1
2    # Detecting heterogeneity in single-cell RNA-Seq data by non-negative matrix factorization
3

4    Xun Zhu[1,2], Travers Ching[1,2], Xinghua Pan[3], Sherman Weissman[3], Lana Garmire[2,*]

5    [1] Molecular Biosciences and Bioengineering Graduate Program, University of Hawaii at

6    Manoa, Honolulu, Hawaii, United States of America

7    [2] Epidemiology Program, University of Hawaii Cancer Center, Honolulu, Hawaii, United

8    States of America

9    [3] Department of Genetics, Yale University, New Haven, Connecticut, United States of

10   America

11

12   [*] Corresponding author.

13   Email address: lgarmire@cc.hawaii.edu

14

15

16

17

18

19

1

20

# **Abstract**

Single-cell RNA-Sequencing (scRNA-Seq) is a cutting edge technology that enables the

understanding of biological processes at an unprecedentedly high resolution. However,

well suited bioinformatics tools to analyze the data generated from this new technology

are still lacking. Here we have investigated the performance of non-negative matrix

factorization (NMF) method to analyze a wide variety of scRNA-Seq data sets, ranging

from mouse hematopoietic stem cells to human glioblastoma data. In comparison to other

unsupervised clustering methods including K-means and hierarchical clustering, NMF

has higher accuracy even when the clustering results of K-means and hierarchical

clustering are enhanced by t-SNE. Moreover, NMF successfully detect the

subpopulations, such as those in a single glioblastoma patient. Furthermore, in

conjugation with the modularity detection method FEM, it reveals unique modules that

are indicative of clinical subtypes. In summary, we propose that NMF is a desirable

method to analyze heterogeneous single-cell RNA-Seq data, and the NMFEM pipeline is

suitable for modularity detection among single-cell RNA-Seq data.

## Introduction

36

37   The advancement of technologies has enabled researchers to separate individual cells

38   from a bulk and sequence their transcriptomes at the single cell level, known as single-

39   cell RNA-Sequencing (scRNA-Seq). This technology has reached an unprecedented fine

40   resolution to reveal the program of gene expression within cells(Kumar et al., 2014). It

41   was used to detect heterogeneity within the cell population, and it has greatly enhanced

42   our understanding of the regulatory programs involved in systems such as

43   glioblastoma(Patel et al., 2014), neuronal cells(Usoskin et al., 2014), or pluripotent stem

44   cells (PSCs)(Kumar et al., 2014). It was also used to delineate cell types and

45   subpopulations in differentiating embryonic cells(Treutlein et al., 2014). Other

46   applications include uncovering multilineage priming processes involved in the initial

47   organogenesis(Brunskill et al., 2014), and substantiating the hypothesis of inter-

48   blastomere differences in 2- and 4-cell mouse embryos(Biase, Cao & Zhong, 2014).

49   Indeed, ScRNA-Seq has already made profound impacts on our understanding of the

50   diversity, complexity, and irregularity of biological activities in cells. It will continue to

51   provide more transformative insights in the near future(Pan, 2014).

52   However, relative to the experimental technology, the bioinformatics tools to analyze

53   scRNA-Seq data are still lagging behind. Given the large amount of noise in the scRNA-

54   Seq data, it is unclear if the tools developed for population-level RNA-Seq differential

55   expression analysis, such as DESeq2(Love, Huber & Anders, 2014) and

56   EdgeR(Robinson, McCarthy & Smyth, 2010), are desirable to identify subpopulations in

57   scRNA-Seq data. Recently, a couple of methods have been reported in the scRNA-Seq

58   analysis domain (Brennecke et al., 2013; McDavid et al., 2013; Kharchenko, Silberstein

3

59   & Scadden, 2014). For example, a statistical variance model based on gamma distribution

60   was developed to account for the high technical noise occurring in scRNA-seq

61   experiments, such that genes with high squared correlation of variations ($CV^2$) relative to

62   mean expression are identified as "significantly differentially expressed" between two

63   conditions(Brennecke et al., 2013). Another Bayesian approach was proposed for

64   scRNA-Seq differential expression analysis, by utilizing a probabilistic model of

65   expression-magnitude distortions that commonly observed in noisy single-cell

66   experiments(Kharchenko, Silberstein & Scadden, 2014). This method later was used for

67   classification of sensory neurons using scRNA-Seq(Usoskin et al., 2014). On the other

68   hand, an R package Monocle was developed recently for single-cell lineage

69   construction(Trapnell et al., 2014). However, it is not clear if all these new methods are

70   suitable for detecting subpopulations in single cells. Moreover, none of the packages

71   mentioned above offers functionalities for modularity identification. For the purpose of

72   network module detection, one has to either use the RNA-Seq transcriptome data as the

73   input for packages such as Module Networks in Genomica(Segal et al., 2003), or use the

74   discovered important genes as seeds to combine with other downstream module detection

75   packages. The fast accumulation of scRNA-Seq data requires new tools to study single-

76   cell transcriptome more efficiently.

77   Previously, NMF has been applied to other areas in computational biology, such as

78   molecular pattern discovery(Brunet et al., 2004), class comparison and prediction(Gao &

79   Church, 2005), cross-platform and cross-species analysis(Tamayo et al., 2007),  and

80   identify subpopulations of cancer patients with mutations in similar network regions.

81   Moreover, NMF has been applied to gene expression profiling studies, in both array(Qi et

4

82    al., 2009) and population-level RNA-Seq platforms(Brunet et al., 2004). Compared to

83    other methods, it showed multiple advantages, such as less sensitivity to a priori selection

84    of genes or initial conditions and the ability to detect context-dependent patterns of gene

85    expression(Rajapakse, Tan & Rajapakse, 2004). Based on these properties, we

86    hypothesize that NMF is less prone to the influence of noise in the scRNA-Seq data, and

87    thus it can detect a group of genes that robustly differentiate single cells from different

88    conditions. In this report, we demonstrate the capabilities of NMF in scRNA-Seq data

89    analysis in these following aspects: (1) accurate clustering of single cells from different

90    conditions in an unsupervised manner; (2) stratification of subpopulations within the

91    same pool of single cells; (3) detection of meaningful genes, pathways and modules

92    associated with differences among populations and subpopulations. We also combine

93    NMF with the modified, seed based module detection tool Functional Epigenetic

94    Modules (FEM)(Jiao, Widschwendter & Teschendorff, 2014), and provide the scientific

95    community with a streamlined modularity detection R package called NMFEM.

## Results

97    The workflow for a typical single-cell analysis using NMF is shown in Fig. 1. Briefly, the

98    pipeline can take raw reads in FastQ files, align and count them to the RefSeq

99    transcriptome, or use raw count data directly as the input matrix. The input data matrix is

100   then subject to quality control and normalization steps. The normalized matrix is operated

101   on by NMF, which clusters the samples into sub-populations and enlists the feature genes

102   that separate the sub-populations. In order to display the insightful biological modules,

103   the feature genes are then used as the seeds for a functional modularity detection

104   algorithm FEM(Jiao, Widschwendter & Teschendorff, 2014), which identifies hotspots in

5

105    the interactome with the scRNA-Seq profiling. We applied this workflow to four scRNA-

106    Seq data sets, varying from mouse hematopoietic stem cells to human glioblastoma

107    primary cancer cells.

## NMF accurately clusters RNA-Seq data from hematopoietic stem cell

## differentiation

110    We first compared the accuracies of NMF in unsupervised clustering, compared to two

111    other commonly used methods: K-means and hierarchical clustering (Hclust) algorithms.

112    We tested these clustering methods on a data set composed of mouse hematopoietic stem

113    cells (HSCs) and stage 1 multipotent progenitor cells (MPP1). These cells were classified

114    using the combined CD62L and CD97 cell surface markers. In order to evaluate the

115    performance of the clustering methods, we removed the cell surface marker based labels.

116    As shown in the PCA plots in Fig. 2A, NMF is the most accurate method, while K-means

117    and hierarchical clustering are much worse. These observations can be quantitatively

118    supported by the results of pairwise Rand measure, a metric that describes the percentage

119    of agreement on a pair of samples belonging to the same group (Fig. 2C). Even though

120    the two cell types are closely related on cell lineage, NMF achieves an overall impressive

121    Rand measure of 83.6% to classify RNA-Seq data by patient ID. In contrast, K-means

122    and hierarchical clustering have much lower Rand measures of 50.6% and 49.7%,

123    respectively (Fig. 2C). Additionally, we plotted the consensus heatmaps of two of the

124    methods — NMF and K-means, which clearly shows the higher accuracy of NMF over

125    K-means (S1 Fig.).

126    Next we investigated the effect of t-SNE modification on NMF, K-means and

127    hierarchical clustering (Fig. 2B). t-SNE is a dimension reduction method that works by

6

128    minimizing the KL-divergence between the distribution of original distances and the

129    distances in the lower-dimensional space. Methods such as K-means are usually

130    conjugated with t-SNE(Van der Maaten & Hinton, 2008) to improve the accuracy of

131    clustering and to be used as a method of visualization in 2-dimensional space(Van der

132    Maaten & Hinton, 2008; Bushati et al., 2011; Junker et al., 2014). However, since NMF

133    is not a distance-based method, applying t-SNE does not improve rather worsen the

134    clustering results of NMF (Fig. 2B and 2C). With the two features extracted by t-SNE,

135    NMF loses its ability to extract meta-genes and to conduct component decomposition, as

136    demonstrated by the clustering accuracy (measured by Rand measure) before and after

137    using t-SNE. On the contrary, K-means and hierarchical clustering have improved

138    accuracies after the application of t-SNE (Fig. 2B and 2C). However, since the

139    differences between HSC vs. MPP1 are very subtle, the ability of t-SNE to improve the

140    clustering accuracy is limited (Fig. 2C).

141    We repeated the same analytical comparisons with another set of dendritic cell

142    differentiation data(Schlitzer et al., 2015), and obtained similar conclusion. That is, NMF

143    has better accuracy than distance-based methods such as K-means and hierarchical

144    clustering, even when the other two methods are boosted by t-SNE (S2 Fig.).

145    **NMF discovers uniquely important genes in mouse embryonic lung**

146    **distal epithelium development**

147    Unlike other conventional differential expression test methods that explicitly model the

148    relationships between the variance and mean in the RNA-Seq data, NMF selects the

149    important genes by Kullback–Leibler divergence (KL-divergence)(Yang et al., 2011).

150    Note, these "important genes" are by no means "differentially expressed (DE) genes", as

7

151    defined by the differential gene expression (DGE) statistical tests. For comparison, we

152    include the recently developed methods for single-cell transcriptome analysis, including

153    Monocle(Trapnell et al., 2014), MAST(McDavid et al., 2013) as well as

154    SCDE(Kharchenko, Silberstein & Scadden, 2014), as well as DESeq2 and EdgeR, two

155    commonly used differential gene selection methods for the bulky RNA-Seq data. We

156    chose another set of mouse embryonic lung distal epithelial cells reported by Treutlein et

157    al.(Treutlein et al., 2014), and focus on the single cells from E14.5 and E16.5 stages,

158    where the RNA-Seq data are so similar that even PCA analysis cannot separate clearly

159    (S3 Fig.). Given that rich experiential knowledge has been accumulated on their

160    developmental process, this dataset allows us to empirically evaluate the results obtained

161    from different RNA-Seq analysis tools.

162    We present the characteristics of "important genes" detected by each method in the MA-

163    plots (Fig. 3). The uniquely identified genes from these methods vary greatly (Fig. 3 and

164    S4 Fig. A). In contrast with all other compared methods, NMF selects genes that are

165    sufficiently expressed in many samples, with a strong preference to select genes around a

166    specific expression level (FPKM 2.740) and but not genes expressed too lowly or too

167    highly (S4 Fig. A). On the other hand, a fair amount of genes selected by MAST, SCDE,

168    and Monocle have very little numerical differences between E14.5 and E16.5 stages. A

169    considerable amount of genes selected by DESeq2 and EdgeR have average low

170    expressions but large variance (Fig. 3). Many of them have zero count in all samples of

171    E16.5 stage. Since lowly expressed genes usually have much higher levels of noise, this

172    suggests that DESeq2 and EdgeR may have detected the expression patterns that are less

173    reliable(Brennecke et al., 2013).

8

174   Such a group of intermediately expressed genes identified by NMF are robust and

175   unlikely a random sample from all expressed genes, since the density distribution of the

176   top 500 genes in NMF per drop-one-out resampling is clearly distinctive from that of

177   random background gene expression (S4 Fig B). The reason that NMF tends to avoid the

178   extremely lowly expressed genes is that KL-divergence intrinsically penalizes lowly

179   expressed genes as $A_{ij}$ can be seen as the weight of $(\log\left(\frac{A_{ij}}{(WH)_{ij}}\right))$ in the formula (see

180   Methods). The lower the original expression level, the weaker that gene can affect the

181   clustering, and thus less likely to be selected as a feature gene by NMF. On the other

182   hand, the highly expressed genes typically have extreme spikes among a few samples,

183   and are also less likely to be selected as feature genes, as the signal linearity of NMF

184   prefers to select genes with consistent expression levels in each cluster.

185   **Important genes selected by NMF yield biologically meaningful modules**

186   We next asked if the important genes detected by NMF convey unique and meaningful

187   biological functions. Towards this, we examined the modularity potentials and used the

188   same number of 500 top genes selected by the eight methods above as the initial seeds for

189   the module detection software FEM(Jiao, Widschwendter & Teschendorff, 2014). FEM

190   is a versatile method that can be adapted to identify hotspots in the interactome with the

191   differential expression profiling, using the seed inputs from external programs including

192   NMF, DESeq2, EdgeR, MAST, SCDE, or Monocle.  We present the results of the top 5

193   most significant modules for each of the eight methods. Within each top module, we

194   conducted Gene Ontology (GO) enrichment analysis and list the top two GO terms

195   (Table 1).

9

196    In comparison, the methods that are established on similar assumptions have higher

197    degrees of agreements on the detected top modules (Table 1) as well as genes in common

198    (S5 Fig.), as expected. For examples, SCDE, MAST and Monocle have more similar

199    results than others; whereas DESeq2 and EdgeR tend to agree to each other better since

200    they were designed for bulky cell RNA-Seq.  Interestingly, all methods except EdgeR,

201    detected that the transcription-related processes play important role from E14.5 to E16.5.

202    NMF finds two unique modules for "mRNA destabilization" (seed gene Pnn) and "rRNA

203    processing" (seed gene exosc4) (Table 1 and Fig. 4).  These results are very interesting as

204    mRNA-destabilizing inflammatory RNA-binding proteins were previous reported to be

205    important in the regulation of miR-155 biogenesis in lung epithelial cells with cystic

206    fibrosis condition(Bhattacharyya et al., 2013). Exosc4 is part of the exosome complex,

207    which has the function of degrading various types of RNA molecules. Since E14.5 cells

208    are prior to sacculation and E16.5 cells are in the early stage of sacculation, the exosc4-

209    centered module may indicate the fast turnover of RNA material associated with the cell

210    growth/apoptosis activities in the process of embryonic lung morphological changes.

211    Additionally, NMF identifies a module related to "G-protein coupled receptor signaling

212    pathway" (seed gene Gna13), which is also shared by DESeq2 and EdgeR methods

213    (Table 1 and Fig. 4). This may indicate active intracellular signal changes during the

214    early phase of embryonic lung epithelial cells. This observation is coherent with another

215    unique module found by NMF, which is related to bone morphogenetic protein (BMP)

216    pathway (seed gene Smad4). BMP pathway previously was verified to have important

217    roles in signal transduction, transcription and adhesion in epithelial bud development,

218    including lung epithelial cells(Jamora et al., 2003).  Moreover, BMPs play important

219    roles in different stem cell systems, including embryonic stem cells(Zhang & Li, 2005).

220    In summary, due to the mechanism of identifying correlated genes rather than genes with

221    numerical differences, NMF is able to extract very unique biological information from

222    different classes of single cells.

223    **NMF identifies tumor sub-populations among a single glioblastoma**

224    **patient**

225    Detecting the subpopulations of single cells within the same bulk is an even subtler

226    problem, in comparison to the issue of accurate clustering of mixed populations. To

227    examine the potential of NMF in this aspect, we next tested the scRNA-Seq data from the

228    five individual glioblastoma patients as reported by Patel, AP et al.(Patel et al., 2014)

229    Interestingly, the consensus clustering results generated from NMF show that among the

230    five patients, only patient MGH28 (Fig. 5A-B) and MGH31 (S6 Fig. A-B) have two

231    distinct subpopulations.

232    To investigate further the characteristics of the two subpopulations in MGH28, we

233    retrieved the top 500 ranked genes that differentiate these two subpopulations and

234    conducted KEGG pathway enrichment analysis on them.  A pathway named "pathogenic

235    Escherichia coli infection" stands out as the most significantly altered pathway between

236    the two subpopulations (FDR < 1E-03) (Fig. 5C). Further examination of this pathway

237    reveals that multiple genes involved in cell mobility are enriched, including ACTG1,

238    ACTB, CTTN, YWHAZ, CDC42, TUBB, RHOA, ROCK, ARPC5, TUBA1A, NCL,

239    TUBA1B, and TUBA1C (Fig. 5D). Glioblastoma is among the most heterogeneous

11

240    tumors in human, and mainly have pro-neuron and mesenchymal phenotypes. The latter

241    is associated with more invasive and infiltrating phenotype. Our results indicate that

242    some cells in patient MGH28 have mesenchymal phenotype. Coincidently, Patel, AP et al

243    also concluded MGH28 as mesenchymal glioblastoma, by comparing the scRNA-Seq

244    signatures to those from TCGA glioblastoma RNA-Seq data(Patel et al., 2014).

245    Interestingly, we also found that patient MGH31 has the same enriched KEGG pathway

246    term of "pathogenic Escherichia coli infection" (S6 Fig. C). Almost all of the important

247    genes in this pathway from patient MGH31 (S6 Fig. D) overlap those from patient

248    MGH28 mentioned above (Fig. 5D). The only exceptions are NCL unique to MGH28,

249    and CDC42 and ROCK2 unique to MGH31. The almost identical genes found in the

250    same pathway that differentiates the subpopulations of both MGH28 and MGH31 suggest

251    that MGH31 may also be classified as mesenchymal glioblastoma, similar to MGH28.

## 252    Discussion and conclusions

253    Due to the high noise levels within scRNA-Seq data(Brennecke et al., 2013), the

254    conventional approaches, which aim to detect numerical differences of gene expression in

255    cell bulks under different conditions, may be limited. Previous applications of NMF to

256    fields such as face reorganization(Rajapakse, Tan & Rajapakse, 2004), image

257    compression(Yuan & Oja, 2005; Monga & Mıhçak, 2007) and sound

258    decomposition(Smaragdis, 2004), have proven successful. Here we propose to utilize

259    NMF as a desirable method for scRNA-Seq analysis. We believe that the pattern based

260    feature extraction ability of NMF can meet the demands to identify genes that signify the

12

261   differences within the noisy scRNA-Seq data. The in-depth analyses on multiple public

262   and private data sets in this study have provided supports from several aspects.

263   We have demonstrated that NMF performs well relative to other popular clustering

264   methods including K-means and hierarchical clustering, even when these methods in

265   comparisons are boosted with t-SNE. Moreover, NMF is capable of identifying

266   subpopulations within the same tumor sample, exemplified by the glioblastoma data here.

267   Through NMF clustering, we found in that patients MGH28 and MGH31 both have a

268   group of genes that can distinguish the single cells into two subpopulations. These genes

269   include actins, tubulins and signaling molecules that can affect cell mobility. Thus we

270   speculate that both MGH28 and MGH31 have mesenchymal phenotypes. The suspected

271   mesenchymal phenotype of MGH28 from scRNA-Seq data alone is directly supported by

272   Patel, AP et al.(Patel et al., 2014), where they used TCGA glioblastoma data and

273   classified MGH28 as the mesenchymal type. On the other hand, the authors could not

274   clearly classified MGH31 as the mesenchymal type, although they suspected two genetic

275   clones from this patient. Here with NMF based subpopulation identification and

276   comparisons of characteristic genes, our analysis confirms the existence of two

277   subpopulations and further, the clinical subtype of MGH31.

278   In summary, we have demonstrated that NMF is a desirable method capable of

279   accomplishing various tasks in scRNA-Seq data analysis, from reclassifying populations

280   of single cells, identifying subpopulations, to revealing meaningful genes, gene sets and

281   modules of biological significance. We expect the new workflow named NMFEM to

282   have wide applications in the field of scRNA-Seq bioinformatics analysis.

13

## Methods

### Data sets

#### Glioblastoma

ScRNA-Seq data were retrieved from the original 875 samples of glioblastoma tumor cells in 5 patients, along with population and cell line controls (GSE57872)(Patel et al., 2014). For NMF, very minimal filtering was employed (filtering steps of other methods are detailed in a later section). First, genes with zero expression across all samples were removed so that 22704 out of 23710 genes (95.8%) remained. Next the smallest number of samples was removed so that at least one gene was expressed across all samples considered, as a quality requirement of DESeq2. As a result, 124 samples with the lowest amount of non-zero expression across all genes are removed, leaving 751 of 875 samples (85.8%).

#### Mouse lung epithelial cells

ScRNA-Seq data were retrieved from the original 201 samples of lung distal epithelial cells of embryonic mouse (GSE52583)(Treutlein et al., 2014). We filtered genes and samples following the sample procedure as in Glioblastoma data set, leaving 16168 out of 23420 genes (69.0%) and 199 out of 201 samples (99.0%).

#### Mouse HSCs and MPP1s

ScRNA-Seq data were extracted from mouse hematopoietic stem cells (HSCs) and early multipotent progenitors (MPP1s). The data were pre-processed into the format of a FPKM expression profile, which include 59 HSCs and 53 MPP1 single cells. We filtered

14

304   genes and samples following the sample procedure as in Glioblastoma data set, leaving

305   12719 out of 21664 genes (58.7%) and 112 out of 112 samples (100.0%).

## Mouse dendritic cells

307   ScRNA-Seq data were extracted from mouse macrophage DC progenitors (MDPs),

308   common DC progenitors (CDPs), and Pre-DCs (GSE60781)(Schlitzer et al., 2015). We

309   used the RPKM table provided by the authors. We filtered genes and samples following

310   the same procedure as in Glioblastoma data set, leaving 15722 out of 29779 genes

311   (52.8%) and 251 out of 251 samples (100.0%).

# Single-cell RNA-Seq analysis

## Read alignment

314   We downloaded the public datasets from NCBI The Gene Expression Omnibus (GEO)

315   database(Edgar, Domrachev & Lash, 2002; Barrett et al., 2013), and retrieved the SRA

316   files from The Sequence Read Archive (SRA)(Leinonen et al., 2011). We used the fastq-

317   dump tool from SRA Toolkit to convert the SRA files into two pair-end FastQ files. We

318   applied FastQC for quality control and Tophat2(Kim et al., 2013) for alignment to the

319   reference genomes. The ready-to-use genome sequences and annotation files were

320   downloaded from Illumina iGenomes page

321   (http://support.illumina.com/sequencing/sequencing_software/igenome.html). For human

322   build hg19 was used, and for mouse genome build mm10 was used(Karolchik et al.,

323   2014).

15

324 **Read Counting**

325 We used featureCounts(Liao, Smyth & Shi, 2014) to map and count the aligned BAM

326 files to the RefSeq transcriptomes from the pre-built packages on Illumina iGenome

327 website above. We used the options to count fragments instead of reads; paired-end

328 distance was checked by featureCounts when assigning fragments to meta-features or

329 features. We only took into account of fragments that have both ends aligned successfully

330 and discarded chimeric fragments. Fragments mapped to multiple locations were counted.

331 The command is "featureCounts -pPBCM --primary -T 6 -a <gtf_file> -o <output_file>

332 <bam_file>".

333 **Normalization of Counts**

334 We used reads per kilo base per million (RPKM) to represent the gene expression level,

335 where the length of each gene was calculated by UCSC RefSeq annotation table, by

336 concatenating all the exons. We normalized the data using DESeq2.

337 # Non-negative Matrix Factorization (NMF)

338 We used the R-package implementation of NMF(Gaujoux & Seoighe, 2010) to perform

339 NMF analysis. NMF is mathematically approximated by: $A \approx WH$, where $A$ ($n$ by $m$) is

340 the matrix representing the scRNA-Seq profile in this report, W is a slim weight matrix

341 ($n$ by $k$, where $n \gg k$), H is a wide matrix ($k$ by $m$, where $m \gg k$), and all three of them

342 are non-negative(Brunet et al., 2004). The column vectors in $W$ are called meta-genes,

343 which are higher-level abstraction of the original gene expression pattern. We used the

344 method "*brunet*" to solve the approximation of $A$, which employs the multiplicative

345 iterative algorithm described by the following rules:

16

346
$$H_{au} \leftarrow H_{au} \frac{\sum_i \frac{W_{ia}V_{iu}}{(WH)_{iu}}}{\sum_k W_{ka}}$$

347
$$W_{ia} \leftarrow W_{ia} \frac{\sum_u \frac{H_{au}A_{iu}}{(WH)_{iu}}}{\sum_v H_{av}}.$$

348 The initialization of $H_{au}$ and $W_{ia}$ was generated as random seed matrices drawn from a

349 uniform distribution within the same range as the entries in the matrix $A$. Since the

350 starting matrices were randomized, we conducted an average of 30 simulations for each

351 NMF run to obtain the consensus clustering results. We used Kullback–Leibler

352 divergence (KL-divergence) as the distance function, as it has significantly better

353 performance theorized in Yang et al.(Yang et al., 2011). The rank ($k$) is chosen by listing

354 the clustering results of all possible $k$'s (usually ranging from 2 to 5, as higher $k$ values

355 requires exponentially more time to run). $k$ is chosen when the best cophenetic

356 correlation coefficient is achieved, as proposed in Brunet et al. 2004(Brunet et al., 2004).

357 NMF package uses the *feature score* to measure the genes for different expression

358 between sample groups, based on a method described in Kim et al.(Kim et al., 2013)

359
$$\text{FeatureScore}(i) = 1 + \frac{1}{\log_2 k} \sum_{q=1}^{k} p(i,q) \log_2 p(i,q),$$

360 where

361
$$p(i,\Omega) = \frac{W(i,\Omega)}{\sum_{q=1}^{k} W(i,q)}.$$

362 The feature score lies between 0 and 1, and is positively related to its factor-specificity.

363 That is, a higher feature score indicates that the gene has more different expression

17

364    patterns between sample groups (phenotypes)(Kim & Park, 2007). We select the top 500

365    genes of NMF based on this feature score.

366    **Other packages used for detecting significant or important genes**

367    We compared a series of computational methods to call "significant genes" with NMF.

368    These methods are divided into three categories.

369    *DE methods for bulky-level RNA-Seq*:  we used two most popular bulky-level RNA-Seq

370    methods: DESeq2 and EdgeR, to compare on the results of DE genes.

371    *DE methods for scRNA-Seq*: three methods were investigated, with default settings of the

372    packages. (1) Monocle: this is a versatile method (V. 1.0.0) that performs differential

373    expression analysis between cell types or states, moreover places cells in order according

374    to their progression through processes such as cell differentiation(Trapnell et al., 2014).

375    (2) SCDE: this package (V 1.2.1) implemented in R is based on Bayesian method, where

376    the individual genes were modeled explicitly as a mixture of the dropout and

377    amplification events by the Poisson model and negative binomial model(Kharchenko,

378    Silberstein & Scadden, 2014). (3) MAST: this method (V 1.0.1) implemented in R was

379    originally used to detected DE genes in qPCR results of single cells. We selected the 500

380    genes with the lowest likelihood ratio test p-value using Hurdle Model provided by the

381    package, as recommended by the authors(McDavid et al., 2013).

382    *Data filtering for other scRNA-Seq methods:*  SCDE model deals with high level noise

383    automatically and requires no filtering as stated by authors. For Monocle and MAST, we

384    first removed the genes of high technical variations using the method as described in

385    Brennecke et al. 2013(Brennecke et al., 2013), then performed filtering steps as instructed

18

386 in each paper. Monocle filters out libraries that contained fewer than 1 million reads in its

387 original report, in the case that reads in some data set do not meet this threshold (such as

388 mouse embryonic lung epithelial cell data), we resorted to no sample filtering to be safe.

389 Additionally, we experimented if introducing t-SNE, a dimension reduction method that

390 was recently successfully applied to scRNA-Seq, would improve the results of NMF. We

391 used the C++ accelerated R-package Rtsne (V 0.10), based on the original C++

392 implementation by van der Maaten et al.(van der Maaten, 2013)

## Module detection package

394 We use Functional Epigenetic Modules (FEM) R package(Jiao, Widschwendter &

395 Teschendorff, 2014) for module detection. FEM utilizes an expansion algorithm based on

396 the z-score of the expression level, by using a list of seed genes as the starting points. It

397 selects the top modules based on p-values calculated by a Monte Carlo method.

398 We modified the source code of the FEM package and changed the process of the seed

399 gene selection. Instead of selecting the seed genes based on the z-score of the expression

400 level, we directly plugged in a list of genes as the seed genes, which were generated from

401 each of the compared method for important gene detection.

## Measuring the performance of unsupervised clustering

## methods

### Label assignments for PCA/t-SNE plots

405 Since multiple assignments of labeling to clusters are possible, for each clustering

406 algorithm we iterated through all possible permutations of labeling and calculated the

19

407     accuracy for each. The one with the best accuracy rate is picked as the *most favorable*

408     *labeling* for the clustering algorithm and is used in plotting its PCA/t-SNE scatter-plots.

## Confusion matrix

410     Confusion matrix $C$ was calculated by the following formula:

411     $$C_{i,j} = |A_i \cap B_j|,$$

412     Where $A_i$ is the set of samples that are labeled as class $i$ according to the correct

413     labelling, and $B_j$ is the set of samples that are labeled as class $j$ in the tested

414     method(Stehman, 1997).

## Chi-square test score

416     Chi-square test score $S_{\chi^2}$ was calculated from the chi-square test p-value $p_{\chi^2}$,

417     $$S_{\chi^2} = \log_{0.05} p_{\chi^2},$$

418     which in turn was calculated by the *chisq.test* function in R(Aguirre & Nikulin,

419     1994).The base of 0.05 was chosen so that a score larger than one indicates that the

420     resulting p-value is significant.

## Pair-wise Rand measure

422     Pair-wise Rand measure for clustering between the test and the reference is defined by

423     $$R = \frac{TP + TN}{TP + FP + FN + TN},$$

424     in which the four quantities $TP$, $FP$, $FN$, and $TN$ are cardinals of the four sets of pairs.

425     $T/F$ means true/false based on the reference, and $P/N$ means positive/negative results

20

426   from the test. Specifically, a positive result ($P$) refers to a pair of samples clustered in the

427   same group by the tested method; a true positive ($TP$) or true negative ($TN$) result

428   represents the case where the agreements between the test and the reference clustering is

429   reached(Rand, 1971).

430   **Modularity detection and pathway Analysis**

431   We used Functional Epigenetic Modules (FEM) package(Jiao, Widschwendter &

432   Teschendorff, 2014) implemented in R for module detection. FEM utilizes SpinGlass

433   algorithm(Reichardt & Bornholdt, 2006) based on the z-score of the expression level, by

434   using a list of seed genes as the starting points. It selects the top modules based on p-

435   values calculated from a Monte Carlo method. We modified the source code of the

436   package to allow seed genes generated from other methods (NMF, DESeq2, EdgeR,

437   SCDE, MAST and Monocle) that detect significant or important genes. In each case, we

438   used top 500 most important genes as the seeds for FEM. We next compared biological

439   meanings of the resulting modules by Gene Ontology (GO) or Kyoto Encyclopedia of

440   Genes and Genomes (KEGG) pathway enrichment analysis, implemented as DAVID

441   Web Service in R(Huang, Sherman & Lempicki, 2008, 2009).

# Data and code availability

443   The Glioblastoma, mouse lung distal epithelial and mouse dendritic cell data are

444   downloaded from GSE57872, GSE52583, and GSE60781. The code used for this

445   package can be found at https://github.com/lanagarmire/NMFEM, and

446   https://github.com/lanagarmire/NMFEM_extra.

# Author contributions

LXG envisioned the project. XZ conducted the data analysis, with assistance from TC.

XP and SW communicated on bioinformatics analysis and provided the mouse HSC and

MPP1 scRNA-Seq data. XZ and LXG wrote the draft. All authors have read, reviewed

and agreed on the manuscript.

# Acknowledgement

# Competing interests

The authors declare that they have no competing interests.

# References

Aguirre N., Nikulin M. 1994. Chi-squared goodness-of-fit test for the family of logistic

   distributions. *Kybernetika* 30:214–222.

Barrett T., Wilhite SE., Ledoux P., Evangelista C., Kim IF., Tomashevsky M., Marshall

   KA., Phillippy KH., Sherman PM., Holko M., Yefanov A., Lee H., Zhang N.,

   Robertson CL., Serova N., Davis S., Soboleva A. 2013. NCBI GEO: archive for

   functional genomics data sets—update. *Nucleic Acids Research* 41 :D991–D995.

   DOI: 10.1093/nar/gks1193.

22

467　　Bhattacharyya S., Kumar P., Tsuchiya M., Bhattacharyya A., Biswas R. 2013. Regulation

468　　　　　of miR-155 biogenesis in cystic fibrosis lung epithelial cells: Antagonistic role of

469　　　　　two mRNA-destabilizing proteins, KSRP and TTP. *Biochemical and Biophysical*

470　　　　　*Research Communications* 433:484–488. DOI: 10.1016/j.bbrc.2013.03.025.

471　　Biase F., Cao X., Zhong S. 2014. Cell fate inclination within 2-cell and 4-cell mouse

472　　　　　embryos revealed by single-cell RNA sequencing. *Genome research*:gr–177725.

473　　Brennecke P., Anders S., Kim JK., Kołodziejczyk AA., Zhang X., Proserpio V., Baying

474　　　　　B., Benes V., Teichmann SA., Marioni JC. 2013. Accounting for technical noise in

475　　　　　single-cell RNA-seq experiments. *Nature methods*.

476　　Brunet J-P., Tamayo P., Golub TR., Mesirov JP. 2004. Metagenes and molecular pattern

477　　　　　discovery using matrix factorization. *Proceedings of the National Academy of*

478　　　　　*Sciences*　101 :4164–4169.

479　　Brunskill EW., Park J-S., Chung E., Chen F., Magella B., Potter SS. 2014. Single cell

480　　　　　dissection of early kidney development: multilineage priming. *Development*

481　　　　　141:3093–3101.

482　　Bushati N., Smith J., Briscoe J., Watkins C. 2011. An intuitive graphical visualization

483　　　　　technique for the interrogation of transcriptome data. *Nucleic acids research*

484　　　　　39:7380–9. DOI: 10.1093/nar/gkr462.

485　　Edgar R., Domrachev M., Lash AE. 2002. Gene Expression Omnibus: NCBI gene

486　　　　　expression and hybridization array data repository. *Nucleic Acids Research*

487　　　　　30 :207–210. DOI: 10.1093/nar/30.1.207.

23

488    Gao Y., Church G. 2005. Improving molecular cancer class discovery through sparse

489         non-negative matrix factorization. *Bioinformatics* 21:3970–3975.

490    Gaujoux R., Seoighe C. 2010. A flexible R package for nonnegative matrix factorization.

491         *BMC bioinformatics* 11:367.

492    Huang DW., Sherman BT., Lempicki RA. 2008. Systematic and integrative analysis of

493         large gene lists using DAVID bioinformatics resources. *Nature protocols* 4:44–57.

494    Huang DW., Sherman BT., Lempicki RA. 2009. Bioinformatics enrichment tools: paths

495         toward the comprehensive functional analysis of large gene lists. *Nucleic acids*

496         *research* 37:1–13.

497    Jamora C., DasGupta R., Kocieniewski P., Fuchs E. 2003. Links between signal

498         transduction, transcription and adhesion in epithelial bud development. *Nature*

499         422:317–322. DOI: 10.1038/nature01458.

500    Jiao Y., Widschwendter M., Teschendorff AE. 2014. A systems-level integrative

501         framework for genome-wide DNA methylation and gene expression data identifies

502         differential gene expression modules under epigenetic control.

503         *Bioinformatics*:btu316.

504    Junker JP., Noël ES., Guryev V., Peterson KA., Shah G., Huisken J., McMahon AP.,

505         Berezikov E., Bakkers J., van Oudenaarden A. 2014. Genome-wide RNA

506         Tomography in the Zebrafish Embryo. *Cell* 159:662–675. DOI:

507         10.1016/j.cell.2014.09.038.

508    Karolchik D., Barber GP., Casper J., Clawson H., Cline MS., Diekhans M., Dreszer TR.,

24

509    Fujita PA., Guruvadoo L., Haeussler M. 2014. The UCSC genome browser

510       database: 2014 update. *Nucleic acids research* 42:D764–D770.

511    Kharchenko P V., Silberstein L., Scadden DT. 2014. Bayesian approach to single-cell

512       differential expression analysis. *Nature methods* 11:740–742.

513    Kim D., Pertea G., Trapnell C., Pimentel H., Kelley R., Salzberg SL. 2013. TopHat2:

514       accurate alignment of transcriptomes in the presence of insertions, deletions and

515       gene fusions. *Genome Biol* 14:R36.

516    Kim H., Park H. 2007. Sparse non-negative matrix factorizations via alternating non-

517       negativity-constrained least squares for microarray data analysis. *Bioinformatics*

518       23:1495–1502.

519    Kumar RM., Cahan P., Shalek AK., Satija R., DaleyKeyser AJ., Li H., Zhang J., Pardee

520       K., Gennert D., Trombetta JJ., Ferrante TC., Regev A., Daley GQ., Collins JJ. 2014.

521       Deconstructing transcriptional heterogeneity in pluripotent stem cells. *Nature*

522       516:56–61. DOI: 10.1038/nature13920.

523    Leinonen R., Sugawara H., Shumway M., Collaboration  on behalf of the INSD. 2011.

524       The Sequence Read Archive. *Nucleic Acids Research* 39:D19–D21. DOI:

525       10.1093/nar/gkq1019.

526    Li Y., Ngom A. 2013. The non-negative matrix factorization toolbox for biological data

527       mining. *Source code for biology and medicine* 8:1–15.

528    Liao Y., Smyth GK., Shi W. 2014. featureCounts: an efficient general purpose program

529       for assigning sequence reads to genomic features. *Bioinformatics* 30:923–930.

25

530     Love MI., Huber W., Anders S. 2014. Moderated estimation of fold change and

531         dispersion for RNA-Seq data with DESeq2. *bioRxiv*.

532     van der Maaten L. 2013. Barnes-hut-sne. *arXiv preprint arXiv:1301.3342*.

533     Van der Maaten L., Hinton G. 2008. Visualizing data using t-SNE. *Journal of Machine*

534         *Learning Research* 9:85.

535     McDavid A., Finak G., Chattopadyay PK., Dominguez M., Lamoreaux L., Ma SS.,

536         Roederer M., Gottardo R. 2013. Data exploration, quality control and testing in

537         single-cell qPCR-based gene expression experiments. *Bioinformatics (Oxford,*

538         *England)* 29:461–7. DOI: 10.1093/bioinformatics/bts714.

539     Monga V., Mıhçak MK. 2007. Robust and secure image hashing via non-negative matrix

540         factorizations. *Information Forensics and Security, IEEE Transactions on* 2:376–

541         390.

542     Pan X. 2014. Single Cell Analysis: From Technology to Biology and Medicine. *Single*

543         *cell biology* 3:106. DOI: 10.4172/2168-9431.1000106.

544     Patel AP., Tirosh I., Trombetta JJ., Shalek AK., Gillespie SM., Wakimoto H., Cahill DP.,

545         Nahed B V., Curry WT., Martuza RL., Louis DN., Rozenblatt-Rosen O., Suvà ML.,

546         Regev A., Bernstein BE. 2014. Single-cell RNA-seq highlights intratumoral

547         heterogeneity in primary glioblastoma. *Science* 344:1396–1401. DOI:

548         10.1126/science.1254257.

549     Qi Q., Zhao Y., Li M., Simon R. 2009. Non-negative matrix factorization of gene

550         expression profiles: a plug-in for BRB-ArrayTools. *Bioinformatics* 25:545–547.

26

551    Rajapakse M., Tan J., Rajapakse J. 2004. Color channel encoding with NMF for face

552        recognition. In: *Image Processing, 2004. ICIP'04. 2004 International Conference*

553        *on*. IEEE, 2007–2010.

554    Rand WM. 1971. Objective criteria for the evaluation of clustering methods. *Journal of*

555        *the American Statistical association* 66:846–850.

556    Reichardt J., Bornholdt S. 2006. Statistical mechanics of community detection. *Physical*

557        *Review E* 74:16110.

558    Robinson MD., McCarthy DJ., Smyth GK. 2010. edgeR: a Bioconductor package for

559        differential expression analysis of digital gene expression data. *Bioinformatics*

560        26:139–140.

561    Schlitzer A., Sivakamasundari V., Chen J., Sumatoh HR Bin., Schreuder J., Lum J.,

562        Malleret B., Zhang S., Larbi A., Zolezzi F. 2015. Identification of cDC1-and cDC2-

563        committed DC progenitors reveals early lineage priming at the common DC

564        progenitor stage in the bone marrow. *Nature immunology* 16:718–728.

565    Segal E., Shapira M., Regev A., Pe'er D., Botstein D., Koller D., Friedman N. 2003.

566        Module networks: identifying regulatory modules and their condition-specific

567        regulators from gene expression data. *Nature genetics* 34:166–176.

568    Smaragdis P. 2004. Non-negative matrix factor deconvolution; extraction of multiple

569        sound sources from monophonic inputs. In: *Independent Component Analysis and*

570        *Blind Signal Separation*. Springer, 494–499.

571    Stehman S V. 1997. Selecting and interpreting measures of thematic classification

27

572      accuracy. *Remote sensing of Environment* 62:77–89.

573      Tamayo P., Scanfeld D., Ebert BL., Gillette MA., Roberts CWM., Mesirov JP. 2007.

574      Metagene projection for cross-platform, cross-species characterization of global

575      transcriptional states. *Proceedings of the National Academy of Sciences* 104:5959–

576      5964.

577      Trapnell C., Cacchiarelli D., Grimsby J., Pokharel P., Li S., Morse M., Lennon NJ., Livak

578      KJ., Mikkelsen TS., Rinn JL. 2014. Pseudo-temporal ordering of individual cells

579      reveals dynamics and regulators of cell fate decisions. *Nature biotechnology* 32:381.

580      Treutlein B., Brownfield DG., Wu AR., Neff NF., Mantalas GL., Espinoza FH., Desai

581      TJ., Krasnow MA., Quake SR. 2014. Reconstructing lineage hierarchies of the distal

582      lung epithelium using single-cell RNA-seq. *Nature* 509:371–375.

583      Usoskin D., Furlan A., Islam S., Abdo H., Lönnerberg P., Lou D., Hjerling-Leffler J.,

584      Haeggström J., Kharchenko O., Kharchenko P V. 2014. Unbiased classification of

585      sensory neuron types by large-scale single-cell RNA sequencing. *Nature*

586      *neuroscience* 18:145–153.

587      Yang Z., Zhang H., Yuan Z., Oja E. 2011. Kullback-Leibler divergence for nonnegative

588      matrix factorization. In: *Artificial Neural Networks and Machine Learning–ICANN*

589      *2011*. Springer, 250–257.

590      Yuan Z., Oja E. 2005. Projective nonnegative matrix factorization for image compression

591      and feature extraction. In: *Image Analysis*. Springer, 333–342.

592      Zhang J., Li L. 2005. BMP signaling and stem cell regulation. *Developmental Biology*

28

593        284:1–11. DOI: 10.1016/j.ydbio.2005.05.009.

594

## Tables

596   **Table 1. Comparison of the top 5 modules selected by FEM with seed genes**

597   **generated by NMF and other differential expression detection methods.** The other

598   compared methods include MAST, SCDE, Monocle, DESeq2 and EdgeR. GO analysis

599   was performed on each module, and the top 2 most enriched GO terms are listed along

600   with their p-values. Connectivity is computed by taking the average of the degree number

601   of all the nodes in the graph. The p-value for each module was calculated by FEM's

602   internal Monte Carlo procedure.

603

# Figure legends

**Fig. 1: The workflow of NMFEM.** The input can be either FastQ files or a raw counts table. If FastQ files are used, they are aligned using TopHat and counted using FeatureCounts (steps shown in brackets). The input or calculated raw counts table are filtered by samples and genes, converted into RPKMs using gene lengths, and normalized by samples. We then run NMF method on them to detect subpopulations, and find the feature genes separating the detected subpopulations. Finally, we feed the feature genes as seed genes in FEM, and generate PPI gene modules that contain highly differentially expressed genes.

**Fig. 2: Comparisons among clustering methods on the HSC vs. MPP1 scRNA-Seq data.**

(A) The PCA scatter-plots of the samples, based on their log normalized expression level. Colors indicate the most favorable labeling that can be assigned to the clustering result generated by each method. The correctly and incorrectly labeled samples are marked by dot (•) and cross (x), respectively. Confusion matrices of the methods in comparison are inserted on the top-right corner of each sub-panel. The closer the matrix is to a diagonal matrix, the more accurate the method is. (B) The scatter-plots of the samples for K-means and hierarchical clustering, after t-SNE based dimension reduction. (C) Rand measures of the methods in comparison, before and after t-SNE. Rand measure ranges from 0 to 1, where a higher value indicates a greater clustering accuracy.

**Fig. 3: MA-plots of significant or important genes defined by different methods.**
Shown are scRNA-Seq data in the mouse lung distal epithelial cell E14.5 vs. E16.5

30

626    samples. The blue color highlights the genes selected as "the most significant" by the

627    corresponding methods. X-axis (A-value) is the mean of the gene expression, and y-axis

628    (M-value) is the difference of the gene expression between E16.5 and E14.5 stages.


629    **Fig. 4: Network of top 5 modules using the seed genes generated by NMF.**

630    Shown are module detection results in the FEM package, using the top 500 most

631    important genes detected by NMF in Fig. 3. ScRNA-Seq data in the mouse lung distal

632    epithelial cell E14.5 vs. E16.5 samples are compared, where the red and green colors

633    indicate up- and down-regulation of genes in E16.5 relative to E14.5, respectively. The

634    top 5 modules are selected by the p-values calculated from the internal Monte-Carlo

635    method in the FEM package (Table 1).


636    **Fig. 5: Using NMF to identify subpopulations in a single glioblastoma tumor from**

637    **patient MGH28.**

638    (A) The consensus heat map generated from NMF. The two subpopulation clusters are

639    the evident 2 red squares, marked out by number 1 and 2. The brightness indicates the

640    confidence level of two subpopulations. (B) The PCA plot of scRNA-Seq samples from

641    patient MGH28, the discovered subpopulations are coded in red and blue colors. (C) The

642    results of KEGG/BioCarta Pathway enrichment analysis. The line of significance (to the

643    right of which meaning the FDR less than 0.05) is shown. (D) The protein interaction

644    diagram of the KEGG pathway "Pathogenic E. Coli infection". The proteins coded by the

645    genes detected by NMF are highlighted yellow, with the gene names marked below.


646


647


31

# Supporting Information

648

649   **S1 Fig.** The consensus map of NMF and K-means methods run on the HSC vs. MPP1

650   dataset. The columns and rows are samples. The brightness indicates the confidence of

651   the method to assign the samples in the same group.

652   **S2 Fig.** (A) comparison of t-SNE two-dimensional scatter-plots of the mouse dendritic

653   cell scRNA-Seq data. Colors indicate the most favorable labeling that can be assigned to

654   the clustering result generated by each method. The correctly and incorrectly labeled

655   samples are marked by dot (•) and cross (x), respectively. (B) Rand measures of the

656   methods in comparison, before and after t-SNE. Rand measure ranges from 0 to 1, where

657   a higher value indicates a greater clustering accuracy.

658   **S3 Fig. PCA plot of the mouse epithelial cell data set.** The groups that are most

659   difficult to separate (E14.5 vs. E16.5) are circled out.

660   **S4 Fig.** (A) The kernel density estimation (KDE) plot showing the frequency of log

661   expression values of "important genes" that separate E14.5 vs. E16.5, as detected by the

662   various methods in comparison. (B) KDE plot of frequency of genes appear in the 71

663   Jackknife runs. For a certain x-value (frequency), a higher y-value (density) means that a

664   higher percentage of genes appear around this frequency among the 71 runs. The blue

665   block is the top 500 genes selected by NMF and the red block is all the genes in the

666   filtered data used by NMF.

667   **S5 Fig. The heatmap of the characteristic genes (E14.5 vs. E16.5) found in common**

668   **pair-wise by the various methods.** The dendrogram at the bottom shows the hierarchical

669    clustering results using the distance measured by the inverse of the number of

670    overlapping genes.

671    **S6 Fig. Using NMF to identify subpopulations in a single glioblastoma tumor from**

672    **Patient MGH31**

673    (A) The consensus heat map generated from NMF. The two subpopulation clusters are

674    the evident 2 red squares, marked out by number 1 and 2. The brightness indicates the

675    confidence level of two subpopulations. (B) The PCA plot of scRNA-Seq samples from

676    patient MGH31, the discovered subpopulations are coded in red and blue colors. (C) The

677    results of KEGG/BioCarta Pathway enrichment analysis. The line of significance (to the

678    right of which meaning the FDR less than 0.05) is shown. (D) The protein interaction

679    diagram of the KEGG pathway "Pathogenic E. Coli infection". The proteins coded by the

680    genes detected by NMF are highlighted yellow, with the gene names marked below.

33

| seed | size | connectivity | p_values | first_term | first_fisher | second_term | second_fisher |
|------|------|-------------|----------|-----------|-------------|------------|--------------|
| | | | | NMF | | | |
| Gna13 | 32 | 4.6875 | 0.004 | G-protein coupled receptor signaling pathway | 1.80E-13 | semaphorin-plexin signaling pathway | 2.50E-13 |
| Med31 | 73 | 8.136986301 | 0.009 | stem cell maintenance | 1.40E-13 | RNA metabolic process | 1.90E-13 |
| Smad4 | 52 | 4.230769231 | 0.017 | BMP signaling pathway | 0.00012 | regulation of BMP signaling pathway | 0.00031 |
| Exosc4 | 42 | 7.857142857 | 0.022 | rRNA catabolic process | 1.10E-16 | rRNA processing | 4.70E-16 |
| Pnn | 14 | 3.857142857 | 0.023 | mRNA destabilization | 0.000028 | RNA destabilization | 0.000059 |
| | | | | MAST | | | |
| Hdac2 | 92 | 5.869565217 | 0 | chromatin organization | 6.10E-29 | negative regulation of nucleic acid-templated transcription | 1.50E-27 |
| Dld | 73 | 8.02739726 | 0.001 | carboxylic acid metabolic process | 1.80E-29 | oxoacid metabolic process | 9.00E-29 |
| Sdhb | 33 | 7.696969697 | 0.006 | aerobic respiration | 3.80E-17 | tricarboxylic acid cycle | 8.10E-17 |
| Ndufv2 | 24 | 7.666666667 | 0.008 | oxidation-reduction process | 0.000000065 | response to protozoan | 0.00024 |
| Twistnb | 46 | 13.13043478 | 0.012 | transcription from RNA polymerase III promoter | 3.70E-14 | nucleobase-containing compound biosynthetic process | 6.10E-13 |
| | | | | SCDE | | | |
| Polr2l | 75 | 12.88 | 0.002 | nucleobase-containing compound biosynthetic process | 2.50E-14 | aromatic compound biosynthetic process | 5.90E-14 |
| Ndufv2 | 24 | 7.666666667 | 0.007 | oxidation-reduction process | 0.000000065 | response to protozoan | 0.00024 |
| Sdhb | 33 | 7.696969697 | 0.008 | aerobic respiration | 3.80E-17 | tricarboxylic acid cycle | 8.10E-17 |
| Ldha | 33 | 7.696969697 | 0.01 | aerobic respiration | 3.80E-17 | tricarboxylic acid cycle | 8.10E-17 |
| Polr2b | 79 | 10.75949367 | 0.014 | nucleobase-containing compound biosynthetic process | 2.60E-18 | transcription, DNA-templated | 4.50E-18 |
| | | | | Monocle | | | |
| Hdac2 | 92 | 5.869565217 | 0 | chromatin organization | 6.10E-29 | negative regulation of nucleic acid-templated transcription | 1.50E-27 |
| Rabgap1 | 10 | 8 | 0.005 | single-organism catabolic process | 0.0014 | cellular catabolic process | 0.0017 |
| Sdhb | 33 | 7.696969697 | 0.006 | aerobic respiration | 3.80E-17 | tricarboxylic acid cycle | 8.10E-17 |
| Twistnb | 46 | 13.13043478 | 0.006 | transcription from RNA polymerase III promoter | 3.70E-14 | nucleobase-containing compound biosynthetic process | 6.10E-13 |
| Ndufv2 | 24 | 7.666666667 | 0.013 | oxidation-reduction process | 0.000000065 | response to protozoan | 0.00024 |
| | | | | DESeq2 | | | |
| Aldh6a1 | 36 | 8.111111111 | 0.003 | aerobic respiration | 1.00E-16 | tricarboxylic acid cycle | 2.00E-16 |
| Gfm2 | 10 | 8 | 0.005 | single-organism catabolic process | 0.0014 | cellular catabolic process | 0.0017 |
| Polr2l | 75 | 12.88 | 0.006 | nucleobase-containing compound biosynthetic process | 2.50E-14 | aromatic compound biosynthetic process | 5.90E-14 |
| Twistnb | 46 | 13.13043478 | 0.006 | transcription from RNA polymerase III promoter | 3.70E-14 | nucleobase-containing compound biosynthetic process | 6.10E-13 |
| Gna13 | 32 | 4.6875 | 0.008 | G-protein coupled receptor signaling pathway | 1.80E-13 | semaphorin-plexin signaling pathway | 2.50E-13 |
| | | | | EdgeR | | | |
| Aldh6a1 | 36 | 8.111111111 | 0.004 | aerobic respiration | 1.00E-16 | tricarboxylic acid cycle | 2.00E-16 |
| Gna13 | 32 | 4.6875 | 0.012 | G-protein coupled receptor signaling pathway | 1.80E-13 | semaphorin-plexin signaling pathway | 2.50E-13 |
| Tpr | 58 | 12.24137931 | 0.016 | proteolysis involved in cellular protein catabolic process | 5.00E-18 | cellular protein catabolic process | 1.30E-17 |
| Thbs1 | 16 | 3.875 | 0.017 | cell adhesion | 0.00001 | biological adhesion | 0.00001 |
| Por | 12 | 7.333333333 | 0.018 | single-organism catabolic process | 0.000018 | cellular catabolic process | 0.000058 |

1 **Table 1.Comparison of the top 5 modules selectedby FEM with seed genes**

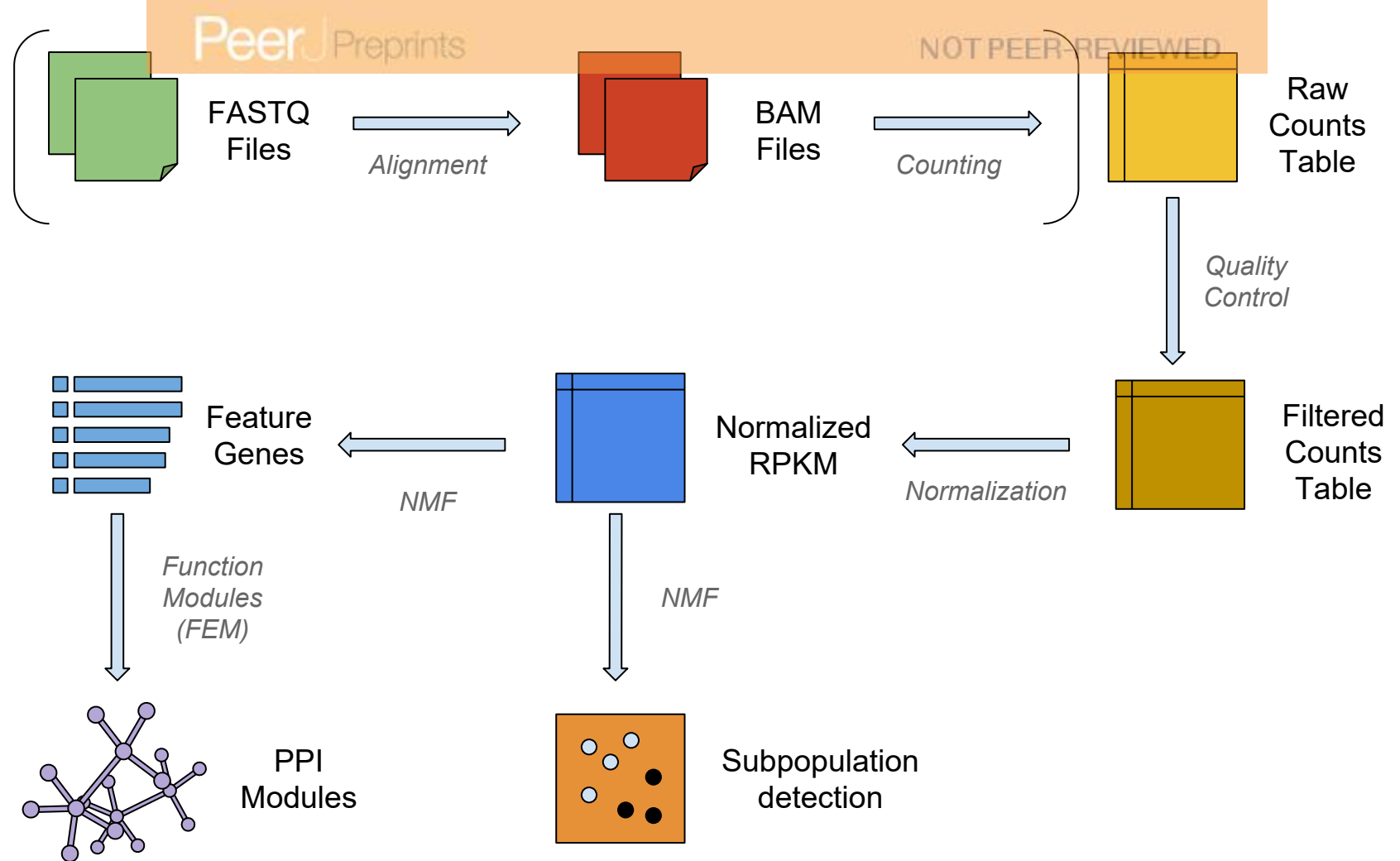2 **generated by NMF and other differential expression detection methods.**The other

3 compared methods include MAST, SCDE, Monocle, DESeq2 and EdgeR. GO analysis

4 was performed on each module, and the top 2 most enriched GO terms are listed along
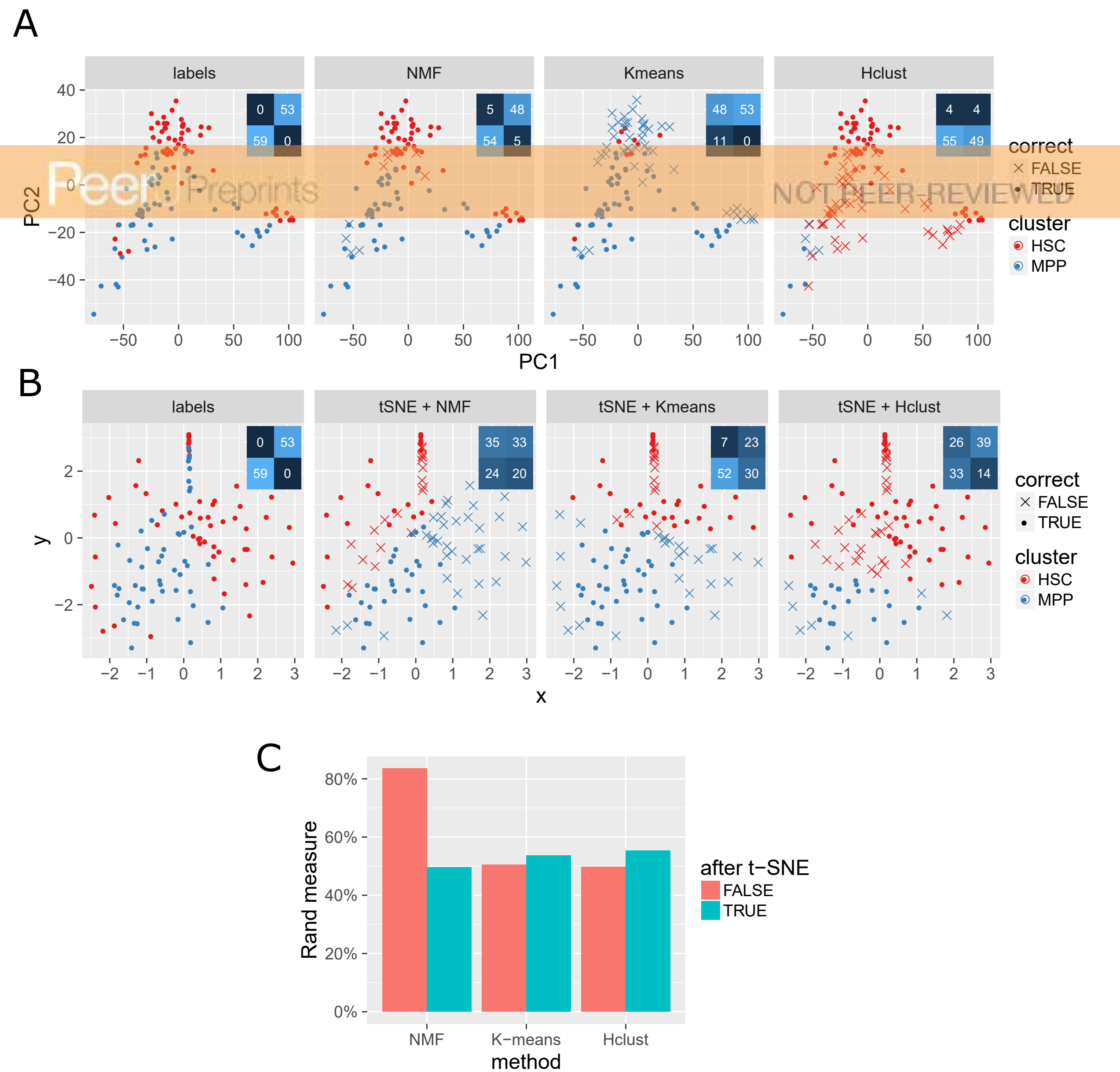
5 with their p-values. Connectivity is computed by taking the average of the degree number

6 of all the nodes in the graph. The p-value for each module was calculated by FEM's

7 internal Monte Carlo procedure.

**Fig. 1: The workflow of NMFEM.** The input can be either FastQ files or a raw counts table. If FastQ files are used, they are aligned using TopHat and counted using FeatureCounts (steps shown in brackets). The input or calculated rawcounts table are filtered by samples and genes, converted into RPKMs using gene lengths, and normalized by samples. We then run NMF method on them to detect subpopulations, and find the feature genes separating the detected subpopulations. Finally,we feed the feature genes as seed genes in FEM, and generate PPI gene modules that contain highly differentially expressed genes.
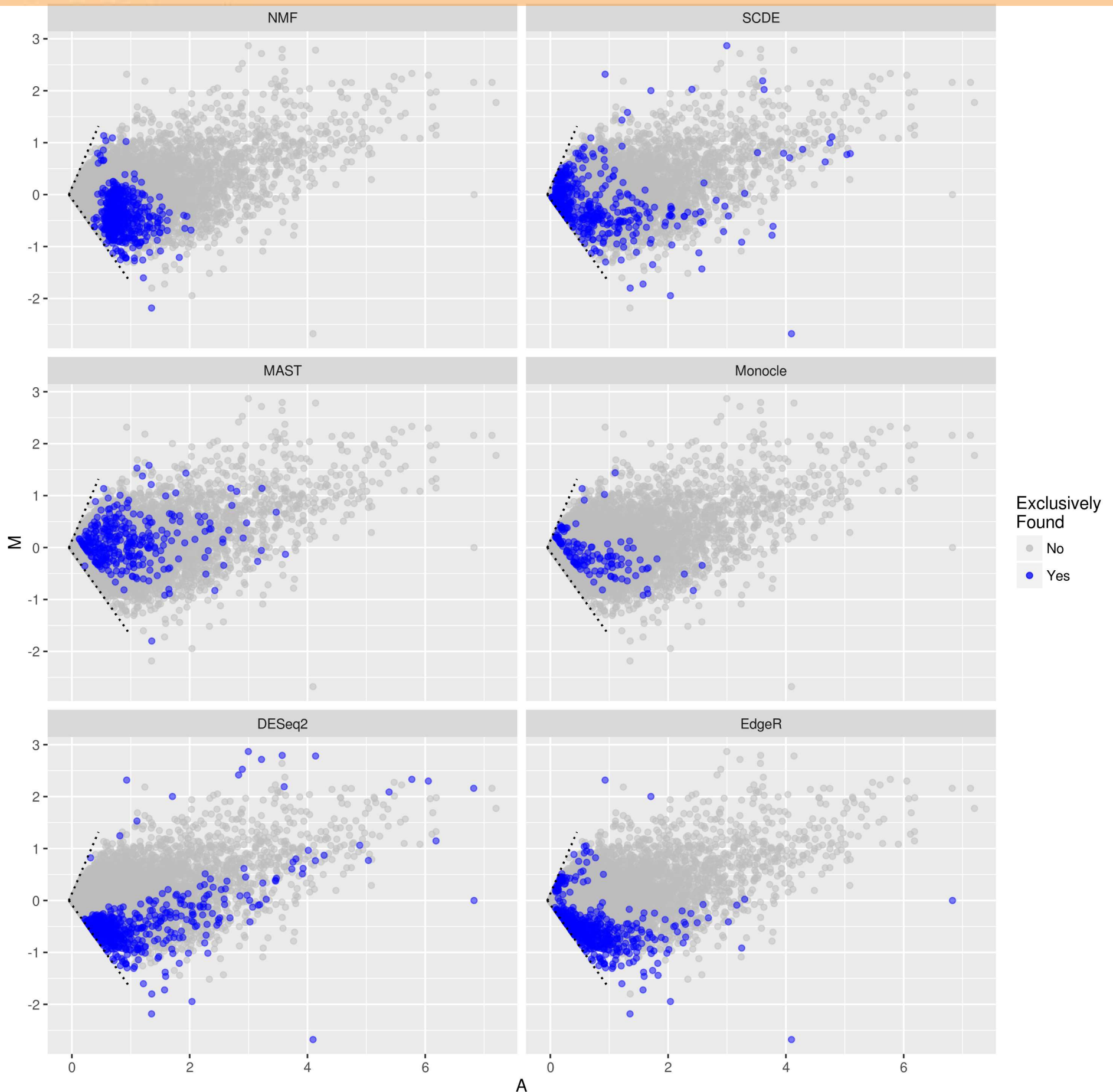
1   **Fig. 2: Comparisons among clustering methods on the HSC vs. MPP1 scRNA-Seq**

2   **data.**

3   (A) The PCA scatter-plots of the samples, based on their log normalized expression level.

4   Colors indicate the most favorable labeling that can be assigned to the clustering result

5   generated by each method. The correctly and incorrectly labeled samples are marked by

6   dot (•) and cross (x), respectively.Confusion matrices of the methods in comparison are

7   inserted on thetop-right corner of each sub-panel. The closer the matrix is to a diagonal

8   matrix, the more accurate the method is. (B) The scatter-plots of the samples for K-means

9   and hierarchical clustering, after t-SNE based dimension reduction. (C) Rand measures of

10  the methods in comparison, before and after t-SNE. Rand measure ranges from 0 to 1,

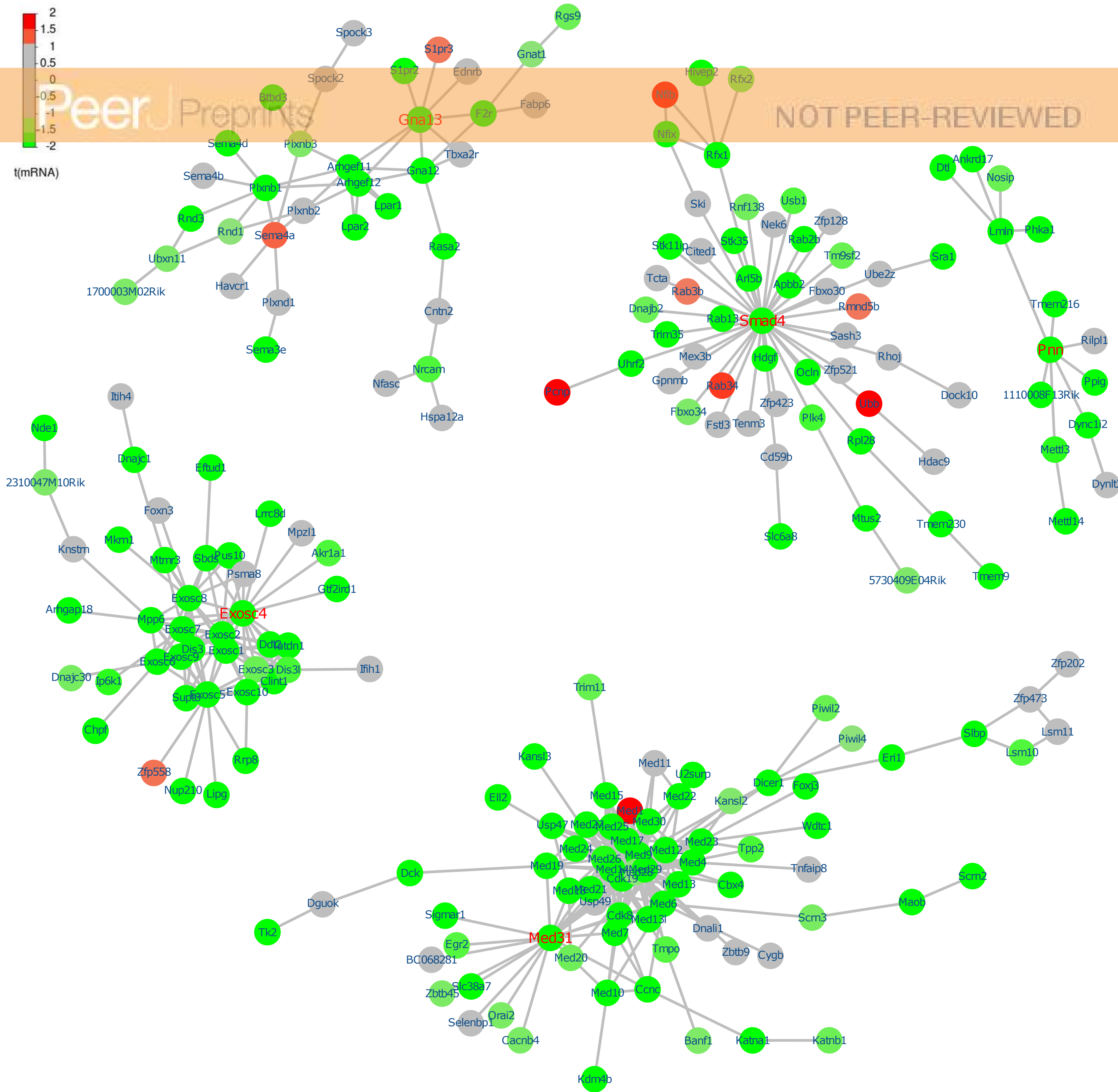11  where a higher value indicates a greater clustering accuracy.

**Fig. 3: MA-plots of significant or important genes defined by different methods.**

Shown are scRNA-Seq data in the mouse lung distal epithelial cell E14.5 vs. E16.5

samples. The blue color highlights the genes selected as "the most significant" by the

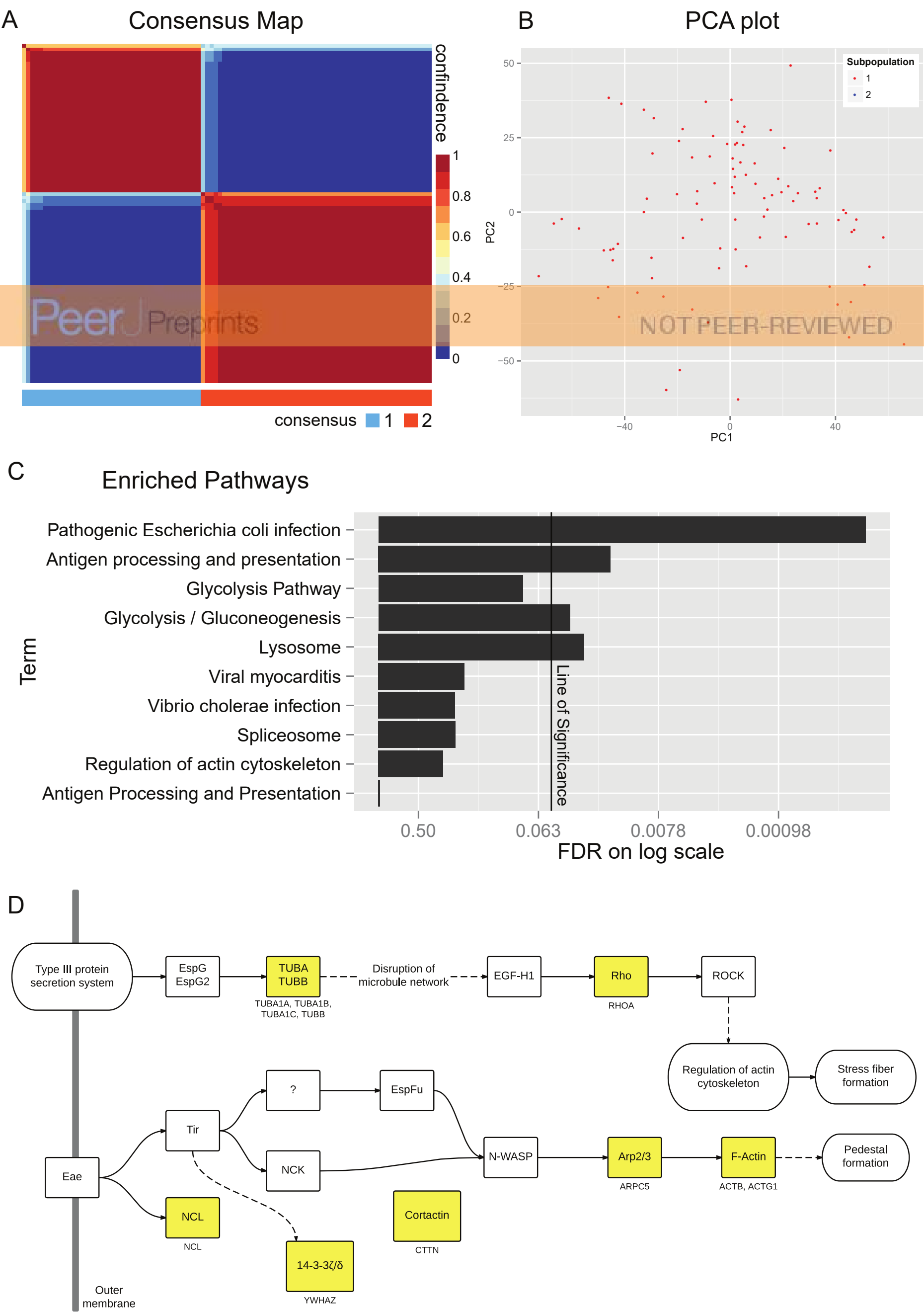corresponding methods. X-axis (A-value) is the mean of the gene expression, and y-axis

(M-value) is the difference of the gene expression between E16.5 and E14.5 stages.

1 **Fig. 4: Network of top 5 modules using the seed genes generated by NMF.**

2 Shown are module detection results in the FEM package, using the top 500 most

3 important genes detected by NMF in Fig. 3. ScRNA-Seq data in the mouse lung distal

4 epithelial cell E14.5 vs. E16.5 samples are compared, where the red and green colors

5 indicate up-and down-regulation of genes in E16.5 relative to E14.5, respectively. The

6 top 5 modules are selected by the p-values calculated from the internal Monte-Carlo

7 method in the FEM package (Table 1).

## A      Consensus Map

## B      PCA plot

## C      Enriched Pathways

## D

**Fig. 5: Using NMF to identify subpopulations in a single glioblastoma tumor from patient MGH28.**

(A) The consensusheat map generated from NMF. The two subpopulation clusters are

the evident 2 red squares, marked out by number 1 and 2. The brightness indicates the

confidence level of two subpopulations. (B) The PCA plot of scRNA-Seq samples from

patient MGH28, the discovered subpopulations are coded in red and blue colors. (C) The

results of KEGG/BioCarta Pathway enrichment analysis. The line of significance (to the

right of which meaning the FDR less than 0.05) is shown. (D) The protein interaction

diagram of the KEGGpathway "Pathogenic E. Coli infection". The proteins coded by the

genes detected by NMF are highlighted yellow, with the gene names marked below.