1 **Pitfalls and opportunities for applying PEER factors in single-cell eQTL analyses**

2

3 Angli Xue[1,2], Seyhan Yazar[1], Drew Neavin[1], Joseph E. Powell[1,3,*]

4

5 [1] Garvan-Weizmann Centre for Cellular Genomics, Garvan Institute of Medical Research, Sydney,

6 NSW, 2010, Australia

7 [2] School of Medical Sciences, University of New South Wales, Sydney, NSW, 2052, Australia

8 [3] UNSW Cellular Genomics Futures Institute, University of New South Wales, Sydney, NSW, 2052,

9 Australia

10 * Correspondence: j.powell@garvan.org.au

11

12 **Abstract**

13 Using latent variables in gene expression data can help correct spurious correlations due to

14 unobserved confounders and increase statistical power for expression Quantitative Trait

15 Loci (eQTL) detection. Probabilistic Estimation of Expression Residuals (PEER) is a widely

16 used statistical method that has been developed to remove unwanted variation and

17 improve eQTL discovery power in bulk RNA-seq analysis. However, its performance has not

18 been largely evaluated in single-cell eQTL data analysis, where it is becoming a commonly

19 used technique. Potential challenges arise due to the structure of single-cell data, including

20 sparsity, skewness, and mean-variance relationship. Here, we show by a series of analyses

21 that this method requires additional quality control and data transformation steps on the

22 pseudo-bulk matrix to obtain valid PEER factors. By using a population-scale single-cell

23 cohort (OneK1K, $N$ = 982), we found that generating PEER factors without further QC or

24 transformation on the pseudo-bulk matrix could result in inferred factors that are highly

25 correlated (Pearson's correlation $r$ = 0.626~0.997). Similar spurious correlations were also

26 found in PEER factors inferred from an independent dataset (induced pluripotent stem cells,

27 $N$ = 31). Optimization of the strategy for generating PEER factors and incorporating the

28 improved PEER factors in the eQTL association model can identify 9.0~23.1% more eQTLs or

29 1.7%~13.3% more eGenes. Sensitivity analysis showed that the pattern of change between

30 the number of eGenes detected and PEER factors fitted varied significantly for different cell

31 types. In addition, using highly variable genes (e.g., top 2000) to generate PEER factors could

32 achieve similar eGenes discovery power as using all genes but save considerable

1    computational resources (~6.2-fold faster). We provide diagnostic guidelines to improve the

2    robustness and avoid potential pitfalls when generating PEER factors for single-cell eQTL

3    association analyses.

4

5    **Keywords:** Single-cell RNA-seq; Pseudo-bulk; Latent variable; PEER factors; Normalisation;

6    eQTL mapping

1    **Background**

2    Inferring latent variables that explain the variation in the gene expression data has been an

3    essential step for expression Quantitative Trait Loci (eQTL) analyses. It can be used to

4    identify the unobserved confounding effects and potential cellular phenotypes (e.g.,

5    transcription factor or pathway activation). Standard methods of inferring latent variables

6    include principal component analysis (PCA)[1], surrogate variable analysis (SVA)[2], and

7    Probabilistic Estimation of Expression Residuals (PEER)[3,4]. PEER is a method that implements

8    a Bayesian framework to estimate the latent variables and jointly learn the contribution to

9    the gene expression variability from genotype, known factors, and hidden factors. The

10   inferred factors (i.e., PEER factors) can be applied to increase the power of eQTL discovery.

11   This method was introduced in 2010 and is widely used in bulk eQTL analyses[5-8], and

12   recently the emerging field of single-cell pseudo-bulk eQTL analysis[9-12].

13

14   As the scale of single-cell RNA-sequencing (scRNA-seq) studies has rapidly grown, eQTL

15   analyses that use pseudo-bulk analysis approaches in scRNA-seq have started to emerge.

16   Pseudo-bulk refers to the aggregation of the gene expression profiling of all cells from one

17   sample into a single pseudo-sample; thus, the data structure will be assimilated into the

18   bulk RNA-sequence data. However, due to the nature of scRNA-seq data structures, the bulk

19   expression matrix and single-cell expression matrix can be very different. There are three

20   main differences between single-cell and bulk RNA data: matrix sparsity, distribution

21   normality or skewness, and mean-variance dependency. First, since the scRNA-seq matrix is

22   sparse and most elements are zero, the pseudo-bulk gene expression matrix also contains

23   many zero values. Second, the underlying distribution of gene expression across cells

24   follows either negative binomial (NB) or zero-inflated NB (ZINB) distributions[13]. Therefore,

25   most inter-individual distributions of mean gene expression in pseudo-bulk are non-normal

26   and heavily right-skewed. Third, mean-variance dependency exists between the intra-

27   individual mean and variance due to the characteristics of the underlying distribution, and

28   such relationships could be retained in the pseudo-bulk data. These features of pseudo-bulk

29   data may violate the assumptions of the PEER method.

30

31   Consequently, the inferred PEER factors could suffer from biases or spurious correlations

32   with each other, which can lead to the problematic interpretation of the factors themselves

1    and compromise the discovery power of pseudo-bulk eQTL association. Moreover, it is

2    unclear how many PEER factors should be fitted in the eQTL association model to maximise

3    the detection power in pseudo-bulk data. Previous bulk eQTL analysis either chose a fixed

4    number[7] or a certain threshold based on the sample size[5,8]. Some studies have run

5    sensitivity tests[5,6,8], but such optimisation has not been systematically evaluated for single-

6    cell data at the population-scale level.

7

8    Here, we identify some common scenarios where pitfalls occur and how they can be

9    avoided with data-driven approaches. To help with the future application of PEER factors to

10   single-cell RNA-sequence data, we propose guidelines for the quality control and scaling of

11   the pseudo-bulk expression matrix, diagnosing and troubleshooting the inferred hidden

12   determinants, and the way to select the optimal number of PEER factors to improve the

13   eQTL discovery.

14

15   **Results and Discussion**

16   Using three independent scRNA-seq datasets, we investigated how PEER factors behave

17   under different quality control scenarios and transformations: one from peripheral blood

18   mononuclear cells (PBMCs, $N$ = 982) and the others from fibroblasts and induced

19   pluripotent stem cells (iPSCs)[10] ($N$ = 79 and 31). The PBMC data were from the OneK1K

20   cohort[12], a population-scale scRNA-seq dataset containing ~1.2 million immune cells

21   collected from 982 donors (**Methods**). This dataset was quality controlled (QC), normalised

22   and variance stabilised at the single-cell level by *sctrasnform*[14], and classified into 14 cell

23   types by *scPred*[15] (**Methods**). To construct a pseudo-bulk expression matrix for each cell

24   type, the gene expression level per individual was calculated as the intra-individual mean

25   counts across cells. We first generated PEER factors while including sex, age, and six

26   genotype PCs as covariates. Using $CD4_{NC}$ cells as an example (**Figure 1,** other cell types

27   shown in **Supp Figure 1**), we observe strong correlations among PEER factors, many of

28   which were nearly equivalent (**Figure 1A**). For instance, while most known covariates are

29   not correlated (Pearson's $r$ = -0.04 ~ 0.06, except -0.13 between PC3 and PC4; **Supp Figure**

30   **2**), the first and second PEER factors have a modest correlation (Pearson's $r$ = 0.20).

31   However, PEER factors 5-7 have pair-wise correlations equal to 1. Although the hidden

32   factor model of PEER allows for non-orthogonal components, the mean of the pair-wise

1  Pearson's *r* across the first 10 PEER factors were all larger than 0.5 in all 14 cell types,
2  suggesting that PEER factors are overfitted.
3  Additionally, we found that the variance explained by the first PEER factor was
4  overwhelmingly more significant than the rest of the PEER factors, where the latter's
5  contributions seem negligible (upper panel in **Figure 1B** and **Supp Figure** 1B). This is
6  consistent with the observation that there is a mean-variance dependency in the pseudo-
7  bulk expression level of each gene (**Figure 1C**). Hence, the highly expressed genes inherently
8  contribute much more variation than other genes. Due to the sparsity in scRNA-seq data,
9  there is a certain proportion of genes whose intra-individual expression is zero (**Figure 1C**);
10  therefore, regardless of the transformation or normalisation methods that are used, the
11  intra-individual distribution of these genes will be strongly right-skewed and violate the
12  normality assumption of PEER (see examples in **Supp Figure 3**).
13
14  To alleviate the impact of these properties, we tested different options in combinations (13
15  options in total) to generate PEER factors: (1) excluding the genes with zero expression in
16  more than a certain % of the individuals for all analyses (i.e., $\pi_0 \geq 0.9$ or 1); (2) log(x+1)
17  transformation; (3) standardisation, which scales the distribution to mean = 0 and standard
18  deviation = 1; (4) Rank-Inverse Normal Transformation (RINT); (5) Selecting the top 2,000
19  highly variable genes (HVGs, ranked by the Fano factor, e.g. variance-to-mean ratio) to
20  generate the PEER factors (**Methods**). The results show that the correlation among PEER
21  factors was still high when genes with high $\pi_0$ were excluded, and gene expression was
22  log(x+1) transformed (options #1-5, **Figure 2A**). Among options #6-11, option #7
23  (standardization + $\pi_0 \geq 0.9$ excluded) and option #11 (log(x+1) + standardization +
24  $\pi_0 \geq 0.9$ excluded) had the lowest mean pair-wise correlation between independent PFs
25  (**Figure 2A** and **Supp Figure 4**). We identified option #11 as the optimal performing
26  approach because the skewness of genes was lower than option #7 (median skewness for all
27  genes is 0.86~3.8 vs 0.90~5.12 across 14 cell types). We also tried to generate PEER factors
28  using the top 2,000 HVGs (options #12-13 in **Figure 2A-B**). The PEER factors generated from
29  all genes are highly correlated with those from the top 2,000 HVGs (**Figure 1D** and **Supp**
30  **Figure 1D**), highlighting that the HVGs could explain most of the variation that was
31  explained when using all the genes and reduce the runtime from 46.2 mins to 7.4 mins on
32  average for different cell types (**Figure 2B**).

1

2     Next, we sought to investigate how PEER factors generated from different strategies affect

3     the power of eQTL discovery. We calculated PEER factors using all genes or the top 2,000

4     HVGs (both pre-excluded genes with $\pi_0 \geq 0.9$) and compared the number of eGenes (at

5     least associated with one significant eQTL) identified when incrementally fitting PEER

6     factors as covariates from 0 to 50. Notably, the pattern of change in the number of eGenes

7     varied across different cell types (**Figure 2**). For CD4$_{NC}$ cells, the number of eGenes

8     continually increased until reaching an asymptote of around 30, while CD4$_{SOX4}$ cells reached

9     a peak between 10-15 and decreased as more factors were included. We also show that the

10    pattern of change in eGenes discovery power was consistent regardless of using all genes or

11    the top 2,000 HVGs (**Figure 2C**), and the gains of made in number of detected eGenes were

12    also similar (**Supp Table 1**). These consistencies reaffirmed that using the top 2,000 HVGs

13    captures most of the latent variation that can be explained by all genes in this dataset. We

14    also compared the number of eQTLs/eGenes when PEER factors were generated without

15    QCs or QC option #11 on the pseudo-bulk matrix. The QC option #11 can identify 9.0~23.1%

16    more eQTLs or 1.7%~13.3% more eGenes at the peak (**Supp Figure 5**). It was also clear that

17    the number of detected eGenes started to drop much earlier when incorporating highly

18    correlated PEER factors (**Supp Figure 5**). Performing these sensitivity analyses in new studies

19    is time-consuming and computationally expensive, especially for large cohorts with many

20    cell types. Our results show that using the top 2000 HVGs to generate PEER factors could

21    achieve similar power in eGenes discovery compared to using all genes while saving

22    significant computational resources (**Figure 2B**). Furthermore, the optimal number of fitted

23    PEER factors is not solely dependent on sample size but on how much variation can be

24    explained. For CD4$_{SOX4}$ cells, the inferred PEER factors did not significantly increase the

25    eGenes detection power in most scenarios (**Figure 2C** and **Supp Table 1**); therefore,

26    selecting the number of PEER factors in eQTL association just based on sample size could be

27    erroneous.

28

29    To expand our exploration into other cell types, we tested the data from Neavin *et al.*[10],

30    who noted that the number of detected eGenes dropped with the incremental increase of

31    PEER factors in the four iPSC clusters but not in the six fibroblast clusters (*Figure S20* in the

32    original paper). Strong correlations among PFs were observed in four iPSC subtypes (after

1   the 4$^{th}$ or 5$^{th}$ PEER factor) but not in fibroblast subtypes (**Supp Figure 6**). In the case of iPSC

2   subtypes, fitting more PEER factors in the eQTL association analysis added more noise,

3   which led to the loss of power. We hypothesise that the difference is due to the sample size

4   since the input expression matrices were already quality controlled using quantile

5   normalisation and z-transformation. There are rules of thumb for the minimum sample size

6   required for factor analysis[16,17], which suggest 3-20 samples per factor. When the sample

7   size is too small, the first several PFs could explain almost all the variations, and there is not

8   enough variation that the additional factors can explain. Thus, the following factors become

9   strongly correlated due to overfitting (observed as very similar or even equivalent weights

10  for certain PEER factors). The sample sizes were 79 for fibroblast and 31 for iPSC; thus, iPSC

11  is more likely to suffer from sample size bias. We validated our hypothesis by down-

12  sampling the fibroblast dataset ($N$ = 31 to match with iPSC; **Methods**). The mean of pair-

13  wise correlations among 10 PEER factors ranged from 0.11 to 0.99 in the six fibroblast

14  subtypes (**Supp Figure 7**), indicating that insufficient sample sizes could result in high

15  correlations among PEER factors even if the expression matrices were well normalised. We

16  also down-sampled the fibroblast clusters to 40 and 50 separately and found a negligible

17  correlation among inferred PEER factors when $N$ = 50 but moderate correlations (0.004-0.39)

18  when $N$ = 40, suggesting that we might need at least five samples per factor in such a

19  dataset.

20

21  Our results demonstrate that generating PEER factors requires careful consideration in

22  single-cell data. The impact of how many PEER factors are included to improve the eGenes

23  discovery power varies across different cell types. We recommend testing the correlation

24  among inferred latent variables (also with the known covariates) and conducting sensitivity

25  analysis to select the optimal number of latent variables to be incorporated in eQTL

26  mapping for each cell type. As we are moving towards the era of identifying single-cell,

27  context-dependent, and dynamic eQTL[18-20], learning latent variables directly from single-cell

28  level data[21,22] and comparing them with those from pseudo-bulk would provide insights into

29  the genetic control of gene expression at a more refined resolution.

30

31  **Conclusions**

1    Applying methods designed for bulk RNA-seq data to pseudo-bulk data could be challenging
2    as the assumptions might not be fully satisfied. This work highlights the pitfalls when
3    learning PEER factors for pseudo-bulk data and presents diagnostic guidelines of performing
4    further QC and normalization on the pseudo-bulk matrix to avoid strong and spurious
5    correlations among the inferred factors. Optimisation for the number of PEER factors
6    included in the eQTL association model should be carried out by a data-driven approach and
7    using highly variable genes to generate PEER factors could achieve similar eGenes discovery
8    power as to using all genes.
9
10    **Methods**
11    Three single-cell datasets were used in this study to explore the performance of the PEER
12    method. The OneK1K consortium[12] is a population-scale single-cell RNA-seq dataset
13    collected in Tasmania, Australia. This cohort includes 982 individuals, each with gene
14    expression profiling for ~1,000 (mean = 1297.0, standard deviation = 23.6) peripheral blood
15    mononuclear cells (PBMCs). The data was quality controlled, normalised and variance
16    stabilised by the *sctransform* method, and classified into 14 cell types (see more details in
17    ref[12]). We further identified two individuals with problematic metrics during the preliminary
18    test of PEER factors (one with a deficient number of cells and the other with abnormal cell
19    composition) and removed them in the primary analysis. The sample sizes for 14 different
20    cell types range from 795 to 980 (**Supp Table 1**). Neavin et al.[10] collected 64,018 fibroblasts
21    from 79 donors and 19,967 iPSC from 31 donors. The fibroblast data were classified into six
22    subtypes and iPSCs into four subtypes. For each subpopulation, the pseudo-bulk was
23    calculated as the mean expression per gene per individual and then quantile-normalised
24    and z-transformed.
25
26    PEER factors are latent variables that can explain the variability in gene expression. The
27    original method[3] was proposed in 2010, and the software[4] was released in 2012. We used
28    the R package 'peer' (v1.0) to generate the PEER factors for the pseudo-bulk data applying
29    max iterations = 2000 and the number of PEER factors = 50. Rank-Inverse Normal
30    Transformation (RINT) was applied to the data by the function *RankNorm()* in the R package
31    '*RNOmni*'[23]. The transformed matrix was standardised to a mean of zero with a unit
32    standard deviation per gene. For analysis using the top 2000 HVGs, a refined gene list (pre-

8

1    excluded genes with $\pi_0 > 0.9$ or mean < 0.001) was ranked by their Fano factor (variance-to-

2    mean ratio) before transformation and scaling. Note that these HVGs are not the same

3    HVGs usually defined in the QC step of the raw expression matrix for single-cell data. The

4    former indicates the genes with high mean variability across individuals, while the latter

5    shows the genes that are highly variable across cells.

6    The eQTL association analysis was performed by Matrix eQTL (v2.3)[24]. We fit sex, age, the

7    first six genotype PCs, and PEER factors as the covariates. The study only included SNPs

8    located in the *cis*-region of the gene within the 1Mb from either upstream or downstream

9    and with minor allele frequency > 5%. A local false discovery rate (LFDR) was calculated to

10    control the false-positive rate for each chromosome tested by the R package '*qvalue*'[25]. An

11    eGene was reported when at least one significant eQTL was found at LFDR < 0.05.

12

13    To investigate whether the strong correlation of PEER factors in iPSC data from Neavin *et*

14    *al.*[10] arose due to the small sample size, we randomly down-sampled the six fibroblast

15    subtypes from 79 to 31 individuals (to match the sample size of the iPSC subtypes) 30 times

16    and then generated PEER factors with these sub-samples. For each sub-sample, pair-wise

17    Pearson's correlations among 10 PEER factors were estimated. A similar down-sampling

18    analysis was also conducted for sample sizes equal to 40 and 50.
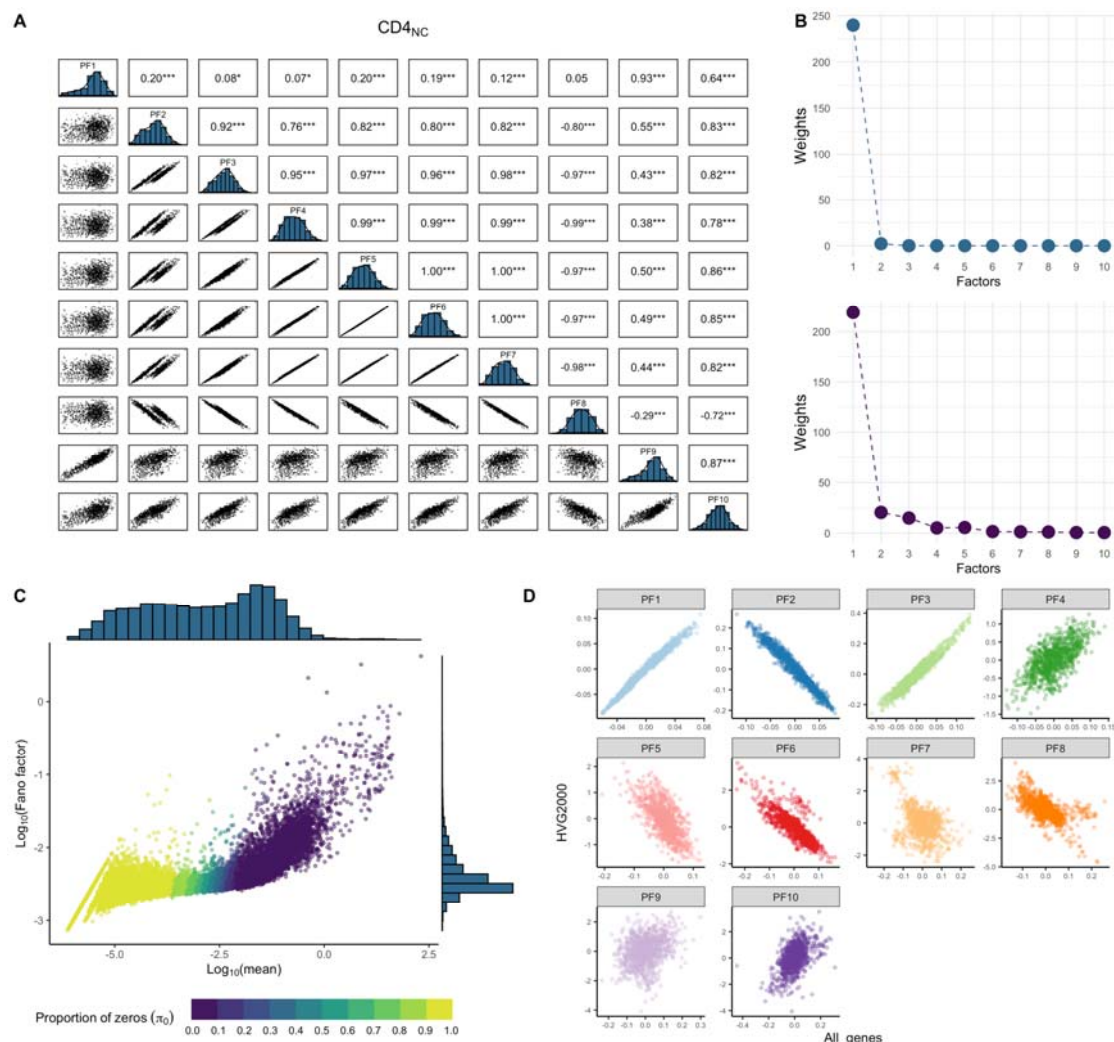
19

20    **Code availability**

21    The analysis code is available on https://github.com/powellgenomicslab/PEER_factors

22

23    **Acknowledgements**

29

**Figure 1. Correlation among inferred PEER factors and global intra-individual mean-variance dependence**. **A,** Pair-wise correlation plot among the first 10 PEER factors generated from pseudobulk expression in $CD4_{NC}$ cells without any quality control (option #1). The upper triangle panel shows the pair-wise estimates of Pearson's correlation, and the bottom triangle panel shows the pair-wise scatter plot between the PEER factors. The diagonal panel shows the distribution of each PEER factor. Significance of correlation test is annotated by * $p$-value $\leq$ 0.05, ** $\leq$ 0.01, *** $\leq$ 0.001. **B,** Diagnostic plot of the factor weights without any further quality control on the pseudo-bulk matrix (option #1, upper panel) and option #11 QC (lower panel). **C,** Relationship between intra-individual pseudo-bulk mean and Fano factor per gene. Both axes are Log10 transformed. The colour of the dots indicates the proportion of zero expression across individuals ( ) for each gene. **D,**

1    scatter plot of first 10 PEER factors generated from all genes against those from top 2000

2    highly variable genes (option #11 vs option #12).

1

2

3  **Figure 2. Performance of different QC options on generation of PEER factors and**

4  **sensitivity test for eGene detection**. **A,** The mean pair-wise correlation among the first 10

5  PEER factors. Each colour and shape represent a specific cell type. **B**, Time to generate 50

6  PEER factors by different quality control options on the pseudo-bulk matrix.  **C**, The x-axis

7  denotes the number of PEER factors fitted as covariates in the association model. The y-axis

8  represents the number of eGenes with at least one eQTL at local FDR < 0.05. The shape of

9  each scatter point indicates whether using all genes or the top 2000 highly variable genes to

1　generate PEER factors (both excluded genes with $\pi_0 \geq 0.9$, log(x+1) transformed and

2　standardised).

3

4　**References**

5

6　1.　Alter, O., Brown, P.O. & Botstein, D. Singular value decomposition for genome-wide

7　　　expression data processing and modeling. *Proc Natl Acad Sci U S A* **97**, 10101-6

8　　　(2000).

9　2.　Leek, J.T. & Storey, J.D. Capturing heterogeneity in gene expression studies by

10　　　surrogate variable analysis. *PLoS Genet* **3**, 1724-35 (2007).

11　3.　Stegle, O., Parts, L., Durbin, R. & Winn, J. A Bayesian framework to account for

12　　　complex non-genetic factors in gene expression levels greatly increases power in

13　　　eQTL studies. *PLoS Comput Biol* **6**, e1000770 (2010).

14　4.　Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of

15　　　expression residuals (PEER) to obtain increased power and interpretability of gene

16　　　expression analyses. *Nature Protocols* **7**, 500-507 (2012).

17　5.　GTEx Consortium. Genetic effects on gene expression across human tissues. *Nature*

18　　　**550**, 204-213 (2017).

19　6.　Bonder, M.J. *et al.* Identification of rare and common regulatory variants in

20　　　pluripotent cells using population-scale transcriptomics. *Nature Genetics* **53**, 313-+

21　　　(2021).

22　7.　Steinberg, J. *et al.* A molecular quantitative trait locus map for osteoarthritis. *Nat*

23　　　*Commun* **12**, 1309 (2021).

24　8.　Ota, M. *et al.* Dynamic landscape of immune cell-specific gene regulation in immune-

25　　　mediated diseases. *Cell* **184**, 3006-3021 e17 (2021).

26　9.　Orozco, L.D. *et al.* Integration of eQTL and a Single-Cell Atlas in the Human Eye

27　　　Identifies Causal Genes for Age-Related Macular Degeneration. *Cell Rep* **30**, 1246-

28　　　1259 e6 (2020).

29　10.　Neavin, D. *et al.* Single cell eQTL analysis identifies cell type-specific genetic control

30　　　of gene expression in fibroblasts and reprogrammed induced pluripotent stem cells.

31　　　*Genome Biol* **22**, 76 (2021).

11.  Cuomo, A.S.E. *et al.* Optimizing expression quantitative trait locus mapping workflows for single-cell studies. *Genome Biol* **22**, 188 (2021).

12.  Yazar, S. *et al.* Single-cell eQTL mapping identifies cell type–specific genetic control of autoimmune disease. *Science* **376**, eabf3041 (2022).

13.  Sarkar, A. & Stephens, M. Separating measurement and expression models clarifies confusion in single-cell RNA sequencing analysis. *Nature Genetics* **53**(2021).

14.  Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biology* **20**, 1--15 (2019).

15.  Alquicira-Hernandez, J., Sathe, A., Ji, H.P., Nguyen, Q. & Powell, J.E. scPred: accurate supervised method for cell-type classification from single-cell RNA-seq data. *Genome Biol* **20**, 264 (2019).

16.  Mundfrom, D.J., Shaw, D.G. & Ke, T.L. Minimum sample size recommendations for conducting factor analyses. *International journal of testing* **5**, 159-168 (2005).

17.  Costello, A.B. & Osborne, J. Best practices in exploratory factor analysis: Four recommendations for getting the most from your analysis. *Practical assessment, research, and evaluation* **10**, 7 (2005).

18.  van der Wijst, M. *et al.* The single-cell eQTLGen consortium. *Elife* **9**(2020).

19.  Schmiedel, B.J. *et al.* Single-cell eQTL analysis of activated T cell subsets reveals activation and cell type-dependent effects of disease-risk variants. *Sci Immunol* **7**, eabm2508 (2022).

20.  Nathan, A. *et al.* Single-cell eQTL models reveal dynamic T cell state dependence of disease loci. *Nature* **606**, 120-128 (2022).

21.  Buettner, F., Pratanwanich, N., McCarthy, D.J., Marioni, J.C. & Stegle, O. f-scLVM: scalable and versatile factor analysis for single-cell RNA-seq. *Genome Biol* **18**, 212 (2017).

22.  Verma, A. & Engelhardt, B.E. A robust nonlinear low-dimensional manifold for single cell RNA-seq data. *BMC Bioinformatics* **21**, 324 (2020).

23.  McCaw, Z. RNOmni: Rank Normal Transformation Omnibus Test; R package version 1.0.0. (2020).

24.  Shabalin, A.A. Matrix eQTL: ultra fast eQTL analysis via large matrix operations. *Bioinformatics* **28**, 1353-8 (2012).

1    25.    Storey JD, B.A., Dabney A, Robinson D. qvalue: Q-value estimation for false discovery

2            rate control; R package version 2.20.0. (2020).

3