1   **Side-by-side analysis of alternative approaches on multi-level RNA-seq data**

2

3   IRINA MOHORIANU[1,2],

4

5   *¹School of Biological Sciences, University of East Anglia, Norwich Research Park, Norwich, NR4 7TJ,*

6   *United Kingdom.*

7   *²School of Computing Sciences, University of East Anglia, Norwich Research Park, Norwich, NR4 7TJ,*

8   *United Kingdom.*

9   **#Corresponding author:** i.mohorianu@gmail.com

10

11   *Running title*: Alternative approaches for multi-level RNA-seq data

12

13

14   **ABSTRACT**

15   **Background**:

16   RNA sequencing (RNA-seq) is widely used for RNA quantification across environmental,

17   biological and medical sciences; it enables the description of genome-wide patterns of

18   expression and the deduction of regulatory interactions and networks. The aim of

19   computational analyses is to achieve an accurate output, i.e. rigorous quantification of

20   genes/transcripts to allow a reliable prediction of differential expression (DE), despite the

21   variable levels of noise and biases present in sequencing data. The evaluation of sequencing

22   quality and normalization are essential components of this process.

23   **Results**:

24   We investigate the discriminative power of existing approaches for the quality checking of

25   mRNA-seq data and also propose additional, quantitative, quality checks. To accommodate

26   the analysis of a nested, multi-level design using data on *D. melanogaster*, we incorporated

27   the sample layout into the analysis. We describe a "subsampling without replacement"-based

28   normalization and identification of DE that accounts for the experimental design i.e. the

29   hierarchy and amplitude of effect sizes within samples. We also evaluate the differential

30   expression call in comparison to existing approaches**.** To assess the broader applicability of

31   these methods, we applied this series of steps to a published set of *H. sapiens* mRNA-seq

32   samples.

33   **Conclusions**:

34   The dataset-tailored methods improved sample comparability and delivered a robust

35   prediction of subtle gene expression changes. Overall, the proposed approach offers the

36   potential to improve key steps in the analysis of RNA-seq data by incorporating the structure

37   and characteristics of biological experiments into the data analysis.

38

39   **Keywords: RNA-seq; quality check; normalization; subsampling normalization;**

40   **identification of differential expression; hierarchical differential expression.**

41 **Background**

42 RNA sequencing (RNA-seq) has revolutionized the field of trascriptomics (Wang et al. 2009;

43 Ozsolak et al. 2011), giving powerful insight into the identity and abundance of RNAs in cells,

44 tissues and whole organisms (Jia et al. 2017). In contrast to the fixed, predefined set of probes

45 used for microarray experiments, RNA-seq generates a diverse set of reads and facilitates

46 analyses of expression level variation for known and unknown RNA transcripts and variants.

47 It also offers the possibility to study additional facets of the transcriptome (Conesa et al. 2016a),

48 such as the (re)annotations of the reference genomes (Torres-Oliva et al. 2016), the

49 identification of alternative splicing events (Trapnell et al. 2012) and variation in abundance

50 across transcripts (Dillies et al. 2013b; Patro et al. 2017).

51 Several bioinformatics methods have been used to analyse the rapidly rising number

52 of RNA-seq datasets (reviewed in (Dillies et al. 2013b; Conesa et al. 2016a; Evans et al.

53 2017)). However, to accommodate the use of RNA-seq in complex experimental designs,

54 there is scope for further developments in: (i) the robust detection of subtle signatures of gene

55 expression (the concordance of which is often very low between different bioinformatics

56 methods, (Rapaport et al. 2013; Roca et al. 2017), (ii) the incorporation of hierarchical

57 experimental designs (Love et al. 2014b; Robinson et al. 2015; Schurch et al. 2016), e.g. from

58 evolutionary experiments (Fang et al. 2011), and (iii) minimising the effect of normalisation on

59 the pattern of DE, especially when differences are subtle (Dillies et al. 2013b). We discuss

60 these themes in the major steps of RNA-seq analysis, below.

61 **Quality Checks** (**QC)**

62 A key step in the analysis of RNA-seq data consists of sample quality checks and the

63 identification, characterization and potential exclusion of sample outliers, e.g. those samples

64 that are compromised due to technical issues (Consortium 2014). Existing tools, such as

65 FastQC (Andrews 2010), SeqMonk (Andrews 2010) or TagCleaner (Schmieder et al. 2010)

66 focus on evaluating the sequencing and the per-base quality. Additional QC may include an

67 analysis of the per-base nucleotide composition and an evaluation of the overall GC content

68   (Risso et al. 2011; DeLuca et al. 2012; Wang et al. 2012). QC procedures focused on the

69   sequencing bias comprise the characterization of k-mer distributions (Hansen et al. 2010) as

70   well as other adapter ligation effects, such as the secondary RNA structure of the insert with

71   the attached adapters (Jayaprakash et al. 2011b; Jackson et al. 2014).

72   Quantitative analysis of sequencing output currently considers measures such as yield,

73   coverage, 3'/5' bias, number of detectable transcripts, strand specificity and read distribution

74   across the genome (Consortium 2014). Additional steps, at the transcript level, include the

75   classification of reads into annotation classes, which can highlight the presence of potential

76   contaminant ncRNAs such as tRNAs and rRNAs (DeLuca et al. 2012; Wang et al. 2012). Such

77   analyses can reveal over-represented classes of sequences, which could be removed in order

78   not to distort the subsequent normalization and lead to changes in the ranking of abundances.

79   In recognition of the multiple sources of variation present in a biological experiment, a

80   QC criterion is based on the Pearson Correlation Coefficient (PCC) between the gene

81   expressions in biological replicates for transcripts that are detected in both samples (McIntyre

82   et al. 2011; Gierlinski et al. 2015). Values of $r^2$ = 0.92-0.98 are generally accepted. If the PCC

83   falls below 0.9 the suggestion is to identify and potentially exclude the problematic samples

84   (Conesa et al. 2016b). However, this criterion may lack discriminatory power; due to the high

85   number of data points (expressed genes, e.g. vector containing >15K genes for *D.*

86   *melanogaster*) correlations will often be very high between *all* samples, regardless of their

87   quality.

88   Additional steps for the quantitative evaluation of samples are possible, but are, as yet,

89   under-utilised. These techniques include analyses of per-sample or per-gene complexities

90   defined as the ratio of non–redundant (NR, unique) to redundant (R, total) reads (Mohorianu

91   et al. 2011a), and similarity comparisons (Jaccard 1901), (Beckers et al. 2017).

92   **Normalization**

93   The next key stage in the analysis of RNA-seq data is the normalization of gene expression

94   levels (Mortazavi et al. 2008b; Conesa et al. 2016a; Roca et al. 2017). Initial reports describing

95   RNA-seq suggested that no normalization method was required (Wang et al. 2009). However,

96      subsequent studies correctly highlighted normalization as a critical step (Bullard et al. 2010;

97      Dillies et al. 2013a). Normalization is designed to transform the distributions of abundances

98      for each sample, without distortion, into distributions that can be compared. A good

99      normalization increases the chances of an accurate call of DE. It accounts for differences in

100     sequencing depths and in biases arising from the library preparation or its sequencing (Li et

101     al. 2014; Li et al. 2015; Lin et al. 2016).

102             However, despite extensive attention from the community (Aanes et al. 2014), there is

103     as yet no clear consensus on whether there is any single optimal normalization method (Dillies

104     et al. 2013b; Roca et al. 2017). Nor is there any general appreciation of the potential

105     magnitude of the consequences of ineffectively normalizing data; their extent depends on the

106     amplitude and distribution of DE, with small gene expression differences being more sensitive

107     than larger ones to the method of normalization. Therefore, particularly for analyses of subtle

108     gene expression differences, it can be important to assess how well the data are normalized

109     by a specific method (Wagner et al. 2012). Such tests are not a routine part of bioinformatics

110     analyses (Evans et al. 2017).

**Identification of Differential Expression**

112     The goal of transcriptomics analyses is the accurate and unbiased identification of expressed

113     genes and genes showing DE between treatments. The majority of existing methods exhibit

114     a good level of overlap in terms of highly differentially expressed genes (Soneson et al. 2013;

115     Roca et al. 2017) . However, their agreement is far less when DE is subtle. Comparative

116     analyses of existing normalization procedures on real and simulated data sets show that only

117     ~50% of significantly differentially expressed genes are identified by all methods (Rapaport et

118     al. 2013; Lin et al. 2016) .

119             ANOVA-based methods are a powerful and extensively applied approach for the

120     analysis of microarray data (Cui et al. 2003). However, such methods are based on *a priori*,

121     to some extent arbitrary, significance thresholds and the type of experiment can greatly

122     influence the expected number of genes showing DE. For example, if the DE frequency

123     distribution is narrow, stringent p-value thresholds can indicate as statistically significant

124    genes that show very small fold change differences. Therefore, the set of DE genes identified

125    by a fixed p-value threshold may not necessarily reflect biologically important facets of the

126    data. Such differences are unlikely to be validated by low throughput methods (Evans et al.

127    2017).

128          Newer methods such as DESeq2 (Love et al. 2014a) and edgeR (Robinson et al. 2010)

129    are based on the negative binomial distribution model for expression levels, using the variance

130    and mean linked by local regression to detect DE genes (DESeq2) and empirical Bayes

131    methods for moderating the degree of over dispersion across transcripts (edgeR). However,

132    existing methods do not easily accommodate inherent, hierarchical experimental design with

133    variable amplitude of DE between the hierarchy levels.

134          To fully encompass this type of experiment (exemplified using the *D. melanogaster*

135    dataset) we incorporated the structure and magnitude of gene expression differences into the

136    analysis, prior to the DE call. We accounted for the specific type of variation in expression for

137    behavioural experiments and structured the analysis framework for the mRNA-seq data (Fig.

138    S1); from QC (including existing and new approaches) to a subsampling without replacement-

139    based normalization and finally a hierarchical approach for the identification of transcripts

140    showing DE. The proposed analysis also features the use of an adjustable, empirically-

141    determined offset to filter out low abundance genes and a DE call using maximal confidence

142    intervals.

143          The adapted methods also performed well in direct comparisons with existing

144    approaches for the analysis of a publicly available *H. sapiens* mRNA-seq dataset. Overall, the

145    methods may offer advantages in the analysis of complex, challenging datasets and are

146    complementary to existing approaches.

147

## 148    RESULTS and DISCUSSION

149    To develop and test the adapted methods we used a *D. melanogaster* dataset in which we

150    tested for subtle effects on gene expression in males exposed to mating rivals (Mohorianu et

151    al. 2017). We compared the output of our pipeline with that of existing methods on the same

152　input data. The description of the steps of this analysis, in line with the approaches described

153　in (Conesa et al. 2016a), are presented in Fig S1.  We also assessed the applicability of the

154　adapted normalization on a publicly available *H. sapiens* mRNA-seq dataset.

155　**Quality Checking**

156　**Stage 1: QC of sequencing quality**

157　The *D. melanogaster* data comprised of 3 replicate mRNA-seq samples of 2 rival treatments

158　(rivals versus no rivals), 2 body parts (Head+Thorax (HT) and Abdomen (A)) and 3 time

159　exposure treatments (2h, 26h or 50h). The first stage of the QC on these data (Fig. S1)

160　focussed on existing approaches (i) the analysis of FastQ quality scores (Andrews 2010), (ii)

161　sequencing depth, (iii) nucleotide (nt) composition / GC content (DeLuca et al. 2012; Wang et

162　al. 2012), (iv) strand bias and (v) proportions of genome and annotation classes - matching

163　reads e.g. mRNAs, t/rRNAs, miRNAs, UTRs, introns, intergenic regions (Conesa et al. 2016a).

164　The FastQ QC indicated good quality reads for all 50nt, though we observed high variability

165　in sequencing depths. Variation in nucleotide content was observed across the first 12nt

166　(Jayaprakash et al. 2011a), but was consistent after that with nucleotide composition of the *D.*

167　*melanogaster* transcriptome. Strand bias was comparable across samples and the proportion

168　of genome-mapping reads was high. Based on these stage 1 quality checks, all samples were

169　retained for further analyses and entered stage 2 QC. The detailed results, supporting the

170　conclusion that the samples were consistent based on these criteria, are presented in Tables

171　S1 and S2 and in (Mohorianu et al. 2017).

172　**Stage 2: Quantitative QC of replicate and sample comparability**

173　*Jaccard/intersection analysis*

174　We computed the Jaccard similarity (Mohorianu et al. 2011b; Beckers et al. 2017) at the gene

175　level, to compare the similarity in expression of the top 1000 most abundant genes present in

176　each sample (Table S3). This measure evaluated what proportion of the most abundantly

177　expressed genes in one sample are also the most abundant in the next sample, and so on.

178　Being calculated on the top most abundant genes it is not biased by different numbers of

179　expressed genes in each sample. In addition, by selecting ~5% of the most abundant genes,

180   this measure is not biased by noise-derived variability in expression. Samples drawn from the

181   same body parts shared > 90% similarity, and between body parts (HT versus A) the similarity

182   dropped to ~50%. Similarity between the experimental (± rivals) treatments was sometimes

183   higher than between replicates, which highlighted the need for alternative approaches for

184   normalizing the gene expression levels.

185

186   ***Complexity analysis***

187   A comparison of complexity (calculated as the ratio of non-redundant (NR), unique, reads to

188   redundant (R), all, reads, on all mRNA matching reads) is an informative measure of the

189   number of unique reads present in each sample, average abundance of reads and

190   subsequently on the average coverage. It is also useful for identifying sample/replicate outliers.

191   A complexity of ~0 would indicate a sample in which all reads were the same and 1 a sample

192   where every single read was different. Sample complexity is influenced by sequencing depth.

193   Samples with high sequencing depth have a lower overall complexity and vice versa. However,

194   samples with comparable sequencing depths, but very different number of unique reads, may

195   highlight incomparable replicates/ samples. To understand how this may influence the

196   accuracy of DE, we calculated the variation in complexity, between replicates, at gene level

197   (Fig. 1). This revealed sizeable differences in complexity for most genes (Fig. 1(A)), which

198   correlate well with the presence of highly variable number of spurious incident reads,

199   especially for low abundance genes; this conclusion was also confirmed by the point-to-point

200   correlation, described in the next section. To reduce this technical variation and normalize the

201   data we tested the effect of subsampling with- (Fig. 1(B)) and without- (Fig. 1(C)) replacement

202   for each set of reads to a fixed total (see below). We observed a reduction in complexity as a

203   result of the subsampling both with and without replacement and an increased similarity

204   between replicates/ samples. However, the subsampling with replacement artificially indicated

205   that the third replicate (R3) is acceptable (Fig. 1(B)), while the subsampling without

206   replacement (Fig 1(C)) maintains the conclusion that the third replicate is an outlier.

207

208    ***Correlation analyses***

209    Correlations were first calculated between gene expression vectors in each sample to assess

210    their comparability (using Pearson (PCC), Spearman (SCC) and Kendall (KCC) correlation

211    coefficients), similarly as in (Gierlinski et al. 2015). Using the raw expression levels,

212    correlations were computed between each sample and every other. Correlations between HT

213    and A samples (in the range of 0.75-0.8) were lower than correlations between same body

214    part samples. This was expected on the basis of HT- and A-specific genes whose expression

215    is restricted to each body part. When only A or HT samples were considered, all correlation

216    coefficients were above 0.95 (Fig. S2). These results showed a high correlation between

217    samples and no sample would be excluded as clear outlier, even though, based on other

218    quantitative QC measures, it would be prudent to do so (see below). Hence, standard

219    correlation metrics may not be sufficiently sensitive to evaluate sample quality.

220

221    ***Point-to-Point Correlation analyses***

222    To gain additional insight we introduced the 'point-to-point' Pearson correlation coefficient

223    (p2pPCC) for each gene, which is a standard Pearson Correlation computed across the whole

224    length of a gene on its expression profile. Although we expect to see alternative splicing events

225    between the different replicates and samples, the number of genes affected by this is expected

226    to be small. Next, we evaluated the distribution of PCC for pairwise sample or replicate

227    comparisons against the corresponding gene abundances. This analysis revealed a higher

228    variability at low in comparison to high abundances and a knock-on effect on the DE among

229    replicates, i.e. more DE at low in comparison to high abundances (Fig. S3).

230    The p2pPCC was also used to determine a variability threshold (noise threshold, denoted as

231    "offset" in the DE call section). This value was calculated, for the whole dataset, once the

232    abundances of the sequencing reads were normalized. Briefly, low abundance genes have a

233    limited number of incident fragments that align to random locations on the gene. As the gene

234    abundance increases, the alignment pattern of the incident fragments starts to resemble the

235    gene model leading to an increase in the p2pPCC. We determined the offset as the gene

236 abundance for which the median of the p2pPCC distribution was > 0.7 between replicates. A

237 similar approach, based on the entropy of strand bias, was implemented for sRNA sequencing

238 in (Beckers et al. 2017).

239 Overall the qualitative QC metrics were informative but lacked discriminatory power.

240 The quantitative QC metrics focused on the comparability at gene level by analyses of

241 complexity and similarity and represent a valuable addition to overall QC.

242 **Subsampling based Normalizations**

243 To attenuate the effect of variation in sequencing depth between replicates and samples

244 described above, we implemented a subsampling (without replacement) normalization on read

245 expression levels adapted from (Li et al. 2013b), enriched with several additional checks on

246 the consistency of the subsample and the similarity to the original sample. The subsampling

247 approaches already in use either focus on the subsampling with replacement option applied

248 on the aligned reads (Gierlinski et al. 2015) or on gene abundances (Robinson et al. 2014),

249 or the subsampling without replacement applied on all reads, without additional per-gene

250 consistency checks (Stupnikov et al. 2015).

251 *Incremental subsampling (without replacement) and bootstrapping-based sample*

252 *checking*

253 First we tested the homogeneity of each sample, as indicated in (Stupnikov et al. 2015). To

254 evaluate the existence of high abundance reads which, due to their higher probability of being

255 selected, could distort the normalized distributions, we conducted a subsampling exercise,

256 from 95% down to 45% (in steps of 5%) of the original redundant reads. Part of the novelty of

257 our approach is to assess the consistency of the sample by checking if the proportion of

258 redundant genome matching reads was affected by the subsampling (Table S4 and Fig S4).

259 Even when the data were subsampled to 45% of the original sequencing depth, the proportion

260 of redundant genome matching reads did not change. However, the complexity of the sample

261 increased and became comparable to other samples with similar number of reads.

262

263         The next novel element was to evaluate the extent to which the data could be

264      subsampled without affecting its structure we calculated the point-to-point PCC on expression

265      levels of the original versus subsampled data from 95% to 40% of the original R set (Fig. S4).

266      This showed that the correlations of abundantly expressed genes remained high over all

267      subsamples, but that the correlation of low abundance transcripts decreased as the proportion

268      of data subsampled dropped (note though that the variability between the original versus

269      subsamples was lower than the variability between the biological replicates). We concluded

270      that the subsampling was effective as it maintained high p2pPCC and strong concordance

271      between the expression levels of the raw versus normalized data (Fig. S4).

272         The number of genes 'lost' due to the exclusion of some low abundance reads was

273      typically <2%. Once we had determined that all samples passed the consistency check, we

274      subsampled every sample to a fixed total of 50M reads and used to check whether

275      subsamples were representative of the original data bootstrapping (Supplemental Material

276      Methods 1). Following this step, one subsample was selected at random for each sample and

277      used in the downstream analysis. This subsampling was efficient at correcting wide variation

278      in read number, complexity differences and minimising the impact of normalization on the

279      original data structure (Fig. 1).

280         The analysis of the distributions of complexity differences between replicates, coupled

281      with the Jaccard similarity analyses applied on the normalized data, was used to identify outlier

282      replicates, which were excluded from subsequent analyses. We classify as an outlier, samples

283      for which the between-replicate similarity was higher than between-sample similarity, as

284      shown by the Jaccard, complexity and p2pPCC analyses. In the 02RH example presented in

285      Fig 1 we showed that the subsampling without replacement correctly highlighted replicate 3

286      as an outlier, while the subsampling with replacement did not. Following this post-

287      normalization QC, we retained two biological replicates for each treatment for the *D.*

288      *melanogaster* data. In general, we advocate the use of as many biological replicates as

289      possible. However, as we show below, the analysis of subtle gene expression even with a

290      limited number of replicates is possible and can be validated.

291    **Subsampling with replacement vs subsampling without replacement**

292    A subsampling with replacement normalisation has previously been proposed (Li, Tibshirani

293    2013b; Robinson, Storey 2014; Gierlinski et al. 2015); a version of the subsampling without

294    replacement was proposed in (Stupnikov et al. 2015). The clear differences, and advantages,

295    of a subsampling based method i.e. the lack of consistency of a scaling factor for all

296    abundances, and thus the advantage over scaling methods, are discussed in (Li, Tibshirani

297    2013b; Stupnikov et al. 2015).

298    The main difference between the two approaches (with- vs without- replacement) is that for

299    the former, the selection probabilities for the reads remains unchanged during the process,

300    whereas for the latter the probabilities of selection are not constant, facilitating the selection

301    of both high and low abundance reads. The abundance range for the later selection is wider

302    and simulates the selection which takes place during the sequencing. The advantage of the

303    former approach is that a higher number of reads can be achieved for a given sample (Li,

304    Tibshirani 2013b) e.g. if a sample with 20M reads needs to be compared with a sample with

305    40M reads, then both samples could be subsampled, with replacement, at 30M reads. The

306    down-side of the approach is that depending on the proportion of selected reads, over-

307    amplification of high abundance reads and exclusion of low abundance reads can occur. This

308    has a knock-on effect in amplifying the expression of abundant fragments and, as a

309    consequence, reducing the expression of genes with low abundance which have fewer, low

310    abundance incident fragments. In addition, the omission of low abundance reads may change

311    the expression profile across transcripts (Fig 2(B)). The disadvantage of the subsampling

312    without replacement is that it has an upper bound, calculated as the minimum sequencing

313    depth between the samples of an experiment. Its advantage is that it renders the samples

314    comparable (as described in the previous section) and allows the identification of potential

315    outliers. In Fig 2(A), we compared the distribution of abundances for the two approaches (with

316    versus without), on the same samples, to the same total. In Fig 2(B), we show the presence

317    plots of a gene with 0.4 difference in complexity between the two approaches. The change in

318    profile is visible for the third replicate, changes in the profile itself are highlighted with arrows.

319

### Calling of DE using a hierarchical approach

320

321 To incorporate the experimental design into the DE call (Anders et al. 2010) we used a simple

322 hierarchical approach for the prediction of DE transcripts (Supplemental Methods 2). The

323 order of levels in the hierarchy was determined based on the amplitude of DE for each factor

324 in the experiment (Fig. S5). For the *D. melanogaster* dataset the highest level in the hierarchy

325 (i.e. that showed most DE) was body part (HT vs A), the second was ± rivals treatment. The

326 distribution of DE between treatments and between replicates overlapped (Figure S5(A,B)),

327 which indicated that the treatment DE was likely to be subtle.

328 We observed direct evidence of the biasing effect of low abundance FC (Fig. 3(A)).

329 For example, using standard FC numerous low abundance DE genes in the HT were in fact a

330 signature from the A body part (e.g. sperm and semen genes are specific to the A body part,

331 but detected as differentially expressed in the HT at low abundance; Fig. 3(A,B)). The RNA-

332 seq technique is highly sensitive and detects these transcripts due to leak through or

333 movement of mRNAs. Unless a correction is applied, the list of DE genes is likely to contain

334 numerous spurious and low-abundance entries. A practical solution was to use offset fold

335 change (OFC), with the offset determined empirically from the data, as described in an earlier

336 section, in preference to FC (Fig. 3A versus B) and to apply the hierarchical DE (Fig. 3C;

337 Supplemental Methods 2). A comparison of the MA plots for all genes versus the A- and the

338 HT-specific genes showed the effect of the hierarchical differential expression (Fig. 3C).

339

### Case study - comparisons with existing approaches

340

341 We next evaluated the output gained from the analysis of the *D. melanogaster* data using our

342 bioinformatics framework, with that obtained from the analysis of the same, original input data

343 (consisting of all 3 replicates for each condition) using DESeq2 (Love et al. 2014a) and edgeR

344 (Robinson, Oshlack 2010; Zhou et al. 2014). Although the DE was conducted on the HT

345 samples, both the HT and AB samples were given as input the edgeR and DEseq2 analyses.

346

347    ***Effect of the normalization***

348    An *a priori* (and necessary, but not sufficient) condition for reliable DE call is good

349    comparability between the distributions of expression levels. To assess whether the

350    normalized expression levels became more similar after the normalization (in particular,

351    whether the replicates became comparable), we compared the distributions of expression

352    levels (log$_2$ scale) for the raw data versus RPM, quantile, subsampling with and without

353    replacement, DEseq2 and edgeR normalizations (Fig. 4). The boxplot of the raw abundances

354    (Fig. 4A) illustrates the variation among the replicates and samples and clearly indicated that

355    normalization was required. The RPM normalization (Fig. 4B) rendered the A and HT

356    distributions comparable to some extent. However, it was difficult to separate the A- or HT-

357    specific genes. In addition, variability between samples was still present (especially for the HT

358    samples). The quantile normalization (Fig. 4C) rendered the distributions comparable, as did

359    the subsampling with and without replacement (Fig. 4D1 and Fig 4D2). DEseq2 performed

360    well (Fig. 4E) – although residual differences in the distributions of the A vs HT samples

361    remained. EdgeR (Fig. 4F) did not effectively equalize the distributions of abundances. We

362    concluded that the subsampling, quantile and DESeq2 normalizations (Fig. 4C,D1, D2,E) were

363    most effective at producing comparable distributions of normalized expression levels for this

364    dataset.

365

366    ***Differences in the DE call between methods***

367    To evaluate the effect of the normalization and hierarchical DE call we compared analyses of

368    the *D. melanogaster* 2h HT and A samples ±rivals (using as input all three biological replicates

369    of the original data) with the output of DEseq2 and edgeR (Fig. 5).

370         The subsampling without replacement normalization and hierarchical DE call (Fig. 5A)

371    showed a relatively low number of up/down-regulated genes with relevant biological functions.

372    The equivalent analysis for DEseq2 (Fig. 5*B*) called many more genes as DE that fell in the

373    region of +/- 0.5 log$_2$ FC (i.e. below the validatable threshold (Morey et al. 2006) (Fig. S5). The

374    analysis using edgeR (Fig. 5C) showed a high frequency of low abundance DE and of leaky

375    genes, which is likely to either represent noise or biological signal of an insufficient magnitude

376    to be captured effectively in the low throughput validation (Fig. S5). The degree of overlap

377    between the three methods (Fig. 5D) revealed a small number of core genes present in the

378    intersection. edgeR and DEseq2 called many more genes as DE and the number of genes

379    uniquely identified by edgeR and DEseq2 was also larger than the number identified in

380    common between the two. These results show that the pipeline chosen will have a strong

381    effect on the biological interpretation of the DE analysis.

382        For the ±rivals comparison for the HT samples, out of the 575 genes that were specific

383    to edgeR, 14 were HT genes (all with max abundance > 50) and 561 were A genes (327 with

384    max abundance > 50 and 234 < 50; Fig. S5). Out of the 578 genes specific to DESeq2, 101

385    were HT genes (100 with abundance > 50) and 477 were A genes (271 with max abundance >

386    50 and 206 < 50; Fig. S5). The predominance of A genes in the DE call supported the use of

387    the hierarchical DE approach. The presence of low abundance genes supported the use of an

388    offset for the calculation of DE.

389        Of some concern was that for genes with a reasonable abundance (> 50) the

390    expression intervals for the ± rivals differences called by DEseq2 and edgeR were

391    close/overlapping, making independent validation using low throughput methods challenging.

392    These results showed evidence of a high number of DE genes called by both or either of

393    DEseq2 and edgeR that would be difficult to validate independently.

394        Reasons for the differences in DE call between methods are likely to result from a

395    propagation of factors throughout the whole analysis. These potentially include low

396    comparability of normalized distributions of expression or the imperfect assessment of

397    hierarchical distribution of DE levels in the experiment. In the calculation of DE by DEseq2,

398    replicate-to-replicate variability is averaged and DE over and above this variation is calculated.

399    The accuracy of such averaging relies on replicates having a low coefficient of variation (CV

400    = standard deviation/mean). However, in the *D. melanogaster* data, in many cases the CV

401    was over 0.25 and in some was > 0.5 (across all abundances, Fig. S6). In the example shown

402    (Figs. S6A-D), there was clearly higher dispersion (CV) at low transcript abundance and

403    consistently high CV across higher abundances. This variation was corrected in our analyses

404    by the subsampling normalization and the use of an OFC (Fig. S6 E,F,G,H, in which dispersion

405    showed minimal variation across transcript abundance and the CV was consistently low

406    (generally < 0.1)). The newer version of DESeq, DESeq2, (Love et al. 2014a) notes the effect

407    of high replicate variation and proposes a shrinkage estimation for dispersions based on

408    empirical Bayes and FC, to improve stability and interpretability of estimates. Our results

409    showed that these changes helped to minimize, but did not fully solve, the overall issue of high

410    variability in transcript abundance.

411

412    ***Comparison of low throughput validated genes with DESeq2 and edgeR outputs***

413    We next investigated whether the set of DE genes identified using the hierarchical approach

414    from our *D. melanogaster* dataset, and validated by qRT-PCR (Mohorianu et al. 2017), were

415    present in the output of DEseq2 and edgeR. Reassuringly, based on DESeq2 our qRT-PCR

416    reference genes were not called DE. Two other genes of interest from the A samples that

417    were validated as DE, had a p-value < 0.05 by DESeq2 (although adjusted p-value > 0.05).

418    One gene of interest validated from the HT had a p-value < 0.05 (but again not according to

419    the adjusted p-value). For DESeq2 the $\log_2$(FC) values were small (0.15 and 0.16,

420    respectively). Hence these genes were not likely to have been selected for further

421    investigation. Based on the edgeR output, our reference genes were also determined as not

422    significantly DE. For the validated A genes of interest (GOI), only one was called marginally

423    DE by edgeR (p = 0.08, FBgn0259998, but with small $\log_2$(FC) = 0.24). For the HT, one GOI,

424    with $\log_2$(FC) = 0.58, was called as significantly DE (the same gene as identified by DESeq2).

425    Another GOI, FBgn0044812, was identified with $\log_2$(FC) = 0.82 yet with a p-value of 0.49

426    from edgeR and therefore would not have been selected (Table S1).

427         Comparing the validated gene set with the output of edgeR and DESeq2, we conclude

428    that some GOIs failed to be identified and therefore the corresponding biological functions

429    (immunity, odorant perception) might have been overlooked.

430    **Which normalization to choose?**

431 An RNA-seq sample is a snapshot of RNA fragments present at a given time, randomly

432 selected according to the RNA abundances, to fill the sequencing space. Due to the stochastic

433 nature of the sequencing process, even technical replicates, at different sequencing depths,

434 do not exhibit a constant scaling factor for all abundances. Also, RNA-seq outputs have

435 varying fits to standard distributions, making it difficult to define "the best" choice. Although the

436 subsampling without replacement normalization was efficient in minimizing the effects of the

437 variable sequencing depth, while preserving a high similarity with the original samples (Fig. 3),

438 we suggest that it is advisable to test different normalizations on mRNA-seq and choose the

439 most appropriate method for the given dataset, on a case-by-case basis (Beckers et al. 2017)

440

**Analysis of human mRNA-seq datasets using subsampling (without replacement)**

**normalization**

443 RNA-seq is expected to have good external validity and produce comparable results when the

444 same RNA is used, across different laboratories. However, a recent study of mRNA-seq

445 conducted on the same human samples (expression level variation in lymphblastoid cell lines)

446 involved the use of the same samples sequenced in two different locations: Yale versus

447 Argonne (Pickrell et al. 2010). Some variation between the results from the different

448 laboratories was observed (Zhou et al. 2014). The authors analysed these data further to

449 explore whether edgeR could reduce the variability between replicates. We tested whether

450 our subsampling normalization could further reduce such variation. To do this, we randomly

451 selected 5 sets of samples (144, 153, 201, 209 and 210) with two replicates each, one from

452 the Yale laboratory source and one from the Argonne source. For these runs, the length of

453 the reads was 36nt for Yale and 46nt for the Argonne-derived data. Since the length of the

454 sequencing read influences the number of unique fragments and the mapping to the reference

455 transcriptome (and, as a result, the gene expression) we trimmed all reads to comparable

456 lengths (35nt) and mapped the reads to the reference human genome using full length, no

457 mis-match or gap criteria and using PatMaN (Prufer et al. 2008). The subsampling was

458    conducted on 7M reads (the number of reads for the smallest sample was 7.1M, and for the

459    largest sample was 8.7M)

460        We created comparable plots to (Zhou et al. 2014) for the data subjected to

461    subsampling normalization. MA plots for the two replicates of each sample showed high

462    reproducibility between runs. The distribution of the coefficient of variation (CV) versus the

463    abundance for the 5 selected sets of samples with two replicates each (one for Yale and one

464    for Argonne) (Fig. S7) showed that the CV for all 5 pairs of samples was consistently (< 0.1)

465    lower than for the analysis of (Zhou et al. 2014)) indicating a very high similarity between the

466    runs. The MA plots on the same sets of two samples showed a high reproducibility between

467    replicates (no genes showing $|log_2(OFC)| > 1$). The genes showing DE were mainly localized

468    in the $2^4$ (16) – $2^6$ (64) range, which is borderline for validation/noise. Together, these analyses

469    showed that: (i) the CV obtained when the subsampling (without replacement) normalization

470    was employed was lower than the CV reported in (Zhou et al. 2014), suggesting that the

471    normalization was tighter, (ii) there was very little DE between replicates, indicating good

472    reproducibility between the sequencing runs.

473        Overall, we conclude that the subsampling, without replacement approach cleared the

474    technical differences between the two runs in the different laboratories and this approach

475    rendered the samples comparable, potentially improving the biological inference.

476

## Conclusion

478    The main conclusion from this study was to emphasise the need to check multiple approaches

479    for the analysis of a dataset and to show that both qualitative and quantitative QC are

480    informative, and the applicability of subsampling (without replacement) -based normalization

481    and hierarchical structuring of the DE call, is efficient in managing variation in read number

482    and differences in sample complexities. In comparison to existing methods, the adapted

483    methods performed well and identified valid candidates that were confirmed using low

484    throughput approaches (Mohorianu et al. 2017). We also successfully applied the subsampling

485    (without replacement) normalization to existing mRNA-seq datasets, used to analyse inter-

486    laboratory variation (Pickrell et al. 2010); the adapted approach proved to be efficient in

487    comparison with existing methods at minimizing potentially confounding sources of variation.

488    Determination of accurate gene expression levels is essential for all mRNA profiles but is also

489    key to successful correlation analysis between mRNAs and sRNAs (Mohorianu et al. 2012;

490    Mohorianu et al. 2013).

491

## 492    **METHODS**

493    **Quality check (QC)**

494    For the mRNA-seq samples, the QC consisted of two stages. Stage 1 comprised of previously

495    described methods (Conesa et al. 2016a) including: (i) the analysis of FastQ quality scores

496    (Andrews 2010), (ii) the total number of reads (sequencing depth) and the read duplication

497    rate, defined as complexity (Mohorianu et al. 2011a), (iii) nucleotide composition relative to

498    the genome and transcriptome of *D. melanogaster*, used to highlight biases such as PCR and

499    ligation bias (Sorefan et al. 2012), (iv) strand bias quantified on CDS incident reads as

500    $|P - 0.5| + |N - 0.5|$, where P and N were the proportion of positive and negative strand read

501    matches, respectively (Mohorianu et al. 2011a) and (v) proportions of reads matching the

502    different genome annotation classes (e.g. mRNAs, t/rRNAs, miRNAs, UTRs, introns,

503    intergenic regions (Conesa et al. 2016a); matching was done on full length reads with no mis-

504    matches or gaps allowed, using PatMaN, (Prufer et al. 2008)). Stage 2 comprised of

505    quantitative approaches, some applied/designed on mRNA-seq data for the first time, which

506    provided an increased insight into sample comparability and enabled us to evaluate the

507    effectiveness of the normalization. The expression level of a gene/ transcript was calculated

508    as the algebraic sum of the raw/normalized abundances of the incident reads (Mortazavi et al.

509    2008a). We examined (i) sample similarity calculated using the Jaccard similarity index

510    (Jaccard 1901) on the top 1000 most abundant genes, and intersection analyses); theseq

511 measures were calculated as the ratio between number of genes found in common to the

512 number of unique genes present in in either samples, (ii) complexities (calculated at gene level

513 and presented as Bland-Altman plots) and (iii) point-to-point PCC between gene expression

514 profiles in different replicates/ samples. The latter were computed on the vector of expression

515 defined for each gene. For all positions *i* on a gene we computed y[i] which is the sum of

516 abundances of fragments incident with position *i*. The point-to-point PCC was computed as

517 the standard PCC on the corresponding vectors from the two samples which were compared.

518 **Normalization**

519 We adapted a normalization procedure based on subsampling (without replacement) (Li et al.

520 2013a); the consistency of the subsample was validated using bootstrapping. The

521 subsampling, without replacement, was done on the redundant set of reads (before genome

522 matching, with the ncRNAs incident reads removed). The proportion of genome matching

523 reads and the variation in gene complexities (coupled with the p2pPCC between the

524 subsamples and the original sample) were used as criteria for consistency of the subsamples.

525 Each sample was first subjected to incremental subsampling in order to investigate the effect

526 on the data structure (complexities, both for non-matching and genome-matching reads) of

527 sampling 95% through to 45% of the data, with successive decreasing steps of 5%. A sample

528 was deemed satisfactory if the proportion of redundant genome matching reads remained

529 constant and the average point-to-point PCC were above 95% as the number of redundant

530 reads was decreased from 95% to 45%. This step represented an empirical determination of

531 the level of subsampling that could be done whilst preserving the original data structure. The

532 second step of the normalization was the subsampling to a fixed total (the minimum

533 sequencing depth of the accepted samples). Samples with low sequencing depths, which

534 would lead to a heavy subsampling for the samples with high read numbers (less than 55%,

535 empirically determined), were treated on a case-by-case basis. A quantile normalization

536 (Bolstad et al. 2003) may be employed after this step to render the distributions fully

537 comparable. The pseudocode is presented in Supplemental Methods 1.

538     Existing procedures which were used for the comparison of the new normalization

539     methods were: scaling normalization (Mortazavi et al. 2008a), for which the scaling total was

540     the mean of the sequencing depths of the compared samples, quantile normalization (Bolstad

541     et al. 2003)  and the normalization approaches from edgeR (Zhou et al. 2014)  and DESeq2

542     (Love et al. 2014a). All were employed using the recommended standard parameters.

543     **Differential Expression call**

544     Existing methods for the DE call are often based on comparing the variability between

545     replicates with the difference between the treatments. However, calculation of variance (or cv)

546     based on a small number of points may often not reflect the true variance of the given

547     gene/transcript (Krzywinski et al. 2013; Blainey et al. 2014; Altman et al. 2015). Moreover,

548     when small numbers of measurements are available, a more conservative approach, which

549     we use here, is to approximate that replicate measurements will fall within the two limits of the

550     maximal interval (Claridge-Chang et al. 2016).

551     The maximal confidence intervals are defined on the minimum and maximum normalized

552     expression levels for the replicated measurements. The amplitude of the DE is calculated on

553     a worst-case scenario, on the proximal ends of the maximal intervals i.e. this method ensures

554     that all points in the treatment are on one side (for up-regulation, above and for down-

555     regulation, below) of the control measurements (Beckers et al. 2017; Collins et al. 2017). As

556     a result of the stringency of this approach all genes called DE using these rules will also be

557     called DE under all statistical tests. In addition, the threshold on the amplitude of the DE (for

558     (Mohorianu et al. 2017) set at 1.5 fold change, in line with the empirical threshold described

559     in Morey et al) prevents the selection of genes with separate but close expression ranges and

560     ensures a higher chance for validation confirmations.

561     DE was calculated using a hierarchical approach and by applying an offset fold change

562     (OFC) method (with offset=20, empirically determined, using the point-to-point PCC, for all

563     replicates within all samples). There were 3 steps to the hierarchical analysis used for the

564  analysis of the *D. melanogaster* transcriptome data. (i) identification of levels for the

565  hierarchical differential expression and the constituent internal classes. For the *D.*

566  *melanogaster* data one 'level' was body part (with HT and A as internal classes) and the other

567  was treatment (with presence or absence of rivals as classes). (ii) the ordering of the

568  hierarchical levels based on the amplitude of differential expression. This was quantified by

569  the width/ spread of the distribution of DE in terms of mean/ median, IQR and min/max values.

570  The amplitude of DE in descending order provided the correct ordering of the levels for the

571  hierarchical DE. (iii) the DE analysis on the proximal ends of the CIs, using OFC (Mohorianu

572  et al. 2011a). The pseudocode is presented in Supplemental Methods 2.

573      The two-step DE procedure (using OFC) consisted of (i) calculation of the list of genes

574  showing DE between body parts, followed by (ii) calculation of the DE between genes in the

575  ± rivals treatment comparisons. Step (i) was conducted on the summed expression levels in

576  the ± rivals pairs (i.e. the ± HT samples combined, and the ± A samples combined, for all time

577  points). The genes were then separated into genes expressed only in HT, only in A, and in

578  both the HT and A. Step (ii) of the DE was then applied on the resulting 3 categories (HT; A;

579  HT+A) using the ± rival condition. We called DE the genes which showed after the second

580  step of the DE described above, of more than 1.5 fold between the treatments (+/- rivals). The

581  DE call as determined by edgeR and DESeq2 were calculated using the default functions and

582  parameters.

583  **DATA ACCESS**

584  *mRNA samples*: (a) *D melanogaster*: males of *D melanogaster* exposed to conspecific rivals

585  (or not) for 3 time periods (GSE55930). (b) *H sapiens*: For the mRNA Human samples, we

586  chose 5 samples from the Pickrell et al. 2010 (Pickrell et al. 2010) study (GSE19480) in order

587  to compare gene expression variation in RNA sequencing between the Argonne and the Yale

588  laboratory sequencing runs. The selected samples were: GSM485369 (NA19144_yale),

589  GSM485380  (NA19144_argonne);  GSM485368  (NA19153_yale),  GSM485383

590  (NA19153_argonne); GSM485367 (NA19201_yale), GSM485381 (NA19201_argonne);

591  GSM485365 (NA19209_yale), GSM485388 (NA19209_argonne); GSM485364

592  (NA19210_yale), GSM485382 (NA19210_argonne). These samples were derived from

593  lymphoblastoid cell lines (LCLs) derived from unrelated individuals from Nigeria (extensively

594  genotyped by the International HapMap Project). The sequencing was done on Illumina GAII,

595  with sequencing reads of 36nt, for the Yale sequencing samples and 46nt for the Argonne

596  sequencing.

597

598  **ACKNOWLEDGEMENTS**

604

605  **COMPETING INTERESTS**

606  The authors declare there are no competing interests.

607 **SUPPLEMENTAL MATERIAL**

608

609 **SUPPLEMENTAL METHODS**

610 **SUPPLEMENTAL METHODS 1 - Subsampling normalization – pseudocode.** A description

611 with details for (1) Incremental subsampling and bootstrapping check for consistency of a

612 sample, and (2) Subsampling to a fixed total.

613 **SUPPLEMENTAL METHODS 2 - Two step (Hierarchical) differential expression (HDE) -**

614 **pseudocode.** A description with technical details for the two step (hierarchical) DE, including

615 the identification of levels in the hierarchy.

616 **SUPPLEMENTAL TABLES**

617 **TABLE S1. Annotation overview for the 02 samples in the *D. melanogaster* dataset.** For

618 each 02 samples (described in Mohorianu et al 2017) we present the number and proportions

619 of reads, matching to the *D. melanogaster* genome (v 6.11) and to the corresponding

620 annotations (exons, introns, 5' and 3' UTRs, ncRNAs and intergenic regions)

621 **TABLE S2. Example of intersection analysis for the 02-A, 02+A, 02-H and 02+H samples**

622 **in the *D. melanogaster* dataset.** Replicates 1 samples 02-A, 02+A, 02-H and 02+H were

623 used to illustrate the proportion of reads mapping simultaneously to pairwise groups of CDSs,

624 exons, 5' and 3' UTRs, introns and intergenic regions. We observed a high proportion of exon

625 matching reads present on 3' and 5' UTR. In the main study we computed expression levels

626 using gene mapping reads.

627 **TABLE S3. Jaccard similarity index on the 02 samples in the *D. melanogaster* dataset**

628 The Jaccard similarity at gene level was computed on the top 1000 most abundant genes in

629 each sample (out of a total of 15 513 genes expressed in at least one sample). As a result, it

630 is not biased by the different number of genes present in each sample. Shown is a 12 by 12

631 matrix of all the original samples compared with each other. Samples are labelled by time

632 point (2h), by ± rivals treatment, by body part (A or HT) and then by replicate number. Each

633    sample tested against itself along the diagonal is therefore 100% similar and shares the top

634    1000 most abundant genes in common. A to A comparisons are shaded in purple, HT to HT

635    comparisons in peach. Samples drawn from the same body parts shared > 90% similarity, and

636    between body parts the similarity dropped to ~50%. Similarity between the ± rivals treatments

637    tended to be higher than between replicates. Two illustrative examples are highlighted, in

638    which ± rivals indices (in red bold) were generally higher than replicate to replicate similarity

639    (blue bold). This highlighted the need for the adapted normalization methods.

640    **TABLE S4. Example of incremental check for subsampling without replacement for**

641    **sample 02EH2 in the *D. melanogaster* dataset.**

642    For sample 02EH2, we present the incremental subsampling, without replacement. To judge

643    whether a sample is consistent, and to determine the consistency threshold, we use the

644    proportion of redundant reads matching to the referenece genome. As a consequence of the

645    incremental subsampling, the complexity increases. A replicate is accepted if it exhibits a

646    similar complexity (and distribution of per-gene complexities) with the other replicates of the

647    same type of sample.

648    **TABLE S5. Results from (A) DEseq2 and (B) edgeR analyses of the *Drosophila***

649    ***melanogaster* qRT-PCR 'validated' gene set**. For the validations we used 3 reference genes

650    and validated 15 A genes and 6 HT genes based on the DE selection using subsampling

651    normalization and hierarchical DE. We investigated whether these genes were called DE by

652    either (A) DESeq2 or (B) edgeR. In Table 1A we present the results for DESeq2, in Table 1B

653    the results for edgeR. For each of the three categories of genes (reference genes, AB genes

654    and HT genes) we show the average of normalized abundances (baseMean for DESeq2 and

655    logCPM for edgeR), the fold change between treatments (log2 FoldChange for DESeq2 and

656    log2FC for edgeR) and the DE p-value and adjusted p-value (used for the DE call).

657

658    # SUPPLEMENTAL FIGURES

659    **FIGURE S1 Analysis framework for mRNA-seq data.**

660     Inputs are shown (sequencing data in FASTQ format, and the corresponding reference

661     genome and transcriptome in FASTA/GFF) and the six main steps: Quality check (QC),

662     alignment, normalization of gene abundances, identification of DE, functional enrichment and

663     finally low-throughput validation.

664     **FIGURE S2**. **Correlation analyses (Pearson, Spearman and Kendall correlation**

665     **coefficients)** between the gene expression levels for the *D. melanogaster* data for (A) all

666     samples, (B) HT samples, (C) A samples. A1, B1, C1 show the PCC; A2, B2, C2 show the

667     SCC; A3, B3 and C3 show the KCC. Each panel shows the distributions of correlation

668     coefficients for all pairwise comparisons. For example, in panel A.1, sample 1 on the x-axis

669     shows the distribution of the n=35 correlation coefficients calculated between the gene

670     expressions in sample 1 compared with gene expressions in all other 35 samples using PCC.

671     The results are presented as a standard boxplots.

672     **FIGURE S3**. **Distribution of point-to-point PCC between gene expression profiles**

673     **against gene expression levels (log$_2$ scale)** for pairwise comparisons for the *D.*

674     *melanogaster* data for the 3 replicates of the 02HT- sample as an example (2h, HT body part,

675     no rivals). Panel a shows replicate 1 vs 2, b replicate 1 vs 3 and C replicate 2 vs 3. Shown are

676     the raw data, prior to normalization. For all replicate comparisons, more variability is

677     consistently observed at lower abundances.

678     **FIGURE S4**. **Point-to-point PCC between the raw and subsampled data** of the *D.*

679     *melanogaster* data. To show the consistency during the subsampling, shown are the point-to-

680     point PCC between the original data and the data incrementally subsampled from 40% to 95%

681     (Panels A to L). On the x-axis is the gene abundance (log$_2$) and on the y-axis the distribution

682     of point-to-point PCCs calculated for each expressed gene.

683     **FIGURE S5**. **Identification of the hierarchy levels for the hierarchical differential**

684     **expression (HDE) analysis based on the distribution of DE for the different classes of**

685     **samples**, i.e. replicates, body parts and ± rivals treatments (for the *D. melanogaster* data).

686    Frequency plots were used to show the distribution of DE between samples. Panel A shows

687    the replicate-replicate DE (blue) and the with/without rivals DE (red) for the abdomen (A)

688    samples. Panel B shows the corresponding data for the HT body part. Panel C shows the

689    distribution of DE for the with/without rivals treatments (blue for HT and green for A samples)

690    and the DE between HT and A (orange).

691    **FIGURE S6**. **Distribution of abundances for the *D. melanogaster* data (for the ± rivals**

692    **treatment DE) for genes identified as DE exclusively by each method**. EdgeR only genes

693    are presented in 5A, DEseq2 only genes in 5B and subsampling normalization only genes in

694    5C. For each gene (FBgn identifier) identified as DE exclusively by each method, the

695    normalized abundance is given for each of the 2h HT and A ± rivals samples. The

696    predominance of leaky genes in the DE calls of edgeR and DESeq2 highlighted the need for

697    the hierarchical DE. The presence of low abundance genes indicated the requirement for an

698    offset for the calculation of the extent of DE.

699    **FIGURE S7**. **Comparison of the coefficient of variation applied on the *D. melanogaster***

700    **data**. On the x-axis is the abundance in $\log_2$ scale, on the y-axis we represent the coefficient

701    of variation (cv), defined as the ratio between the standard deviation and the mean. For clarity,

702    the distributions are represented as standard boxplots. The upper panels (A,B,C,D) show the

703    cv for the original data for A samples, without rivals, A samples with rivals, HT samples without

704    rivals and HT samples with rivals, respectively. The lower panels (E,F,G,H) give the CV for

705    the same samples, after the subsampling normalization. The horizontal lines indicate 0.5 and

706    0,25 cv, to ease visualization. It is clear that the subsampling normalization reduced the

707    variance between the replicates to < 0.25 cv across most abundances (Panels E-H), whereas

708    the cv was much higher across all abundances for the raw data (Panels A-D).

709    **FIGURE S8**. **Analysis of the effect of the subsampling normalization on technical**

710    **(laboratory-laboratory) variation in mRNA-seq for the human mRNA-seq data in**

711    **(Pickrell et al. 2010)**. In the upper plots we show the coefficient of variation (CV) obtained

712 after the subsampling normalization for 5 sequencing pairs (each pair consisted of a Yale

713 laboratory run compared to an Argonne run: A1,A2 = Sample 144; B1,B2 = Sample 153;

714 C1,C2 = Sample 201; D1,D2 = Sample 209; E1,E2 = Sample 210). Shown is the CV against

715 abundance (log$_2$ scale). For all comparisons we achieved lower CVs in comparison to the

716 Zhou et al (2014) analysis of these sample data. In red we represent the CV of these data

717 obtained using edgeR, in blue the CV using DESeq2. It is evident that our subsampling

718 normalization achieved lower CV across all abundances in comparison to both edgeR and

719 DRseq2. Based on these distributions we conclude that the samples from the different

720 laboratories can be rendered comparable using the subsampling approach, i.e. the

721 subsampling normalization removed the technical differences between the two different

722 laboratory runs. In the lower panels, we present the MA plots, after the subsampling

723 normalization, for the same pairs of samples. The tightness of these plots (all falling within

724 ±0.5 OFC) supports the conclusion that the subsampling has rendered these samples derived

725 from sequencing in different laboratories highly comparable.

726

727

728

**REFERENCES**

Aanes H, Winata C, Moen LF, Ostrup O, Mathavan S, Collas P, Rognes T, Alestrom P. 2014. Normalization of RNA-sequencing data from samples with varying mRNA levels. *PLoS One* 9:e89158.

Altman N, Krzywinski M. 2015. Points of significance: Sources of variation. *Nat Methods* 12:5-6.

Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome Biol* 11:R106.

Andrews S. 2010. FastQC: A quality control tool for high throughput sequence data. *abraham Bioinformatics*.

Beckers ML, Mohorianu I, Stocks MB, Applegate C, Dalmay T, Moulton V. 2017. Comprehensive processing of high throughput small RNA sequencing data including quality checking, normalization and differential expression analysis using the UEA sRNA Workbench. *RNA*.

Blainey P, Krzywinski M, Altman N. 2014. Points of significance: replication. *Nat Methods* 11:879-880.

Bolstad BM, Irizarry RA, Astrand M, Speed TP. 2003. A comparison of normalization methods for high density oligonucleotide array data based on variance and bias. *Bioinformatics* 19:185-193.

Bullard JH, Purdom E, Hansen KD, Dudoit S. 2010. Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics* 11:94.

Claridge-Chang A, Assam PN. 2016. Estimation statistics should replace significance testing. *Nat Methods* 13:108-109.

Collins DH, Mohorianu I, Beckers M, Moulton V, Dalmay T, Bourke AF. 2017. MicroRNAs Associated with Caste Determination and Differentiation in a Primitively Eusocial Insect. *Sci Rep* 7:45674.

Conesa A, Madrigal P, Tarazona S, et al. 2016a. A survey of best practices for RNA-seq data analysis. *Genome Biol* 17:13.

Conesa A, Madrigal P, Tarazona S, et al. 2016b. A survey of best practices for RNA-seq data analysis. *Genome biology* 17:13.

Consortium SM-I. 2014. A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium. *Nat Biotechnol* 32:903-914.

Cui X, Churchill GA. 2003. Statistical tests for differential expression in cDNA microarray experiments. *Genome biology* 4:210.

DeLuca DS, Levin JZ, Sivachenko A, Fennell T, Nazaire MD, Williams C, Reich M, Winckler W, Getz G. 2012. RNA-SeQC: RNA-seq metrics for quality control and process optimization. *Bioinformatics* 28:1530-1532.

Dillies MA, Rau A, Aubert J, et al. 2013a. A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis. *Briefings in bioinformatics* 14:671-683.

Dillies MA, Rau A, Aubert J, et al. 2013b. A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis. *Brief Bioinform* 14:671-683.

Evans C, Hardin J, Stoebel DM. 2017. Selecting between-sample RNA-Seq normalization methods from the perspective of their assumptions. *Brief Bioinform*.

Fang Z, Cui X. 2011. Design and validation issues in RNA-seq experiments. *Brief Bioinform* 12:280-287.

Gierlinski M, Cole C, Schofield P, et al. 2015. Statistical models for RNA-seq data derived from a two-condition 48-replicate experiment. *Bioinformatics* 31:3625-3630.

Hansen KD, Brenner SE, Dudoit S. 2010. Biases in Illumina transcriptome sequencing caused by random hexamer priming. *Nucleic acids research* 38:e131.

Jaccard P. 1901. Étude comparative de la distribution florale dans une portion des Alpes et des Jura. *Bulletin de la Société Vaudoise des Sciences Naturelles*:547-579.

783    Jackson TJ, Spriggs RV, Burgoyne NJ, Jones C, Willis AE. 2014. Evaluating bias-reducing protocols
784         for RNA sequencing library preparation. *BMC Genomics* 15:569.
785    Jayaprakash AD, Jabado O, Brown BD, Sachidanandam R. 2011a. Identification and remediation
786         of biases in the activity of RNA ligases in small-RNA deep sequencing. *Nucleic Acids Res*
787         39:e141.
788    Jayaprakash AD, Jabado O, Brown BD, Sachidanandam R. 2011b. Identification and remediation
789         of biases in the activity of RNA ligases in small-RNA deep sequencing. *Nucleic acids*
790         *research* 39:e141.
791    Jia B, Xu S, Xiao G, Lamba V, Liang F. 2017. Learning gene regulatory networks from next
792         generation sequencing data. *Biometrics*.
793    Krzywinski M, Altman N. 2013. Significance, P values and t-tests. *Nat Methods* 10:1041-1042.
794    Li J, Tibshirani R. 2013a. Finding consistent patterns: a nonparametric approach for identifying
795         differential expression in RNA-Seq data. *Statistical methods in medical research* 22:519-
796         536.
797    Li J, Tibshirani R. 2013b. Finding consistent patterns: a nonparametric approach for identifying
798         differential expression in RNA-Seq data. *Stat Methods Med Res* 22:519-536.
799    Li P, Piao Y, Shon HS, Ryu KH. 2015. Comparing the normalization methods for the differential
800         analysis of Illumina high-throughput RNA-Seq data. *BMC Bioinformatics* 16:347.
801    Li S, Labaj PP, Zumbo P, et al. 2014. Detecting and correcting systematic variation in large-scale
802         RNA sequencing data. *Nat Biotechnol* 32:888-895.
803    Lin Y, Golovnina K, Chen ZX, Lee HN, Negron YL, Sultana H, Oliver B, Harbison ST. 2016.
804         Comparison of normalization and differential expression analyses using RNA-Seq data
805         from 726 individual Drosophila melanogaster. *BMC Genomics* 17:28.
806    Love MI, Huber W, Anders S. 2014a. Moderated estimation of fold change and dispersion for
807         RNA-seq data with DESeq2. *Genome biology* 15:550.
808    Love MI, Huber W, Anders S. 2014b. Moderated estimation of fold change and dispersion for
809         RNA-seq data with DESeq2. *Genome Biol* 15:550.
810    McIntyre LM, Lopiano KK, Morse AM, Amin V, Oberg AL, Young LJ, Nuzhdin SV. 2011. RNA-seq:
811         technical variability and sampling. *BMC Genomics* 12:293.
812    Mohorianu, II, Bretman A, Smith DT, Fowler E, Dalmay T, Chapman T. 2017. Genomic responses
813         to socio-sexual environment in male Drosophila melanogaster exposed to conspecific
814         rivals. *RNA*.
815    Mohorianu I, Lopez-Gomollon S, Schwach F, Dalmay T, Moulton V. 2012. FiRePat - Finding
816         Regulatory Patterns between sRNAs and Genes. *Wiley Interdisc. Rew.: Data Mining and*
817         *Knowledge Discovery* 2:273-284.
818    Mohorianu I, Schwach F, Jing R, Lopez-Gomollon S, Moxon S, Szittya G, Sorefan K, Moulton V,
819         Dalmay T. 2011a. Profiling of short RNAs during fleshy fruit development reveals stage-
820         specific sRNAome expression patterns. *The Plant journal : for cell and molecular biology*
821         67:232-246.
822    Mohorianu I, Schwach F, Jing R, Lopez-Gomollon S, Moxon S, Szittya G, Sorefan K, Moulton V,
823         Dalmay T. 2011b. Profiling of short RNAs during fleshy fruit development reveals stage-
824         specific sRNAome expression patterns. *Plant J* 67:232-246.
825    Mohorianu I, Stocks MB, Wood J, Dalmay T, Moulton V. 2013. CoLIde: a bioinformatics tool for
826         CO-expression-based small RNA Loci Identification using high-throughput sequencing
827         data. *RNA biology* 10:1221-1230.
828    Morey JS, Ryan JC, Van Dolah FM. 2006. Microarray validation: factors influencing correlation
829         between oligonucleotide microarrays and real-time PCR. *Biological procedures online*
830         8:175-193.
831    Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. 2008a. Mapping and quantifying
832         mammalian transcriptomes by RNA-Seq. *Nature methods* 5:621-628.
833    Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. 2008b. Mapping and quantifying
834         mammalian transcriptomes by RNA-Seq. *Nat Methods* 5:621-628.
835    Ozsolak F, Milos PM. 2011. RNA sequencing: advances, challenges and opportunities. *Nature*
836         *reviews. Genetics* 12:87-98.

837 Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. 2017. Salmon provides fast and bias-aware
838      quantification of transcript expression. *Nat Methods* 14:417-419.
839 Pickrell JK, Marioni JC, Pai AA, et al. 2010. Understanding mechanisms underlying human gene
840      expression variation with RNA sequencing. *Nature* 464:768-772.
841 Prufer K, Stenzel U, Dannemann M, Green RE, Lachmann M, Kelso J. 2008. PatMaN: rapid
842      alignment of short sequences to large databases. *Bioinformatics* 24:1530-1531.
843 Rapaport F, Khanin R, Liang Y, Pirun M, Krek A, Zumbo P, Mason CE, Socci ND, Betel D. 2013.
844      Comprehensive evaluation of differential gene expression analysis methods for RNA-seq
845      data. *Genome Biol* 14:R95.
846 Risso D, Schwartz K, Sherlock G, Dudoit S. 2011. GC-content normalization for RNA-Seq data.
847      *BMC Bioinformatics* 12:480.
848 Robinson DG, Storey JD. 2014. subSeq: determining appropriate sequencing depth through
849      efficient read subsampling. *Bioinformatics* 30:3424-3426.
850 Robinson DG, Wang JY, Storey JD. 2015. A nested parallel experiment demonstrates differences
851      in intensity-dependence between RNA-seq and microarrays. *Nucleic Acids Res* 43:e131.
852 Robinson MD, Oshlack A. 2010. A scaling normalization method for differential expression
853      analysis of RNA-seq data. *Genome biology* 11:R25.
854 Roca CP, Gomes SI, Amorim MJ, Scott-Fordsmand JJ. 2017. Variation-preserving normalization
855      unveils blind spots in gene expression profiling. *Sci Rep* 7:42460.
856 Schmieder R, Lim YW, Rohwer F, Edwards R. 2010. TagCleaner: Identification and removal of
857      tag sequences from genomic and metagenomic datasets. *BMC bioinformatics* 11:341.
858 Schurch NJ, Schofield P, Gierlinski M, et al. 2016. How many biological replicates are needed in
859      an RNA-seq experiment and which differential expression tool should you use? *RNA*
860      22:839-851.
861 Soneson C, Delorenzi M. 2013. A comparison of methods for differential expression analysis of
862      RNA-seq data. *BMC Bioinformatics* 14:91.
863 Sorefan K, Pais H, Hall AE, Kozomara A, Griffiths-Jones S, Moulton V, Dalmay T. 2012. Reducing
864      ligation bias of small RNAs in libraries for next generation sequencing. *Silence* 3:4.
865 Stupnikov A, Glazko GV, Emmert-Streib F. 2015. Effects of subsampling on characteristics of
866      RNA-seq data from triple-negative breast cancer patients. *Chin J Cancer* 34:427-438.
867 Torres-Oliva M, Almudi I, McGregor AP, Posnien N. 2016. A robust (re-)annotation approach to
868      generate unbiased mapping references for RNA-seq-based analyses of differential
869      expression across closely related species. *BMC Genomics* 17:392.
870 Trapnell C, Roberts A, Goff L, et al. 2012. Differential gene and transcript expression analysis of
871      RNA-seq experiments with TopHat and Cufflinks. *Nature protocols* 7:562-578.
872 Wagner GP, Kin K, Lynch VJ. 2012. Measurement of mRNA abundance using RNA-seq data:
873      RPKM measure is inconsistent among samples. *Theory Biosci* 131:281-285.
874 Wang L, Wang S, Li W. 2012. RSeQC: quality control of RNA-seq experiments. *Bioinformatics*
875      28:2184-2185.
876 Wang Z, Gerstein M, Snyder M. 2009. RNA-Seq: a revolutionary tool for transcriptomics. *Nature*
877      *reviews. Genetics* 10:57-63.
878 Zhou X, Lindsay H, Robinson MD. 2014. Robustly detecting differential expression in RNA
879      sequencing data using observation weights. *Nucleic acids research* 42:e91.

880

881    **FIGURES**

882

883    **FIGURE 1 Distributions of complexities, calculated at the gene level, on the *D.***

884    ***melanogaster* mRNA-seq data.**

885    Transcript abundances (x-axis, $\log_2$ scale) are plotted against the absolute difference in

886    complexities, i.e. the non-redundant/redundant (NR/R) ratio (y-axis) for all genes and all

887    biological replicate comparisons (a = replicate 1 vs 2; b = replicate 1 vs 3; c = replicate 2 vs 3.

888    Example data shown are for the three original replicates of the 02+H (2 hours, rivals present,

889    head thorax) samples. The differences in complexities were calculated on the raw data (top

890    row), and on the data after subsampling normalization with replacement (middle row) and

891    without replacement (bottom row). Red horizontal lines indicate 0.05 and 0.1 differences in

892    complexity. Before subsampling the complexity differences were frequently > 0.1. The

893    subsampling approaches (with or without replacement) rendered the biological replicates

894    more comparable and reduced the complexity differences to < 0.1 across all transcript

895    abundances. In addition, the subsampling without replacement maintained the conclusion that

896    the third replicate (R3) was problematic, whereas the subsampling with replacement masked

897    this conclusion.

898    **FIGURE 2 Comparison of results obtained using the subsampling with or without**

899    **replacement.** On the top row we present the MA plots on the gene expression levels,

900    normalized using either the with- or without- replacement approaches, for the three replicates

901    of the 02+H sample. Although the variability between the two approaches is contained within

902    the +/- 0.5 $\log_2$(OFC), we observe a higher variability in expression for the low abundance

903    genes. In subplot (B) we present the presence plots for the gene FBgn0033865 for each of

904    the three replicates (the individual panels) obtained using either the subsampling without

905    replacement (black solid line) or subsampling with replacement (red solid line). The arrows

906    indicate the regions where the two approaches provide different answers. The arrow indicating

907    the first exon of the gene highlights the difference observed for the third replicate (02+H, rep3).

908 **FIGURE 3 Distribution of DE as calculated by using fold change (FC) versus offset fold**

909 **change (OFC) and the effect of incorporating hierarchical DE (HDE).** Shown are MA plots

910 (x-axis showing gene abundance ($log_2$), y-axis indicating FC/OFC for replicate-to-replicate

911 comparisons for the 2h samples. Panels A1, B1, C1 show 02A- comparisons, panels A2, B2,

912 C2 for 02A+ samples, A3, B3, C3 for 02HT- and A4, B4, C4 for 02HT+ samples (Sample

913 codes: 02 = 2h of exposure, A = abdomen, HT = head-thorax, + = with rivals, - = without rivals).

914 Panel a shows the distribution of DE calculated using FC, showing how the low abundance

915 genes distort the distribution of DE. Panel b shows the DE distribution using OFC (offset=20).

916 Here most of the low abundance genes were excluded. Panel c shows the DE distribution

917 following hierarchical DE analysis using OFC for A- and HT-specific genes highlighting the

918 elimination of low abundance, potentially spurious, DE. The red horizontal lines denote 0 $log_2$

919 FC/OFC and the blue lines ± 0.5 $log_2$ FC/OFC.

920 **FIGURE 4 Comparison of expression distributions resulting from different**

921 **normalization methods.** Shown are standard boxplots of normalized gene expressions. On

922 the x-axis are the different samples (e.g. 02A-1 = 2h time point abdomen body part, no rivals,

923 replicate 1) and the on the y-axis the $log_2$ gene expression. Panel A shows the raw expression

924 levels, B the RPM normalization to a fixed total of 50M reads, C the quantile normalization,

925 D1 the subsampling (with replacement) normalization to a total of 50M reads D2, the

926 subsampling without replacement, E the DESeq2 normalization and F the edgeR

927 normalization. Effective normalization (e.g. C and D) is observed when the distributions
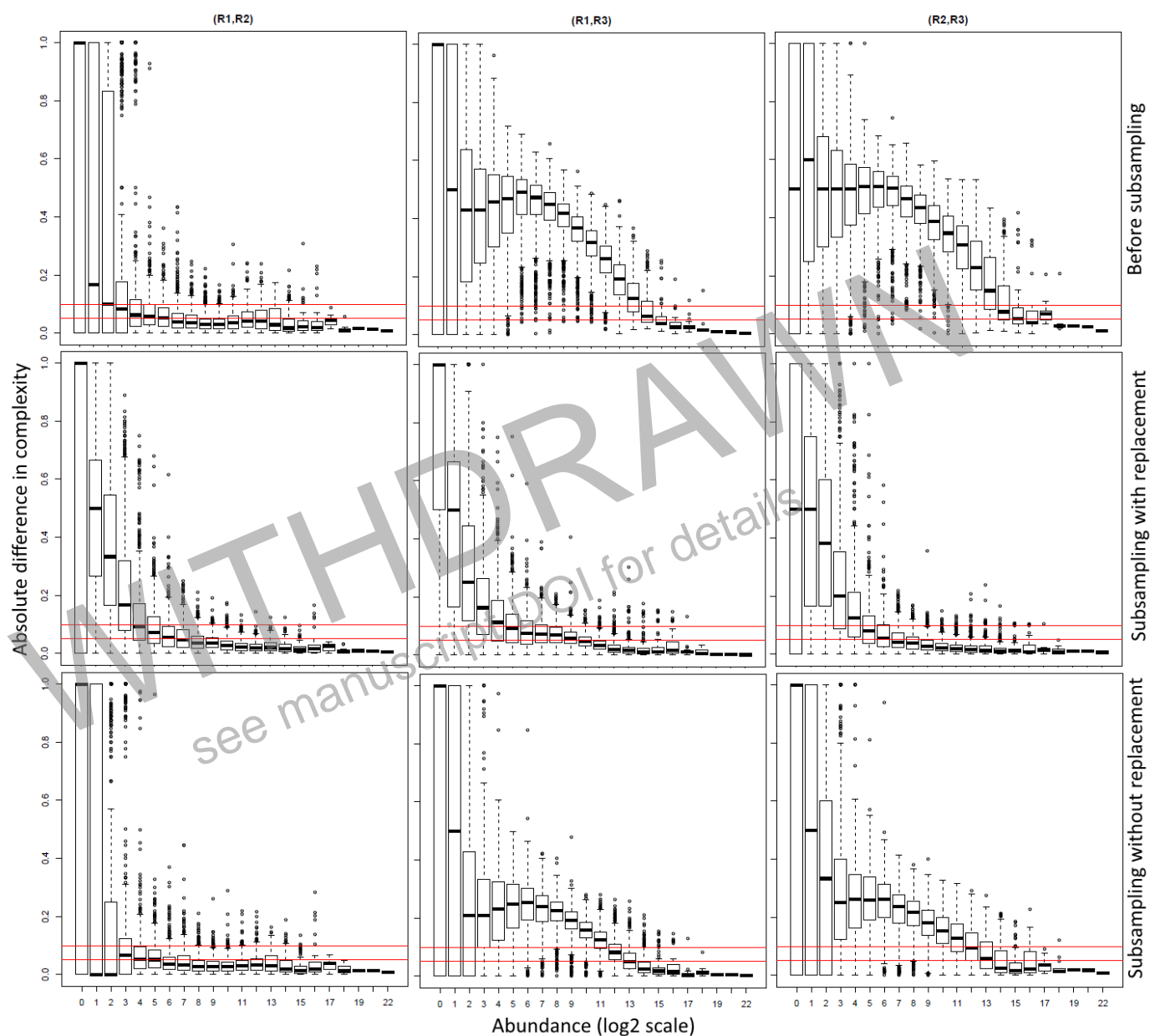
928 become most comparable.

929 **FIGURE 5 Comparison of distribution of DE obtained using the subsampling**

930 **normalization and HDE, DESeq2 and edgeR.** MA plots, with x-axis showing $log_2$ average

931 abundances against OFC with an offset of 20 (Panel A) and FC (Panels B and C). The

932 example shown ism for the 02HT ± rivals DE comparison. The red line indicates 0 $log_2$ FC/OFC

933 and the blue lines ±0.5 $log_2$ FC/OFC. Red data points represent the genes 'called' differentially

934 expressed by each of the methods. Panel A shows the results for subsampling normalization
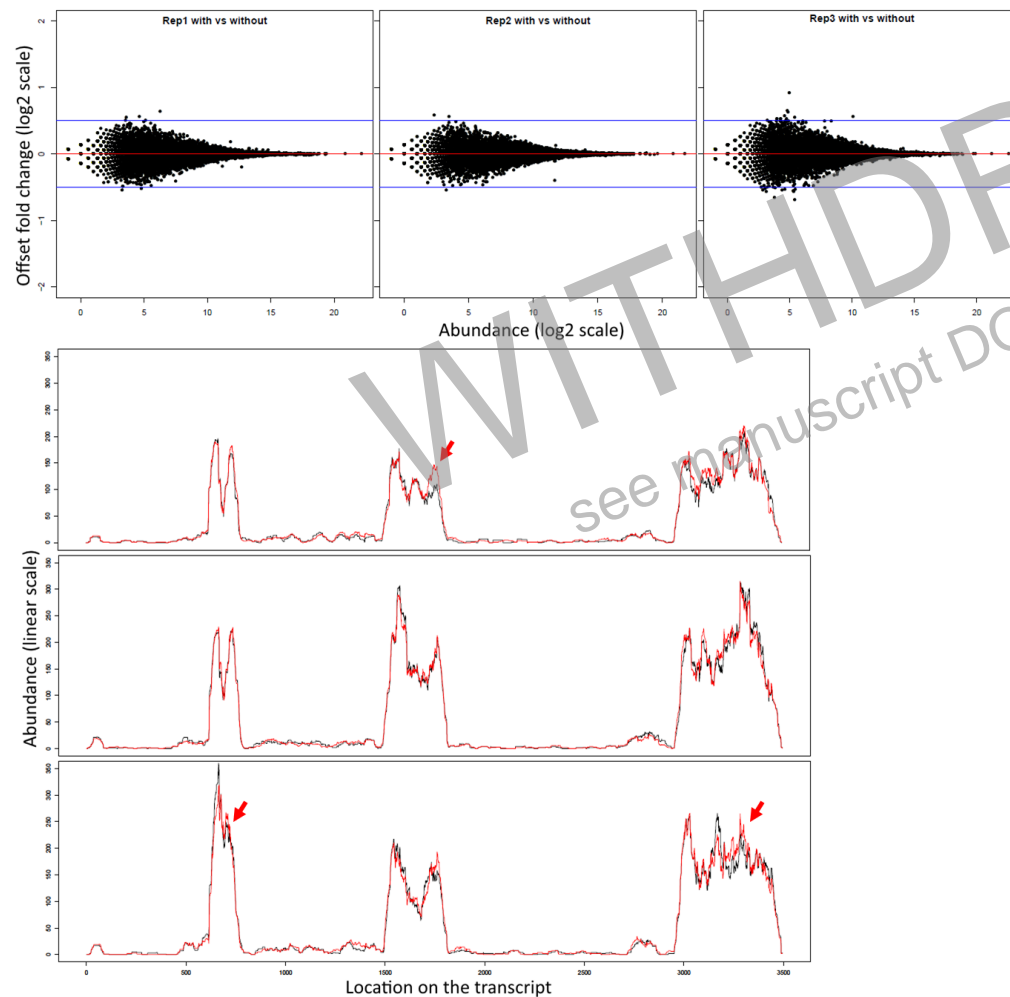
935    with DE calculated using the hierarchical approach, Panel B for DEseq2 and Panel C for

936    edgeR. Panel d shows a Venn diagram identifying the number of differentially expressed

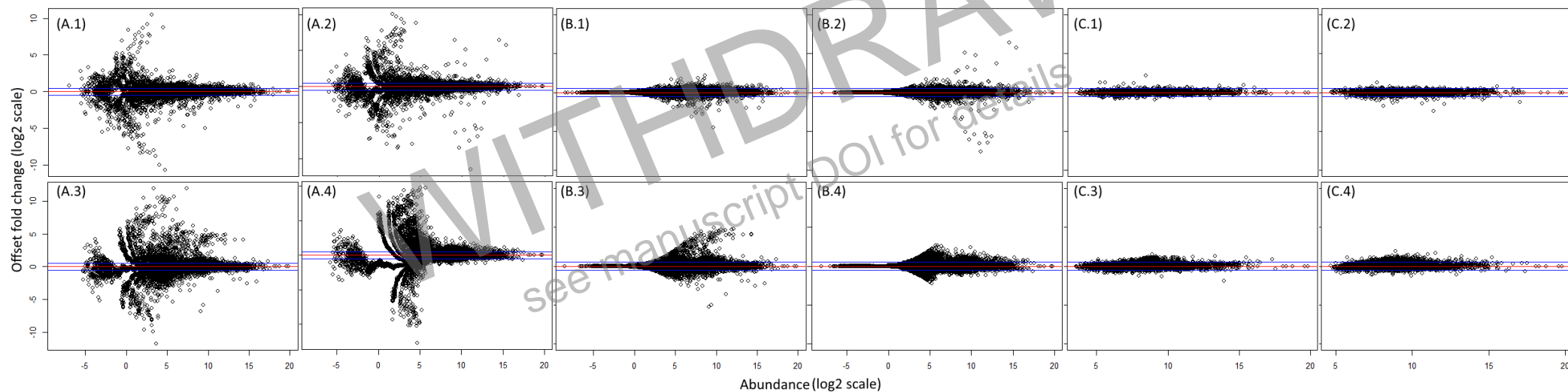937    genes identified by two or more methods versus uniquely by each.

938

WITHDRAWN

see manuscript DOI for details

939     **FIGURE 1**



940

941 **FIGURE 2**
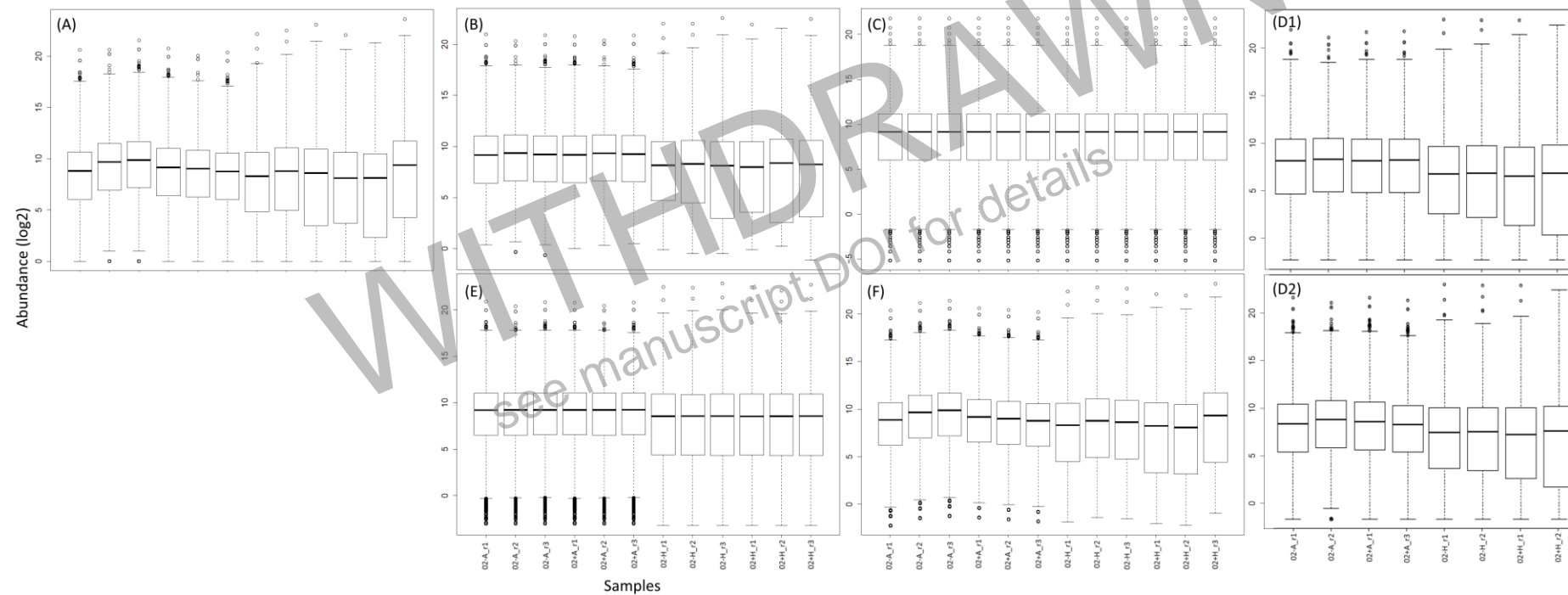


942

943 **FIGURE 3**
944
945
946



947
948
949
950
951
952
953
954
955
956
957
958

959 **FIGURE 4**
960



961
962
963
964
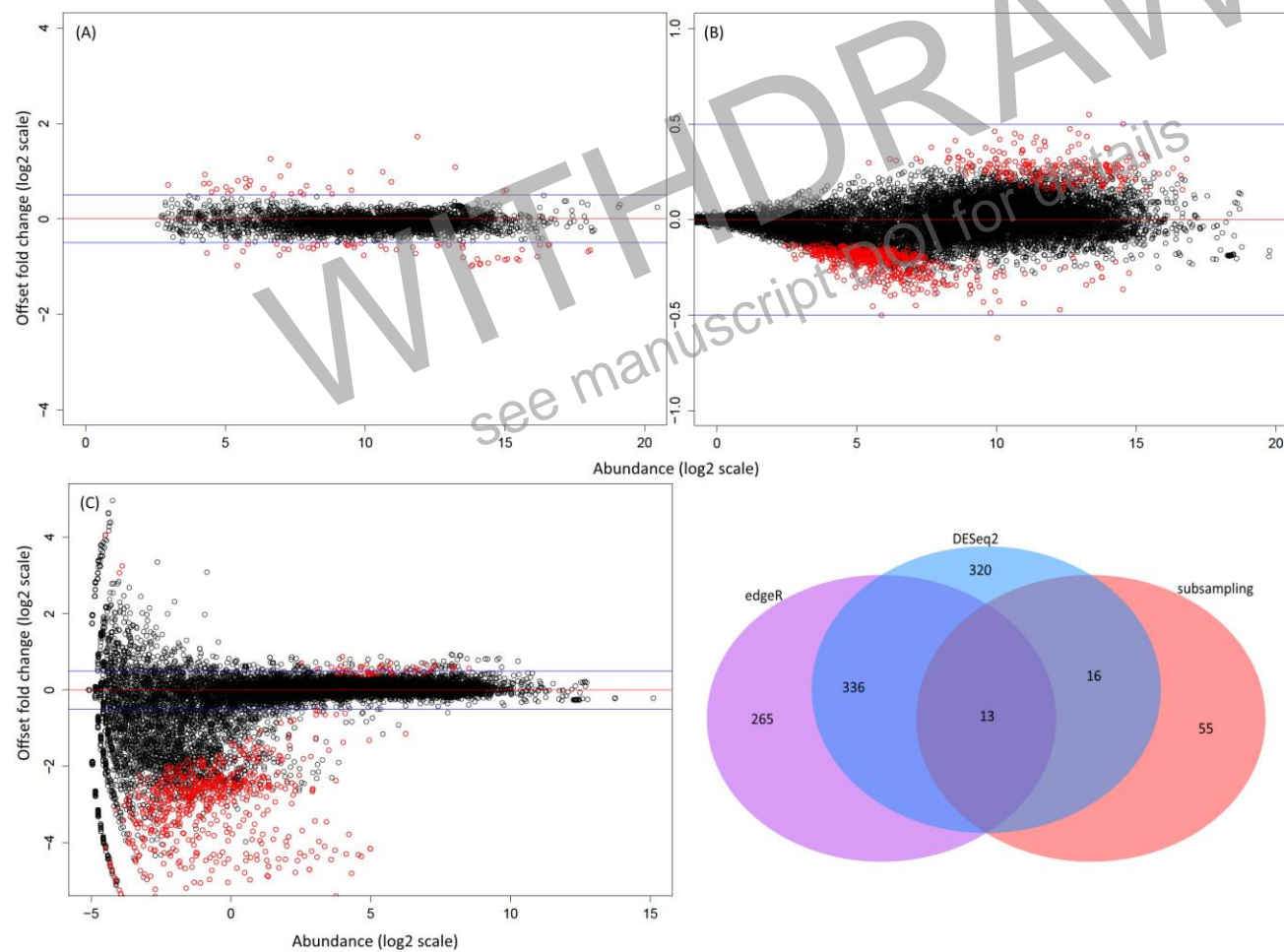
965 **FIGURE 5**

966

WITHDRAWN

see manuscript DOI for details