# A comparative study of deconvolution methods for RNA-seq data under a dynamic testing landscape

Haijing Jin[1] and Zhandong Liu*[2,3]

1.  Graduate Program in Quantitative and Computational Biosciences, Baylor College of Medicine, United States

2.  Jan and Dan Duncan Neurological Research Institute at Texas Children's Hospital, Houston, United States

3.  Department of Pediatrics, Baylor College of Medicine, Houston, United States

## Abstract

Deconvolution analyses have been widely used to track compositional alternations of cell-types in gene expression data. Even though numerous novel methods have been developed in recent years, researchers are still having difficulty selecting optimal deconvolution methods due to the lack of comprehensive benchmarks relative to the newly developed methods. To systematically reveal the pitfalls and challenges of deconvolution analyses, we studied the impact of several technical and biological factors such as simulation model, quantification unit, component number, weight matrix, and unknown content by constructing three benchmarking frameworks that cover comparative analysis of 11 popular deconvolution methods under 1,766 conditions. We hope this study can provide new

20    insights to researchers for future application, standardization, and development of

21    deconvolution tools on RNA-seq data.

## Background

23        Deconvolution refers to a process that separates a heterogeneous mixture signal into

24    its constituent components. In the biomedical field, researchers have been using

25    deconvolution methods to derive cell-type-specific signals[1–3] from heterogeneous mixture

26    data. Cellular composition information is crucial for developing sophisticated diagnostic

27    techniques as it enables researchers to track each cellular component's contribution during

28    disease progressions[4]. Although some experimental approaches like fluorescence-activated

29    cell sorting(FACS), immunohistochemistry(IHC), and single-cell RNA-seq can derive cell-

30    type proportion data[3], all these approaches are either restricted by its throughput or remain

31    too costly and laborious for large-scale clinical applications. By far, deconvolution is

32    recognized as the most cost-effect approach to derive cell-type proportion data from

33    heterogenous biospecimens and has the potential to bring a considerable improvement in

34    the speed and scale of cell-type-specific clinical diagnosis.

35        By January 2018, there have been around 50 deconvolution methods developed[2] and

36    researchers are now facing the challenge of selecting the right method for deconvolution

37    analysis. In a methodological paper, authors usually compared the method of their own to

38    a chosen set of published methods and arrived at the conclusion that their method was the

39    best. However, only a limited number of deconvolution methods and biological conditions

40    were considered in these comparisons. Moreover, different research groups applied

41    inconsistent testing frameworks with different simulation strategies, evaluation metrics, and

42    cell-type annotations, making it difficult for researchers to determine the optimal method

43    for the deconvolution analysis.  For a fair and comprehensive comparison of deconvolution

44    applications in complex biological systems, an independent benchmarking is in need[5].

45    Previously, Sturm *et al.*[3] and Cobos *et al.*[6] performed quantitative evaluations of reference-

46    based and marker-based deconvolution methods on RNA-seq data. Sturm *et al.*[3] focused on

47    spill-over effects, minimal detection fraction, and background predictions and suggested

48    removing non-specific signature genes to improve deconvolution accuracy. Cobos *et*

49    *al.*[6] focused on the impact of different normalization strategies, reference platforms, marker

50    gene selection strategies, and missing cellular components in the reference. Compared with

51    previous benchmarks, our study focuses on technical and biological factors caused by varied

52    experimental mixture conditions such as mixture noise levels, quantification unit selection,

53    cellular component number, weight matrix property, and unknown cellular contents.  We

54    also studied the major factors that determine an evaluation framework, such as simulation

55    model selection, evaluation metric selection, and measurement scale selection. Our work

56    carefully examined the joint impact of different technical parameters and biological design

57    factors to provide an insightful reference guide for mixture condition determination and

58    deconvolution method selection.

59        There are three types of benchmarking frameworks for the evaluation of

60    deconvolution methods: *in silico* framework[7,8], *in vitro* framework[9], and *in vivo* framework[10]

61    (Supplementary Table 1). The *in vivo* testing framework mainly rely on indirect performance

62    assessment and usually cannot derive a definite conclusion of the method's performance.

63    Only a few in vivo benchmarking datasets[3] have coupled FACS results. Nevertheless, these

3

64    benchmarking datasets are often restricted by limited cell types and sample numbers[3,8] . The

65    *in vitro* testing framework where mixtures are generated in the tube with predefined mixing

66    compositions also suffers from limited cell types and sample numbers. Moreover, most *in*

67    *vitro* testing frameworks applied 'orthogonal' weights, leading to over-optimistic

68    performance assessment. The *in silico* testing framework uses RNA-seq profiles from

69    purified biological samples as primary building blocks and generates heterogeneous mixing

70    samples by *in silico* mixing procedures. Among all three benchmarking frameworks, we

71    selected the *in silico* testing framework to systematically explore the impact of different

72    biological and technical factors, which require large amounts of benchmarking datasets

73    under controlled and finely tuned multi-factor testing environments.

74        To provide a reliable reference for the application and development of deconvolution

75    methods, we compared 11 deconvolution methods (Figure 1b and Supplementary Table 3).

76    To establish benchmarking frameworks that mimic application scenarios of more

77    complicated and diverse biological systems, we designed three sets of benchmarking

78    frameworks that mimic up to 1,766 biological conditions with varying noise levels, library

79    sizes, cellular component numbers, weight matrix properties, simulation models, and

80    proportions of unknown contents (Figure 1a, Supplementary Table 2). To determine the

81    impact of evaluation frameworks, we performed comparisons under different simulation

82    models and measurement scales with two sets of evaluation metrics: correlation (Pearson's

83    Correlation Coefficient) and mAD (Mean Absolute Deviation)(Methods). Compared with

84    previous benchmarks, we applied more flexible and sophisticated simulation strategies to

85    create mixtures covering dynamic conditions, which enable us to investigate the tipping

86    point where each method deteriorates. Moreover, we studied the impact of commonly
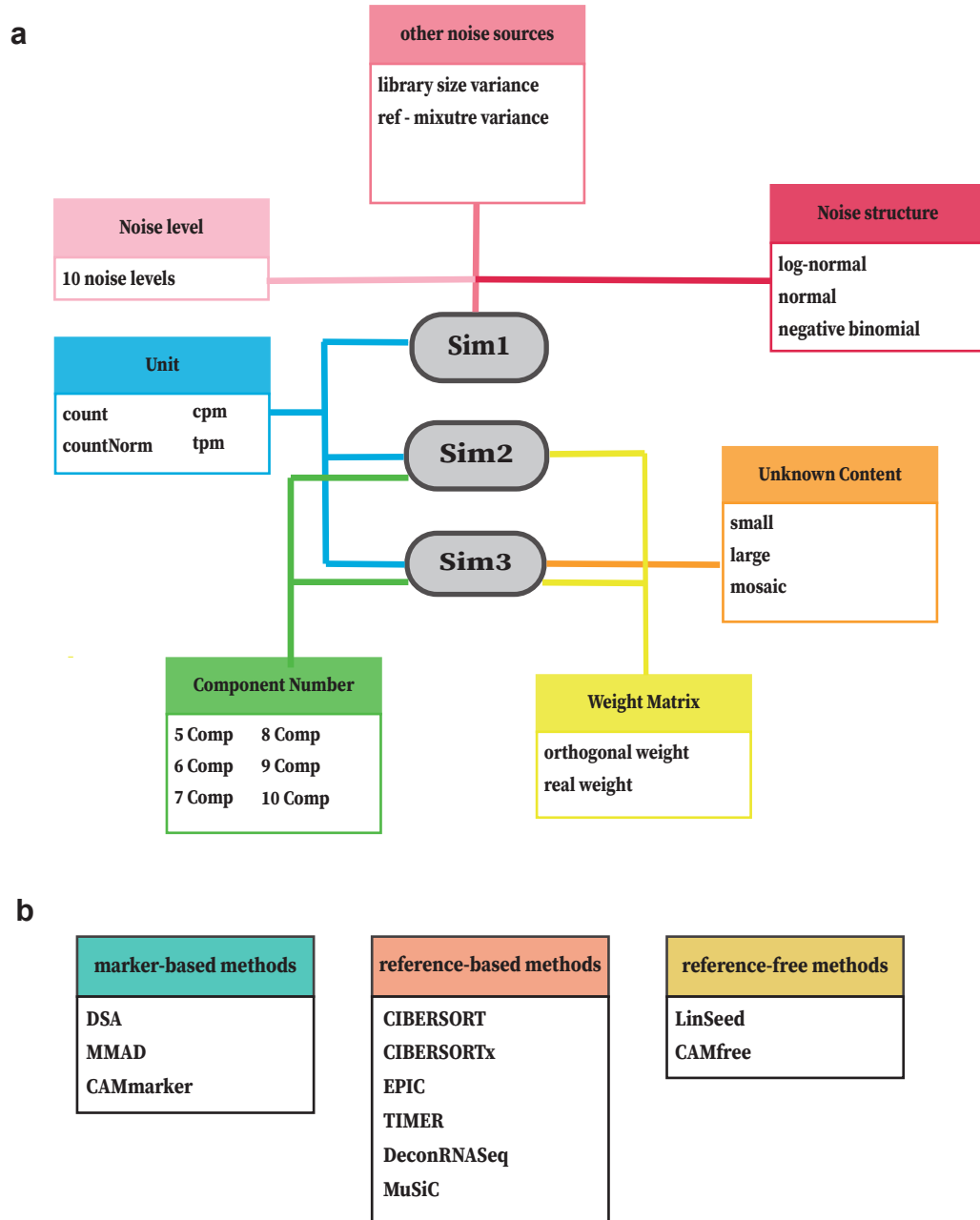
4

87   applied simulation strategies, and by comparison to the real mixture data, we derived

88   improved simulation strategies that can generate more complex and yet authentic

89   simulation data. Our results provide a dynamic testing landscape that allows the user to

90   select the right method that performs well in the targeted experimental condition.

## Results

**Using simulation to generate diverse deconvolution testing environments**

93   We designed three benchmarking frameworks to test the performance of

94   deconvolution methods under multiple application scenarios. Each framework was designed

95   to study the impact of specific technical and biological factors on deconvolution analysis

96   (Figure 1a). The first benchmarking framework (Sim1) was designed to reveal the impact of

97   the noise structure under diverse noise levels. The second benchmarking framework (Sim2)

98   was designed to reveal the impact of cellular component numbers and weight matrix

99   properties. The third benchmarking framework (Sim3) was designed to reveal the impact of

100  unknown biological contents and measurement scales.

101  In an *in silico* benchmarking framework, a deconvolution testing environment

102  consists of mixture data, reference data, ground truths, and testing methods. Mixture data

103  refers to heterogeneous gene expression profiles for deconvolution. Reference data refers to

104  homogeneous cell-type-specific data, which is used to guide the deconvolution process.

**a**



**b**

**Fig.1| Overview of *in silico* testing frameworks and methods categorization**

**a,** Three benchmarking frameworks were constructed to investigate the impact of seven factors that affect deconvolution analysis: noise level, noise structure, other noise sources, quantification unit, unknown content, component number, and weight matrix. **b,** 11 deconvolution methods are tested and have been categorized based on the required reference input: marker-based, reference-based, and reference-free.

111    Ground truths refer to the real mixing proportions of constituent cell types in the mixture

112    data. The accuracy of deconvolution methods can be assessed by comparing estimated

113    proportions to the ground truths. Reference data can vary based on the required input of the

114    tested deconvolution method. In this study, we classified eleven deconvolution methods

115    according to the required reference data in the following categories: marker-based,

116    reference-based, and reference-free (Figure 1b, Supplementary Table 3). Marker-based

117    methods such as DSA[11], MMAD[12], and CAMmarker[13] use marker gene lists to guide the

118    deconvolution analysis. Reference-based methods such as CIBERSORT[7], CIBERSORTx[8],

119    EPIC[14], TIMER[10], DeconRNASeq[15], and MuSiC[16] use cell-type-specific gene expression

120    profiles. Except for MuSiC[16], nearly all reference-based methods require signature gene lists

121    as an additional input. MuSiC[16] implements weighted non-negative least squares regression

122    (W-NNLS) and does not require any pre-determined gene sets. Finally, reference-free

123    methods such as LinSeed[17] and CAMfree[13] do not require any external references.
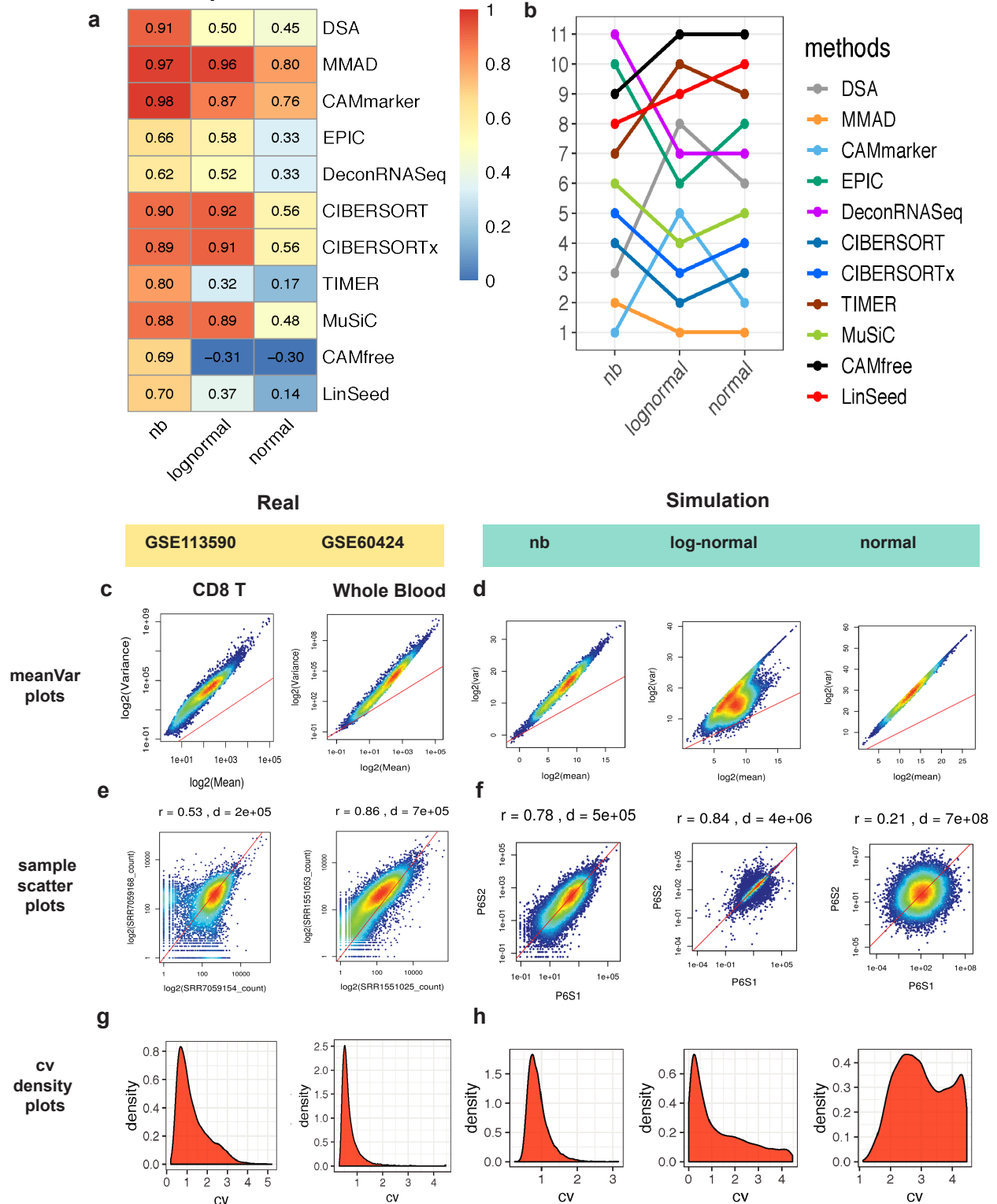
124    **Selection of simulation model affects the deconvolution evaluation**

125    The benchmarking framework Sim1_simModel is designed to learn the impact of

126    noise structure under different noise levels (Fig. 1a, Methods). To understand the impact of

127    noise structure, we simulated noise based on three simulation models: normal, log-normal,

128    and negative binomial (nb). All these simulation models have been applied in previous

129    publications[7,15,17–19] to generate *in silico* mixing expression profiles. For each simulation

130    model, we generated ten levels of noise to evaluate the robustness of deconvolution methods

131    to the magnitude of noise(Supplementary Fig. 1a). To ensure the generality of our conclusion

132    across different datasets and account for reference-mixture variance, we performed

133    repeated mixture simulation with three independent blood datasets and created nine testing

134   environments with different mixture-reference pairs (Methods, Supplementary Table 2 and

135   Supplementary Table 4).

136       For the noise level, consistent with previous findings, we observed that the accuracies

137   of the deconvolution methods decreased as the noise level increased, which was exhibited

138   as decreasing correlation (Supplementary Fig. 3) and increasing mAD (Supplementary Fig.

139   4) values. We also noticed that the impact of the RNA-seq quantification unit is trivial

140   (Supplementary Fig. 3 and 4) and thus selected the most commonly used unit tpm for

141   remaining illustrations of testing results in Sim1_simModel. Unless specifically indicated (as

142   in Sim1_libSize), all results in this study are from mixture data with the tpm unit.

143       To reveal the impact of the simulation models, we averaged evaluation metrics across

144   noise levels and generated summarized evaluation heatmaps ($11 \times 3$) where row index

145   number 11 indicates the number of methods and column index number 3 indicates the

146   number of simulation models. Based on the summarized evaluation heatmaps of correlation

147   (Fig. 2a) and mAD (Supplementary Fig. 5a), we observed that the selection of the simulation

148   model strongly affected evaluation results. For instance, methods like DSA[11], TIMER[10], and

149   CAMfree's[13] rankings were all relatively higher in the negative binomial group in both

150   correlation (Fig. 2b) and mAD (Supplementary Fig. 5b) metrics when comparing with

151   evaluations from normal and log-normal groups. The above phenomenon indicated that the

152   performance of some deconvolution methods is underestimated due to the underlying

153   simulation model.

**Fig.2| Evaluation results of Sim1_simModel and noise structure comparisons between real and simulated data**

157 **a,** Heatmap of summarized evaluation results based on the Pearson's correlation coefficients and **b,** rankings

158 of tested deconvolution methods in the Sim1_simModel. In each heatmap, row indexes refer to the tested

159 methods and column indexes refer to the simulation models (negative binomial, log-normal, and normal). **c,d,**

160 Mean-variance plots of (**c**) real and (**d**) simulated data. (r: Spearman's correlation coefficient, d: Euclidean

161 distance) **e,f,** sample-sample scatter plots of (**e**) real and (**f**) simulated data. **g,h,** Density plots of CV (Coefficient

162 of variation) of (**g**) real and (**d**) simulated data. (Real data are derived from GSE113590 and GSE60424 and

163 Supplementary Figure 6 and 7 contain detailed variance analysis results for each dataset) (All simulated data

164 in Figure 2 are based on simulations derived from GSE51984 with the P6 noise level.) (Results in **a** and **b** are

165 in tpm unit, results in **c-f** are in count unit)

### The negative binomial model recapitulates noise structures of real data

167 In the Sim1_simModel, we found that the noise structure is the main factor obscuring

168 deconvolution performance assessment (Fig. 2a and b, Supplementary Fig. 5). To identify the

169 simulation model that best recapitulates the essential characteristics of real data, we

170 performed noise structure comparisons between real and simulated data by mean-variance

171 plots, sample-sample scatter plots and coefficient of variance (CV) density plots.

172 We used the mean-variance plots to study the overall trend of variance along with the

173 gene expression level in both real and simulated data (P6 noise level) (Fig. 2c and d). As

174 expected, we observed that the variance and mean value of counts follow a linear trend in

175 the log space with a clear overdispersion phenomenon, which is typical to the RNA-seq

176 data[20](Fig. 2c). However, in the simulation group, only the simulations generated from the

177 negative binomial and normal models showed a similar mean-variance trend to the trend

178 observed in the real data (Fig. 2d).

179 Next, we used sample-sample scatter plots to study the concordance trend of gene

180 expression profiles(Fig. 2e and f). In real data, we observed that lowly expressed genes

10

181  exhibited larger relative deviances to the diagonal reference line (y = x) than highly

182  expressed genes (Fig. 2e). This phenomenon indicates larger uncertainties in quantifying

183  RNA molecules with lower abundance. In the simulation group, only simulation data from

184  the negative binomial model recapitulated higher deviances of lowly expressed genes (Fig.

185  2f).

186  We also compared the magnitude of noise between the real and simulated data. In the

187  real data, the sample-sample Spearman's correlation values range from 0.53 to 0.99 while

188  the sample-sample Euclidean distances fluctuate around the order of $10^4 \sim 10^5$

189  (Supplementary Fig.6 a and b and Supplementary Fig. 7 a and b). In three tested simulation

190  models, only the negative binomial model was capable of generating simulated profiles with

191  comparable sample-sample correlation (0.57 – 0.98) and Euclidean distance (around the

192  order of$10^4 \sim 10^5$) to the real datasets (Supplementary Figure 8) while maintaining mean-

193  variance trend with overdispersion phenomenon (Supplementary Fig. 9).

194  We compared the density curve of CV (coefficient variation) values in real and

195  simulated data (Fig. 2g and h). Real data exhibited a unimodal bell-shaped curve, indicating

196  that most of the genes had low to moderate levels of CV (Fig. 2g). In the simulation group,

197  only simulations derived from the negative binomial model maintained the unimodal bell-

198  shaped curve throughout all noise levels (Fig. 2h). CV density distributions of normal and

199  log-normal simulation models showed density curves that were skewed towards the high CV

200  value from noise level P6 to P10, which indicating unauthentic noise

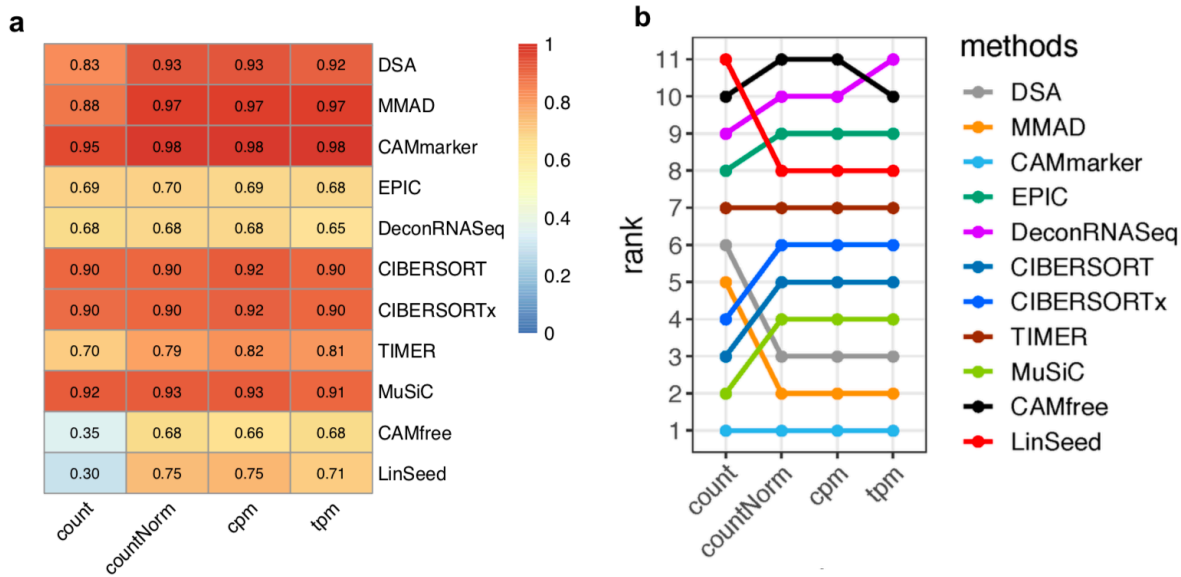201  structure(Supplementary Fig. 10b).

202    In conclusion, the negative binomial simulation model, which successfully

203    recapitulates the mean-variance trend, sample-sample concordance, the density of CV,

204    presents the most similar noise structure to the real data. The negative binomial model also

205    kept the magnitude of noise at comparable levels to the real data and thus should be

206    considered as the most appropriate simulation model for generating *in silico* mixtures for

207    deconvolution benchmarking.

208    **Library size normalization is required to ensure the deconvolution accuracy**

209    In this benchmarking framework, we focused on the impact of RNA-seq quantification

210    units with mixtures that varied in their library sizes (Supplementary Fig. 1b). To reveal bias

211    caused by varied library sizes, we designed Sim1_libSize in which every mixture comprised

212    of samples with varied library sizes (first 10 samples with 12M reads, and remaining 10

213    samples with 24M reads), and our results indicate using quantification units normalized by

214    library sizes can mitigate the bias caused by library size variation (Fig. 3a, Supplementary

215    Fig. 11a). We summarized evaluation results across all 10 noise levels and generated

216    evaluation heatmaps with dimensions 11 by 4 where 11 indicates the number of methods

217    and 4 indicates the number of quantification units being tested.

218    We observed that three methods, CIBERSORT[7], CIBERSORTx[8], and MuSiC[16], which

219    implemented normalization procedures, showed decent performance ($r \geq 0.9, mAD \leq 0.1$)

220    regardless of the selected quantification unit (Fig. 3a, Supplementary Fig. 11a). Six methods

221    (DSA[11], MMAD[12], CAMmarker[13], TIMER[10], CAMfree[13], and LinSeed[17]) showed improved

222    accuracy after library size normalization (Fig. 3a, Supplementary Fig. 11a).

223

12

224

**Fig.3| Evaluation results of Sim1_libSize**

**a,** Heatmap of summarized evaluation results based on the Pearson's correlation coefficients and **b,** rankings of tested deconvolution methods. In each heatmap, row indexes refer to the tested methods and column indexes refer to the quantification units (count, countNorm, cpm, and tpm).

229

230    Contradicting to the Sim1_simModel (Supplementary Fig.3 and 4), we observed that

231    the choice of quantification unit had a high impact on Sim1_libSize, which was reflected by

232    discrepant rankings of tested methods (Supplementary Fig. 3b and 11b). As the only

233    difference between the two benchmarking frameworks was the library size, we deduced that

234    the inconsistent performance over different quantification units was due to the library size

235    variation in the mixture dataset.   We thus suggest researchers applying RNA-seq

236    quantification units that are normalized by library sizes to mitigate the bias caused by varied

237    library sizes unless indicated by the author of the method(MuSiC[16]) to use the count unit.

**Impact of cellular component number and weight matrix on deconvolution analysis**

To investigate the joint impact of the cellular component number and weight matrix property, we designed the benchmarking framework Sim2 with six gradients of component number ranging from 5 to 10 and two types of weight matrices: 'orthog' and 'real' (Supplementary Fig. 2a and Supplementary Table 2 and 4). The 'orthog' weight matrix was generated by minimizing the condition number, and the 'real' weight matrix is constructed based on whole blood immune cell proportions in the real biological samples[21](Methods). We discarded the CAMfree[13] method in Sim2 due to the poor scalability of CAMfree[13] on mixtures with large component numbers.

We found that nearly all deconvolution methods achieved higher accuracies with the 'orthog' weight matrices (Fig. 4a) than the 'real' weight matrices, indicating that the mathematical property of the weight matrix has a significant impact on deconvolution analysis. In the mixtures with five components (Comp 5), eight methods (DSA[11], MMAD[12], CAMmarker[13], EPIC[14], CIBERSORT[7], CIBERSORTx[8], MuSiC[16], and LinSeed[17]) exhibited high accuracy levels($r \geq 0.95, mAD \leq 0.05$) in the 'orthog' group (Fig. 4a and Supplementary Fig. 12a) while only three of those eight methods (CIBERSORT[7], CIBERSORTx[8], and MuSiC[16]) in the 'real' group achieved the same level of accuracy (Fig. 4b and Supplementary Fig. 12b).

In addition to the impact of the weight matrix selection, cellular component numbers also affect deconvolution accuracy. In both 'orthog' and 'real' groups, the majority of methods exhibited poorer performance as cellular component number increasing (Fig. 4 a,b and Supplementary Fig. 12). It is also worth noting that none of the tested deconvolution
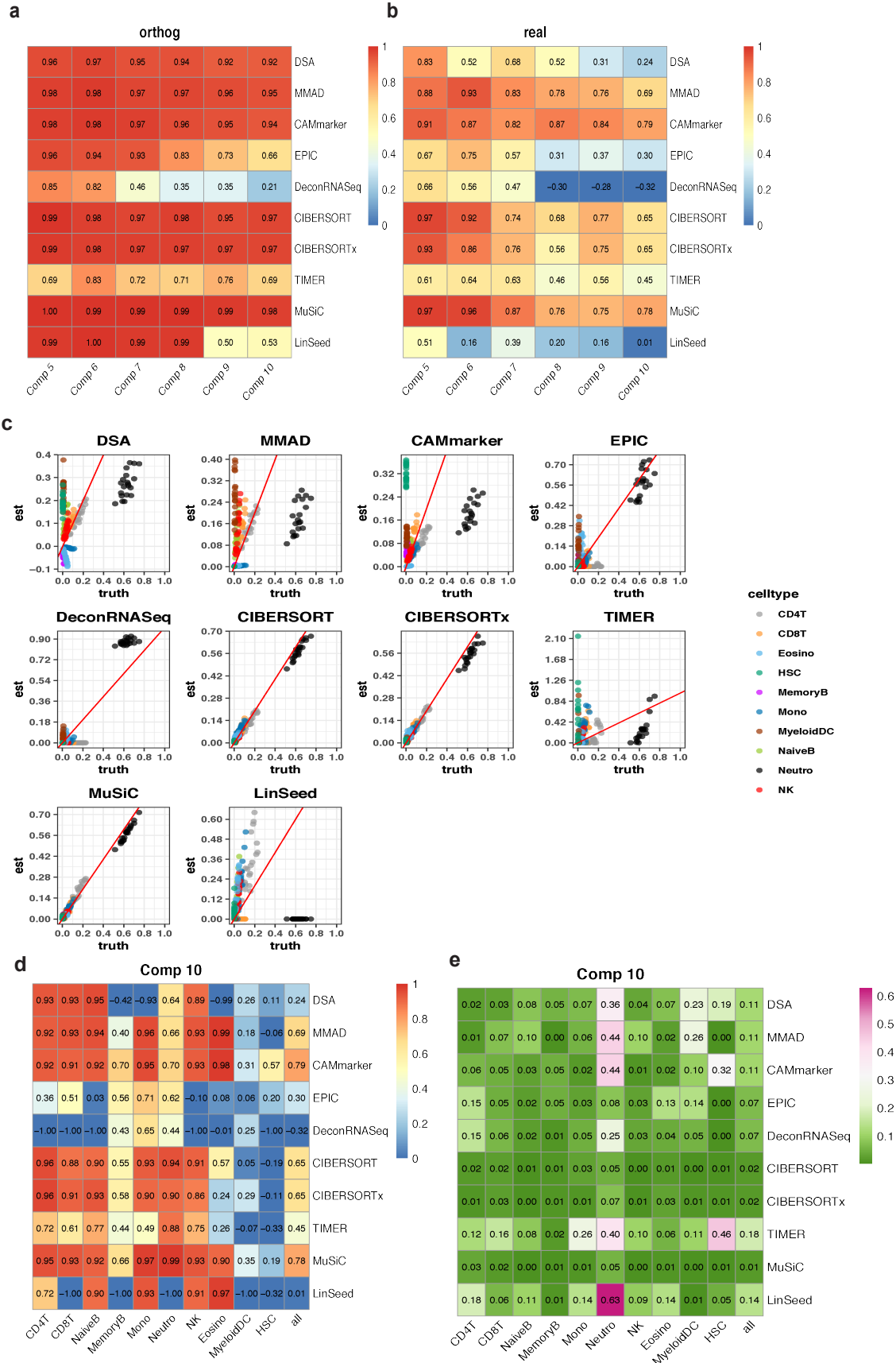
14

259  methods showed a correlation larger than 0.9 with mixtures consist of large cellular

260  component numbers (Comp 7 to Comp 10) in the 'real' group (Fig. 4b).

261      To further investigate the performance of deconvolution methods with large

262  component numbers, we explored the accuracies of mixtures with 10 cellular components

263  and the 'real' weight matrix by drawing scatters plots of estimations and ground truths (data

264  corresponds to the last column of Fig. 4b and Supplementary Fig. 12b). Surprisingly, we

265  found that the correlation evaluation metric, which was considered as the golden standard

266  for the evaluation of deconvolution methods, cannot reflect the deviance of estimations from

267  ground truths (Fig. 4c). However, the deviance of estimation can be reflected by another

268  evaluation metric mAD (Supplementary Fig. 12). For instance, MMAD[12] and CAMmarker[13]

269  performed relatively well on the correlation evaluation metric ($r \geq 0.65$, Fig. 4b),  but both

270  methods had mAD values larger than 0.1, indicating large estimation deviance

271  (Supplementary Fig. 12b). Consistent with the results from scatter plots (Fig. 4c), we found

272  that the best performers were CIBERSORT[7], CIBERSORTx[8], and MuSiC[16]. All three methods

273  achieved high accuracies on both correlation evaluation metric ($r \geq 0.65$) (Supplementary

274  Fig. 4b) and mAD evaluation metric ($mAD \leq 0.02$) (Supplementary Fig. 12b) in the Comp

275  10 mixture with 'real' weight matrix.

276      To understand the impact of each cellular component on deconvolution analysis, we

277  drew  evaluation  heatmaps  with  cell-type-specific  correlation  and  mAD  values

278  (Supplementary Fig. 13, 14). Based on the evaluation heatmap of mixtures with ten cellular

279  components and the 'real' weight matrix, which is the most complicated *in silico* mixture set

280  in the Sim2 benchmark framework, we identified three best performers: CIBERSORT[7],

281     CIBEERSORTx[8], and MuSiC[16] (Fig. 4 d and e). First, we found that all three methods correctly

282     estimated major cellular components ($r \geq 0.85, mAD \leq 0.05$), such as Neutrophils, CD4T,

283     and CD8T in the respective mixtures. Second, while all three methods failed to estimate the

284     linear trend of proportions of rare cell subpopulations that occupies less than 1% in the

285     mixture, such as Myeloid DC and HSC (Hematopoietic Stem Cells) ($r: -0.19 \sim 0.35$ ), they

286     correctly identified them as minor components and did not attribute the percentages of

287     other cell types to these rare cell populations ($mAD: 0 \sim 0.01$). Moreover, because none of

288     the tested deconvolution methods showed good accuracies in both correlation and mAD

289     metrics with Myeloid DC and HSC (Figure 4 d and e), we concluded that none of the currently

290     developed deconvolution methods could not reliably estimate some rare cellular

291     populations that have proportions less than 1%. Finally, we also discovered that marker-

292     gene based methods like DSA[11], MMAD[12], and CAMmarker[13] showed high mAD values

293     (Figure 4d and e), indicating larger deviances in their estimations in the major

294     components($mAD: 0.36 \sim 0.44$)(Fig. 4e).

295          By inspecting cell-type-specific evaluation results of 'real' weight matrices across 6

296     component gradients, we found that introducing rare cellular components MyeloidDC in the

297     Comp 7 mixture caused the deterioration of deconvolution performance, which might be due

298     to the close relationship between MyeloidDC to the monocytes[22]. However, introducing

299     relatively distinct HSC in the Comp 8 mixture further exacerbated the performance

300     deterioration (Supplementary Figures 13 and 14, 'real' group). Therefore, we concluded that

301     the deterioration of deconvolution performance on mixtures with large component number

302     is due to the confounding effect of the highly correlated cellular component and the rare

303     cellular component in the mixture dataset.

305    **Fig.4| Evaluation results of Sim2**

306    **a,b,** Heatmaps of summarized evaluation results based on the Pearson's correlation coefficients with (**a**)

307    'orthog' weight matrix and (**b**) real weight matrix. In each heatmap, row indexes refer to the tested methods

308    and column indexes refer to the cellular component numbers. **c,** Scatter plots of estimated weights vs. ground

309    truths of mixtures with 10 cellular components. **d,e,** Cell-type specific evaluation metrics of mixtures consist of

310    10 cellular components based on (**d**) Pearson's correlation coefficient and (**e**) Mean absolute deviance.

311

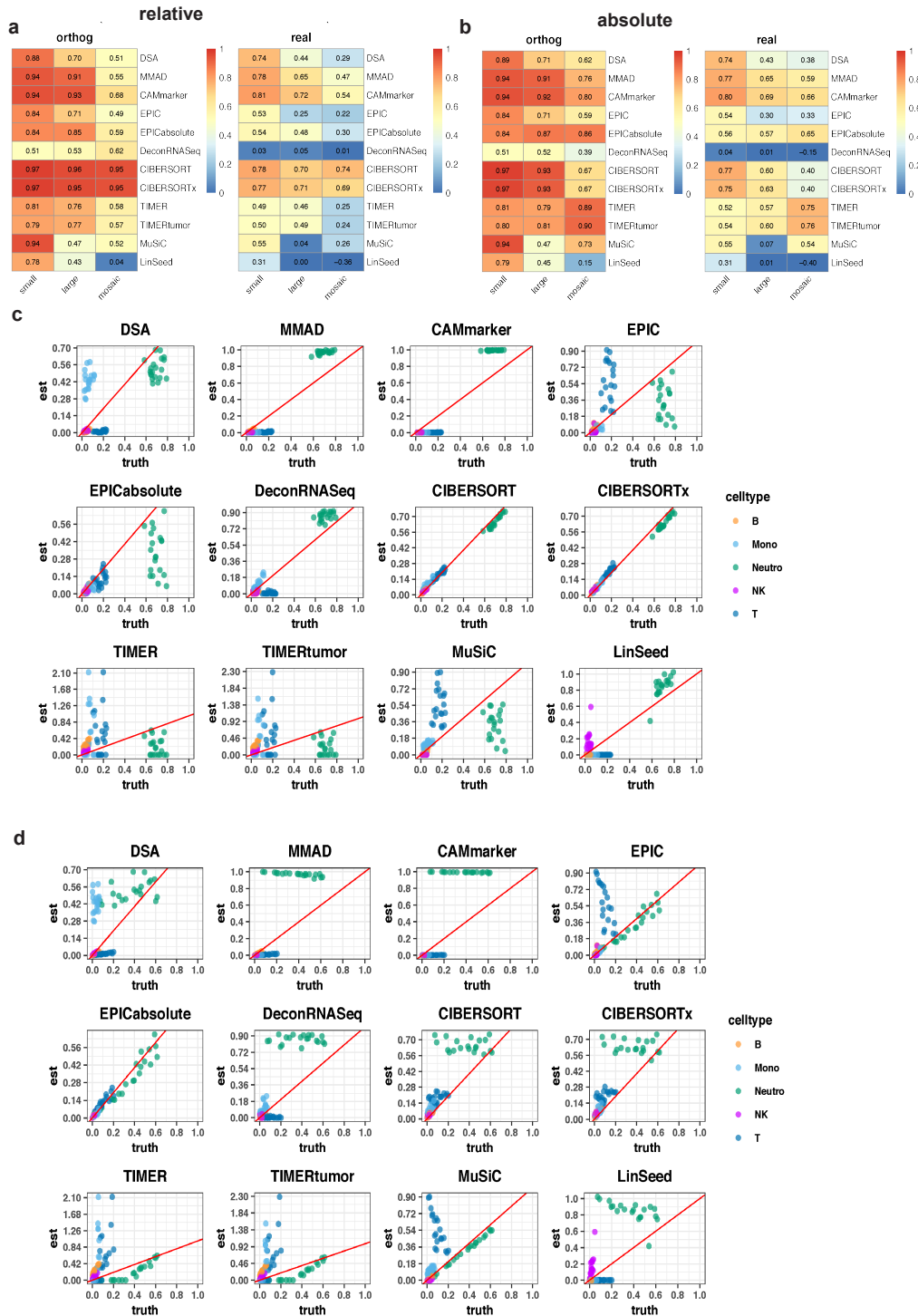312    **Impact of tumor content on deconvolution analysis**

313        Unknown biological content, such as tumor content, is another major factor that

314    influences deconvolution analysis for several reasons.  First, unknown content could be

315    treated as a source of noise unless explicitly modeled by deconvolution methods[7,14]. Second,

316    unknown content is not counted in the estimated cell-type proportions and violates the sum-

317    to-one assumption applied by the majority of deconvolution methods[2,9].

318        To study the impact of unknown biological content on deconvolution analysis, we

319    designed a benchmarking framework that contains mixtures with three sets of tumor spike-

320    ins: the 'small' group refers to mixtures with low levels of tumor spike-ins (0 – 20%), the

321    'large' group refers to mixtures with high levels of tumor spike-ins (70 – 90%), and the

322    'mosaic' group refers to mixtures with more dynamic levels of tumor spike-ins (5% - 95%).

323    Tumor spike-ins were introduced to the 12 mixture sets generated in the Sim2 framework

324    to analyze the joint impact of the component numbers, weight matrix properties, and

325    unknown biological contents (Supplementary Fig. 2b, Methods).  In the performance

326    assessment step, we used two sets of ground truths to derive evaluation results that

18

327    represent different measurement scales (Supplementary Table 5, Methods). The first set of

328    ground truths used the absolute proportions of immune cell types and led to 'absolute'

329    deconvolution accuracy. The second set of ground truths used the relative proportions of

330    immune cells and led to 'relative' deconvolution accuracy. In this set of analyses, we

331    considered additional settings of deconvolution methods that were relevant to the tumor

332    content. Thus, we evaluated eleven methods and two specific method settings

333    TIMERtumor[10] and EPICabsolute[14], which are tailored for deconvolution analysis with

334    unknown tumor contents (Methods, Supplementary Table 3).

335        Our results indicated the weight matrix property as the leading factor that affected

336    deconvolution accuracy because the 'orthog' group presented higher accuracies throughout

337    all deconvolution methods and tumor content conditions (Fig. 5a, b and Supplementary Fig.

338    15). In addition to the weight matrix property, we found that the size of tumor content also

339    affected deconvolution accuracy as we observed deconvolution methods performed better

340    on mixtures with smaller tumor content (Fig. 5a, b and Supplementary Fig. 15). Moreover,

341    we found that all methods showed inconsistent performance with the 'mosaic' mixture group

342    when evaluated on different measurement scales (Fig. 5a, b and Supplementary Fig. 15). For

343    instance, in the 'mosaic' column, CIBERSORT[7] and CIBERSORTx[8] showed higher accuracies

344    ($r: 0.69 \sim 0.95$, $mAD: 0.03$) in the relative measurement scale (Fig. 5a and Supplementary Fig.

345    15a) than in the absolute measurement scale ($r: 0.4 \sim 0.97$, $mAD: 0.06 \sim 0.07$) (Fig. 5b and

346    Supplementary Fig. 15b). Methods like DSA[11], MMAD[12], CAMmarker[13], EPIC[14],

347    EPICabsolute[14], TIMER[10], TIMERtumor[10], and MuSiC[16] showed higher accuracies in the

348    absolute measurement scale ($r: 0.33 \sim 0.9$, $mAD: 0.21$) (Fig. 5b and Supplementary Fig. 15b)

349    than  in  the  relative  measurement  scale  ( $r: 0.22 \sim 0.68$ , $mAD: 0.17$ )  in  the  'mosaic'

350    column(Fig. 5a and Supplementary Fig. 15a).



351

352    **Fig.5| Evaluation results of Sim3**

353    **a,b,** Heatmaps of summarized evaluation metric based on Pearson's correlation coefficients on the (**a**) relative

354    measurement scale and (**b**) absolute measurement scale. In each heatmap, row indexes refer to the tested

355    methods and column indexes refer to the types of tumor spike-ins (small, large, and mosaic). **c,d,** Scatter plots

356    of estimated weights vs. ground truths of mixtures consist of 5 cellular components and mosaic tumor spike-

357    ins. (**c**) estimated weights vs. relative ground truth (**d**) estimated weights vs. absolute ground truth.

358

359        To further investigate the performance of deconvolution methods under the cell-type

360    resolution, we drew scatter plots of estimations from 5 Comp mixtures with 'mosaic' tumor

361    spike-ins and 'real' weight matrix (Fig. 5 c,d). In the relative measurement scale, CIBERSORT[7]

362    and CIBERSORTx[8] were the top performers and achieved high accuracy ($r \geq 0.95, mAD \leq$

363    $0.05$) (Fig. 5c and Supplementary Fig. 16). However, in the absolute measurement scale,

364    EPICabsolute[14] was the top performer and correctly estimated the absolute immune cell

365    proportions ($r \geq 0.95$, $mAD \leq 0.05$ ) (Fig. 5d and Supplementary Fig. 17). Based on

366    inconsistent evaluation results from two measurement scales, we suggest researchers pay

367    attention to the impact of measurement scales when performing deconvolution analysis on

368    mixtures with unknown contents.

369        Next, we checked the robustness of the three best performers in terms of component

370    number and tumor content in the 'real' weight matrix group. The robustness of CIBERSORT[7]

371    and CIBERSORTx[8]'s performance to the component number is high on the mAD evaluation

372    metric (mAD: $0.02 \sim 0.05$) in the relative measurement scale (Supplementary Fig. 16b).

373    EPICabsolute[14] also showed high robustness to the component number on the mAD

374    evaluation metric ($mAD$: $0.02 \sim 0.07$ ) in the absolute measurement scale(Supplementary

21

375    Fig. 17b). We found that having a larger variance in tumor content will increase the accuracy

376    of EPICabsolute[14], as we observed that with mosaic tumor spike-ins, EPICabsolute achieved

377    higher accuracies (r: 0.31~0.95, mAD: 0.02~0.05) than other tumor spike-in groups(r:

378    0.17~0.84, mAD: 0.02~0.07) (Supplementary Fig. 17) in the absolute scale. Consistent with

379    the observation in Sim2, we observed decreasing accuracies of CIBERSORT[7], CIBERSORTx[8],

380    and EPICabsolute[14] with the increasing component number (Supplementary Fig. 16a and

381    Supplementary Fig. 17a), and we deduced this phenomenon is due to the difficulty of current

382    deconvolution methods estimating rare subpopulations and closely related cell-types.

383         Our results revealed the impact of unknown biological content on deconvolution

384    analysis. We found both size (large vs. small spike-ins) and variance (large vs. mosaic spike-

385    ins) of unknown content affected deconvolution analysis. We also observed a discrepancy in

386    performance evaluation when used different measurement scales. In the relative scale, we

387    concluded CIBERSORT[7] and CIBERSORTx[8] were the top performers, while in the absolute

388    scale, EPICabsolute[14] was the top performer.

## Discussion

390         In this study, we designed three *in silico* benchmarking frameworks to systematically

391    explore the impact of several biological and technical factors. We identified top-performing

392    deconvolution methods for each framework and clearly illustrated the strengths and limits

393    of these tested methods under different application scenarios. Moreover, we offered several

394    strategies to mitigate systematic biases caused by different technical and biological factors

395    such as varied library sizes, simulation models, and cellular compositions.

396       In the first framework (Sim1), we explored the impact of noise structure under

397       different noise levels. We identified CAMmarker, MMAD, DSA, and CIBERSORT as the best

398       performers since these methods showed high accuracy and high robustness to diverse noise

399       levels. For the noise structure, we identified the negative binomial as the best simulation

400       model that captures the essential characteristics of real data. In the second framework

401       (Sim2), we explored the impact of the cellular component number and the weight matrix

402       property. We identified CIBERSORT, CIBERSORTx, and MuSiC as top-performers since these

403       two methods achieved high accuracies across a gradient of cellular component numbers with

404       both 'orthog' and 'real' weight matrices. We also found all marker-gene based methods

405       exhibited larger estimation deviances from ground truths, this type of estimation biases is

406       reflected in the scatter plots and can be quantitatively captured by the mAD evaluation

407       metric, indicating the necessity of using mAD as an auxiliary evaluation metric for

408       deconvolution performance assessment. In the third framework (Sim3), we explored the

409       impact of unknown biological content and measurement scales. In the relative measurement

410       scale, CIBERSORT and CIBERSORTx were the best performers. In the absolute measurement

411       scale, EPICabsolute was the best performer. Our analysis also illustrated different evaluation

412       results under the absolute and relative measurement scale, which have been overlooked in

413       the previous deconvolution benchmarks.

414       Based on the observations in this benchmark, we give the following suggestions for

415       best practices of deconvolution analysis and evaluations. For the *in silico* benchmarking data

416       generation, we suggest researchers 1) Use the negative binomial model as the primary

417       simulation model for *in silico* mixture data generation. 2) Referencing real biological

418       composition data when building weight matrices. 3) Consider at least two evaluation metrics.

419     One is used for checking linear concordance between estimation and ground truth, and the

420     other one is used for checking estimation deviances. 4) In the context of unknown biological

421     content, beware of the influence caused by different measurement scales(absolute vs.

422     relative). 5) Constructing multi-factor conditions on a large scale to ensure the robustness

423     and comprehensiveness of the benchmark.

424     For deconvolution analysis, we suggest researchers 1) Use the quantification unit

425     (countNorm, cpm, or tpm) that is normalized by library sizes. 2) Check for the compositional

426     information from previous publications. When the targeted tissue type has a relatively stable

427     composition over several samples, consider using deconvolution methods that are robust to

428     non-orthog weight matrices such as CIBERSORT, CIBERSORTx, and MuSiC. When an

429     unknown cellular component is expected (i.e., tumor sample) and the researcher needs to

430     derive absolute proportion, consider methods like EPIC, which is specifically tailored for

431     deconvolution with unknown content. 3) When referencing benchmark paper to select the

432     optimal method, beware of different technical factors that might derive different estimation

433     accuracies such as the resolution of analysis(number of cellular components), the variance

434     of proportions across samples(weight matrix property), reference selection, evaluation

435     metric selection, and measurement scale selection.

436     In addition to the suggestions mentioned above, previous benchmark publications

437     also clarified the impact of signature matrices[1], multicollinearity issue[7], spill-over effects[3,23]

438     caused by missing cellular components in the reference, minimal detection fraction[3],

439     background predictions[3], marker/signature gene selection[4,6], the variance between

440     reference and mixture sources[4]. Some deconvolution methods like CIBERSORT, CIBERSORTx,

441     and MuSiC can derive both cell-type-specific expression and composition signals. However,

442   by far, all independent deconvolution benchmark studies have been focused on the accuracy

443   of compositional information[3,6]. More benchmarks that derive accuracies of cell-type-

444   specific expression estimation are still in need.

445        For the future advancement of deconvolution analysis on RNA-seq data, we suggest

446   more efforts be put into the refinement of simulation models to generate more authentic *in*

447   *silico* testing environments that mimic diverse application scenarios. The weight matrix

448   property was revealed as the most important factor affecting deconvolution analysis in this

449   study and have been overlooked by the community. Therefore, more studies on the cellular

450   compositional information and its corresponding effects on deconvolution analysis are still

451   in need. Devotions on improving in silico benchmark generation strategy could further

452   enhance the efficiency of deconvolution method development and enable a wide range of

453   clinical applications.

## Methods

**Data processing:**

456        Raw SRA files were downloaded from the GEO repository, processed by SRA Toolkit

457   (2.10.0)[24], and reads were aligned to the hard masked human reference GRCh38 (v95) using

458   alignment tool STAR (2.6.1)[25], and quantification was performed with RSEM (1.3.1)[26] with

459   default parameter settings. Quantification matrices with the count, tpm, and fpkm units

460   were loaded into R (3.6.1)[27] for feature ID transformation, duplication removal, and low-

461   abundant gene removal. For low-abundant gene removal, we relied on two parameters:

462   minimum sample threshold (GSE113590[28] - 4, other datasets - 5) and minimum expression

463   threshold (10 counts, 1 tpm, and 1fpkm). For instance, the filtering parameter (5, 10) is used

464    to retain genes with more than 10 counts in at least 5 samples. GSE113590 only has 4

465    samples per cellular category, and we set the minimum sample thresholds as 4. In the Sim1,

466    we performed filtering independently on each dataset with a minimum sample threshold set

467    at 5. For Sim2 and Sim3, we first concatenated samples into one matrix and then performed

468    filtering with a minimum sample threshold set at 10. For the information of datasets involved

469    in Sim1, Sim2, and Sim3, please refer to Supplementary Table 4.

470    **Marker gene selection:**

471        For the marker gene selection, we selected genes that are highly expressed in the

472    targeted cell-type and lowly expressed in other cell-types. The expression threshold is set at

473    the 80th percentile for high expression (the targeted group) and 50th percentile for low

474    expression (other groups). Ideally, it would be nice if all samples pass the criteria; however,

475    to successfully derive marker genes with a larger number of cellular components, we

476    gradually relaxed the threshold (the percentage of samples pass the criteria, initial value p =

477    0.95) by a step parameter (default value s = 0.03) until there are at least two marker genes

478    determined.

479    **Signature gene selection:**

480        We performed differential expression testing on all cell-type pairs (all combinations

481    of 2 elements) using DESeq2[29]. Then we selected genes with $p_{adj} \leq 0.01$    and

482    $log2FoldChange \geq 10$.

483    **Benchmarking framework construction:**

26

484    Three benchmarking frameworks are constructed to study the impact of different

485    technical and biological factors on deconvolution analysis (Figure 1). We created simulated

486    mixture data M (N by J) by multiplying signature gene profiles S (N by K) to the predefined

487    weight matrix W (K by J). Here, N is the number of genes, J is the number of samples, and K

488    is the number of cellular components. The noise term $\varepsilon$ is used to model sample to sample

489    variability where the value of $\varepsilon$ determines the noise level.

490    $$M \ = \ S \times W \ + \ \varepsilon$$

491    **Sim1:** In the Sim1, we aimed at understanding the impact of noise from different

492    aspects such as noise structure and noise level. Sim1 consists of two sub frameworks:

493    Sim1_simModel and Sim1_libSize, where Sim1_simModel focuses on the noise structure, and

494    Sim1_libSize focuses on noise caused by varied library sizes.

495    **Sim1_simModel:** In this benchmarking framework, we mainly focused on the impact

496    of the simulation model that was used to generate noise. We selected three models for this

497    study, which are the normal, log-normal, and negative binomial models. For each simulation

498    model, we generated ten levels of noise where the magnitude of the noise is controlled by a

499    corresponding variance term in each model.

500    **Normal model:**

501    $$M \ = \ 2^{(log2(S \times W) + N(0, \sigma \times p_t))}$$

502

503    **Log-normal model:**

504    $$M = S \times W \ + \ 2^{N(0, \sigma \times p_t)}$$

27

505

506

507    In both Log-normal and Normal simulation models, the level of noise is controlled by

508    the product of a constant variance parameter $\sigma$ and a perturbation level parameter $p_t$. In this

509    study, we set $\sigma$ to 10 based on previous publications[7] and set $p_t$ as a length-10-vector (0, 0.1,

510    0.2, ... , 0.9).

511    **Negative binomial model:**

512    $$\mu_0 = r \times L_j$$

513    $$\mu_j = Gamma(shape = \frac{1}{\sigma^2}, scale = \frac{\mu_0}{shape})$$

514    $$\sigma = \left(1.8 \times p_t + \frac{1}{\sqrt{\mu_0}}\right) \times \delta \; where \; \delta \sim e^{N(0,0.25)}$$

515

516    $$v_j = Poisson(\mu_j)$$

517    We followed the simulation process suggested by Law *et al.*[19] and used $p_t$ to control

518    the noise level for simulation. r is a vector of genomic feature proportions, $L_j$ is the library

519    size and, $\mu_0$ is the expected gene expression in the simulation. In the negative binomial model,

520    two layers of variance are added from the Gamma distribution and Poisson distribution. We

521    derived sample gene expression vector $\mu_j$ from Gamma sampling to model biological

522    variance. In the Gamma distribution, the variance is determined by shape parameter $\sigma$. We

523    used $p_t$, a length-10 vector (0.1, 0.2, ..., 0.9, 1), to regulate the value of $\sigma$ to control the noise

28

524     level in the negative binomial simulation. Then we performed Poisson sampling to model

525     technical variance and get the final simulated expression vector.

526     To ensure the universality of our conclusion on different datasets, we applied the

527     Sim1 framework on 3 blood datasets to generate reference and *in silico* mixtures

528     (Supplementary Fig.1). Different from previous studies that concatenate samples derived

529     from different datasets, we generated 3 sets of simulated mixtures and 3 sets of references

530     independently. And then used combinations of mixtures and references to generate 9

531     replicated testing environments for each noise level. For one testing environment, there are

532     9 (3 times 3) deconvolution results from which 6 of them have mixture-reference pairs

533     derived from different sources. For simplicity, we only presented the averaged performance

534     across 9 mixture-reference pairs, but the impact of mixture-reference variance is considered

535     in this analysis. Above mentioned mixture-reference variance modeled in Sim1 is named as

536     other noise sources in Supplementary Table 2.

537     To understand the impact of quantification units over different application scenarios,

538     we generated simulations of the most commonly used RNA-seq quantification units: count,

539     countNorm, cpm, and tpm.

540
$$cpm_{i,j} = \frac{Count_{i,j}}{\sum_i Count_{i,j}} \times 10^6$$

541
$$tpm_{i,j} = \frac{Count_{i,j}}{L_{i,j}} \times \left(\frac{1}{\sum_i \frac{Count_{i,j}}{L_{i,j}}}\right) \times 10^6$$

29

542    Here j is the index of the sample and i is the index of the gene. cpm is normalized by

543    library size. countNorm is acquired from cpm units with every value rounded to the integer.

544    tpm is normalized by both library size and feature-length.

545    **Sim1_libSize**: In this testing framework, we mainly focused on bias derived from

546    varied library sizes. We first simulated mixtures based on the negative binomial model with

547    the lowest level of noise in Sim1_simModel ($p_1$ perturbation level). The library size variation

548    is controlled by the library size parameter $L_j$ in the negative binomial model. For every

549    simulation dataset that consists of 20 simulated profiles, we set the library size of the first

550    ten samples as 12 million reads and the remaining ten samples as 24 million reads

551    (Supplementary Fig. 1b).

552    **Sim2:** In this benchmarking framework, we studied the impact of cellular component

553    numbers and the mathematical property of the weight matrix (Supplementary Fig.2a).

554    Mixtures are generated based on the negative binomial model with the $p_1$ level noise. For

555    component number, we generated six sets of mixtures from 5 components up to 10

556    components. For the weight matrix, we generated two sets of weight matrix: orthog and real.

557    **Weight simulations:**

558    '**Orthog**' refers to the idealized weight matrix with a small condition number, which

559    provides a relatively optimal mathematical condition for deconvolution analysis. We first

560    simulated 1000 matrices (K by J) by randomly sampling weights from a uniform distribution

561    and then rescaled sampled weights so that for each mixture sample, all components sum to

562    1. Among 1000 proportion matrices, we picked the one weight matrix that has the smallest

563    condition number. '**Real**' refers to the weight matrix that mimics immune cell compositions

564　in the real whole blood sample. We generated weights based on uniform distribution with

565　min and max value defined based on previous observations of whole blood samples[21] and

566　then rescaled weights so that all components sum to 1.

567　**Sim3:** In this benchmarking framework, we studied the impact of unknown biological

568　content and measurement scales (Supplementary Fig.2b). To study unknown biological

569　content, we generated mixtures with tumor content spike-ins. In total, we created three sets

570　of tumor spike-ins: small, large, and mosaic. Tumor proportions are sampled from uniform

571　distributions and only differ in parameters used to set minimum and maximum values in the

572　sampling. 'Small' tumor spike-ins are sampled within the range 0-0.2, 'large' tumor spike-ins

573　are sampled within the range 0.7-0.9, and 'mosaic' tumor spike-ins are sampled within the

574　range 0.05-0.95. We then added three sets of tumor spike-in proportions to the weight

575　matrices generated in the Sim2 and rescaled them to have proportions of all components

576　sum to 1. After defining weights, we performed *in silico* mixing in the count unit and then

577　normalized it to other quantification units. To study the impact of the measurement scale,

578　we generated two sets of evaluations where one used absolute proportions of immune

579　components as the ground truth and the other used relative proportions of immune

580　components as the ground truth. The toy example of the absolute measurement scale and

581　the relative measurement scale is in Supplementary Table 5.

**Assessment of deconvolution performance**

583　J is the total number of mixture samples in a dataset and j is the sample index. $x_j$ is the

584　estimated proportion of sample j and $y_j$ is the ground truth of sample j. When a

585    deconvolution returns NA values, we directly assign highest penalty for the evaluation

586    metrics: r = -1, and mAD = 1.

587    **Pearson Correlation Coefficient (r):**

588

$$\frac{\sum_{j=1}^{J}(x_j - \bar{x})(y_j - \bar{y})}{\sqrt{\sum_{j=1}^{J}(x_j - \bar{x})^2 \sum_{j=1}^{J}(y_j - \bar{y})^2}}$$

589    **Mean Absolute Deviance (mAD):**

590

$$\frac{\sum_{j=1}^{J}|x_j - y_j|}{J}$$

591

592    **Datasets description:**

593        **1. GSE60424**[30] **-** Consists of 134 RNA-seq profiles of 6 immune cell types and whole

594         blood from both healthy donors and donors with five immune-associated diseases.

595        **2. GSE113590**[28] **–** Consists of 32 CD8 T cell RNA-seq profiles from peripheral blood,

596         colorectal tumor samples, and lung tumor samples.

597        **3. GSE64655**[31] **-** Consists of 56 RNA-seq profiles of 6 immune cell types and peripheral

598         blood from two vaccinated donors.

599        **4. GSE51984**[32] **–** Consists of 24 RNA-seq profiles of 5 immune cell types and total white

600         blood cells from healthy donors

601        **5. GSE115736**[33] **–** Consists of 42 RNA-seq profiles of 12 immune cell types from healthy

602         donors.

603        **6. GSE118490**[34] **–** HCT116 profiles (unknown tumor content in Sim3)

32

## Data and code availability

All data and codes are available in the https://github.com/LiuzLab/paper_deconvBenchmark under MIT license.

## Author contributions

H.J. designed, planned, and conducted data analysis and wrote the manuscript.

Z.L. supervised the analysis and wrote the manuscript. All authors read and approved the final manuscript.

## Competing interests

The authors declare no competing interests.

## References

1. Vallania, F. *et al.* Leveraging heterogeneity across multiple datasets increases cell-mixture deconvolution accuracy and reduces biological and technical biases. *Nat. Commun.* **9**, (2018).

2. Avila Cobos, F., Vandesompele, J., Mestdagh, P. & De Preter, K. Computational deconvolution of transcriptomics data from mixed cell populations. *Bioinformatics* 1–11 (2018). doi:10.1093/bioinformatics/bty019

3. Sturm, G. *et al.* Comprehensive evaluation of transcriptome-based cell-type quantification methods for immuno-oncology. *Bioinformatics* **35**, i436–i445 (2019).

4. Schelker, M. *et al.* Estimation of immune cell content in tumour tissue using single-cell RNA-seq data. *Nat. Commun.* **8**, 2032 (2017).

624    5.    Weber, L. M. *et al.* Essential guidelines for computational method benchmarking. *arXiv*

625          1–12 (2018).

626    6.    Cobos, F. A., Alquicira-Hernandez, J., Powell, J., Mestdagh, P. & De Preter, K.

627          Comprehensive benchmarking of computational deconvolution of transcriptomics

628          data. (2020). doi:10.1101/2020.01.10.897116

629    7.    Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression

630          profiles. *Nat. Methods* **12**, 1–10 (2015).

631    8.    Newman, A. M. *et al.* Determining cell type abundance and expression from bulk

632          tissues with digital cytometry. *Nat. Biotechnol.* **37**, 773–782 (2019).

633    9.    Mohammadi, S., Zuckerman, N., Goldsmith, A. & Grama, A. A Critical Survey of

634          Deconvolution Methods for Separating Cell Types in Complex Tissues. *Proc. IEEE* **105**,

635          340–366 (2017).

636    10.   Li, B. *et al.* Comprehensive analyses of tumor immunity: Implications for cancer

637          immunotherapy. *Genome Biol.* **17**, 1–16 (2016).

638    11.   Zhong, Y., Wan, Y.-W., Pang, K., Chow, L. M. L. & Liu, Z. Digital sorting of complex tissues

639          for cell type-specific gene expression profiles. *BMC Bioinformatics* **14**, 89 (2013).

640    12.   Liebner, D. A., Huang, K. & Parvin, J. D. MMAD: Microarray microdissection with

641          analysis of differences is a computational tool for deconvoluting cell type-specific

642          contributions from tissue samples. *Bioinformatics* **30**, 682–689 (2014).

643    13.   Chen, L. CAMTHC: Convex Analysis of Mixtures for Tissue Heterogeneity

644          Characterization. (2019).

645    14.    Racle, J., de Jonge, K., Baumgaertner, P., Speiser, D. E. & Gfeller, D. Simultaneous

646             enumeration of cancer and immune cell types from bulk tumor gene expression data.

647             *Elife* **6**, 1–25 (2017).

648    15.    Gong, T. & Szustakowski, J. D. DeconRNASeq: A statistical framework for

649             deconvolution of heterogeneous tissue samples based on mRNA-Seq data.

650             *Bioinformatics* **29**, 1083–1085 (2013).

651    16.    Wang, X., Park, J., Susztak, K. & Zhang, N. R. Bulk tissue cell type deconvolution with

652             multi-subject single-cell expression reference. *Nat. Commun.* doi:10.1038/s41467-

653             018-08023-x

654    17.    Zaitsev, K., Bambouskova, M., Swain, A. & Artyomov, M. N. based on linearity of

655             transcriptional signatures. *Nat. Commun.* 1–16 doi:10.1038/s41467-019-09990-5

656    18.    Zappia, L., Phipson, B. & Oshlack, A. Splatter: Simulation of single-cell RNA sequencing

657             data. *Genome Biol.* **18**, 1–15 (2017).

658    19.    Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: Precision weights unlock linear model

659             analysis tools for RNA-seq read counts. *Genome Biol.* **15**, R29 (2014).

660    20.    Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: A Bioconductor package for

661             differential expression analysis of digital gene expression data. *Bioinformatics* **26**,

662             139–140 (2009).

663    21.    Inc., S. T. Frequencies of Cell Types in Human Peripheral Blood. (2017).

664    22.    Calderon, D. *et al.* Landscape of stimulation-responsive chromatin across diverse

665             human immune cells. *Nat. Genet.* **51**, 1494–1505 (2019).

666    23.    Aran, D., Hu, Z. & Butte, A. J. xCell: Digitally portraying the tissue cellular heterogeneity

667            landscape. *Genome Biol.* **18**, 1–14 (2017).

668    24.    Leinonen, R., Sugawara, H. & Shumway, M. The sequence read archive. *Nucleic Acids

669            Res.* **39**, 2010–2012 (2011).

670    25.    Dobin, A. *et al.* STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21

671            (2013).

672    26.    Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with

673            or without a reference genome. *BMC Bioinformatics* **12**, 323 (2011).

674    27.    R Core Team. R: A Language and Environment for Statistical Computing. (2019).

675    28.    Simoni, Y. *et al.* Bystander CD8+ T cells are abundant and phenotypically distinct in

676            human tumour infiltrates. *Nature* **557**, 575–579 (2018).

677    29.    Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion

678            for RNA-seq data with DESeq2. *Genome Biol.* **15**, 1–21 (2014).

679    30.    Linsley, P. S., Speake, C., Whalen, E. & Chaussabel, D. Copy number loss of the interferon

680            gene cluster in melanomas is linked to reduced T cell infiltrate and poor patient

681            prognosis. *PLoS One* **9**, (2014).

682    31.    Hoek, K. L. *et al.* A cell-based systems biology assessment of human blood to monitor

683            immune responses after influenza vaccination. *PLoS One* **10**, 1–24 (2015).

684    32.    Pabst, C. *et al.* GPR56 identifies primary human acute myeloid leukemia cells with high

685            repopulating potential in vivo. *Blood* **127**, 2018–2027 (2016).

686  33.  Choi, J. *et al.* Haemopedia RNA-seq: A database of gene expression during

687       haematopoiesis in mice and humans. *Nucleic Acids Res.* **47**, D780–D785 (2019).

688  34.  Wagner, S. *et al.* Suppression of interferon gene expression overcomes resistance to

689       MEK inhibition in KRAS-mutant colorectal cancer. *Oncogene* **38**, 1717–1733 (2019).

690

691