.07</i>	<i>59</i>	<i>GCNT1</i>	<i>0.06</i>	<i>37</i>

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>TMEM86A</i>	<i>0.07</i>	<i>138</i>	<i>CDC20</i>	<i>0.06</i>	<i>58</i>	<i>EN2</i>	<i>0.06</i>	<i>36</i>
<i>HMBS</i>	<i>0.06</i>	<i>137</i>	<i>DPP4</i>	<i>0.06</i>	<i>57</i>	<i>SPINK1</i>	<i>0.06</i>	<i>35</i>
<i>ERG5'</i>	<i>0.06</i>	<i>136</i>	<i>TFDP1</i>	<i>0.06</i>	<i>56</i>	<i>PTPRC</i>	<i>0.06</i>	<i>34</i>
<i>DNAH5</i>	<i>0.06</i>	<i>135</i>	<i>AR exon 9</i>	<i>0.06</i>	<i>55</i>	<i>ANKRD34B</i>	<i>0.05</i>	<i>33</i>
<i>MSMB</i>	<i>0.05</i>	<i>134</i>	<i>MYOF</i>	<i>0.05</i>	<i>54</i>	<i>IGFBP3</i>	<i>0.05</i>	<i>32</i>
<i>MFSD2A</i>	<i>0.05</i>	<i>133</i>	<i>RNF157</i>	<i>0.05</i>	<i>53</i>	<i>UPK2</i>	<i>0.05</i>	<i>31</i>
<i>SERPINB5</i>	<i>0.05</i>	<i>132</i>	<i>AMH</i>	<i>0.05</i>	<i>52</i>	<i>AURKA</i>	<i>0.05</i>	<i>30</i>
<i>P712P</i>	<i>0.05</i>	<i>131</i>	<i>GABARAPL2</i>	<i>0.05</i>	<i>51</i>	<i>SNCA</i>	<i>0.05</i>	<i>29</i>
<i>CAMKK2</i>	<i>0.05</i>	<i>130</i>	<i>FOLH1</i>	<i>0.05</i>	<i>50</i>	<i>CACNA1D</i>	<i>0.05</i>	<i>28</i>
<i>TMEM47</i>	<i>0.04</i>	<i>129</i>	<i>ANPEP</i>	<i>0.05</i>	<i>49</i>	<i>LASS1</i>	<i>0.05</i>	<i>27</i>
<i>PPFIA2</i>	<i>0.04</i>	<i>128</i>	<i>EN2</i>	<i>0.05</i>	<i>48</i>	<i>GAPDH</i>	<i>0.05</i>	<i>26</i>
<i>ITGBL1</i>	<i>0.04</i>	<i>127</i>	<i>ST6GALNAC1</i>	<i>0.04</i>	<i>47</i>	<i>CAMKK2</i>	<i>0.04</i>	<i>25</i>
<i>MNX1</i>	<i>0.04</i>	<i>126</i>	<i>ANKRD34B</i>	<i>0.04</i>	<i>46</i>	<i>B2M</i>	<i>0.04</i>	<i>24</i>
<i>RIOK3</i>	<i>0.04</i>	<i>125</i>	<i>AMACR</i>	<i>0.04</i>	<i>45</i>	<i>ERG3' exons 4-5</i>	<i>0.04</i>	<i>23</i>
<i>GJB1</i>	<i>0.04</i>	<i>124</i>	<i>ERG5'</i>	<i>0.03</i>	<i>44</i>	<i>SLC12A1</i>	<i>0.04</i>	<i>22</i>
<i>TWIST1</i>	<i>0.04</i>	<i>123</i>	<i>CP</i>	<i>0.03</i>	<i>43</i>	<i>ITPR1</i>	<i>0.04</i>	<i>21</i>
<i>SRSF3</i>	<i>0.03</i>	<i>122</i>	<i>EIF2D</i>	<i>0.03</i>	<i>42</i>	<i>MAPK8IP2</i>	<i>0.03</i>	<i>20</i>
<i>AGR2</i>	<i>0.03</i>	<i>121</i>	<i>MCM7</i>	<i>0.03</i>	<i>41</i>	<i>SRSF3</i>	<i>0.03</i>	<i>19</i>
<i>PPAP2A</i>	<i>0.03</i>	<i>120</i>	<i>Met</i>	<i>0.03</i>	<i>40</i>	<i>ISX</i>	<i>0.03</i>	<i>18</i>
<i>PPP1R12B</i>	<i>0.03</i>	<i>119</i>	<i>DNAH5</i>	<i>0.02</i>	<i>39</i>	<i>FOLH1</i>	<i>0.03</i>	<i>17</i>
<i>STEAP4</i>	<i>0.03</i>	<i>118</i>	<i>SIM2 short</i>	<i>0.02</i>	<i>38</i>	<i>EIF2D</i>	<i>0.03</i>	<i>16</i>
<i>MYOF</i>	<i>0.03</i>	<i>117</i>	<i>SMAP1 exons 7-8</i>	<i>0.02</i>	<i>37</i>	<i>CDC20</i>	<i>0.03</i>	<i>15</i>
<i>STEAP2</i>	<i>0.03</i>	<i>116</i>	<i>AGR2</i>	<i>0.02</i>	<i>36</i>	<i>GABARAPL2</i>	<i>0.02</i>	<i>14</i>
<i>IGFBP3</i>	<i>0.03</i>	<i>115</i>	<i>NLRP3</i>	<i>0.02</i>	<i>35</i>	<i>MXI1</i>	<i>0.02</i>	<i>13</i>
<i>CCDC88B</i>	<i>0.03</i>	<i>114</i>	<i>KLK2</i>	<i>0.02</i>	<i>34</i>	<i>AMH</i>	<i>0.02</i>	<i>12</i>
<i>SLC4A1.S</i>	<i>0.03</i>	<i>113</i>	<i>AR exons 4-8</i>	<i>0.02</i>	<i>33</i>	<i>TBP</i>	<i>0.01</i>	<i>11</i>
<i>SULF2</i>	<i>0.03</i>	<i>112</i>	<i>MAK</i>	<i>0.02</i>	<i>32</i>	<i>PDLIM5</i>	<i>0.01</i>	<i>10</i>
<i>DLX1</i>	<i>0.03</i>	<i>111</i>	<i>TMEM47</i>	<i>0.02</i>	<i>31</i>	<i>ARHGEF25</i>	<i>0.01</i>	<i>9</i>

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>SEC61A1</i>	<i>0.03</i>	<i>110</i>	<i>CDKN3</i>	<i>0.02</i>	<i>30</i>	<i>ACTR5</i>	<i>0.01</i>	<i>8</i>
<i>PECI</i>	<i>0.03</i>	<i>109</i>	<i>RPS11</i>	<i>0.02</i>	<i>29</i>	<i>NLRP3</i>	<i>0.01</i>	<i>7</i>
<i>HIST1H2BG</i>	<i>0.03</i>	<i>108</i>	<i>PPAP2A</i>	<i>0.02</i>	<i>28</i>	<i>CD10</i>	<i>0.01</i>	<i>6</i>
<i>LASS1</i>	<i>0.03</i>	<i>107</i>	<i>PALM3</i>	<i>0.02</i>	<i>27</i>	<i>TERF2IP</i>	<i>0.005</i>	<i>5</i>
<i>NLRP3</i>	<i>0.02</i>	<i>106</i>	<i>RP11_97012.7</i>	<i>0.02</i>	<i>26</i>	<i>ANPEP</i>	<i>0.004</i>	<i>4</i>
<i>SULT1A1</i>	<i>0.02</i>	<i>105</i>	<i>CAMK2N2</i>	<i>0.02</i>	<i>25</i>	<i>MIC1</i>	<i>0.004</i>	<i>3</i>
<i>ACTR5</i>	<i>0.02</i>	<i>104</i>	<i>PECI</i>	<i>0.02</i>	<i>24</i>	<i>CASKIN1</i>	<i>0.003</i>	<i>2</i>
<i>MDK</i>	<i>0.02</i>	<i>103</i>	<i>FDPS</i>	<i>0.02</i>	<i>23</i>	<i>SACMIL</i>	<i>5.46x10<sup>-17</sup></i>	<i>1</i>
<i>SLC12A1</i>	<i>0.02</i>	<i>102</i>	<i>ARHGEF25</i>	<i>0.02</i>	<i>22</i>			
<i>TMEM45B</i>	<i>0.02</i>	<i>101</i>	<i>HOXC4</i>	<i>0.02</i>	<i>21</i>			
<i>MAK</i>	<i>0.02</i>	<i>100</i>	<i>MARCH5</i>	<i>0.01</i>	<i>20</i>			
<i>SIRT1</i>	<i>0.02</i>	<i>99</i>	<i>TBP</i>	<i>0.01</i>	<i>19</i>			
<i>MAPK8IP2</i>	<i>0.02</i>	<i>98</i>	<i>ABCB9</i>	<i>0.01</i>	<i>18</i>			
<i>MCTP1</i>	<i>0.02</i>	<i>97</i>	<i>B2M</i>	<i>0.01</i>	<i>17</i>			
<i>AATF</i>	<i>0.02</i>	<i>96</i>	<i>ALAS1</i>	<i>0.01</i>	<i>16</i>			
<i>RAB17</i>	<i>0.02</i>	<i>95</i>	<i>DLX1</i>	<i>0.01</i>	<i>15</i>			
<i>MEMO1</i>	<i>0.02</i>	<i>94</i>	<i>BTG2</i>	<i>0.01</i>	<i>14</i>			
<i>PALM3</i>	<i>0.02</i>	<i>93</i>	<i>PCSK6</i>	<i>0.01</i>	<i>13</i>			
<i>TRPM4</i>	<i>0.02</i>	<i>92</i>	<i>SSTR1</i>	<i>0.01</i>	<i>12</i>			
<i>SMIMI</i>	<i>0.02</i>	<i>91</i>	<i>STEAP2</i>	<i>0.01</i>	<i>11</i>			
<i>ABCB9</i>	<i>0.02</i>	<i>90</i>	<i>CLU</i>	<i>0.01</i>	<i>10</i>			
<i>MIR146A</i>	<i>0.02</i>	<i>89</i>	<i>LBH</i>	<i>0.01</i>	<i>9</i>			
<i>IMPDH2</i>	<i>0.02</i>	<i>88</i>	<i>MIATNB</i>	<i>0.01</i>	<i>8</i>			
<i>MGAT5B</i>	<i>0.02</i>	<i>87</i>	<i>COL10A1</i>	<i>0.01</i>	<i>7</i>			
<i>DPP4</i>	<i>0.02</i>	<i>86</i>	<i>COL9A2</i>	<i>0.01</i>	<i>6</i>			
<i>MIR4435_1HG</i>	<i>0.02</i>	<i>85</i>	<i>OGT</i>	<i>0.01</i>	<i>5</i>			
<i>CACNAID</i>	<i>0.01</i>	<i>84</i>	<i>MEX3A</i>	<i>0.01</i>	<i>4</i>			
<i>CDC20</i>	<i>0.01</i>	<i>83</i>	<i>GOLM1</i>	<i>0.01</i>	<i>3</i>			



## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>RPS10</i>	<b>0.01</b>	<b>82</b>	<i>CASKINI</i>	<b>0.004</b>	<b>2</b>			
<i>CASKINI</i>	<b>0.01</b>	<b>81</b>	<i>BRAF</i>	<b>0.002</b>	<b>1</b>			
<i>Met</i>	<b>0.01</b>	<b>80</b>						
<i>SPON2</i>	<b>0.01</b>	<b>79</b>						
<i>TERF2IP</i>	<b>0.01</b>	<b>78</b>						
<i>HIST1H1E</i>	<b>0.01</b>	<b>77</b>						
<i>GAPDH</i>	<b>0.01</b>	<b>76</b>						
<i>AURKA</i>	<b>0.01</b>	<b>75</b>						
<i>NKAIN1</i>	<b>0.01</b>	<b>74</b>						
<i>PVT1</i>	<b>0.01</b>	<b>73</b>						
<i>STOM</i>	<b>0.01</b>	<b>72</b>						
<i>VPS13A</i>	<b>0.01</b>	<b>71</b>						
<i>AMH</i>	<b>0.01</b>	<b>70</b>						
<i>COL9A2</i>	<b>0.01</b>	<b>69</b>						
<i>AMACR</i>	<b>0.01</b>	<b>68</b>						
<i>SIM2 short</i>	<b>0.01</b>	<b>67</b>						
<i>CD10</i>	<b>0.01</b>	<b>66</b>						
<i>FDPS</i>	<b>0.01</b>	<b>65</b>						
<i>MMP26</i>	<b>0.01</b>	<b>64</b>						
<i>MXI1</i>	<b>0.01</b>	<b>63</b>						
<i>ARHGEF25</i>	<b>0.01</b>	<b>62</b>						
<i>IFT57</i>	<b>0.01</b>	<b>61</b>						
<i>KLK2</i>	<b>0.01</b>	<b>60</b>						
<i>HOXC4</i>	<b>0.01</b>	<b>59</b>						
<i>KLK4</i>	<b>0.01</b>	<b>58</b>						
<i>MED4</i>	<b>0.01</b>	<b>57</b>						
<i>RPLP2</i>	<b>0.01</b>	<b>56</b>						
<i>CDKN3</i>	<b>0.01</b>	<b>55</b>						

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>CDC37L1</i>	<i>0.01</i>	<i>54</i>						
<i>MMP11</i>	<i>0.01</i>	<i>53</i>						
<i>AR exons 4-8</i>	<i>0.01</i>	<i>52</i>						
<i>RPS11</i>	<i>0.01</i>	<i>51</i>						
<i>SMAP1 exons 7-8</i>	<i>0.01</i>	<i>50</i>						
<i>FOLH1</i>	<i>0.01</i>	<i>49</i>						
<i>GOLM1</i>	<i>0.01</i>	<i>48</i>						
<i>PTN</i>	<i>0.01</i>	<i>47</i>						
<i>HIST3H2A</i>	<i>0.01</i>	<i>46</i>						
<i>ERG3' exons 4-5</i>	<i>0.01</i>	<i>45</i>						
<i>TERT</i>	<i>0.01</i>	<i>44</i>						
<i>MEX3A</i>	<i>0.01</i>	<i>43</i>						
<i>SYNM</i>	<i>0.01</i>	<i>42</i>						
<i>B2M</i>	<i>0.01</i>	<i>41</i>						
<i>SChLAPI</i>	<i>0.01</i>	<i>40</i>						
<i>RP11_97O12.7</i>	<i>0.01</i>	<i>39</i>						
<i>RPL18A</i>	<i>0.01</i>	<i>38</i>						
<i>GABARAPL2</i>	<i>0.01</i>	<i>37</i>						
<i>HIST1H1C</i>	<i>0.01</i>	<i>36</i>						
<i>BRAF</i>	<i>0.01</i>	<i>35</i>						
<i>SNORA20</i>	<i>0.01</i>	<i>34</i>						
<i>OR52A2</i>	<i>0.01</i>	<i>33</i>						
<i>ANPEP</i>	<i>0.01</i>	<i>32</i>						
<i>PSTPIP1</i>	<i>0.01</i>	<i>31</i>						
<i>RPL23AP53</i>	<i>0.01</i>	<i>30</i>						
<i>COL10A1</i>	<i>0.01</i>	<i>29</i>						
<i>SSTR1</i>	<i>0.01</i>	<i>28</i>						
<i>LBH</i>	<i>0.005</i>	<i>27</i>						

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ITPRI</i>	<i>0.005</i>	<i>26</i>						
<i>TFDPI</i>	<i>0.005</i>	<i>25</i>						
<i>CAMK2N2</i>	<i>0.004</i>	<i>24</i>						
<i>TBP</i>	<i>0.004</i>	<i>23</i>						
<i>PTPRC</i>	<i>0.004</i>	<i>22</i>						
<i>ZNF577</i>	<i>0.004</i>	<i>21</i>						
<i>MARCH5</i>	<i>0.003</i>	<i>20</i>						
<i>ALAS1</i>	<i>0.003</i>	<i>19</i>						
<i>HPRT</i>	<i>0.003</i>	<i>18</i>						
<i>OGT</i>	<i>0.003</i>	<i>17</i>						
<i>KLK3 exons 1-2</i>	<i>0.002</i>	<i>16</i>						
<i>MCM7</i>	<i>0.002</i>	<i>15</i>						
<i>VAX2</i>	<i>0.002</i>	<i>14</i>						
<i>SSPO</i>	<i>0.002</i>	<i>13</i>						
<i>BTG2</i>	<i>0.002</i>	<i>12</i>						
<i>MIC1</i>	<i>0.002</i>	<i>11</i>						
<i>NEAT1</i>	<i>0.002</i>	<i>10</i>						
<i>MKi67</i>	<i>0.002</i>	<i>9</i>						
<i>MIATNB</i>	<i>0.002</i>	<i>8</i>						
<i>EIF2D</i>	<i>0.002</i>	<i>7</i>						
<i>SACMIL</i>	<i>1.60x10<sup>-17</sup></i>	<i>6</i>						
<i>SNCA</i>	<i>1.08x10<sup>-17</sup></i>	<i>5</i>						
<i>PCSK6</i>	<i>7.99x10<sup>-18</sup></i>	<i>4</i>						
<i>HIST1H2BF</i>	<i>1.78x10<sup>-18</sup></i>	<i>3</i>						
<i>KLK3 exons 2-3</i>	<i>4.44x10<sup>-19</sup></i>	<i>2</i>						
<i>PDLIM5</i>	<i>0</i>	<i>1</i>						

**6.19 High Risk Vs CB**

Supplementary Table 37 Transcripts that have significant differential expression (using glm and MWU tests) between clinically benign and high-risk cancer samples in the baseline normalized NanoString data.

<i>Transcript</i>	<i>MWU</i>	<i>glm</i>		<i>Log<sub>2</sub>(FC)</i>	
	<i>p-value</i>	<i>Adjusted p-value</i>	<i>p-value</i>		<i>Adjusted p-value</i>
<i>HOXC6</i>	<b>0.0002</b>	<b>0.0299</b>	<b>0.004</b>	<b>0.6711</b>	<b>2</b>
<i>ERG3' exons 6-7</i>	<b>6.21x10<sup>-06</sup></b>	<b>0.001</b>	<b>0.0371</b>	<b>0.9942</b>	<b>1.6</b>
<i>TDRD</i>	<b>0.0011</b>	<b>0.1558</b>	<b>0.0333</b>	<b>0.9942</b>	<b>1.5</b>
<i>TMPRSS2-ERG</i>	<b>0.0004</b>	<b>0.0668</b>	<b>0.0386</b>	<b>0.9942</b>	<b>1.3</b>
<i>B4GALNT4</i>	<b>2.88x10<sup>-05</sup></b>	<b>0.0048</b>	<b>0.0409</b>	<b>0.9942</b>	<b>1.2</b>
<i>SLC43A1</i>	<b>0.002</b>	<b>0.2897</b>	<b>0.0117</b>	<b>0.9942</b>	<b>1.2</b>
<i>CADPS</i>	<b>6.70x10<sup>-05</sup></b>	<b>0.011</b>	<b>0.02</b>	<b>0.9942</b>	<b>1.1</b>
<i>CLIC2</i>	<b>0.0002</b>	<b>0.0386</b>	<b>0.0087</b>	<b>0.9942</b>	<b>1</b>
<i>HPN</i>	<b>0.0008</b>	<b>0.1258</b>	<b>0.0092</b>	<b>0.9942</b>	<b>0.9</b>
<i>LASS1</i>	<b>0.0011</b>	<b>0.1558</b>	<b>0.0103</b>	<b>0.9942</b>	<b>0.9</b>
<i>MAPK8IP2</i>	<b>0.0148</b>	<b>1</b>	<b>0.0336</b>	<b>0.9942</b>	<b>0.9</b>
<i>SFRP4</i>	<b>0.0013</b>	<b>0.1919</b>	<b>0.0155</b>	<b>0.9942</b>	<b>0.9</b>
<i>CKAP2L</i>			<b>0.0392</b>	<b>0.9942</b>	<b>0.9</b>
<i>CDKN3</i>			<b>0.0326</b>	<b>0.9942</b>	<b>0.9</b>
<i>ANKRD34B</i>	<b>0.0054</b>	<b>0.7002</b>	<b>0.0368</b>	<b>0.9942</b>	<b>0.8</b>
<i>ERG3' exons 4-5</i>	<b>0.0037</b>	<b>0.5042</b>	<b>0.0434</b>	<b>0.9942</b>	<b>0.8</b>
<i>APOC1</i>	<b>0.0002</b>	<b>0.0386</b>	<b>0.0055</b>	<b>0.9103</b>	<b>0.7</b>
<i>ERG5'</i>	<b>0.0077</b>	<b>0.9678</b>			<b>0.7</b>
<i>MMP25</i>	<b>0.0045</b>	<b>0.5959</b>	<b>0.0162</b>	<b>0.9942</b>	<b>0.7</b>
<i>AMH</i>	<b>0.0108</b>	<b>1</b>			<b>0.6</b>
<i>CCDC88B</i>	<b>0.0007</b>	<b>0.1014</b>	<b>0.0104</b>	<b>0.9942</b>	<b>0.6</b>
<i>FOLH1</i>	<b>0.0108</b>	<b>1</b>	<b>0.027</b>	<b>0.9942</b>	<b>0.6</b>
<i>ISX</i>	<b>0.0026</b>	<b>0.3638</b>	<b>0.0234</b>	<b>0.9942</b>	<b>0.6</b>
<i>MCTP1</i>	<b>0.0001</b>	<b>0.0228</b>	<b>0.0098</b>	<b>0.9942</b>	<b>0.6</b>
<i>SIM2 long</i>	<b>0.0002</b>	<b>0.0386</b>	<b>0.0124</b>	<b>0.9942</b>	<b>0.6</b>
<i>SRSF3</i>	<b>0.0234</b>	<b>1</b>	<b>0.0225</b>	<b>0.9942</b>	<b>0.6</b>
<i>ANPEP</i>	<b>0.0234</b>	<b>1</b>	<b>0.0324</b>	<b>0.9942</b>	<b>0.5</b>
<i>GCNT1</i>	<b>0.0054</b>	<b>0.7002</b>	<b>0.0209</b>	<b>0.9942</b>	<b>0.5</b>
<i>MFS2A</i>	<b>0.0017</b>	<b>0.2364</b>	<b>0.0216</b>	<b>0.9942</b>	<b>0.5</b>
<i>NLRP3</i>	<b>0.0202</b>	<b>1</b>	<b>0.0277</b>	<b>0.9942</b>	<b>0.5</b>
<i>SLC12A1</i>	<b>0.0064</b>	<b>0.8317</b>	<b>0.0267</b>	<b>0.9942</b>	<b>0.5</b>
<i>SULF2</i>	<b>0.0007</b>	<b>0.1014</b>	<b>0.0141</b>	<b>0.9942</b>	<b>0.5</b>
<i>TMEM86A</i>	<b>0.0005</b>	<b>0.0823</b>	<b>0.022</b>	<b>0.9942</b>	<b>0.5</b>
<i>AATF</i>	<b>0.0007</b>	<b>0.1014</b>	<b>0.0108</b>	<b>0.9942</b>	<b>0.4</b>
<i>CAMKK2</i>	<b>0.0007</b>	<b>0.1014</b>	<b>0.0128</b>	<b>0.9942</b>	<b>0.4</b>
<i>CDC20</i>	<b>0.0202</b>	<b>1</b>	<b>0.0343</b>	<b>0.9942</b>	<b>0.4</b>
<i>EN2</i>	<b>0.0091</b>	<b>1</b>			<b>0.4</b>

9: APPENDICES

<i>ARHGEF25</i>	<b>0.0464</b>	<b>1</b>			<b>0.3</b>
<i>AURKA</i>	<b>0.0234</b>	<b>1</b>			<b>0.3</b>
<i>CD10</i>	<b>0.0127</b>	<b>1</b>	<b>0.0481</b>	<b>0.9942</b>	<b>0.3</b>
<i>HPRT</i>	<b>0.0025</b>	<b>0.3503</b>	<b>0.0255</b>	<b>0.9942</b>	<b>0.3</b>
<i>MEX3A</i>	<b>0.0045</b>	<b>0.5959</b>	<b>0.0496</b>	<b>0.9942</b>	<b>0.3</b>
<i>MIC1</i>	<b>0.031</b>	<b>1</b>			<b>0.3</b>
<i>PTPRC</i>	<b>0.0013</b>	<b>0.1919</b>	<b>0.0116</b>	<b>0.9942</b>	<b>0.3</b>
<i>RIOK3</i>	<b>5.63x10<sup>-05</sup></b>	<b>0.0093</b>	<b>0.0237</b>	<b>0.9942</b>	<b>0.3</b>
<i>SEC61A1</i>	<b>0.0031</b>	<b>0.4184</b>	<b>0.0341</b>	<b>0.9942</b>	<b>0.3</b>
<i>SIRT1</i>	<b>0.0025</b>	<b>0.3503</b>	<b>0.0206</b>	<b>0.9942</b>	<b>0.3</b>
<i>SNCA</i>	<b>0.027</b>	<b>1</b>			<b>0.3</b>
<i>ACTR5</i>	<b>0.0127</b>	<b>1</b>			<b>0.2</b>
<i>CACNA1D</i>	<b>0.0202</b>	<b>1</b>	<b>0.0356</b>	<b>0.9942</b>	<b>0.2</b>
<i>CASKIN1</i>	<b>0.0464</b>	<b>1</b>			<b>0.2</b>
<i>EIF2D</i>	<b>0.0077</b>	<b>0.9678</b>	<b>0.0297</b>	<b>0.9942</b>	<b>0.2</b>
<i>GABARAPL2</i>	<b>0.0127</b>	<b>1</b>	<b>0.036</b>	<b>0.9942</b>	<b>0.2</b>
<i>ITPR1</i>	<b>0.0008</b>	<b>0.1258</b>	<b>0.0091</b>	<b>0.9942</b>	<b>0.2</b>
<i>MAK</i>	<b>0.0148</b>	<b>1</b>			<b>0.2</b>
<i>MIR4435 1HG</i>	<b>0.0007</b>	<b>0.1014</b>	<b>0.0169</b>	<b>0.9942</b>	<b>0.2</b>
<i>MXI1</i>	<b>0.0054</b>	<b>0.7002</b>	<b>0.0206</b>	<b>0.9942</b>	<b>0.2</b>
<i>NEAT1</i>	<b>0.0002</b>	<b>0.0299</b>	<b>0.0039</b>	<b>0.6541</b>	<b>0.2</b>
<i>PDLIM5</i>	<b>0.0464</b>	<b>1</b>			<b>0.2</b>
<i>TBP</i>	<b>0.0202</b>	<b>1</b>			<b>0.2</b>
<i>B2M</i>	<b>0.0008</b>	<b>0.1258</b>	<b>0.0203</b>	<b>0.9942</b>	<b>0.1</b>
<i>GAPDH</i>	<b>0.0031</b>	<b>0.4184</b>	<b>0.0111</b>	<b>0.9942</b>	<b>0.1</b>
<i>SACM1L</i>	<b>0.0234</b>	<b>1</b>			<b>0.1</b>
<i>TERF2IP</i>	<b>0.0045</b>	<b>0.5959</b>	<b>0.0197</b>	<b>0.9942</b>	<b>0.1</b>
<i>IGFBP3</i>	<b>0.0464</b>	<b>1</b>			<b>-0.3</b>
<i>SPINK1</i>	<b>0.0077</b>	<b>0.9678</b>			<b>-0.4</b>
<i>UPK2</i>	<b>0.0202</b>	<b>1</b>			<b>-0.8</b>

Supplementary Table 38 Transcripts that have significant differential expression (using glm and MWU tests) between clinically benign and high-risk cancer samples in the *KLK2* ratio NanoString data.

<i>Transcript</i>	<i>MWU</i>		<i>glm</i>		<i>Log<sub>2</sub>(FC)</i>
	<i>p-value</i>	<i>Adjusted p-value</i>	<i>p-value</i>	<i>Adjusted p-value</i>	
<i>TMPRSS2:ERG</i>	<b>0.004</b>	<b>0.68</b>	<b>0.028</b>	<b>1.000</b>	<b>0.25</b>
<i>ERG 3' exons 6-7</i>	<b>0.000</b>	<b>0.07</b>	<b>0.008</b>	<b>1.000</b>	<b>0.25</b>
<i>HOXC6</i>	<b>4.28E-05</b>	<b>0.01</b>			<b>0.25</b>
<i>TDRD</i>	<b>0.001</b>	<b>0.09</b>	<b>0.017</b>	<b>1.000</b>	<b>0.24</b>
<i>SLC43A1</i>	<b>0.002</b>	<b>0.27</b>	<b>0.022</b>	<b>1.000</b>	<b>0.21</b>
<i>CADPS</i>	<b>0.007</b>	<b>1</b>			<b>0.18</b>
<i>B4GALNT4</i>	<b>0.002</b>	<b>0.33</b>	<b>0.035</b>	<b>1.000</b>	<b>0.17</b>
<i>ERG 5'</i>	<b>0.027</b>	<b>1</b>			<b>0.16</b>
<i>SLC12A1</i>	<b>0.013</b>	<b>1</b>			<b>0.15</b>
<i>ERG 3' exons 4-5</i>	<b>0.046</b>	<b>1</b>	<b>0.050</b>	<b>1.000</b>	<b>0.14</b>

9: APPENDICES

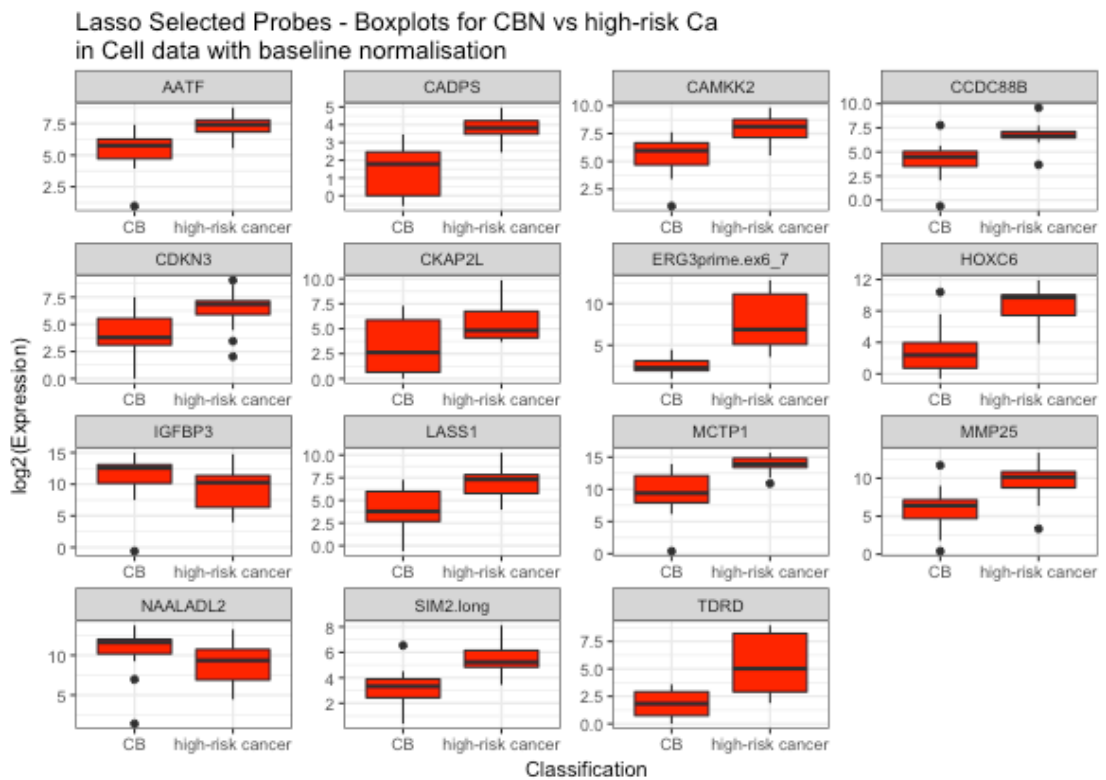
LASS1	0.015	1	0.029	1.000	0.13
CLIC2	0.004	0.59	0.027	1.000	0.13
HPN	0.003	0.49	0.032	1.000	0.11
ISX	0.017	1			0.11
APOC1	0.009	1			0.09
TMEM86A	0.017	1			0.09
PCA3	0.003	0.49	0.015	1.000	0.08
CCDC88B	0.027	1			0.08
SFRP4	0.031	1			0.08
MCTP1	0.020	1			0.08
SIM2 long	0.001	0.14	0.028	1.000	0.08
FOLH1	0.008	1			0.07
CAMKK2	0.036	1			0.05
SEC61A1	0.046	1			0.05
GCNT1	0.027	1	0.043	1.000	0.04

**Supplementary Table 39** Transcripts that have significant differential expression (using glm and MWU tests) between clinically benign and high-risk cancer samples in the HK normalised NanoString data.

Transcript	MWU	glm			Log <sub>2</sub> (FC)
	p-value	Adjusted p-value	p-value	Adjusted p-value	
HOXC6	0.0005	0.0882	0.0059	0.9765	1.6
ERG3' exons 6-7	0.0013	0.2186	0.0266	0.9765	1.4
TDRD	0.0031	0.4948	0.0272	0.9765	1.1
TMPRSS2:ERG fusion	0.0094	1	0.033	0.9765	1.1
ST6GALNAC1	0.0037	0.5969	0.0168	0.9765	-1
SLC43A1	0.0013	0.2186	0.0197	0.9765	0.9
B4GALNT4	0.0202	1			0.8
HPN	0.0077	1	0.0314	0.9765	0.8
CADPS	0.0145	1	0.0326	0.9765	0.7
CCDC88B	0.031	1	0.0482	0.9765	0.7
SPINK1	0.0007	0.1115	0.0092	0.9765	-0.7
UPK2	0.0054	0.8564	0.0237	0.9765	-0.7
CLIC2	0.0108	1	0.0278	0.9765	0.6
LASS1	0.0202	1	0.0451	0.9765	0.6
GJB1	0.0108	1	0.0197	0.9765	-0.6
IGFBP3	0.0464	1			-0.6
NAALADL2	0.0031	0.4948	0.0133	0.9765	-0.6
SERPINB5	0.0054	0.8564	0.0199	0.9765	-0.6
ISX	0.0288	1			0.5
MMP25	0.0407	1			0.5
GCNT1	0.0356	1	0.0446	0.9765	0.4
MCTP1	0.0464	1			0.4
SIM2 long	0.0234	1	0.0317	0.9765	0.4
PALM3	0.0108	1	0.0394	0.9765	-0.4

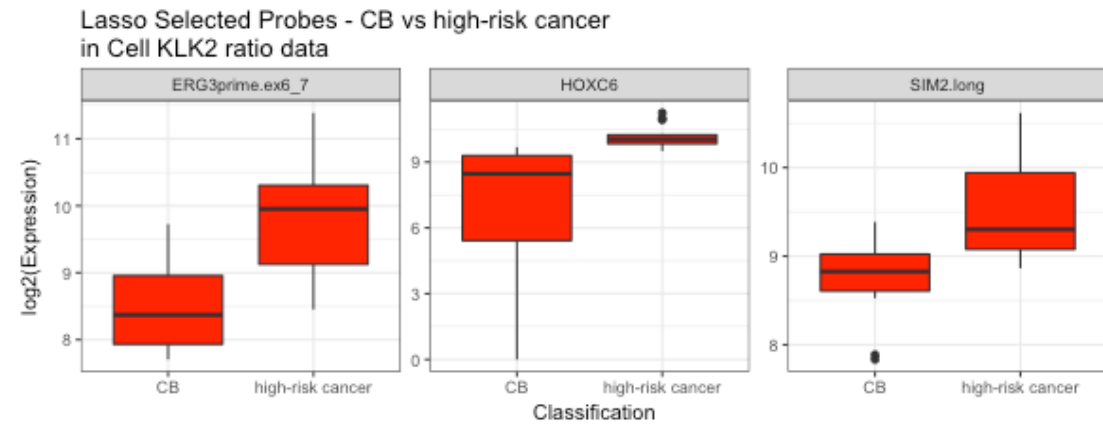
9: APPENDICES

<i>APOC1</i>	<b>0.0173</b>	<b>1</b>			<b>0.3</b>
<i>MSMB</i>	<b>0.0356</b>	<b>1</b>			<b>-0.3</b>
<i>PPAP2A</i>	<b>0.0077</b>	<b>1</b>			<b>-0.3</b>
<i>RAB17</i>	<b>0.0356</b>	<b>1</b>			<b>-0.3</b>
<i>RPS10</i>	<b>0.0464</b>	<b>1</b>			<b>-0.3</b>
<i>SPON2</i>	<b>0.0464</b>	<b>1</b>			<b>-0.3</b>
<i>STEAP2</i>	<b>0.0173</b>	<b>1</b>			<b>-0.3</b>
<i>VAX2</i>	<b>0.0109</b>	<b>1</b>			<b>-0.3</b>
<i>TMEM86A</i>	<b>0.0464</b>	<b>1</b>			<b>0.2</b>
<i>IFT57</i>	<b>0.027</b>	<b>1</b>			<b>-0.2</b>
<i>PTN</i>	<b>0.031</b>	<b>1</b>	<b>0.0491</b>	<b>0.9765</b>	<b>-0.2</b>

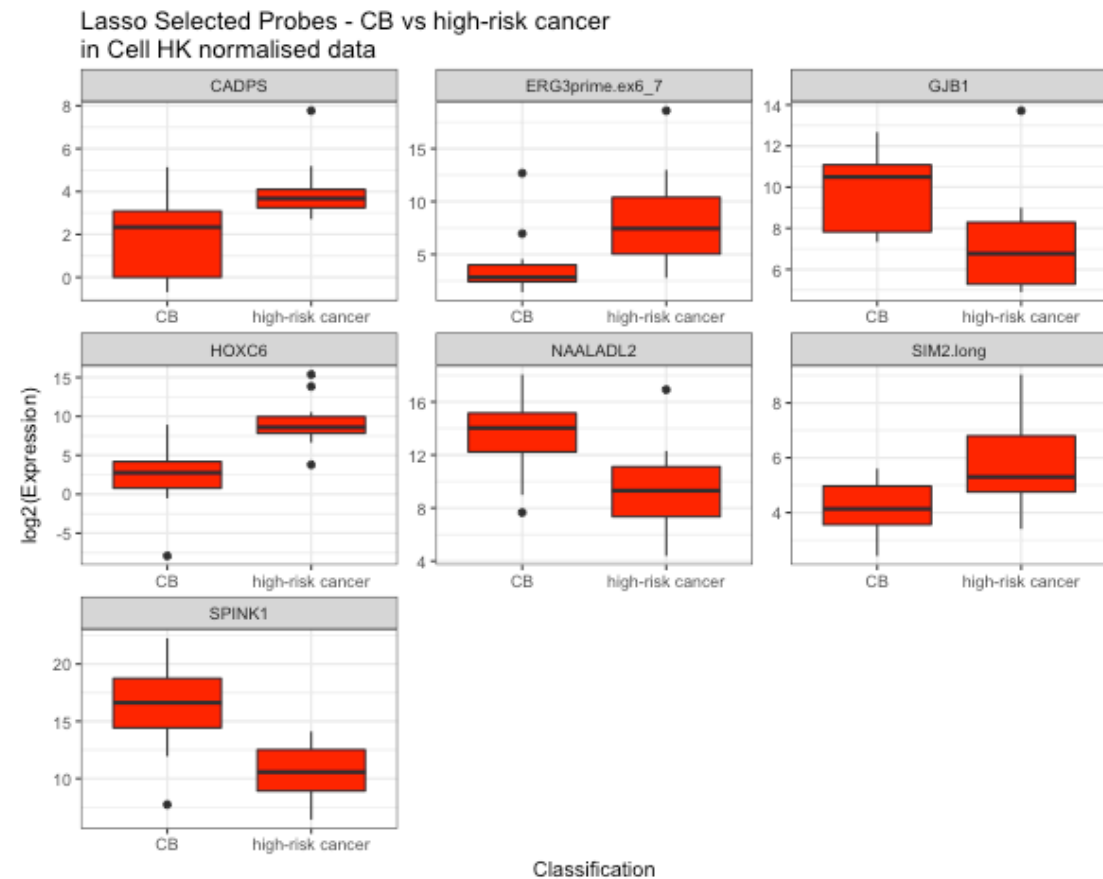


Supplementary Figure 16 Boxplots showing all of the Lasso selected probes in CB Vs. high-risk cancer models in baseline normalised cell data.

## 9: APPENDICES



**Supplementary Figure 17** Boxplots showing all of the Lasso selected probes in CB Vs. high-risk cancer models in *KLK2* ratio cell data.



**Supplementary Figure 18** Boxplots showing all of the Lasso selected probes in CB Vs. high-risk cancer models in HK normalised cell data.



9: APPENDICES

Supplementary Table 40 Random Forest results for HR-Ca vs CBN.

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 51)</i>			<i>Transcripts identified by Mann Whitney U (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>SPINK1</i>	<b>0.70</b>	<b>167</b>	<i>CADPS</i>	<b>0.67</b>	<b>51</b>	<i>ERG3' exons 6-7</i>	<b>0.76</b>	<b>65</b>
<i>CADPS</i>	<b>0.69</b>	<b>166</b>	<i>ERG3' exons 6-7</i>	<b>0.66</b>	<b>50</b>	<i>CCDC88B</i>	<b>0.53</b>	<b>64</b>
<i>NAALADL2</i>	<b>0.67</b>	<b>165</b>	<i>CCDC88B</i>	<b>0.51</b>	<b>49</b>	<i>CADPS</i>	<b>0.51</b>	<b>63</b>
<i>HOXC6</i>	<b>0.64</b>	<b>164</b>	<i>RIOK3</i>	<b>0.34</b>	<b>48</b>	<i>B4GALNT4</i>	<b>0.29</b>	<b>62</b>
<i>HPN</i>	<b>0.57</b>	<b>163</b>	<i>HOXC6</i>	<b>0.27</b>	<b>47</b>	<i>HOXC6</i>	<b>0.27</b>	<b>61</b>
<i>ERG3' exons 6-7</i>	<b>0.49</b>	<b>162</b>	<i>B4GALNT4</i>	<b>0.26</b>	<b>46</b>	<i>TMPRSS2:ERG fusion</i>	<b>0.24</b>	<b>60</b>
<i>UPK2</i>	<b>0.47</b>	<b>161</b>	<i>TMPRSS2:ERG fusion</i>	<b>0.26</b>	<b>45</b>	<i>RIOK3</i>	<b>0.19</b>	<b>59</b>
<i>TMPRSS2:ERG fusion</i>	<b>0.43</b>	<b>160</b>	<i>SIM2 long</i>	<b>0.25</b>	<b>44</b>	<i>SIM2 long</i>	<b>0.19</b>	<b>58</b>
<i>CP</i>	<b>0.42</b>	<b>159</b>	<i>MIR4435_IHG</i>	<b>0.22</b>	<b>43</b>	<i>MIR4435_IHG</i>	<b>0.18</b>	<b>57</b>
<i>TDRD</i>	<b>0.39</b>	<b>158</b>	<i>SFRP4</i>	<b>0.21</b>	<b>42</b>	<i>NEAT1</i>	<b>0.18</b>	<b>56</b>
<i>MNX1</i>	<b>0.39</b>	<b>157</b>	<i>HPRT</i>	<b>0.17</b>	<b>41</b>	<i>AATF</i>	<b>0.15</b>	<b>55</b>
<i>PCA3</i>	<b>0.36</b>	<b>156</b>	<i>APOC1</i>	<b>0.16</b>	<b>40</b>	<i>SIRT1</i>	<b>0.13</b>	<b>54</b>
<i>ST6GALNAC1</i>	<b>0.34</b>	<b>155</b>	<i>AATF</i>	<b>0.16</b>	<b>39</b>	<i>APOC1</i>	<b>0.12</b>	<b>53</b>
<i>TMCC2</i>	<b>0.34</b>	<b>154</b>	<i>NEAT1</i>	<b>0.15</b>	<b>38</b>	<i>HPRT</i>	<b>0.11</b>	<b>52</b>
<i>SIM2 long</i>	<b>0.34</b>	<b>153</b>	<i>TMEM86A</i>	<b>0.13</b>	<b>37</b>	<i>MMP25</i>	<b>0.11</b>	<b>51</b>
<i>APOC1</i>	<b>0.34</b>	<b>152</b>	<i>MEX3A</i>	<b>0.13</b>	<b>36</b>	<i>TDRD</i>	<b>0.10</b>	<b>50</b>
<i>CKAP2L</i>	<b>0.31</b>	<b>151</b>	<i>SLC43A1</i>	<b>0.11</b>	<b>35</b>	<i>MCTP1</i>	<b>0.09</b>	<b>49</b>
<i>PPFIA2</i>	<b>0.29</b>	<b>150</b>	<i>HPN</i>	<b>0.11</b>	<b>34</b>	<i>TMEM86A</i>	<b>0.09</b>	<b>48</b>
<i>SERPINB5</i>	<b>0.27</b>	<b>149</b>	<i>SEC61A1</i>	<b>0.10</b>	<b>33</b>	<i>CLIC2</i>	<b>0.09</b>	<b>47</b>
<i>TMEM45B</i>	<b>0.27</b>	<b>148</b>	<i>SIRT1</i>	<b>0.10</b>	<b>32</b>	<i>SFRP4</i>	<b>0.09</b>	<b>46</b>
<i>AGR2</i>	<b>0.26</b>	<b>147</b>	<i>CLIC2</i>	<b>0.09</b>	<b>31</b>	<i>ERG5'</i>	<b>0.09</b>	<b>45</b>
<i>EN2</i>	<b>0.25</b>	<b>146</b>	<i>GCNT1</i>	<b>0.09</b>	<b>30</b>	<i>MEX3A</i>	<b>0.09</b>	<b>44</b>

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 51)</i>			<i>Transcripts identified by Mann Whitney U (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ISX</i>	<b>0.25</b>	<b>145</b>	<i>MCTP1</i>	<b>0.08</b>	<b>29</b>	<i>SLC43A1</i>	<b>0.08</b>	<b>43</b>
<i>GCNT1</i>	<b>0.25</b>	<b>144</b>	<i>MFSD2A</i>	<b>0.08</b>	<b>28</b>	<i>MFSD2A</i>	<b>0.08</b>	<b>42</b>
<i>MFSD2A</i>	<b>0.24</b>	<b>143</b>	<i>TDRD</i>	<b>0.08</b>	<b>27</b>	<i>SEC61A1</i>	<b>0.07</b>	<b>41</b>
<i>DNAH5</i>	<b>0.24</b>	<b>142</b>	<i>SULF2</i>	<b>0.06</b>	<b>26</b>	<i>MAK</i>	<b>0.07</b>	<b>40</b>
<i>SFRP4</i>	<b>0.24</b>	<b>141</b>	<i>PTPRC</i>	<b>0.06</b>	<b>25</b>	<i>HPN</i>	<b>0.07</b>	<b>39</b>
<i>SLC43A1</i>	<b>0.24</b>	<b>140</b>	<i>GAPDH</i>	<b>0.06</b>	<b>24</b>	<i>SULF2</i>	<b>0.07</b>	<b>38</b>
<i>B4GALNT4</i>	<b>0.24</b>	<b>139</b>	<i>ISX</i>	<b>0.06</b>	<b>23</b>	<i>GCNT1</i>	<b>0.06</b>	<b>37</b>
<i>PTN</i>	<b>0.24</b>	<b>138</b>	<i>ANKRD34B</i>	<b>0.05</b>	<b>22</b>	<i>EN2</i>	<b>0.06</b>	<b>36</b>
<i>GJB1</i>	<b>0.23</b>	<b>137</b>	<i>MMP25</i>	<b>0.05</b>	<b>21</b>	<i>SPINK1</i>	<b>0.06</b>	<b>35</b>
<i>MMP25</i>	<b>0.23</b>	<b>136</b>	<i>ITPR1</i>	<b>0.04</b>	<b>20</b>	<i>PTPRC</i>	<b>0.06</b>	<b>34</b>
<i>Timp4</i>	<b>0.22</b>	<b>135</b>	<i>CACNA1D</i>	<b>0.04</b>	<b>19</b>	<i>ANKRD34B</i>	<b>0.05</b>	<b>33</b>
<i>RIOK3</i>	<b>0.21</b>	<b>134</b>	<i>MXI1</i>	<b>0.04</b>	<b>18</b>	<i>IGFBP3</i>	<b>0.05</b>	<b>32</b>
<i>MDK</i>	<b>0.21</b>	<b>133</b>	<i>SRSF3</i>	<b>0.04</b>	<b>17</b>	<i>UPK2</i>	<b>0.05</b>	<b>31</b>
<i>CLU</i>	<b>0.20</b>	<b>132</b>	<i>LASS1</i>	<b>0.03</b>	<b>16</b>	<i>AURKA</i>	<b>0.05</b>	<b>30</b>
<i>LASS1</i>	<b>0.20</b>	<b>131</b>	<i>B2M</i>	<b>0.03</b>	<b>15</b>	<i>SNCA</i>	<b>0.05</b>	<b>29</b>
<i>MMP11</i>	<b>0.20</b>	<b>130</b>	<i>SLC12A1</i>	<b>0.03</b>	<b>14</b>	<i>CACNA1D</i>	<b>0.05</b>	<b>28</b>
<i>ERG3' exons 4-5</i>	<b>0.19</b>	<b>129</b>	<i>GABARAPL2</i>	<b>0.03</b>	<b>13</b>	<i>LASS1</i>	<b>0.05</b>	<b>27</b>
<i>VAX2</i>	<b>0.19</b>	<b>128</b>	<i>ERG3' exons 4-5</i>	<b>0.03</b>	<b>12</b>	<i>GAPDH</i>	<b>0.05</b>	<b>26</b>
<i>SPON2</i>	<b>0.18</b>	<b>127</b>	<i>EIF2D</i>	<b>0.03</b>	<b>11</b>	<i>CAMKK2</i>	<b>0.04</b>	<b>25</b>
<i>PPAP2A</i>	<b>0.18</b>	<b>126</b>	<i>MAPK8IP2</i>	<b>0.02</b>	<b>10</b>	<i>B2M</i>	<b>0.04</b>	<b>24</b>
<i>TMEM47</i>	<b>0.17</b>	<b>125</b>	<i>FOLH1</i>	<b>0.02</b>	<b>9</b>	<i>ERG3' exons 4-5</i>	<b>0.04</b>	<b>23</b>
<i>CLIC2</i>	<b>0.17</b>	<b>124</b>	<i>CAMKK2</i>	<b>0.02</b>	<b>8</b>	<i>SLC12A1</i>	<b>0.04</b>	<b>22</b>
<i>SLC12A1</i>	<b>0.17</b>	<b>123</b>	<i>ANPEP</i>	<b>0.02</b>	<b>7</b>	<i>ITPR1</i>	<b>0.04</b>	<b>21</b>
<i>COL9A2</i>	<b>0.17</b>	<b>122</b>	<i>CDC20</i>	<b>0.01</b>	<b>6</b>	<i>MAPK8IP2</i>	<b>0.03</b>	<b>20</b>
<i>ANKRD34B</i>	<b>0.17</b>	<b>121</b>	<i>CKAP2L</i>	<b>0.01</b>	<b>5</b>	<i>SRSF3</i>	<b>0.03</b>	<b>19</b>
<i>TWIST1</i>	<b>0.17</b>	<b>120</b>	<i>CDKN3</i>	<b>0.01</b>	<b>4</b>	<i>ISX</i>	<b>0.03</b>	<b>18</b>
<i>SSPO</i>	<b>0.17</b>	<b>119</b>	<i>TERF2IP</i>	<b>0.01</b>	<b>3</b>	<i>FOLH1</i>	<b>0.03</b>	<b>17</b>

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 51)</i>			<i>Transcripts identified by Mann Whitney U (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>MYOF</i>	<b>0.17</b>	<b>118</b>	<i>CD10</i>	<b>0.01</b>	<b>2</b>	<i>EIF2D</i>	<b>0.03</b>	<b>16</b>
<i>CCDC88B</i>	<b>0.16</b>	<b>117</b>	<i>NLRP3</i>	<b>0.01</b>	<b>1</b>	<i>CDC20</i>	<b>0.03</b>	<b>15</b>
<i>SIM2 short</i>	<b>0.16</b>	<b>116</b>				<i>GABARAPL2</i>	<b>0.02</b>	<b>14</b>
<i>DLX1</i>	<b>0.16</b>	<b>115</b>				<i>MXI1</i>	<b>0.02</b>	<b>13</b>
<i>CAMKK2</i>	<b>0.16</b>	<b>114</b>				<i>AMH</i>	<b>0.02</b>	<b>12</b>
<i>IGFBP3</i>	<b>0.15</b>	<b>113</b>				<i>TBP</i>	<b>0.01</b>	<b>11</b>
<i>IFT57</i>	<b>0.15</b>	<b>112</b>				<i>PDLIM5</i>	<b>0.01</b>	<b>10</b>
<i>MMP26</i>	<b>0.15</b>	<b>111</b>				<i>ARHGEF25</i>	<b>0.01</b>	<b>9</b>
<i>SNORA20</i>	<b>0.15</b>	<b>110</b>				<i>ACTR5</i>	<b>0.01</b>	<b>8</b>
<i>RNF157</i>	<b>0.14</b>	<b>109</b>				<i>NLRP3</i>	<b>0.01</b>	<b>7</b>
<i>TMEM86A</i>	<b>0.14</b>	<b>108</b>				<i>CD10</i>	<b>0.01</b>	<b>6</b>
<i>MSMB</i>	<b>0.14</b>	<b>107</b>				<i>TERF2IP</i>	<b>0.01</b>	<b>5</b>
<i>P712P</i>	<b>0.14</b>	<b>106</b>				<i>ANPEP</i>	<b>0.00</b>	<b>4</b>
<i>PALM3</i>	<b>0.14</b>	<b>105</b>				<i>MIC1</i>	<b>0.00</b>	<b>3</b>
<i>SLC4A1.S</i>	<b>0.14</b>	<b>104</b>				<i>CASKIN1</i>	<b>0.00</b>	<b>2</b>
<i>MAPK8IP2</i>	<b>0.14</b>	<b>103</b>				<i>SACMIL</i>	<b>0.00</b>	<b>1</b>
<i>MCTP1</i>	<b>0.14</b>	<b>102</b>						
<i>ERG5'</i>	<b>0.14</b>	<b>101</b>						
<i>FOLH1</i>	<b>0.14</b>	<b>100</b>						
<i>AMH</i>	<b>0.13</b>	<b>99</b>						
<i>SEC61A1</i>	<b>0.13</b>	<b>98</b>						
<i>AR.ex9</i>	<b>0.13</b>	<b>97</b>						
<i>ABCB9</i>	<b>0.13</b>	<b>96</b>						
<i>MIR146A</i>	<b>0.13</b>	<b>95</b>						
<i>RPS11</i>	<b>0.12</b>	<b>94</b>						
<i>RAB17</i>	<b>0.12</b>	<b>93</b>						
<i>OR52A2</i>	<b>0.12</b>	<b>92</b>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 51)</i>			<i>Transcripts identified by Mann Whitney U (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ACTR5</i>	<b>0.11</b>	<b>91</b>						
<i>RPS10</i>	<b>0.11</b>	<b>90</b>						
<i>Met</i>	<b>0.11</b>	<b>89</b>						
<i>OGT</i>	<b>0.11</b>	<b>88</b>						
<i>STEAP2</i>	<b>0.11</b>	<b>87</b>						
<i>MEX3A</i>	<b>0.11</b>	<b>86</b>						
<i>ITGBL1</i>	<b>0.11</b>	<b>85</b>						
<i>PECI</i>	<b>0.10</b>	<b>84</b>						
<i>SSTR1</i>	<b>0.10</b>	<b>83</b>						
<i>HIST1H1E</i>	<b>0.10</b>	<b>82</b>						
<i>HIST1H2BG</i>	<b>0.10</b>	<b>81</b>						
<i>MGAT5B</i>	<b>0.10</b>	<b>80</b>						
<i>SULF2</i>	<b>0.10</b>	<b>79</b>						
<i>HMBS</i>	<b>0.10</b>	<b>78</b>						
<i>MAK</i>	<b>0.10</b>	<b>77</b>						
<i>AR exons 4-8</i>	<b>0.10</b>	<b>76</b>						
<i>SMAP1 exons 7-8</i>	<b>0.10</b>	<b>75</b>						
<i>CDC37L1</i>	<b>0.09</b>	<b>74</b>						
<i>RPLP2</i>	<b>0.09</b>	<b>73</b>						
<i>AMACR</i>	<b>0.09</b>	<b>72</b>						
<i>NEAT1</i>	<b>0.09</b>	<b>71</b>						
<i>STEAP4</i>	<b>0.09</b>	<b>70</b>						
<i>MED4</i>	<b>0.09</b>	<b>69</b>						
<i>AURKA</i>	<b>0.09</b>	<b>68</b>						
<i>NKAIN1</i>	<b>0.08</b>	<b>67</b>						
<i>GOLM1</i>	<b>0.08</b>	<b>66</b>						
<i>CD10</i>	<b>0.08</b>	<b>65</b>						

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 51)</i>			<i>Transcripts identified by Mann Whitney U (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ZNF577</i>	<b>0.08</b>	<b>64</b>						
<i>GAPDH</i>	<b>0.08</b>	<b>63</b>						
<i>KLK3 exons 1-2</i>	<b>0.08</b>	<b>62</b>						
<i>IMPDH2</i>	<b>0.08</b>	<b>61</b>						
<i>MXI1</i>	<b>0.08</b>	<b>60</b>						
<i>RPL23AP53</i>	<b>0.08</b>	<b>59</b>						
<i>FDPS</i>	<b>0.08</b>	<b>58</b>						
<i>ALAS1</i>	<b>0.08</b>	<b>57</b>						
<i>PPP1R12B</i>	<b>0.08</b>	<b>56</b>						
<i>PCSK6</i>	<b>0.08</b>	<b>55</b>						
<i>NLRP3</i>	<b>0.08</b>	<b>54</b>						
<i>MCM7</i>	<b>0.07</b>	<b>53</b>						
<i>DPP4</i>	<b>0.07</b>	<b>52</b>						
<i>ARHGEF25</i>	<b>0.07</b>	<b>51</b>						
<i>SRSF3</i>	<b>0.07</b>	<b>50</b>						
<i>STOM</i>	<b>0.07</b>	<b>49</b>						
<i>PTPRC</i>	<b>0.07</b>	<b>48</b>						
<i>VPSI3A</i>	<b>0.07</b>	<b>47</b>						
<i>CACNA1D</i>	<b>0.07</b>	<b>46</b>						
<i>ANPEP</i>	<b>0.07</b>	<b>45</b>						
<i>MIC1</i>	<b>0.07</b>	<b>44</b>						
<i>CAMK2N2</i>	<b>0.06</b>	<b>43</b>						
<i>AATF</i>	<b>0.06</b>	<b>42</b>						
<i>KLK4</i>	<b>0.06</b>	<b>41</b>						
<i>HIST1H1C</i>	<b>0.06</b>	<b>40</b>						
<i>TRPM4</i>	<b>0.06</b>	<b>39</b>						
<i>KLK3 exons 2-3</i>	<b>0.06</b>	<b>38</b>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 51)</i>			<i>Transcripts identified by Mann Whitney U (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>PVT1</i>	<i>0.06</i>	<i>37</i>						
<i>BTG2</i>	<i>0.06</i>	<i>36</i>						
<i>TERT</i>	<i>0.06</i>	<i>35</i>						
<i>SIRT1</i>	<i>0.06</i>	<i>34</i>						
<i>HPRT</i>	<i>0.06</i>	<i>33</i>						
<i>MIATNB</i>	<i>0.05</i>	<i>32</i>						
<i>KLK2</i>	<i>0.05</i>	<i>31</i>						
<i>MEMO1</i>	<i>0.05</i>	<i>30</i>						
<i>RPL18A</i>	<i>0.05</i>	<i>29</i>						
<i>COL10A1</i>	<i>0.05</i>	<i>28</i>						
<i>RP11_97012.7</i>	<i>0.05</i>	<i>27</i>						
<i>GABARAPL2</i>	<i>0.05</i>	<i>26</i>						
<i>LBH</i>	<i>0.04</i>	<i>25</i>						
<i>MKi67</i>	<i>0.04</i>	<i>24</i>						
<i>EIF2D</i>	<i>0.04</i>	<i>23</i>						
<i>SULT1A1</i>	<i>0.04</i>	<i>22</i>						
<i>HOXC4</i>	<i>0.04</i>	<i>21</i>						
<i>CDC20</i>	<i>0.04</i>	<i>20</i>						
<i>HIST3H2A</i>	<i>0.04</i>	<i>19</i>						
<i>CDKN3</i>	<i>0.04</i>	<i>18</i>						
<i>CASKIN1</i>	<i>0.03</i>	<i>17</i>						
<i>MARCH5</i>	<i>0.03</i>	<i>16</i>						
<i>BRAF</i>	<i>0.03</i>	<i>15</i>						
<i>HIST1H2BF</i>	<i>0.03</i>	<i>14</i>						
<i>PSTPIP1</i>	<i>0.03</i>	<i>13</i>						
<i>ITPR1</i>	<i>0.03</i>	<i>12</i>						
<i>TFDPI</i>	<i>0.03</i>	<i>11</i>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 51)</i>			<i>Transcripts identified by Mann Whitney U (n = 65)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>TERF2IP</i>	<b>0.03</b>	<b>10</b>						
<i>TBP</i>	<b>0.03</b>	<b>9</b>						
<i>MIR4435 1HG</i>	<b>0.03</b>	<b>8</b>						
<i>SYNM</i>	<b>0.03</b>	<b>7</b>						
<i>SACMIL</i>	<b>0.03</b>	<b>6</b>						
<i>SChLAPI</i>	<b>0.03</b>	<b>5</b>						
<i>SNCA</i>	<b>0.02</b>	<b>4</b>						
<i>SMIMI</i>	<b>0.02</b>	<b>3</b>						
<i>PDLIM5</i>	<b>0.02</b>	<b>2</b>						
<i>B2M</i>	<b>0.01</b>	<b>1</b>						

9: APPENDICES

**Supplementary Table 41 Random Forest results for comparing high-risk cancer samples with clinically benign samples in *KLK2* factorised cell data.**

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 12)</i>			<i>Transcripts identified by Mann Whitney U (n = 25)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>HOXC6</i>	<b>1.66</b>	<b>166</b>	<i>CLIC2</i>	<b>0.80</b>	<b>12</b>	<i>HOXC6</i>	<b>1.73</b>	<b>25</b>
<i>TDRD</i>	<b>0.33</b>	<b>165</b>	<i>TDRD</i>	<b>0.70</b>	<b>11</b>	<i>FOLH1</i>	<b>0.48</b>	<b>24</b>
<i>SLC43A1</i>	<b>0.30</b>	<b>164</b>	<i>SLC43A1</i>	<b>0.68</b>	<b>10</b>	<i>CADPS</i>	<b>0.39</b>	<b>23</b>
<i>FOLH1</i>	<b>0.30</b>	<b>163</b>	<i>SIM2 long</i>	<b>0.65</b>	<b>9</b>	<i>TDRD</i>	<b>0.36</b>	<b>22</b>
<i>CADPS</i>	<b>0.27</b>	<b>162</b>	<i>ERG3' exons 6-7</i>	<b>0.58</b>	<b>8</b>	<i>SIM2 long</i>	<b>0.35</b>	<b>21</b>
<i>SIM2 long</i>	<b>0.26</b>	<b>161</b>	<i>PCA3</i>	<b>0.56</b>	<b>7</b>	<i>CLIC2</i>	<b>0.33</b>	<b>20</b>
<i>CLIC2</i>	<b>0.24</b>	<b>160</b>	<i>B4GALNT4</i>	<b>0.45</b>	<b>6</b>	<i>SLC43A1</i>	<b>0.30</b>	<b>19</b>
<i>ERG3' exons 6-7</i>	<b>0.21</b>	<b>159</b>	<i>HPN</i>	<b>0.40</b>	<b>5</b>	<i>ERG3' exons 6-7</i>	<b>0.25</b>	<b>18</b>
			<b><i>TMPRSS2:ERG</i></b>					
<i>HPN</i>	<b>0.19</b>	<b>158</b>	<i>fusion</i>	<b>0.34</b>	<b>4</b>	<i>PCA3</i>	<b>0.24</b>	<b>17</b>
<i>PCA3</i>	<b>0.14</b>	<b>157</b>	<i>GCNT1</i>	<b>0.32</b>	<b>3</b>	<i>APOC1</i>	<b>0.22</b>	<b>16</b>
<i>B4GALNT4</i>	<b>0.13</b>	<b>156</b>	<i>LASS1</i>	<b>0.23</b>	<b>2</b>	<i>SLC12A1</i>	<b>0.22</b>	<b>15</b>
<i>APOC1</i>	<b>0.09</b>	<b>155</b>	<i>ERG3' exons 4-5</i>	<b>0.12</b>	<b>1</b>	<i>B4GALNT4</i>	<b>0.20</b>	<b>14</b>
<i>NAALADL2</i>	<b>0.09</b>	<b>154</b>				<i>GCNT1</i>	<b>0.11</b>	<b>13</b>
						<b><i>TMPRSS2:ERG</i></b>		
<i>SLC12A1</i>	<b>0.08</b>	<b>153</b>				<i>fusion</i>	<b>0.10</b>	<b>12</b>
<i>LASS1</i>	<b>0.07</b>	<b>152</b>				<i>TMEM86A</i>	<b>0.10</b>	<b>11</b>
<i>TMEM86A</i>	<b>0.07</b>	<b>151</b>				<i>SEC61A1</i>	<b>0.10</b>	<b>10</b>
<i>GCNT1</i>	<b>0.07</b>	<b>150</b>				<i>HPN</i>	<b>0.10</b>	<b>9</b>
<i>ISX</i>	<b>0.05</b>	<b>149</b>				<i>CCDC88B</i>	<b>0.09</b>	<b>8</b>
<i>HIST1H2BG</i>	<b>0.05</b>	<b>148</b>				<i>ISX</i>	<b>0.09</b>	<b>7</b>
<i>DLX1</i>	<b>0.05</b>	<b>147</b>				<i>MCTP1</i>	<b>0.08</b>	<b>6</b>
<i>Timp4</i>	<b>0.04</b>	<b>146</b>				<i>ERG5'</i>	<b>0.07</b>	<b>5</b>
<i>CAMKK2</i>	<b>0.04</b>	<b>145</b>				<i>LASS1</i>	<b>0.07</b>	<b>4</b>
<i>TMPRSS2:ERG fusion</i>	<b>0.04</b>	<b>144</b>				<i>ERG3' exons 4-5</i>	<b>0.06</b>	<b>3</b>



9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 12)</i>			<i>Transcripts identified by Mann Whitney U (n = 25)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>GJB1</i>	<b>0.04</b>	<b>143</b>				<b>SFRP4</b>	<b>0.03</b>	<b>2</b>
<i>TMEM45B</i>	<b>0.04</b>	<b>142</b>				<b>CAMKK2</b>	<b>0.03</b>	<b>1</b>
<i>HIST1H1C</i>	<b>0.04</b>	<b>141</b>						
<i>SMIMI</i>	<b>0.04</b>	<b>140</b>						
<i>MMP25</i>	<b>0.03</b>	<b>139</b>						
<i>VAX2</i>	<b>0.03</b>	<b>138</b>						
<i>UPK2</i>	<b>0.03</b>	<b>137</b>						
<i>SULT1A1</i>	<b>0.03</b>	<b>136</b>						
<i>ABCB9</i>	<b>0.03</b>	<b>135</b>						
<i>SEC61A1</i>	<b>0.03</b>	<b>134</b>						
<i>RNF157</i>	<b>0.03</b>	<b>133</b>						
<i>CKAP2L</i>	<b>0.03</b>	<b>132</b>						
<i>AR exons 4-8</i>	<b>0.03</b>	<b>131</b>						
<i>AURKA</i>	<b>0.03</b>	<b>130</b>						
<i>IGFBP3</i>	<b>0.02</b>	<b>129</b>						
<i>P712P</i>	<b>0.02</b>	<b>128</b>						
<i>SIM2 short</i>	<b>0.02</b>	<b>127</b>						
<i>SFRP4</i>	<b>0.02</b>	<b>126</b>						
<i>GOLM1</i>	<b>0.02</b>	<b>125</b>						
<i>SPINK1</i>	<b>0.02</b>	<b>124</b>						
<i>ERG3' exons 4-5</i>	<b>0.02</b>	<b>123</b>						
<i>CD10</i>	<b>0.02</b>	<b>122</b>						
<i>ERG5'</i>	<b>0.02</b>	<b>121</b>						
<i>MGAT5B</i>	<b>0.02</b>	<b>120</b>						
<i>STEAP2</i>	<b>0.02</b>	<b>119</b>						
<i>ANKRD34B</i>	<b>0.02</b>	<b>118</b>						
<i>CP</i>	<b>0.02</b>	<b>117</b>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 12)</i>			<i>Transcripts identified by Mann Whitney U (n = 25)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>SLC4A1.S</i>	<b>0.02</b>	<b>116</b>						
<i>MNX1</i>	<b>0.02</b>	<b>115</b>						
<i>ST6GALNAC1</i>	<b>0.02</b>	<b>114</b>						
<i>LBH</i>	<b>0.02</b>	<b>113</b>						
<i>COL9A2</i>	<b>0.02</b>	<b>112</b>						
<i>NKAIN1</i>	<b>0.02</b>	<b>111</b>						
<i>SRSF3</i>	<b>0.02</b>	<b>110</b>						
<i>SERPINB5</i>	<b>0.02</b>	<b>109</b>						
<i>KLK3 exons 2-3</i>	<b>0.02</b>	<b>108</b>						
<i>PPP1R12B</i>	<b>0.01</b>	<b>107</b>						
<i>ACTR5</i>	<b>0.01</b>	<b>106</b>						
<i>SPON2</i>	<b>0.01</b>	<b>105</b>						
<i>SULF2</i>	<b>0.01</b>	<b>104</b>						
<i>RPL23AP53</i>	<b>0.01</b>	<b>103</b>						
<i>CAMK2N2</i>	<b>0.01</b>	<b>102</b>						
<i>CDC37L1</i>	<b>0.01</b>	<b>101</b>						
<i>HIST1H2BF</i>	<b>0.01</b>	<b>100</b>						
<i>MIR146A</i>	<b>0.01</b>	<b>99</b>						
<i>TERT</i>	<b>0.01</b>	<b>98</b>						
<i>SACMIL</i>	<b>0.01</b>	<b>97</b>						
<i>ALAS1</i>	<b>0.01</b>	<b>96</b>						
<i>OR52A2</i>	<b>0.01</b>	<b>95</b>						
<i>HIST3H2A</i>	<b>0.01</b>	<b>94</b>						
<i>RPS11</i>	<b>0.01</b>	<b>93</b>						
<i>KLK3 exons 1-2</i>	<b>0.01</b>	<b>92</b>						
<i>NLRP3</i>	<b>0.01</b>	<b>91</b>						
<i>TMEM47</i>	<b>0.01</b>	<b>90</b>						

## 9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 12)</i>			<i>Transcripts identified by Mann Whitney U (n = 25)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>MEX3A</i>	<i>0.01</i>	<i>89</i>						
<i>MKi67</i>	<i>0.01</i>	<i>88</i>						
<i>RIOK3</i>	<i>0.01</i>	<i>87</i>						
<i>PSTPIP1</i>	<i>0.01</i>	<i>86</i>						
<i>BRAF</i>	<i>0.01</i>	<i>85</i>						
<i>SSPO</i>	<i>0.01</i>	<i>84</i>						
<i>MDK</i>	<i>0.01</i>	<i>83</i>						
<i>ITGBL1</i>	<i>0.01</i>	<i>82</i>						
<i>AMACR</i>	<i>0.01</i>	<i>81</i>						
<i>VPS13A</i>	<i>0.01</i>	<i>80</i>						
<i>RAB17</i>	<i>0.01</i>	<i>79</i>						
<i>MIC1</i>	<i>0.01</i>	<i>78</i>						
<i>PPAP2A</i>	<i>0.01</i>	<i>77</i>						
<i>KLK4</i>	<i>0.01</i>	<i>76</i>						
<i>SNORA20</i>	<i>0.01</i>	<i>75</i>						
<i>PECI</i>	<i>0.01</i>	<i>74</i>						
<i>PTN</i>	<i>0.01</i>	<i>73</i>						
<i>RPS10</i>	<i>0.01</i>	<i>72</i>						
<i>MFSD2A</i>	<i>0.01</i>	<i>71</i>						
<i>CACNA1D</i>	<i>0.01</i>	<i>70</i>						
<i>PALM3</i>	<i>0.01</i>	<i>69</i>						
<i>MCTP1</i>	<i>0.01</i>	<i>68</i>						
<i>CCDC88B</i>	<i>0.01</i>	<i>67</i>						
<i>AMH</i>	<i>0.01</i>	<i>66</i>						
<i>STOM</i>	<i>0.01</i>	<i>65</i>						
<i>AGR2</i>	<i>0.01</i>	<i>64</i>						
<i>DNAH5</i>	<i>0.01</i>	<i>63</i>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 12)</i>			<i>Transcripts identified by Mann Whitney U (n = 25)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>HOXC4</i>	<i>0.01</i>	<i>62</i>						
<i>TWIST1</i>	<i>0.01</i>	<i>61</i>						
<i>PDLIM5</i>	<i>0.01</i>	<i>60</i>						
<i>AATF</i>	<i>0.01</i>	<i>59</i>						
<i>PVT1</i>	<i>0.004</i>	<i>58</i>						
<i>B2M</i>	<i>0.004</i>	<i>57</i>						
<i>HPRT</i>	<i>0.004</i>	<i>56</i>						
<i>DPP4</i>	<i>0.004</i>	<i>55</i>						
<i>RPLP2</i>	<i>0.004</i>	<i>54</i>						
<i>MEMO1</i>	<i>0.004</i>	<i>53</i>						
<i>MSMB</i>	<i>0.004</i>	<i>52</i>						
<i>PPFIA2</i>	<i>0.004</i>	<i>51</i>						
<i>COL10A1</i>	<i>0.004</i>	<i>50</i>						
<i>ZNF577</i>	<i>0.004</i>	<i>49</i>						
<i>TRPM4</i>	<i>0.004</i>	<i>48</i>						
<i>MIATNB</i>	<i>0.004</i>	<i>47</i>						
<i>SChLAP1</i>	<i>0.004</i>	<i>46</i>						
<i>GAPDH</i>	<i>0.004</i>	<i>44.5</i>						
<i>RPL18A</i>	<i>0.004</i>	<i>44.5</i>						
<i>TMCC2</i>	<i>0.003</i>	<i>43</i>						
<i>MCM7</i>	<i>0.003</i>	<i>42</i>						
<i>NEAT1</i>	<i>0.003</i>	<i>41</i>						
<i>HIST1H1E</i>	<i>0.003</i>	<i>40</i>						
<i>CLU</i>	<i>0.002</i>	<i>39</i>						
<i>MYOF</i>	<i>0.002</i>	<i>38</i>						
<i>BTG2</i>	<i>0.002</i>	<i>37</i>						
<i>ITPR1</i>	<i>0.002</i>	<i>36</i>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 12)</i>			<i>Transcripts identified by Mann Whitney U (n = 25)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>CDC20</i>	<b>0.002</b>	<b>35</b>						
<i>STEAP4</i>	<b>0.002</b>	<b>34</b>						
<i>Met</i>	<b>0.002</b>	<b>33</b>						
<i>EN2</i>	<b>0.002</b>	<b>32</b>						
<i>SMAP1 exons 7-8</i>	<b>0.002</b>	<b>31</b>						
<i>SSTR1</i>	<b>0.002</b>	<b>30</b>						
<i>MAPK8IP2</i>	<b>0.002</b>	<b>29</b>						
<i>MAK</i>	<b>0.002</b>	<b>28</b>						
<i>GABARAPL2</i>	<b>0.002</b>	<b>27</b>						
<i>CASKIN1</i>	<b>0.002</b>	<b>26</b>						
<i>MED4</i>	<b>0.002</b>	<b>25</b>						
<i>IFT57</i>	<b>0.002</b>	<b>24</b>						
<i>AR.ex9</i>	<b>0.002</b>	<b>23</b>						
<i>TFDPI</i>	<b>2.17604x10<sup>-17</sup></b>	<b>22</b>						
<i>RP11_97012.7</i>	<b>2.13163x10<sup>-17</sup></b>	<b>21</b>						
<i>CDKN3</i>	<b>2.04281x10<sup>-17</sup></b>	<b>19.5</b>						
<i>HMBS</i>	<b>2.04281x10<sup>-17</sup></b>	<b>19.5</b>						
<i>SNCA</i>	<b>1.77636x10<sup>-17</sup></b>	<b>18</b>						
<i>ARHGEF25</i>	<b>1.73195x10<sup>-17</sup></b>	<b>17</b>						
<i>OGT</i>	<b>1.59872x10<sup>-17</sup></b>	<b>16</b>						
<i>MXI1</i>	<b>1.55431x10<sup>-17</sup></b>	<b>15</b>						
<i>MARCH5</i>	<b>1.46549x10<sup>-17</sup></b>	<b>14</b>						
<i>MMP11</i>	<b>1.42109x10<sup>-17</sup></b>	<b>13</b>						
<i>TERF2IP</i>	<b>1.33227x10<sup>-17</sup></b>	<b>12</b>						
<i>SYNM</i>	<b>1.24345x10<sup>-17</sup></b>	<b>11</b>						
<i>ANPEP</i>	<b>1.19904x10<sup>-17</sup></b>	<b>9.5</b>						
<i>IMPDH2</i>	<b>1.19904x10<sup>-17</sup></b>	<b>9.5</b>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 12)</i>			<i>Transcripts identified by Mann Whitney U (n = 25)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>MIR4435 IHG</i>	<i>1.02141x10<sup>-17</sup></i>	<i>8</i>						
<i>MMP26</i>	<i>9.32587x10<sup>-18</sup></i>	<i>6</i>						
<i>PCSK6</i>	<i>9.32587x10<sup>-18</sup></i>	<i>6</i>						
<i>TBP</i>	<i>9.32587x10<sup>-18</sup></i>	<i>6</i>						
<i>FDPS</i>	<i>8.43769x10<sup>-18</sup></i>	<i>4</i>						
<i>EIF2D</i>	<i>7.10543x10<sup>-18</sup></i>	<i>3</i>						
<i>SIRT1</i>	<i>3.9968x10<sup>-18</sup></i>	<i>2</i>						
<i>PTPRC</i>	<i>3.55271x10<sup>-18</sup></i>	<i>1</i>						

Supplementary Table 42 Random Forest results when comparing clinically benign samples to high risk cancer samples using the *RPLP2* and *TWIST1* normalised data.

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 20)</i>			<i>Transcripts identified by polr (n = 35)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>CADPS</i>	<i>0.59</i>	<i>167</i>	<i>HOXC6</i>	<i>0.87</i>	<i>20</i>	<i>SPINK1</i>	<i>0.69</i>	<i>35</i>
<i>SPINK1</i>	<i>0.58</i>	<i>166</i>	<i>CADPS</i>	<i>0.77</i>	<i>19</i>	<i>HOXC6</i>	<i>0.68</i>	<i>34</i>
<i>HOXC6</i>	<i>0.50</i>	<i>165</i>	<i>SPINK1</i>	<i>0.77</i>	<i>18</i>	<i>CADPS</i>	<i>0.63</i>	<i>33</i>
<i>ST6GALNAC1</i>	<i>0.36</i>	<i>164</i>	<i>ERG3' exons 6-7</i>	<i>0.47</i>	<i>17</i>	<i>ST6GALNAC1</i>	<i>0.27</i>	<i>32</i>
<i>VAX2</i>	<i>0.27</i>	<i>163</i>	<i>ST6GALNAC1</i>	<i>0.39</i>	<i>16</i>	<i>VAX2</i>	<i>0.26</i>	<i>31</i>
<i>ERG3' exons 6-7</i>	<i>0.19</i>	<i>162</i>	<i>NAALADL2</i>	<i>0.34</i>	<i>15</i>	<i>ERG3' exons 6-7</i>	<i>0.25</i>	<i>30</i>
<i>NAALADL2</i>	<i>0.18</i>	<i>161</i>	<i>SLC43A1</i>	<i>0.25</i>	<i>14</i>	<i>B4GALNT4</i>	<i>0.23</i>	<i>29</i>
<i>SLC43A1</i>	<i>0.15</i>	<i>160</i>	<i>CLIC2</i>	<i>0.24</i>	<i>13</i>	<i>NAALADL2</i>	<i>0.22</i>	<i>28</i>
<i>HPN</i>	<i>0.13</i>	<i>159</i>	<i>TDRD</i>	<i>0.23</i>	<i>12</i>	<i>PPAP2A</i>	<i>0.20</i>	<i>27</i>
<i>PPAP2A</i>	<i>0.12</i>	<i>158</i>	<i>UPK2</i>	<i>0.22</i>	<i>11</i>	<i>SLC43A1</i>	<i>0.19</i>	<i>26</i>
<i>UPK2</i>	<i>0.12</i>	<i>157</i>	<i>PTN</i>	<i>0.21</i>	<i>10</i>	<i>HPN</i>	<i>0.16</i>	<i>25</i>
<i>TDRD</i>	<i>0.12</i>	<i>156</i>	<i>SERPINB5</i>	<i>0.20</i>	<i>9</i>	<i>UPK2</i>	<i>0.15</i>	<i>24</i>
<i>PTN</i>	<i>0.11</i>	<i>155</i>	<i>PALM3</i>	<i>0.17</i>	<i>8</i>	<i>CLIC2</i>	<i>0.14</i>	<i>23</i>

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 20)</i>			<i>Transcripts identified by polr (n = 35)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>TMPRSS2:ERG fusion</i>	<b>0.10</b>	<b>154</b>	<i>GJB1</i>	<b>0.16</b>	<b>7</b>	<i>APOC1</i>	<b>0.14</b>	<b>22</b>
<i>IFT57</i>	<b>0.10</b>	<b>153</b>	<i>HPN</i>	<b>0.16</b>	<b>6</b>	<i>TDRD</i>	<b>0.13</b>	<b>21</b>
<i>SERPINB5</i>	<b>0.09</b>	<b>152</b>	<i>TMPRSS2:ERG fusion</i>	<b>0.15</b>	<b>5</b>	<i>ISX</i>	<b>0.13</b>	<b>20</b>
<i>SFRP4</i>	<b>0.08</b>	<b>151</b>	<i>CCDC88B</i>	<b>0.13</b>	<b>4</b>	<i>IFT57</i>	<b>0.12</b>	<b>19</b>
<i>GJB1</i>	<b>0.08</b>	<b>150</b>	<i>SIM2 long</i>	<b>0.13</b>	<b>3</b>	<i>TMEM86A</i>	<b>0.12</b>	<b>18</b>
<i>CLIC2</i>	<b>0.08</b>	<b>149</b>	<i>LASS1</i>	<b>0.11</b>	<b>2</b>	<i>MSMB</i>	<b>0.12</b>	<b>17</b>
<i>MAPK8IP2</i>	<b>0.08</b>	<b>148</b>	<i>GCNT1</i>	<b>0.09</b>	<b>1</b>	<i>SERPINB5</i>	<b>0.11</b>	<b>16</b>
<i>B4GALNT4</i>	<b>0.07</b>	<b>147</b>				<i>TMPRSS2:ERG fusion</i>	<b>0.11</b>	<b>15</b>
<i>ERG5'</i>	<b>0.06</b>	<b>146</b>				<i>MMP25</i>	<b>0.11</b>	<b>14</b>
<i>COL9A2</i>	<b>0.06</b>	<b>145</b>				<i>SPON2</i>	<b>0.09</b>	<b>13</b>
<i>APOC1</i>	<b>0.06</b>	<b>144</b>				<i>STEAP2</i>	<b>0.08</b>	<b>12</b>
<i>SIM2 long</i>	<b>0.06</b>	<b>143</b>				<i>PALM3</i>	<b>0.08</b>	<b>11</b>
<i>PECI</i>	<b>0.06</b>	<b>142</b>				<i>GCNT1</i>	<b>0.08</b>	<b>10</b>
<i>ISX</i>	<b>0.06</b>	<b>141</b>				<i>RAB17</i>	<b>0.08</b>	<b>9</b>
<i>MSMB</i>	<b>0.06</b>	<b>140</b>				<i>PTN</i>	<b>0.08</b>	<b>8</b>
<i>CP</i>	<b>0.06</b>	<b>139</b>				<i>SIM2 long</i>	<b>0.07</b>	<b>7</b>
<i>TMEM86A</i>	<b>0.05</b>	<b>138</b>				<i>RPS10</i>	<b>0.07</b>	<b>6</b>
<i>MNX1</i>	<b>0.05</b>	<b>137</b>				<i>GJB1</i>	<b>0.07</b>	<b>5</b>
<i>PALM3</i>	<b>0.05</b>	<b>136</b>				<i>CCDC88B</i>	<b>0.07</b>	<b>4</b>
<i>IGFBP3</i>	<b>0.04</b>	<b>135</b>				<i>MCTP1</i>	<b>0.05</b>	<b>3</b>
<i>ANKRD34B</i>	<b>0.04</b>	<b>134</b>				<i>LASS1</i>	<b>0.05</b>	<b>2</b>
<i>LASS1</i>	<b>0.04</b>	<b>133</b>				<i>IGFBP3</i>	<b>0.05</b>	<b>1</b>
<i>CCDC88B</i>	<b>0.04</b>	<b>132</b>						
<i>TMCC2</i>	<b>0.04</b>	<b>131</b>						
<i>GCNT1</i>	<b>0.04</b>	<b>130</b>						
<i>FOLH1</i>	<b>0.04</b>	<b>129</b>						

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 20)</i>			<i>Transcripts identified by polr (n = 35)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>DLX1</i>	<i>0.03</i>	<i>128</i>						
<i>MMP25</i>	<i>0.03</i>	<i>127</i>						
<i>RAB17</i>	<i>0.03</i>	<i>126</i>						
<i>RPL18A</i>	<i>0.03</i>	<i>125</i>						
<i>MDK</i>	<i>0.03</i>	<i>124</i>						
<i>RPS10</i>	<i>0.03</i>	<i>123</i>						
<i>EN2</i>	<i>0.03</i>	<i>122</i>						
<i>RIOK3</i>	<i>0.03</i>	<i>121</i>						
<i>MFSD2A</i>	<i>0.03</i>	<i>120</i>						
<i>KLK3 exons 2-3</i>	<i>0.03</i>	<i>119</i>						
<i>CKAP2L</i>	<i>0.03</i>	<i>118</i>						
<i>PCA3</i>	<i>0.03</i>	<i>117</i>						
<i>PPFIA2</i>	<i>0.02</i>	<i>116</i>						
<i>MCTP1</i>	<i>0.02</i>	<i>115</i>						
<i>MYOF</i>	<i>0.02</i>	<i>114</i>						
<i>RNF157</i>	<i>0.02</i>	<i>113</i>						
<i>CDC37L1</i>	<i>0.02</i>	<i>112</i>						
<i>AMACR</i>	<i>0.02</i>	<i>111</i>						
<i>Timp4</i>	<i>0.02</i>	<i>110</i>						
<i>CDC20</i>	<i>0.02</i>	<i>109</i>						
<i>SEC61A1</i>	<i>0.02</i>	<i>108</i>						
<i>STEAP2</i>	<i>0.02</i>	<i>107</i>						
<i>SRSF3</i>	<i>0.02</i>	<i>106</i>						
<i>STOM</i>	<i>0.02</i>	<i>105</i>						
<i>SPON2</i>	<i>0.02</i>	<i>104</i>						
<i>MKi67</i>	<i>0.02</i>	<i>103</i>						
<i>SMIMI</i>	<i>0.02</i>	<i>102</i>						
<i>ITGBL1</i>	<i>0.02</i>	<i>101</i>						



9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 20)</i>			<i>Transcripts identified by polr (n = 35)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>HOXC4</i>	<b>0.02</b>	<b>100</b>						
<i>PPP1R12B</i>	<b>0.01</b>	<b>99</b>						
<i>KLK4</i>	<b>0.01</b>	<b>98</b>						
<i>ACTR5</i>	<b>0.01</b>	<b>97</b>						
<i>CLU</i>	<b>0.01</b>	<b>96</b>						
<i>AR exon 9</i>	<b>0.01</b>	<b>95</b>						
<i>RP11_97012.7</i>	<b>0.01</b>	<b>94</b>						
<i>CAMKK2</i>	<b>0.01</b>	<b>93</b>						
<i>TRPM4</i>	<b>0.01</b>	<b>92</b>						
<i>MIR146A</i>	<b>0.01</b>	<b>91</b>						
<i>SIRT1</i>	<b>0.01</b>	<b>90</b>						
<i>GOLM1</i>	<b>0.01</b>	<b>89</b>						
<i>SLC4A1.S</i>	<b>0.01</b>	<b>88</b>						
<i>ZNF577</i>	<b>0.01</b>	<b>87</b>						
<i>RPS11</i>	<b>0.01</b>	<b>86</b>						
<i>PTPRC</i>	<b>0.01</b>	<b>85</b>						
<i>NLRP3</i>	<b>0.01</b>	<b>84</b>						
<i>TMEM47</i>	<b>0.01</b>	<b>83</b>						
<i>CACNA1D</i>	<b>0.01</b>	<b>82</b>						
<i>HMBS</i>	<b>0.01</b>	<b>81</b>						
<i>ABCB9</i>	<b>0.01</b>	<b>80</b>						
<i>PVT1</i>	<b>0.01</b>	<b>79</b>						
<i>SSPO</i>	<b>0.01</b>	<b>78</b>						
<i>ITPR1</i>	<b>0.01</b>	<b>77</b>						
<i>KLK3 exons 1-2</i>	<b>0.01</b>	<b>76</b>						
<i>STEAP4</i>	<b>0.01</b>	<b>75</b>						
<i>PCSK6</i>	<b>0.01</b>	<b>74</b>						
<i>AURKA</i>	<b>0.01</b>	<b>73</b>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 20)</i>			<i>Transcripts identified by polr (n = 35)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>TBP</i>	<i>0.01</i>	<i>72</i>						
<i>SChLAP1</i>	<i>0.01</i>	<i>71</i>						
<i>VPS13A</i>	<i>0.01</i>	<i>70</i>						
<i>NKAIN1</i>	<i>0.01</i>	<i>69</i>						
<i>MIATNB</i>	<i>0.01</i>	<i>68</i>						
<i>FDPS</i>	<i>0.01</i>	<i>67</i>						
<i>OR52A2</i>	<i>0.01</i>	<i>66</i>						
<i>RPL23AP53</i>	<i>0.01</i>	<i>65</i>						
<i>HIST1H2BF</i>	<i>0.01</i>	<i>64</i>						
<i>CAMK2N2</i>	<i>0.01</i>	<i>63</i>						
<i>DPP4</i>	<i>0.01</i>	<i>62</i>						
<i>SMAP1 exons 7-8</i>	<i>0.01</i>	<i>61</i>						
<i>HIST1H1C</i>	<i>0.01</i>	<i>60</i>						
<i>ALAS1</i>	<i>0.01</i>	<i>59</i>						
<i>TMEM45B</i>	<i>0.005</i>	<i>58</i>						
<i>TWIST1</i>	<i>0.005</i>	<i>57</i>						
<i>HIST1H1E</i>	<i>0.005</i>	<i>56</i>						
<i>MMP26</i>	<i>0.004</i>	<i>55</i>						
<i>SNCA</i>	<i>0.004</i>	<i>54</i>						
<i>BRAF</i>	<i>0.004</i>	<i>53</i>						
<i>GABARAPL2</i>	<i>0.004</i>	<i>52</i>						
<i>RPLP2</i>	<i>0.004</i>	<i>51</i>						
<i>MIR4435 1HG</i>	<i>0.004</i>	<i>50</i>						
<i>ERG3' exons 4-5</i>	<i>0.004</i>	<i>49</i>						
<i>AMH</i>	<i>0.004</i>	<i>48</i>						
<i>ANPEP</i>	<i>0.004</i>	<i>47</i>						
<i>SACMIL</i>	<i>0.004</i>	<i>46</i>						
<i>AGR2</i>	<i>0.004</i>	<i>45</i>						

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 20)</i>			<i>Transcripts identified by polr (n = 35)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>MEX3A</i>	<i>0.004</i>	<i>44</i>						
<i>MXII</i>	<i>0.004</i>	<i>43</i>						
<i>MARCH5</i>	<i>0.004</i>	<i>42</i>						
<i>CD10</i>	<i>0.003</i>	<i>41</i>						
<i>EIF2D</i>	<i>0.003</i>	<i>40</i>						
<i>ARHGEF25</i>	<i>0.003</i>	<i>39</i>						
<i>NEAT1</i>	<i>0.003</i>	<i>38</i>						
<i>IMPDH2</i>	<i>0.003</i>	<i>37</i>						
<i>Met</i>	<i>0.002</i>	<i>36</i>						
<i>PSTPIP1</i>	<i>0.002</i>	<i>35</i>						
<i>P712P</i>	<i>0.002</i>	<i>34</i>						
<i>DNAH5</i>	<i>0.002</i>	<i>33</i>						
<i>MAK</i>	<i>0.002</i>	<i>32</i>						
<i>SIM2 short</i>	<i>0.002</i>	<i>31</i>						
<i>SYNM</i>	<i>0.002</i>	<i>30</i>						
<i>MCM7</i>	<i>0.002</i>	<i>29</i>						
<i>TERT</i>	<i>0.002</i>	<i>28</i>						
<i>AR exons 4-8</i>	<i>0.002</i>	<i>27</i>						
<i>PDLIM5</i>	<i>0.002</i>	<i>26</i>						
<i>B2M</i>	<i>0.002</i>	<i>25</i>						
<i>COL10A1</i>	<i>0.002</i>	<i>24</i>						
<i>LBH</i>	<i>2.53x10<sup>-17</sup></i>	<i>23</i>						
<i>SULF2</i>	<i>2.18x10<sup>-17</sup></i>	<i>22</i>						
<i>MMP11</i>	<i>2.13x10<sup>-17</sup></i>	<i>20.5</i>						
<i>SULT1A1</i>	<i>2.13x10<sup>-17</sup></i>	<i>20.5</i>						
<i>MGAT5B</i>	<i>1.95x10<sup>-17</sup></i>	<i>19</i>						
<i>CASKIN1</i>	<i>1.87x10<sup>-17</sup></i>	<i>18</i>						
<i>SLC12A1</i>	<i>1.82x10<sup>-17</sup></i>	<i>17</i>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 20)</i>			<i>Transcripts identified by polr (n = 35)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>HIST3H2A</i>	<i>1.78x10<sup>-17</sup></i>	<i>15.5</i>						
<i>MIC1</i>	<i>1.78x10<sup>-17</sup></i>	<i>15.5</i>						
<i>KLK2</i>	<i>1.69x10<sup>-17</sup></i>	<i>14</i>						
<i>BTG2</i>	<i>1.64x10<sup>-17</sup></i>	<i>13</i>						
<i>SNORA20</i>	<i>1.60x10<sup>-17</sup></i>	<i>11.5</i>						
<i>TFDP1</i>	<i>1.60x10<sup>-17</sup></i>	<i>11.5</i>						
<i>HIST1H2BG</i>	<i>1.24x10<sup>-17</sup></i>	<i>10</i>						
<i>HPRT</i>	<i>1.15x10<sup>-17</sup></i>	<i>8.5</i>						
<i>SSTR1</i>	<i>1.15x10<sup>-17</sup></i>	<i>8.5</i>						
<i>AATF</i>	<i>1.02x10<sup>-17</sup></i>	<i>7</i>						
<i>TERF2IP</i>	<i>9.77x10<sup>-18</sup></i>	<i>6</i>						
<i>CDKN3</i>	<i>9.33x10<sup>-18</sup></i>	<i>5</i>						
<i>MED4</i>	<i>8.88x10<sup>-18</sup></i>	<i>4</i>						
<i>MEMO1</i>	<i>8.44x10<sup>-18</sup></i>	<i>3</i>						
<i>GAPDH</i>	<i>7.99x10<sup>-18</sup></i>	<i>2</i>						
<i>OGT</i>	<i>5.33x10<sup>-18</sup></i>	<i>1</i>						

**6.20 CB- L-I-H Trend**

**Supplementary Table 43** Transcripts that have significant expression trend (using polr and glm) across clinically benign, low-risk, intermediate-risk and high-risk cancer samples in the baseline normalised cell NanoString data.

<i>Transcript</i>	<i>Glm p-value</i>	<i>glm Adjusted p-value</i>	<i>Polr p-value</i>	<i>Polr adjusted p-value</i>
<i>AATF</i>	0.0002	0.035	0.0023	0.3326
<i>ACTR5</i>	0.0114	0.9827	0.0204	0.9931
<i>AMH</i>	0.0387	0.9827		
<i>ANKRD34B</i>	0.0033	0.4275	0.0065	0.878
<i>ANPEP</i>	0.0007	0.1031	0.0033	0.4854
<i>APOC1</i>	2.90x10 <sup>-06</sup>	0.0005	0.0001	0.0246
<i>ARHGEF25</i>	0.026	0.9827	0.0251	0.9931
<i>AURKA</i>	0.0066	0.7921	0.0232	0.9931
<i>B2M</i>	0.0022	0.2821	0.0092	0.9931
<i>B4GALNT4</i>	7.83x10 <sup>-06</sup>	0.0013	9.95x10 <sup>-05</sup>	0.0166
<i>BTG2</i>	0.0121	0.9827	0.0398	0.9931
<i>CACNA1D</i>	0.003	0.3935	0.0184	0.9931
<i>CADPS</i>	0.0001	0.019	0.0017	0.261
<i>CAMKK2</i>	5.44x10 <sup>-05</sup>	0.0087	0.0012	0.1813
<i>CCDC88B</i>	0.0008	0.1059	0.0043	0.6079
<i>CD10</i>	0.0011	0.158	0.0047	0.6525
<i>CDC20</i>	0.018	0.9827	0.041	0.9931
<i>CDKN3</i>	0.0218	0.9827	0.0342	0.9931
<i>CKAP2L</i>	0.0063	0.7631	0.0163	0.9931
<i>CLIC2</i>	0.0012	0.1699	0.0071	0.9579
<i>COL9A2</i>	0.0273	0.9827	0.031	0.9931
<i>DPP4</i>	0.0385	0.9827		
<i>EIF2D</i>	0.0038	0.4774	0.0197	0.9931
<i>EN2</i>	0.0141	0.9827	0.0189	0.9931
<i>ERG3' exons 4-5</i>	0.0065	0.7819	0.0083	0.9931
<i>ERG3' exons 6-7</i>	5.87x10 <sup>-06</sup>	0.001	0.0001	0.0191
<i>FDPS</i>	0.0231	0.9827		
<i>FOLH1</i>	0.0031	0.3935	0.0055	0.7553
<i>GABARAPL2</i>	0.0039	0.4932	0.0317	0.9931
<i>GAPDH</i>	0.0004	0.0542	0.0022	0.3326
<i>GCNT1</i>	0.0002	0.025	0.0008	0.1318
<i>HIST1H2BF</i>	0.0265	0.9827	0.0481	0.9931
<i>HIST1H2BG</i>	0.0092	0.9827	0.026	0.9931
<i>HOXC6</i>	1.85x10 <sup>-05</sup>	0.003	0.0001	0.0191
<i>HPN</i>	6.51x10 <sup>-05</sup>	0.0102	0.0007	0.1079
<i>HPRT</i>	0.0015	0.1973	0.0107	0.9931
<i>ISX</i>	3.91x10 <sup>-05</sup>	0.0063	0.0004	0.066
<i>ITPRI</i>	0.0016	0.2177	0.0118	0.9931
<i>LASS1</i>	0.0002	0.0331	0.0019	0.2854
<i>MAK</i>	0.0482	0.9827		

## 9: APPENDICES

<i>Transcript</i>	<i>Glm p-value</i>	<i>glm Adjusted p-value</i>	<i>Polr p-value</i>	<i>Polr adjusted p-value</i>
<i>MAPK8IP2</i>	<b>0.036</b>	<b>0.9827</b>		
<i>MCTP1</i>	<b>4.10x10<sup>-06</sup></b>	<b>0.0007</b>	<b>0.0003</b>	<b>0.0481</b>
<i>MED4</i>	<b>0.0335</b>	<b>0.9827</b>		
<i>MEMO1</i>	<b>0.0368</b>	<b>0.9827</b>		
<i>MEX3A</i>	<b>0.0046</b>	<b>0.5676</b>	<b>0.0073</b>	<b>0.9736</b>
<i>MFSD2A</i>	<b>0.0003</b>	<b>0.0379</b>	<b>0.0035</b>	<b>0.4969</b>
<i>MGAT5B</i>	<b>0.0254</b>	<b>0.9827</b>	<b>0.0352</b>	<b>0.9931</b>
<i>MIC1</i>	<b>0.0484</b>	<b>0.9827</b>	<b>0.0359</b>	<b>0.9931</b>
<i>MIR146A</i>	<b>0.021</b>	<b>0.9827</b>		
<i>MIR4435 IHG</i>	<b>0.0018</b>	<b>0.2367</b>	<b>0.0131</b>	<b>0.9931</b>
<i>MMP11</i>	<b>0.0231</b>	<b>0.9827</b>	<b>0.0497</b>	<b>0.9931</b>
<i>MMP25</i>	<b>0.0002</b>	<b>0.0277</b>	<b>0.0014</b>	<b>0.2178</b>
<i>MMP26</i>	<b>0.0317</b>	<b>0.9827</b>	<b>0.032</b>	<b>0.9931</b>
<i>MXII</i>	<b>0.0006</b>	<b>0.0855</b>	<b>0.0057</b>	<b>0.7873</b>
<i>NEAT1</i>	<b>0.0002</b>	<b>0.0265</b>	<b>0.0003</b>	<b>0.047</b>
<i>NLRP3</i>	<b>0.0012</b>	<b>0.1632</b>	<b>0.0042</b>	<b>0.6027</b>
<i>PCA3</i>	<b>0.0154</b>	<b>0.9827</b>		
<i>PDLIM5</i>	<b>0.0101</b>	<b>0.9827</b>	<b>0.0468</b>	<b>0.9931</b>
<i>PSTPIP1</i>	<b>0.0149</b>	<b>0.9827</b>	<b>0.0217</b>	<b>0.9931</b>
<i>PTPRC</i>	<b>0.0003</b>	<b>0.0504</b>	<b>0.0036</b>	<b>0.5209</b>
<i>RIOK3</i>	<b>6.64x10<sup>-06</sup></b>	<b>0.0011</b>	<b>0.0003</b>	<b>0.0515</b>
<i>RPL18A</i>	<b>0.04</b>	<b>0.9827</b>		
<i>RPS11</i>	<b>0.0436</b>	<b>0.9827</b>		
<i>SACMIL</i>	<b>0.0239</b>	<b>0.9827</b>	<b>0.0396</b>	<b>0.9931</b>
<i>SEC61A1</i>	<b>0.001</b>	<b>0.1335</b>	<b>0.0074</b>	<b>0.9864</b>
<i>SFRP4</i>	<b>0.0054</b>	<b>0.6665</b>	<b>0.0235</b>	<b>0.9931</b>
<i>SIM2 long</i>	<b>0.0004</b>	<b>0.061</b>	<b>0.0022</b>	<b>0.3237</b>
<i>SIM2 short</i>	<b>0.0139</b>	<b>0.9827</b>	<b>0.0253</b>	<b>0.9931</b>
<i>SIRT1</i>	<b>0.0018</b>	<b>0.2367</b>	<b>0.0176</b>	<b>0.9931</b>
<i>SLC12A1</i>	<b>0.0317</b>	<b>0.9827</b>		

9: APPENDICES

**Supplementary Table 44** Transcripts that have significant expression trend (using polr and glm) across clinically benign, low-risk, intermediate-risk and high-risk cancer samples in the *KLK2* ratio cell NanoString data.

<i>Transcript</i>	<i>Glm p-value</i>	<i>glm Adjusted p-value</i>	<i>Polr p-value</i>	<i>Polr adjusted p-value</i>
<i>ANKRD34B</i>	<b>0.0112</b>	<b>0.998</b>	<b>0.0374</b>	<b>0.994</b>
<i>APOC1</i>	<b>0.0178</b>	<b>0.998</b>		
<i>B4GALNT4</i>	<b>0.0216</b>	<b>0.998</b>		
<i>CADPS</i>	<b>0.0052</b>	<b>0.8154</b>	<b>0.0112</b>	<b>0.994</b>
<i>CAMKK2</i>	<b>0.0385</b>	<b>0.998</b>		
<i>CCDC88B</i>	<b>0.0175</b>	<b>0.998</b>		
<i>CKAP2L</i>	<b>0.0046</b>	<b>0.731</b>	<b>0.0187</b>	<b>0.994</b>
<i>CLIC2</i>	<b>0.0191</b>	<b>0.998</b>	<b>0.0397</b>	<b>0.994</b>
<i>ERG3' exons 4-5</i>	<b>0.0275</b>	<b>0.998</b>	<b>0.0243</b>	<b>0.994</b>
<i>ERG3' exons 6-7</i>	<b>2.18 x10<sup>-05</sup></b>	<b>0.0036</b>	<b>0.0002</b>	<b>0.0283</b>
<i>FOLH1</i>	<b>0.0124</b>	<b>0.998</b>	<b>0.0348</b>	<b>0.994</b>
<i>GCNT1</i>	<b>0.0066</b>	<b>0.998</b>	<b>0.0104</b>	<b>0.994</b>
<i>HOXC6</i>	<b>1.36x10<sup>-05</sup></b>	<b>0.0023</b>	<b>0.0002</b>	<b>0.0406</b>
<i>HPN</i>	<b>0.0027</b>	<b>0.4419</b>	<b>0.0081</b>	<b>0.994</b>
<i>ISX</i>	<b>0.0022</b>	<b>0.3607</b>	<b>0.0367</b>	<b>0.994</b>
<i>LASS1</i>	<b>0.005</b>	<b>0.7872</b>	<b>0.0209</b>	<b>0.994</b>
<i>MAPK8IP2</i>	<b>0.0215</b>	<b>0.998</b>		
<i>MCTP1</i>	<b>0.0262</b>	<b>0.998</b>	<b>0.0467</b>	<b>0.994</b>
<i>MEX3A</i>	<b>0.0227</b>	<b>0.998</b>	<b>0.0412</b>	<b>0.994</b>
<i>MFSD2A</i>	<b>0.0103</b>	<b>0.998</b>		
<i>MIR146A</i>	<b>0.0494</b>	<b>0.998</b>		
<i>MMP25</i>	<b>0.0274</b>	<b>0.998</b>	<b>0.0463</b>	<b>0.994</b>
<i>NLRP3</i>	<b>0.0386</b>	<b>0.998</b>		
<i>PCA3</i>	<b>0.0182</b>	<b>0.998</b>		
<i>PSTPIP1</i>	<b>0.0274</b>	<b>0.998</b>		
<i>RIOK3</i>	<b>0.0386</b>	<b>0.998</b>		
<i>SEC61A1</i>	<b>0.0337</b>	<b>0.998</b>		
<i>SFRP4</i>	<b>0.0156</b>	<b>0.998</b>	<b>0.0494</b>	<b>0.994</b>
<i>SIM2 long</i>	<b>0.0031</b>	<b>0.5021</b>	<b>0.0056</b>	<b>0.9028</b>
<i>SLC43A1</i>	<b>0.0199</b>	<b>0.998</b>		
<i>SNORA20</i>	<b>0.045</b>	<b>0.998</b>		
<i>SULF2</i>	<b>0.0153</b>	<b>0.998</b>		
<i>TDRD</i>	<b>0.0002</b>	<b>0.0263</b>	<b>0.0011</b>	<b>0.1757</b>
<i>TMCC2</i>	<b>0.0331</b>	<b>0.998</b>		
<i>TMEM86A</i>	<b>0.0079</b>	<b>0.998</b>	<b>0.0184</b>	<b>0.994</b>
<i>TMPRSS2:ERG</i>	<b>6.86 x10<sup>-05</sup></b>	<b>0.0112</b>	<b>0.0007</b>	<b>0.1136</b>

9: APPENDICES

**Supplementary Table 45** Transcripts that have significant expression trend (using polr and glm) across clinically benign, low-risk, intermediate-risk and high-risk cancer samples in the HK normalised cell NanoString data.

<i>Transcript</i>	<i>Glm p-value</i>	<i>glm Adjusted p-value</i>	<i>Polr p-value</i>	<i>Polr adjusted p-value</i>
<i>CADPS</i>	<b>0.0213</b>	<b>0.9941</b>		
<i>CLIC2</i>	<b>0.0333</b>	<b>0.9941</b>		
<i>EN2</i>			<b>0.0463</b>	<b>0.9994</b>
<i>ERG 3' exons 6-7</i>	<b>0.0043</b>	<b>0.6877</b>	<b>0.0098</b>	<b>0.9994</b>
<i>FOLH1</i>	<b>0.0174</b>	<b>0.9941</b>	<b>0.0191</b>	<b>0.9994</b>
<i>GJB1</i>	<b>0.0215</b>	<b>0.9941</b>		
<i>HOXC6</i>	<b>4.54 x10<sup>-6</sup></b>	<b>0.0008</b>	<b>6.37x10<sup>-05</sup></b>	<b>0.0106</b>
<i>LASS1</i>	<b>0.0287</b>	<b>0.9941</b>		
<i>MEX3A</i>	<b>0.0243</b>	<b>0.9941</b>	<b>0.0337</b>	<b>0.9994</b>
<i>MSMB</i>	<b>0.0334</b>	<b>0.9941</b>		
<i>NAALADL2</i>	<b>0.0018</b>	<b>0.2913</b>	<b>0.0098</b>	<b>0.9994</b>
<i>PALM3</i>	<b>0.027</b>	<b>0.9941</b>	<b>0.0461</b>	<b>0.9994</b>
<i>SERPINB5</i>	<b>0.0162</b>	<b>0.9941</b>	<b>0.0425</b>	<b>0.9994</b>
<i>SIM2 long</i>	<b>0.0032</b>	<b>0.5147</b>	<b>0.0043</b>	<b>0.7056</b>
<i>SLC43A1</i>	<b>0.0011</b>	<b>0.1895</b>	<b>0.006</b>	<b>0.978</b>
<i>ST6GALNAC1</i>	<b>0.0049</b>	<b>0.7755</b>	<b>0.0179</b>	<b>0.9994</b>
<i>TDRD</i>	<b>0.0012</b>	<b>0.2024</b>	<b>0.0034</b>	<b>0.564</b>
<i>TMEM86A</i>	<b>0.0107</b>	<b>0.9941</b>	<b>0.0337</b>	<b>0.9994</b>
<i>TMPRSS2:ERG fusion</i>	<b>0.004</b>	<b>0.6414</b>	<b>0.0127</b>	<b>0.9994</b>
<i>UPK2</i>	<b>0.0028</b>	<b>0.4609</b>	<b>0.0077</b>	<b>0.9994</b>

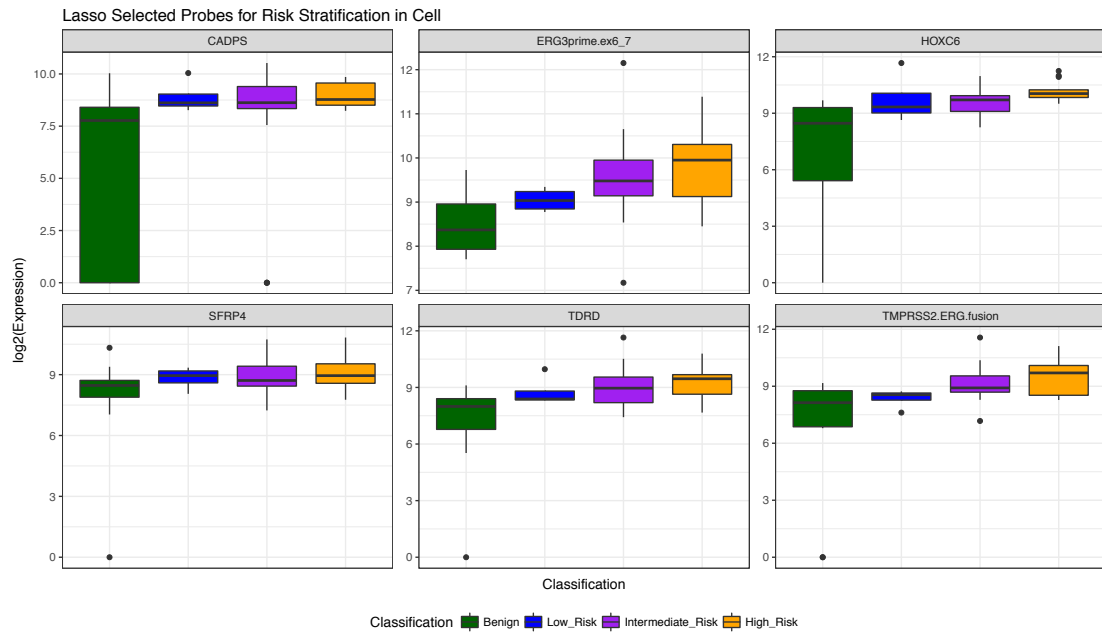


## 9: APPENDICES

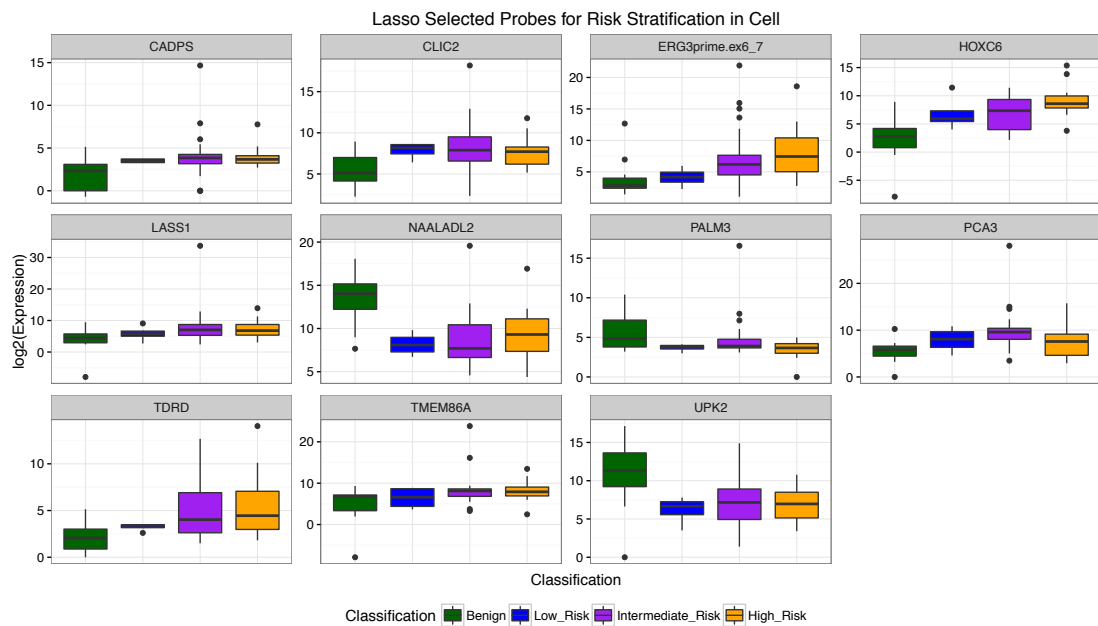


**Supplementary Figure 19** Boxplots of the Lasso identified transcripts for modeling between clinically benign, low-risk, intermediate-risk and high-risk cancer categories in the baseline normalised data.

## 9: APPENDICES



**Supplementary Figure 20** Boxplots showing the Lasso selected transcripts for CB-L-I-H trend in *KLK2* ratio data.



**Supplementary Figure 21** Transcripts selected by Lasso for showing trend of expression levels across clinical categories: clinically benign, low-risk, intermediate-risk and high-risk cancer in the HK normalised cell data.

9: APPENDICES

**Supplementary Table 46 Random Forest rankings for three subsets of transcripts (all 167, the 87 chosen by glm, and the 70 chosen by polr\*), for groups CB, L, I and H. (\*All 70 transcripts identified by polr are were also common to those identified by glm) in the baseline normalization data.**

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 70)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ERG3' exons 6-7</i>	<b>6.60</b>	<b>167</b>	<i>ERG3' exons 6-7</i>	<b>8.10</b>	<b>87</b>	<i>ERG3' exons 6-7</i>	<b>5.00</b>	<b>70</b>
<i>TMPRSS2:ERG fusion</i>	<b>4.31</b>	<b>166</b>	<i>APOC1</i>	<b>3.84</b>	<b>86</b>	<i>TMPRSS2:ERG fusion</i>	<b>4.76</b>	<b>69</b>
<i>RIOK3</i>	<b>2.55</b>	<b>165</b>	<i>SPINK1</i>	<b>2.94</b>	<b>85</b>	<i>NEAT1</i>	<b>3.46</b>	<b>68</b>
<i>NEAT1</i>	<b>2.32</b>	<b>164</b>	<i>CCDC88B</i>	<b>2.58</b>	<b>84</b>	<i>RIOK3</i>	<b>3.08</b>	<b>67</b>
<i>CADPS</i>	<b>1.74</b>	<b>163</b>	<i>CADPS</i>	<b>2.36</b>	<b>83</b>	<i>APOC1</i>	<b>2.26</b>	<b>66</b>
<i>APOC1</i>	<b>1.51</b>	<b>162</b>	<i>B4GALNT4</i>	<b>2.09</b>	<b>82</b>	<i>SIM2 long</i>	<b>2.14</b>	<b>65</b>
<i>SIM2 long</i>	<b>1.51</b>	<b>161</b>	<i>CAMKK2</i>	<b>2.06</b>	<b>81</b>	<i>CCDC88B</i>	<b>1.84</b>	<b>64</b>
<i>SPINK1</i>	<b>1.42</b>	<b>160</b>	<i>GAPDH</i>	<b>1.79</b>	<b>80</b>	<i>MCTP1</i>	<b>1.61</b>	<b>63</b>
<i>MCTP1</i>	<b>1.34</b>	<b>159</b>	<i>CKAP2L</i>	<b>1.45</b>	<b>79</b>	<i>GCNT1</i>	<b>1.55</b>	<b>62</b>
<i>CCDC88B</i>	<b>1.11</b>	<b>158</b>	<i>TDRD</i>	<b>1.24</b>	<b>78</b>	<i>CADPS</i>	<b>1.49</b>	<b>61</b>
<i>MFSD2A</i>	<b>1.02</b>	<b>157</b>	<i>CP</i>	<b>1.19</b>	<b>77</b>	<i>HOXC6</i>	<b>1.49</b>	<b>60</b>
<i>CAMKK2</i>	<b>1.01</b>	<b>156</b>	<i>ISX</i>	<b>1.07</b>	<b>76</b>	<i>MFSD2A</i>	<b>1.45</b>	<b>59</b>
<i>GCNT1</i>	<b>1.00</b>	<b>155</b>	<i>CD10</i>	<b>1.01</b>	<b>75</b>	<i>SPINK1</i>	<b>1.43</b>	<b>58</b>
<i>MXI1</i>	<b>0.90</b>	<b>154</b>	<i>HPN</i>	<b>0.99</b>	<b>74</b>	<i>CAMKK2</i>	<b>1.40</b>	<b>57</b>
<i>TMEM86A</i>	<b>0.89</b>	<b>153</b>	<i>UPK2</i>	<b>0.97</b>	<b>73</b>	<i>SIRT1</i>	<b>1.27</b>	<b>56</b>
<i>HOXC6</i>	<b>0.86</b>	<b>152</b>	<i>SLC4A1 S</i>	<b>0.93</b>	<b>72</b>	<i>SULT1A1</i>	<b>1.18</b>	<b>55</b>
<i>CP</i>	<b>0.86</b>	<b>151</b>	<i>PCA3</i>	<b>0.89</b>	<b>71</b>	<i>CKAP2L</i>	<b>1.14</b>	<b>54</b>
<i>SLC43A1</i>	<b>0.83</b>	<b>150</b>	<i>CACNAID</i>	<b>0.89</b>	<b>70</b>	<i>LASS1</i>	<b>0.93</b>	<b>53</b>
<i>SFRP4</i>	<b>0.82</b>	<b>149</b>	<i>AATF</i>	<b>0.83</b>	<b>69</b>	<i>TMEM86A</i>	<b>0.91</b>	<b>52</b>
<i>B4GALNT4</i>	<b>0.78</b>	<b>148</b>	<i>SNORA20</i>	<b>0.81</b>	<b>68</b>	<i>SLC43A1</i>	<b>0.88</b>	<b>51</b>
<i>SIRT1</i>	<b>0.72</b>	<b>147</b>	<i>IGFBP3</i>	<b>0.75</b>	<b>67</b>	<i>AURKA</i>	<b>0.87</b>	<b>50</b>
<i>MIR4435 IHG</i>	<b>0.72</b>	<b>146</b>	<i>ANKRD34B</i>	<b>0.70</b>	<b>66</b>	<i>HPRT</i>	<b>0.85</b>	<b>49</b>
<i>CKAP2L</i>	<b>0.71</b>	<b>145</b>	<i>CLIC2</i>	<b>0.69</b>	<b>65</b>	<i>MIR4435 IHG</i>	<b>0.80</b>	<b>48</b>
<i>SULT1A1</i>	<b>0.71</b>	<b>144</b>	<i>SMAP1 exons 7-</i>	<b>0.69</b>	<b>64</b>	<i>B4GALNT4</i>	<b>0.76</b>	<b>47</b>

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 70)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
			<b>8</b>					
<i>LASS1</i>	<b>0.67</b>	<b>143</b>	<i>DLX1</i>	<b>0.66</b>	<b>63</b>	<i>CD10</i>	<b>0.72</b>	<b>46</b>
<i>MYOF</i>	<b>0.56</b>	<b>142</b>	<i>AURKA</i>	<b>0.65</b>	<b>62</b>	<i>MIC1</i>	<b>0.69</b>	<b>45</b>
<i>HPN</i>	<b>0.54</b>	<b>141</b>	<i>MYOF</i>	<b>0.64</b>	<b>61</b>	<i>MMP25</i>	<b>0.69</b>	<b>44</b>
<i>PCA3</i>	<b>0.54</b>	<b>140</b>	<i>ERG3' exons 4-5</i>	<b>0.62</b>	<b>60</b>	<i>SFRP4</i>	<b>0.58</b>	<b>43</b>
<i>SNORA20</i>	<b>0.52</b>	<b>139</b>	<i>AMH</i>	<b>0.59</b>	<b>59</b>	<i>ANKRD34B</i>	<b>0.57</b>	<b>42</b>
<i>SLC4A1 S</i>	<b>0.50</b>	<b>138</b>	<i>TERF2IP</i>	<b>0.58</b>	<b>58</b>	<i>SRSF3</i>	<b>0.56</b>	<b>41</b>
<i>CD10</i>	<b>0.49</b>	<b>137</b>	<i>AMACR</i>	<b>0.58</b>	<b>57</b>	<i>HIST1H2BF</i>	<b>0.51</b>	<b>40</b>
<i>SMAP1 exons 7-8</i>	<b>0.49</b>	<b>136</b>	<i>ERG5'</i>	<b>0.55</b>	<b>56</b>	<i>SULF2</i>	<b>0.50</b>	<b>39</b>
<i>MMP25</i>	<b>0.45</b>	<b>135</b>	<i>MAK</i>	<b>0.48</b>	<b>55</b>	<i>TDRD</i>	<b>0.50</b>	<b>38</b>
<i>AMH</i>	<b>0.42</b>	<b>134</b>	<i>ANPEP</i>	<b>0.46</b>	<b>54</b>	<i>CDKN3</i>	<b>0.49</b>	<b>37</b>
<i>TDRD</i>	<b>0.41</b>	<b>133</b>	<i>DPP4</i>	<b>0.44</b>	<b>53</b>	<i>EN2</i>	<b>0.48</b>	<b>36</b>
<i>SULF2</i>	<b>0.40</b>	<b>132</b>	<i>NAALADL2</i>	<b>0.44</b>	<b>52</b>	<i>MXI1</i>	<b>0.46</b>	<b>35</b>
<i>OR52A2</i>	<b>0.39</b>	<b>131</b>	<i>MEX3A</i>	<b>0.43</b>	<b>51</b>	<i>STEAP4</i>	<b>0.46</b>	<b>34</b>
<i>MIC1</i>	<b>0.39</b>	<b>130</b>	<i>EN2</i>	<b>0.43</b>	<b>50</b>	<i>ITPR1</i>	<b>0.45</b>	<b>33</b>
<i>HPRT</i>	<b>0.38</b>	<b>129</b>	<i>RNF157</i>	<b>0.41</b>	<b>49</b>	<i>ERG3' exons 4-5</i>	<b>0.45</b>	<b>32</b>
<i>HIST1H1E</i>	<b>0.37</b>	<b>128</b>	<i>DNAH5</i>	<b>0.40</b>	<b>48</b>	<i>SNCA</i>	<b>0.44</b>	<b>31</b>
<i>MAPK8IP2</i>	<b>0.37</b>	<b>127</b>	<i>HIST1H2BG</i>	<b>0.39</b>	<b>47</b>	<i>MEX3A</i>	<b>0.44</b>	<b>30</b>
<i>UPK2</i>	<b>0.37</b>	<b>126</b>	<i>ACTR5</i>	<b>0.38</b>	<b>46</b>	<i>PSTPIP1</i>	<b>0.44</b>	<b>29</b>
<i>GAPDH</i>	<b>0.36</b>	<b>125</b>	<i>TFDP1</i>	<b>0.37</b>	<b>45</b>	<i>GAPDH</i>	<b>0.41</b>	<b>28</b>
<i>AURKA</i>	<b>0.33</b>	<b>124</b>	<i>CDKN3</i>	<b>0.37</b>	<b>44</b>	<i>CACNAID</i>	<b>0.40</b>	<b>27</b>
<i>ANKRD34B</i>	<b>0.32</b>	<b>123</b>	<i>MCM7</i>	<b>0.37</b>	<b>43</b>	<i>ISX</i>	<b>0.39</b>	<b>26</b>
<i>STOM</i>	<b>0.29</b>	<b>122</b>	<i>CDC20</i>	<b>0.36</b>	<b>42</b>	<i>TMCC2</i>	<b>0.38</b>	<b>25</b>
<i>Timp4</i>	<b>0.29</b>	<b>121</b>	<i>AR exon 9</i>	<b>0.34</b>	<b>41</b>	<i>PTPRC</i>	<b>0.38</b>	<b>24</b>
<i>DLX1</i>	<b>0.28</b>	<b>120</b>	<i>PALM3</i>	<b>0.33</b>	<b>40</b>	<i>AATF</i>	<b>0.38</b>	<b>23</b>
<i>EN2</i>	<b>0.27</b>	<b>119</b>	<i>SSTR1</i>	<b>0.31</b>	<b>39</b>	<i>MMP26</i>	<b>0.38</b>	<b>22</b>
<i>RPL23AP53</i>	<b>0.27</b>	<b>118</b>	<i>AGR2</i>	<b>0.31</b>	<b>38</b>	<i>CLIC2</i>	<b>0.36</b>	<b>21</b>
<i>ITPR1</i>	<b>0.26</b>	<b>117</b>	<i>ABCB9</i>	<b>0.30</b>	<b>37</b>	<i>COL9A2</i>	<b>0.34</b>	<b>20</b>

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 70)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ISX</i>	<i>0.26</i>	<i>116</i>	<i>CAMK2N2</i>	<i>0.29</i>	<i>36</i>	<i>TERF2IP</i>	<i>0.33</i>	<i>19</i>
<i>RNF157</i>	<i>0.25</i>	<i>115</i>	<i>CLU</i>	<i>0.28</i>	<i>35</i>	<i>MMP11</i>	<i>0.31</i>	<i>18</i>
<i>TMEM45B</i>	<i>0.24</i>	<i>114</i>	<i>BTG2</i>	<i>0.27</i>	<i>34</i>	<i>ACTR5</i>	<i>0.30</i>	<i>17</i>
<i>FOLH1</i>	<i>0.24</i>	<i>113</i>	<i>ALAS1</i>	<i>0.27</i>	<i>33</i>	<i>ANPEP</i>	<i>0.30</i>	<i>16</i>
<i>CDKN3</i>	<i>0.24</i>	<i>112</i>	<i>NLRP3</i>	<i>0.27</i>	<i>32</i>	<i>SACMIL</i>	<i>0.30</i>	<i>15</i>
<i>AATF</i>	<i>0.23</i>	<i>111</i>	<i>TMCC2</i>	<i>0.26</i>	<i>31</i>	<i>HIST1H2BG</i>	<i>0.30</i>	<i>14</i>
<i>SLC12A1</i>	<i>0.23</i>	<i>110</i>	<i>OGT</i>	<i>0.25</i>	<i>30</i>	<i>PDLIM5</i>	<i>0.28</i>	<i>13</i>
<i>CACNA1D</i>	<i>0.23</i>	<i>109</i>	<i>EIF2D</i>	<i>0.24</i>	<i>29</i>	<i>NLRP3</i>	<i>0.27</i>	<i>12</i>
<i>SACMIL</i>	<i>0.23</i>	<i>108</i>	<i>ARHGEF25</i>	<i>0.24</i>	<i>28</i>	<i>HPN</i>	<i>0.27</i>	<i>11</i>
<i>IGFBP3</i>	<i>0.22</i>	<i>107</i>	<i>FOLH1</i>	<i>0.24</i>	<i>27</i>	<i>MGAT5B</i>	<i>0.26</i>	<i>10</i>
<i>ACTR5</i>	<i>0.21</i>	<i>106</i>	<i>LBH</i>	<i>0.23</i>	<i>26</i>	<i>FOLH1</i>	<i>0.26</i>	<i>9</i>
<i>MEX3A</i>	<i>0.20</i>	<i>105</i>	<i>TMEM47</i>	<i>0.22</i>	<i>25</i>	<i>EIF2D</i>	<i>0.21</i>	<i>8</i>
<i>DPP4</i>	<i>0.20</i>	<i>104</i>	<i>ST6GALNAC1</i>	<i>0.22</i>	<i>24</i>	<i>SIM2 short</i>	<i>0.20</i>	<i>7</i>
<i>SRSF3</i>	<i>0.20</i>	<i>103</i>	<i>B2M</i>	<i>0.22</i>	<i>23</i>	<i>CDC20</i>	<i>0.20</i>	<i>6</i>
<i>AR.ex9</i>	<i>0.20</i>	<i>102</i>	<i>Met</i>	<i>0.21</i>	<i>22</i>	<i>BTG2</i>	<i>0.19</i>	<i>5</i>
<i>ANPEP</i>	<i>0.19</i>	<i>101</i>	<i>RP11_97012.7</i>	<i>0.19</i>	<i>21</i>	<i>GABARAPL2</i>	<i>0.19</i>	<i>4</i>
<i>VAX2</i>	<i>0.19</i>	<i>100</i>	<i>PPAP2A</i>	<i>0.19</i>	<i>20</i>	<i>SEC61A1</i>	<i>0.16</i>	<i>3</i>
<i>ERG5'</i>	<i>0.19</i>	<i>99</i>	<i>COL9A2</i>	<i>0.18</i>	<i>19</i>	<i>B2M</i>	<i>0.15</i>	<i>2</i>
<i>COL9A2</i>	<i>0.19</i>	<i>98</i>	<i>COL10A1</i>	<i>0.18</i>	<i>18</i>	<i>ARHGEF25</i>	<i>0.15</i>	<i>1</i>
<i>AGR2</i>	<i>0.18</i>	<i>97</i>	<i>GOLM1</i>	<i>0.17</i>	<i>17</i>			
<i>CLIC2</i>	<i>0.18</i>	<i>96</i>	<i>PECI</i>	<i>0.16</i>	<i>16</i>			
<i>ABCB9</i>	<i>0.18</i>	<i>95</i>	<i>AR exons 4-8</i>	<i>0.15</i>	<i>15</i>			
<i>DNAH5</i>	<i>0.18</i>	<i>94</i>	<i>CDC37L1</i>	<i>0.15</i>	<i>14</i>			
<i>MSMB</i>	<i>0.18</i>	<i>93</i>	<i>SIM2 short</i>	<i>0.15</i>	<i>13</i>			
<i>STEAP4</i>	<i>0.17</i>	<i>92</i>	<i>TBP</i>	<i>0.14</i>	<i>12</i>			
<i>SSTR1</i>	<i>0.17</i>	<i>91</i>	<i>MIATNB</i>	<i>0.14</i>	<i>11</i>			
<i>HIST1H2BF</i>	<i>0.16</i>	<i>90</i>	<i>PCSK6</i>	<i>0.14</i>	<i>10</i>			
<i>HMBS</i>	<i>0.15</i>	<i>89</i>	<i>MARCH5</i>	<i>0.13</i>	<i>9</i>			

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 70)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>SChLAP1</i>	<b>0.15</b>	<b>88</b>	<i>BRAF</i>	<b>0.13</b>	<b>8</b>			
<i>LBH</i>	<b>0.15</b>	<b>87</b>	<i>HOXC4</i>	<b>0.13</b>	<b>7</b>			
<i>PPFIA2</i>	<b>0.15</b>	<b>86</b>	<i>FDPS</i>	<b>0.12</b>	<b>6</b>			
<i>NLRP3</i>	<b>0.15</b>	<b>85</b>	<i>KLK2</i>	<b>0.12</b>	<b>5</b>			
<i>PTPRC</i>	<b>0.15</b>	<b>84</b>	<i>RPS11</i>	<b>0.11</b>	<b>4</b>			
<i>SPON2</i>	<b>0.14</b>	<b>83</b>	<i>CASKIN1</i>	<b>0.11</b>	<b>3</b>			
<i>AMACR</i>	<b>0.14</b>	<b>82</b>	<i>GABARAPL2</i>	<b>0.10</b>	<b>2</b>			
<i>MCM7</i>	<b>0.14</b>	<b>81</b>	<i>STEAP2</i>	<b>0.08</b>	<b>1</b>			
<i>MIR146A</i>	<b>0.14</b>	<b>80</b>						
<i>KLK4</i>	<b>0.14</b>	<b>79</b>						
<i>ALAS1</i>	<b>0.13</b>	<b>78</b>						
<i>MMP26</i>	<b>0.13</b>	<b>77</b>						
<i>PPAP2A</i>	<b>0.13</b>	<b>76</b>						
<i>MNX1</i>	<b>0.13</b>	<b>75</b>						
<i>ERG3' exons 4-5</i>	<b>0.13</b>	<b>74</b>						
<i>CDC37L1</i>	<b>0.13</b>	<b>73</b>						
<i>NAALADL2</i>	<b>0.13</b>	<b>72</b>						
<i>PSTPIP1</i>	<b>0.12</b>	<b>71</b>						
<i>SSPO</i>	<b>0.11</b>	<b>70</b>						
<i>EIF2D</i>	<b>0.11</b>	<b>69</b>						
<i>CDC20</i>	<b>0.11</b>	<b>68</b>						
<i>CLU</i>	<b>0.11</b>	<b>67</b>						
<i>PALM3</i>	<b>0.11</b>	<b>66</b>						
<i>KLK3 exons 2-3</i>	<b>0.11</b>	<b>65</b>						
<i>TRPM4</i>	<b>0.11</b>	<b>64</b>						
<i>MKi67</i>	<b>0.11</b>	<b>63</b>						
<i>TERF2IP</i>	<b>0.10</b>	<b>62</b>						
<i>HOXC4</i>	<b>0.10</b>	<b>61</b>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 70)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>COL10A1</i>	<i>0.10</i>	<i>60</i>						
<i>RP11_97012.7</i>	<i>0.10</i>	<i>59</i>						
<i>SEC61A1</i>	<i>0.10</i>	<i>58</i>						
<i>RAB17</i>	<i>0.10</i>	<i>57</i>						
<i>NKAIN1</i>	<i>0.10</i>	<i>56</i>						
<i>MDK</i>	<i>0.10</i>	<i>55</i>						
<i>SNCA</i>	<i>0.09</i>	<i>54</i>						
<i>MGAT5B</i>	<i>0.09</i>	<i>53</i>						
<i>VPS13A</i>	<i>0.09</i>	<i>52</i>						
<i>MED4</i>	<i>0.09</i>	<i>51</i>						
<i>ARHGEF25</i>	<i>0.09</i>	<i>50</i>						
<i>MAK</i>	<i>0.09</i>	<i>49</i>						
<i>PPP1R12B</i>	<i>0.09</i>	<i>48</i>						
<i>TBP</i>	<i>0.09</i>	<i>47</i>						
<i>SERPINB5</i>	<i>0.08</i>	<i>46</i>						
<i>GJB1</i>	<i>0.08</i>	<i>45</i>						
<i>BTG2</i>	<i>0.08</i>	<i>44</i>						
<i>MEMO1</i>	<i>0.08</i>	<i>43</i>						
<i>HIST3H2A</i>	<i>0.08</i>	<i>42</i>						
<i>TERT</i>	<i>0.08</i>	<i>41</i>						
<i>PVT1</i>	<i>0.08</i>	<i>40</i>						
<i>TFDP1</i>	<i>0.07</i>	<i>39</i>						
<i>P712P</i>	<i>0.07</i>	<i>38</i>						
<i>ZNF577</i>	<i>0.07</i>	<i>37</i>						
<i>Met</i>	<i>0.07</i>	<i>36</i>						
<i>OGT</i>	<i>0.07</i>	<i>35</i>						
<i>AR exons 4-8</i>	<i>0.07</i>	<i>34</i>						
<i>ITGBL1</i>	<i>0.07</i>	<i>33</i>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 70)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>GOLM1</i>	<i>0.07</i>	<i>32</i>						
<i>FDPS</i>	<i>0.07</i>	<i>31</i>						
<i>MIATNB</i>	<i>0.07</i>	<i>30</i>						
<i>B2M</i>	<i>0.07</i>	<i>29</i>						
<i>RPS10</i>	<i>0.06</i>	<i>28</i>						
<i>ST6GALNAC1</i>	<i>0.06</i>	<i>27</i>						
<i>RPL18A</i>	<i>0.06</i>	<i>26</i>						
<i>IMPDH2</i>	<i>0.06</i>	<i>25</i>						
<i>SMIMI</i>	<i>0.05</i>	<i>24</i>						
<i>HIST1H2BG</i>	<i>0.05</i>	<i>23</i>						
<i>TMCC2</i>	<i>0.05</i>	<i>22</i>						
<i>STEAP2</i>	<i>0.05</i>	<i>21</i>						
<i>RPS11</i>	<i>0.05</i>	<i>20</i>						
<i>IFT57</i>	<i>0.05</i>	<i>19</i>						
<i>BRAF</i>	<i>0.05</i>	<i>18</i>						
<i>TWIST1</i>	<i>0.05</i>	<i>17</i>						
<i>CAMK2N2</i>	<i>0.05</i>	<i>16</i>						
<i>SIM2 short</i>	<i>0.05</i>	<i>15</i>						
<i>MMP11</i>	<i>0.04</i>	<i>14</i>						
<i>HIST1H1C</i>	<i>0.04</i>	<i>13</i>						
<i>PCSK6</i>	<i>0.04</i>	<i>12</i>						
<i>PECI</i>	<i>0.04</i>	<i>11</i>						
<i>PDLIM5</i>	<i>0.04</i>	<i>10</i>						
<i>MARCH5</i>	<i>0.04</i>	<i>9</i>						
<i>CASKIN1</i>	<i>0.04</i>	<i>8</i>						
<i>TMEM47</i>	<i>0.04</i>	<i>7</i>						
<i>RPLP2</i>	<i>0.04</i>	<i>6</i>						
<i>KLK2</i>	<i>0.03</i>	<i>5</i>						



9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 87)</i>			<i>Transcripts identified by polr (n = 70)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>GABARAPL2</i>	<b>0.03</b>	<b>4</b>						
<i>PTN</i>	<b>0.03</b>	<b>3</b>						
<i>KLK3 exons 1-2</i>	<b>0.03</b>	<b>2</b>						
<i>SYNM</i>	<b>0.01</b>	<b>1</b>						

## 9: APPENDICES

**Supplementary Table 47 Random Forest results for trend across clinical categories: CBN-L-I-H in *KLK2* factorised data.**

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 36)</i>			<i>Transcripts identified by polr (n = 20)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>PCA3</i>	<b>1.19</b>	<b>166</b>	<i>PCA3</i>	<b>3.36</b>	<b>36</b>	<i>HOXC6</i>	<b>3.58</b>	<b>20</b>
<i>HOXC6</i>	<b>1.14</b>	<b>165</b>	<i>HOXC6</i>	<b>2.56</b>	<b>35</b>	<i>ERG3' exons 6-7</i>	<b>2.98</b>	<b>19</b>
<i>TMPRSS2:ERG fusion</i>	<b>0.89</b>	<b>164</b>	<i>ERG3' exons 6-7</i>	<b>1.85</b>	<b>34</b>	<i>TMPRSS2:ERG fusion</i>	<b>2.97</b>	<b>18</b>
<i>ERG3' exons 6-7</i>	<b>0.88</b>	<b>163</b>	<i>TMPRSS2:ERG fusion</i>	<b>1.81</b>	<b>33</b>	<i>FOLH1</i>	<b>2.30</b>	<b>17</b>
<i>SLC12A1</i>	<b>0.59</b>	<b>162</b>	<i>FOLH1</i>	<b>1.42</b>	<b>32</b>	<i>TMEM86A</i>	<b>2.09</b>	<b>16</b>
<i>NAALADL2</i>	<b>0.59</b>	<b>161</b>	<i>TDRD</i>	<b>1.24</b>	<b>31</b>	<i>HPN</i>	<b>2.01</b>	<b>15</b>
<i>APOC1</i>	<b>0.53</b>	<b>160</b>	<i>APOC1</i>	<b>1.20</b>	<b>30</b>	<i>CKAP2L</i>	<b>1.97</b>	<b>14</b>
<i>FOLH1</i>	<b>0.50</b>	<b>159</b>	<i>SLC43A1</i>	<b>1.16</b>	<b>29</b>	<i>GCNT1</i>	<b>1.82</b>	<b>13</b>
<i>CP</i>	<b>0.49</b>	<b>158</b>	<i>TMEM86A</i>	<b>1.15</b>	<b>28</b>	<i>CADPS</i>	<b>1.79</b>	<b>12</b>
<i>OR52A2</i>	<b>0.49</b>	<b>157</b>	<i>GCNT1</i>	<b>1.11</b>	<b>27</b>	<i>TDRD</i>	<b>1.77</b>	<b>11</b>
<i>SIM2 long</i>	<b>0.49</b>	<b>156</b>	<i>CKAP2L</i>	<b>1.09</b>	<b>26</b>	<i>MMP25</i>	<b>1.71</b>	<b>10</b>
<i>TDRD</i>	<b>0.47</b>	<b>155</b>	<i>SIM2 long</i>	<b>1.08</b>	<b>25</b>	<i>SIM2 long</i>	<b>1.54</b>	<b>9</b>
<i>PALM3</i>	<b>0.45</b>	<b>154</b>	<i>HPN</i>	<b>1.07</b>	<b>24</b>	<i>CLIC2</i>	<b>1.52</b>	<b>8</b>
<i>SERPINB5</i>	<b>0.44</b>	<b>153</b>	<i>B4GALNT4</i>	<b>1.07</b>	<b>23</b>	<i>ISX</i>	<b>1.48</b>	<b>7</b>
<i>AR exons 4-8</i>	<b>0.41</b>	<b>152</b>	<i>CADPS</i>	<b>1.02</b>	<b>22</b>	<i>ANKRD34B</i>	<b>1.47</b>	<b>6</b>
<i>TMEM86A</i>	<b>0.40</b>	<b>151</b>	<i>MAPK8IP2</i>	<b>0.95</b>	<b>21</b>	<i>MCTP1</i>	<b>1.46</b>	<b>5</b>
<i>MSMB</i>	<b>0.40</b>	<b>150</b>	<i>ANKRD34B</i>	<b>0.91</b>	<b>20</b>	<i>LASS1</i>	<b>1.44</b>	<b>4</b>
<i>MDK</i>	<b>0.39</b>	<b>149</b>	<i>MMP25</i>	<b>0.90</b>	<b>19</b>	<i>SFRP4</i>	<b>1.40</b>	<b>3</b>
<i>CKAP2L</i>	<b>0.38</b>	<b>148</b>	<i>SEC61A1</i>	<b>0.90</b>	<b>18</b>	<i>MEX3A</i>	<b>1.30</b>	<b>2</b>
<i>DLX1</i>	<b>0.37</b>	<b>147</b>	<i>LASS1</i>	<b>0.87</b>	<b>17</b>	<i>ERG3' exons 4-5</i>	<b>1.08</b>	<b>1</b>
<i>HPN</i>	<b>0.36</b>	<b>146</b>	<i>CLIC2</i>	<b>0.85</b>	<b>16</b>			
<i>SLC43A1</i>	<b>0.36</b>	<b>145</b>	<i>SULF2</i>	<b>0.82</b>	<b>15</b>			
<i>STEAP2</i>	<b>0.35</b>	<b>144</b>	<i>TMCC2</i>	<b>0.80</b>	<b>14</b>			
<i>IGFBP3</i>	<b>0.34</b>	<b>143</b>	<i>CCDC88B</i>	<b>0.77</b>	<b>13</b>			

## 9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 36)</i>			<i>Transcripts identified by polr (n = 20)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>SPON2</i>	<i>0.34</i>	<i>142</i>	<i>MEX3A</i>	<i>0.75</i>	<i>12</i>			
<i>LASS1</i>	<i>0.34</i>	<i>141</i>	<i>CAMKK2</i>	<i>0.74</i>	<i>11</i>			
<i>TMEM47</i>	<i>0.34</i>	<i>140</i>	<i>SNORA20</i>	<i>0.72</i>	<i>10</i>			
<i>AGR2</i>	<i>0.33</i>	<i>139</i>	<i>SFRP4</i>	<i>0.72</i>	<i>9</i>			
<i>CADPS</i>	<i>0.32</i>	<i>138</i>	<i>MFS2A</i>	<i>0.69</i>	<i>8</i>			
<i>MMP11</i>	<i>0.31</i>	<i>137</i>	<i>MCTP1</i>	<i>0.68</i>	<i>7</i>			
<i>GJB1</i>	<i>0.30</i>	<i>136</i>	<i>ERG3' exons 4-5</i>	<i>0.66</i>	<i>6</i>			
<i>SSTR1</i>	<i>0.30</i>	<i>135</i>	<i>NLRP3</i>	<i>0.62</i>	<i>5</i>			
<i>TMCC2</i>	<i>0.30</i>	<i>134</i>	<i>PSTPIP1</i>	<i>0.59</i>	<i>4</i>			
<i>AMACR</i>	<i>0.30</i>	<i>133</i>	<i>MIR146A</i>	<i>0.58</i>	<i>3</i>			
<i>B4GALNT4</i>	<i>0.29</i>	<i>132</i>	<i>RIOK3</i>	<i>0.54</i>	<i>2</i>			
<i>SULF2</i>	<i>0.29</i>	<i>131</i>	<i>ISX</i>	<i>0.45</i>	<i>1</i>			
<i>GCNT1</i>	<i>0.29</i>	<i>130</i>						
<i>ZNF577</i>	<i>0.28</i>	<i>129</i>						
<i>ANKRD34B</i>	<i>0.28</i>	<i>128</i>						
<i>HIST1H2BG</i>	<i>0.27</i>	<i>127</i>						
<i>SPINK1</i>	<i>0.27</i>	<i>126</i>						
<i>MMP25</i>	<i>0.27</i>	<i>125</i>						
<i>HIST3H2A</i>	<i>0.26</i>	<i>124</i>						
<i>TRPM4</i>	<i>0.26</i>	<i>123</i>						
<i>SLC4A1.S</i>	<i>0.25</i>	<i>122</i>						
<i>SULT1A1</i>	<i>0.25</i>	<i>121</i>						
<i>CDKN3</i>	<i>0.25</i>	<i>120</i>						
<i>Timp4</i>	<i>0.25</i>	<i>119</i>						
<i>ST6GALNAC1</i>	<i>0.25</i>	<i>118</i>						
<i>SNORA20</i>	<i>0.25</i>	<i>117</i>						
<i>EN2</i>	<i>0.25</i>	<i>116</i>						
<i>AR exon 9</i>	<i>0.24</i>	<i>115</i>						

## 9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 36)</i>			<i>Transcripts identified by polr (n = 20)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ITGBL1</i>	<i>0.24</i>	<i>114</i>						
<i>UPK2</i>	<i>0.24</i>	<i>113</i>						
<i>MKi67</i>	<i>0.24</i>	<i>112</i>						
<i>SChLAP1</i>	<i>0.24</i>	<i>111</i>						
<i>AMH</i>	<i>0.23</i>	<i>110</i>						
<i>MCTP1</i>	<i>0.23</i>	<i>109</i>						
<i>SFRP4</i>	<i>0.23</i>	<i>108</i>						
<i>MFSD2A</i>	<i>0.23</i>	<i>107</i>						
<i>SIM2 short</i>	<i>0.23</i>	<i>106</i>						
<i>PPP1R12B</i>	<i>0.23</i>	<i>105</i>						
<i>TERT</i>	<i>0.23</i>	<i>104</i>						
<i>RAB17</i>	<i>0.22</i>	<i>103</i>						
<i>NKAIN1</i>	<i>0.22</i>	<i>102</i>						
<i>SMIMI</i>	<i>0.22</i>	<i>101</i>						
<i>P712P</i>	<i>0.22</i>	<i>100</i>						
<i>ERG3' exons 4-5</i>	<i>0.22</i>	<i>99</i>						
<i>PECI</i>	<i>0.22</i>	<i>98</i>						
<i>ERG5'</i>	<i>0.22</i>	<i>97</i>						
<i>VAX2</i>	<i>0.22</i>	<i>96</i>						
<i>CLIC2</i>	<i>0.22</i>	<i>95</i>						
<i>RNF157</i>	<i>0.21</i>	<i>94</i>						
<i>CDC37L1</i>	<i>0.21</i>	<i>93</i>						
<i>CCDC88B</i>	<i>0.21</i>	<i>92</i>						
<i>CLU</i>	<i>0.20</i>	<i>91</i>						
<i>MIC1</i>	<i>0.20</i>	<i>90</i>						
<i>TMEM45B</i>	<i>0.20</i>	<i>89</i>						
<i>MNX1</i>	<i>0.20</i>	<i>88</i>						
<i>ISX</i>	<i>0.20</i>	<i>87</i>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 36)</i>			<i>Transcripts identified by polr (n = 20)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>HIST1H1C</i>	<b>0.19</b>	<b>86</b>						
<i>KLK4</i>	<b>0.18</b>	<b>85</b>						
<i>LBH</i>	<b>0.18</b>	<b>84</b>						
<i>COL10A1</i>	<b>0.18</b>	<b>83</b>						
<i>MED4</i>	<b>0.18</b>	<b>82</b>						
<i>HIST1H2BF</i>	<b>0.18</b>	<b>81</b>						
<i>PPAP2A</i>	<b>0.18</b>	<b>80</b>						
<i>ABCB9</i>	<b>0.17</b>	<b>79</b>						
<i>STOM</i>	<b>0.17</b>	<b>78</b>						
<i>DNAH5</i>	<b>0.17</b>	<b>77</b>						
<i>DPP4</i>	<b>0.17</b>	<b>76</b>						
<i>MMP26</i>	<b>0.17</b>	<b>75</b>						
<i>HOXC4</i>	<b>0.16</b>	<b>74</b>						
<i>MGAT5B</i>	<b>0.16</b>	<b>73</b>						
<i>MIR146A</i>	<b>0.16</b>	<b>72</b>						
<i>PCSK6</i>	<b>0.16</b>	<b>71</b>						
<i>CAMKK2</i>	<b>0.16</b>	<b>70</b>						
<i>MARCH5</i>	<b>0.15</b>	<b>69</b>						
<i>RPL23AP53</i>	<b>0.15</b>	<b>68</b>						
<i>IMPDH2</i>	<b>0.15</b>	<b>67</b>						
<i>HPRT</i>	<b>0.15</b>	<b>66</b>						
<i>ACTR5</i>	<b>0.15</b>	<b>65</b>						
<i>MAPK8IP2</i>	<b>0.15</b>	<b>64</b>						
<i>SNCA</i>	<b>0.15</b>	<b>63</b>						
<i>SYNM</i>	<b>0.15</b>	<b>62</b>						
<i>PSTPIP1</i>	<b>0.15</b>	<b>61</b>						
<i>CACNA1D</i>	<b>0.14</b>	<b>60</b>						
<i>PVT1</i>	<b>0.14</b>	<b>59</b>						

## 9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 36)</i>			<i>Transcripts identified by polr (n = 20)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>HMBS</i>	<i>0.14</i>	<i>58</i>						
<i>SACM1L</i>	<i>0.14</i>	<i>57</i>						
<i>KLK3 exons 2-3</i>	<i>0.14</i>	<i>56</i>						
<i>COL9A2</i>	<i>0.14</i>	<i>55</i>						
<i>SRSF3</i>	<i>0.14</i>	<i>54</i>						
<i>KLK3 exons 1-2</i>	<i>0.13</i>	<i>53</i>						
<i>RPS10</i>	<i>0.13</i>	<i>52</i>						
<i>NLRP3</i>	<i>0.13</i>	<i>51</i>						
<i>RP11_97012.7</i>	<i>0.13</i>	<i>50</i>						
<i>PPFIA2</i>	<i>0.13</i>	<i>49</i>						
<i>SMAP1 exons 7-8</i>	<i>0.13</i>	<i>48</i>						
<i>MAK</i>	<i>0.13</i>	<i>47</i>						
<i>AATF</i>	<i>0.13</i>	<i>46</i>						
<i>CDC20</i>	<i>0.13</i>	<i>45</i>						
<i>MXII</i>	<i>0.13</i>	<i>44</i>						
<i>SSPO</i>	<i>0.13</i>	<i>43</i>						
<i>MEX3A</i>	<i>0.13</i>	<i>42</i>						
<i>MCM7</i>	<i>0.12</i>	<i>41</i>						
<i>PDLIM5</i>	<i>0.12</i>	<i>40</i>						
<i>OGT</i>	<i>0.12</i>	<i>39</i>						
<i>GOLM1</i>	<i>0.12</i>	<i>38</i>						
<i>MYOF</i>	<i>0.12</i>	<i>37</i>						
<i>VPSI3A</i>	<i>0.12</i>	<i>36</i>						
<i>CASKIN1</i>	<i>0.12</i>	<i>35</i>						
<i>RPS11</i>	<i>0.11</i>	<i>34</i>						
<i>RIOK3</i>	<i>0.11</i>	<i>33</i>						
<i>B2M</i>	<i>0.11</i>	<i>32</i>						
<i>FDPS</i>	<i>0.11</i>	<i>31</i>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 36)</i>			<i>Transcripts identified by polr (n = 20)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>ANPEP</i>	<i>0.11</i>	<i>30</i>						
<i>CD10</i>	<i>0.11</i>	<i>29</i>						
<i>RPL18A</i>	<i>0.11</i>	<i>28</i>						
<i>ITPRI</i>	<i>0.11</i>	<i>27</i>						
<i>SEC61A1</i>	<i>0.11</i>	<i>26</i>						
<i>EIF2D</i>	<i>0.11</i>	<i>25</i>						
<i>TFDP1</i>	<i>0.10</i>	<i>24</i>						
<i>TWIST1</i>	<i>0.10</i>	<i>23</i>						
<i>MEMO1</i>	<i>0.10</i>	<i>22</i>						
<i>RPLP2</i>	<i>0.10</i>	<i>21</i>						
<i>HIST1H1E</i>	<i>0.10</i>	<i>20</i>						
<i>Met</i>	<i>0.10</i>	<i>19</i>						
<i>GABARAPL2</i>	<i>0.10</i>	<i>18</i>						
<i>AURKA</i>	<i>0.10</i>	<i>17</i>						
<i>MIATNB</i>	<i>0.10</i>	<i>16</i>						
<i>ALAS1</i>	<i>0.09</i>	<i>15</i>						
<i>PTN</i>	<i>0.09</i>	<i>14</i>						
<i>STEAP4</i>	<i>0.09</i>	<i>13</i>						
<i>GAPDH</i>	<i>0.09</i>	<i>12</i>						
<i>TERF2IP</i>	<i>0.08</i>	<i>11</i>						
<i>IFT57</i>	<i>0.08</i>	<i>10</i>						
<i>MIR4435 1HG</i>	<i>0.08</i>	<i>9</i>						
<i>TBP</i>	<i>0.08</i>	<i>8</i>						
<i>BRAF</i>	<i>0.07</i>	<i>7</i>						
<i>BTG2</i>	<i>0.07</i>	<i>6</i>						
<i>CAMK2N2</i>	<i>0.07</i>	<i>5</i>						
<i>ARHGEF25</i>	<i>0.07</i>	<i>4</i>						
<i>NEAT1</i>	<i>0.06</i>	<i>3</i>						

9: APPENDICES

<i>All Transcripts (n = 166)</i>			<i>Transcripts identified by glm (n = 36)</i>			<i>Transcripts identified by polr (n = 20)</i>		
<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>	<i>Transcript</i>	<i>IncNodePurity</i>	<i>Rank</i>
<i>PTPRC</i>	<i>0.05</i>	<i>2</i>						
<i>SIRT1</i>	<i>0.05</i>	<i>1</i>						



9: APPENDICES

**Supplementary Table 48** Random Forest results for CB, low-risk, intermediate-risk and high-risk cancer trend using the *RPLP2* and *TWIST1* normalised data.

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 19)</i>			<i>Transcripts identified by polr (n = 15)</i>		
<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>
<i>PCA3</i>	<b>1.20</b>	<b>167</b>	<i>HOXC6</i>	<b>3.24</b>	<b>19</b>	<i>HOXC6</i>	<b>3.53</b>	<b>15</b>
<i>HOXC6</i>	<b>0.81</b>	<b>166</b>	<i>NAALADL2</i>	<b>2.78</b>	<b>18</b>	<i>NAALADL2</i>	<b>3.35</b>	<b>14</b>
<i>CP</i>	<b>0.79</b>	<b>165</b>	<i>ERG3' exons 6-7</i>	<b>2.75</b>	<b>17</b>	<i>TMPRSS2:ERG fusion</i>	<b>2.84</b>	<b>13</b>
<i>PALM3</i>	<b>0.70</b>	<b>164</b>	<i>PALM3</i>	<b>2.47</b>	<b>16</b>	<i>UPK2</i>	<b>2.76</b>	<b>12</b>
<i>ERG3' exons 6-7</i>	<b>0.67</b>	<b>163</b>	<i>UPK2</i>	<b>2.34</b>	<b>15</b>	<i>PALM3</i>	<b>2.75</b>	<b>11</b>
<i>NAALADL2</i>	<b>0.61</b>	<b>162</b>	<i>TMPRSS2:ERG fusion</i>	<b>2.28</b>	<b>14</b>	<i>ERG3' exons 6-7</i>	<b>2.64</b>	<b>10</b>
<i>UPK2</i>	<b>0.60</b>	<b>161</b>	<i>ST6GALNAC1</i>	<b>2.11</b>	<b>13</b>	<i>SIM2 long</i>	<b>2.53</b>	<b>9</b>
<i>TMPRSS2:ERG fusion</i>	<b>0.57</b>	<b>160</b>	<i>TMEM86A</i>	<b>2.06</b>	<b>12</b>	<i>TMEM86A</i>	<b>2.47</b>	<b>8</b>
<i>OR52A2</i>	<b>0.52</b>	<b>159</b>	<i>CADPS</i>	<b>2.00</b>	<b>11</b>	<i>TDRD</i>	<b>2.44</b>	<b>7</b>
<i>SPINK1</i>	<b>0.51</b>	<b>158</b>	<i>SIM2 long</i>	<b>1.97</b>	<b>10</b>	<i>ST6GALNAC1</i>	<b>2.43</b>	<b>6</b>
<i>VAX2</i>	<b>0.49</b>	<b>157</b>	<i>SERPINB5</i>	<b>1.79</b>	<b>9</b>	<i>EN2</i>	<b>2.30</b>	<b>5</b>
<i>TDRD</i>	<b>0.48</b>	<b>156</b>	<i>GJB1</i>	<b>1.76</b>	<b>8</b>	<i>SERPINB5</i>	<b>2.14</b>	<b>4</b>
<i>CKAP2L</i>	<b>0.44</b>	<b>155</b>	<i>TDRD</i>	<b>1.75</b>	<b>7</b>	<i>FOLH1</i>	<b>1.94</b>	<b>3</b>
<i>CADPS</i>	<b>0.43</b>	<b>154</b>	<i>LASS1</i>	<b>1.66</b>	<b>6</b>	<i>SLC43A1</i>	<b>1.85</b>	<b>2</b>
<i>HPN</i>	<b>0.42</b>	<b>153</b>	<i>CLIC2</i>	<b>1.47</b>	<b>5</b>	<i>MEX3A</i>	<b>1.68</b>	<b>1</b>
<i>AMH</i>	<b>0.42</b>	<b>152</b>	<i>SLC43A1</i>	<b>1.43</b>	<b>4</b>			
<i>HMBS</i>	<b>0.38</b>	<b>151</b>	<i>MSMB</i>	<b>1.40</b>	<b>3</b>			
<i>APOC1</i>	<b>0.37</b>	<b>150</b>	<i>FOLH1</i>	<b>1.22</b>	<b>2</b>			
<i>PPAP2A</i>	<b>0.37</b>	<b>149</b>	<i>MEX3A</i>	<b>1.18</b>	<b>1</b>			
<i>TMEM47</i>	<b>0.37</b>	<b>148</b>						
<i>LASS1</i>	<b>0.34</b>	<b>147</b>						
<i>SIM2 long</i>	<b>0.34</b>	<b>146</b>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 19)</i>			<i>Transcripts identified by polr (n = 15)</i>		
<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>
<i>ST6GALNAC1</i>	<b>0.33</b>	<b>145</b>						
<i>ISX</i>	<b>0.33</b>	<b>144</b>						
<i>DLX1</i>	<b>0.32</b>	<b>143</b>						
<i>MAPK8IP2</i>	<b>0.31</b>	<b>142</b>						
<i>AR.ex9</i>	<b>0.30</b>	<b>141</b>						
<i>MKi67</i>	<b>0.30</b>	<b>140</b>						
<i>TMEM45B</i>	<b>0.30</b>	<b>139</b>						
<i>TMEM86A</i>	<b>0.29</b>	<b>138</b>						
<i>PTN</i>	<b>0.29</b>	<b>137</b>						
<i>TERT</i>	<b>0.28</b>	<b>136</b>						
<i>EN2</i>	<b>0.28</b>	<b>135</b>						
<i>B4GALNT4</i>	<b>0.28</b>	<b>134</b>						
<i>CAMKK2</i>	<b>0.28</b>	<b>133</b>						
<i>ERG5'</i>	<b>0.27</b>	<b>132</b>						
<i>IGFBP3</i>	<b>0.27</b>	<b>131</b>						
<i>GCNT1</i>	<b>0.27</b>	<b>130</b>						
<i>MMP11</i>	<b>0.27</b>	<b>129</b>						
<i>AGR2</i>	<b>0.27</b>	<b>128</b>						
<i>MFSD2A</i>	<b>0.26</b>	<b>127</b>						
<i>SFRP4</i>	<b>0.26</b>	<b>126</b>						
<i>NKAIN1</i>	<b>0.26</b>	<b>125</b>						
<i>MDK</i>	<b>0.26</b>	<b>124</b>						
<i>DNAH5</i>	<b>0.26</b>	<b>123</b>						
<i>Timp4</i>	<b>0.25</b>	<b>122</b>						
<i>SLC4A1.S</i>	<b>0.24</b>	<b>121</b>						
<i>SPON2</i>	<b>0.24</b>	<b>120</b>						
<i>GJB1</i>	<b>0.24</b>	<b>119</b>						
<i>KLK4</i>	<b>0.24</b>	<b>118</b>						

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 19)</i>			<i>Transcripts identified by polr (n = 15)</i>		
<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>
<i>PPP1R12B</i>	<i>0.24</i>	<i>117</i>						
<i>AMACR</i>	<i>0.24</i>	<i>116</i>						
<i>SLC43A1</i>	<i>0.24</i>	<i>115</i>						
<i>TMCC2</i>	<i>0.24</i>	<i>114</i>						
<i>HOXC4</i>	<i>0.23</i>	<i>113</i>						
<i>ANKRD34B</i>	<i>0.23</i>	<i>112</i>						
<i>SERPINB5</i>	<i>0.23</i>	<i>111</i>						
<i>SChLAP1</i>	<i>0.23</i>	<i>110</i>						
<i>SLC12A1</i>	<i>0.23</i>	<i>109</i>						
<i>MMP25</i>	<i>0.23</i>	<i>108</i>						
<i>CLU</i>	<i>0.23</i>	<i>107</i>						
<i>TWIST1</i>	<i>0.23</i>	<i>106</i>						
<i>MYOF</i>	<i>0.22</i>	<i>105</i>						
<i>Met</i>	<i>0.22</i>	<i>104</i>						
<i>MARCH5</i>	<i>0.22</i>	<i>103</i>						
<i>MIR146A</i>	<i>0.22</i>	<i>102</i>						
<i>FOLH1</i>	<i>0.21</i>	<i>101</i>						
<i>CCDC88B</i>	<i>0.21</i>	<i>100</i>						
<i>COL9A2</i>	<i>0.21</i>	<i>99</i>						
<i>HIST1H2BG</i>	<i>0.20</i>	<i>98</i>						
<i>MNX1</i>	<i>0.20</i>	<i>97</i>						
<i>PCSK6</i>	<i>0.20</i>	<i>96</i>						
<i>AATF</i>	<i>0.20</i>	<i>95</i>						
<i>SMIM1</i>	<i>0.20</i>	<i>94</i>						
<i>PDLIM5</i>	<i>0.20</i>	<i>93</i>						
<i>HPRT</i>	<i>0.20</i>	<i>92</i>						
<i>ACTR5</i>	<i>0.19</i>	<i>91</i>						
<i>KLK3 exons 2-3</i>	<i>0.19</i>	<i>90</i>						

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 19)</i>			<i>Transcripts identified by polr (n = 15)</i>		
<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>
<i>HIST1H2BF</i>	<b>0.19</b>	<b>89</b>						
<i>ZNF577</i>	<b>0.19</b>	<b>88</b>						
<i>SRSF3</i>	<b>0.19</b>	<b>87</b>						
<i>ANPEP</i>	<b>0.19</b>	<b>86</b>						
<i>CLIC2</i>	<b>0.18</b>	<b>85</b>						
<i>MAK</i>	<b>0.18</b>	<b>84</b>						
<i>RIOK3</i>	<b>0.18</b>	<b>83</b>						
<i>SIRT1</i>	<b>0.18</b>	<b>82</b>						
<i>SMAP1 exons 7-8</i>	<b>0.18</b>	<b>81</b>						
<i>VPS13A</i>	<b>0.18</b>	<b>80</b>						
<i>PPF1A2</i>	<b>0.18</b>	<b>79</b>						
<i>ERG3' exons 4-5</i>	<b>0.17</b>	<b>78</b>						
<i>IMPDH2</i>	<b>0.17</b>	<b>77</b>						
<i>IFT57</i>	<b>0.17</b>	<b>76</b>						
<i>GOLM1</i>	<b>0.17</b>	<b>75</b>						
<i>LBH</i>	<b>0.17</b>	<b>74</b>						
<i>TFDP1</i>	<b>0.17</b>	<b>73</b>						
<i>CDKN3</i>	<b>0.17</b>	<b>72</b>						
<i>ITGBL1</i>	<b>0.17</b>	<b>71</b>						
<i>RP11_97O12.7</i>	<b>0.17</b>	<b>70</b>						
<i>BTG2</i>	<b>0.17</b>	<b>69</b>						
<i>CACNA1D</i>	<b>0.16</b>	<b>68</b>						
<i>HIST1H1C</i>	<b>0.16</b>	<b>67</b>						
<i>MIC1</i>	<b>0.16</b>	<b>66</b>						
<i>CASKIN1</i>	<b>0.16</b>	<b>65</b>						
<i>CDC37L1</i>	<b>0.16</b>	<b>64</b>						
<i>PECI</i>	<b>0.16</b>	<b>63</b>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 19)</i>			<i>Transcripts identified by polr (n = 15)</i>		
<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>
<i>MMP26</i>	<b>0.16</b>	<b>62</b>						
<i>MCTP1</i>	<b>0.16</b>	<b>61</b>						
<i>MGAT5B</i>	<b>0.15</b>	<b>60</b>						
<i>HIST3H2A</i>	<b>0.15</b>	<b>59</b>						
<i>TRPM4</i>	<b>0.15</b>	<b>58</b>						
<i>HIST1H1E</i>	<b>0.15</b>	<b>57</b>						
<i>RNF157</i>	<b>0.15</b>	<b>56</b>						
<i>ARHGEF25</i>	<b>0.15</b>	<b>55</b>						
<i>SNORA20</i>	<b>0.14</b>	<b>54</b>						
<i>STEAP2</i>	<b>0.14</b>	<b>53</b>						
<i>MEX3A</i>	<b>0.14</b>	<b>52</b>						
<i>CD10</i>	<b>0.14</b>	<b>51</b>						
<i>RAB17</i>	<b>0.14</b>	<b>50</b>						
<i>MCM7</i>	<b>0.14</b>	<b>49</b>						
<i>PTPRC</i>	<b>0.14</b>	<b>48</b>						
<i>PSTPIP1</i>	<b>0.14</b>	<b>47</b>						
<i>SULF2</i>	<b>0.14</b>	<b>46</b>						
<i>SSTR1</i>	<b>0.14</b>	<b>45</b>						
<i>SACMIL</i>	<b>0.14</b>	<b>44</b>						
<i>RPLP2</i>	<b>0.13</b>	<b>43</b>						
<i>KLK3 exons 1-2</i>	<b>0.13</b>	<b>42</b>						
<i>KLK2</i>	<b>0.13</b>	<b>41</b>						
<i>P712P</i>	<b>0.13</b>	<b>40</b>						
<i>SIM2 short</i>	<b>0.13</b>	<b>39</b>						
<i>MSMB</i>	<b>0.13</b>	<b>38</b>						
<i>SEC61A1</i>	<b>0.13</b>	<b>37</b>						
<i>AR.ex4 8</i>	<b>0.13</b>	<b>36</b>						
<i>SULT1A1</i>	<b>0.13</b>	<b>35</b>						

## 9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 19)</i>			<i>Transcripts identified by polr (n = 15)</i>		
<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>
<i>SSPO</i>	<b>0.13</b>	<b>34</b>						
<i>OGT</i>	<b>0.13</b>	<b>33</b>						
<i>ALAS1</i>	<b>0.12</b>	<b>32</b>						
<i>RPL23AP53</i>	<b>0.12</b>	<b>31</b>						
<i>STEAP4</i>	<b>0.12</b>	<b>30</b>						
<i>SYNM</i>	<b>0.12</b>	<b>29</b>						
<i>COL10A1</i>	<b>0.12</b>	<b>28</b>						
<i>AURKA</i>	<b>0.12</b>	<b>27</b>						
<i>ABCB9</i>	<b>0.12</b>	<b>26</b>						
<i>NEAT1</i>	<b>0.12</b>	<b>25</b>						
<i>PVT1</i>	<b>0.12</b>	<b>24</b>						
<i>RPS11</i>	<b>0.12</b>	<b>23</b>						
<i>DPP4</i>	<b>0.12</b>	<b>22</b>						
<i>SNCA</i>	<b>0.12</b>	<b>21</b>						
<i>CAMK2N2</i>	<b>0.11</b>	<b>20</b>						
<i>STOM</i>	<b>0.10</b>	<b>19</b>						
<i>RPL18A</i>	<b>0.10</b>	<b>18</b>						
<i>MED4</i>	<b>0.10</b>	<b>17</b>						
<i>GABARAPL2</i>	<b>0.10</b>	<b>16</b>						
<i>RPS10</i>	<b>0.10</b>	<b>15</b>						
<i>FDPS</i>	<b>0.10</b>	<b>14</b>						
<i>CDC20</i>	<b>0.10</b>	<b>13</b>						
<i>MXI1</i>	<b>0.10</b>	<b>12</b>						
<i>ITPR1</i>	<b>0.09</b>	<b>11</b>						
<i>TBP</i>	<b>0.09</b>	<b>10</b>						
<i>MIR4435 1HG</i>	<b>0.09</b>	<b>9</b>						
<i>TERF2IP</i>	<b>0.09</b>	<b>8</b>						
<i>BRAF</i>	<b>0.09</b>	<b>7</b>						

9: APPENDICES

<i>All Transcripts (n = 167)</i>			<i>Transcripts identified by glm (n = 19)</i>			<i>Transcripts identified by polr (n = 15)</i>		
<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>	<i>Transcript</i>	<i>MeanDecreaseGini</i>	<i>Rank</i>
<i>MIATNB</i>	<b>0.09</b>	<b>6</b>						
<i>NLRP3</i>	<b>0.08</b>	<b>5</b>						
<i>EIF2D</i>	<b>0.08</b>	<b>4</b>						
<i>GAPDH</i>	<b>0.07</b>	<b>3</b>						
<i>B2M</i>	<b>0.07</b>	<b>2</b>						
<i>MEMO1</i>	<b>0.06</b>	<b>1</b>						

**6.21 Cell vs EV fraction**

Supplementary Table 49 The 129 transcripts that are significantly (post multiple testing correction) different between the cell and microvesicular fraction.

<i>Transcript</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>Log2 Fold Change</i>
<i>NEAT1</i>	1.77E-16	2.94E-14	-0.88
<i>PTPRC</i>	1.77E-16	2.94E-14	-1.97
<i>MMP25</i>	2.17E-16	3.58E-14	-1.43
<i>SULF2</i>	2.48E-16	4.07E-14	-1.69
<i>HIST1H1C</i>	2.83E-16	4.59E-14	0.25
<i>MCTP1</i>	2.83E-16	4.59E-14	-1.09
<i>IFT57</i>	2.93E-16	4.66E-14	0.43
<i>CCDC88B</i>	2.93E-16	4.66E-14	-1.27
<i>STOM</i>	2.93E-16	4.66E-14	-1.73
<i>MFSD2A</i>	3.24E-16	5.12E-14	-1.66
<i>B2M</i>	3.46E-16	5.44E-14	-0.36
<i>PSTPIP1</i>	3.96E-16	6.17E-14	-1.44
<i>APOC1</i>	6.72E-16	1.03E-13	-1.07
<i>NLRP3</i>	6.72E-16	1.03E-13	-1.64
<i>MIR4435 IHG</i>	6.94E-16	1.06E-13	-0.43
<i>MSMB</i>	7.18E-16	1.09E-13	0.29
<i>KLK2</i>	8.19E-16	1.24E-13	0.55
<i>AR.ex4 8</i>	9.33E-16	1.40E-13	0.51
<i>KLK3.ex2 3</i>	1.03E-15	1.53E-13	0.51
<i>KLK4</i>	1.10E-15	1.63E-13	0.41
<i>PTN</i>	1.17E-15	1.73E-13	0.76
<i>BTG2</i>	1.34E-15	1.95E-13	-0.34
<i>DPP4</i>	1.52E-15	2.20E-13	0.47
<i>CLIC2</i>	1.52E-15	2.20E-13	-1.35
<i>STEAP2</i>	1.57E-15	2.25E-13	0.52
<i>MIR146A</i>	1.74E-15	2.47E-13	-0.96
<i>PECI</i>	2.04E-15	2.88E-13	0.30
<i>IMPDH2</i>	2.65E-15	3.70E-13	0.30
<i>RPLP2</i>	3.77E-15	5.24E-13	0.11
<i>P712P</i>	6.29E-15	8.69E-13	0.67
<i>TWIST1</i>	8.38E-15	1.15E-12	0.34
<i>RP11 97O12.7</i>	8.65E-15	1.18E-12	0.22
<i>RPS11</i>	9.51E-15	1.27E-12	0.10
<i>TMEM86A</i>	9.51E-15	1.27E-12	-1.19
<i>MAK</i>	1.39E-14	1.85E-12	-1.15
<i>ZNF577</i>	1.43E-14	1.89E-12	0.38
<i>PPAP2A</i>	1.68E-14	2.20E-12	0.29
<i>HIST1H2BF</i>	2.16E-14	2.80E-12	0.24
<i>TERT</i>	2.22E-14	2.87E-12	0.51
<i>05-Mar</i>	3.04E-14	3.89E-12	0.26
<i>PCA3</i>	7.66E-14	9.73E-12	0.58
<i>SERPINB5</i>	9.78E-14	1.23E-11	0.79
<i>NKAIN1</i>	1.07E-13	1.34E-11	0.57



## 9: APPENDICES

<i>Transcript</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>Log2 Fold Change</i>
<i>SMIMI</i>	<i>1.69E-13</i>	<i>2.09E-11</i>	<i>0.62</i>
<i>SSPO</i>	<i>1.85E-13</i>	<i>2.27E-11</i>	<i>0.38</i>
<i>FOLH1</i>	<i>2.28E-13</i>	<i>2.78E-11</i>	<i>0.51</i>
<i>RPL18A</i>	<i>2.57E-13</i>	<i>3.11E-11</i>	<i>0.17</i>
<i>FDPS</i>	<i>3.27E-13</i>	<i>3.92E-11</i>	<i>0.19</i>
<i>AMACR</i>	<i>3.37E-13</i>	<i>4.01E-11</i>	<i>0.50</i>
<i>OR52A2</i>	<i>3.57E-13</i>	<i>4.22E-11</i>	<i>0.77</i>
<i>MIC1</i>	<i>4.15E-13</i>	<i>4.85E-11</i>	<i>-0.67</i>
<i>GABARAPL2</i>	<i>6.86E-13</i>	<i>7.95E-11</i>	<i>0.15</i>
<i>PDLIM5</i>	<i>7.06E-13</i>	<i>8.12E-11</i>	<i>0.19</i>
<i>RPS10</i>	<i>7.27E-13</i>	<i>8.29E-11</i>	<i>0.14</i>
<i>KLK3.ex1 2</i>	<i>7.71E-13</i>	<i>8.71E-11</i>	<i>0.46</i>
<i>PPFIA2</i>	<i>2.07E-12</i>	<i>2.32E-10</i>	<i>0.65</i>
<i>SEC61A1</i>	<i>2.40E-12</i>	<i>2.66E-10</i>	<i>-0.46</i>
<i>MNX1</i>	<i>3.48E-12</i>	<i>3.82E-10</i>	<i>0.40</i>
<i>CD10</i>	<i>4.24E-12</i>	<i>4.63E-10</i>	<i>0.30</i>
<i>NAALADL2</i>	<i>5.18E-12</i>	<i>5.59E-10</i>	<i>0.38</i>
<i>CAMK2N2</i>	<i>6.68E-12</i>	<i>7.14E-10</i>	<i>0.53</i>
<i>TFDP1</i>	<i>9.35E-12</i>	<i>9.91E-10</i>	<i>0.18</i>
<i>Met</i>	<i>9.62E-12</i>	<i>1.01E-09</i>	<i>-0.88</i>
<i>SIM2.long</i>	<i>1.68E-11</i>	<i>1.74E-09</i>	<i>0.65</i>
<i>COL10A1</i>	<i>2.09E-11</i>	<i>2.15E-09</i>	<i>-0.59</i>
<i>SSTR1</i>	<i>3.16E-11</i>	<i>3.22E-09</i>	<i>0.24</i>
<i>CP</i>	<i>4.49E-11</i>	<i>4.54E-09</i>	<i>-0.99</i>
<i>PCSK6</i>	<i>6.05E-11</i>	<i>6.05E-09</i>	<i>0.40</i>
<i>Timp4</i>	<i>6.74E-11</i>	<i>6.67E-09</i>	<i>0.61</i>
<i>VAX2</i>	<i>1.09E-10</i>	<i>1.06E-08</i>	<i>0.36</i>
<i>CACNA1D</i>	<i>1.09E-10</i>	<i>1.06E-08</i>	<i>0.19</i>
<i>HOXC6</i>	<i>1.12E-10</i>	<i>1.07E-08</i>	<i>0.81</i>
<i>SPON2</i>	<i>2.40E-10</i>	<i>2.28E-08</i>	<i>0.34</i>
<i>AMH</i>	<i>2.60E-10</i>	<i>2.44E-08</i>	<i>0.30</i>
<i>ARHGEF25</i>	<i>6.56E-10</i>	<i>6.10E-08</i>	<i>0.58</i>
<i>EIF2D</i>	<i>6.90E-10</i>	<i>6.35E-08</i>	<i>0.12</i>
<i>SchLAPI</i>	<i>8.67E-10</i>	<i>7.89E-08</i>	<i>0.67</i>
<i>GJB1</i>	<i>1.01E-09</i>	<i>9.08E-08</i>	<i>0.49</i>
<i>AURKA</i>	<i>1.17E-09</i>	<i>1.04E-07</i>	<i>-0.35</i>
<i>HIST3H2A</i>	<i>2.08E-09</i>	<i>1.83E-07</i>	<i>0.37</i>
<i>RAB17</i>	<i>2.47E-09</i>	<i>2.15E-07</i>	<i>0.38</i>
<i>HMBS</i>	<i>2.80E-09</i>	<i>2.41E-07</i>	<i>0.25</i>
<i>MKi67</i>	<i>3.85E-09</i>	<i>3.27E-07</i>	<i>-1.18</i>
<i>DNAH5</i>	<i>5.02E-09</i>	<i>4.22E-07</i>	<i>0.57</i>
<i>CKAP2L</i>	<i>8.93E-09</i>	<i>7.41E-07</i>	<i>-0.50</i>
<i>CASKIN1</i>	<i>1.05E-08</i>	<i>8.65E-07</i>	<i>0.24</i>
<i>SULT1A1</i>	<i>1.08E-08</i>	<i>8.75E-07</i>	<i>-0.18</i>
<i>MXII</i>	<i>1.57E-08</i>	<i>1.24E-06</i>	<i>0.13</i>
<i>ITPR1</i>	<i>1.57E-08</i>	<i>1.24E-06</i>	<i>-0.14</i>
<i>MMP11</i>	<i>1.94E-08</i>	<i>1.51E-06</i>	<i>0.30</i>
<i>HPRT</i>	<i>3.16E-08</i>	<i>2.43E-06</i>	<i>0.18</i>
<i>SIM2.short</i>	<i>3.30E-08</i>	<i>2.51E-06</i>	<i>0.35</i>
<i>PALM3</i>	<i>3.62E-08</i>	<i>2.72E-06</i>	<i>0.31</i>

## 9: APPENDICES

<i>Transcript</i>	<i>p-value</i>	<i>Adjusted p-value</i>	<i>Log2 Fold Change</i>
<i>AGR2</i>	<b>4.06E-08</b>	<b>3.00E-06</b>	<b>0.32</b>
<i>SYNM</i>	<b>4.87E-08</b>	<b>3.55E-06</b>	<b>0.54</b>
<i>MDK</i>	<b>1.42E-07</b>	<b>1.02E-05</b>	<b>0.21</b>
<i>EN2</i>	<b>2.66E-07</b>	<b>1.89E-05</b>	<b>0.36</b>
<i>MED4</i>	<b>3.15E-07</b>	<b>2.21E-05</b>	<b>0.09</b>
<i>RNF157</i>	<b>3.58E-07</b>	<b>2.47E-05</b>	<b>0.58</b>
<i>MGAT5B</i>	<b>7.01E-07</b>	<b>4.76E-05</b>	<b>0.28</b>
<i>LBH</i>	<b>1.01E-06</b>	<b>6.80E-05</b>	<b>0.28</b>
<i>IGFBP3</i>	<b>1.06E-06</b>	<b>6.98E-05</b>	<b>-0.56</b>
<i>TMEM45B</i>	<b>1.30E-06</b>	<b>8.42E-05</b>	<b>-0.28</b>
<i>HOXC4</i>	<b>1.82E-06</b>	<b>0.0001</b>	<b>0.37</b>
<i>CLU</i>	<b>2.99E-06</b>	<b>0.0002</b>	<b>0.61</b>
<i>SNCA</i>	<b>2.99E-06</b>	<b>0.0002</b>	<b>0.15</b>
<i>MYOF</i>	<b>4.76E-06</b>	<b>0.0003</b>	<b>0.12</b>
<i>CDC37L1</i>	<b>5.24E-06</b>	<b>0.0003</b>	<b>0.11</b>
<i>GOLM1</i>	<b>7.66E-06</b>	<b>0.0005</b>	<b>0.39</b>
<i>SACMIL</i>	<b>1.11E-05</b>	<b>0.0006</b>	<b>0.11</b>
<i>SFRP4</i>	<b>1.27E-05</b>	<b>0.0007</b>	<b>0.35</b>
<i>ERG3prime.ex4 5</i>	<b>2.31E-05</b>	<b>0.001</b>	<b>0.64</b>
<i>LASS1</i>	<b>2.31E-05</b>	<b>0.001</b>	<b>-0.46</b>
<i>B4GALNT4</i>	<b>2.71E-05</b>	<b>0.001</b>	<b>-0.45</b>
<i>MEX3A</i>	<b>2.86E-05</b>	<b>0.002</b>	<b>0.39</b>
<i>STEAP4</i>	<b>4.06E-05</b>	<b>0.002</b>	<b>-0.11</b>
<i>HPN</i>	<b>4.66E-05</b>	<b>0.002</b>	<b>0.21</b>
<i>MAPK8IP2</i>	<b>4.74E-05</b>	<b>0.002</b>	<b>-0.45</b>
<i>TRPM4</i>	<b>7.91E-05</b>	<b>0.004</b>	<b>0.37</b>
<i>ANPEP</i>	<b>9.04E-05</b>	<b>0.004</b>	<b>-0.18</b>
<i>TERF2IP</i>	<b>9.50E-05</b>	<b>0.004</b>	<b>0.04</b>
<i>SRSF3</i>	<b>9.99E-05</b>	<b>0.005</b>	<b>-0.19</b>
<i>HIST1H1E</i>	<b>0.0002</b>	<b>0.007</b>	<b>0.09</b>
<i>ANKRD34B</i>	<b>0.0003</b>	<b>0.012</b>	<b>-0.35</b>
<i>PPP1R12B</i>	<b>0.0003</b>	<b>0.015</b>	<b>0.12</b>
<i>HIST1H2BG</i>	<b>0.0005</b>	<b>0.021</b>	<b>0.15</b>
<i>CDC20</i>	<b>0.0009</b>	<b>0.039</b>	<b>0.22</b>
<i>TMEM47</i>	<b>0.0011</b>	<b>0.043</b>	<b>0.61</b>
<i>SMAP1.ex7 8</i>	<b>0.0013</b>	<b>0.049</b>	<b>0.14</b>