

Submitted to the *Annals of Applied Statistics*  
arXiv: arXiv:0000.0000

# CONTROLLING FOR CONFOUNDING EFFECTS IN SINGLE CELL RNA SEQUENCING STUDIES USING BOTH CONTROL AND TARGET GENES

BY MENGJIE CHEN<sup>\*</sup> AND XIANG ZHOU<sup>†</sup>

*University of North Carolina and University of Michigan*

Single cell RNA sequencing (scRNAseq) technique is becoming increasingly popular for unbiased and high-resolucional transcriptome analysis of heterogeneous cell populations. Despite its many advantages, scRNAseq, like any other genomic sequencing technique, is susceptible to the influence of confounding effects. Controlling for confounding effects in scRNAseq data is a crucial step for accurate downstream analysis. Here, we present a novel statistical method, which we refer to as scPLS, for robust and accurate inference of confounding effects. scPLS takes advantage of the fact that genes in a scRNAseq study can often be naturally classified into two sets: a control set of genes that are free of effects of the predictor variables and a target set of genes that are of primary interest. By modeling the two sets of genes jointly using the partial least squares regression, scPLS is capable of making full use of the data to improve the inference of confounding effects. scPLS is closely related to and bridges between two existing subcategories of methods, and enjoys robust performance across a range of application scenarios. To accompany our method, we also develop a new, block-wise expectation maximization algorithm for scalable inference. Our algorithm is an order of magnitude faster than a standard one, making scPLS applicable to hundreds of cells and hundreds of thousands of genes. With extensive simulations and comparisons with other methods, we demonstrate the effectiveness of scPLS. Finally, we apply scPLS to analyze two scRNAseq data sets to illustrate its benefits in removing technical confounding effects as well as for removing cell cycle effects.

**1. Introduction.** Single-cell RNA sequencing (scRNAseq) has emerged as a powerful tool in genomics. While the traditional RNA sequencing, known as the bulk RNAseq, measures gene expression levels averaged across many different cells in a sample of potentially heterogeneous cell population, scRNAseq can measure gene expression levels directly at the single

---

<sup>\*</sup>Departments of Biostatistics and Genetics, University of North Carolina, Chapel Hill, NC 27599. Email: mengjie@email.unc.edu.

<sup>†</sup>Department of Biostatistics, Center for Statistical Genetics, University of Michigan, Ann Arbor, MI 48109. Email: xzhousph@umich.edu.

*Keywords and phrases:* single cell RNA sequencing, partial least squares, confounding factors

cell resolution. As a result, scRNAseq is less influenced by the variation of cell type and cell composition across different samples – a major confounding in the analyses of bulk RNAseq studies. Because of this benefit and its high resolution, scRNAseq provides unprecedented insights into many basic biological questions that are previously difficult to address. For example, scRNAseq has been applied to classify novel cell subtypes [49, 55] and cellular states [16, 31], reconstruct cell lineage and quantify progressive gene expression during development [47, 46, 9, 53], perform spatial mapping and re-localization [1, 39], identify differentially expressed genes and gene expression modulars [41, 21, 27], and investigate the genetic basis of gene expression variation by detecting heterogenic allelic specific expressions [4, 8].

Like any other genomic sequencing experiment, scRNAseq studies are influenced by many factors that can introduce unwanted variation in the sequencing data and confound the down-stream analysis [43]. Due to low capture efficiency and low amount of input material, such unwanted variation are exacerbated in scRNAseq experiments [50]. Indeed, adjusting for confounding factors in scRNAseq data has been shown to be crucial for accurate estimation of gene expression levels and successful down-stream analysis [13, 20, 22, 43, 50]. However, depending on the source, adjusting for confounding factors in scRNAseq can be non-trivial. Some confounding effects, such as read sampling noise and drop-out events, are direct consequences of low sequencing-depth, which are random in nature and can be readily addressed by probabilistic modeling using existing statistical methods [20, 13, 22, 10, 36]. Other confounding effects are inherent to a particular experimental protocol and can cause amplification bias, but can be easily mitigated by using new protocols [14]. Yet other confounding effects are due to observable batches and can be adjusted for by including batch labels and technician ids as covariates or dealt with other statistical methods [18, 51]. However, many confounding factors are hidden and are difficult or even impossible to measure. Common hidden confounding factors include various technical artifacts during library preparation and sequencing, and unwanted biological confounders such as cell cycle status. These hidden confounding factors can cause systematic bias, are notoriously difficult to control for, and are the focus of the present study.

To effectively infer and control for hidden confounding factors in scRNAseq studies, we develop a novel statistical method, which we refer to as scPLS. scPLS is specifically designed in the unsupervised setting where the predictor variables are not known *a priori* (e.g. cell clustering problems). scPLS takes advantage of the fact that genes in a scRNAseq study can often

be naturally classified into two sets: a control set of genes that are free of effects of the predictor variables and a target set of genes that are of primary interest. By modeling the two sets of genes jointly using the partial least squares regression, scPLS is capable of making full use of the data to improve the inference of confounding factors. scPLS is closely related to and bridges between two existing subcategories of methods: a subcategory of methods (e.g. PCA [37, 28, 34, 44] and LMM [17, 19, 29]) that treat control and target genes in the same fashion, and another subcategory of methods (e.g. RUV [37, 15] and scLVM [6]) that use control genes alone for inferring confounding factors. By bridging between the two subcategories of methods, scPLS enjoys robust performance across a range of application scenarios. scPLS is also computationally efficient: with a new block-wise expectation maximization (EM) algorithm, it is scalable to thousands of cells and tens of thousands of genes. Using simulations and two real data applications, we show how scPLS can be used to remove confounding effects and enable accurate down-stream analysis in scRNAseq studies. Our method is implemented as a part of the Citrus project and is freely available at: <http://chenmengjie.github.io/Citrus/>.

The paper is organized as follows. In Section 2 we provide a brief review of existing statistical methods for removing confounding effects and describe how scPLS is related to and motivated from these methods. In Section 3, we provide a methodological description of the scPLS model, with inference details provided in Section 4. In Section 5 we present comparisons between scPLS and several existing methods using simulations. In Section 6, we apply scPLS to two real scRNAseq data sets to remove technical confounding effects or cell cycle effects. Finally, we conclude the paper with a summary and discussion in Section 7.

**2. Review of Previous Methods.** Many statistical methods have been developed in sequencing- and array-based genomic studies to infer hidden confounding factors and control for hidden confounding effects. Based on their targeted application, these statistical methods can be generally classified into two categories.

The first category of methods are supervised and application-specific: these methods are designed to infer the confounding factors in the presence of a *known* predictor variable, and to correct for the confounding effects without removing the effects of the predictor variable. For example, scientists are often interested in identifying genes that are differentially expressed between two pre-determined treatment conditions or that are associated with a measured predictor variable of interest. To remove the confounding effects

in these applications, methods, include SVA [28], sparse regression models [45, 54], and, more recently, RUV [12, 11], are developed. Although these application-specific methods are widely applied in many genomics studies, their usage is naturally restricted to cases where the primary variable of interest is known. The application-specific methods become inconvenient in cases where there are multiple variables of interest (e.g. in eQTL mapping problems). They also become inapplicable when the primary variable of interest is not observed (e.g. in clustering problems).

Our scPLS belongs to the second category of unsupervised methods, which are designed to infer the confounding factors without knowing or using the predictor variable of interest. Notable applications of unsupervised methods in scRNAseq studies include cell type clustering and classification [49, 55, 16, 31, 47, 46, 9, 53]. Existing unsupervised statistical methods can be further classified into two subcategories. The first subcategory of methods treat all genes in the same fashion and use all of them to infer the confounding factors. For example, the principal component analysis (PCA) or the factor model extracts the principal components or factors from all genes as surrogates for the confounding factors [37, 28, 34, 44]. The inferred factors are treated as covariates whose effects are further removed from gene expression levels before downstream analyses. Similarly, the linear mixed models (LMMs) construct a sample relatedness matrix based on all genes to capture the influence of the confounding factors [17, 19, 29]. The relatedness matrix are then included in the downstream analyses to control for the confounding effects. In contrast, the second subcategory of unsupervised methods are recently developed to take advantage of a set of control genes for inferring the confounding factors [6, 37]. These methods divide genes into two sets: a control set of genes that are known to be free of effects of interest *a priori* and a target set of genes that are of primary interest. Unlike the first subcategory, the second subcategory of methods treat the two gene sets differently in inferring the confounding factors: the confounding factors are only inferred from the control set, and are then used to remove the confounding effects in the target genes for subsequent downstream analysis. For example, scRNAseq studies often add ERCC spike-in controls prior to the PCR amplification and sequencing steps. The spike-in controls can be used to capture the hidden confounding technical factors associated with the experimental procedures, which are further used to remove technical confounding effects (e.g. reverse transcription or PCR amplification confounding effects) from the target genes [17]. Similarly, most scRNAseq studies include a set of control genes that are known to have varying expression levels across cell cycles. These cell cycle genes can be used to capture the unmeasured cell

cycle status of each cell, which are further used to remove cell cycle effects in the target genes [6]. Prominent methods in the second subcategory include the unsupervised version of RUV [37, 15] and scLVM [6].

The two subcategories of unsupervised methods use different strategies to infer the confounding factors. Therefore, these two sets of methods are expected to perform well in different settings. Specifically, the first subcategory of methods have the advantage of using information contained in all genes to accurately infer the confounding effects. However, when the predictor variable of interest influences a large number of genes, then this subcategory of methods may incorrectly remove the primary effects of interest. On the other hand, the second subcategory of methods infer confounding factors only from the control genes and are thus not prone to mistakenly removing the primary effects of interest. However, these methods overlook one important fact – that the hidden confounding factors not only influence the control genes but also the target genes, i.e. the exact reason that we need to remove such confounding effects in the first place. Because the confounding factors influence both control and target genes, using control genes alone to infer the confounding factors can be suboptimal as it fails to use the information from target genes.

To more effectively infer and control for hidden confounding factors in scRNAseq studies, we develop a novel statistical method, which we refer to as scPLS. scPLS bridges between the two subcategories of unsupervised methods and effectively includes each as a special case. Like the first subcategory of methods, scPLS models both control and target genes jointly to infer the confounding factors. Like the second subcategory of methods, scPLS is capable of taking advantage of a control set to guide the inference of confounding factors. scPLS builds upon the partial least squares regression model and relies on a key modeling assumption that only target genes contain the primary effects of interest or other systematic biological variations. By incorporating such systematic variations in the target genes only, we can jointly model both control and target genes to infer the confounding effects while avoiding mis-removing the primary effects of interest. Therefore, scPLS has the potential to make full use of the data to improve the inference of confounding factors and the removal of confounding effects.

**3. Details of scPLS.** We provide modeling details for scPLS here. Our scPLS is generally applicable to both sequencing- and array-based genomic studies, but we focus on its application in scRNAseq. The scRNAseq data resembles that of the bulk RNAseq data and consists of a gene expression matrix on  $n$  cells and  $p + q$  genes. We consider dividing the genes into

two sets: a control set that contains  $q$  control genes and a target set that contains  $p$  genes of primary interest. The control genes are selected based on the purpose of the analysis. For example, the control set would contain ERCC spike-ins if we want to remove technical confounding factors, and would contain cell cycle genes if we want to remove cell cycle effects. We use the following partial least squares regression to jointly model both control and target genes:

$$(3.1) \quad \mathbf{x}_i = \mathbf{\Lambda}_x \mathbf{z}_i + \boldsymbol{\epsilon}_{xi}, \boldsymbol{\epsilon}_{xi} \sim \text{MVN}(0, \boldsymbol{\Psi}_{xi})$$

$$(3.2) \quad \mathbf{y}_i = \mathbf{\Lambda}_y \mathbf{z}_i + \mathbf{\Lambda}_u \mathbf{u}_i + \boldsymbol{\epsilon}_{yi}, \boldsymbol{\epsilon}_{yi} \sim \text{MVN}(0, \boldsymbol{\Psi}_{yi})$$

where for  $i$ 'th individual cell,  $\mathbf{x}_i$  is a  $q$ -vector of expression levels for  $q$  control genes;  $\mathbf{y}_i$  is a  $p$ -vector of expression levels for  $p$  target genes;  $\mathbf{z}_i$  is  $k_z$ -vector of unknown confounding factors that affect both control and target genes; the coefficients of the confounding factors are represented by the  $q$  by  $k_z$  loading matrix  $\mathbf{\Lambda}_x$  for the control genes and the  $p$  by  $k_z$  loading matrix  $\mathbf{\Lambda}_y$  for the target genes;  $\mathbf{u}_i$  is a  $k_u$ -vector of unknown factors in the target genes and potentially represents the predictors of interest or other structured variations (see below);  $\mathbf{\Lambda}_u$  is a  $p$  by  $k_u$  loading matrix;  $\boldsymbol{\epsilon}_{xi}$  is a  $q$ -vector of idiosyncratic error with covariance  $\boldsymbol{\Psi}_{xi} = \text{diag}(\sigma_{x1}^2, \dots, \sigma_{xq}^2)$ ;  $\boldsymbol{\epsilon}_{yi}$  is a  $p$ -vector of idiosyncratic error with covariance  $\boldsymbol{\Psi}_{yi} = \text{diag}(\sigma_{y1}^2, \dots, \sigma_{yp}^2)$ ; MVN denotes the multivariate normal distribution. We assume  $\mathbf{z}_i \sim \text{MVN}(0, \mathbf{I})$  and  $\mathbf{u}_i \sim \text{MVN}(0, \mathbf{I})$ . We model transformed data instead of the raw read counts. We also assume that the expression levels of each gene have been centered to have mean zero, which allows us to ignore the intercept.

scPLS includes two types of unknown latent factors. The first set of factors,  $\mathbf{z}_i$ , represents the unknown confounding factors that affect both control and target genes. The effects of  $\mathbf{z}_i$  on the control and target genes are captured in the loading matrices  $\mathbf{\Lambda}_x$  and  $\mathbf{\Lambda}_y$ , respectively. We call  $\mathbf{z}_i$  the confounding factors throughout the text, and we aim to remove the confounding effects  $\mathbf{\Lambda}_y \mathbf{z}_i$  from the target genes. The second set of factors,  $\mathbf{u}_i$ , aims to capture a low dimensional structure of the expression level of  $p$  target genes. The factors  $\mathbf{u}_i$  can represent the unknown predictor variables of interest, specific experimental perturbations, gene signatures or other intermediate factors that coordinately regulate a set of genes. Therefore, the factors  $\mathbf{u}_i$  can be interpreted as cell subtypes, treatment status, transcription factors or regulators of biological pathways in different studies [7, 35, 30, 3, 33]. Although  $\mathbf{u}_i$  could be of direct biological interest in many data sets, we do not explicitly examine the inferred  $\mathbf{u}_i$  here. Rather, we view modeling  $\mathbf{u}_i$  in the target genes as a way to better capture the complex variance structure there and to facilitate the precise estimation of confounding factors  $\mathbf{z}_i$ . For

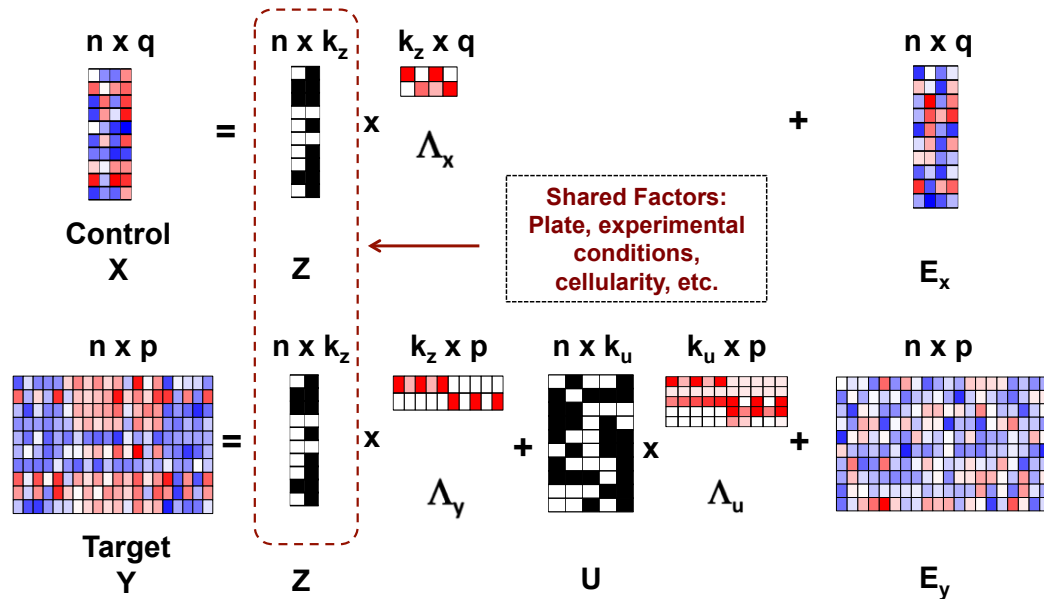


FIG 1. Illustration of scPLS. We model the expression level of genes in the control set ( $\mathbf{X}$ ) and genes in the target set ( $\mathbf{Y}$ ) jointly. Both control and target genes are affected by the common confounding factors ( $\mathbf{Z}$ ) with effects  $\Lambda_x$  and  $\Lambda_y$  in the two gene sets, respectively. The target genes are also influenced by biological factors ( $\mathbf{U}$ ) with effects  $\Lambda_u$ . The biological factors represent intermediate factors that coordinately regulate a set of genes, and are introduced to better capture the complex variance structure in the target genes.  $\mathbf{E}_x$  and  $\mathbf{E}_y$  represent residual errors. scPLS aims to remove the confounding effects  $\mathbf{Z}\Lambda_y$  in the target genes.

simplicity, we call  $\mathbf{u}_i$  the biological factors throughout the text, though we note that  $\mathbf{u}_i$  could well represent non-biological processes such as treatment or environmental effects. Thus, the expression levels of the control genes can be described by a linear combination of the confounding factors  $\mathbf{z}_i$  and residual errors; the expression levels of the target genes can be described by a linear combination of the confounding factors  $\mathbf{z}_i$ , the biological factors  $\mathbf{u}_i$  and residual errors. For both types of confounding factors, we are interested in inferring the factor effects  $\Lambda_y \mathbf{z}_i$  and  $\Lambda_u \mathbf{u}_i$  rather than the individual factors  $\mathbf{z}_i$  and  $\mathbf{u}_i$ . Therefore, unlike in standard factor models, we are not concerned with the identifiability of the factors. Figure 1 shows an illustration of scPLS.

scPLS is closely related to the two subcategories of unsupervised methods described in Section 2. Specifically, without the biological effects term  $\Lambda_u \mathbf{u}_i$ , scPLS effectively reduces to the first subcategory of methods that treat all

genes in the same fashion for inferring the confounding factors. Without the Equation 3.2 term, scPLS effectively reduces to the second subcategory of methods that use only control genes for inference. (Note that, after inferring the confounding factors  $\mathbf{z}_i$  from Equation 3.1, the second subcategory of methods still use a reduced version of Equation 3.2 without the biological effects term  $\Lambda_u \mathbf{u}_i$  to remove the confounding effects.) By including both modeling terms, scPLS can robustly control for confounding effects across a range of scenarios. Therefore, scPLS provides a flexible modeling framework that effectively includes the two subcategories of unsupervised methods as special cases and has the potential to outperform these previous methods.

**4. EM Algorithms for scPLS.** We develop an expectation-maximization (EM) algorithm for inference in scPLS. Specifically, we first initialize the factor loading matrices  $(\Lambda_x, \Lambda_y, \Lambda_u)$  based on sequential single value decompositions on the gene expression matrices  $(\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_q), \mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_p))$  (Algorithm 1). Afterwards, we treat the latent factors  $(\mathbf{w}_i = (\mathbf{z}_i^T, \mathbf{u}_i^T)^T)$  as missing data, use an iterative procedure to compute the expectation of the factors conditional on each individual cell data  $(\mathbf{v}_i = (\mathbf{x}_i^T, \mathbf{y}_i^T)^T)$  in turn in the E-step, and then update the factor loading matrices  $(\Lambda = (\Lambda_x, \Lambda_y, \Lambda_u))$  by merging information across all individuals in the M-step (Algorithm 2). We list the EM algorithm below, with detailed derivation provided in Appendix A.

---

**Algorithm 1:** Initializer of EM algorithms for scPLS

---

**Input:** Data matrices  $\mathbf{X}$ ,  $\mathbf{Y}$ , and the number of latent factors  $k_z$  and  $k_u$ .

**Output:**  $\Lambda^{(0)}$ , the initial value for  $\Lambda$ .

Apply SVD on  $\mathbf{X}$ , obtain  $\mathbf{U}, \mathbf{D}, \mathbf{V}$ ;

Calculate  $\mathbf{Z} = \mathbf{U}_{(k_u)} \mathbf{D}_{(k_z)}^{1/2}$  and standardize the elements in  $\mathbf{Z}$  to have mean 0 and variance 1;

Use least squares to estimate  $\Lambda_x^{(0)} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X}$  and

$$\Lambda_y^{(0)} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{Y};$$

Obtain the residuals of  $\mathbf{X}$  after removing the confounding effects, or

$$\mathbf{R} = \mathbf{X} - \mathbf{Z} \Lambda_x^{(0)};$$

Similarly, apply SVD on  $\mathbf{R}$ , obtain  $\mathbf{U}', \mathbf{D}', \mathbf{V}'$ ;

Calculate  $\mathbf{S} = \mathbf{U}'_{(k_u)} \mathbf{D}'_{(k_u)}{}^{1/2}$  and standardize elements in  $\mathbf{S}$  so that all elements have mean 0 and variance 1;

Use least squares to estimate  $\Lambda_u^{(0)} = (\mathbf{S}^T \mathbf{S})^{-1} \mathbf{S}^T \mathbf{R}$ ;

---



---

**Algorithm 2:** Naive EM algorithm for scPLS

---

**Input:** Data  $\mathbf{w}$ .

**Output:**  $\hat{\mathbf{v}}, \hat{\Lambda}$ .

Initialize  $\Lambda^{(0)}$  using Algorithm 1 ;

Initialize  $\psi^{(0)} = \mathbf{I}$  ;

**E step:** Compute  $E(\mathbf{v}_i|\mathbf{w}_i)^{(t)}$  and  $E(\mathbf{v}_i\mathbf{v}_i^T|\mathbf{w}_i)^{(t)}$ , given  $\Lambda^{(t)}, \psi^{(t)}$ ;

**M step:**

$$(\Lambda_x^T)^{(t+1)} = \left( \sum_{i=1}^n \mathbf{x}_i (E(\mathbf{z}_i|\mathbf{w}_i)^T)^{(t)} \right) \left( \sum_{i=1}^n E(\mathbf{z}_i\mathbf{z}_i^T|\mathbf{w}_i)^{(t)} \right)^{-1} ;$$

$$(\Lambda_y^T)^{(t+1)} = \left( \sum_{i=1}^n \mathbf{y}_i (E(\mathbf{z}_i|\mathbf{w}_i)^T)^{(t)} - \sum_{i=1}^n (\Lambda_u^T)^{(t)} E(\mathbf{u}_i\mathbf{z}_i^T|\mathbf{w}_i)^{(t)} \right) \left( \sum_{i=1}^n E(\mathbf{z}_i\mathbf{z}_i^T|\mathbf{w}_i)^{(t)} \right)^{-1} ;$$

$$(\Lambda_u^T)^{(t+1)} = \left( \sum_{i=1}^n \mathbf{y}_i (E(\mathbf{u}_i|\mathbf{w}_i)^T)^{(t)} - \sum_{i=1}^n (\Lambda_y^T)^{(t+1)} E(\mathbf{z}_i\mathbf{u}_i^T|\mathbf{w}_i)^{(t)} \right) \left( \sum_{i=1}^n E(\mathbf{u}_i\mathbf{u}_i^T|\mathbf{w}_i)^{(t)} \right)^{-1} ;$$

$$\Lambda^{(t+1)} = \begin{pmatrix} \Lambda_x^{(t+1)} & \Lambda_y^{(t+1)} \\ \mathbf{0} & \Lambda_u^{(t+1)} \end{pmatrix} ;$$

$$\psi^{(t+1)} = \frac{1}{n} \text{diag} \left\{ \sum_{i=1}^n (\mathbf{w}_i\mathbf{w}_i^T - (\Lambda_x^T)^{(t+1)} E(\mathbf{v}_i|\mathbf{w}_i)^{(t)} \mathbf{w}_i^T) \right\} ;$$

Stop when  $\|(\Lambda^T)^{(t+1)} \Lambda^{(t+1)} - (\Lambda^T)^{(t)} \Lambda^{(t)}\|_F^2$  is below a threshold;

---

We refer to the above algorithm (Algorithm 2) as the naive EM algorithm. The naive EM algorithm is computationally expensive: it scales quadratically with the number of genes and linearly with the number of cells/samples. To improve the computational speed, we develop a new EM-in-chunks algorithm (Algorithm 3). Our algorithm is based on the observation that the expression levels of the target genes are determined by the same set of underlying factors and that these factors can be estimated accurately even with a small subset set of target genes. This allows us to randomly divide target genes into dozens of chunks, compute the expectation of the factors in each chunk separately in the E-step, and then average these expectations across chunks. With the averaged expectations, we then update the factor loading matrices in the M-step. Thus, our new algorithm modifies the E-step in the naive algorithm and becomes  $K$  times faster than the naive one, where  $K$  is the number of chunks. Simulations (detailed in Section 5) show that our EM-in-chunks algorithm yields almost comparable results to the naive EM algorithm with respect to estimation errors, but can be close to an order of magnitude faster (Table 1). Therefore, we apply the EM-in-chunks algorithm

with chunk size 500 throughout the rest of the paper.

---

**Algorithm 3:** EM-in-chunks algorithm for scPLS

---

**Input:** Data  $W$ .

**Output:**  $\hat{V}$ ,  $\hat{\Lambda}$ .

Initialize  $\Lambda^{(0)}$  using Algorithm 2 ;

Initialize  $\psi^{(0)} = \mathbf{I}$  ;

Initialize  $E(\mathbf{v}_i|\mathbf{w}_i)^{(0)}$  and  $E(\mathbf{v}_i\mathbf{v}_i^T|\mathbf{w}_i)^{(0)}$  using E step in Algorithm 1 ;

**M step:**

$$((\Lambda_x^k)^T)^{(t+1)} = \left( \sum_{i=1}^n \mathbf{x}_i (E(\mathbf{z}_i|\mathbf{w}_i)^T)^{(t)} \right) \left( \sum_{i=1}^n E(\mathbf{z}_i\mathbf{z}_i^T|\mathbf{w}_i)^{(t)} \right)^{-1} ;$$

$$((\Lambda_y^k)^T)^{(t+1)} = \left( \sum_{i=1}^n \mathbf{y}_i^k (E(\mathbf{z}_i|\mathbf{w}_i)^T)^{(t)} - \right.$$

$$\left. \sum_{i=1}^n ((\Lambda_u^k)^T)^{(t)} E(\mathbf{u}_i\mathbf{z}_i^T|\mathbf{w}_i)^{(t)} \right) \left( \sum_{i=1}^n E(\mathbf{z}_i\mathbf{z}_i^T|\mathbf{w}_i)^{(t)} \right)^{-1} ;$$

$$((\Lambda_u^k)^T)^{(t+1)} = \left( \sum_{i=1}^n \mathbf{y}_i^k (E(\mathbf{u}_i|\mathbf{w}_i)^T)^{(t)} - \right.$$

$$\left. \sum_{i=1}^n ((\Lambda_y^k)^T)^{(t+1)} E(\mathbf{z}_i\mathbf{u}_i^T|\mathbf{w}_i)^{(t)} \right) \left( \sum_{i=1}^n E(\mathbf{u}_i\mathbf{u}_i^T|\mathbf{w}_i)^{(t)} \right)^{-1} ;$$

$$(\Lambda^k)^{(t+1)} = \begin{pmatrix} (\Lambda_x^k)^{(t+1)} & (\Lambda_y^k)^{(t+1)} \\ \mathbf{0} & (\Lambda_u^k)^{(t+1)} \end{pmatrix} ;$$

$$(\psi_x^k)^{(n+1)} = \frac{1}{n} \text{diag} \{ \sum_{i=1}^n (\mathbf{w}_i\mathbf{w}_i^T - ((\Lambda_x^k)^T)^{(t+1)} E(\mathbf{v}_i|\mathbf{w}_i)^{(t)} \mathbf{w}_i^T) \} ;$$

**E step;**

**for**  $k = 1$  **to**  $K$  **do**

    | Compute  $E(\mathbf{z}_i^k|\mathbf{z}_i^k)$  and  $E(\mathbf{z}_i^k(\mathbf{z}_i^k)^T|\mathbf{w}_i^k)$ , given  $\Lambda^k, \psi^k$ ;

**end**

Average among  $K$  chunks and obtain  $E(\mathbf{z}_i|\mathbf{w}_i) = \frac{1}{K} E(\mathbf{z}_i^k|\mathbf{w}_i^k)$ ,

$$E(\mathbf{z}_i\mathbf{z}_i^T|\mathbf{w}_i) = \frac{1}{K} E(\mathbf{z}_i^k(\mathbf{z}_i^k)^T|\mathbf{w}_i^k);$$

Iterate between M and E step until last cycle;

Given  $E(\mathbf{z}_i|\mathbf{w}_i)$  and  $E(\mathbf{z}_i\mathbf{z}_i^T|\mathbf{w}_i)$  from the last cycle, the final estimate of  $\Lambda$  and  $\psi$  are calculated using one M step in Algorithm 1 ;

---

Finally, we use the Bayesian information criterion (BIC) to determine the number of confounding factors  $k_z$  and the number of biological factors  $k_u$ . Specifically, we evaluate the likelihood on a grid of  $k_z$  (1 to 3) and  $k_u$  values (1 to 10) and choose the optimal combination that minimizes the BIC. After estimating the model parameters on the optimal set of  $k_z$  and  $k_u$ , we use the residuals  $\hat{\mathbf{y}}_i = \mathbf{y}_i - \hat{\Lambda}_y \hat{\mathbf{z}}_i$  as the de-noised values for subsequent analysis. Note that the residuals are only free of the confounding effects  $\Lambda_y \mathbf{z}_i$  but still contain the biological effects  $\Lambda_u \mathbf{u}_i$ .

**5. Simulations.** We performed a simulation study to compare scPLS with other methods. Specifically, we simulated gene expression levels for 50

TABLE 1

Comparison of the naive EM algorithm and the EM-in-chunks algorithm in terms of accuracy and speed. The EM-in-chunks algorithm uses either a chunk size of 500 genes or a chunk size of 1,000 genes. Accuracy is measured by the estimation error of the loading matrix in terms of the normalized Frobenius norm (i.e.  $\sqrt{\|\Lambda_x - \hat{\Lambda}_x\|_F/n}$ ). Speed is measured by CPU time in minutes. Standard deviations across 10 replicates are listed inside parenthesis.  $s$ : number of genes per chunk.  $n$ : the number of cells.  $p$ : the number of genes in the target set. The number of genes in the control set is  $q = 50$  in all simulations.

$n$	$p$	Naive EM		EM-in-chunks ( $s = 1,000$ )		EM-in-chunks ( $s = 500$ )	
		Accuracy	CPU time	Accuracy	CPU time	Accuracy	CPU time
200	2000	67.29 (5.33)	66.77 (12.72)	73.32 (6.09)	31.70 (7.95)	75.6 (6.52)	15.37 (3.28)
200	4000	135.07(10.48)	190.06 (23.30)	144.00 (13.38)	61.76 (15.43)	148.57 (14.11)	26.39 (5.04)
400	2000	72.96(5.58)	107.05 (63.66)	66.98 (5.15)	63.66 (18.39)	53.48 (4.61)	26.87 (3.00)
400	4000	95.5 (7.41)	296.16 (18.55)	101.8(9.46)	121.73 (23.65)	105.05 (9.97)	39.52 (3.84)

control genes and 1,000 target genes for 200 cells. These 200 cells come from two equal-sized groups, representing two treatment conditions or two sub-cell types. Among the 1,000 target genes, only 100 of them are differentially expressed (DE) between the two groups and thus represent the signature of the two groups. The effect sizes of the DE genes were simulated from a normal distribution, and we scaled the effects further so that the group label explains a fixed percentage of phenotypic variation (PVE) in expression levels in the DE genes (ranging from 1% to 20%, with 1% increment). In addition to the group effects, we set  $k_z = 2, k_u = 5$  and simulated each element of  $\mathbf{z}_i$  and  $\mathbf{u}_i$  from a standard normal distribution. We simulated each element of  $\Lambda_x$  from  $N(-0.25, \sigma_l^2)$  and each element of  $\Lambda_y$  from  $N(0.25, \sigma_l^2)$ . Note that  $\Lambda_x$  and  $\Lambda_y$  were simulated differently to capture the fact that the effect sizes of the confounding factors could be different for control and target genes. We simulated each element of  $\Lambda_u$  from  $N(0, \sigma_b^2)$ . We simulated each element of  $\epsilon_{xi}$  and  $\epsilon_{yi}$  from a standard normal distribution. We set  $\sigma_l^2 = 0.4$  and  $\sigma_b^2 = 0.6$  to ensure that, in non-DE genes, the confounding factors  $\mathbf{z}_i$  explain 20% PVE in either the control or the target genes; the biological factors  $\mathbf{u}_i$  explain 30% PVE of the target genes; and the residual errors to explain the rest of PVE. After we simulated gene expression levels, we further converted these continuous values into count data by using a Poisson distribution: the final observation for  $i$ th cell and  $j$ th gene  $c_{ij}$  is from  $c_{ij} \sim \text{Poi}(N \exp(\mu + w_{ij}))$ , with  $w_{ij}$  being the continuous gene expression levels simulated above and  $N = 500000, \mu = \log(10/500000)$ . Note that, because of the residual errors, the resulting count data are over-dispersed with respect to a Poisson distribution.

We considered three different simulation scenarios. In scenario I, the confounding factors  $\mathbf{z}_i$  are independent of group labels. In scenario II, the con-

founding factors are correlated with group labels. To simulate correlated data, we simulated each element of  $\mathbf{z}_i$  from  $N(0, 1)$  if the corresponding sample belongs to the first group, but from  $N(-0.25, 1)$  if the corresponding sample belongs to the second group. Finally, we also considered a scenario III where there is no biological factor (i.e. data were simulated effectively under the PCA modeling assumption and all genes could be used to infer the confounding factors). We performed 10 simulation replicates for each scenario.

We compared our method to four different methods: (1) PCA and (2) LMM (implemented in GEMMA [56, 57]) all genes used to infer the confounding effects; while (3) RUVseq (version 1.2.0); which we simply refer to as RUV in the following text) and (4) scLVM (version 0.99.1) only control genes used to infer the confounding effects. We used default settings in all the above methods. We used the count data directly for RUV and used log transformed data (i.e.  $\log(c_{ij} + 1)$ ) for all other methods. For PCA and RUV, we set the number of latent factors to be the true number (i.e. 2). Such number is determined automatically by the software itself for scLVM, and is not needed for LMM. Our goal on the simulated data is twofold: we want to identify these differentially expressed genes and to classify the 200 cells into two groups. Therefore, we compared the performance of various methods based on two criteria: the power to identify the DE genes and the power to classify cells into two groups. We permuted group labels to construct an empirical null and compared methods based on either power given 5% false discovery rate (FDR) for identifying DE genes.

It is useful to point out that the three simulation scenarios are designed to highlight the hybrid nature of scPLS. In particular, in the presence of biological factors (i.e. scenarios I and II), the methods that use all genes to remove confounding effects, such as PCA and LMM, may incorrectly remove the primary effects of interest. Therefore, we would expect RUV and scLVM to outperform PCA and LMM in scenarios I and II. In the absence of biological factors (i.e. scenarios III), the methods that only use control genes to remove confounding effects, such as RUV and scLVM, may fail to utilize all information contained in the data. Therefore, we would expect PCA and LMM to outperform RUV and scLVM in scenario III. Because of the hybrid nature of scPLS, we would expect it to perform well across all scenarios.

Simulation results confirm our expectations. Specifically, in scenario I (Figure 2a), scPLS outperforms the other four methods in identifying DE genes across a range of PVEs. Among the rest of the four methods, RUV and scLVM outperform PCA and LMM. Similarly, in scenario II (Figure 2b),

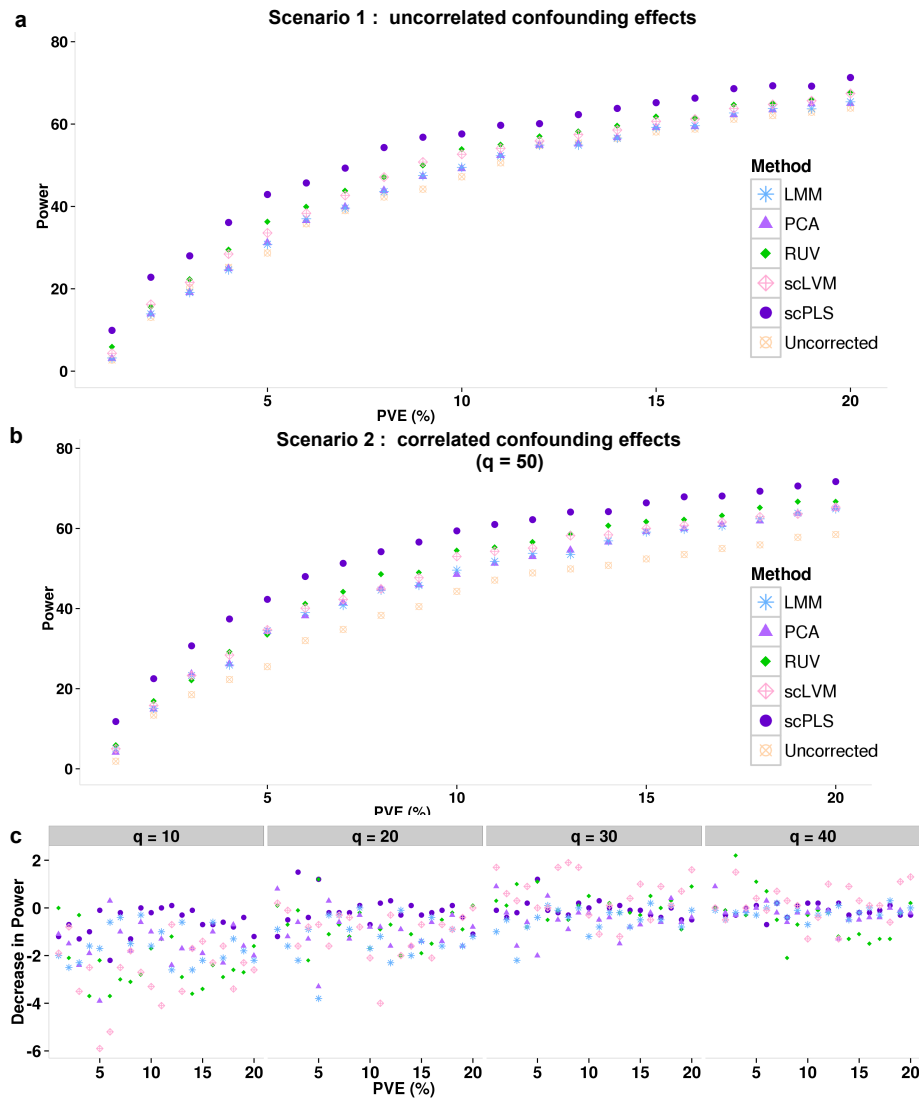


FIG 2. Method comparison in simulations. Identifying differentially expressed genes using scPLS-corrected data achieves higher power than using LMM-, PCA-, RUV- and scLVM-corrected data or uncorrected data in both scenario I (a) and scenario II (b) across a range of effect sizes. Power is evaluated at an empirical false discovery (FDR) rate of 0.05 and is averaged across ten simulation replicates. x-axis shows the effect sizes, which are measured as the percentage of phenotypic variation (PVE) in expression levels explained by the group label (ranges from 1% to 20%). (c) Sensitivity analysis shows that, compared with other methods, scPLS has the least percentage reduction in power (y-axis) when a smaller subset of control genes are used ( $q = 10, 20, 30$  or  $40$  instead of  $50$ ).

scPLS performs the best, followed by RUV and scLVM. PCA and LMM perform the worst. Compared with RUV and scLVM, scPLS is also more robust with respect to the number of control genes used in the analysis (Figure 2c). In particular, because scPLS does not completely rely on the information contained in the control genes, it achieves good performance even if we only use a much smaller subset of control genes. In contrast, the performance of RUV and scLVM compromises more quickly when a reduced number of control genes is used (especially when using  $q = 10$ ). The higher power of scPLS to detect DE genes in scenario I and II also translates to a better performance of classifying single cells (Figure 3a). To quantify the classification performance, we applied the support vector machine (SVM) to classify the cells. We performed a five-fold cross-validation, training SVM with 80% of the samples and evaluating the prediction accuracy with the rest of the samples. All methods achieve similarly high power in the easiest case when PVE is greater than 10%. However scPLS outperforms the other four methods when PVE is low and the classification task is difficult. For example, in scenario I, when  $PVE = 1\%$ , scPLS achieves an average accuracy of 77% across 10 replicates, while LMM, PCA, RUV and scLVM achieve 71.2%, 71.2%, 73.5%, 70.5%, respectively.

On the other hand, in scenario III, scPLS performs as well as PCA and LMM, and all these three methods outperform RUV and scLVM (Figure 3b). Importantly, scPLS is not sensitive with respect to the number of biological factors used in fitting the model, and achieves similar power for a range of reasonable  $k_u$  values when the truth is 0 in scenario III (Figure 3c). As it is often unknown whether a low-rank structural variation exists in a real data set, our simulation suggests that we can always include the biological factors  $\mathbf{u}_i$  in the model even in the absence of such factors.

Therefore, the simulation results highlight the hybrid nature of scPLS. scPLS works robustly well across a range of scenarios while the other two subcategories of methods work preferentially well only in scenarios that most favor their modeling assumptions.

**6. Real Data Applications.** Next, we applied scPLS to two real data sets. The first dataset is used to demonstrate the effectiveness of scPLS in removing the technical confounding effects by using ERCC spike-ins. Removing technical confounding effects is a common and important task in transcriptome analysis. The second dataset is used to demonstrate the effectiveness of scPLS in removing cell cycle effects by using a known set of cell cycle genes. Removing cell cycle effects can reveal gene expression heterogeneity that is otherwise obscured.

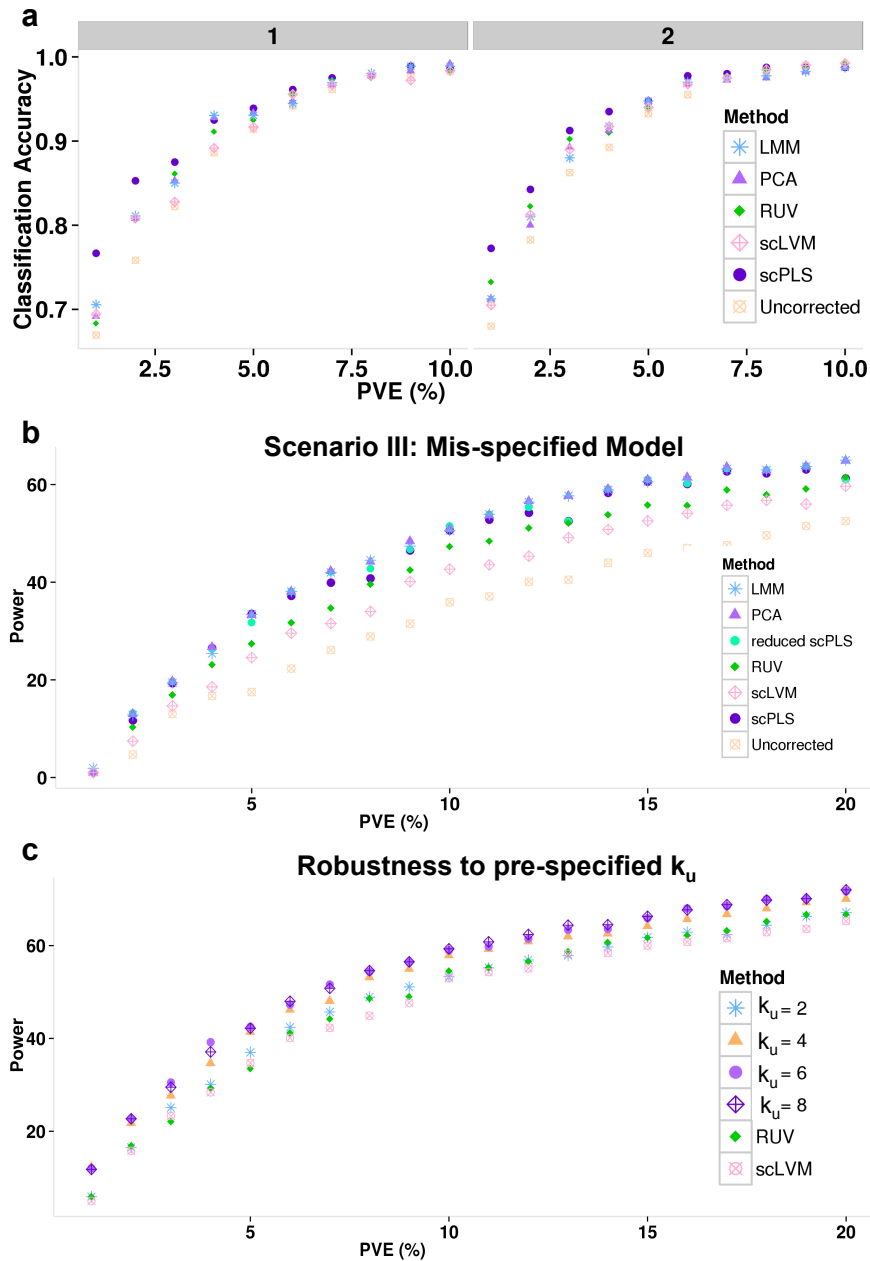


FIG 3. Method comparison in simulations (continued). (a) scPLS corrected expression data can be used to better classify cells into the two known clusters than LMM-, PCA-, RUV- and scLVM-corrected data or uncorrected data in both scenario I and scenario II across a range of effect sizes. Classification is based on support vector machine (SVM) with five-fold cross-validation. Accuracy is computed as the mean percentage of true positives in the test set across replicates. (b) Identifying differentially expressed genes using scPLS-corrected data achieves similar power as using LMM-, PCA-corrected data, all of which are more powerful than RUV- and scLVM-corrected or uncorrected data in scenario III across a range of effect sizes. (c) scPLS is robust with respect to  $k_u$ , as the power to identify differentially expressed genes remains similar when a different number of biological factors is used ( $k_u=2, 4, 6, 8$ ) in scenario III where in turn  $k_u = 0$ . x-axis shows the effect sizes, which are measured as the percentage of phenotypic variation (PVE) in expression levels explained by the group label (ranges from 1% to 20%).

6.1. *Removing Technical Confounding Factors.* The first dataset consists of 119 mouse embryonic stem cells (mESCs), including 74 mESCs cultured in a two-inhibitor (2i) medium and 45 mESCs cultured in a serum medium [13]. We obtained the raw UMI counts data directly from the authors and the data contains measurements for 92 ERCC spike-ins and 23,459 genes. Due to the low coverage of this dataset (median coverage equals one), we filtered out lowly expressed genes and selected only genes that have at least five counts in more than a third of the cells. This filtering step resulted in a total of 17 ERCC spike-ins that were used as the controls and 2,772 genes that were used as the targets. As in the simulations, we log transformed the count data and centered the transformed values for scPLS, PCA, LMM and scLVM. We used the count data for RUV. In this data, scPLS infers  $k_z = 1$  confounding factors and  $k_u = 1$  biological factors. In the target genes, the confounding factors and structured biological factors explain a median of 20% and 19% of gene expression variance, respectively. The PVE by the confounding and biological factors can be as high as 86.2% and 76.9%, respectively, in the target genes.

We applied scPLS and the other four methods to remove confounding effects in the data. To compare the performance of different methods in the real data, we performed a clustering analysis. We reason that if method is effective in removing confounding effects, then the corrected data from the method could be used to separate the mESCs into the two known clusters (i.e. 2i medium vs serum medium). For the clustering analysis, we applied the k-means method, an unsupervised method, with the number of clusters set to two, on uncorrected data and data corrected by different methods. Consistent with our simulations, scPLS outperforms all other methods based on a variety of clustering performance measurements (Table 2).

6.2. *Removing Cell Cycle Effects.* Our method can also be used to remove cell cycle effects. To demonstrate its effectiveness there, we applied scPLS and several other methods to a second dataset [6]. This dataset contains gene expression measurements on 9,570 genes from 182 embryonic stem cells (ESCs) with pre-determined cell-cycle phases (G1, S and G2M). The uncorrected data we obtained are already pre-processed by the original study to remove the technical effects and are thus continuous. Therefore, we did not apply RUV here. To remove cell cycle effects, we used 629 annotated cell-cycle genes as controls and the other genes as targets. scPLS infers  $k_z = 1$  cell cycle confounding factors, and  $k_u = 1$  biological factors. These factors explain a median of 0.4% and 0.1% of gene expression variance, respectively. The PVE by cell cycle factors and biological factors can



CONTROLLING FOR CONFOUNDING EFFECTS IN SCRNASQ 17

TABLE 2

A number of clustering measurements show that the corrected data by scPLS can be used to better reveal the two known clusters than the other four methods in the first data set.

A *k*-means algorithm (with two clusters) is applied to the uncorrected data and data corrected by different methods. Clustering performance is measured by Rand Index, sensitivity, specificity, positive predictive value (PPV), negative predictive value (NPV), mutual information (Mutual) or adjusted Rand Index (RandIndexAdj). Blue color labels the best performer by each criterion. All performance measurements are averaged across 10 runs and are multiplied by a factor of 100.

	RandIndex	Sensitivity	Specificity	PPV	NPV	Mutual	RandIndexAdj
uncorrected	74	50	50	53	48	0	0
scPLS	84	70	69	71	68	14	39
RUV	76	55	45	52	48	0	0
LMM	74	50	51	53	49	1	2
PCA	74	50	51	53	49	1	2
scLVM	73	50	50	53	48	0	1

be as high as 7% and 2%, respectively. We visualized the uncorrected data and scPLS corrected data on a PCA plot (Figure 4). In the uncorrected data, there is a clear separation of cells according to cell-cycle stage. Such separation of cells is not observed in the corrected data, indicating that the cell cycle related expression signature is effectively removed.

We compared scPLS and the other three methods in their effectiveness in removing cell cycle effects. Following the original study [6], we evaluated method performance based on the following criteria. Specifically, we computed for each gene the proportion of expression variance explained by the cell cycle factor. We denote this quantity as PVE<sub>i</sub>, which stands for inferred PVE. Because the cell-cycle stage of each cell had been experimentally determined in this data set, we further computed the variance explained by the true cell cycle labels. We denote this quantity as PVE<sub>t</sub>, which stands for true PVE. For scPLS, PVE<sub>i</sub> and PVE<sub>t</sub> are highly correlated ( $r^2 = 0.94$ ), demonstrating the efficacy of scPLS. The correlation remains the same whether we use the full control set or with a subset of 300 controls. The correlation between PVE<sub>i</sub> and PVE<sub>t</sub> in scPLS is slightly higher, with statistical significance, than scLVM ( $r^2 = 0.92$ ; p-value  $< 10^{-16}$  comparing scPLS vs scLVM), LMM ( $r^2 = 0.92$ ; p-value  $< 10^{-16}$  comparing scPLS vs LMM), and PCA ( $r^2 = 0.92$ ; p-value  $< 10^{-16}$  comparing scPLS vs PCA). In addition, as an alternative measurement, the median of the absolute difference between PVE<sub>i</sub> and PVE<sub>t</sub> across genes from scPLS, scLVM, LMM and PCA are 0.018, 0.023, 0.019 and 0.019, respectively, again supporting a small advantage of scPLS. Therefore, the results suggest that scPLS works slightly better than the other three methods, though all methods work reasonably

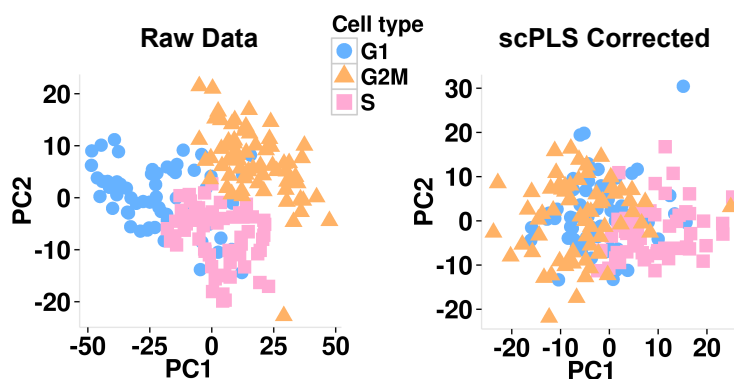


FIG 4. PCA plots for the uncorrected data and scPLS corrected data in the second dataset. In the uncorrected data, there is a clear separation of cells by cell-cycle stage. Such separation of cells is no longer observed in the scPLS corrected data.

well in removing cell cycle effects in this data set (which is consistent with the low variance explained by the confounding factors).

**7. Discussion.** We have presented scPLS for removing hidden confounding effects in scRNAseq studies. scPLS models both control and target genes jointly to infer the confounding factors and shows robust performance across a range of application scenarios. With simulations and applications to two real data sets, we have demonstrated its effectiveness for removing technical confounding effects or cell cycle effects in scRNAseq studies.

Although we have focused on its applications to scRNAseq studies, scPLS can be readily applied to other genomic sequencing studies. For instance, our method can be used to remove confounding effects from gene expression levels in bulk RNAseq studies [48] or from methylation levels in bisulfite sequencing studies [25]. The main requirement of our method is a set of pre-specified control genes that are measured together with the target genes in the sequencing studies. It is often straightforward to obtain such control genes. For example, many scRNAseq studies include a set of ERCC spike-in controls that could be used to model and remove technical confounding effects [17]. Even when such ERCC spike-in controls are not present or when they are unreliable [37], we can select a known set of house-keeping genes as controls to remove technical confounding [37]. Similarly, we can use a set of known cell cycle genes to remove cell cycle effects. Importantly, the performance of scPLS is robust to the number of genes included in the control set

and yields comparable results even when a much smaller number of control genes is used. This is because scPLS not only uses information from control genes but also relies on information from target genes. Insensitivity to the control set makes scPLS especially suited to removing confounding factors in studies where a control set is not clearly defined. Because of its effectiveness and robustness, we expect scPLS to be useful in removing confounding effects in a wide variety of sequencing studies.

One important feature of scPLS is that it includes a low-rank component to model the structured biological variation often observed in real data. By decomposing the (residual) gene expression variation into a low-rank structured component that is likely to be contributed by a sparse set of biological factors, and an unstructured component that reflects the remaining variation, scPLS can better model the residual error structure for accurate inference of confounding effects. Although here we have focused on using the biological factors to better infer the confounding effects, we note that the low-rank biology factors themselves could be of direct interest. In fact, low-rank factors inferred from many data sets using standard factor models have been linked to important biological pathways or transcription factors [7, 35, 30, 3, 33]. Inferring the biological factors using scPLS is not feasible at the moment, however: because of model identifiability, scPLS can only be used to infer the biological effects (i.e.  $\Lambda_u \mathbf{u}_i$ ) but not the biological factors (i.e.  $\mathbf{u}_i$ ). That said, additional assumptions can be made on the structure of the factors or the factor loading matrices to make factor inference possible [52]. For example, we could impose sparsity assumptions on the low-rank factors to facilitate the inference of a parsimonious set of biological factors. Exploring the use of biological factors in scPLS is an interesting avenue for future research.

Like many other methods for scRNAseq [5] or bulk [24, 38] RNAseq studies, scPLS requires a data transformation step that converts the count data into quantitative expression data. Different transformation methods can affect the interpretation of the data and are advantageous in different situations [43]. Because scPLS does not rely on a particular transformation procedure, scPLS can be paired with any transformation methods to take advantage of their benefits. One potential disadvantage of scPLS is that it does not model raw count data directly. However, despite the count nature of sequencing data, it has been shown that there is often a limited advantage of modeling the raw read counts directly, at least for RNAseq studies [42, 40]. Statistical methods that convert and model the quantitative expression data have been shown to be robust [24, 38] and most large scale bulk RNAseq studies in recent years have used transformed data instead of count data

[23, 34, 2, 32]. However, we note that, unlike bulk RNAseq studies, single cell RNAseq data often come with low read depth. In low read depth cases, modeling count data while accounting for over-dispersion or dropout events in single cell RNAseq studies may have added benefits [20, 50]. Therefore, extending our framework to modeling count data [26, 58] is another promising avenue for future research.

#### APPENDIX A: EM ALGORITHMS FOR SCPLS

To derive the EM algorithm, we first integrate out the latent variables  $\mathbf{z}_i$  and  $\mathbf{u}_i$

$$(A.1) \quad P(\mathbf{x}_i | \Lambda_x, \psi_x) = MVN(0, \psi_x + \Lambda_x^T \Lambda_x),$$

$$(A.2) \quad P(\mathbf{y}_i | \Lambda_y, \Lambda_u, \psi_y) = MVN(0, \psi_y + \Lambda_y^T \Lambda_y + \Lambda_u^T \Lambda_u).$$

The latent variable  $\mathbf{y}_i$  and  $\mathbf{z}_i$  follow a joint normal distribution

$$(A.3) \quad \begin{pmatrix} \mathbf{x}_i \\ \mathbf{z}_i \end{pmatrix} \sim MVN \left( \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \psi_x + \Lambda_x^T \Lambda_x & \Lambda_x \\ \Lambda_x^T & \mathbf{I} \end{pmatrix} \right).$$

Denoting  $\Lambda = \begin{pmatrix} \Lambda_x & \Lambda_y \\ \mathbf{0} & \Lambda_u \end{pmatrix}$ ,  $\mathbf{v}_i = \begin{pmatrix} \mathbf{z}_i \\ \mathbf{u}_i \end{pmatrix}$ , and  $\psi = \begin{pmatrix} \psi_x & \mathbf{0} \\ \mathbf{0} & \psi_y \end{pmatrix}$ , we can re-write  $\mathbf{w}_i = \begin{pmatrix} \mathbf{x}_i \\ \mathbf{y}_i \end{pmatrix}$  as  $\mathbf{w}_i = \Lambda^T \mathbf{v}_i + \psi$ . The variables  $\mathbf{w}_i$  and  $\mathbf{v}_i$  then follow a joint normal distribution

$$(A.4) \quad \begin{pmatrix} \mathbf{w}_i \\ \mathbf{v}_i \end{pmatrix} \sim MVN \left( \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \psi_y & \mathbf{0} \\ \mathbf{0} & \psi_x \\ & \Lambda^T & \Lambda \\ & & & \mathbf{I} \end{pmatrix} \right).$$

We view the latent factors  $\mathbf{v}_i$  as the missing data. In the E step, we calculate the expectation of the log likelihood function for complete data. The expectation is taken with respect to the conditional distribution of  $\mathbf{v}_i$  given  $\mathbf{w}_i$

$$\begin{aligned} E(\log l(\mathbf{v}, \mathbf{w}) | \mathbf{w}) &= -\frac{1}{2} \sum_{i=1}^n E[\mathbf{v}_i^T \Lambda \psi^{-1} \Lambda^T \mathbf{v}_i - 2\mathbf{v}_i^T \Lambda \psi^{-1} \mathbf{w}_i | \mathbf{w}_i] - \frac{n}{2} \log |\psi| - \frac{1}{2} \sum_{i=1}^n \mathbf{w}_i^T \psi^{-1} \mathbf{w}_i \\ &= -\frac{1}{2} \sum_{i=1}^n E[\text{tr}(\Lambda \psi (\Lambda \Lambda^T)^T \mathbf{v}_i \mathbf{v}_i^T) | \mathbf{w}_i] + \sum_{i=1}^n E(\mathbf{v}_i | \mathbf{w}_i)^T \Lambda \psi^{-1} \mathbf{w}_i - \frac{n}{2} \log |\psi| - \frac{1}{2} \sum_{i=1}^n \mathbf{w}_i^T \psi^{-1} \mathbf{w}_i. \end{aligned}$$

In the M step, we maximize the above expectation. To do so, we take derivatives of the log-likelihood function with respect to  $\Lambda_x$ ,  $\Lambda_y$  and  $\Lambda_u$ , and obtain

$$\begin{aligned}\frac{\partial E\log l}{\partial \Lambda_x} &= (\Lambda_x^{-1}) \sum_{i=1}^n \psi_x^{-1} \Lambda_x^T E(\mathbf{z}_i \mathbf{z}_i^T | \mathbf{w}_i) - \sum_{i=1}^n \psi_x^{-1} \mathbf{x}_i E(\mathbf{z}_i | \mathbf{w}_i)^T, \\ \frac{\partial E\log l}{\partial \Lambda_y} &= (\Lambda_y^{-1}) \sum_{i=1}^n \psi_y^{-1} \Lambda_y^T E(\mathbf{z}_i \mathbf{z}_i^T | \mathbf{w}_i) + \sum_{i=1}^n \psi_y^{-1} \Lambda_u^T E(\mathbf{u}_i \mathbf{z}_i^T | \mathbf{w}_i) - \sum_{i=1}^n \psi_y^{-1} \mathbf{y}_i E(\mathbf{z}_i | \mathbf{w}_i)^T, \\ \frac{\partial E\log l}{\partial \Lambda_u} &= (\Lambda_u^{-1}) \sum_{i=1}^n \psi_y^{-1} \Lambda_u^T E(\mathbf{u}_i \mathbf{u}_i^T | \mathbf{w}_i) + \sum_{i=1}^n \psi_y^{-1} \Lambda_y^T E(\mathbf{z}_i \mathbf{u}_i^T | \mathbf{w}_i) - \sum_{i=1}^n \psi_y^{-1} \mathbf{y}_i E(\mathbf{u}_i | \mathbf{w}_i)^T,\end{aligned}$$

where the conditional expectations are

$$(A.9) \quad E(\mathbf{v}_i | \mathbf{w}_i) = \Lambda(\psi + \Lambda^T \Lambda)^{-1} \mathbf{w}_i,$$

$$(A.10) \quad \text{Var}(\mathbf{v}_i | \mathbf{w}_i) = \mathbf{I} - \Lambda(\psi + \Lambda^T \Lambda)^{-1} \Lambda^T$$

$$(A.11) \quad E(\mathbf{v}_i \mathbf{v}_i^T | \mathbf{w}_i) = \text{Var}(\mathbf{v}_i | \mathbf{w}_i) + E(\mathbf{v}_i | \mathbf{w}_i) E(\mathbf{v}_i | \mathbf{w}_i)^T.$$

The above equations form the basis of our EM algorithms.

#### ACKNOWLEDGEMENTS

MC was supported by the National Institutes of Health (NIH) grants R01 GM105785, R01 CA082659 and P01 CA142538. XZ was supported by NIH grants R01HL117626 (PI Abecasis) and R21ES024834 (PI Pierce), and a grant from the Foundation for the National Institutes of Health through the Accelerating Medicines Partnership (BOEH15AMP, co-PIs Boehnke and Abecasis). We thank Dr. Dominic Grun for providing the raw read counts of the first dataset.

#### REFERENCES

- [1] ACHIM, K., PETTIT, J. B., SARAIVA, L. R., GAVRIOUCHKINA, D., LARSSON, T., ARENDT, D. and MARIONI, J. C. (2015). High-throughput spatial mapping of single-cell RNA-seq data to tissue of origin. *Nat Biotechnol* **33** 503-9.
- [2] BATTLE, A., MOSTAFAVI, S., ZHU, X., POTASH, J. B., WEISSMAN, M. M., MCCORMICK, C., HAUDENSCHILD, C. D., BECKMAN, K. B., SHI, J., MEI, R., URBAN, A. E., MONTGOMERY, S. B., LEVINSON, D. F. and KOLLER, D. (2014). Characterizing the genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals. *Genome Res* **24** 14-24.
- [3] BLUM, Y., LE MIGNON, G., LAGARRIGUE, S. and CAUSEUR, D. (2010). A factor model to analyze heterogeneity in gene expression. *BMC Bioinformatics* **11** 368.
- [4] BOREL, C., FERREIRA, P. G., SANTONI, F., DELANEAU, O., FORT, A., POPADIN, K. Y., GARIERI, M., FALCONNET, E., RIBAU, P., GUIPPONI, M., PADIOLEAU, I., CARNINCI, P., DERMITZAKIS, E. T. and ANTONARAKIS, S. E. (2015).

- Biased allelic expression in human primary fibroblast single cells. *Am J Hum Genet* **96** 70-80.
- [5] BRENECKE, P., ANDERS, S., KIM, J. K., KOŁODZIEJCZYK, A. A., ZHANG, X. W., PROSERPIO, V., BAYING, B., BENES, V., TEICHMANN, S. A., MARIONI, J. C. and HEISLER, M. G. (2013). Accounting for technical noise in single-cell RNA-seq experiments. *Nature Methods* **10** 1093-1095.
- [6] BUETTNER, F., NATARAJAN, K. N., CASALE, F. P., PROSERPIO, V., SCIALDONE, A., THEIS, F. J., TEICHMANN, S. A., MARIONI, J. C. and STEGLE, O. (2015). Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat Biotechnol* **33** 155-60.
- [7] CARVALHO, C. M., CHANG, J., LUCAS, J. E., NEVINS, J. R., WANG, Q. L. and WEST, M. (2008). High-Dimensional Sparse Factor Modeling: Applications in Gene Expression Genomics. *Journal of the American Statistical Association* **103** 1438-1456.
- [8] DENG, Q., RAMSKOLD, D., REINIUS, B. and SANDBERG, R. (2014). Single-cell RNA-seq reveals dynamic, random monoallelic gene expression in mammalian cells. *Science* **343** 193-6.
- [9] DURRUTHY-DURRUTHY, R., GOTTLIEB, A., HARTMAN, B. H., WALDHAUS, J., LASKE, R. D., ALTMAN, R. and HELLER, S. (2014). Reconstruction of the mouse otocyst and early neuroblast lineage at single-cell resolution. *Cell* **157** 964-78.
- [10] FINAK, G., MCDAVID, A., YAJIMA, M., DENG, J., GERSUK, V., SHALEK, A. K., SLICHTER, C. K., MILLER, H. W., MCEL RATH, M. J., PR LIC, M., LINSLEY, P. S. and GOTTARDO, R. (2015). MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol* **16** 278.
- [11] GAGNON-BARTSCH, J. A., JACOB, L. and SPEED, T. P. (2013). Removing unwanted variation from high dimensional data with negative controls. Technical Report.
- [12] GAGNON-BARTSCH, J. A. and SPEED, T. P. (2012). Using control genes to correct for unwanted variation in microarray data. *Biostatistics* **13** 539-552.
- [13] GRUN, D., KESTER, L. and VAN OUDENAARDEN, A. (2014). Validation of noise models for single-cell transcriptomics. *Nat Methods* **11** 637-40.
- [14] ISLAM, S., ZEISEL, A., JOOST, S., LA MANNO, G., ZAJAC, P., KASPER, M., LONNERBERG, P. and LINNARSSON, S. (2014). Quantitative single-cell RNA-seq with unique molecular identifiers. *Nat Methods* **11** 163-6.
- [15] JACOB, L., GAGNON-BARTSCH, J. A. and SPEED, T. P. (2015). Correcting gene expression data when neither the unwanted variation nor the factor of interest are observed. *Biostatistics* **17** 16-28.
- [16] JAITIN, D. A., KENIGSBERG, E., KEREN-SHAUL, H., ELEFANT, N., PAUL, F., ZARETSKY, I., MILDNER, A., COHEN, N., JUNG, S., TANAY, A. and AMIT, I. (2014). Massively parallel single-cell RNA-seq for marker-free decomposition of tissues into cell types. *Science* **343** 776-9.
- [17] JIANG, L., SCHLESINGER, F., DAVIS, C. A., ZHANG, Y., LI, R., SALIT, M., GINGERAS, T. R. and OLIVER, B. (2011). Synthetic spike-in standards for RNA-seq experiments. *Genome Res* **21** 1543-51.
- [18] JOHNSON, W. E., LI, C. and RABINOVIC, A. (2007). Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* **8** 118-127.
- [19] KANG, H. M., YE, C. and ESKIN, E. (2008). Accurate discovery of expression quantitative trait loci under confounding from spurious and genuine regulatory hotspots. *Genetics* **180** 1909-1925.
- [20] KHARCHENKO, P. V., SILBERSTEIN, L. and SCADDEN, D. T. (2014). Bayesian ap-

- proach to single-cell differential expression analysis. *Nature Methods* **11** 740-U184.
- [21] KIM, K. T., LEE, H. W., LEE, H. O., KIM, S. C., SEO, Y. J., CHUNG, W., EUM, H. H., NAM, D. H., KIM, J., JOO, K. M. and PARK, W. Y. (2015a). Single-cell mRNA sequencing identifies subclonal heterogeneity in anti-cancer drug responses of lung adenocarcinoma cells. *Genome Biol* **16** 127.
- [22] KIM, J. K., KOLODZIEJCZYK, A. A., ILLICIC, T., TEICHMANN, S. A. and MARIONI, J. C. (2015b). Characterizing noise structure in single-cell RNA-seq distinguishes genuine from technical stochastic allelic expression. *Nat Commun* **6** 8687.
- [23] LAPPALAINEN, T., SAMMETH, M., FRIEDLANDER, M. R., T HOEN, P. A., MONLONG, J., RIVAS, M. A., GONZALEZ-PORTA, M., KURBATOVA, N., GRIEBEL, T., FERREIRA, P. G., BARANN, M., WIELAND, T., GREGER, L., VAN ITERSON, M., ALMLOF, J., RIBECA, P., PULYAKHINA, I., ESSER, D., GIGER, T., TIKHONOV, A., SULTAN, M., BERTIER, G., MACARTHUR, D. G., LEK, M., LIZANO, E., BUERMANS, H. P., PADIOLEAU, I., SCHWARZMAYR, T., KARLBERG, O., ONGEN, H., KILPINEN, H., BELTRAN, S., GUT, M., KAHLEM, K., AMSTISLAVSKIY, V., STEGLE, O., PIRINEN, M., MONTGOMERY, S. B., DONNELLY, P., MCCARTHY, M. I., FLICEK, P., STROM, T. M., GEUVADIS, C., LEHRACH, H., SCHREIBER, S., SUDBRACK, R., CARRACEDO, A., ANTONARAKIS, S. E., HASLER, R., SYVANEN, A. C., VAN OMMEN, G. J., BRAZMA, A., MEITINGER, T., ROSENSTIEL, P., GUIGO, R., GUT, I. G., ESTIVILL, X. and DERMITZAKIS, E. T. (2013). Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* **501** 506-11.
- [24] LAW, C. W., CHEN, Y., SHI, W. and SMYTH, G. K. (2014). voom: Precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol* **15** R29.
- [25] LEA, A. J., TUNG, J. and ZHOU, X. (2015). A Flexible, Efficient Binomial Mixed Model for Identifying Differential DNA Methylation in Bisulfite Sequencing Data. *PLoS Genet* **11** e1005650.
- [26] LEE, S., CHUGH, P. E., SHEN, H., EBERLE, R. and DITTMER, D. P. (2013). Poisson factor models with applications to non-normalized microRNA profiling. *Bioinformatics* **29** 1105-11.
- [27] LEE, M. C., LOPEZ-DIAZ, F. J., KHAN, S. Y., TARIQ, M. A., DAYN, Y., VASKE, C. J., RADENBAUGH, A. J., KIM, H. J., EMERSON, B. M. and POURMAND, N. (2014). Single-cell analyses of transcriptional heterogeneity during drug tolerance transition in cancer cells by RNA sequencing. *Proc Natl Acad Sci U S A* **111** E4726-35.
- [28] LEEK, J. T. and STOREY, J. D. (2007). Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genet* **3** 1724-35.
- [29] LISTGARTEN, J., KADIE, C., SCHADT, E. E. and HECKERMAN, D. (2010). Correction for hidden confounders in the genetic analysis of gene expression. *Proc Natl Acad Sci USA* **107** 16465-16470.
- [30] LUCAS, J. E., KUNG, H. N. and CHI, J. T. (2010). Latent factor analysis to discover pathway-associated putative segmental aneuploidies in human cancers. *PLoS Comput Biol* **6** e1000920.
- [31] MACOSKO, E. Z., BASU, A., SATIJA, R., NEMESH, J., SHEKHAR, K., GOLDMAN, M., TIROSH, I., BIALAS, A. R., KAMITAKI, N., MARTERSTECK, E. M., TROMBETTA, J. J., WEITZ, D. A., SANES, J. R., SHALEK, A. K., REGEV, A. and MCCARROLL, S. A. (2015). Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* **161** 1202-14.
- [32] MONTGOMERY, S. B., SAMMETH, M., GUTIERREZ-ARCELUS, M., LACH, R. P., INGLE, C., NISBETT, J., GUIGO, R. and DERMITZAKIS, E. T. (2010). Transcriptome genetics using second generation sequencing in a Caucasian population. *Nature* **464**

- 773-7.
- [33] PARTS, L., STEGLE, O., WINN, J. and DURBIN, R. (2011). Joint genetic analysis of gene expression data with inferred cellular phenotypes. *PLoS Genet* **7** e1001276.
  - [34] PICKRELL, J. K., MARIONI, J. C., PAI, A. A., DEGNER, J. F., ENGELHARDT, B. E., NKADORI, E., VEYRIERAS, J. B., STEPHENS, M., GILAD, Y. and PRITCHARD, J. K. (2010). Understanding mechanisms underlying human gene expression variation with RNA sequencing. *Nature* **464** 768-72.
  - [35] POURNARA, I. and WERNISCH, L. (2007). Factor analysis for gene regulatory networks and transcription factor activity profiles. *BMC Bioinformatics* **8** 61.
  - [36] REINIUS, B. and SANDBERG, R. (2015). Random monoallelic expression of autosomal genes: stochastic transcription and allele-level regulation. *Nat Rev Genet* **16** 653-64.
  - [37] RISSO, D., NGAI, J., SPEED, T. P. and DUDOIT, S. (2014). Normalization of RNA-seq data using factor analysis of control genes or samples. *Nat Biotechnol* **32** 896-902.
  - [38] RITCHIE, M. E., Phipson, B., Wu, D., Hu, Y., Law, C. W., Shi, W. and Smyth, G. K. (2015). limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res* **43** e47.
  - [39] SATIJA, R., FARRELL, J. A., GENNERT, D., SCHIER, A. F. and REGEV, A. (2015). Spatial reconstruction of single-cell gene expression data. *Nat Biotechnol* **33** 495-502.
  - [40] SEYEDNASROLLAH, F., LAIHO, A. and ELO, L. L. (2015). Comparison of software packages for detecting differential expression in RNA-seq studies. *Brief Bioinform* **16** 59-70.
  - [41] SHALEK, A. K., SATIJA, R., SHUGA, J., TROMBETTA, J. J., GENNERT, D., LU, D. N., CHEN, P. L., GERTNER, R. S., GAUBLomme, J. T., YOSEF, N., SCHWARTZ, S., FOWLER, B., WEAVER, S., WANG, J., WANG, X. H., DING, R. H., RAYCHOWDHURY, R., FRIEDMAN, N., HACOHEM, N., PARK, H., MAY, A. P. and REGEV, A. (2014). Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. *Nature* **510** 363-+.
  - [42] SONESON, C. and DELORENZI, M. (2013). A comparison of methods for differential expression analysis of RNA-seq data. *BMC Bioinformatics* **14** 91.
  - [43] STEGLE, O., TEICHMANN, S. A. and MARIONI, J. C. (2015). Computational and analytical challenges in single-cell transcriptomics. *Nat Rev Genet* **16** 133-45.
  - [44] STEGLE, O., PARTS, L., DURBIN, R. and WINN, J. (2010). A Bayesian framework to account for complex non-genetic factors in gene expression levels greatly increases power in eQTL studies. *PLoS Comput Biol* **6** e1000770.
  - [45] SUN, Y., ZHANG, N. R. and OWEN, A. B. (2012). Multiple hypothesis testing adjusted for latent variables, with an application to the AGEMAP gene expression data. *Annals of Applied Statistics* **6** 1664-1688.
  - [46] TANG, F., BARBACIORU, C., BAO, S., LEE, C., NORDMAN, E., WANG, X., LAO, K. and SURANI, M. A. (2010). Tracing the derivation of embryonic stem cells from the inner cell mass by single-cell RNA-Seq analysis. *Cell Stem Cell* **6** 468-78.
  - [47] TREUTLEIN, B., BROWNFIELD, D. G., WU, A. R., NEFF, N. F., MANTALAS, G. L., ESPINOZA, F. H., DESAI, T. J., KRASNOW, M. A. and QUAKE, S. R. (2014). Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. *Nature* **509** 371-5.
  - [48] TUNG, J., ZHOU, X., ALBERTS, S. C., STEPHENS, M. and GILAD, Y. (2015). The genetic architecture of gene expression levels in wild baboons. *Elife* **4**.
  - [49] USOSKIN, D., FURLAN, A., ISLAM, S., ABDO, H., LONNERBERG, P., LOU, D., HJERLING-LEFFLER, J., HAEGGSTROM, J., KHARCHENKO, O., KHARCHENKO, P. V., LINNARSSON, S. and ERNFORS, P. (2015). Unbiased classification of sensory neuron types by large-scale single-cell RNA sequencing. *Nat Neurosci* **18** 145-53.



- [50] VALLEJOS, C. A., MARIONI, J. C. and RICHARDSON, S. (2015). BASiCS: Bayesian Analysis of Single-Cell Sequencing Data. *PLoS Comput Biol* **11** e1004333.
- [51] WALKER, W. L., LIAO, I. H., DONALD L. GILBERT, K. S. P. C. E. M. L. L. BRENDA WONG and SHARP, F. R. (2008). Empirical Bayes accomodation of batch-effects in microarray data using identical replicate reference samples: application to RNA expression profiling of blood from Duchenne muscular dystrophy patients. *BMC Genomics* **9** 494.
- [52] WEST, M. (2003). Bayesian factor regression models in the "Large p, Small n" paradigm. *Bayesian Statistics* **7** 733-742.
- [53] XUE, Z., HUANG, K., CAI, C., CAI, L., JIANG, C. Y., FENG, Y., LIU, Z., ZENG, Q., CHENG, L., SUN, Y. E., LIU, J. Y., HORVATH, S. and FAN, G. (2013). Genetic programs in human and mouse early embryos revealed by single-cell RNA sequencing. *Nature* **500** 593-7.
- [54] YANG, C., WANG, L., ZHANG, S. and ZHAO, H. (2013). Accounting for non-genetic factors by low-rank representation and sparse regression for eQTL mapping. *Bioinformatics* **29** 1026-1034.
- [55] ZEISEL, A., MUNOZ-MANCHADO, A. B., CODELUPPI, S., LONNERBERG, P., LA MANNO, G., JUREUS, A., MARQUES, S., MUNGUBA, H., HE, L., BETSHOLTZ, C., ROLNY, C., CASTELO-BRANCO, G., HJERLING-LEFFLER, J. and LINNARSSON, S. (2015). Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* **347** 1138-42.
- [56] ZHOU, X. and STEPHENS, M. (2012). Genome-wide efficient mixed-model analysis for association studies. *Nat Genet* **44** 821-4.
- [57] ZHOU, X. and STEPHENS, M. (2014). Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nat Methods* **11** 407-9.
- [58] ZHOU, M., HANNAH, L., DUNSON, D. and CARIN, L. (2012). Beta-negative binomial process and Poisson factor analysis. *Artificial Intelligence and Statistics* **22** 1462-1471.

DEPARTMENT OF BIostatISTICS  
DEPARTMENT OF GENETICS  
UNIVERSITY OF NORTH CAROLINA  
CHAPEL HILL, NC 27599  
E-MAIL: [mengjie@email.unc.edu](mailto:mengjie@email.unc.edu)

DEPARTMENT OF BIostatISTICS  
CENTER FOR STATISTICAL GENETICS  
UNIVERSITY OF MICHIGAN  
ANN ARBOR, MI 48109  
E-MAIL: [xzhousph@umich.edu](mailto:xzhousph@umich.edu)