

1 **NormExpression: an R package**

2 **to normalize gene expression data using evaluated methods**

3
4 Zhenfeng Wu^{12§}, Weixiang Liu³⁴, Haishuo Ji²⁵, Deshui Yu²,

5 Hua Wang², Liu Lin², Jishou Ruan^{16*}, Shan Gao^{25*}

6
7 1. School of Mathematical Sciences, Nankai University, Tianjin 300071, P.R.China.

8 2. College of Life Sciences, Nankai University, Tianjin, Tianjin 300071, P.R.China.

9 3. School of Biomedical Engineering, Health Science Center, Shenzhen University,
10 Shenzhen 518060, P.R.China.

11 4. Guangdong Provincial Key Laboratory of Biomedical Measurements and Ultrasound
12 Imaging, Shenzhen 518060, P.R.China.

13 5. Institute of Statistics, Nankai University, Tianjin 300071, P.R.China.

14 6. State Key Laboratory of Medicinal Chemical Biology, Nankai University, Tianjin
15 300071, P.R.China.

16
17
18
19 § These authors contributed equally to this paper.

20 * The corresponding authors.

21 SG: gao_shan@mail.nankai.edu.cn

22 JR: jsruan@nankai.edu.cn

23

24 **Abstract**

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

Data normalization is a crucial step in the gene expression analysis as it determines the validity of its downstream analyses. Although many metrics has been designed to evaluate the relative success of these methods, the results by different metrics did not show consistency. Based on the previous work, we designed a new metric named Area Under normalized CV threshold Curve (AUCVC) to evaluate 13 commonly used normalization methods and achieved consistency in our evaluation results using both bulk RNA-seq and scRNA-seq data from the same library construction protocol. These gene expression data, normalization methods and evaluation metrics have been included in an R package named NormExpression. NormExpression provides a framework for researchers to select normalization methods with a fast and simple way to evaluate different methods, particularly some data-driven methods or their own methods.

Keyword: gene expression; normalization; evaluation; R package; scRNA-seq

54 Introduction

55 Global gene expression analysis provides quantitative information about the
56 population of RNA species in cells and tissues [1]. High-throughput technologies to
57 measure global gene expression levels started with Serial Analysis of Gene Expression
58 method (SAGE) and are widely used with microarray and RNA-seq [2]. Recently, single-
59 cell RNA sequencing (scRNA-seq) has been used to simultaneously measure the expression
60 levels of genes from a single-cell and to provide a higher resolution of cellular differences
61 than bulk RNA-seq, which can only produce an expression value for each gene by
62 averaging its expression levels across a large population of cells [3]. Gene expression raw
63 data from these high-throughput technologies must be normalized to remove technical
64 variation so that meaningful biological comparisons can be made. Data normalization is a
65 crucial step in the gene expression analysis as it determines the validity of its downstream
66 analyses. Although the significance of gene expression data normalization has been
67 demonstrated [4], how to successfully select a normalization method is still a controversial
68 problem, particularly for scRNA-seq data.

69 Basically, two classes of methods are available to normalize gene expression data.
70 They are the control-based normalization and the average-bulk normalization. The former
71 class of methods assumes the total expression level summed over a small group of genes is
72 approximately the same across all the samples. The latter class of methods assumes most of
73 genes are not Differentially Expressed (DE) genes across all the samples. The control-based
74 normalization often uses RNA from a group of internal control genes (*e.g.* housekeeping
75 genes) or external spike-in RNA (*e.g.* ERCC RNA [5]), while the average-bulk
76 normalization is more commonly used for their universality. Five average-bulk
77 normalization methods designed to normalize bulk RNA-seq data are library size, median of
78 the ratios of observed counts that is also referred to as the DESeq method [6], Relative Log
79 Expression (RLE), upperquartile (UQ) and Trimmed Mean of M values (TMM) [7].
80 Recently, three new methods have been introduced as Total Ubiquitous (TU), Network
81 Centrality Scaling (NCS) and Evolution Strategy (ES) with best performance among 15
82 tested methods [8].

83 Although many metrics has been designed to evaluate the relative success of these
 84 methods, the results by different metrics did not show consistency. In 2013, Gustavo *et al.*
 85 designed two novel and mutually independent metrics to evaluate 15 normalization methods
 86 and achieved consistent results using bulk RNA-seq data [8]. Based on their work, we
 87 designed a new metric named Area Under normalized CV threshold Curve (AUCVC) and
 88 tested it using both bulk RNA-seq and scRNA-seq data from the same library construction
 89 protocol. As a result, the evaluation by both our metric AUCVC and their metrics achieved
 90 consistency. On the other hand, with many new normalization methods developed,
 91 researchers need a fast and simple way to evaluate different methods, particularly some
 92 data-driven methods or their own methods rather than obtain information from published
 93 evaluation results, which could have bias or mistakes, *e.g.* misunderstanding of RLE, UQ
 94 and TMM methods [9]. To satisfy this demand, we developed an R package
 95 NormExpression to include gene expression data, normalization methods and evaluation
 96 metrics used in this study and provide a framework for researchers to evaluate and select
 97 methods for the normalization of their gene expression data.

A.

$$\begin{bmatrix} x_{11} & \cdots & x_{1j} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots & & \vdots \\ x_{i1} & \cdots & x_{ij} & \cdots & x_{in} \\ \vdots & & \vdots & \ddots & \vdots \\ x_{m1} & \cdots & x_{mj} & \cdots & x_{mn} \end{bmatrix} \times \begin{bmatrix} f_1 \\ \vdots \\ f_j \\ \vdots \\ f_n \end{bmatrix}$$

B.

For all methods: $f_j = f_j / \exp\left[\frac{1}{n} \sum_{j=1}^n \log(f_j)\right]$

Library size: $f_j = 10^6 / N_j^*$

DESeq(RLE): $f_j = 1/s_j$

$$s_j = \exp\left[\text{median}\left(\log(x_{ij}) - \frac{1}{n} \sum_{k=1}^n \log(x_{ik})\right)\right]$$

Upperquartile: $f_j = 10^6 / N_j s_j$ $s_j = Q_3(x_{ij} / N_j)$

98

99

Figure 1

100 **Results**

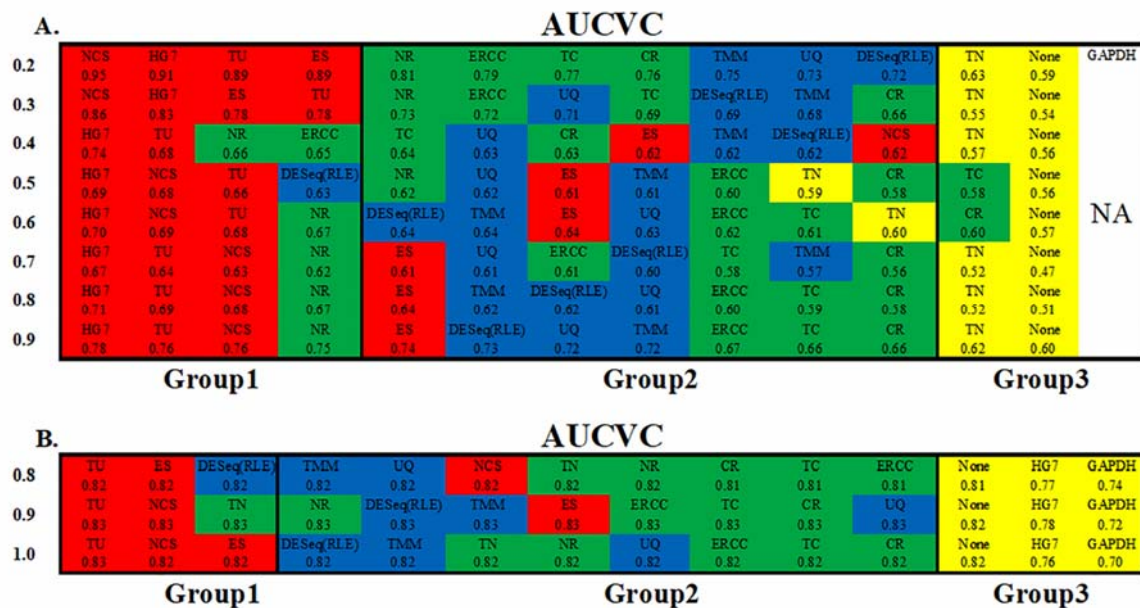
101 **Normalization factor and library size**

102 In total, 13 normalization methods (**Materials and Methods**) have been included in
103 the R package NormExpression. They are Housekeeping Genes (HG7), External RNA
104 Control Consortium (ERCC), Total Read Number (TN), Total Read Count (TC), Cellular
105 RNA (CR), Nuclear RNA (NR), the ratios of observed counts (DESeq), Relative Log
106 Expression (RLE), upperquartile (UQ), Trimmed Mean of M values (TMM), Total
107 Ubiquitous (TU), Network Centrality Scaling (NCS) and Evolution Strategy (ES). Currently,
108 all the commonly used methods are used to normalize a raw gene expression matrix (n
109 samples by m genes) by the multiplication of a factor to each column of it and produce a
110 normalized gene expression matrix (**Figure 1A**). This factor is named as normalization
111 factor in the package NormExpression or scaling factor in TU, NCS and ES methods. In
112 NormExpression, the reciprocal of normalization factor is named as library size (**Figure**
113 **1B**), which is also named as size factor in the Bioconductor package DESeq [6]. Definitions
114 of normalization factor and size factor in the Bioconductor package edgeR [7] are different
115 from the definition of normalization factor in NormExpression and the definition of size
116 factor in DESeq. RLE, UQ and TMM in edgeR produce normalization factors to adjust
117 library sizes, which should be used to calculate the Counts Per Million (CPM) for the
118 normalization of gene expression data and CPM should be calculated by one million
119 multiplying reciprocals of adjusted library sizes (**Figure 1B**). However, edgeR provides a
120 function named calcNormFactors to produce normalization factors for library-size
121 adjustment, which have been wrongly used for the normalization of gene expression data in
122 many studies [9]. Since ES, HG7, ERCC, TN, TC, CR, NR and TU produce normalization
123 factors by the estimation of library sizes as CPM, their normalization factors are amplified
124 by one million for a uniform representation (**Figure 1B**) in NormExpression. DESeq, RLE,
125 UQ and TMM have been modified to ignore zero values for both scRNA-seq and bulk
126 RNA-seq data and the resulting normalization factors need be further normalized by their
127 geometric mean values (**Figure 1B**). UQ and TMM use library sizes estimated by NR. After
128 modification, RLE is identical to DESeq. We verified that these modifications did not
129 change the evaluation or normalization results.

130

131 Evaluation of 13 normalization methods

132 In the previous study [8], Gustavo *et al.* had quantified success of normalization
 133 methods by the number of uniform genes (**Materials and Methods**) and used the
 134 Coefficient of Variation (CV) cutoff 0.25 to determine the number of uniform genes for
 135 each method. This metric was designed based on the theory that the relative values among
 136 different normalization methods were quite stable, although the absolute number of uniform
 137 genes depended on the cutoff value. However, it is almost impossible to determine a CV
 138 cutoff for scRNA-seq data since the CV in scRNA-seq data has a much more large dynamic
 139 range than that in bulk RNA-seq data.



140

141

Figure 2.

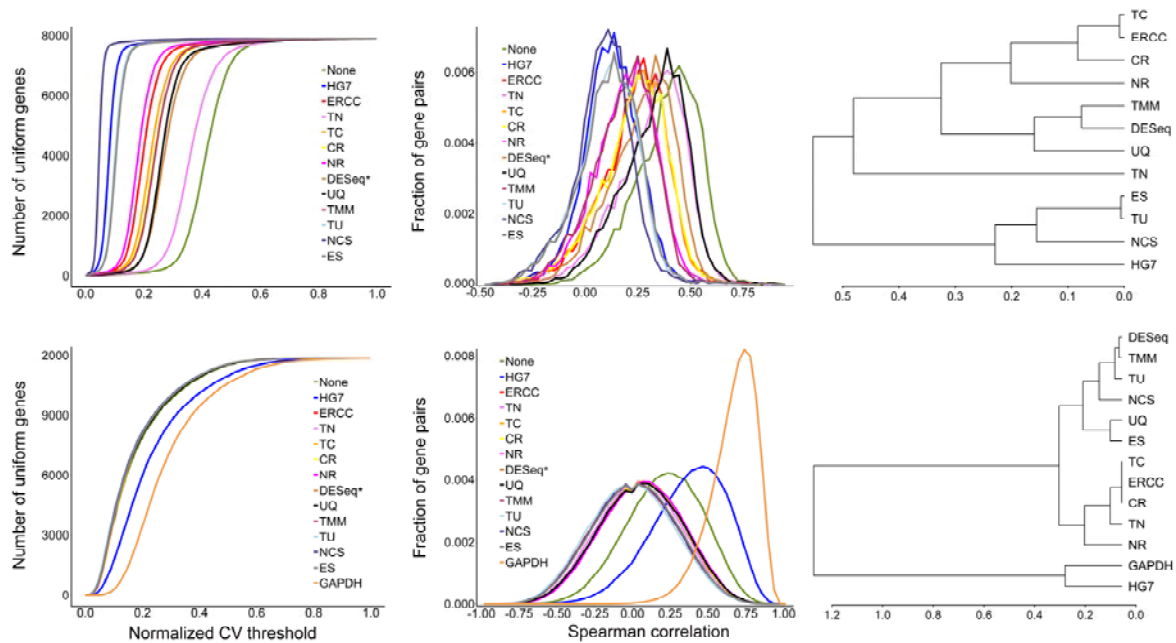
142 Inspired by Area Under the receiver operating characteristic Curve (AUC) [10], we
 143 designed a new metric named Area Under normalized CV threshold Curve (AUCVC) to
 144 evaluate 13 normalization methods using one scRNA-seq dataset scRNA663 and one bulk
 145 RNA-seq dataset bkrNA18 (**Materials and Methods**). A single housekeeping gene
 146 GAPDH was also used for comparison in the evaluation of normalization methods using
 147 bulk RNA-seq data, but it was not available for that using scRNA-seq data due to zero
 148 counts of GAPDH in many samples. Parameter grid of non-zero ratio (**Materials and**
 149 **Methods**) from 0.2 to 0.9 for scRNA663 and from 0.8 to 1 for bkrNA18 was used to

150 produce AUCVC values of all methods (**Figure 2**). For each non-zero ratio, TU used the
151 maximum AUCVC, which had been determined by testing all possible combinations of
152 presence rate, lower and upper cutoffs (**Materials and Methods**) at 5% resolution. The
153 presence rate was tested from 0.2 to 0.6 for scRNA663 and set 1 for bkRNA18. The lower
154 cutoff was tested from 5% to 40% and the upper cutoff was tested from 60% to 95%. In
155 addition, the calculation only considered each combination of lower and upper cutoffs
156 which produced ubiquitous genes (**Materials and Methods**) more than 1,00 for scRNA663
157 and more than 1,000 for bkRNA18. For each non-zero ratio, NCS and ES used the
158 ubiquitous genes produced by the TU method, when it achieved the maximum AUCVC.
159 The raw gene expression matrix (None) was also used to produce AUCVC values for
160 comparison.

161 The evaluation results using both scRNA663 and bkRNA18 achieved consistency that
162 all the normalization methods were classified into three groups (**Figure 2**) based on their
163 AUCVC values sorted in descending order. The first group including TU, NCS and ES
164 achieved the best performances using both scRNA663 and bkRNA18. The second group
165 including ERCC, TC, CR, NR, DESeq, RLE, UQ and TMM achieved medial performances
166 using both scRNA663 and bkRNA18. In the second group, ERCC, TC, CR and NR
167 outperformed DESeq, RLE, UQ and TMM using scRNA663, while DESeq, RLE, UQ and
168 TMM outperformed ERCC, TC, CR and NR using bkRNA18. The third group achieved the
169 poorest performances, including TN and None for scRNA663 (**Figure 2A**) and HG7,
170 GAPDH and None for bkRNA18 (**Figure 2B**). HG7 and GAPDH achieved the poorest
171 performances using bkRNA18, which suggested that a predefined set of housekeeping
172 genes could not be appropriate guides for data normalization of bulk RNA-seq data.
173 However, it could be coincidental that HG7 was classified into the first group using
174 scRNA663. TN outperformed the second group of methods using bkRNA18 but was
175 outperformed by the second group of methods using scRNA663.

176 The evaluation results by the medians of Spearman Correlation Coefficients (SCCs)
177 (**Materials and Methods**) and the cluster analysis results were also consistent with the
178 evaluation results by AUCVC. Generally speaking, a normalization method with a higher
179 AUCVC value produced a lower median of Spearman Correlation Coefficients (SCCs)

180 between normalized expression profiles of ubiquitous gene pairs using both scRNA-seq and
 181 bulk RNA-seq data. The hierarchical clustering result showed that 13 methods had been
 182 classified into the same groups (**Figure 3CF**) by SCCs between normalization factor pairs
 183 as those (**Figure 2AB**) by AUCVC. By our new designed metric AUCVC, TU, NCS and
 184 ES were evaluated as the best normalization methods using both scRNA-seq and bulk
 185 RNA-seq data, which enhanced the discovery using only bulk RNA-seq data in the previous
 186 study [8]. Since the non-zero ratio 0.2 allowed the maximum number of uniform genes for
 187 calculation, we presented this snapshot of evaluation results to show the consistency of the
 188 evaluation results using both scRNA663 (**Figure 3ABC**) and bkRNA18 (**Figure 3DEF**).



190 **Figure 3.**

189

190

191 **Implementation and availability**

192 The gene expression data (scRNA663 and bkRNA18), normalization methods and
 193 evaluation metrics (AUCVC and SCCs) have been included in the R package
 194 NormExpression. All the functions except the NCS and ES methods have been
 195 implemented in R programs [2] for their running on R platforms of any version. The NCS
 196 and ES methods had been implemented in Perl programs on the Linux system by Gustavo *et*
 197 *al.* [8] but they need be installed with many Perl modules. We have modified them into a
 198 stand-alone program (**Supplementary file 2**).

199 A quick evaluation is usually started with 10 normalization methods, which are HG7,
200 ERCC (if available), TC, CR, NR, DESeq, RLE, UQ, TMM and TU. The quick evaluation
201 produce AUCVC values of 10 methods and the raw gene expression matrix for users to
202 evaluate and select methods. NCS and ES are not included in 10 methods, since they have a
203 similar performance of TU but are much more time consuming. The non-zero ratio and
204 presence rate can be set to 1 to calculate AUCVC for bulk RNA-seq data, while they need
205 be set to appropriate values (default 0.2) for scRNA-seq data to avoid parameter grid. Based
206 on our experiences, **both non-zero ratio and presence rate need be set to the values to ensure**
207 **that both the product of the sample number multiplying non-zero ratio and that of the**
208 **sample number multiplying presence rate are larger than 100 for scRNA-seq data.**

209

210 **Materials and Methods**

211 **Datasets**

212 In the previous study by Lin Liu *et al.* (SRA: SRP113436), 663 single-cell samples and
213 18 bulk samples had been sequenced using the Smart-seq2 scRNA-seq protocol. In this
214 study, we built a scRNA-seq dataset including 653 single cells from colon tumor tissues and
215 10 single cells from distal tissues (>10 cm) as control. We also built a bulk RNA-seq dataset
216 including nine samples from colon tumor tissues and nine samples from distal tissues.
217 Samples with total read number less than 288,289 were removed in the data filtering step.
218 The cleaning and quality control of both scRNA-seq and bulk RNA-seq data were
219 performed using the pipeline Fastq_clean [15] that was optimized to clean the raw reads
220 from Illumina platforms. Using the software STAR [11] v2.5.2b, we aligned all the cleaned
221 scRNA-seq and bulk RNA-seq reads to the human genome GRCh38/hg38 and the
222 expression levels of 57,992 annotated genes (57,955 nuclear genes and 37 mitochondrial
223 genes) were quantified. Non-polyA RNAs were not discarded to test the robustness of
224 normalization methods, although the Smart-seq2 protocol theoretically had only captured
225 polyA RNAs. In addition, the expression levels of 92 ERCC RNA and the long non-coding
226 RNA (lncRNA) MDL1 in human mitochondrial [12] were also quantified. ERCC RNA had
227 been spiked into 208 single-cell samples before library construction, the expression levels

228 of ERCC RNA in other 455 single-cell samples and 18 bulk samples were simulated by
229 linear regression. Finally, two datasets were named scRNA663 (58085×663) and
230 bkRNA18 (58085×18) and included into the R package NormExpression.

231

232 **Normalization methods**

233 All 13 methods in the package NormExpression are HG7, ERCC, TN, TC, CR, NR,
234 DESeq, RLE, UQ, TMM, TU, NCS and ES. HG7, ERCC, TN, TC, CR, NR and TU are
235 based on a set of pre-selected genes and each of these methods uses the gene expression
236 level summed over these pre-selected genes in a sample as the library size (**Figure 1B**) to
237 calculate the normalization factor. HG7 includes seven genes (UBC, HMBS, TBP, GAPDH,
238 HPRT1, RPL13A and ACTB), which had been used to achieve the best evaluation result
239 among those using all possible combinations of tested housekeeping genes in the previous
240 study by Gustavo *et al.* [8]. ERCC is a set of commonly used spike-in RNA consisting of 92
241 polyadenylated transcripts with short 3' polyA tails but without 5' caps [5]. The pre-selected
242 genes used by HG7, ERCC, and TU are seven housekeeping genes, 92 ERCC RNA and
243 ubiquitous genes (described below), respectively. NR only counts reads which have been
244 aligned to nuclear genomes, while CR counts reads which have been aligned to both nuclear
245 and mitochondrial genomes. The library size estimated by TC is equal to that estimated by
246 CR plus that estimated by ERCC. TN uses the number of all reads which can be aligned to
247 ERCC RNA, nuclear and mitochondrial genomes.

248 The DESeq method was obtained from the Bioconductor package DESeq [6] and
249 modified to process scRNA-seq data. RLE, UQ and TMM were obtained from the
250 Bioconductor package edgeR [7] and modified to process scRNA-seq data. TU, NCS and
251 ES were obtained from the previous study by Gustavo *et al.* [8]. Since TU sums counts of
252 all ubiquitous genes as the library size to calculate the normalization factor, a process to
253 select ubiquitous genes (describe below) has been integrated into the TU method. TU
254 maximizes AUCVC instead of the number of resulting uniform genes to select ubiquitous
255 genes in the R package NormExpression.

256

257 **Uniform genes and ubiquitous genes**

258 A uniform gene was defined if the Coefficient of Variation (CV, **Formula 1**) of its
259 post-normalization expression levels across all samples was not more than a cutoff.
260 Ubiquitous genes were defined as the intersection of a trimmed sets of all samples [8]. This
261 trimmed set of genes were selected for each sample by 1) excluding genes with zero values,
262 2) sorting the non-zero genes by expression level in that sample, and 3) removing the upper
263 and lower ends of the sample-specific expression distribution. Gustavo *et al.* determined the
264 upper and lower cutoffs by testing all possible combinations of lower and upper cutoffs at
265 5% resolution to maximize the number of resulting uniform genes using one bulk RNA-seq
266 dataset [8]. The size of a scRNA-seq dataset is usually very large, which could result in a
267 very small or even empty set of ubiquitous genes, since the number of ubiquitous genes
268 depends on the sizes of datasets. To select ubiquitous genes using scRNA-seq data, we
269 defined a parameter named presence rate, which required that one selected ubiquitous gene
270 must appear in at least a proportion of the trimmed sets.

271

272 **Evaluation metrics**

273 In the previous study [8], Gustavo *et al.* designed two novel and mutually independent
274 metrics, which were the number of uniform genes and Spearman Correlation Coefficients
275 (SCCs) between expression profiles of gene pairs. Two basic ideas to support these two
276 evaluation metrics are successful normalization methods increase the number of uniform
277 genes and decrease the correlation between the expression profiles of gene pairs. In this
278 study, we designed a new metric AUCVC instead of the number of uniform genes to
279 evaluate normalization methods. We randomly selected 1,000,000 ubiquitous gene pairs to
280 calculate the medians of SCCs. Then, we compared the evaluation results of all the
281 normalization methods by the medians of SCCs with those by AUCVC.

282 AUCVC is created by plotting the number of uniform genes (y-axis) at each
283 normalized CV (**Formula 2**) threshold (x-axis). To determine the number of uniform genes
284 using scRNA-seq data containing a high frequency of zeros, **we only considered genes with**
285 **non-zero expression values divided by the sample number not less than a threshold, which**
286 **was designed as a parameter non-zero ratio.** Since a high or a low normalized CV threshold
287 produces more false or less uniform genes, it is reasonable to consider the overall

288 performance of each method at various threshold settings instead of that at one specific
289 threshold setting. In formula 1 and 2, symbols have the same meanings as those in figure 1
290 and n^* does not count zero elements in each sample.

291

292

$$293 \quad CV_i = \left(\sqrt{\frac{1}{n} \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2} \right) / \bar{x}_i, \quad \bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ij} \quad (1)$$

$$294 \quad \text{Normalized } CV_i = \left\{ CV_i - \min_i(CV_i) \right\} / \left\{ \max_i(CV_i) - \min_i(CV_i) \right\}$$
$$CV_i = \left(\sqrt{\frac{1}{n^*} \sum_{j=1}^{n^*} (\log_2(x_{ij}) - \bar{x}_i)^2} \right) / \bar{x}_i, \quad \bar{x}_i = \frac{1}{n^*} \sum_{j=1}^{n^*} \log_2(x_{ij}), \quad x_{ij} > 0 \quad (2)$$

295 **Conclusion and Discussion**

296

297 **Acknowledgments**

298 We thank Professor Sui Huang and Dr. Joseph Zhou from Institute for Systems
299 Biology (ISB) their hosts on Shan Gao' visiting to ISB.

300

301 **Funding**

302 This work was supported by grants from by National Key Research and
303 Development Program of China (2016YFC0502304-03) to Defu Chen and Fundamental
304 Research Funds for the Central Universities (for Nankai University) to Shan Gao.

305

306 **Competing interests**

307 Non-financial competing interests are claimed in this study.

308

309

310 Authors' contributions

311 SG conceived this project. SG and JR supervised this project. ZW, WL and SL
312 performed programming. ZW, DY and HJ analyzed the data. WD prepared all the figures,
313 tables and additional files. SG drafted the main manuscript. XX and XX revised the
314 manuscript. **All authors have read and approved the manuscript.**

315
316
317
318

319 REFERENCES

320

- 321 1. Lovén, J., Orlando, D.A., Sigova, A.A., Lin, C.Y., Rahl, P.B., Burge, C.B., Levens,
322 D.L., Lee, T.I., and Young, R.A., *Revisiting global gene expression analysis*. Cell,
323 2012. **151**(3): 476-482.
- 324 2. Gao, S., Ou, J., and Xiao, K., *R language and Bioconductor in bioinformatics*
325 *applications(Chinese Edition)*. 2014, Tianjin: Tianjin Science and Technology
326 Translation Publishing Ltd.
- 327 3. Tao, H. and Gao, S., *Computational Systems Biology: Methods and Protocols*. 2018,
328 New York: Springer.
- 329 4. Bullard, J.H., Purdom, E., Hansen, K.D., and Dudoit, S., *Evaluation of statistical*
330 *methods for normalization and differential expression in mRNA-Seq experiments*.
331 BMC Bioinformatics, 2010. **11**(1): 94.
- 332 5. Jiang, L., Schlesinger, F., Davis, C.A., Zhang, Y., Li, R., Salit, M., Gingeras, T.R.,
333 and Oliver, B., *Synthetic spike-in standards for RNA-seq experiments*. Genome
334 Research, 2011. **21**(9): 1543.
- 335 6. Anders, S. and Huber, W., *Differential expression analysis for sequence count data*.
336 Genome biology, 2010. **11**(10): R106.
- 337 7. Robinson, M.D., McCarthy, D.J., and Smyth, G.K., *edgeR: a Bioconductor package*
338 *for differential expression analysis of digital gene expression data*. Bioinformatics,
339 2010. **26**(1): 139-140.
- 340 8. Glusman, G., Caballero, J., Robinson, M., Kutlu, B., and Hood, L., *Optimal scaling*
341 *of digital transcriptomes*. Plos One, 2013. **8**(11): e77885.
- 342 9. Li, P., Piao, Y., Shon, H.S., and Ryu, K.H., *Comparing the normalization methods*
343 *for the differential analysis of Illumina high-throughput RNA-Seq data*. BMC
344 Bioinformatics, 2015. **16**(1): 347.
- 345 10. Gao, S., Zhang, N., Duan, G.Y., Yang, Z., Ruan, J.S., and Zhang, T., *Prediction of*
346 *function changes associated with single-point protein mutations using support*
347 *vector machines (SVMs)*. Human Mutation, 2009. **30**(8): 1161-6.
- 348 11. Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P.,
349 Chaisson, M., and Gingeras, T.R., *STAR: ultrafast universal RNA-seq aligner*.
350 Bioinformatics, 2013. **29**(1): 15-21.

351 12. Gao, S., Tian, X., Chang, H., Sun, Y., Wu, Z., Cheng, Z., Dong, P., Zhao, Q., Ruan,
352 J., and Bu, W., *Two novel lncRNAs discovered in human mitochondrial DNA using*
353 *PacBio full-length transcriptome data*. Mitochondrion, 2017.
354

355 **Figure legends**

356 **Figure 1. Normalization factor and library size**

357 (A). A raw gene expression matrix can be transformed into a normalized gene expression
358 matrix by the multiplication of a factor f_j to each column. Each column represents the
359 expression levels of all genes from a sample and each row represents the expression levels
360 of a gene across all samples. (B). HG7, ERCC, TN, TC, CR, NR and TU use library sizes
361 N_j^* to calculate normalization factors. N_j represents the library size estimated by TC.
362 DESeq, RLE, UQ and TMM have been modified in NormExpression to ignore zero values
363 and the resulting normalization factors need be further normalized by their geometric mean
364 values. After modification, RLE is identical to DESeq. Q3 means that about 75% of genes
365 in the j th sample have expression levels below Q3 and about 25% have those above Q3. For
366 all methods, log represents the natural logarithm.

367

368

369 **Figure 2. Parameter grid to evaluate normalization methods**

370 Parameter grid of non-zero ratio from 0.2 to 0.9 for scRNA663 and from 0.8 to 1 for
371 bkRNA18 was used to produce AUCVC values All the normalization methods were
372 classified into three groups based on their AUCVC values sorted in descending order using
373 one scRNA-seq dataset scRNA663 (A) and one bulk RNA-seq dataset bkRNA18 (B).

374

375 **Figure 3. Consistency in the evaluation results by different metrics**

376 A normalization method with a higher AUCVC value produced a lower median of
377 Spearman Correlation Coefficients (SCCs) between normalized expression profiles of
378 ubiquitous gene pairs using both scRNA-seq (AB) and bulk RNA-seq data (DE). The
379 hierarchical clustering result showed that 13 methods had been classified into the same
380 groups (CF) by SCCs between normalization factor pairs as those (Figure 2AB) by
381 AUCVC.

382