

1 **Comprehensive processing of high throughput small RNA sequencing**
2 **data including quality checking, normalization and differential**
3 **expression analysis using the UEA sRNA Workbench**

4
5
6 Matthew Beckers^{1*}, Irina Mohorianu^{2,1*}, Matthew Stocks^{1*},
7 Christopher Applegate¹, Tamas Dalmay², Vincent Moulton^{1#}

8 ¹School of Computing Sciences, University of East Anglia, Norwich, NR4 7TJ, UK

9 ²School of Biological Sciences, University of East Anglia, Norwich, NR4 7TJ, UK

10
11 #corresponding author.

12 Tel: +44 1603 592607; Fax: +44 1603 593345;

13 Email: v.moulton@uea.ac.uk

14 Present Address: School of Computing Sciences, University of East Anglia,
15 Norwich, NR4 7TJ, United Kingdom

16 *joint first authors

17
18 Running title: sRNA-seq processing with the UEA sRNA Workbench

19 **ABSTRACT**

20 Recently High Throughput Sequencing (HTS) has revealed compelling details about the small RNA
21 (sRNA) population in eukaryotes. These 20-25 nt non-coding RNAs can influence gene expression
22 by acting as guides for the sequence-specific regulatory mechanism known as RNA silencing. The
23 increase in sequencing depth and number of samples per project enables a better understanding of
24 the role sRNAs play by facilitating the study of expression patterns. However, the intricacy of the
25 biological hypotheses coupled with a lack of appropriate tools often leads to inadequate mining of the
26 available data and thus, an incomplete description of the biological mechanisms involved.

27 To enable a comprehensive study of differential expression in sRNA datasets we present a new
28 interactive pipeline that guides researchers through the various stages of data pre-processing and
29 analysis. This includes various tools, some of which we specifically developed for sRNA analysis, for
30 quality checking and normalization of sRNA samples as well as tools for the detection of differentially
31 expressed sRNAs and identification of the resulting expression patterns.

32 The pipeline is available within the UEA sRNA Workbench, a user-friendly software package for the
33 processing of sRNA datasets. We demonstrate the use of the pipeline on a *H. sapiens* dataset;
34 additional examples on a *B. terrestris* dataset and on an *A. thaliana* dataset are described in the
35 supplementary information. A comparison with existing approaches is also included, which
36 exemplifies some of the issues that need to be addressed for sRNA analysis, and how the new
37 pipeline may be used to do this.

38 **Keywords:**

39 high-throughput sequencing (HTS), microRNA (miRNA), small RNA (sRNA), UEA sRNA Workbench,
40 quality checking, normalization, differential expression

41 **Introduction**

42 RNA silencing is known to play a key role in the fine-tuning of gene expression in eukaryotes
43 (Brodersen et al. 2006). The process is mediated by a set of RNA molecules referred to as small
44 RNAs (sRNAs). Well-known examples of sRNAs include microRNAs (miRNAs) (Bartel 2009; Voinnet
45 2009) and small interfering RNAs (siRNAs) (Carthew et al. 2009; Meister 2013). These RNA
46 fragments are excised by Dicer/Dicer-like proteins from double stranded RNA precursors deriving
47 either from single stranded RNAs with a hairpin-like secondary structure, the miRNAs (Zhu et al.
48 2013), or long double stranded RNA created by a polymerase, the siRNAs (Chen 2012). The sRNAs
49 target and subsequently silence genes and thus play an important role in gene regulation (Lippman
50 et al. 2004; Omidvar et al. 2015), defence against pathogens (Szittyá et al. 2010; Donaszi-Ivanov et
51 al. 2013) and general maintenance of the genome (Molnar et al. 2007; Mohorianu et al. 2011).

52 For most molecular biology experiments, an important question is how the observed phenotype or
53 inherent differences (e.g. time or organ/tissue series) are reflected in the variation in expression of
54 sRNAs, commonly referred to as differential expression analysis or DE analysis (Mohorianu et al.
55 2010; Garber et al. 2011; Oszolak et al. 2011; Xu et al. 2014). Identification of DE sequences consists
56 of several distinct stages: first, the quality of the data is investigated to identify (and potentially
57 exclude) samples containing artefacts such as over-representing biases originating from sequencing
58 inaccuracies (Sorefan et al. 2012; Raabe et al. 2014) or introduced from the handling of the original
59 biological sample. Second, the reads are annotated to determine which categories of sRNAs are
60 present. Finally, the expression levels in the samples are normalised to improve the comparability
61 between samples and, subsequently, to refine the accuracy of DE predictions (McCormick et al. 2011;
62 Dillies et al. 2013).

63 Bioinformatics methods developed for DE analysis have thus far largely focused on analysing
64 messenger RNA (mRNA) data, first from microarray experiments and now, more commonly, from
65 RNA-seq (mRNA-seq) datasets (Rapaport et al. 2013; Soneson et al. 2013). Many of these
66 approaches devised for each stage of a DE analysis are transferable to sRNA datasets (see table 1).

67 However, there are a number of conceptual differences between sRNA microarrays, which capture a
68 small number of known sequences (mainly miRNAs), and sRNA-seq, which capture a wider variety
69 known and novel sRNAs (usually, in excess of 100k unique reads). Similar differences in number of
70 quantified transcripts are also observed between the output of mRNA-seq experiments and sRNA-
71 seq output. More specifically, for mRNA studies, the expression levels of the reads are aggregated
72 into a gene abundance (Mortazavi et al. 2008) whereas each sRNA sequence contributes individually
73 to the distribution of abundances (McCormick et al. 2011; Studholme 2012). Because of this, the
74 resulting distributions are different both in shape; mRNA-seq abundances have a Gaussian-like
75 distribution whereas sRNA-seq abundances follow an exponential-like distribution, and in number of
76 points; thousands of genes compared to millions of unique sRNAs (Barquist et al. 2015). In addition,
77 sRNA-seq data has higher ratio of noise (random degradation products) to signal (genuine sRNAs);
78 due to the nature of sRNA-seq processing the median of sRNA abundances lies within the noise
79 range (Vidal et al. 2013). This implies that existing methodologies for microarrays or mRNA-seq DE
80 analyses are applicable but not always appropriate for sRNA-seq datasets (McCormick et al. 2011;
81 Gupta et al. 2012; Lohse et al. 2012; Vidal et al. 2013). Therefore it is important to develop tools that
82 address the specific characteristics of sRNA-seq datasets and their analysis to complement those
83 currently used for mRNA-seq analysis.

84 A common approach for HTS data analysis is to group several tools into a pipeline. As well as
85 providing the ability to tailor pipelines to individual experiments, this enables researchers to configure
86 the distinct stages of the analysis as required (Davis et al. 2013). After the setup is complete the (likely
87 lengthy) procedure can be executed without the need for further input from the user. Currently there
88 are several mRNA-seq pipelines available, such as DESeq/DESeq2 (Anders et al. 2013; Love et al.
89 2014) or edgeR (Zhou et al. 2014) that can be configured to handle, to some extent, the various
90 stages of a sRNA DE pipeline as well (see Table 1). However, none of these cover the entire analysis
91 of an sRNA dataset.

92 Here we present a comprehensive, interactive processing pipeline for the analysis of sRNA-seq
93 datasets included as part of the UEA small RNA Workbench (Moxon et al 2008; Stocks et al. 2012).

94 The pipeline summarizes approaches for quality checking (Mohorianu et al. 2011; Axtell 2013),
95 normalization (Dillies et al. 2013) and identification of expression-derived patterns (Lopez-Gomollon
96 et al. 2012; Mohorianu et al. 2013). To enable the user to compare sRNA-seq libraries and indicate
97 the level of confidence to place on predictions made during downstream analysis, we also provide a
98 series of diagnostic plots used throughout the pipeline to assess the characteristics and overall quality
99 of the samples. Users can also evaluate different normalization methods in order to decide which
100 approach is suitable for their dataset. In addition, we present a confidence interval (CI)-based
101 approach (Lopez-Gomollon et al. 2012) to summarise the magnitude and direction of fold changes,
102 for each sRNA. On an *H. sapiens* dataset, described in the main text, we demonstrate how this can
103 be extended to multiple comparisons that can be used to group sequences with similar patterns
104 across the whole experiment.

105

106 **RESULTS**

107 In this section, we illustrate the features of our pipeline on a publically available dataset in *H. sapiens*,
108 GSE47532 (Barrett et al. 2013; Camps et al. 2014) to highlight its use to identify characteristics and
109 diagnose problems in real data. Additional examples are presented in Supplementary information 1
110 (example on a *B. terrestris* dataset) and in Supplementary information 3 (example on an *A. thaliana*
111 dataset). The impact of the number of samples or available memory (RAM) on the runtime is
112 discussed in Supplementary information 2.

113 **Workflows and implementation details.**

114 The pipeline is part of the existing UEA small RNA Workbench package (Stocks et al. 2012) which
115 provides a user friendly environment designed for all users regardless of computing experience. The
116 latest version of the workbench also facilitates the chaining together of multiple tools within a workflow.
117 This allows each distinct part of a pipeline to be fully configured prior to runtime forgoing the need for
118 many separate programs that require interlinked inputs/outputs. For example, given a set of sRNA

119 samples, a workflow for the identification of DE sRNAs could consist of the quality checking of the
120 samples, the normalization of expression levels, the identification of differentially expressed,
121 annotated reads and the overview of resulting expression patterns – a diagram illustrating this series
122 of steps is presented in Fig. 1a. Within the workbench interface, the workflow (Fig. 1c) consists of
123 multiple user configurable nodes that represent the various stages in the analysis.

124 A standard pipeline takes as input sequence data in FASTA format with the adapters trimmed. The
125 files can be generated using the adapter removal tool (Stocks et al. 2012) which also allows users to
126 process samples created using the HD sequencing protocol (Sorefan et al. 2012). The next step is
127 the configuration of the workflow using the setup wizard. The first stage is to organise the
128 data/samples in a manner that reflects the original wet lab experimental design. The sample hierarchy
129 is represented as a tree diagram where leaf nodes represent the replicates and the parents represent
130 the individual samples (Fig. 1b). Users then provide a reference genome and an (optional) GFF file,
131 corresponding to the genome build, which will be used for the annotation stage. If an annotation file
132 is provided, users can then choose which annotations are relevant for the analysis.

133 After configuring the sample files, users can choose to begin the workflow immediately or enter each
134 stage of the workflow and change the configurable parameters, as necessary. In addition, during the
135 workflow, users can mark problematic replicates (resulting from the first stage of quality checking) or
136 individual size classes for removal, then select up to six normalization methods to be investigated.
137 The quality check reports are then recreated on the normalised data and can be inspected. Next, the
138 user can select the method that best corrects the data artefacts based on the nuanced characteristics
139 of the dataset's expression distributions.

140 The quality check, normalization, and DE steps are computationally intensive and pose significant
141 demands on both processor and in particular memory (RAM). To counteract this issue, we developed
142 a series of back end improvements, which enable users with a wide range of computing hardware to
143 use the pipeline. More specifically, we employed disk solutions based on relational database
144 management interfaced with a Java front end and interacted via a JavaScript GUI (which is also used

145 to display resulting graphs and tabular results). However, as the use of disk for runtime storage and
146 calculations can have significant impacts on processing time, a RAM only version of the software is
147 also available for users with access to high-end computing hardware.

148 **Quality Checking**

149 To illustrate the quality check stage of the pipeline, initial checks on a *H. sapiens* dataset (H data)
150 were conducted both before and after aligning reads to the reference genome. The first step of the
151 pipeline is to evaluate the overall features of the data being analysed. The sequencing quality of
152 individual sRNA-seq samples is initially assessed based the positional nucleotide composition. Next,
153 the total library size (redundant count) and the total number of unique sequences (non-redundant)
154 count are compared across libraries to assess the variation in sequencing depth. The size class
155 distributions for both redundant and non-redundant reads (Fig 2a1 and 2a2) can indicate abundant or
156 otherwise important sRNA classes early on in the analysis, or identify issues with the sequencing or
157 mapping of certain size classes. The distribution of complexities, defined as the ratio of redundant to
158 non-redundant reads, provides an approximation of the number and abundances of reads in each
159 size class (Fig 2a3). Complexity values that are close to 1 indicate a highly diverse set of low abundant
160 sequences whereas lower complexity values are caused by a smaller set of highly abundant
161 sequences (Mohorianu et al. 2011). For the H dataset we observe a peak in the redundant count
162 distribution at 22-23nt and a sharp and focused decrease in complexity (Figure 2a1 and 2a3). This
163 indicates the presence of a few highly abundant sRNAs for these particular lengths. We also notice
164 that one replicate of the H32 condition contains more unique reads than the other samples for sizes
165 lower than 22nt, and that there is a markedly higher complexity for an H16 replicate across the lower
166 and higher range of size classes, indicating an over representation of read variants.

167 The qualitative replication analysis is conducted through the replicate versus replicate scatter plots
168 and MA plots/Bland-Altman plots (Bland et al. 1986), Fig 2b, with similar characteristics and
169 interpretation to those on microarray data (Bolstad et al. 2003; McCormick et al. 2011; Dillies et al.
170 2013); for the latter each dot corresponds to a gene, in this context each dot represents an sRNA.

171 This comparative analysis can be extended to higher levels (such as at the sample or treatment level)
172 and it should be reviewed again using the normalised expression levels. For the H dataset, this
173 analysis indicated a high consistency for the H32 and H48 replicates and reduced agreement between
174 the H16 replicates. Supporting the initial observation, the most dispersed size-separated fold changes
175 are those found between the replicates of H16 (Fig 2b). Low Jaccard indices generated in the second
176 report indicate that these replicates have poor comparability caused by large differences in both the
177 sequence count distribution and sequence composition of the first replicate (Fig 2c). Since there are
178 only two replicates per treatment and there is no objective approach for choosing one of the two, this
179 plot indicates that the H16 treatment should probably be excluded from further analysis. The other
180 treatments show a high similarity between replicates, with very few fold changes greater than an
181 absolute log₂ fold change of one at higher average expression levels. Although treatment H32 shows
182 a slight skew towards positive fold changes caused by a higher sequencing depth in the second
183 replicate, the pipeline can be used to correct this issue at the normalization stage.

184 The percentage of genome-matching reads is also calculated for both redundant and non-redundant
185 sequences and across size classes (Fig 2a). In addition to examining the entire sRNA population in
186 a dataset, all quality checks described so far can also be calculated and compared visually across
187 individual annotations of interest. These include miRNAs, other ncRNAs (such as tRNAs, rRNAs or
188 snoRNAs), protein-coding genes and repeat/ transposable elements depending on available
189 annotation information (Mortazavi et al. 2008; Xu et al. 2014). These analyses indicate a high
190 proportion of reads in these samples are likely to be miRNAs.

191 **Normalization**

192 The next step in the pipeline is the normalization of the expression levels. In the normalization node
193 we incorporate several existing methods for normalization, with additional features that we have
194 developed especially for sRNA datasets. For scaling-based methods, the normalization total
195 influences the subsequent DE call; ideally it should not be much lower than the original number of
196 reads. For example, if the scaling is done at 1M for samples with >10M reads then all the expression

197 levels will be reduced and DE may be hidden. Alternatively, if scaling is done at 10M for samples with
198 <1M reads, then DE could be artificially be generated. An appropriate normalization total therefore
199 lies in the same range as the sample totals (the average and median options are presented as
200 alternatives). Other options are rank-based quantile normalization adapted for sRNA-seq data
201 (Bolstad et al. 2003) and subsampling normalization (Li et al. 2012).

202 The analysis of the H dataset highlights a common issue with normalization where two replicates are
203 sequenced with different overall depths (mainly due to the characteristics of the sequencing platform
204 employed). To evaluate which method(s) are suitable for this dataset, we tested all six normalization
205 techniques described in the methods section. Figure 3 illustrates the size separated distributions of
206 differential expression which can be used to evaluate the suitability of each normalization method.
207 Fold-changes between replicates should be minimal and produce a distribution centred on zero, after
208 normalization. Whilst the TMM, DESeq2, and quantile methods appear to centre the distributions of
209 all size classes, the total count, subsampling, and upper quartile methods do not improve on the
210 comparability of the distributions. This suggests that for the H data, either TMM, DESeq or quantile
211 normalization should be chosen as normalization approaches.

212 **Differential Expression comparison with existing approaches**

213 We exemplify the DE analysis on the H dataset for two comparisons: N00/H32 and H32/H48; the left
214 hand side is considered the reference sample.

215 We compared our results, obtained using the offset fold change, in log₂ scale (LOFC) and confidence
216 interval (CI) pattern approach – described in the methods, with two of the most widely used tools for
217 detecting DE reads, DESeq2 (Love et al. 2014) and edgeR (Zhou et al. 2014). Both approaches
218 control for false positives by estimating dispersions and weighting fold changes based on these
219 dispersion estimates. For the DESeq2 and edgeR analyses we used a significance cut-off of 0.05.
220 For the method implemented in the workbench, we applied a threshold of 1 LOFC (both for U and D
221 patterns) to call sequences as DE. This was selected based on empirical evidence that a sequence
222 with a log₂ fold change of 1 can be detectable on a northern blot or via qPCR (Morey et al. 2006).

223 The KL divergence curves generated from the H dataset used for determining the appropriate offsets
224 are shown in Figure 4. We also assessed the dependence of the offset on the number of strand bias
225 bins and length of the alignment window. In the H dataset the number of strand bias bins heavily
226 affected the resulting offset up to 100 bins, after which point the KL curve remains unchanged which
227 resulted in an offset biased towards the lower end of abundance levels. The offset was also affected
228 by alignment window length and can vary erratically when using the raw measures; however, we
229 utilize a LOESS smoothing function (Cleaveland 1979) to produce a more stable offset.

230 For the N00 vs H32 comparison of the H dataset, 427 sequences were called DE by all methods (Fig
231 5b). DESeq2 and edgeR both predicted 241 sequences not called DE by the LOFC method; DESeq2
232 returned 110 differentially expressed sequences that edgeR and our method did not find significant
233 and edgeR predicted 15 differentially expressed sequences which were not captured using the other
234 methods. Based on the MA plot (Fig 5a) we observe that the abundance and/or offset fold-change of
235 these specific calls low. These artefacts can be identified and evaluated on a case-by-case basis by
236 using the LOFC and the CI approach. In addition, we present the expression levels of the 4 reads
237 identified exclusively using the LOFC approach (Fig 5c).

238 **DISCUSSION**

239 We have described a sRNA processing pipeline, part of the UEA sRNA Workbench, that includes
240 steps for quality checking, normalization, and identification of DE sRNAs considering the unique
241 characteristics of sRNA-seq datasets. To achieve a better understanding of these datasets, the
242 pipeline generates a set of diagnostic plots, which can be used initially to review the raw sequencing
243 quality of the replicates and then to assess the effect that different normalization techniques have of
244 the abundance distributions. The use of a suitable normalization is essential for reducing false positive
245 predictions; however no single normalization technique can be invariably applied to all sRNA-seq
246 datasets. To evaluate which approach is appropriate for a given dataset (i.e. by rendering the samples
247 comparable from most (preferably all) quality check angles) we encourage the user to investigate their
248 using the revised quality check plots.

249 When identifying DE transcripts in HTS data it is important to take into account the level of noise, a
250 quantity that increases with the depth of sequencing. To account for this, we have implemented a
251 user-friendly tool for the identification of a suitable offset, which estimates the abundance range of
252 the reads lacking sRNA characteristics (e.g. specific size), taking into account the sequencing depth
253 and the characteristics of the sRNA population present in the samples. We compared the results of
254 our DE analysis (LOFC) to that of DESeq2 and edgeR DE packages to determine the level of overlap
255 between other methods and our own. In lieu of a p-value threshold to assess DE genes, which often
256 reports large numbers of significant genes often with a low difference in expression, we used a cut-
257 off of 1 LOFC to filter the reported sequences. The cut-off can, however, be user defined in order to
258 reduce or increase the number of reported sequences. Importantly, the ranking of sequences by
259 LOFC is not populated with high but insignificant fold changes.

260 To further accommodate the variability between replicates we use CIs created over normalized
261 replicate expression levels which produce more stringent lists of DE sequences between treatments.
262 The method is also extended to multiple conditions by using pattern-based grouping of the sequences
263 (Fig 6). The method is not only suitable for (ordered) time-series datasets, but can also be applied to
264 other types of comparative experiments such as wild type versus multiple treatments or cross tissue
265 comparison. Grouping DE sequences allows users to quickly view sets of sRNAs that follow the same
266 pattern of expression throughout the experiment.

267 During our analyses we observed that problematic datasets arise when whole size classes are
268 affected by a condition, causing a high rate of DE for a large proportion of the sRNAs e.g. RNAi
269 mutants which cause the exclusion of a whole class of sRNAs or virus infections which produce a
270 large set of viral siRNAs in the infected samples (Szittyá et al. 2010). To our knowledge no current
271 normalization is able to correct for such experiments, and further approaches will need to be
272 developed to provide an appropriate normalization solution to this kind of data.

273 In conclusion, we have described a user-friendly pipeline for sRNA DE analysis which allows the
274 evaluation of a variety of techniques to identify the most suitable approach for a given dataset. The

275 workbench includes both established approaches and tools that we have specifically developed for
276 sRNA sequence analysis and facilitates a coherent and informed analysis through linking the different
277 aspects into a workflow. The UEA sRNA Workbench and the pipeline design devised for the data
278 analysis may prove to be a valuable resource facilitating the expansion of our knowledge of sRNAs,
279 especially for the study of new or less well characterised classes of sRNAs.

280 **MATERIALS**

281 To illustrate the use of the pipeline we use a *H. sapiens* dataset referred to as "H" data (publicly
282 available on Gene Expression Omnibus (GEO) under accession number GSE47532). This is an
283 experiment on the effects of hypoxic conditions on MCF7 cells (Camps et al. 2014), organised into a
284 time series of four points, each with two biological replicates, Normoxia (N00), Hypoxia at 16 hours
285 (H16), Hypoxia at 32 hours (H32), and Hypoxia at 48 hours (H48). The additional examples presented
286 in the supplementary information are based on publicly available *B. terrestris* data (GSE64512)
287 consisting of two samples, with four biological replicates each (Sadd et al. 2012) and a publicly
288 available *A. thaliana* data (GSE35562, GSM1178880 to GSM1178882 for the wild-type and
289 GSM1178883 to GSM1178885 for the Hen1-8 mutant) consisting of two samples, with three biological
290 replicates each (Zhai et al. 2013).

291 **METHODS**

292 In this section we describe the methodology and software underpinning the new pipeline; the main
293 workflow diagram is presented in the diagram in Figure 2a.

294 **Quality checking**

295 The sequencing quality of individual sRNA-seq samples is assessed based on properties such as the
296 positional nucleotide composition (Consortium 2014), sequencing depth and the number of unique
297 sequences present in a sample (Rajagopalan et al. 2006). The accuracy of expression replication is
298 evaluated by comparing, qualitatively and quantitatively, the abundances of reads between replicates
299 (Mapleson et al. 2014). The quantitative analysis includes the study of size-class separated

300 distributions of abundances and complexities, defined as the ratio of unique (non-redundant) to total
301 (redundant) reads and the Jaccard similarity index on the top 500 most abundant reads (Mohorianu
302 et al. 2011; Jaccard 1901). The qualitative comparison is conducted through the replicate versus
303 replicate scatter plots and MA plots /Bland-Altman plots (Bland, Altman 1986). We also assess the
304 stability of distributions of fold changes between replicate libraries for each size class presented in
305 \log_2 scale. An appropriate similarity between the compared replicates/ samples is indicated by tight
306 distributions, symmetric on 0 \log_2 fold change with no deviations for any particular size classes (Yang
307 et al. 2002; Mohorianu et al. 2011). The percentage of genome-matching reads is calculated for both
308 redundant and non-redundant sequences and across size classes. Selected annotations for which
309 similar checks are performed include miRNAs, other ncRNAs (such as tRNAs, rRNAs or snoRNAs),
310 protein-coding genes and repeat/ transposable elements depending on available annotation
311 information (Xu et al. 2014; Omidvar et al. 2015).

312 Abundance distributions of reads in each sample are plotted in a series of boxplots (McCormick et al.
313 2011; Dillies et al. 2013). However, due to the high proportion of low abundance reads characteristic
314 to sRNA-seq data these distributions for all reads are often uninformative. To counter this, we break
315 the data into abundance ranges of user defined length (referred to as abundance windows) and
316 assess the comparability of the sample distributions within each window.

317 **Normalization**

318 The aim of the normalization of the expression levels is to minimize the technical variation between
319 replicates and treatments which is not biologically relevant e.g. sequencing errors and biases or
320 artefacts from the RNA itself (Sorefan et al. 2012; Raabe et al. 2014) since DE predictions are only
321 considered reliable when the variability between replicates is lower than the differences between the
322 treatments. In the Normalization component of the pipeline, we incorporate several existing methods
323 for normalization (scaling-based, rank-based and statistical), with additional features, adapted for
324 sRNA datasets. Scaling normalizations, based on the identification of a scaling factor which brings
325 the total number of reads to an a priori fixed total include: the reads per million (RPM)/ reads per total

326 (RPT) method (Mortazavi et al. 2008) for which the total abundance of all reads in a sample is
327 considered, upper quartile normalization (Bullard et al. 2010) for which only the reads with
328 abundances in the upper quartile are considered, the trimmed mean of M-values (TMM) (Anders et
329 al. 2010) and DESeq (Anders, Huber 2010).

330 Quantile normalization (Bolstad et al. 2003), originally designed for microarray experiments, is also
331 included as an option in the pipeline. This method imposes the same distribution of ranks over all
332 sequences in the dataset. We adapted this method to sRNA sequencing data by adding two extra
333 conditions: (1) if, within a sample, two or more reads have the same abundance before normalization,
334 they are assigned the same abundance after normalization which is the average of the normalised
335 abundances. (2) If a read is not present in the original sample (abundance=0) then it is assigned an
336 expression level of 0 in its normalised version.

337 We also include a subsampling-based normalization which is an adapted version of the method
338 described in (Li et al. 2012). Our method is based on sampling reads (without replacement) to the
339 minimum library size (for all samples that pass the quality check). It consists of two steps: (1) to ensure
340 that the distribution of abundances are consistent within a sample, the sampling is conducted for
341 decreasing proportions until the sample's distribution has significantly changed or the lowest sample
342 size has been reached; (2) a subsample of reads with a fixed total is selected repeatedly and, using
343 bootstrapping, the variability of the subsamples is tested. If the variability is low, a random sample
344 (representative for the distribution, i.e. not an outlier) is selected.

345 **Differential Expression call**

346 To identify DE sequences between conditions/treatments, the pipeline includes a confidence interval
347 (CI) based approach (Lopez-Gomollon et al. 2012; Mohorianu et al. 2013). For each sequence, in
348 each condition, a CI is calculated over replicate expressions using either Chebyshev's intervals
349 calculated from the mean and the standard deviation (Singh et al. 2006) or the minimum and
350 maximum expression levels if only two replicates are available. For a selected comparison between
351 a reference and observed condition, the direction of DE and its amplitude are also calculated. A

352 directional descriptor from the set {up (U), down (D), straight (S)} is assigned to each sequence as
353 follows: S is used if the CIs overlap, U indicates that the observed CI is higher than the reference, and
354 D indicates the opposite result. The issue of performing pairwise comparisons with sample counts
355 greater than two can then be addressed by forming patterns using the {U,D,S} descriptors. This allows
356 sorting and filtering of sequences that result in potentially relevant/interesting expression changes
357 throughout the course of the experiment.

358 The amplitude of the difference in expression between conditions is considered on proximate
359 extremes (the closest ends of the neighbouring CI) of the reference and observed CIs and is only
360 calculated on sequences that have been assigned an U or D descriptor. The amplitude is calculated
361 using the \log_2 Offset Fold Change (LOFC) method previously described in (Mohorianu et al. 2011;
362 Mohorianu et al. 2013). The offset prevents low abundance variation from being included in the
363 significant DE distribution. The aim of the offset-approach is to reduce the number of false positives
364 from low abundance sequences and to allow fold change values to be used directly when assessing
365 the relative significance of differentially expressed sequences.

366 To determine an appropriate offset for a dataset, the pipeline can be used to estimate the abundance
367 level around which the majority of noise-related reads lie. Previous studies have observed that low
368 abundance regions/loci have a high strand bias (derived from the reduced number of reads), but loci
369 within the noise to signal range have no preferred strand bias (Mohorianu et al. 2011). Based on this
370 observation, the method assigns sRNAs to windows of a set length along the genome reference and
371 the total expression and strand bias is then calculated for each window. For all expression levels, the
372 distribution of strand biases is compared to a random uniform distribution using the Kullback-Leibler
373 (KL) divergence measure (Kullback et al. 1951). We define the noise to signal threshold (the offset)
374 as the value for which the global minimum of the KL divergence distribution is reached. The
375 distribution is smoothed by a LOESS function (Cleaveland 1979) to prevent expression level outliers
376 from giving a local minimum. Expression levels lower than this threshold tend to have a higher
377 divergence from a uniform strand bias due to low number of incident reads, and expression levels

378 that are higher than the threshold have an increasing divergence measure due to biologically relevant
379 reads.

380 **Availability**

381 The workbench and all the supporting data and tutorials are freely available from the website
382 <http://srna-workbench.cmp.uea.ac.uk>. The licence is a custom licence written for the UEA sRNA
383 Workbench and can be found in the Workbench installation directory or by visiting the following web
384 link <http://srna-workbench.cmp.uea.ac.uk/wp-content/uploads/2016/11/sRNA-WorkbenchEULA.pdf>.
385 There are no restrictions on use other than requiring citations to specific papers when conducting
386 research with the software; specific details can be found on the website.

387 **ACKNOWLEDGEMENTS**

388 The authors would like to thank the members of the Dalmy and Moulton labs for constructive
389 discussions and suggestions.

390 **FUNDING**

391 This study was supported by the Biotechnology and Biological Sciences Research Council: grant numbers
392 BB/L003139/1, BB/L021269/1.

393

394 **References**

- 395 Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome Biol*
396 11:R106.
- 397 Anders S, McCarthy DJ, Chen Y, Okoniewski M, Smyth GK, Huber W, Robinson MD. 2013. Count-
398 based differential expression analysis of RNA sequencing data using R and Bioconductor. *Nat Protoc*
399 8:1765-1786.
- 400 Axtell MJ. 2013. ShortStack: comprehensive annotation and quantification of small RNA genes. *RNA*
401 19:740-751.
- 402 Barquist L, Vogel J. 2015. Accelerating Discovery and Functional Analysis of Small RNAs with New
403 Technologies. *Annu Rev Genet* 49:367-394.
- 404 Barrett T, Wilhite SE, Ledoux P, et al. 2013. NCBI GEO: archive for functional genomics data sets--
405 update. *Nucleic Acids Res* 41:D991-995.
- 406 Bartel DP. 2009. MicroRNAs: target recognition and regulatory functions. *Cell* 136:215-233.
- 407 Bland JM, Altman DG. 1986. Statistical methods for assessing agreement between two methods of
408 clinical measurement. *Lancet* 1:307-310.
- 409 Bolstad BM, Irizarry RA, Astrand M, Speed TP. 2003. A comparison of normalization methods for high
410 density oligonucleotide array data based on variance and bias. *Bioinformatics* 19:185-193.
- 411 Brodersen P, Voinnet O. 2006. The diversity of RNA silencing pathways in plants. *Trends Genet*
412 22:268-280.
- 413 Bullard JH, Purdom E, Hansen KD, Dudoit S. 2010. Evaluation of statistical methods for normalization
414 and differential expression in mRNA-Seq experiments. *BMC Bioinformatics* 11:94.
- 415 Camps C, Saini HK, Mole DR, et al. 2014. Integrated analysis of microRNA and mRNA expression
416 and association with HIF binding reveals the complexity of microRNA expression regulation under
417 hypoxia. *Mol Cancer* 13:28.
- 418 Carthew RW, Sontheimer EJ. 2009. Origins and Mechanisms of miRNAs and siRNAs. *Cell* 136:642-
419 655.
- 420 Chen X. 2012. Small RNAs in development - insights from plants. *Curr Opin Genet Dev* 22:361-367.
- 421 Cleveland WS. 1979. Robust Locally Weighted Regression and Smoothing catterplots. *Journal of*
422 *the American Statistical Association* 74:829-836.
- 423 Consortium SM-I. 2014. A comprehensive assessment of RNA-seq accuracy, reproducibility and
424 information content by the Sequencing Quality Control Consortium. *Nat Biotechnol* 32:903-914.
- 425 Davis MP, van Dongen S, Abreu-Goodger C, Bartonicek N, Enright AJ. 2013. Kraken: a set of tools
426 for quality control and analysis of high-throughput sequence data. *Methods* 63:41-49.
- 427 Dillies MA, Rau A, Aubert J, et al. 2013. A comprehensive evaluation of normalization methods for
428 Illumina high-throughput RNA sequencing data analysis. *Brief Bioinform* 14:671-683.
- 429 Donaszi-Ivanov A, Mohorianu I, Dalmay T, Powell PP. 2013. Small RNA analysis in Sindbis virus
430 infected human HEK293 cells. *PLoS One* 8:e84070.
- 431 Garber M, Grabherr MG, Guttman M, Trapnell C. 2011. Computational methods for transcriptome
432 annotation and quantification using RNA-seq. *Nat Methods* 8:469-477.

- 433 Gupta V, Markmann K, Pedersen CN, Stougaard J, Andersen SU. 2012. shortran: a pipeline for small
434 RNA-seq data analysis. *Bioinformatics* 28:2698-2700.
- 435 Jaccard P. 1901. Etude comparative de la distribution florale dans une portion des Alpes et de Jura.
436 *Bulletin de la Societe Vaudoise des Sciences Naturelles*.
- 437 Kullback S, Leibler RA. 1951. On Information and Sufficiency. *The Annals of Mathematical Statistics*
438 22:79-86.
- 439 Li J, Witten DM, Johnstone IM, Tibshirani R. 2012. Normalization, testing, and false discovery rate
440 estimation for RNA-sequencing data. *Biostatistics* 13:523-538.
- 441 Lippman Z, Martienssen R. 2004. The role of RNA interference in heterochromatic silencing. *Nature*
442 431:364-370.
- 443 Lohse M, Bolger AM, Nagel A, Fernie AR, Lunn JE, Stitt M, Usadel B. 2012. RobiNA: a user-friendly,
444 integrated software solution for RNA-Seq-based transcriptomics. *Nucleic Acids Res* 40:W622-627.
- 445 Lopez-Gomollon S, Mohorianu I, Szittyá G, Moulton V, Dalmay T. 2012. Diverse correlation patterns
446 between microRNAs and their targets during tomato fruit development indicates different modes of
447 microRNA actions. *Planta* 236:1875-1887.
- 448 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for RNA-seq
449 data with DESeq2. *Genome Biol* 15:550.
- 450 Mapleson D, Mohorianu I, Pais H, Stocks MB, Folkes L, Moulton V. 2014. Processing Large-scale
451 Small RNA datasets in silico. *Next-generation Sequencing: Current Technologies and Applications:*
452 *Caister Academic Press*.
- 453 McCormick KP, Willmann MR, Meyers BC. 2011. Experimental design, preprocessing, normalization
454 and differential expression analysis of small RNA sequencing experiments. *Silence* 2:2.
- 455 Meister G. 2013. Argonaute proteins: functional insights and emerging roles. *Nat Rev Genet* 14:447-
456 459.
- 457 Mohorianu I, Moulton V. 2010. Revealing biological information using data structuring and automated
458 learning. *Recent Pat DNA Gene Seq* 4:181-191.
- 459 Mohorianu I, Schwach F, Jing R, Lopez-Gomollon S, Moxon S, Szittyá G, Sorefan K, Moulton V,
460 Dalmay T. 2011. Profiling of short RNAs during fleshy fruit development reveals stage-specific
461 sRNAome expression patterns. *Plant J* 67:232-246.
- 462 Mohorianu I, Stocks MB, Wood J, Dalmay T, Moulton V. 2013. CoLlde: a bioinformatics tool for CO-
463 expression-based small RNA Loci Identification using high-throughput sequencing data. *RNA Biol*
464 10:1221-1230.
- 465 Molnar A, Schwach F, Studholme DJ, Thuenemann EC, Baulcombe DC. 2007. miRNAs control gene
466 expression in the single-cell alga *Chlamydomonas reinhardtii*. *Nature* 447:1126-1129.
- 467 Morey JS, Ryan JC, Van Dolah FM. 2006. Microarray validation: factors influencing correlation
468 between oligonucleotide microarrays and real-time PCR. *Biol Proced Online* 8:175-193.
- 469 Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. 2008. Mapping and quantifying mammalian
470 transcriptomes by RNA-Seq. *Nat Methods* 5:621-628.
- 471 Moxon S, Schwach F, Dalmay T, Maclean D, Studholme DJ, Moulton V. 2008. A toolkit for analysing
472 large-scale plant small RNA datasets. *Bioinformatics*. 1;24(19):2252-3. doi: 10.1093/bioinformatics/

- 473 btn428.Omidvar V, Mohorianu I, Dalmay T, Fellner M. 2015. Identification of miRNAs with potential
474 roles in regulation of anther development and male-sterility in 7B-1 male-sterile tomato mutant. *BMC*
475 *Genomics* 16:878.
- 476 Oszolak F, Milos PM. 2011. RNA sequencing: advances, challenges and opportunities. *Nat Rev*
477 *Genet* 12:87-98.
- 478 Raabe CA, Tang TH, Brosius J, Rozhdestvensky TS. 2014. Biases in small RNA deep sequencing
479 data. *Nucleic Acids Res* 42:1414-1426.
- 480 Rajagopalan R, Vaucheret H, Trejo J, Bartel DP. 2006. A diverse and evolutionarily fluid set of
481 microRNAs in *Arabidopsis thaliana*. *Genes Dev* 20:3407-3425.
- 482 Rapaport F, Khanin R, Liang Y, Pirun M, Krek A, Zumbo P, Mason CE, Socci ND, Betel D. 2013.
483 Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data.
484 *Genome Biol* 14:R95.
- 485 Sadd BM, Barribeau SM, Bloch G, de Graaf DC et al. The genomes of two key bumblebee species
486 with primitive eusocial organization. *Genome Biol* 2015 Apr 24;16:76
- 487 Singh A, Maichle R, Lee S. 2006. On the computation of 95% upper confidence limit of the unknown
488 population mean based upon data sets with below detection limit observations. US Environmental
489 Protection Agency, Office of Research and Development.
- 490 Sonesson C, Delorenzi M. 2013. A comparison of methods for differential expression analysis of RNA-
491 seq data. *BMC Bioinformatics* 14:91.
- 492 Sorefan K, Pais H, Hall AE, Kozomara A, Griffiths-Jones S, Moulton V, Dalmay T. 2012. Reducing
493 ligation bias of small RNAs in libraries for next generation sequencing. *Silence* 3:4.
- 494 Stocks MB, Moxon S, Mapleson D, Woolfenden HC, Mohorianu I, Folkes L, Schwach F, Dalmay T,
495 Moulton V. 2012. The UEA sRNA workbench: a suite of tools for analysing and visualizing next
496 generation sequencing microRNA and small RNA datasets. *Bioinformatics* 28:2059-2061.
- 497 Szittyá G, Moxon S, Pantaleo V, Toth G, Rusholme Pilcher RL, Moulton V, Burgyan J, Dalmay T.
498 2010. Structural and functional analysis of viral siRNAs. *PLoS Pathog* 6:e1000838.
- 499 Studholme DJ 2012. Deep sequencing of small RNAs in plants: applied bioinformatics. *Brief Funct*
500 *Genomics*. 11(1):71-85. doi: 10.1093/bfgp/elr039. Vidal EA, Moyano TC, Krouk G, Katari MS,
501 Tanurdzic M, McCombie WR, Coruzzi GM, Gutierrez RA. 2013. Integrated RNA-seq and sRNA-seq
502 analysis identifies novel nitrate-responsive genes in *Arabidopsis thaliana* roots. *BMC Genomics*
503 14:701.
- 504 Voinnet O. 2009. Origin, biogenesis, and activity of plant microRNAs. *Cell* 136:669-687.
- 505 Xu P, Mohorianu I, Yang L, Zhao H, Gao Z, Dalmay T. 2014. Small RNA profile in moso bamboo root
506 and leaf obtained by high definition adapters. *PLoS One* 9:e103590.
- 507 Yang YH, Dudoit S, Luu P, Lin DM, Peng V, Ngai J, Speed TP. 2002. Normalization for cDNA
508 microarray data: a robust composite method addressing single and multiple slide systematic variation.
509 *Nucleic Acids Res* 30:e15.
- 510 Zhai J, Zhao Y, Simon SA, Huang S et al. Plant microRNAs display differential 3' truncation and tailing
511 modifications that are ARGONAUTE1 dependent and conserved across species. *Plant Cell* 2013
512 Jul;25(7):2417-28.

513 Zhou X, Lindsay H, Robinson MD. 2014. Robustly detecting differential expression in RNA
514 sequencing data using observation weights. *Nucleic Acids Res* 42:e91.

515 Zhu H, Zhou Y, Castillo-Gonzalez C, et al. 2013. Bidirectional processing of pri-miRNAs with branched
516 terminal loops by *Arabidopsis* Dicer-like1. *Nat Struct Mol Biol* 20:1106-1115.

517

518

519 **TABLE AND FIGURES LEGENDS**

520 **Table 1. A summary of current tools designed for RNA-seq analysis, which can be applied for**
521 **sRNA-seq analyses.** For each tool we present the type of expected input (e.g. mRNA-seq, sRNA-
522 seq, etc), the availability of quality checks, analysis of the nucleotide distributions and possibility of
523 adapter trimming. Additional features include the evaluation of size class distributions and MA or
524 scatter plots. Higher level checks such as the annotation of reads, normalization of abundances and
525 differential expression calls are also reviewed.

526 **Figure 1. Overview of the analysis pipeline and the input of the Differential expression**
527 **workflow implemented in the UEA sRNA Workbench.** (a) Diagram showing the steps of
528 the pipeline, including the pre-processing, alignment to the reference genome and available
529 annotations, quality checking of the raw and processed data, normalization and differential
530 expression call; (b) hierarchical representation of the input data obtained using the input
531 wizard (c) the user interface for a workflow containing Quality Checks, Normalization and
532 Differential Expression call; each node can configured individually.

533 **Figure 2. Quality checks for the H dataset.** (a) The characterization of reads within a sample can
534 be obtained by creating the size class distributions for redundant (a.1) and non-redundant
535 (a.2) reads. Next, the ratio of unique to total reads can be investigated using the complexity
536 distribution (a.3). Lastly, the proportions of genome matching reads for redundant (a.4) and
537 non-redundant (a.5) reads highlights the quality of the sRNA library. (b) MA plots on the raw
538 abundances (prior to any normalization or filtering) for evaluating the reproducibility of the
539 replicates. On the x-axis we represent the average abundance between the replicates; on the
540 y-axis we represent the fold changes. Good samples show low variability with the increase of
541 abundance (e.g. N00, H32 and H48); problematic samples are characterized by high variability
542 between replicates (e.g. H16). (c) Jaccard similarity indexes computed on the top 1000 most
543 abundant reads. These indicate a high reproducibility between the N00, H32 and H48

544 replicates (in excess of 0.8) and a low reproducibility for the H16 replicates (0.62).
 545 Interestingly, the second H16 replicates is more similar to the first replicate in H32 time point.

546 **Figure 3. Evaluation of the appropriateness of the normalization methods on the H dataset.** For
 547 each sample and for each set of replicates we represent the fold change distributions (y-axis)
 548 for each individual size class (x-axis). Based on the assumption that no significant differences
 549 are expected between replicates, a suitable normalization is one which brings all distributions
 550 on the 0 line (in log2 scale, this corresponds to equal values in both replicates). For the H
 551 dataset, the TMM, DeSeq2 and the adapted quantile normalization fulfil this criterion for all
 552 samples.

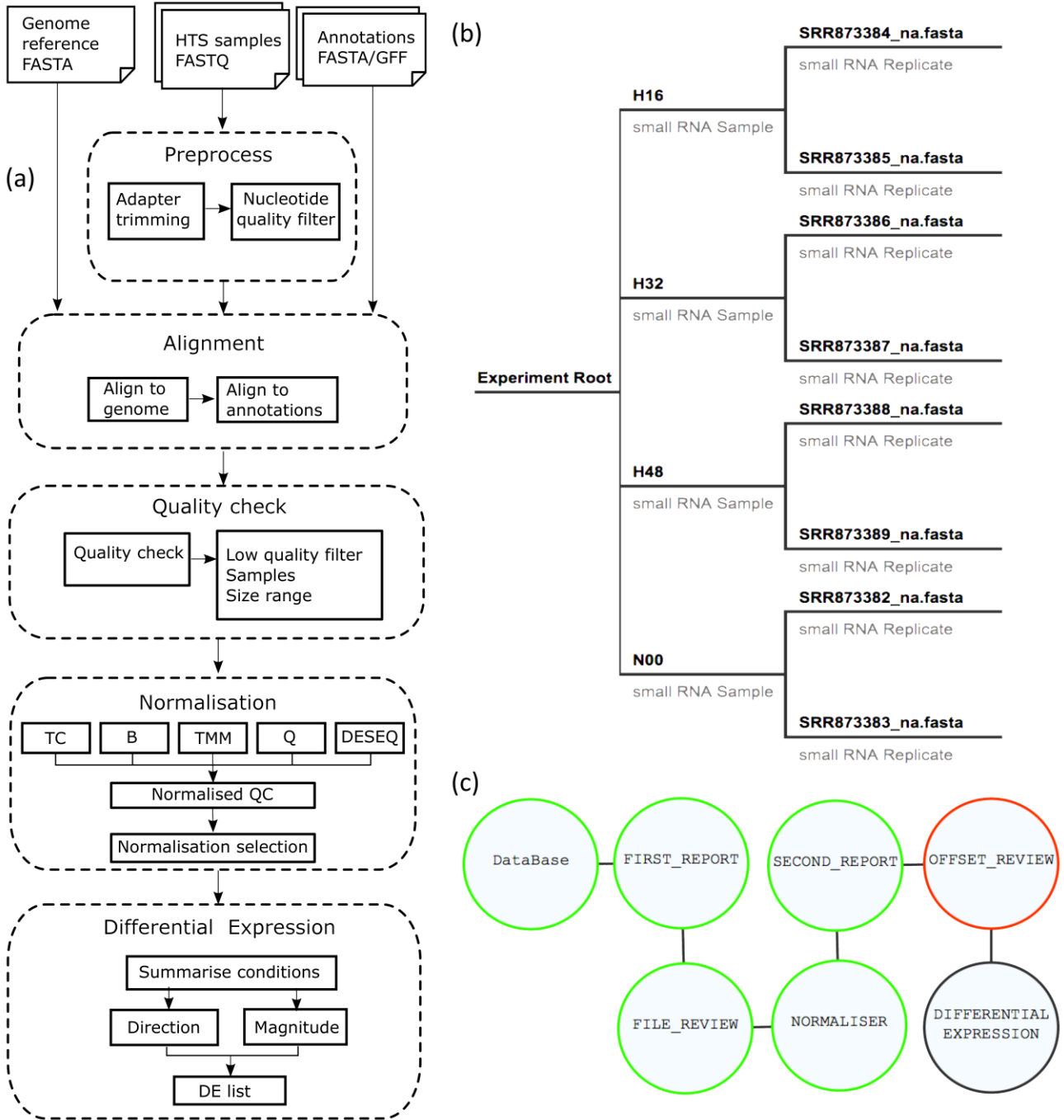
553 **Figure 4. Identification of an offset for each sample in the H dataset using Kullback-Leibler**
 554 **divergence to compare the strand bias distributions of reads to a random uniform**
 555 **distribution, in windows of various lengths.** This analysis was done on windows of length
 556 1000nt (parameter which can be modified from the GUI), for each of the two replicates (_1
 557 and _2) of the three accepted samples, N00, H32 and H48. On the x-axis we show the
 558 abundances of the considered windows (the abundance within a window is the algebraic sum
 559 of the abundances of all incident reads); on the y-axis we represent the value of the KL
 560 divergence. The grey line indicates the unsmoothed KL divergence values and the blue line
 561 shows the divergence values smoothed by loess (span=0.3). The offset for each sample is
 562 determined as the minimum of the smoothed divergence. The offset for the whole dataset is
 563 the overall minimum of these values, for this dataset this value was determined to be 42.

564 **Figure 5. Assessment of three approaches used for the identification differentially expressed**
 565 **reads applied on the N00 vs H32 comparison (H dataset).** (a) MA plot created using the
 566 normalized expression levels (TMM method, see figure 3). On the x-axis we represent the
 567 average abundance; on the y-axis we represent the log2(OFC). The colour of the dots
 568 indicates whether the reads were called DE by both edgeR and DESeq2 (orange), exclusively
 569 by edgeR (blue), or exclusively by DESeq2 (red). Reads accepted as DE using the LOFC

570 approach are those outside the dotted lines. (b) Venn diagram showing the number of reads
571 called DE using the LOFC, edgeR and DESeq2 methods. (c) distributions of expression levels
572 (represented as maximal intervals) for the 4 sequences called DE exclusively by the LOFC
573 method.

574 **Figure 6. Clusters of reads sharing similar patterns** (only the clusters with more than 15 entries
575 were presented; the SS cluster was excluded, since the vast majority of the reads are not
576 expected to be differentially expressed between treatments). The U and D descriptors were
577 assigned to reads for which the LOFC on the proximal ends of the maximal expression
578 intervals was in excess of 1. Each line corresponds to the averaged expression profile, on the
579 two available replicates, for one sRNA; the red lines are used to highlight miRNA expression
580 profiles. The boxplot inter-quartile ranges (IQRs) are used to highlight the distributions of
581 expression in each time point and underline the pattern.

582

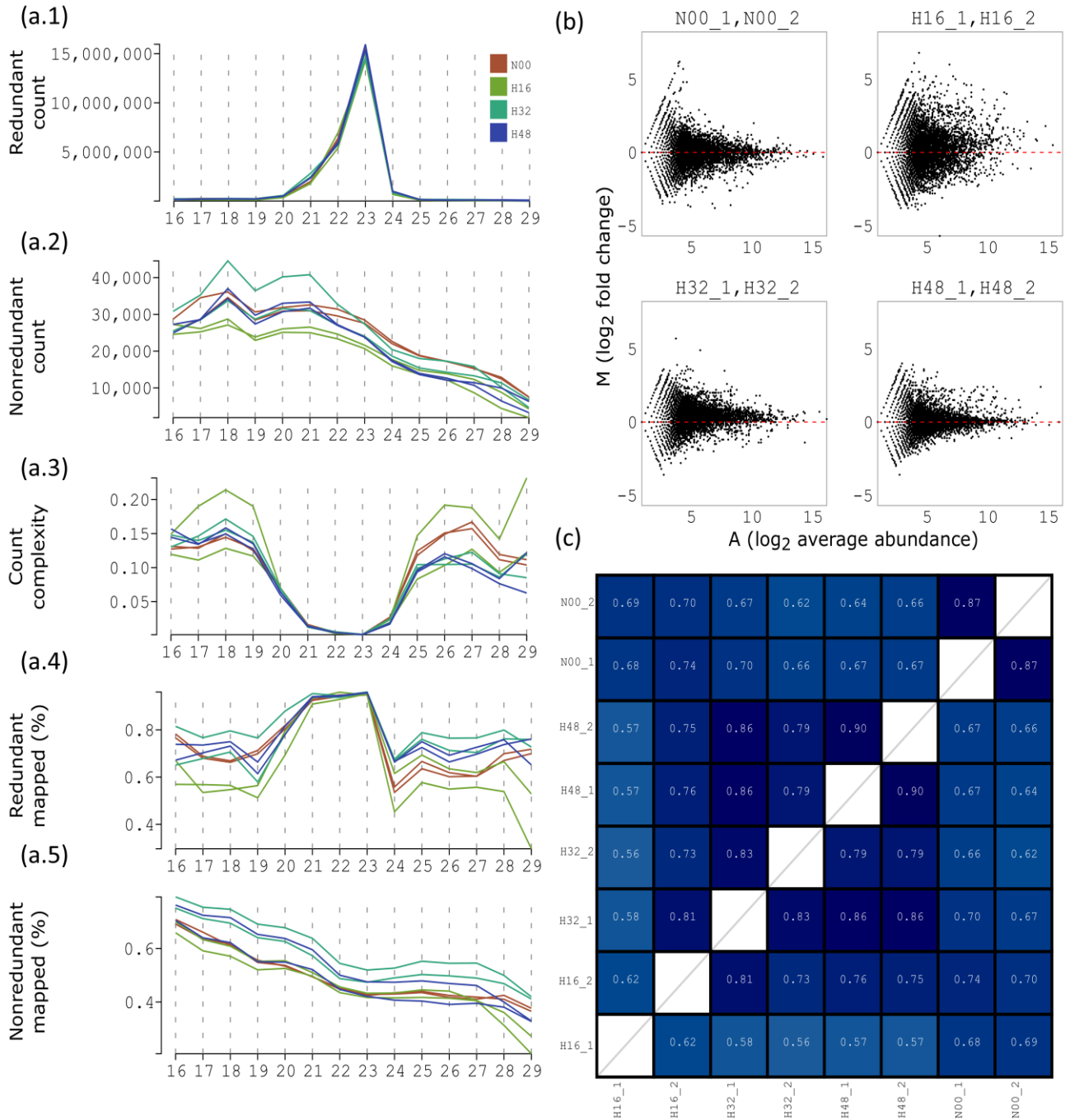


583

584 Figure 1.

585

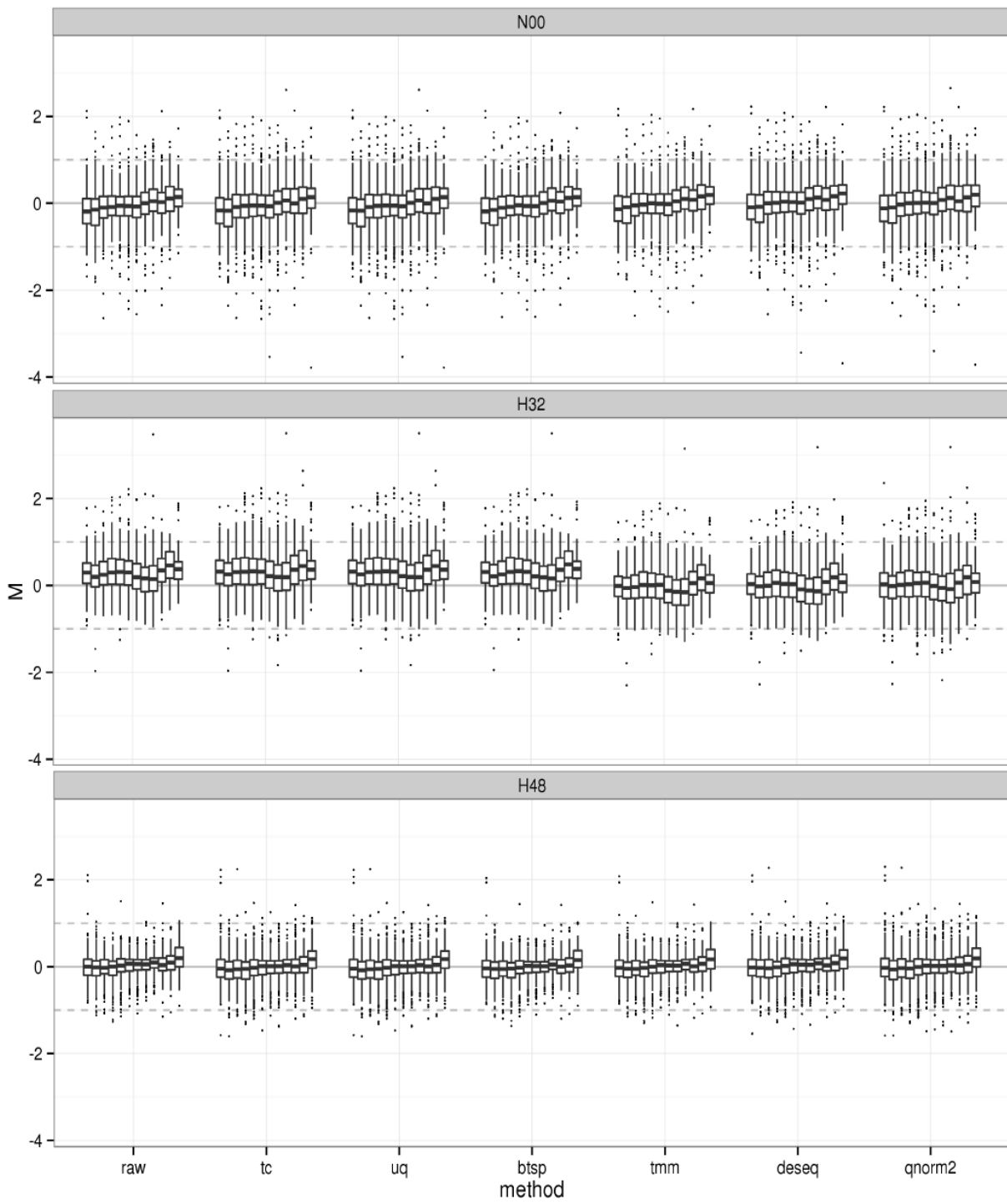
586



587

588 Figure 2

589

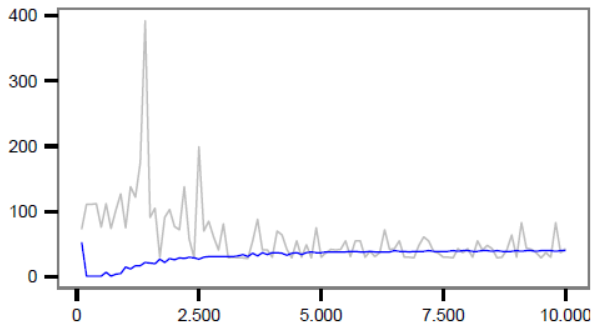


590

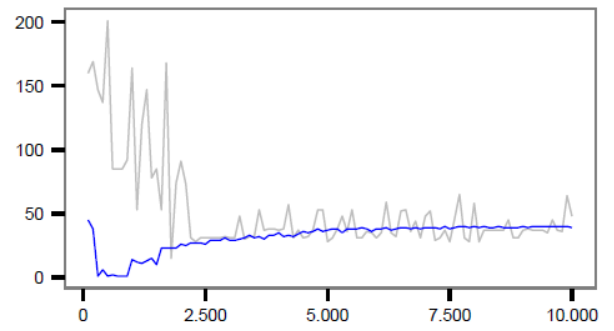
591 Figure 3

592

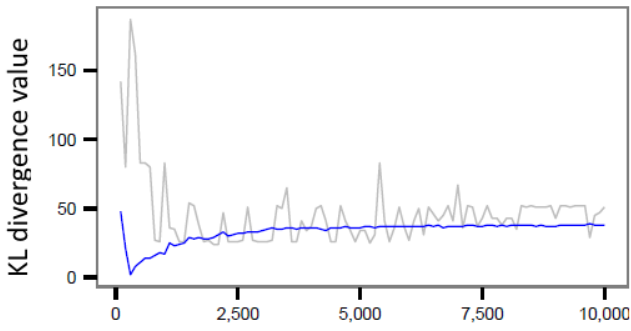
N00_1



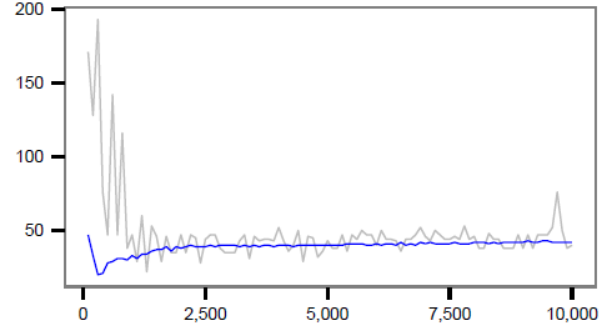
N00_2



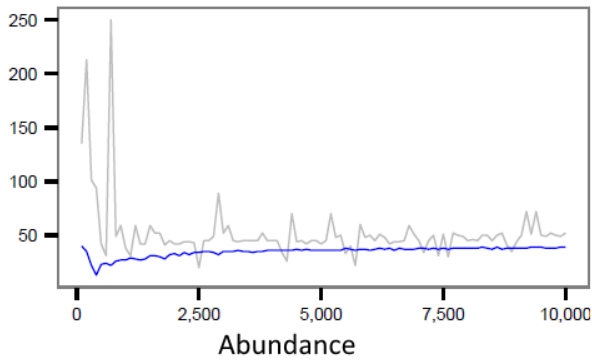
H32_1



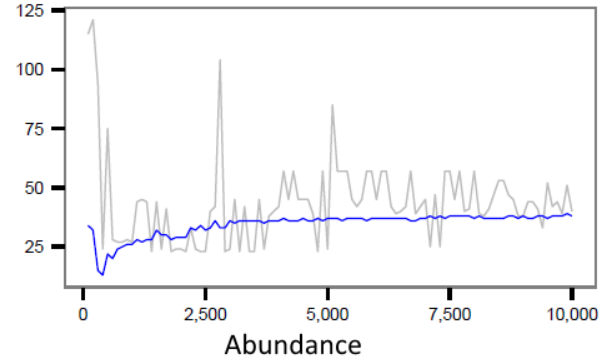
H32_2



H48_1



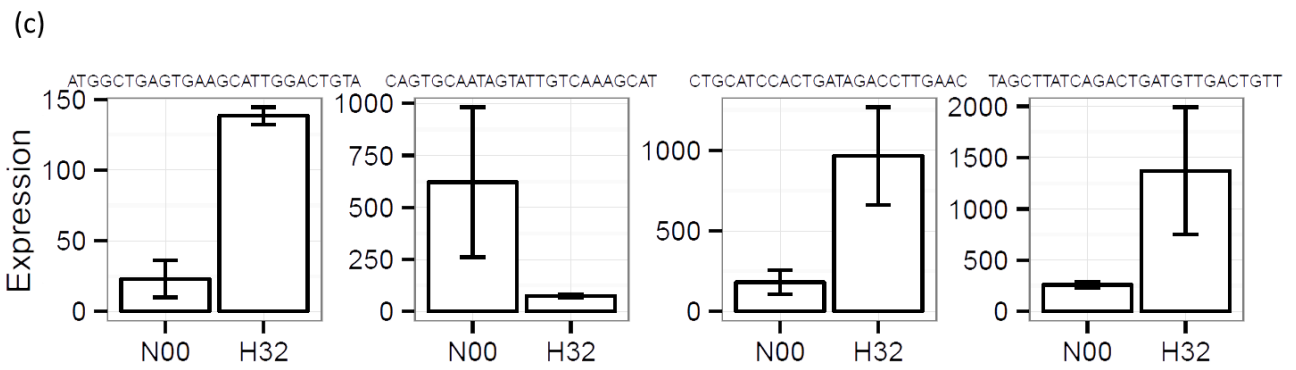
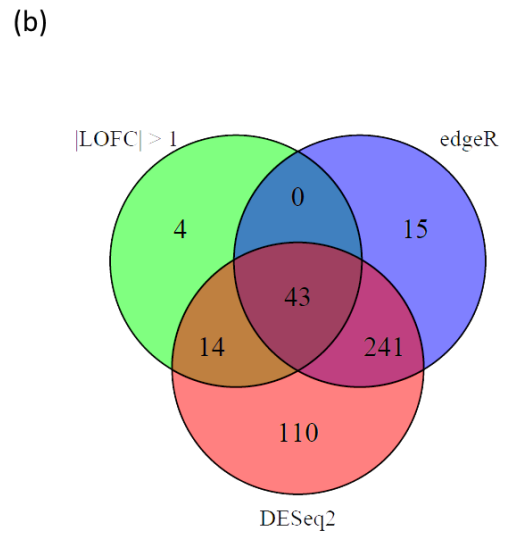
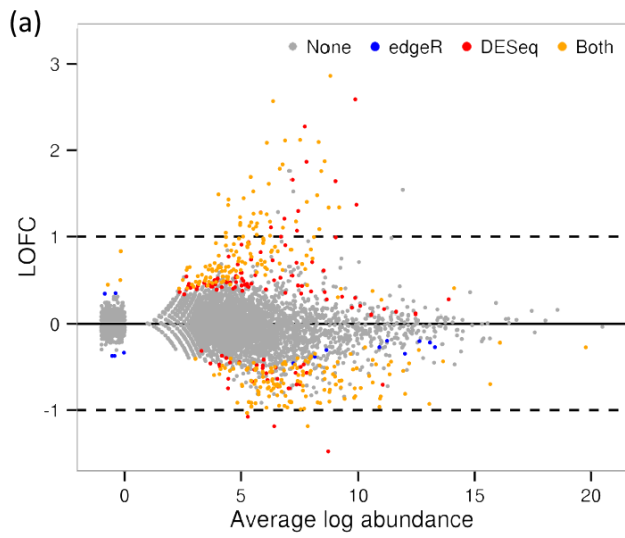
H48_2



593

594 Figure 4

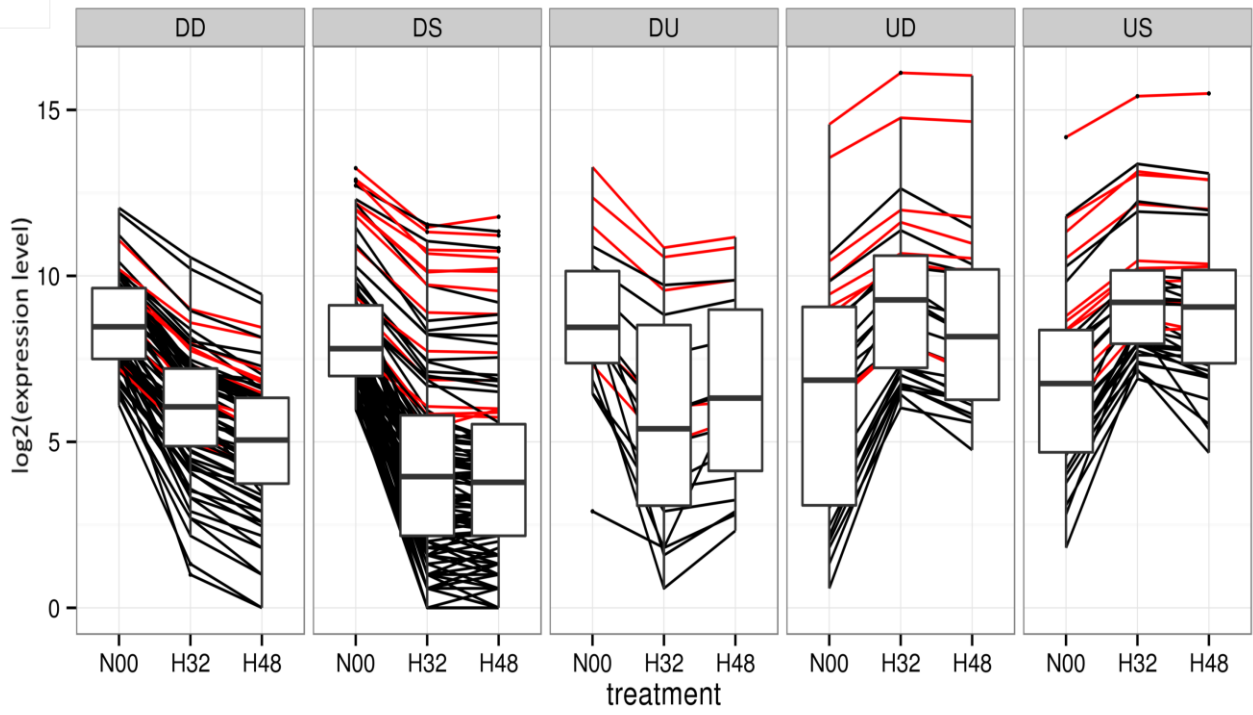
595



596

597 Figure 5

598



599

600 Figure 6

Tool	Format	DataType	Fastq QC	Nt freq	Adp trim	Size class	Annotation	MA/scatter	Norm	DE	Ref
DESeq	R library	RNA-seq	N	N	N	N	N	Y	DESeq	Y	Love et al 2014
edgeR	R library	RNA-seq	N	N	N	N	N	Y	TMM	Y	Zhou et al 2014
baySeq	R library	RNA-seq	N	N	N	N	N	N	Quantile	Y	Hardcastle et al 2010
RSEQtools	Software	mRNA-seq	N	N	N	N	Y	N	RPKM	Y	Habegger et al 2011
DARIO	web	ncRNA-seq	N	N	N	Y	Y	N	-	N	Fasold et al 2011
Cyber-T	Web	RNA-seq	N	N	N	N	N	N	Logarithmic,VSN	Y	Kayala et al 2012
ncPRO-seq	Software	sRNA-seq	Y	Y	N	Y	Y	N	-	N	Chen et al 2012
Shortran	Software	sRNA-seq	N	N	Y	N	N	N	Total count	Y	Gupta et al 2012
RobiNA	Software	RNA-seq	Y	Y	Y	N	N	N	RPKM	DeSeq/edgeR	Lohse et al 2012
omiRas	Web	miRNA-seq	Y	N	Y	N	N	N	DESeq	DeSeq	Muller et al 2013
Kraken	Software	RNA-seq	Y	Y	Y	N	N	N	-	N	David et al 2013
TCC	R library	RNA-seq	N	N	N	N	N	N	DEGES/TbT	Multiple	Sun et al 2013
sRNAtoolbox	Web	sRNA-seq	N	N	Y	N	N	N	edgeR, NOIseq	edgeR, NOIseq	Rueda et al 2015
UEA sRNA Workbench	Software	sRNA-seq	Y	Y	Y	Y	Y	Y	RPM, quantile, subsampling, DESeq, TMM	Y	Stocks et al 2012

601

602 Main Table 1