

1 **A Novel Stratification Framework for Predicting Outcome in Patients**  
2 **with Prostate Cancer**

3 Bogdan-Alexandru Luca<sup>1,2,\*</sup>, Vincent Moulton<sup>2,#</sup>, Christopher Ellis<sup>1,2</sup>, Dylan R Edwards<sup>1</sup>,  
4 Colin Campbell<sup>3</sup>, Rosalin A Cooper<sup>4</sup>, Jeremy Clark<sup>1</sup>, Daniel S Brewer<sup>1,5,\*,#</sup>, & Colin S  
5 Cooper<sup>1,#,§</sup>

6  
7 <sup>1</sup> Norwich Medical School, University of East Anglia, Norwich Research Park, Norwich,  
8 Norfolk, UK.

9 <sup>2</sup> School of Computing Sciences, University of East Anglia, Norwich Research Park,  
10 Norwich, Norfolk, UK.

11 <sup>3</sup> Intelligent Systems Laboratory, University of Bristol, Bristol, UK.

12 <sup>4</sup> Department of Pathology, University Hospital Southampton NHS Foundation Trust,  
13 Southampton, UK.

14 <sup>5</sup> The Earlham Institute, Norwich Research Park, Norwich, Norfolk, UK.

15 \* Contributed equally to this work.

16 # Jointly supervised this work.

17

18 <sup>§</sup> **Corresponding author.**

19 Professor Colin Cooper (colin.cooper@uea.ac.uk), University of East Anglia, Norwich  
20 Research Park, Norwich, NR4 7UG, UK

21

22 **Running title: prostate cancer stratification**

1 **Abstract**

2 **Background**

3 Unsupervised learning methods such as Hierarchical Cluster Analysis are commonly used  
4 for the analysis of genomic platform data. Unfortunately, such approaches ignore the well  
5 documented heterogeneous composition of prostate cancer samples. Our aim is to use more  
6 sophisticated analytical approaches to deconvolute the structure of prostate cancer  
7 transcriptome data providing novel clinically actionable information for this disease.

8 **Methods**

9 We apply an unsupervised model called Latent Process Decomposition (LPD), which can  
10 handle heterogeneity within individual cancer samples, to genome-wide expression data  
11 from eight prostate cancer clinical series including 1,785 malignant samples with the clinical  
12 endpoints of PSA failure and metastasis.

13 **Results**

14 We show that PSA failure is correlated with the level of an expression signature called  
15 DESNT (HR = 1.52, 95% CI = [1.36, 1.7],  $P = 9.0 \times 10^{-14}$ , Cox model) and that patients with a  
16 majority DESNT signature have an increased metastatic risk ( $X^2$ -test,  $P = 0.0017$ , and  
17  $P = 0.0019$ ). Additionally, we develop a stratification framework that incorporates DESNT  
18 and identifies three novel molecular subtypes of prostate cancer.

19 **Conclusions**

20 These results highlight the importance of using more complex approaches for the analysis of  
21 genomic data, may assist drug targeting, and have allowed the construction of a nomogram  
22 combining DESNT with other clinical factors for use in clinical management.

23

24

1

## 2 **Background**

3 Driven by technological advances and decreased costs, a plethora of genomic datasets now  
4 exist. This is illustrated by the availability of expression data from over 1.3 million samples  
5 from the Gene Expression Omnibus [1] and DNA sequence data on 25,000 cases from the  
6 International Cancer Genome Consortium [2]. Such datasets have been used as the raw  
7 material for the discovery of disease sub-classes using a variety of mathematical approaches.  
8 Hierarchical clustering, k-means clustering, and self-organising maps have been applied to  
9 expression datasets leading, for example, to the discovery of five molecular breast cancer  
10 types (Basal, Luminal A, Luminal B, ERBB2-overexpressing, and Normal-like) [3]. The  
11 inherent shortcoming of this type of approach is the implicit assumption of sample  
12 assignment to a particular cluster or group. Such analyses are in complete contrast to the  
13 well documented heterogeneous composition of most individual cancer samples [4,5].

14

15 Unsupervised analysis of prostate cancer transcriptome profiles using the above approaches  
16 have failed to identify robust disease categories that have distinct clinical outcomes [6,7].  
17 Noting that prostate cancer samples derived from genome wide studies frequently harbour  
18 multiple cancer lineages, and can have intra-tumour variations in genetic compositions [8-  
19 10], we applied an unsupervised learning method called Latent Process Decomposition  
20 (LPD) [11] that can take into account the issue of heterogeneity of composition within  
21 individual cancer samples. By heterogeneity we mean that an individual cancer sample can  
22 be made up of several different components that each has distinct properties We had  
23 previously used Latent Process Decomposition: (i) to confirm the presence of the basal and

1 ERBB2 overexpressing subtypes in breast cancer transcriptome datasets [12]; (ii) to  
2 demonstrate that data from the MammaPrint breast cancer recurrence assay would be  
3 optimally analyzed using four separate prognostic categories [12]; and (iii) to show that  
4 patients with advanced prostate cancer can be stratified into two clinically distinct categories  
5 based on expression profiles in blood [13]. LPD (closely related to Latent Dirichlet  
6 Allocation) is a mixed membership model in which the expression profile for a single cancer  
7 is represented as a combination of underlying latent (hidden) signatures. Each latent  
8 signature has a representative gene expression pattern. A given sample can be represented  
9 over a number of these underlying functional states, or just one such state. The appropriate  
10 number of signatures to use is determined using the LPD algorithm by maximising the  
11 probability of the model given the data.

12

13 The application of LPD to prostate cancer transcriptome datasets led to the discovery of an  
14 expression pattern, called DESNT, that was observed in all datasets examined [14]. Cancer  
15 samples were assigned as DESNT when this pattern was more common than any other  
16 signature, and this designation was associated with poor outcomes independently of other  
17 clinical parameters including Gleason, Clinical stage and PSA. In the current paper we test  
18 whether the presence of even a small proportion of the DESNT cancer signature confers  
19 poor outcome and use LPD to develop a new prostate cancer stratification framework.

20

## 21 **Materials and Methods**

### 22 **Transcriptome datasets**

23 Eight publically available transcriptome microarray datasets derived from prostatectomy  
24 samples from men with prostate cancer were used and are referred to as: Memorial Sloan  
25 Kettering Cancer Centre (MSKCC) [7], CancerMap [14], CamCap [6], Stephenson [15], TCGA

1 [16], Klein [17], Erho [18], and Karnes [19]. There were 1785 samples from primary  
2 malignant tissue and 173 from normal tissue (Table 1). MSKCC also had data from 19  
3 metastatic cancer samples. The CamCap dataset was produced by combining two Illumina  
4 HumanHT-12 V4.0 expression beadchip datasets (GEO: [GSE70768](#) and [GSE70769](#)) obtained  
5 from two prostatectomy series (Cambridge and Stockholm) [6]. The original CamCap [6]  
6 and CancerMap [14] datasets have 40 patients in common and thus 20 of the common  
7 samples were excluded at random from each dataset. Each Affymetrix Exon microarray  
8 dataset was normalised using the RMA algorithm [20] implemented in the Affymetrix  
9 Expression Console software. For CamCap and Stephenson previous normalised values  
10 were used. For the TCGA dataset, the counts per gene previously calculated were used [16]  
11 and transformed using the variance stabilising transformation implemented in the DESeq2  
12 package[21]. For the CamCap and CancerMap datasets the ERG gene alterations had been  
13 scored by fluorescence in situ hybridization [6,14]. Only probes corresponding to genes  
14 measured by all platforms were retained. The ComBat algorithm from the sva R package  
15 and quantile transformation, was used to mitigate series-specific effects. Flow diagrams  
16 presenting each of the analyses performed in the current study, with the datasets used, are  
17 shown in the Supplementary Materials. The ethical approvals obtained for each dataset are  
18 listed in the original publications.

19

## 20 **Latent Process Decomposition (LPD)**

21 LPD [11,12] is an unsupervised Bayesian approach which breaks down (decomposes) each  
22 sample into component sub-elements (signatures). Each signature is a representative gene  
23 expression pattern. LPD is able to classify complex data based on the relative representation  
24 of these signatures in each sample. LPD can objectively assess the most likely number of  
25 signatures. We assessed the hold-out validation log-likelihood of the data computed at

1 various number of signatures and used a combination of both the uniform (equivalent to a  
2 maximum likelihood approach) and non-uniform (missed approach point) priors to choose  
3 the number of signatures. For input, each dataset was reduced to probes that detect the 500  
4 genes with the greatest variance across the MSKCC dataset. For robustness, LPD is run 100  
5 times with different seeds, for each dataset. Out of the 100 runs we selected the run with the  
6 survival log-rank  $p$ -value closest to the mode as a representative run that was used for  
7 subsequent analysis.

8

### 9 **OAS-LPD (One Added Sample LPD)**

10 The OAS-LPD algorithm is a modified version of the LPD algorithm in which new sample(s)  
11 are decomposed into LPD signatures, without retraining the model (i.e. without re-  
12 estimating the model parameters  $\mu_{gk}$ ,  $\sigma^2_{gk}$ , and  $\alpha$  in Rogers et al. [11]). Only the variational  
13 parameters  $Q_{kga}$  and  $\gamma_{ak}$ , corresponding to the new sample(s), are iteratively updated until  
14 convergence, according to Eq. (6) and Eq. (7) from Rogers et al. 2005 [11]. LPD as presented  
15 by Rogers et al. [11] was first applied to the MSKCC dataset of 131 cancer and 29 normal  
16 samples, as described above. The model parameters  $\mu_{gk}$ ,  $\sigma^2_{gk}$ , and  $\alpha$ , corresponding to the  
17 representative LPD run, were then used to classify additional expression profiles from all  
18 datasets, one sample at a time. A detailed description is provided in the Supplementary  
19 Methods.

20

### 21 **Statistical tests**

22 All statistical tests were performed in R version 3.3.1. For characterisation of signatures, each  
23 sample was assigned to the signature that had the largest gamma ( $\gamma$ ) value for that sample.

### 24 Correlations

1 Pearson correlations between the expression profiles between the MSKCC and CancerMap  
2 were calculated for each of the eight signatures: (i) for each gene we select one  
3 corresponding probe at random; (ii) for each probe we transformed its distribution across all  
4 samples to a standard normal distribution; (iii) the mean expression for each gene across the  
5 samples assigned to signature  $j$  (gene subgroup mean) in each dataset was determined; (iv)  
6 the Pearson's correlation between the gene subgroup mean expression profile in MSKCC vs  
7 the gene subgroup mean expression profile in CancerMap is calculated for each signature.

#### 8 Differentially expressed and methylated features

9 Differentially expressed probesets were identified for each signature using a moderated  $t$ -  
10 test implemented in the limma R package (Benjamin-Hocberg false discovery rate  $< 0.05$ ,  
11 differentially expressed in at least 50/100 runs; samples assigned to the signature vs the  
12 rest).

13 Thus differential methylation was assigned at the probe level. Hypo and hypermethylated  
14 genes that are predictive of transcription were identified using the methylMix R package  
15 (functionally differentially methylated in at least 50/100 runs) using genes that are found to  
16 be differentially expressed in that signature as input. Datasets where there were  $< 10$  samples  
17 assigned to a signature were removed from the identification of intersection genes for that  
18 signature.

#### 19 Survival analyses and nomogram

20 Survival analyses were performed using Cox proportional hazards models, the log-rank test,  
21 and Kaplan-Meier estimator, with biochemical recurrence after prostatectomy as the end  
22 point. For nomogram construction, the Cox proportional hazards model was fitted on the  
23 meta-dataset obtained by combining MSKCC, CancerMap, and Stephenson datasets, and  
24 validated on CamCap, using the rms R package. The Gleason grade was divided into  $< 7$ ,  
25  $3+4$ ,  $4+3$ ,  $> 7$ , the pathological stage in T1-T2 vs T3-T4, while DESNT percentage and PSA

1 were considered continuous covariates. The missing values for the predictors were imputed  
2 using the flexible additive models with predictive mean matching, implemented in the  
3 Hmisc R package. The linearity of the continuous covariates was assessed using the  
4 Martingale residuals [22]. The lack of collinearity between covariates was determined by  
5 calculating the variance inflation factors (VIF) (VIF values between 1.04 and 3.01) [23]. All  
6 covariates met the Cox proportional hazards assumption, as determined by the Schoenfeld  
7 residuals. The internal validation and calibration of the Cox model were performed by  
8 bootstrapping the training dataset 1,000 times. The calibration of the model was estimated  
9 by comparing the predicted and observed survival probabilities at five years. For comparing  
10 the discrimination accuracy of two non-nested Cox models the U-statistic calculated by the  
11 Hmisc rcorr.cens function was used.

#### 12 Detecting over-representation of genomic features

13 Mutated cancer genes identified by the Cancer Genome Atlas Research Network (2015) [16]  
14 were examined at the sample level. The under-/over-representation of these features in  
15 samples assigned to a particular LPD signature was determined using the  $\chi^2$  independence  
16 test.

#### 17 Pathway over-representation analysis and signature correlation analysis

18 The GO biological process annotations were tested for over-representation (or under-  
19 representation) in the lists of differentially expressed genes in each signature, using  
20 clusterProfiler version 3.4.4. For a given pathway and a given sample the pathway activation  
21 score was calculated as indicated in Levine et al. [24]. Using the complete combined dataset  
22 of all 8 datasets, Z-scores were calculated for each sample for each of the 17,697 MSigDB v6.0  
23 gene sets. These were correlated with DESNT  $\gamma$  values, and the top 20 sets with the highest  
24 absolute Pearson's correlation were selected. The resulting  $p$ -values from pathway over-  
25 representation analysis were adjusted for multiple testing using the false discovery rate.

1

## 2 **Results**

### 3 **Presence of DESNT signature as a continuous variable is associated with poor** 4 **clinical outcome**

5 In our previous studies, LPD detected between three and eight underlying signatures (also  
6 called processes) in expression microarray datasets collected from prostate cancer samples  
7 after prostatectomy [14]. Decomposition of the MSKCC dataset [7] gave eight signatures  
8 [14]. Fig. 1a illustrates the proportion of the DESNT expression signature identified in each  
9 MSKCC sample, with individual cancer samples being assigned as a “DESNT cancer” when  
10 the DESNT signature was the most abundant as shown in Fig. 1a and Fig. 1c. Based on PSA  
11 failure, patients with DESNT cancer always exhibited poorer outcome relative to other  
12 cancer samples in the same dataset [14]. The implication is that it is the presence of regions  
13 of cancer containing the DESNT signature conferred poor outcome. If this idea is correct, we  
14 would predict that cancer samples containing a smaller contribution of DESNT signature,  
15 such as those shown in Fig. 1b for the MSKCC dataset, should also exhibit poorer outcome.

16 To increase the power to test this prediction we combined transcriptome data from the  
17 MSKCC [7], CancerMap [14], Stephenson [15], and CamCap [6] studies (n=503). There was a  
18 significant association with PSA recurrence when the proportion of expression assigned to  
19 the DESNT signature was treated as a continuous variable (HR = 1.52, 95% CI=[1.36, 1.7],  
20  $P = 9.0 \times 10^{-14}$ , Cox proportional hazard regression model). Outcome became worse as the  
21 proportion of DESNT signature increased. For illustrative purposes cancer samples were  
22 divided into four groups based on the proportion of DESNT, with 47.4% of cancer samples  
23 contained at least some DESNT cancer (proportion greater than 0.001; Fig. 2a). PSA failure

1 free survival at 60 months is 82.5%, 67.4%, 59.5% and 44.9% for the proportion of DESNT  
2 signature being <0.001, 0.001 to 0.3, 0.3 to 0.6, and >0.6, respectively (Fig. 2b).

### 3 **Nomogram for DESNT predicting PSA failure**

4 The proportion of DESNT cancer was combined with other clinical variables (Gleason grade,  
5 PSA levels, pathological stage, and the surgical margins status) in a Cox proportional  
6 hazards model and fitted to a combined dataset of 318 cancer samples (MSKCC,  
7 CancerMap, and Stephenson); CamCap cancer samples (n=185) were used for external  
8 validation. The proportion of DESNT was an independent predictor of worse clinical  
9 outcome (HR = 1.33, 95% CI=[1.14, 1.56],  $P = 3.0 \times 10^{-4}$ ,) along with Gleason grade=4+3  
10 (HR = 2.43, 95% CI=[1.10, 5.37],  $P = 2.7 \times 10^{-2}$ ), Gleason grade>7 (HR = 5.05, 95% CI=[2.35,  
11 10.89],  $P < 1 \times 10^{-4}$ ), and positive surgical margins (HR = 1.65, 95% CI=[1.07, 2.56],  $P = 2.2 \times 10^{-2}$ )  
12 (Fig. S1: Supplementary Figure 1). PSA level and pathological stage were below the  
13 threshold of statistical significance ( $P = 0.09$ , HR = 1.14, 95% CI=[0.97, 1.34]) and ( $P = 0.055$ ,  
14 HR = 1.51, 95% CI=[0.99, 2.31]) respectively. At internal validation, the Cox model obtained  
15 a 1,000 bootstrap-corrected C-index of 0.747, and at external validation a C-index of 0.795.  
16 Using this model, a nomogram was constructed for use of DESNT cancer information in  
17 conjunction with clinical variables to predict the risk of biochemical recurrence at one, three,  
18 five, and, seven years following prostatectomy (Fig. 2c, Fig. S1).

### 19 **LPD algorithm for detecting the presence of DESNT cancer in individual samples**

20 The ability of LPD to detect structure is likely to be dependent on sample size, cohort  
21 composition, disease severity range and data quality. We observed optimal decompositions  
22 varying between three and eight underlying signatures in different datasets [14]. When we  
23 examined the two datasets that had an optimal eight underlying signatures (MSKCC and  
10

1 CancerMap) we noted a striking relationship: based on correlations of expression profiles;  
2 all eight of the LPD signatures appeared to be common (Fig. S2;  $R^2 > 0.5$ ). To provide a more  
3 consistent classification framework where the number of classes did not vary between  
4 datasets, we therefore used the MSKCC dataset and its decomposition into eight distinct  
5 signatures as a reference for identifying categories of prostate cancer type.

6 LPD is a computer intensive procedure and analyses can take days to run on a high-  
7 performance computing cluster. This would restrict ease of DESNT detection for clinical  
8 implementation. We therefore developed a variant of LPD called One Added Sample-LPD  
9 (OAS-LPD), where data from a single additional cancer sample could be decomposed into  
10 signatures, following normalisation, without repeating the entire LPD procedure. LPD  
11 model parameters [11] were first derived by decomposition of the MSKCC dataset into eight  
12 signatures. These signature parameters were then used as a framework for decomposition of  
13 additional data from single samples, selected in this case from a dataset, or in future from a  
14 patient undergoing assessment in the clinic. To test this procedure, we applied OAS-LPD  
15 individually to cancer samples from MSKCC, CancerMap, Stephenson, and CamCap (Fig.  
16 S3) and repeated Cox regression analysis and nomogram construction. Proportion of DESNT  
17 ( $P = 0.0011$ , HR = 1.53, 95% CI=[1.19, 1.98]), Gleason=4+3 ( $P = 0.0061$ , HR = 2.83, 95%  
18 CI=[1.35, 5.96]), Gleason>7 ( $P < 1 \times 10^{-4}$ , HR = 5.39, 95% CI=[2.54, 11.44]) and surgical margin  
19 status ( $P = 0.0015$ , HR = 2.00, 95% CI=[1.30, 3.07]) remained independent predictors of  
20 clinical outcome (Fig. S4). Notably the performance of the Cox model (internal validation C-  
21 index=0.742; external validation C-index=0.786) was not significantly different to that of the  
22 original separate dataset Cox model (train dataset  $Z = -0.65$ , two-tailed  $P = 0.52$ ; validation  
23 dataset  $Z = 0.89$ , two-tailed  $P = 0.38$ ; U-statistic) and the nomogram (Fig. S5) had almost an  
24 identical presentation of parameters to that shown in Fig. 2c. This observation is consistent

1 the high degree of correlation between LPD and OAS-LPD DESNT gamma values across the  
2 MSKCC, CancerMap, Stephenson, and CamCap datasets  $P= 2.39 \times 10^{-110}$ )

### 3 **New categories of prostate cancer**

4 We wished to determine whether LPD signatures were characterized by particular clinical or  
5 molecular features indicating that they represented distinct categories of prostate cancer.  
6 OAS-LPD using the MSKCC derived model of gene signatures was applied to all datasets  
7 ( $n=1958$ , Table 1) and each sample was assigned to the signature that was the most  
8 abundant. Samples from non-cancerous (benign) prostate tissue were more frequently  
9 assigned to LPD2, LPD4, and LPD8 than to the other groups ( $P < 0.05$ ,  $\chi^2$  test, Fig. S3, Table  
10 S1). When datasets with linked clinical data were combined (MSKCC, CancerMap,  
11 Stephenson, CamCap, Fig. 3a-c) primary cancers assigned to DESNT had worse outcome  
12 ( $P = 3.4 \times 10^{-14}$ , log-rank test, DESNT assigned samples vs the rest) while those assigned to  
13 LPD4 had improved outcome ( $P = 0.0081$ , log-rank test, LPD4 assigned samples vs the rest)  
14 as judged by PSA failure. Cancer samples with ERG-alterations assigned to signature LPD3  
15 also exhibited better outcome ( $P < 0.05$ ; log-rank test, comparison to all other ETS positive  
16 cancer samples) in all three datasets where ERG status was available (Fig. 4b-d).

17 To gain information about the new LPD categories we examined the distribution of genetic  
18 alterations in the decomposition of the TGCA dataset [16] (Fig. 4a), LPD3 cancer samples  
19 had over-representation of ETS and *PTEN* gene alterations, and under-representation of  
20 *CDH1* and *SPOP* gene alterations ( $P < 0.05$ ,  $\chi^2$  test, Table 2). LPD5 cancer samples exhibited  
21 exactly the reverse pattern of genetic alteration: there was under-repression of ETS and  
22 *PTEN* gene alterations and over-representation of *SPOP* and *CHD1* alterations (Table 2). The  
23 statistically different distribution of ETS-gene alterations in samples assigned to LPD3, and

1 LPD5, observed in the TCGA dataset were confirmed in the CamCap and CancerMap  
2 dataset (Table 2). In summary we have identified three additional prostate cancer categories  
3 that have altered genetic and/or clinical associations: LPD3, LPD4, and LPD5 (Fig. 5) and  
4 that may be relevant for drug targeting.

## 5 **Altered patterns of gene expression and DNA methylation**

6 We examined samples assigned to each OAS-LPD signature for genes with significantly  
7 altered expression levels in all eight datasets ( $P < 0.05$  after FDR correction, samples in LPD  
8 group vs all other LPD categories from the same dataset, Supplementary data 1). LPD3  
9 cancers samples exhibited seven commonly overexpressed genes including ERG, GHR, and  
10 HDAC1. Pathway analysis suggested the involvement of Stat3 gene signalling (Fig. S6a,  
11 Supplementary data 2). LPD5 exhibited 47 significantly overexpressed gene and 13 under-  
12 expressed genes. Many of the genes had established roles in fatty acid metabolism and the  
13 control of secretion (Fig. S6b). LPD6-cancers and LPD8 cancers had failed to exhibit  
14 statistically significant changes in genetic alteration or clinical outcome in the current study  
15 but did have characteristic altered patterns of gene expression (Fig. S6c,e). The five genes  
16 commonly overexpressed in LPD6 cancers suggested involvement in metal ion homeostasis.  
17 30 genes were overexpressed and 36 genes under expressed in in LPD8 cancers including  
18 several genes involved in extracellular matrix organisation. Cross referencing differential  
19 methylation data available for the TCGA dataset with gene associated with each LPD group  
20 indicated that many expression changes may be explained, at least in part, by changes in  
21 DNA methylation (Fig. 5, Fig. S7, Supplementary data 3).

## 22 **DESNT as a signature of metastasis**

1 The MSKCC study includes data from 19 metastatic cancer samples. For each metastatic  
2 sample, DESNT was the most abundant signature when OAS-LPD was applied (Fig. 3d).  
3 Two of the studied datasets (MSKCC and Erho) had publicly available annotations  
4 indicating that the patients from which primary cancer expression profiles were examined  
5 had progressed to develop metastasis after prostatectomy (Fig. S3). From nine cancer  
6 patients developing metastasis in the MSKCC dataset five occurred from samples in which  
7 the DESNT signature is most common ( $\chi^2$ -test,  $P = 0.0017$ ) and of 212 cancer patients  
8 developing metastases in the Erho dataset 50 were from DESNT cancers ( $\chi^2$ -test,  $P = 0.0019$ )  
9 (Fig. S8). From these studies we concluded that DESNT cancers have an increased risk of  
10 developing metastasis, consistent with the higher risk of PSA failure. For the Erho dataset,  
11 membership of LPD1 was associated with lower risk of metastasis ( $\chi^2$ -test,  $P = 0.026$ , Fig.  
12 S8).

13 To further investigate the underlying nature of DESNT cancer we used the transcriptome  
14 profile for each primary prostate cancer sample to investigate associations with the 17,697  
15 signatures and pathways annotated in the MSigDB database. The top 20 signatures where  
16 expression was associated with proportion of DESNT are shown in Table S2. The third most  
17 significant correlation was to genes downregulated in metastatic prostate cancer. This  
18 resulting data gives additional clues to the underlying biology of DESNT cancer including  
19 associations with genes altered in ductal breast cancer, in stem cells and during FGFR1  
20 signaling.

21

## 22 **Discussion**

23 We have confirmed a key prediction of the DESNT cancer model by demonstrating that the  
24 presence of a small proportion of the DESNT cancer signature confers poorer outcome. The

1 proportion of DESNT signature can be considered a continuous variable such that as DESNT  
2 cancer content increases outcome became worse. This observation led to the development of  
3 nomograms for estimating PSA failure at three years, five years, and seven years following  
4 prostatectomy. The result provides an extension of previous studies in which nomograms  
5 incorporating Gleason score, Stage, and PSA value have been used to predict outcome  
6 following surgery [25].

7  
8 The match between the eight underlying signatures detected for the MSKCC and  
9 CancerMap datasets was used as the basis for developing a novel classification framework  
10 for prostate cancer. A new algorithm called OAS-LPD was developed to allow rapid  
11 assessment of the presence of the signatures in individual cancer samples. In total four  
12 clinically and or genetically distinct subgroups were identified (DESNT, LPD3, LPD4, and  
13 LPD5, Fig. 5). The functional significance of the new disease groupings, for example in  
14 determining drug sensitivity, remains to be established. However, with the use of OAS-LPD  
15 it will be possible to undertake assessments of the response of patients in each of the groups  
16 DESNT, LPD3, LPD4, and LPD5 to drug treatments. There is limited overlap between the  
17 new classification and previously proposed subgroups based on genetic alterations [16,26-  
18 29].

19 Multiplatform data (expression, mutation, and methylation data from each cancer sample)  
20 are available for many cancer types, for example from The Cancer Genome Atlas. This has  
21 prompted the development of additional methods for sub-class discovery that can combine  
22 information from different platforms including the copula mixed model [30], Bayesian  
23 consensus clustering [31] and the iCluster model [16]. These approaches can suffer from the  
24 problem of sample assignment to a particular cluster or group, and the failure to take into

1 consideration the heterogeneous composition of individual cancer samples. These  
2 observations highlight the need to develop methods similar to LPD that can be applied to  
3 multiplatform data.

4

5 An important issue for patients diagnosed with prostate cancer is that clinical outcome is  
6 highly variable and precise prediction of the course of disease progression at the time of  
7 diagnosis is not possible [32]. In some studies, the use of population PSA screening can  
8 reduce mortality from prostate cancer by up to 21% [33]. However many, if not most,  
9 prostate cancers that are currently detected by PSA screening are clinically insignificant [34].

10 Over-diagnosis of clinically insignificant prostate cancer is a major issue and is set to  
11 increase still further [35]. There is therefore an urgent need for the identification of cancer  
12 categories that are associated with clinically aggressive or indolent disease to allow the  
13 targeting of radical therapies to the men that need them. For breast cancer unsupervised  
14 hierarchical clustering of transcriptome data resulted in a classification system that is  
15 routinely used to guide the management and treatment of this disease. Here we established  
16 a novel classification framework for the analysis of prostate cancer that has its origins in  
17 unsupervised analyses of transcriptome data. In future studies we plan to analyse the utility  
18 of DESNT and other LPD processes (particularly LPD3, LPD4 and LPD5) in managing  
19 prostate cancer patients, including predicting response to drug treatment. This will be  
20 performed through the assessment of LPD status in the contexts of established clinical trials.  
21 For evaluation we would plan to use each LPD assignment (eg DESNT, LPD3, LPD4 and  
22 LPD5) as a continuous variable as illustrated here by the development of a nomograms for  
23 the use of DESNT in predicting PSA failure. In conclusion our results highlight the

1 importance of devising and using more sophisticated approaches for the analysis of genomic  
2 datasets from all biological systems.

3

#### 4 **Additional Information**

#### 5 **Acknowledgements**

6 The research presented in this paper was carried out on the High Performance Computing  
7 Cluster supported by the Research and Specialist Computing Support service at the  
8 University of East Anglia. We thank Shea Connell for useful comments and suggestions on  
9 the manuscript.

#### 10 **Authors' Contributions**

11 CSC, DSB, and VM were involved in funding acquisition and supervised the project. CSC,  
12 DSB, B-AL, VM were involved in conceptualisation and planned the data analysis. B-AL,  
13 CE, DSB performed the majority of the analyses and investigations, with additional analysis  
14 and insight provided by VM, DRE, CC, RAC, JC, and CSC. B-AL, CE, CC, DSB, and CSC  
15 were involved in developing the methodology for the project. CSC, DSB, and B-AL wrote  
16 the original draft of the manuscript. All authors reviewed and edited the manuscript. All  
17 authors read and approved the final manuscript.

#### 18 **Ethics approval and consent to participate**

19 All data were from other publications. The ethical approvals obtained for each dataset are  
20 listed in the original publications.

#### 21 **Consent for publication**

22 All authors consent to publication.

#### 23 **Data availability**

1 The datasets analysed during the current study are available (Table 1). The majority are  
2 available from the Gene Expression Omnibus repository:

- 3 • MSKCC[7] : <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE21034>
- 4 • CancerMap[14] : <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE94767>
- 5 • Klein[17] : <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE62667>
- 6 • CamCap[6] : <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE70768> and  
7 <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE70769>
- 8 • Erho[18] : <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE46691>
- 9 • Karnes[19] : <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE62116>
- 10 • Stephenson[15] : Data available from the corresponding author of this paper.
- 11 • TCGA[16]: Data available from the TCGA Data Portal  
12 <https://portal.gdc.cancer.gov/projects/TCGA-PRAD>

### 13 **Competing Interests**

14 C.S.C., D.S.B., B-A L., and V.M. are co-inventors on a patent application from the University  
15 of East Anglia on the detection of DESNT prostate cancer.

### 16 **Funding**

17 This work was funded by the Bob Champion Cancer Trust, The Masonic Charitable  
18 Foundation successor to The Grand Charity, The King Family, The Stephen Hargrave  
19 Foundation, and The University of East Anglia. We acknowledge support from Movember,  
20 from Prostate Cancer UK, Callum Barton, and from The Andy Ripley Memorial Fund.

21

### 22 **References**

- 23 1. Edgar R, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression  
24 and hybridization array data repository. *Nucleic Acids Res.* 2002;30(1):207–10.

- 1 2. Consortium ICG, Anderson W, Artez A, Bell C, Bernabé RR, Bhan MK, et al.  
2 International network of cancer genome projects. *Nature*. 2010;464(7291):993–8.
- 3 3. Sorlie T, Tibshirani R, Parker J, Hastie T, Marron JS, Nobel A, et al. Repeated  
4 observation of breast tumor subtypes in independent gene expression data sets. *Proc*  
5 *Natl Acad Sci U S A*. 2003 Jul 8;100(14):8418–23.
- 6 4. Blanco-Calvo M, Concha Á, Figueroa A, Garrido F, Valladares-Ayerbes M. Colorectal  
7 cancer classification and cell heterogeneity: A systems oncology approach. *Int J Mol*  
8 *Sci*. 2015;16(6):13610–32.
- 9 5. Polyak K. Heterogeneity in breast cancer Review series introduction Heterogeneity in  
10 breast cancer. *J Clin Invest*. 2011;121(10):3786.
- 11 6. Ross-Adams H, Lamb ADD, Dunning MJJ, Halim S, Lindberg J, Massie CMM, et al.  
12 Integration of copy number and transcriptomics provides risk stratification in  
13 prostate cancer: A discovery and validation cohort study. *EBioMedicine*.  
14 2015;2(9):1133–44.
- 15 7. Taylor BS, Schultz N, Hieronymus H, Gopalan A, Xiao Y, Carver BS, et al. Integrative  
16 Genomic Profiling of Human Prostate Cancer. *Cancer Cell*. 2010 Jun 13;18(1):11–22.
- 17 8. Cooper CS, Eeles R, Wedge DC, Van Loo P, Gundem G, Alexandrov LB, et al.  
18 Analysis of the genetic phylogeny of multifocal prostate cancer identifies multiple  
19 independent clonal expansions in neoplastic and morphologically normal prostate  
20 tissue. *Nat Genet*. 2015 Mar 2;47(4):367–72.
- 21 9. Boutros PC, Fraser M, Harding NJ, de Borja R, Trudel D, Lalonde E, et al. Spatial  
22 genomic heterogeneity within localized, multifocal prostate cancer. *Nat Genet*. 2015  
23 May 25;47(7):736–45.
- 24 10. Tsourlakis M-C, Stender A, Quaas A, Kluth M, Wittmer C, Haese A, et al.  
25 Heterogeneity of ERG expression in prostate cancer: a large section mapping study of

- 1 entire prostatectomy specimens from 125 patients. *BMC Cancer*. 2016;16(1):641.
- 2 11. Rogers S, Girolami M, Campbell C, Breitling R. The latent process decomposition of  
3 cDNA microarray data sets. *IEEE/ACM Trans Comput Biol Bioinforma*.  
4 2005;2(2):143-56.
- 5 12. Carrivick L, Rogers S, Clark J, Campbell C, Girolami M, Cooper C. Identification of  
6 prognostic signatures in breast cancer microarray data using Bayesian techniques. *J R*  
7 *Soc Interface*. 2006;3(8):367-81.
- 8 13. Olmos D, Brewer D, Clark J, Danila DC, Parker C, Attard G, et al. Prognostic value of  
9 blood mRNA expression signatures in castration-resistant prostate cancer: a  
10 prospective, two-stage study. *Lancet Oncol*. 2012 Oct 8;2045(12):1-11.
- 11 14. Luca B, Brewer DS, Edwards DR, Edwards S, Whitaker HC, Merson S, et al. DESNT:  
12 A Poor Prognosis Category of Human Prostate Cancer. *Eur Urol Focus*. 2018 Dec  
13 6;4(6):842-50.
- 14 15. Stephenson AJ, Smith A, Kattan MW, Satagopan J, Reuter VE, Scardino PT, et al.  
15 Integration of gene expression profiling and clinical variables to predict prostate  
16 carcinoma recurrence after radical prostatectomy. *Cancer*. 2005;104(2):290-8.
- 17 16. Network CGAR, Cancer Genome Atlas Research Network. The Molecular Taxonomy  
18 of Primary Prostate Cancer. *Cell*. 2015 Nov 5;163(4):1011-25.
- 19 17. Klein EA, Yousefi K, Haddad Z, Choeurng V, Buerki C, Stephenson AJ, et al. A  
20 genomic classifier improves prediction of metastatic disease within 5 years after  
21 surgery in node-negative high-risk prostate cancer patients managed by radical  
22 prostatectomy without adjuvant therapy. *Eur Urol*. 2015;67(4):778-86.
- 23 18. Erho N, Crisan A, Vergara IA, Mitra AP, Ghadessi M, Buerki C, et al. Discovery and  
24 validation of a prostate cancer genomic classifier that predicts early metastasis  
25 following radical prostatectomy. *PLoS One*. 2013;8(6):e66855.

- 1 19. Karnes RJ, Bergstralh EJ, Davicioni E, Ghadessi M, Buerki C, Mitra AP, et al.  
2 Validation of a Genomic Classifier that Predicts Metastasis Following Radical  
3 Prostatectomy in an At Risk Patient Population. *J Urol*. 2013;190(6):2047–53.
- 4 20. Irizarry RA, Hobbs B, Collin F, Beazer-Barclay YD, Antonellis KJ, Scherf U, et al.  
5 Exploration, normalization, and summaries of high density oligonucleotide array  
6 probe level data. *Biostatistics*. 2003;4(2):249–64.
- 7 21. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for  
8 RNA-seq data with DESeq2. *Genome Biol*. 2014;15(12):550.
- 9 22. Therneau TM, GRAMBSCH PM, Fleming TR. Martingale-based residuals for survival  
10 models. *Biometrika*. 1990 Mar 1;77(1):147–60.
- 11 23. Hair JF, Black WC, Babin BJ, Anderson RE, Tatham RL. *Multivariate data analysis*.  
12 1998.
- 13 24. Levine DM, Haynor DR, Castle JC, Stepaniants SB, Pellegrini M, Mao M, et al.  
14 Pathway and gene-set activation measurement from mRNA expression data: the  
15 tissue distribution of human pathways. *Genome Biol*. 2006;7(10):R93.
- 16 25. Shariat SF, Kattan MW, Vickers AJ, Karakiewicz PI, Scardino PT. Critical review of  
17 prostate cancer predictive tools. *Future Oncol*. 2009;5(10):1555–84.
- 18 26. Attard G, Clark J, Ambroisine L, Fisher G, Kovacs G, Flohr P, et al. Duplication of the  
19 fusion of TMPRSS2 to ERG sequences identifies fatal human prostate cancer.  
20 *Oncogene*. 2008;27(3):253–63.
- 21 27. Reid AHM, Attard G, Ambroisine L, Fisher G, Kovacs G, Brewer D, et al. Molecular  
22 characterisation of ERG, ETV1 and PTEN gene loci identifies patients at low and high  
23 risk of death from prostate cancer. *Br J Cancer*. 2010;102(4):678–84.
- 24 28. Mosquera JM, Beltran H, Park K, MacDonald TY, Robinson BD, Tagawa ST, et al.  
25 Concurrent AURKA and MYCN Gene Amplifications Are Harbingers of Lethal

- 1 TreatmentRelated Neuroendocrine Prostate Cancer. *Neoplasia*. 2013;15(1):1-IN4.
- 2 29. Rodrigues LU, Rider L, Nieto C, Romero L, Karimpour-Fard A, Loda M, et al.  
3 Coordinate loss of MAP3K7 and CHD1 promotes aggressive prostate cancer. *Cancer*  
4 *Res*. 2015 Mar 15;75(6):1021-34.
- 5 30. Rey M, Roth V. Copula Mixture Model for Dependency-seeking Clustering. 2012;
- 6 31. Lock EF, Dunson DB. Bayesian consensus clustering. *Bioinformatics*.  
7 2013;29(20):2610-6.
- 8 32. Buyyounouski MK, Pickles T, Kestin LL, Allison R, Williams SG. Validating the  
9 Interval to Biochemical Failure for the Identification of Potentially Lethal Prostate  
10 Cancer. *J Clin Oncol*. 2016;30(15):1857-63.
- 11 33. Schröder FH, Hugosson J, Roobol MJ, Tammela TLJ, Zappa M, Nelen V, et al.  
12 Screening and prostate cancer mortality: results of the European Randomised Study  
13 of Screening for Prostate Cancer (ERSPC) at 13 years of follow-up. *Lancet*.  
14 2014;384(9959):2027-35.
- 15 34. Etzioni R, Gulati R, Mallinger L, Mandelblatt J. Influence of study features and  
16 methods on overdiagnosis estimates in breast and prostate cancer screening. *Ann*  
17 *Intern Med*. 2013 Jun 4;158(11):831-8.
- 18 35. Parker C, Emberton M. Screening for prostate cancer appears to work, but at what  
19 cost? *BJU Int*. 2009;104(3):290-2.

20

21

## 1 **Legends**

2 **Figure 1.** LPD decomposition of the MSKCC dataset. (a) DESNT bar chart from the LPD  
3 decomposition of the MSKCC dataset [14]. showing the number ID assigned to 23 examples  
4 samples that had some amount of DESNT signature. (b,c) Pie charts showing the relative  
5 proportions of the eight LPD signatures in 23 example samples. DESNT is in red, other LPD  
6 signatures are represented by different colours as indicated in the key. The number next to  
7 each pie chart indicates which cancer it represents from the bar chart above. Individual  
8 cancer samples were assigned as a “DESNT cancer” when the DESNT signature was the  
9 most abundant; examples are shown in the right-hand box (c, ‘DESNT’). Many other cancer  
10 samples contained a smaller proportion of DESNT cancer and were associated with a poor  
11 outcome: examples shown in larger box (b, ‘SOME DESNT’).

12

13 **Figure 2.** Stratification of prostate cancer samples based on the percentage of DESNT cancer  
14 present. For these analyses the data from the MSKCC, CancerMap, CamCap, and  
15 Stephenson datasets were combined (n=503). (a) Plot showing the proportion of DESNT  
16 signature in each cancer sample and the division into four groups of increasing DESNT.  
17 Group 1 samples have a proportion of less than 0.001 of the DESNT signature. (b) Kaplan-  
18 Meier plot showing the Biochemical Recurrence (BCR) free survival based on proportion of  
19 DESNT cancer present as determined by LPD. Number of cancer patients in each Group are  
20 indicated (bottom right) and the number of PCR failures in each group are show in  
21 parentheses. The definition of Groups 1-4 is shown in Fig. 2a. Cancer samples with  
22 proportions up to 0.3 DESNT (Group 2) exhibited poorer clinical outcome ( $X^2$ -test,  $P = 0.011$ )  
23 compared to cancer samples lacking DESNT ( $<0.001$ ). Cancer samples with the  
24 intermediate (0.3 to 0.6) and high ( $>0.6$ ) proportions of DESNT also exhibited significantly  
25 worse outcome (respectively  $P = 2.6 \times 10^{-5}$  and  $P = 8.3 \times 10^{-9}$  compare to cancer samples  
23

1 lacking DESNT. The combined log-rank  $P = 1.3 \times 10^{-8}$ ). (c) Nomogram model developed to  
2 predict PSA free survival at one, three, five, and seven years using proportion of DESNT.  
3 Assessing a single patient each clinical variable has a corresponding point score (top scales).  
4 The point scores for each variable are added to produce a total points score for each patient.  
5 The predicted probability of PSA free survival at one, three, five, and seven years can be  
6 determined by drawing a vertical line from the total points score to the probability scales  
7 below.

8

9 **Figure 3.** Prediction of clinical outcome according to OAS-LPD group. (a-c) Kaplan-Meier  
10 plots showing PSA free survival outcomes for the cancer patients assigned to LPD groups in  
11 analyses of the combine MSKCC, CancerMap, CamCap, and Stephenson datasets: (a)  
12 comparison of all LPD groups (LPD7 is DESNT); (b) cancer patients assigned to LPD4  
13 compared to patients assigned to all other LPD groups; (c) cancer patients assigned to  
14 DESNT (LPD7) compared to cancers assigned to all other LPD groups. (d) OAS-LPD  
15 signature assignment proportions for the 19 metastatic tissue samples reported as part of the  
16 MSKCC dataset. In all cases DESNT (LPD7) was the dominant expression signature  
17 detected.

18

19 **Figure 4.** (a) OAS-LPD sub-groups in The Cancer Genome Atlas Dataset (n=333). Cancer  
20 samples were assigned to subgroups based on the most prominent signature as detected by  
21 OAS-LPD. The types of genetic alteration are shown for each gene (mutations, fusions,  
22 deletions, and over-expression). Clinical parameters including biochemical recurrence (BCR)  
23 are represented at the bottom together with groups for iCluster, methylation, somatic copy  
24 number alteration (SVNA), and messenger RNA (mRNA)[16]. Comparison of the frequency  
25 of genetic alterations present in each subgroup are shown in Table 2. (b-d) Kaplan-Meier

1 plots showing PSA free survival outcomes for ETS-rearrangement positive cancers in LPD3  
2 compared to all other ETS-positive cancers for the CancerMap, CamCap, and TCGA  
3 datasets.

4

5 **Figure 5.** A classification framework for prostate cancer. Based on the analyses of genetic  
6 and clinical correlations we consider that there is good evidence for the existence of LPD3,  
7 LPD4, and LPD5 as separate cancer categories, moderate evidence of the existence of LPD6  
8 and LPD8 (based on alteration of expression only), and weak evidence for LPD1. The  
9 methylation column list all genes that exhibit differential expression and that also contain at  
10 least one locus that is differentially methylated.

11

12

**Table 1**

Dataset	Primary	Normal	Type	Platform	Citation
MSKCC [7]	131	29	FF	Affymetrix Exon 1.0 ST v2	Taylor <i>et al.</i> 2010
CancerMap [14]	137	17	FF	Affymetrix Exon 1.0 ST v2	Luca <i>et al.</i> 2017
Stephenson [15]	78	11	FF	Affymetrix U133A	Stephenson <i>et al.</i> 2005
Klein [17]	182	0	FFPE	Affymetrix Exon 1.0 ST v2	Klein <i>et al.</i> 2015
CamCap [6]	147	73	FF	Illumina HT12 v4.0 BeadChip	Ross-Adams <i>et al.</i> 2015
TCGA [16]	333	43	FF	Illumina HiSeq 2000 RNA-Seq v2	TCGA network 2015
Erho [18]	545	0	FFPE	Affymetrix Exon 1.0 ST v2	Erho <i>et al.</i> 2013
Karnes [19]	232	0	FFPE	Affymetrix Exon 1.0 ST v2	Karnes <i>et al.</i> 2013

**Table 1.** Transcriptome datasets. The MSKCC study additionally reported expression profiles from 19 metastatic cancers. The ethical approvals obtained for each dataset are listed in the original publications.

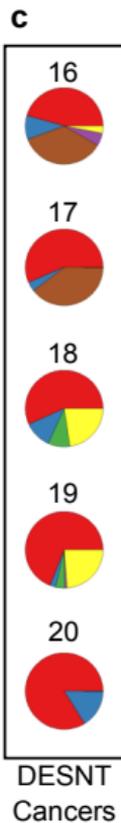
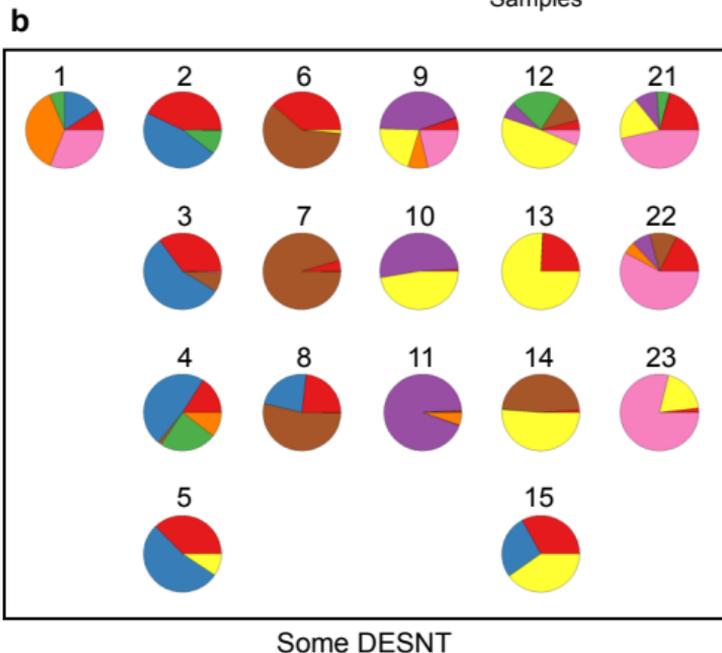
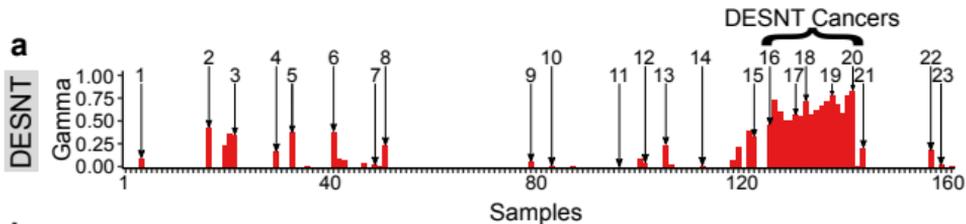
**Table 2**

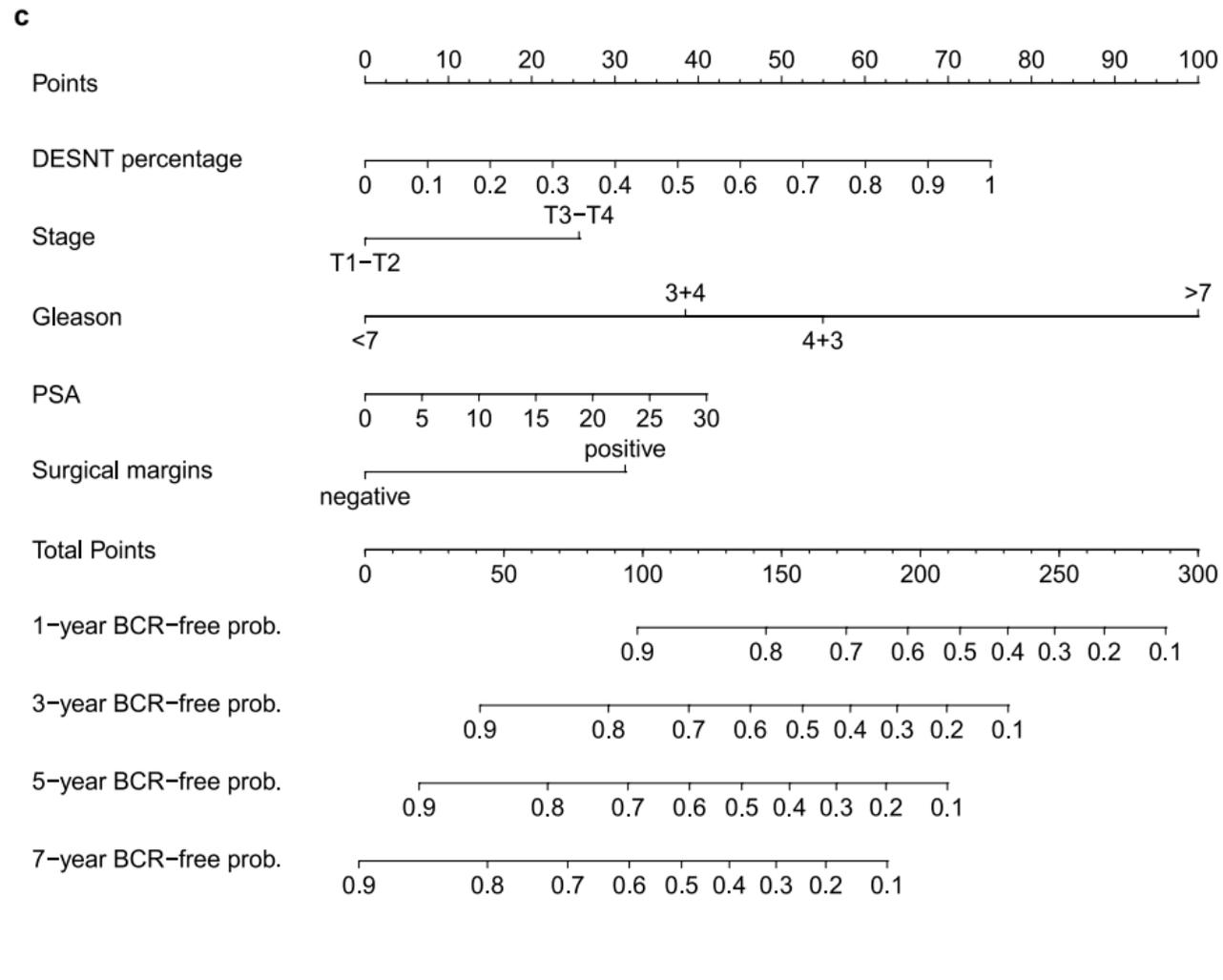
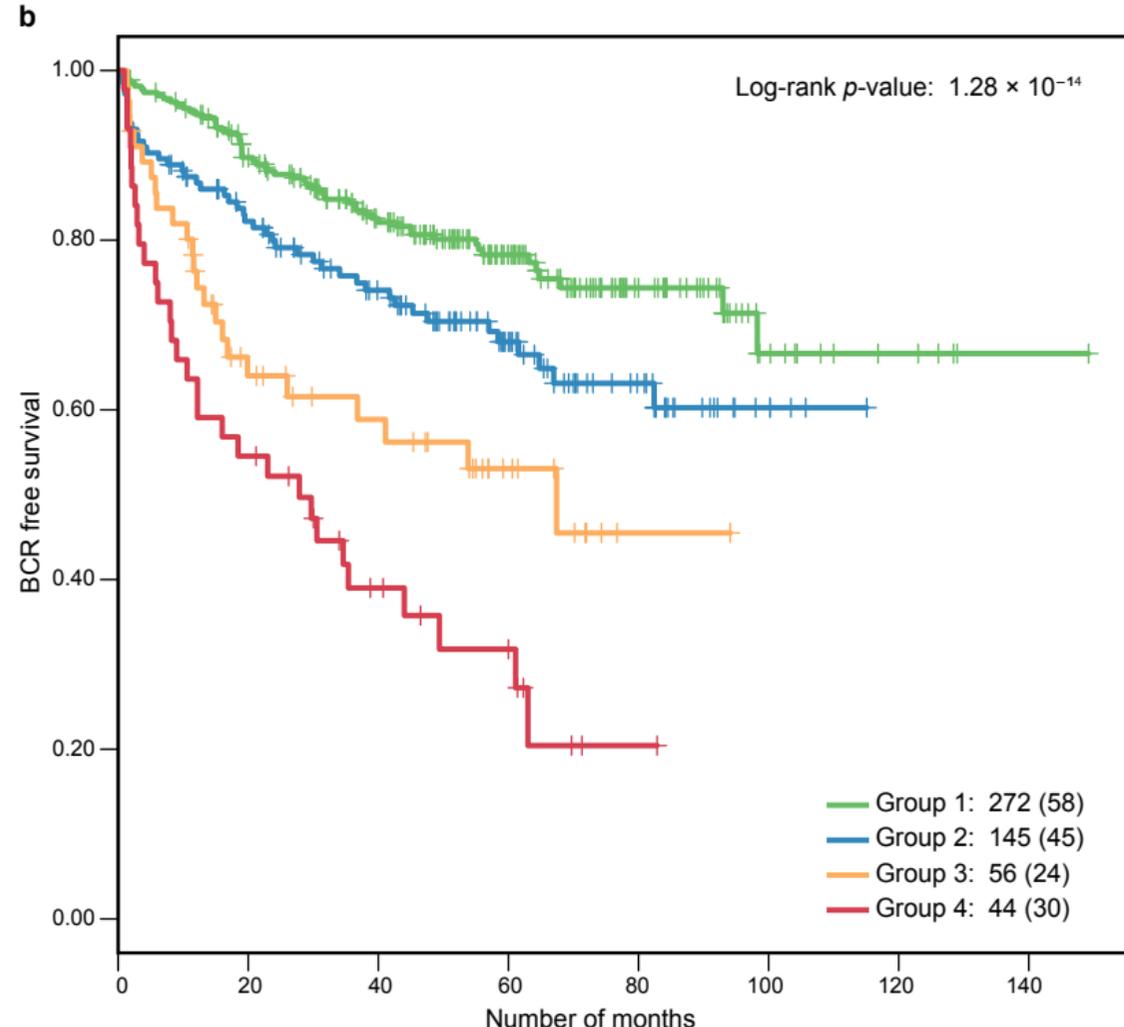
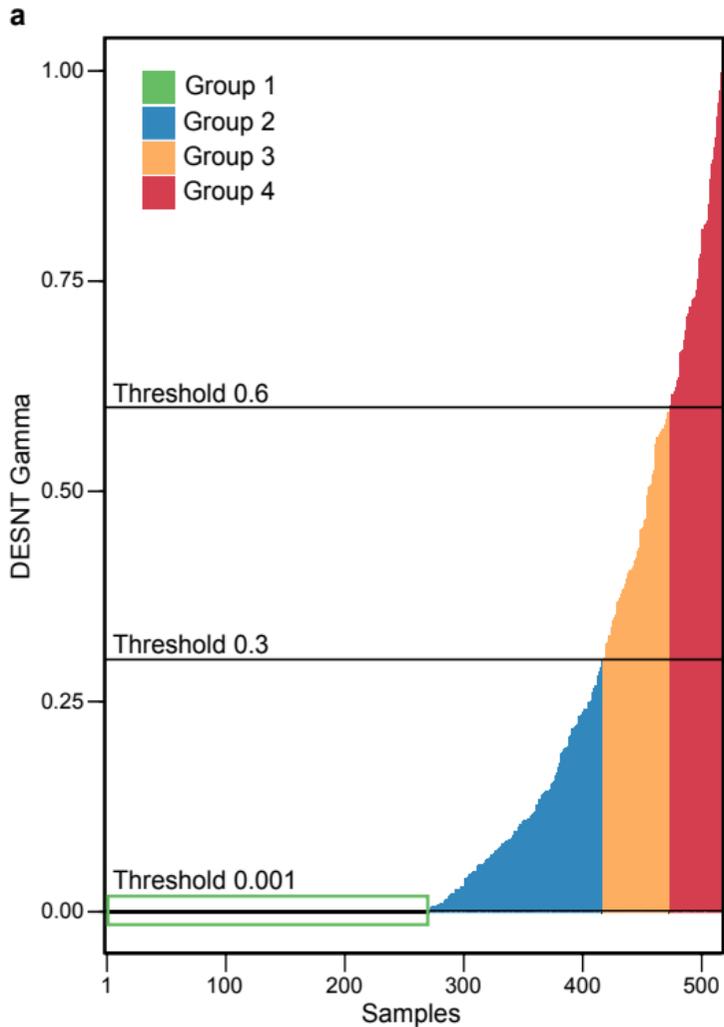
	TCGA			CancerMap			CamCap		
	ETS-	ETS+	$\chi^2$ P-val	ERG-	ERG+	$\chi^2$ P-val	ERG-	ERG+	$\chi^2$ P-val
LPD1	8	3	0.0588	13	4	0.0851	0	3	0.235
LPD2	4	8	0.827	3	3	1	0	2	0.467
LPD3	9	67	1.45x10 <sup>-08</sup>	5	15	0.00977	4	17	0.00299
LPD4	14	21	1	14	15	0.619	1	2	0.987
LPD5	65	5	2.20x10 <sup>-16</sup>	19	1	0.000180	34	0	1.15x10 <sup>-11</sup>
LPD6	13	22	0.802	5	5	1	2	4	0.657
DESNT	13	66	1.17x10 <sup>-06</sup>	6	15	0.0207	9	24	0.00274
LPD8	9	6	0.193	8	4	0.540	4	1	0.371

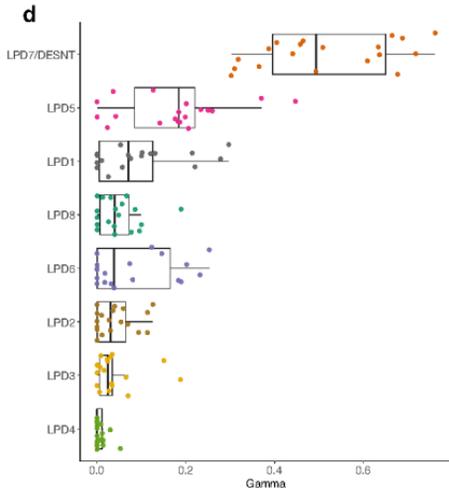
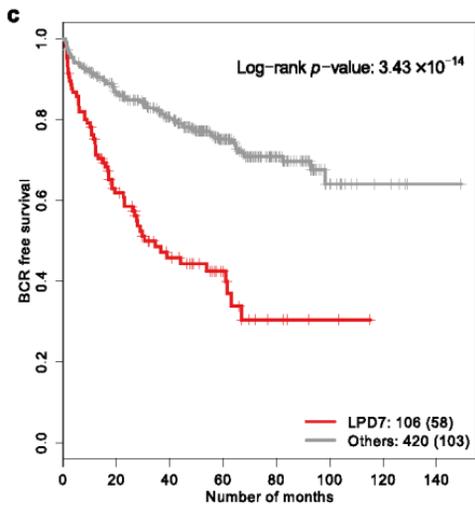
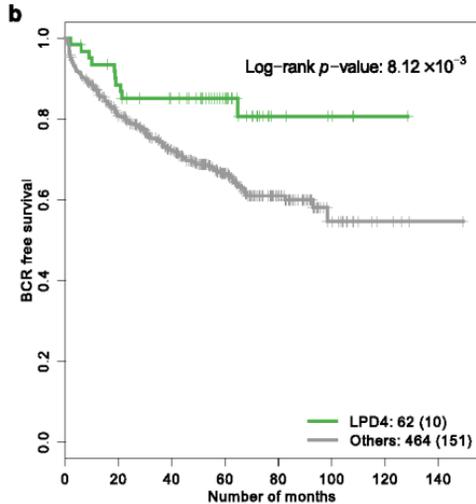
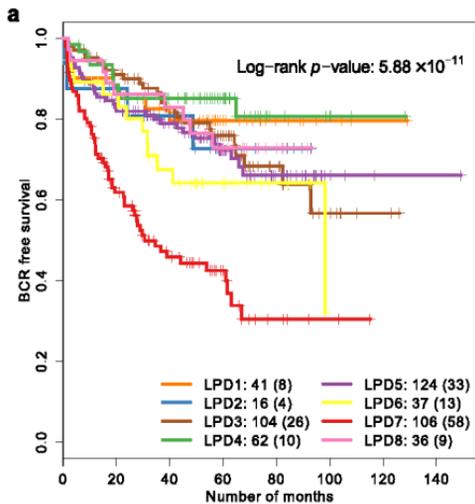
  

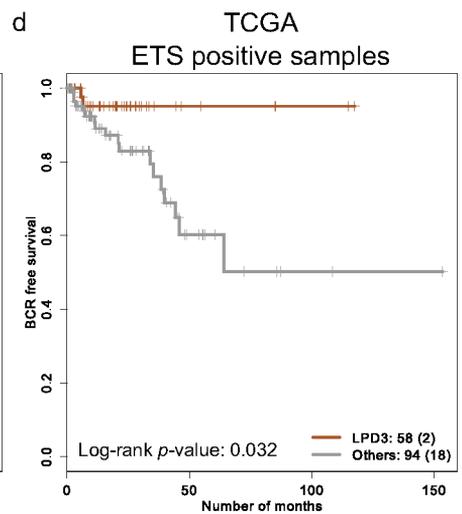
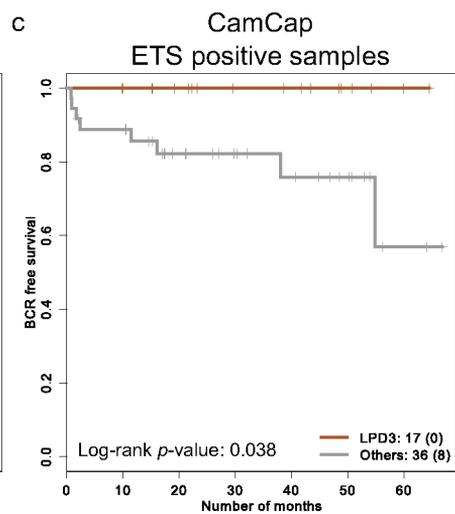
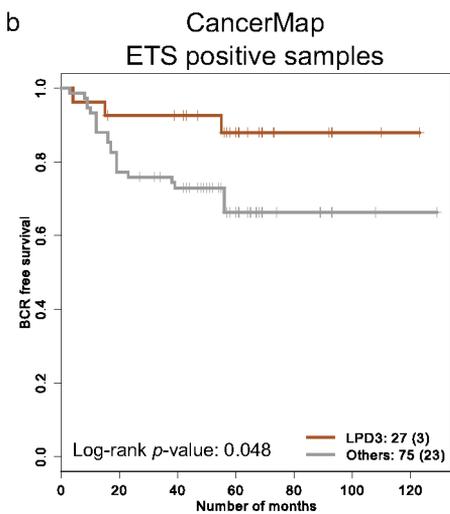
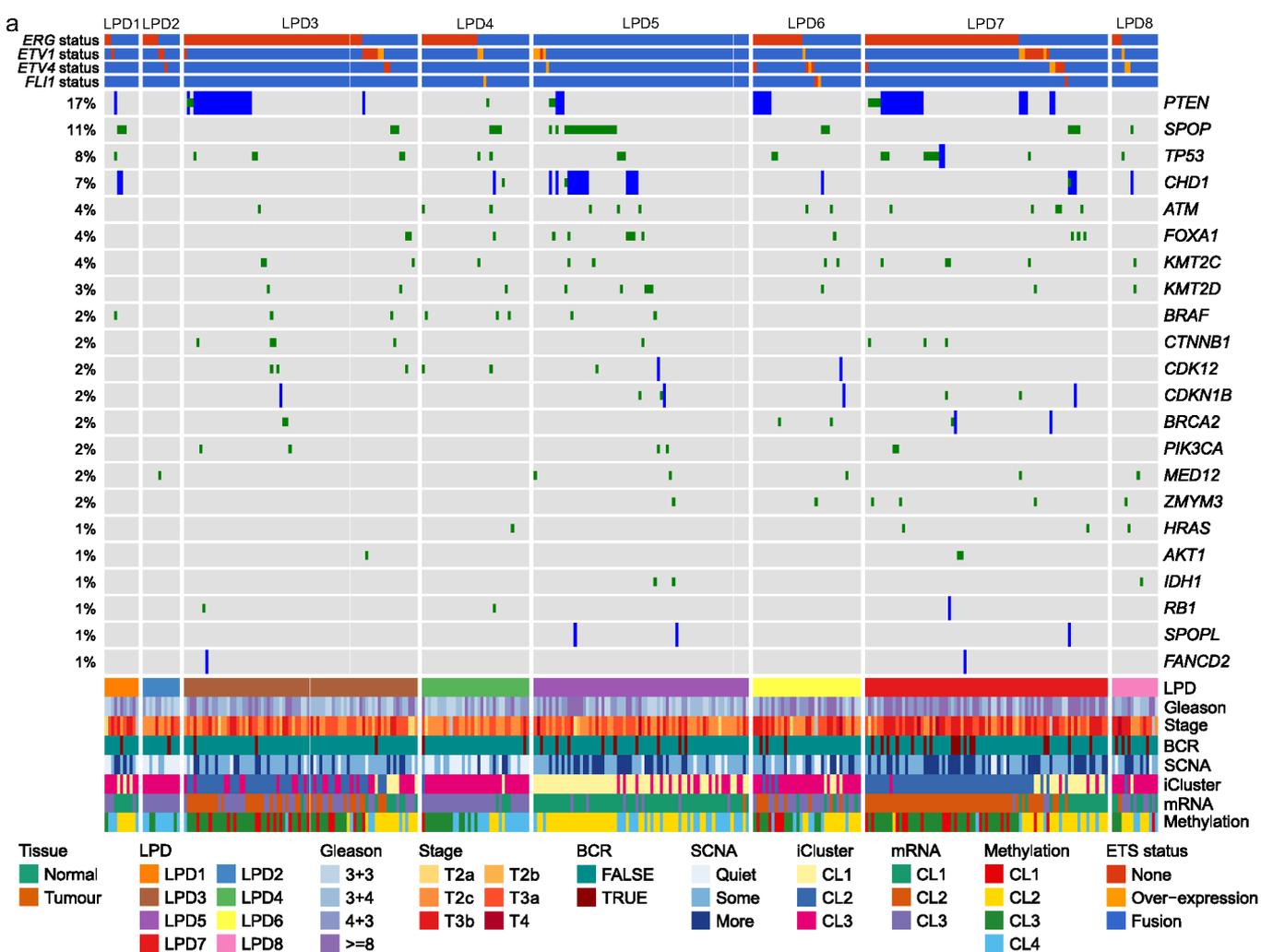
	<i>PTEN</i>			<i>SPOP</i>			<i>CHD1</i>		
	Non-homdel	Homdel	$\chi^2$ P-val	Non-mut	Mut	$\chi^2$ P-val	Non-homdel	Homdel	$\chi^2$ P-val
LPD1	10	1	0.896	8	3	0.213	9	2	0.309
LPD2	12	0	0.284	12	0	0.436	12	0	0.756
LPD3	55	21	0.000894	73	3	0.0400	76	0	0.0211
LPD4	35	0	0.0174	31	4	1	34	1	0.603
LPD5	67	3	0.00830	51	19	4.46x10 <sup>-06</sup>	57	13	7.69x10 <sup>-06</sup>
LPD6	29	6	0.903	32	3	0.825	34	1	0.603
DESNT	60	19	0.0167	75	4	0.0795	76	3	0.432
LPD8	15	0	0.195	14	1	0.889	14	1	1

**Table 2.** Correlation of OAS-LPD subgroups with genetic alterations in The Cancer Genome Atlas Dataset. Statistically significant differences are highlighted in grey.









	CLINICAL	GENE MUTATIONS	EXPRESSION	METHYLATION
LPD1	Less metastases in Erho <i>et al.</i> dataset.	No current evidence.		No current evidence.
LPD2 – NORMAL LIKE	Frequently contains normal samples. (Over-represented in 4/5 datasets)	No current evidence	<ul style="list-style-type: none"> <li>▲ 2 genes ▼ 0 genes</li> <li>▲ Structural molecule activity</li> <li>▲ Protein-glutamine gamma-glutamyl transferase activity</li> </ul>	<i>KRT13, TGM4</i>
LPD3 – STAT SIGNALING	<i>ERG+</i> cancers have better outcome. (2/3 datasets)	<ul style="list-style-type: none"> <li>▲ <i>ERG/ETS</i> ▲ <i>PTEN</i></li> <li>▼ <i>SPOP</i> ▼ <i>CHD1</i></li> </ul>	<ul style="list-style-type: none"> <li>▲ 7 genes ▼ 0 genes</li> <li>▲ Regulation of STAT proteins</li> <li>▲ Regulation of insulin secretion</li> <li>▲ Second-messenger-mediated signalling</li> </ul>	<i>CSGALNACT1, ERG, GHR, GUCY1A3, HDAC1, ITPR3, PLA2G7</i>
LPD4 – IMPROVED SURVIVAL	Improved survival based on PSA failure. Frequently contains normal samples. (Over-represented in 2/5 datasets)	▼ <i>PTEN</i>		No current evidence.
LPD5 – SPOP/CHD1	No current evidence.	<ul style="list-style-type: none"> <li>▼ <i>ERG/ETS</i> ▼ <i>PTEN</i></li> <li>▲ <i>SPOP</i> ▲ <i>CHD1</i></li> </ul>	<ul style="list-style-type: none"> <li>▲ 43 genes ▼ 13 genes</li> <li>▲ Fatty acid metabolism</li> <li>▲ Branched-chain aminoacid metabolism</li> <li>▲ Positive regulation of secretion</li> </ul>	<i>ABHD2, ACAD8, ACLY, ALCAM, ALDH6A1, ALOX15B, ARHGEF7, AUH, BBS4, C1orf115, CAMKK2, COG5, CPEB3, CYP2J2, DHRS3, DHX32, EHHADH, ELOVL2, ERG, EXTL2, F3, FAM111A, GATA3, GLUD1, GNMT, HES1, HPGD, KHDRBS3, LAMB2, LAMC2, MIPEP, MON1B, NANS, NAT1, NCAPD3, PDE8B, PPFIBP2, PTK7, PTPN13, PTPRM, RAB27A, REPS2, RFX3, SCIN, SLC1A1, SLC4A4, SMPDL3A, SORL1, STXBP6, SYTL2, TBPL1, TFF3, TRIM29, TUBB2A, YIPF1, ZNF516</i>
LPD6 – TRANSITION METAL HOMEOSTASIS	No current evidence.	No current evidence.	<ul style="list-style-type: none"> <li>▲ 9 genes ▼ 0 genes</li> <li>▲ Transition metal ion homeostasis</li> <li>▲ Viral genome replication</li> <li>▲ Humoral immune response</li> </ul>	<i>CCL2, CFB, CFTR, CXCL2, IFI16, LCN2, LTF, LXN, TFRC</i>
LPD7 – DESNT POOR PROGNOSIS	Independent predictor of poor outcome based on PSA failure. More metastases in Erho <i>et al.</i> and MSKCC datasets. Dominant signature in metastatic cancers.	▲ <i>ERG/ETS</i> ▲ <i>PTEN</i>	<ul style="list-style-type: none"> <li>▲ 2 genes ▼ 49 genes</li> <li>▼ Cell-substrate adhesion</li> <li>▼ Muscle contraction</li> <li>▼ Cell junction organisation</li> </ul>	<i>ACTG2, ACTN1, ADAMTS1, ANPEP, ARMCX1, AZGP1, C7, CD44, CHRDL1, CNN1, CRISPLD2, CSRP1, CYP27A1, CYR61, DES, EGR1, ETS2, F5, FBLN1, FERMT2, FHL2, FLNA, FXYD6, FZD7, ITGA5, ITM2C, JAM3, JUN, KHDRBS3, LMOD1, LPHN2, MT1M, MYH11, MYL9, NFIL3, PARM1, PCP4, PDK4, PLAGL1, RAB27A, SERPINF1, SNAI2, SORBS1, SPARCL1, SPOCK3, SYNM, TAGLN, TCEAL2, TGFB3, TPM2, VCL</i>
LPD8 – NORMAL LIKE	Contains normal samples. (Over-represented in 5/5 datasets)	No current evidence.	<ul style="list-style-type: none"> <li>▲ 29 genes ▼ 37 genes</li> <li>▲ Extracellular matrix organisation</li> <li>▲ Extracellular structure organisation</li> </ul>	<i>ABCC4, ACAT2, ARHGFE6, ATP8A1, AXL, CANT1, CD83, CDH1, COL15A1, DCXR, DHCR24, DHRS7, DPYSL3, EPB41L3, FAM174B, FAM189A2, FBN1, FCHSD2, FHL1, FKBP4, FOXA1, FXYD5, GNAO1, GOLM1, GPX3, GTF3C1, HPN, IFI16, IRAK3, ITGA5, KIF5C, KLK3, LAPTM5, MAP7, MBOAT2, MFAP4, MFG8, MIOS, MLPH, MMP2, MYO5C, NEDD4L, PART1, PARVA, PDIA5, PIGH, PLEKH01, PLSCR4, PMEPA1, PRSS8, RFTN1, SAMD4A, SAMSN1, SEC23B, SERPINF1, SLC43A1, SPDEF, SPINT2, STEAP4, TMPRSS2, TRPM8, TSPAN1, VCAM1, WIPF1, XBP1, ZYX</i>