# Contrastive latent variable modeling with application to case-control sequencing experiments

Andrew Jones[1], F. William Townes[1], Didong Li[1,2], and Barbara E. Engelhardt[1,3]

[1]*Department of Computer Science, Princeton University*
[2]*Department of Biostatistics, University of California, Los Angeles*
[3]*Center for Statistics and Machine Learning, Princeton University*

February 16, 2021

## Abstract

High-throughput RNA-sequencing (RNA-seq) technologies are powerful tools for understanding cellular state. Often it is of interest to quantify and summarize changes in cell state that occur between experimental or biological conditions. Differential expression is typically assessed using univariate tests to measure gene-wise shifts in expression. However, these methods largely ignore changes in transcriptional correlation. Furthermore, there is a need to identify the low-dimensional structure of the gene expression shift to identify collections of genes that change between conditions. Here, we propose contrastive latent variable models designed for count data to create a richer portrait of differential expression in sequencing data. These models disentangle the sources of transcriptional variation in different conditions, in the context of an explicit model of variation at baseline. Moreover, we develop a model-based hypothesis testing framework that can test for global and gene subset-specific changes in expression. We test our model through extensive simulations and analyses with count-based gene expression data from perturbation and observational sequencing experiments. We find that our methods can effectively summarize and quantify complex transcriptional changes in case-control experimental sequencing data.

## 1 Introduction

High-throughput RNA-sequencing technologies have emerged as useful tools for understanding transcriptional patterns. Recently, single-cell RNA-sequencing (scRNA-seq) technologies have allowed for investigation of these patterns at the level of individual cells. Together, these sequencing technologies have revealed new insights about a range of biological questions, from how cell types differ from one another to how cells respond to therapeutic drugs. In many scRNA-seq and bulk RNA-seq experiments, there are gene expression readouts

for two or more experimental conditions or biological traits, such as tumor versus normal (Young et al., 2018; Kinker et al., 2020), drug exposure versus placebo exposure (Srivastava and Yanagihara, 2010; McFarland et al., 2020), or ventilated versus non-ventilated lung tissue (Consortium et al., 2020). In these cases, it is of interest to understand how the transcription levels differ in a *foreground dataset* (collected from the treatment condition) relative to a *background dataset* (collected from the control condition). These changes are traditionally identified using methods for differential gene expression, which estimate the average shift in expression levels between conditions.

Most differential expression methods for RNA-seq and scRNA-seq compute the univariate change between conditions for each gene separately (Kharchenko et al., 2014; Finak et al., 2015; Qiu et al., 2017). In scRNA-seq, this amounts to treating each cell as an independent sample from one of two distributions over expression state; then each gene is considered marginally. These analyses use the repeated samples to quantify the uncertainty in the estimate of differential expression.

However, these methods for differential expression ignore a fundamental benefit of collecting scRNA-seq data over bulk RNA-sequencing: The ability to quantify population-level variation within a single sample across all of the cells. Population variation exists in pools of single cells, even from the same tissue. This gives us an opportunity to find differences in variation among genes, not just each gene itself. Furthermore, it is believed that these transcriptional changes follow a low-dimensional structure. (Dixit et al., 2016; McFarland et al., 2020; Becht et al., 2019). In particular, cells' fixed energy budgets and strongly correlated and interacting gene networks constrain the types of structural changes in gene expression between conditions. In other words, changes in expression levels can be described in fewer dimensions than there are genes — a phenomenon that most differential expression methods fail to capture. This low-dimensional projection can then be studied to find groups of genes that change similarly (Xia et al., 2015), and to quantify gene covariation in a low-dimensional representation (Townes et al., 2019; Ding et al., 2018). There is a need for methods that can robustly identify the low-dimensional structure of variation in gene expression, and how this structure differs across experimental and biological conditions.

To address this gap in methodology, we develop a family of probabilistic models — contrastive Poisson latent variable models (CPLVMs) — that are designed to estimate the low-dimensional structure of the transcriptional response to perturbations as measured using sequencing technologies. Existing contrastive dimension reduction methods have shown promise for understanding the global shifts in variation between multiple conditions (Zou et al., 2013; Severson et al., 2019; Abid et al., 2018). But these methods fall short in a number of ways. First, they typically assume normally-distributed data, and do not treat count data from sequencing methods appropriately. Second, these methods are not designed to properly normalize count-based sequencing profiles. Finally, these approaches do not provide a hypothesis testing framework for testing differential expression across multiple genes.

Our contrastive Poisson latent variable model (CPLVM) bridges the gaps between ex-

isting contrastive methods and the needs of sequencing experiments by addressing these issues. We use a Poisson data likelihood that accounts for the count-based data produced by sequencing technologies. We show that these methods can identify changes in experimental and biological conditions that standard differential expression methods are not able to detect. Using our model, we decompose two-condition scRNA-seq and RNA-seq data into a small set of interpretable, nonnegative factors. Moreover, we build a hypothesis testing framework that accommodates hypotheses at varying scales, from testing for global shifts in gene expression across conditions to testing for correlated changes to a small set of candidate genes.

In this paper, we first describe the CPLVM in the context of related work and motivated by experiments in scRNA-seq. We then demonstrate the utility of our model through extensive simulations and experiments with multiple gene expression datasets. Using case-control scRNA-seq readouts of cells exposed to genetic and chemical perturbations (Dixit et al., 2016; McFarland et al., 2020), we show that the CPLVM can identify structure that is specific to the foreground data. Furthermore, we show that these methods can identify changes across case-control data that standard differential expression methods are not able to detect. In addition, we apply the CPLVM to RNA-seq measurements from coronary artery tissue of donors with heart disease and healthy donors (Consortium et al., 2020). We also show that the CPLVM hypothesis testing framework identifies experiments in which a specific, structured change in a small group of genes has been observed.

## 2 Methods

### 2.1 Problem definition

We motivate the problem using a generic scRNA-seq experiment, but our models can be applied to other count-based sequencing technologies as well. A scRNA-seq experiment with two conditions yields a set of unique molecular identifier (UMI) counts for cells from each condition. In this paper, we call measurements from the control condition the *background data* and measurements from the treatment condition the *foreground data.*

Suppose there are $n$ cells measured in the background condition, and $m$ cells measured in the foreground condition, with gene expression measured across $p$ (total) genes. We denote the data in matrix form as $\mathbf{Y} \in \mathbb{N}_0^{p \times n}$ and $\mathbf{X} \in \mathbb{N}_0^{p \times m}$, which contain the UMI counts for each cell and gene in the background and foreground data, respectively.

The column vectors $\mathbf{y}_i \in \mathbb{N}_0^p$ and $\mathbf{x}_j \in \mathbb{N}_0^p$ denote UMI counts across the $p$ genes for cell $i = \{1, \ldots, n\}$ or $j = \{1, \ldots, m\}$ from their respective conditions.

In this study, we are concerned with characterizing the transcriptional structure that exists in the foreground data $\mathbf{X}$ but not in $\mathbf{Y}$, as well as identifying the structure that is shared between the conditions. Decomposing these sources of variation into interpretable, low-dimensional structure is critical to understanding the effects of a treatment or different biological condition on cell state, regulation, and dynamics.

3

## 2.2 Related work

Several families of methods have been developed to characterize the changes in gene expression between experimental or biological conditions. In this section, we outline several of these approaches.

### 2.2.1 Differential expression methods

The most common approaches for quantifying transcriptional changes in bulk and scRNA-seq data are differential expression methods. In general, these approaches compute univariate estimates of the change in expression for each gene between conditions. We review several approaches below; see Wang et al. (2019) for a thorough review and benchmarking of differential expression methods in the context of scRNA-seq.

Most differential expression methods use linear models or generalized linear models (GLMs) to estimate the change in gene expression. For example, Multi-Input Multi-Output Single-Cell Analysis (MIMOSCA) — which was developed specifically for the setting of genetic perturbation experiments in scRNA-seq — uses a linear model with a Gaussian noise assumption (Dixit et al., 2016). Specifically, it assumes the following model for the log-transformed and normalized UMI counts:

$$\log_2\left(\frac{\mathbf{y}_i}{n_i}s + c\right) = \boldsymbol{\beta}_0 + \boldsymbol{\epsilon} \tag{1}$$

$$\log_2\left(\frac{\mathbf{x}_j}{n_j}s + c\right) = \boldsymbol{\beta}_1 + \boldsymbol{\beta}_0 + \boldsymbol{\epsilon} \tag{2}$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I) \tag{3}$$

where $n_i = \sum_k y_{ik}$ is the total number of counts in cell $i$ ($n_j$ is similarly defined), $s$ is a constant multiplicative factor, and $c$ is a "pseudocount" added to avoid taking $\log_2(0)$ (typically, $c = 1$). The coefficient vector $\boldsymbol{\beta}_1 \in \mathbb{R}^p$ then captures the average fold-change in gene expression for a single gene between the conditions, and can be tested directly for significance.

GLMs more flexibly capture non-Gaussian data likelihoods. For example, in the context of scRNA-seq count data, a popular choice is the Poisson likelihood. In this case, the UMI count for gene $k$ in each cell is assumed to be a draw from a Poisson distribution, whose rate parameter is a transformation of the linear predictor:

$$y_{ik} \sim \text{Poisson}\left(\frac{n_i}{s}g^{-1}(\beta_{0k})\right) \quad i = 1, \ldots, n \tag{4}$$

$$x_{jk} \sim \text{Poisson}\left(\frac{n_j}{s}g^{-1}(\beta_{0k} + \beta_{1k})\right) \quad j = 1, \ldots, m. \tag{5}$$

The canonical link function $g(\cdot)$ for a Poisson likelihood, $\log(\cdot)$, is typically used in the Poisson setting. Several existing methods, such as single-cell differential expression (SCDE Kharchenko

4

et al. 2014), use a Poisson GLM to identify differential expression across conditions in scRNA-seq data. Closely related to the Poisson GLM, a common approach is to allow for overdispersion by using a negative binomial likelihood, which is equivalent to a gamma-Poisson mixture (Love et al., 2014; Robinson et al., 2010; Hafemeister and Satija, 2019). One method uses a zero-inflated negative binomial to model dropout events (Miao et al., 2018).

Other distributional assumptions have also been proposed for differential expression in sequencing data. One approach, Model-based Analysis of Single-cell Transcriptomics (MAST), uses a hurdle model combined with a Gaussian likelihood (Finak et al., 2015). Another approach, scDD, uses a Dirichlet process mixture of Gaussians to model the potentially multimodal expression across cells and computes Bayes factors to quantify differential expression (Korthauer et al., 2016). Other nonparametric approaches have also been proposed, including using Earth Mover's Distance to quantify expression changes (Nabavi et al., 2016) and Cramér-von Mises and Kolmogorov-Smirnov hypothesis tests (Delmans and Hemberg, 2016).

While differential expression methods have proven to be reliable for identifying changes that occur in individual genes across conditions, they typically ignore any correlation structure between genes. This is an important limitation, as gene expression has been shown to have substantial correlation structure (Stuart et al., 2003), and identifying changes in this structure across conditions is of great interest.

### 2.2.2 Two-sample covariance comparison methods

Another related line of work has focused on identifying differences in the covariance between two conditions. Most commonly, these approaches rely on a hypothesis test to decide whether covariance matrices $\Sigma_X, \Sigma_Y$ are different:

$$H_0: \quad \Sigma_X = \Sigma_Y, \qquad H_1: \quad \Sigma_X \neq \Sigma_Y.$$

There exist a number of such tests in the setting of low-dimensional data, including modified multivariate generalizations of Levene's test (O'Brien, 1992) and the commonly-used likelihood ratio test (Anderson, 1958). But high-dimensional data pose a greater challenge. Since these existing tests were designed for small numbers of features $p$ relative to sample size $n$, they have poor statistical power when applied to high-dimensional data where $p \gg n$, and in some cases are not even well-defined (Cai et al., 2013).

Some covariance inequality tests have been designed to address the problem of high-dimensional data by using estimators of the distance between covariance matrices based on the Frobenius norm (Li et al., 2012; Srivastava and Yanagihara, 2010). However, these tests have low power to detect the true effect when the differences between the covariance matrices are sparse.

Two techniques have emerged to address both the issue of high-dimensional data and the possibility of a small number of entries driving the differences between two covariance

matrices. The first approach uses a test statistic based on the largest standardized difference between the two covariance matrices' entries (Cai et al., 2013). The second approach uses a Gaussian graphical model framework to infer the differential network structure (Xia et al., 2015). Other related approaches consider building Gaussian graphical models from precision matrices and identifying network edges that are differentially identified across two conditions (Glass et al., 2013).

More recently, a hypothesis testing approach that is robust in the setting of high-dimensionality and low sample size was developed using the strongly spiked eigenvalue (SSE) model (Aoshima and Yata, 2018; Ishii et al., 2019). The SSE model assumes the first eigenvalues $\lambda_{1x}, \lambda_{1y}$ of the covariance matrices $\Sigma_x, \Sigma_y$ are "strongly spiked" relative to the subsequent eigenvalues, in the sense of

$$\liminf_{p \to \infty} \left\{ \frac{\lambda_{1i}}{\text{tr}(\Sigma_i^2)} \right\} > 0.$$

The authors show that the SSE assumption is reasonable in many high-dimensional settings, especially when $p \gg n$. They derived the limiting distributions for test statistics under on this model, as well as the power and size of the accompanying hypothesis tests. They found that the SSE model had greater statistical power compared to previous models that assumed more diffuse eigenvalue spectra. However, methods for hypothesis testing neglect the goals of exploratory data analysis, including identifying interpretable, low-dimensional factors that explain the changes in covariance structure across conditions.

### 2.2.3 Contrastive dimension reduction methods

Contrastive dimension reduction methods estimate low-dimensional changes in variation between conditions. In particular, these methods aim to identify variation that exists in the foreground data but not in the background data. Furthermore, they typically assume that the variation in each condition can be explained by a small number of latent dimensions.

As one of the first steps in this direction, a framework for contrastive learning in mixture models was proposed (Zou et al., 2013). This approach assumes that the background and foreground data are generated from a set of mixture distributions, some of which are shared between the two conditions, and some of which are exclusive to each condition. Specifically, given a set of mixture parameters $\{\mu_\ell\}, w_\ell\}_{\ell=1}^L$, condition-specific mixture weights $\{w_\ell^{\text{b}}\}_{\ell=1}^L, \{w_\ell^{\text{f}}\}_{\ell=1}^L$, and three disjoint index sets $A, B, C \subseteq [L]$, it assumes $\mathbf{Y}$ and $\mathbf{X}$ are drawn from a set of mixtures:

$$p(\mathbf{y}_i; \{\mu_\ell, w_\ell\}_{\ell \in A \cup B}) = \sum_{\ell \in A \cup B} w_\ell^{\text{b}} f(\mathbf{y}_i | \mu_\ell), \quad i = 1, \ldots, n$$

$$p(\mathbf{x}_j; \{\mu_\ell, w_\ell\}_{\ell \in B \cup C}) = \sum_{\ell \in B \cup C} w_\ell^{\text{f}} f(\mathbf{x}_j | \mu_\ell), \quad j = 1, \ldots, m.$$

6

Note that the mixture components indexed by $B$ are shared between the conditions, while those indexed by $A$ and $C$ are unique to the background and foreground, respectively. This framework encompasses several general models, including topic models, such as Latent Dirichlet Allocation (LDA). The authors were primarily interested in estimating the foreground-specific model parameters, $\{\mu_\ell, w_\ell^{\mathrm{f}}\}_{\ell \in A}$. Their inference approach relies on a tensor decomposition method for estimating the foreground-specific latent components without estimating the background or shared components.

As an important special case of this contrastive learning framework, contrastive principal component analysis (CPCA) was derived explicitly, which extends the classical PCA method (Abid et al., 2018). Specifically, given the sample covariance matrices of the background and foreground conditions, $\widehat{\Sigma}_x$ and $\widehat{\Sigma}_y$, the objective function of CPCA seeks to find a unit vector $\mathbf{v}$ that maximizes the variance in the foreground and minimizes the variance in the background:

$$\mathrm{argmax}_{\mathbf{v} \in \mathbb{S}^{p-1}} \left\{ \mathbf{v}^\top \widehat{\Sigma}_x \mathbf{v} - \gamma \mathbf{v}^\top \widehat{\Sigma}_y \mathbf{v} \right\}.$$

Here, $\gamma \geq 0$ is a tuning parameter controlling the relative influence of the background data. When $\gamma = 0$, this model reduces to PCA on the foreground data. This problem can be solved analytically: the top $k$ "contrastive principal components" correspond to the $k$ eigenvectors $\left[ \mathbf{u}_1, \ldots, \mathbf{u}_k \right]$, which represent the top $k$ eigenvalues $\lambda_1 \geq \cdots \geq \lambda_k$ of the differential covariance:

$$C = \widehat{\Sigma}_x - \gamma \widehat{\Sigma}_y.$$

The authors showed that CPCA accurately recovers structure that is unique to the foreground data, and CPCA is able to identify heterogeneous responses in two-condition gene expression data (Abid et al., 2017, 2018).

A sparse version of CPCA was recently developed, which allows for greater interpretability of the estimated components, especially in high-dimensional settings (Boileau et al., 2020). Building off of sparse PCA (Zou et al., 2006), which uses element-wise $\ell_1$ regularization to encourage zeros in the loadings matrix, the authors propose an estimation procedure that alternates between estimating the principal components and the sparse loadings matrix. They demonstrate the behavior of sparse CPCA on a series of gene and protein expression datasets.

Most closely related to our work, probabilistic counterparts to CPCA have been proposed. The contrastive latent variable model (CLVM, Severson et al. 2019) captures structure that is unique to the foreground data, as well as structure shared between the conditions. In particular, the shared variation is described by a set of latent variables $\{\mathbf{z}_i^{\mathrm{b}}\}_{i=1}^n$ and $\{\mathbf{z}_j^{\mathrm{f}}\}_{j=1}^m$, and the foreground-specific variation is captured by another set of latent variables

$\{\mathbf{t}_j\}_{j=1}^m$. Using Gaussian likelihoods and priors, the CLVM has the following form:

$$\mathbf{y}_i|\mathbf{z}_i^{\mathrm{b}} \sim \mathcal{N}(\mathbf{S}\mathbf{z}_i^{\mathrm{b}} + \boldsymbol{\mu}^{\mathrm{b}}, \sigma^2 I)$$
$$\mathbf{x}_j|\mathbf{z}_j^{\mathrm{f}}, \mathbf{t}_j \sim \mathcal{N}(\mathbf{S}\mathbf{z}_j^{\mathrm{f}} + \mathbf{W}\mathbf{t}_j + \boldsymbol{\mu}^{\mathrm{f}}, \sigma^2 I)$$
$$\mathbf{z}_i^{\mathrm{b}} \sim \mathcal{N}(\mathbf{0}, I), \quad \mathbf{z}_j^{\mathrm{f}} \sim \mathcal{N}(\mathbf{0}, I), \quad \mathbf{t}_j \sim \mathcal{N}(\mathbf{0}, I).$$

Here, $\mathbf{S} \in \mathbb{R}^{p \times k_1}$ and $\mathbf{W} \in \mathbb{R}^{p \times k_2}$ are loadings matrices that map from the latent dimensions $k_1$ and $k_2$ to the data feature dimension $p$. Through experiments with gene expression and image data, the authors showed that the CLVM can disentangle low-dimensional latent structure that is shared between the two conditions and structure that is specific to the foreground.

Another model-based contrastive method, probabilistic contrastive principal component analysis (PCPCA, Li et al. 2020) was developed as a direct generalization of CPCA and probabilistic PCA. PCPCA provides a simple estimation procedure based on a likelihood ratio and was shown to be robust to noise and missing data. In applications, PCPCA was successful in identifying subgroup structure in case-control gene expression data. Although these probabilistic models have many advantages over previous approaches, both CLVM and PCPCA assume Gaussian errors, which is not ideal for modeling count-based expression profiles.

While contrastive dimension reduction methods have shown promise for analyzing two-condition datasets, there remains a need to adapt these methods to the setting of sequencing data, where observations are counts of RNA sequence fragments mapping to genes across the genome. Moreover, there is a substantial need to provide a common framework for both factor analysis and hypothesis testing in these models when it is useful to quantify the statistical significance of changes in the covariance structure of expression across cases and controls in an experimental setting.

## 3 Contrastive Poisson latent variable models for scRNA-seq

In this study, we develop a family of contrastive Poisson latent variable models (CPLVMs) that are designed to capture variation and covariation among count data that are unique to the foreground condition, as well as variation and covariation that are shared between the foreground and background. Furthermore, we build a hypothesis testing framework that quantifies support for structured changes in variation across conditions. Throughout, we rely on principled probabilistic modeling of count data, rather than data transformations and Gaussian models. In the context of sequencing data, our model explicitly accounts for the count-based nature of expression profiles, while decomposing case-control data into a small set of interpretable factors.

In the following sections, we first describe the CPLVM. Then we explain our inference procedure for the CPLVM, and we develop the corresponding hypothesis testing framework.

## 3.1 CPLVM definition

As above, let $\mathbf{Y} \in \mathbb{N}_0^{p \times n}$ and $\mathbf{X} \in \mathbb{N}_0^{p \times m}$ be the count matrices for $p$ genes and $n, m$ cells for the background and foreground conditions, respectively.

The CPLVM assumes that transcription variation in a sequencing experiment with multiple conditions can be described by a small set of latent factors. In particular, the model assumes that the variation shared between the conditions is described by a set of $k_1$-dimensional latent variables, $\{\mathbf{z}_i^{\mathrm{b}}\}_{i=1}^n$ and $\{\mathbf{z}_j^{\mathrm{f}}\}_{j=1}^m$. Furthermore, we assume that the foreground-specific variation is captured by another set of $k_2$-dimensional latent variables $\{\mathbf{t}_j\}_{j=1}^m$. To describe the mapping between these latent spaces and the data, we introduce loadings matrices $\mathbf{S} \in \mathbb{R}_+^{k_1 \times p}$ and $\mathbf{W} \in \mathbb{R}_+^{k_2 \times p}$, which map to the data space of dimension $p$ from the shared latent space of dimension $k_1$ and the foreground-specific latent space of dimension $k_2$.

To account for varying numbers of total counts between cells and experimental conditions, we include size factors $\alpha_i^{\mathrm{b}}$ and $\alpha_j^{\mathrm{f}}$ for each cell in each condition; these terms control for technical variation in the total number of counts per cell. Furthermore, to account for shifts in each gene's mean counts between conditions, we also include gene-specific multiplicative scale parameters $\boldsymbol{\delta} \in \mathbb{R}_+^p$. These terms are analogous to the gene-wise additive intercept terms in a linear model, but we constrain them to be nonnegative in this model. Thus, for gene $k$, $0 < \widehat{\delta}_k < 1$ indicates that there is lower relative expression of gene $k$ in the background cells, while $\widehat{\delta}_k > 1$ indicates higher relative expression of gene $k$ in the background cells.

The full generative model for the nonnegative CPLVM is then

$$\mathbf{y}_i | \mathbf{z}_i \sim \mathrm{Poisson}\left(\alpha_i^{\mathrm{b}} \boldsymbol{\delta} \odot \left(\mathbf{S}^\top \mathbf{z}_i^{\mathrm{b}}\right)\right) \tag{6}$$

$$\mathbf{x}_j | \mathbf{z}_j, \mathbf{t}_j \sim \mathrm{Poisson}\left(\alpha_j^{\mathrm{f}} \left(\mathbf{S}^\top \mathbf{z}_j^{\mathrm{f}} + \mathbf{W}^\top \mathbf{t}_j\right)\right) \tag{7}$$

$$z_{il}^{\mathrm{b}} \sim \mathrm{Gamma}(\gamma_1, \beta_1), \quad z_{jl}^{\mathrm{f}} \sim \mathrm{Gamma}(\gamma_2, \beta_2), \quad t_{jd} \sim \mathrm{Gamma}(\gamma_3, \beta_3), \tag{8}$$

$$W_{kd} \sim \mathrm{Gamma}(\gamma_4, \beta_4), \quad S_{jl} \sim \mathrm{Gamma}(\gamma_5, \beta_5), \quad \boldsymbol{\delta} \sim \mathrm{LogNormal}(0, \mathbf{I}_p), \tag{9}$$

where $l \in \{1, \ldots, k_1\}$, $d \in \{1, \ldots, k_2\}$, and $\odot$ represents a Hadamard (element-wise) product. Following previous work (Lopez et al., 2018), we place log normal priors on the size factors $\alpha_i, \alpha_j$ with parameters given by the empirical mean and variance of the log total counts for each cell.

## 3.2 Stochastic variational inference for the CPLVM

For a given experiment, we are interested in estimating the posterior distribution of $\mathbf{Z}^{\mathrm{b}}, \mathbf{Z}^{\mathrm{f}}, \mathbf{T}, \mathbf{S}$, and $\mathbf{W}$ given the data $\{\mathbf{Y}, \mathbf{X}\}$. Since the true posterior is intractable, we use a variational approximation.

Specifically, we perform approximate posterior inference on the latent variables $\mathbf{Z}^{\mathrm{b}}, \mathbf{Z}^{\mathrm{f}}, \mathbf{T}$, the loadings matrices $\mathbf{S}, \mathbf{W}$, and the mean-shift parameter $\boldsymbol{\delta}$ using a mean-field variational approximation. In other words, we approximate the true posterior distribution

$p(\mathbf{Z}^{\mathrm{b}}, \mathbf{Z}^{\mathrm{f}}, \mathbf{T}, \mathbf{S}, \mathbf{W}, \boldsymbol{\delta}|\mathbf{Y}, \mathbf{X})$ with a variational posterior distribution $q$ that fully factorizes:

$$q(\mathbf{Z}^{\mathrm{b}}, \mathbf{Z}^{\mathrm{f}}, \mathbf{T}, \mathbf{S}, \mathbf{W}) = q(\mathbf{Z}^{\mathrm{b}})q(\mathbf{Z}^{\mathrm{f}})q(\mathbf{T})q(\mathbf{S})q(\mathbf{W})q(\boldsymbol{\delta}).$$

For numerical stability and speed, we specify each of these variational distributions to be log normal for the CPLVM:

$$q(\mathbf{z}_i) = \mathrm{LogNormal}(\boldsymbol{\mu}_1, \sigma_1^2 I_{k_1}), \quad q(\mathbf{z}_j) = \mathrm{LogNormal}(\boldsymbol{\mu}_2, \sigma_2^2 I_{k_1}), \quad q(\mathbf{t}_j) = \mathrm{LogNormal}(\boldsymbol{\mu}_3, \sigma_3^2 I_{k_2}),$$
$$q(\mathbf{S}_l) = \mathrm{LogNormal}(\boldsymbol{\mu}_4, \sigma_4^2 I_p), \quad q(\mathbf{W}_d) = \mathrm{LogNormal}(\boldsymbol{\mu}_5, \sigma_5^2 I_p), \quad q(\boldsymbol{\delta}) = \mathrm{LogNormal}(\boldsymbol{\mu}_6, \sigma_6^2 I_p).$$

We perform approximate inference by minimizing the Kullback-Leibler (KL) divergence between the true posterior and the approximate posterior with respect to the variational parameters. This is equivalent to maximizing the lower bound on the log marginal likelihood of the data, known as the evidence lower bound (ELBO):

$$\log p(\mathbf{Y}, \mathbf{X}) \geq \mathbb{E}_{\mathcal{Z} \sim q(\mathcal{Z})} \left[ \log \frac{q(\mathcal{Z})}{p(\mathbf{Y}, \mathbf{X}, \mathcal{Z})} \right], \tag{10}$$

where $\mathcal{Z} = \{\mathbf{Z}^{\mathrm{b}}, \mathbf{Z}^{\mathrm{f}}, \mathbf{T}, \mathbf{S}, \mathbf{W}, \boldsymbol{\delta}\}$.

We use stochastic gradient descent to minimize the negative ELBO (Hoffman et al., 2013). We define and fit the variational model using TensorFlow probability (Dillon et al., 2017).

### 3.3 Contrastive generalized latent variable model

We also develop a second CPLVM that allows the factors to be negative by leveraging the exponential family distributions. In this case, we use a log-link function to transform the linear predictors to $\mathbb{R}_+$, similar to a generalized linear model (GLM). We call this model the contrastive generalized latent variable model (CGLVM). Here, in place of the multiplicative scale terms $\boldsymbol{\delta}$ in the CPLVM, we use additive coefficient terms $\boldsymbol{\mu}^{\mathrm{f}}, \boldsymbol{\mu}^{\mathrm{b}} \in \mathbb{R}^p$, similar to a traditional GLM.

$$\mathbf{y}_i|\mathbf{z}_i \sim \mathrm{Poisson}\left(\exp\left\{\left(\mathbf{S}^\top \mathbf{z}_i^{\mathrm{b}} + \boldsymbol{\mu}^{\mathrm{f}} + \log \alpha_i^{\mathrm{b}}\right)\right\}\right) \tag{11}$$

$$\mathbf{x}_j|\mathbf{z}_j, \mathbf{t}_j \sim \mathrm{Poisson}\left(\exp\left\{\left(\mathbf{S}^\top \mathbf{z}_j^{\mathrm{f}} + \mathbf{W}^\top \mathbf{t}_j + \boldsymbol{\mu}^{\mathrm{b}} + \log \alpha_j^{\mathrm{f}}\right)\right\}\right) \tag{12}$$

$$\mathbf{z}_i^{\mathrm{b}} \sim \mathcal{N}(\mathbf{0}, I), \quad \mathbf{z}_j^{\mathrm{f}} \sim \mathcal{N}(\mathbf{0}, I), \quad \mathbf{t}_j \sim \mathcal{N}(\mathbf{0}, I) \tag{13}$$

$$\mathbf{W}_d \sim \mathcal{N}(0, \mathbf{I}_p), \quad \mathbf{S}_l \sim \mathcal{N}(0, \mathbf{I}_p), \quad \boldsymbol{\mu}^{\mathrm{b}}, \boldsymbol{\mu}^{\mathrm{f}} \sim \mathcal{N}(0, \mathbf{I}_p) \tag{14}$$

where $l \in \{1, \ldots, k_1\}$, $d \in \{1, \ldots, k_2\}$. We place log normal priors on the size factors $\alpha_i^{\mathrm{b}}, \alpha_j^{\mathrm{f}}$, similar to those for the CPLVM.

|            | Contrastive | LVM | Background | Orthogonal | Counts | Nonnegative |
|:----------:|:-----------:|:---:|:----------:|:----------:|:------:|:-----------:|
| PCA        |             |     |            | x          |        |             |
| PPCA       |             | x   |            | x          |        |             |
| NMF        |             |     |            |            |        | x           |
| CPCA       | x           |     |            | x          |        |             |
| PCPCA      | x           | x   |            | x          |        |             |
| CLVM       | x           | x   | x          |            |        |             |
| CGLVM (ours) | x         | x   | x          |            | x      |             |
| CPLVM (ours) | x         | x   | x          |            | x      | x           |

Table 1: **Overview of related dimension reduction methods.** *Contrastive*: Method directly models the contrast between two datasets. *LVM*: Method has a latent variable model formulation. *Background*: Method includes an explicit model of of the background data. *Orthogonal*: Method constrains factors to be orthogonal to one another. *Counts*: Method directly accounts for count-based data. *Nonnegative*: Method has nonnegative factors and loadings.

For inference in the CGLVM, we again use a variational approach. Here, we specify the variational distributions as multivariate Gaussians:

$$q(\mathbf{z}_i) = \mathcal{N}(\boldsymbol{\mu}_1, \sigma_1^2 I_{k_1}), \quad q(\mathbf{z}_j) = \mathcal{N}(\boldsymbol{\mu}_2, \sigma_2^2 I_{k_1}), \quad q(\mathbf{t}_j) = \mathcal{N}(\boldsymbol{\mu}_3, \sigma_3^2 I_{k_2}),$$
$$q(\mathbf{S}_l) = \mathcal{N}(\boldsymbol{\mu}_4, \sigma_4^2 I_p), \quad q(\mathbf{W}_d) = \mathcal{N}(\boldsymbol{\mu}_5, \sigma_5^2 I_p), \quad q(\boldsymbol{\mu}^{\mathrm{b}}), q(\boldsymbol{\mu}^{\mathrm{f}}) = \mathcal{N}(\boldsymbol{\mu}_6, \sigma_6^2 I_p).$$

Similar to the CPLVM, we apply stochastic variational inference to optimize the ELBO.

### 3.4   Hypothesis testing with CPLVMs

In addition to exploratory data analysis using the model defined above, the CPLVM framework allows us to test hypotheses about whether the transcriptional structure is altered between conditions. Specifically, these tests quantify the extent to which the model's goodness of fit is improved when including the foreground-specific latent variables in addition to the shared latent variables.

The Bayesian framework of our model allows for model comparison using Bayes factors (Goodman, 1999). Bayes factors compare the ratio of data log likelihoods between an alternative model $\mathcal{M}_1$ and a null model $\mathcal{M}_0$, integrating over model parameters. Specifically, the Bayes factor is the ratio of model evidence (or marginal likelihoods):

$$\mathrm{BF} = \frac{p(\mathbf{Y}, \mathbf{X} | \mathcal{M}_1)}{p(\mathbf{Y}, \mathbf{X} | \mathcal{M}_0)} = \frac{\int_{\Theta_1} p(\mathbf{Y}, \mathbf{X} | \theta_1, \mathcal{M}_1) p(\theta_1 | \mathcal{M}_1) d\theta_1}{\int_{\Theta_0} p(\mathbf{Y}, \mathbf{X} | \theta_0, \mathcal{M}_0) p(\theta_0 | \mathcal{M}_0) d\theta_0}.$$

In practice, the log-Bayes factor is often used, and the null hypothesis is rejected if it surpasses some threshold $\tau$:

$$\text{Reject } H_0 \Leftrightarrow \log p(\mathbf{Y}, \mathbf{X}|\mathcal{M}_1) - \log p(\mathbf{Y}, \mathbf{X}|\mathcal{M}_0) > \tau.$$

Here, we implicitly assume equal prior weight on the null and alternative hypotheses, $p(\mathcal{M}_0) = p(\mathcal{M}_1)$, which we find to be well-calibrated in our numerical experiments.

Selecting a proper value for $\tau$ depends on the application area (Kass and Raftery, 1995). In practical settings, often $\tau$ is chosen based on a frequentist calibration of the hypothesis test.

In general, computing the model evidence requires solving an intractable integral, in turn making Bayes factors difficult to estimate. Following previous work (Lopez et al., 2018), we address this issue by approximating the model evidence with the ELBO (Equation (10)), which is a lower bound on the evidence, leading to ELBO-based Bayes factors (EBFs). We note that the tightness of the lower bound on the evidence depends on a number of modeling choices — for example, the choice of variational families and parameter initialization for stochastic VI — and the gap between the ELBO and the evidence could be different for the numerator and the denominator. However, we find these EBFs to be reliable in a number of simulations for our models. Thus, the general form of the CPLVM hypothesis test is the following:

$$\text{Reject } H_0 \Leftrightarrow \text{ELBO}_{\mathcal{M}_1} - \text{ELBO}_{\mathcal{M}_0} > \tau. \tag{15}$$

Defining the null and alternative models depends on the hypothesis of interest. Here, we consider two types of hypotheses for our model: *global* hypotheses and *gene set* hypotheses.

### 3.4.1 Global hypothesis test for changes in gene covariance structure

In this paper, a global hypothesis test is one that considers changes in expression across all genes that occur between conditions. This type of test is useful for assessing the effect of interventions that are expected to impact the expression of many genes, and the covariance structure among those genes.

For these hypotheses, we propose constructing the null model by removing the latent variables specific to the foreground data $\mathbf{t}_j$. Hence, the null model for the global hypothesis is

$$\mathbf{y}_i|\mathbf{z}_i \sim \text{Poisson}\left(\alpha_i^{\text{b}}\boldsymbol{\delta} \odot \mathbf{S}\mathbf{z}_i^{\text{b}}\right) \tag{16}$$

$$\mathbf{x}_j|\mathbf{z}_j \sim \text{Poisson}\left(\alpha_j^{\text{f}}\left(\mathbf{S}\mathbf{z}_j^{\text{f}}\right)\right) \tag{17}$$

$$\mathbf{z}_i^{\text{b}} \sim \text{Gamma}(\gamma, \beta_1), \quad \mathbf{z}_j^{\text{f}} \sim \text{Gamma}(\gamma, \beta_2). \tag{18}$$

Intuitively, the null model assumes that the latent structure in both matrices can be captured using a single, shared latent space $S$, and the samples from both matrices are projected onto that latent space using $\mathbf{z}_j^{\text{b}}$ and $\mathbf{z}_j^{\text{f}}$ for background and foreground data, respectively.

The alternative hypothesis is that there is structure specific to the foreground data across all genes. For the alternative model, we use the full CPLVM defined in Equations (6)-(9). This model includes the latent variables specific to the foreground data $\mathbf{t}_j$.

### 3.4.2 Gene-set hypothesis test for foreground changes to a subset of genes

To test for changes in variation in foreground data relative to background data involving only a subset of genes, we propose a gene set hypothesis test. This test is useful to quantify support for changes to joint expression within specific known gene modules that are unique to the foreground matrix.

Suppose we would like to test for differential variation in a set of $G$ genes indexed by $\{\ell_1, \ldots, \ell_G\}$. In this case, the null hypothesis is encoded in a model constructed by constraining rows of $\mathbf{W}$ indexed by $\{\ell_1, \ldots, \ell_G\}$ to be zero. Intuitively, this model assumes no change in variation specific to that gene set in the foreground matrix relative to the background matrix. To be precise, let $\mathbf{Q}$ be a $G \times p$ binary matrix which, when multiplied with $\mathbf{W}$, only takes rows corresponding to genes in the gene set. Specifically, for $g \in \{1, \ldots, G\}$ and $k \in \{1, \ldots, p\}$,

$$\mathbf{Q}_{gk} = \begin{cases} 1 & k = \ell_g \\ 0 & \text{otherwise.} \end{cases}$$

Then the null model is

$$\mathbf{y}_i | \mathbf{z}_i \sim \text{Poisson}\left( \alpha_i^{\mathrm{b}} \boldsymbol{\delta} \odot \mathbf{S} \mathbf{z}_i^{\mathrm{b}} \right) \tag{19}$$

$$\mathbf{x}_j | \mathbf{z}_j \sim \text{Poisson}\left( \alpha_j^{\mathrm{f}} \left( \mathbf{S} \mathbf{z}_j^{\mathrm{f}} + \mathbf{W} \mathbf{t}_j \right) \right) \quad \text{s.t.} \quad \mathbf{Q}\mathbf{W} = \mathbf{0}_G, \tag{20}$$

where $\mathbf{0}_G$ is the zero vector of length $G$. We specify the alternative hypothesis in a model that includes the full CPLVM as described in (6)-(9).

A gene set hypothesis will test whether there are changes in the covariance structure specific to a subset of genes that are prespecified. To define gene sets of interest, one could use established gene set collections, such as the MSigDB sets (Liberzon et al., 2015), as we show in our experiments below.

## 4  Simulation results

In this section, we evaluate the performance of the CGLVM and CPLVM using synthetic data. We compare the performance of these models to four related state-of-the-art methods for dimension reduction (PCA, NMF, CPCA, and PCPCA), a related linear model for differential expression discovery, and a two sample test for differential expression testing.

## 4.1 Visualizing CGLVM and CPLVM latent spaces

In order to demonstrate the behavior of the CGLVM and CPLVM, we first fit the model on simple two-dimensional simulated count data. In this dataset, the foreground data is made up of two subgroups, while the background data is homogeneous. Specifically, to generate count data with a prespecified covariance matrix $\Sigma$, we use a Gaussian copula with a Poisson likelihood. In particular, we generate the foreground samples $\{\mathbf{x}_i\}_{i=1}^n$ and background samples $\{\mathbf{y}_j\}_{j=1}^m$ as

$$\mathbf{x}_i = F_\lambda^{-1}(\Phi(z_i))$$
$$\mathbf{y}_j = F_\lambda^{-1}(\Phi(z_j))$$

where $F_\lambda^{-1}$ is the inverse CDF of a Poisson with parameter $\lambda$, $\Phi$ is the standard normal CDF, and $z_i, z_j \sim \mathcal{N}(0, \Sigma)$. We then shift the subgroups to give the background a mean of $\binom{14}{14}$ and the foreground subgroups means of $\binom{18}{10}$ and $\binom{10}{18}$. Here, we set $\Sigma = \left(\begin{smallmatrix} 2.7 & 2.6 \\ 2.6 & 2.7 \end{smallmatrix}\right)$, $\lambda = 10$, and $n = m = 200$.

We fit the CGLVM and CPLVM, as well as the related methods listed above, on this dataset. For the CGLVM, we use a single latent dimension for the shared and foreground-specific compartments, $k_1 = k_2 = 1$. For the CPLVM, we set $k_1 = 1, k_2 = 2$ in order to allow the nonnegative factors to discover negative associations.

We find that, in both the CGLVM and CPLVM, the subspaces defined by $\mathbf{S}$ and $\mathbf{W}$ picked up on the directions of shared and foreground-specific variation, respectively (Figure 1e, f). The other contrastive methods, CPCA and PCPCA, were able to detect the axis of variation unique to the foreground data, but these approaches do not have an explicit model for the background data (Figure 1c, d). Finally, PCA and NMF — which do not model the contrast between the conditions — are unable to identify the axis that separates the two foreground subgroups (Figure 1a, b).

This result suggests that the CPLVM is able to disentangle these two sources of variation: those that are shared between conditions, and those that are unique to the foreground.

## 4.2 Discovering heterogeneous responses

We next examined whether the CPLVM discovers the latent structure of a dataset in which there is a heterogeneous response across samples in the foreground data. To study this behavior, we generate a synthetic count dataset from the CPLVM model (6)-(9). We set the parameters such that all background cells have the same latent state, but each foreground cell is drawn from one of two unique latent states. Specifically, we set the data dimension $p = 100$, and the latent dimensions as $k_1 = k_2 = 2$. Furthermore, we set $\beta_1 = \beta_2 = \beta_4 = \beta_5 = 1$ and $\gamma_1 = \gamma_2 = \gamma_4 = \gamma_5 = 1$. For half of the foreground cells, we sampled $t_{j1} \sim \text{Gamma}(1,1)$ and $t_{j2} \sim \text{Gamma}(1,.01)$. For the other half, we sampled $t_{j1} \sim \text{Gamma}(1,.01)$ and $t_{j2} \sim \text{Gamma}(1,1)$.

Figure 1: **Illustration of the CPLVM with toy data.** Related dimension reduction methods applied to a toy dataset in which the foreground data contains two subgroups. (a) PCA ($k = 1$), (b) NMF ($k = 2$), (c) CPCA ($k = 1$), (d) PCPCA ($k = 1$), (e) CGLVM (ours, $k_1 = k_2 = 1$), (f) CPLVM (ours, $k_1 = 1, k_2 = 2$). In each, we plot the one-dimensional line defined by each column of the estimated loadings matrix from each method.

We fit the CPLVM on this dataset and examine its estimated latent projections of the foreground cells (Figure 2d). For comparison, we also visualize the latent projections of these cells under PCA and CPCA (Figure 2b, c). To quantify whether the two foreground subgroups are preserved in the latent space, we compute the silhouette score for the latent variables relative to the true cluster identities. As benchmarks, we also compute the silhouette score for PCA, NMF, CPCA, and CGLVM (Figure 2e).

We found that the CPLVM's latent space was able to recover the structure of the two subgroups in the foreground (Figure 2d). In contrast, PCA and CPCA were not able to capture this two-cluster response as clearly (Figure 2b, c). Moreover, the cluster analysis revealed that the CPLVM was better able to retain the subgroup structure compared to PCA, NMF, CPCA, and CGLVM (Figure 2e).

To further test the goodness-of-fit of our models, we quantified their ability to recover

Figure 2: **Cluster identification in simulated data with the CPLVM.** The foreground data was generated from two subgroups of samples. (a) The true underlying foreground-specific latent variables. (b) PCA does not separate the two clusters. (c) CPCA shows an improvement over PCA, but still has overlap in subgroups in the reconstructed data. (d) The CPLVM is able to capture the difference between the subgroups, as well as better preserving nearest-neighbor relationships. (e) Silhouette scores computed on the foreground latent variables for competing methods.

relationships between samples in their latent spaces. To do this, we generated count data from a small set of latent factors. We then fit four models — PCA, CPCA, CGLVM, and CPLVM — and computed the distance between each method's recovered latent variables and the true latent variables. Specifically, we computed the Wasserstein distance between normalized pairwise distance matrices of the simulated and estimated latent variables. We repeated this ten times for each method. We found that the CGLVM and CPLVM outperformed PCA and CPCA in reconstruction error (Figure 3a, b). The relative performance of the CGLVM and CPLVM was even more noticeable in the background samples, likely due to the CGLVM and CPLVM's explicit models of the background, which PCA and CPCA do not have. In both metrics, the CPLVM outperforms the CGLVM slightly. These results suggest that the CPLVM captures variation in foreground count data relative to background count data, and enables subgroup discovery within the foreground data.

## 4.3 Estimating the latent dimension

Next, we asked whether the CPLVM could estimate the dimension of the generative low-dimensional space. To test this, we generated data from the CPLVM with $k_1 = k_2 = 5$. We then fit a series of CPLVMs, each with a different latent dimension ranging from $k_1 = k_2 = 1$ to $k_1 = k_2 = 15$ with $k_1 = k_2$ in all cases. We measured the quality of each model's fit to the data by computing the ELBO for each fit, repeating this procedure ten times for each latent dimension. We found that the ELBO peaked near the true latent dimension $k = 5$ (Figure 3b) and that the ELBO was lower for models with over- or under-constrained latent spaces (corresponding to latent dimensions that were higher or lower than the true dimension, respectively). This result suggests that the CPLVM can recover the true complexity of the variation in the data, and that the ELBO can be used as a reasonable measure of the model's

Figure 3: **Simulation experiments with the CPLVM.** We fit our contrastive models to data generated from a small set of shared and foreground-specific latent variables. (a) The average Wasserstein-2 distance between the estimated and true pairwise distances between samples in the foreground for each method. (b) Same as (a), but for background samples. (c) The ELBO for our CPLVM with a range of latent dimensions. The true latent structure of the simulated data is shown by the vertical dotted line. Vertical lines show 95% confidence intervals.

fit to the data.

## 4.4 Hypothesis testing

Next, we examined whether the CPLVM hypothesis testing framework detects changes in variation between conditions. We consider two types of changes in variation that are found in scientific data: Global shifts in variation across all features, and changes to variation specific to a subset of features (Chandrasekaran et al., 2009; Leek and Storey, 2008).

### 4.4.1 Global hypothesis tests

To evaluate the global hypothesis testing framework, we generated three datasets: one "alternative" dataset simulating true global change in variation between conditions and two "null" datasets simulating no change between conditions. The alternative dataset — which we call the *perturbed dataset* — was drawn from the alternative model defined in (6)-(8) such that there was substantial change in variation across most genes. The first null dataset — which we call the *unperturbed null dataset* — was drawn from the null CPLVM in Equations (16)-(18) such that there was no change in variation between conditions. The second null dataset — which we call the *shuffled null dataset* — was constructed from the samples in the *perturbed dataset* by randomly reassigning cells to the background and foreground conditions. These datasets allow us to calibrate the Bayes factors for a truly alternative dataset relative to two truly null datasets. The shuffled dataset is intended to emulate a real-world scenario, in which calibration relative to a true null is not possible.

We computed EBFs for a global hypothesis test for each of the three datasets. For each dataset, we fit the null and alternative models defined in Equations (16)-(18) and

Figure 4: **Hypothesis testing on simulated data with the CPLVM.** (a) Global hypothesis testing with data generated from a null model (left box), shuffled data approximating truly null data (middle box), and data generated from an alternative model (right box). (b) Gene set hypothesis tests with data in which only variation among genes in gene set one has been altered between conditions (indicated by the red box).

Equations (6)-(8), respectively, and computed the EBFs as in Equation (15). We repeated this procedure ten times.

We found that the EBFs for the *perturbed dataset* were all substantially above zero. These EBFs were also higher than the EBFs for either of the truly null datasets, indicating a consistently higher lower bound on the model evidence for the alternative model on the *perturbed dataset* (Figure 4a). The EBFs for the *unperturbed null dataset* were all below zero, implying that the model evidence did not favor the alternative model in this case. The *shuffled null dataset* showed Bayes factors that were between the other two datasets, but distinct from them both. Using the same datasets, we found that the CGLVM was similarly well-calibrated (Supplementary Fig 1). This implies that, for global hypothesis testing, the *shuffled null dataset* can be used as the empirical null to calibrate EBFs in practice.

To assess the reliability of the hypothesis testing framework, we quantified how frequently the CPLVM correctly rejected the null hypothesis. To do this, we classified each Bayes factor as "accept $H_0$" or "reject $H_0$" for a range of thresholds $\tau_1, \ldots, \tau_t$, where

$$\text{Reject } H_0 \Leftrightarrow \text{ELBO}_{\mathcal{M}_1} - \text{ELBO}_{\mathcal{M}_0} > \tau_i. \tag{21}$$

Using this decision rule, we then estimated the true positive rate (TPR) and false positive rate (FPR) for each threshold $\tau_i$:

$$\text{TPR}_i = \mathbb{P}(\text{reject } H_0 | H_1)$$
$$\text{FPR}_i = \mathbb{P}(\text{reject } H_0 | H_0).$$

To be precise, the *TPR* is the probability of correctly rejecting the null hypothesis (also called the statistical power), and the *FPR* is the probability of incorrectly rejecting the null hypothesis.

To quantify how the CPLVM performs under these metrics, we generated data from the CPLVM (6)-(9). In particular, we sampled data with three different data dimensions, $p \in \{10, 100, 1000\}$, creating 50 datasets for each value of $p$. For each setting of $p$, we then computed EBFs for the corresponding datasets. For each dataset, we created a corresponding negative control, or "null", dataset by shuffling the foreground or background labels of the samples. Finally, at a range of thresholds $\tau$, we accepted or rejected the null hypothesis for each dataset based on the decision rule in (21). Finally, we computed the TPR and FPR for each value of $\tau$, and we computed the corresponding ROC curves (Figure 5).

For comparison, we computed the same metrics for a competing two-sample covariance matrix test (Cai et al., 2013). This approach tests whether the foreground and background covariance matrices are equal,

$$H_0 : \Sigma_x = \Sigma_y \quad \text{versus} \quad H_1 : \Sigma_x \neq \Sigma_y.$$

This procedure computes a test statistic,

$$M_n = \max_{1 \leq i \leq j \leq p} \frac{(\widehat{\sigma}_{kl}^{\mathrm{f}} - \widehat{\sigma}_{kl}^{\mathrm{b}})^2}{\widehat{\theta}_{kl}^{\mathrm{f}}/n + \widehat{\theta}_{kl}^{\mathrm{b}}/m}$$

where $\widehat{\sigma}_{kl}^{\mathrm{f}}$ and $\widehat{\sigma}_{kl}^{\mathrm{b}}$ are the covariance between features $k$ and $l$ in the foreground and background, respectively, and $\widehat{\theta}_{kl}^{\mathrm{f}}$ and $\widehat{\theta}_{kl}^{\mathrm{b}}$ are the variance of the covariance elements,

$$\widehat{\theta}_{kl}^{\mathrm{f}} = \frac{1}{n} \sum_{i=1}^{n} \left[ (X_{ki} - \bar{X}_k)(X_{li} - \bar{X}_l) - \widehat{\sigma}_{kl}^{\mathrm{f}} \right]^2$$

$$\widehat{\theta}_{kl}^{\mathrm{b}} = \frac{1}{m} \sum_{i=1}^{n} \left[ (Y_{kj} - \bar{Y}_k)(Y_{lj} - \bar{Y}_l) - \widehat{\sigma}_{kl}^{\mathrm{b}} \right]^2.$$

Here, $\bar{X}, \bar{Y} \in \mathbb{R}^p$ are vectors of sample means. Based on the limiting distribution of $M_n$, the decision rule for this test at level $\alpha$ is

$$\text{Reject } H_0 \Leftrightarrow M_n \geq q_\alpha + 4 \log p - \log \log p$$

where $q_\alpha$ is the $1 - \alpha$ quantile of the Type I extreme value distribution (Gumbel distribution) with cumulative distribution function $F(x) = \exp\left(-\frac{1}{\sqrt{8\pi}} \exp(-x/2)\right)$.

The ROC curves for these data settings show that the CPLVM test consistently outperforms Cai's two-sample covariance test (Figure 5). For the CPLVM test, we found that the TPR and FPR remained strong across each data setting, performing perfectly for $p \in \{100, 1000\}$. For the datasets with $p = 10$, the CPLVM test did not achieve perfect

Figure 5: **Benchmarking the global hypothesis test.** Using simulated data with varying data dimensionalities ($p \in \{10, 100, 1000\}$), we computed ROC curves based on the CPLVM's rejection or acceptance of the null (orange curves). For comparison, we computed the same metrics for Cai's two-sample covariance test that relies on explicitly computing the full sample covariance matrix (blue curves, Cai et al. 2013).

TPR and FPR, but still performed well above random. In contrast, the two-sample covariance test consistently performed worse than the CPLVM, and indeed performed no better than random guessing with $p = 10$ (Figure 5a). Moreover, the two-sample covariance test had substantially lower TPRs and FPRs than the CPLVM for $p \in \{100, 1000\}$.

These results demonstrate that the CPLVM hypothesis testing framework is able to reliably detect global changes in variation between conditions. Furthermore, the analysis suggests that shuffling cell condition labels is a viable strategy for calibrating the EBFs.

### 4.4.2   Gene set hypothesis tests

To test the gene set hypothesis testing framework, we created ten synthetic genes sets, each made up of 25 genes. We arbitrarily designated the first gene set — whose genes are indexed by $\{1, \ldots, 25\}$ — as the *perturbed gene set*. In other words, this gene set was chosen to show substantial change in variation between the background and foreground conditions. We simulated data for these genes using the CPLVM (6)-(9). All other gene sets were designed to not show substantial variation between conditions. We call these truly null gene sets the *unperturbed gene sets*, and we simulated these with the CPLVM corresponding to the null gene set hypothesis (Equations (19)-(20)). We also included 250 genes that did not belong to a gene set, which were also simulated from the null gene set model. This led to a total of $p = 500$ genes, half of which belong to gene sets.

To calibrate the gene set EBFs, we estimated an empirical null distribution of EBFs by creating gene sets with randomly assigned genes. In particular, for the 250 genes belonging to gene sets, we randomly reassigned each of them to synthetic gene sets of size 25, repeating this 50 times to create 50 new synthetic, shuffled gene sets. These gene sets — which we

call *shuffled null gene sets* — are useful because the true null distribution of EBFs is not available in practice.

For each gene set, we fit the null and alternative models described by (19)-(20) and (6)-(8), respectively, and computed EBFs for each model. We repeated this experiment five times, with each iteration yielding ten EBFs (one for each gene set).

We found that the *perturbed gene set* showed consistently higher EBFs than all other gene sets (Figure 4b). Furthermore, all EBFs for the *perturbed gene set* were positive, while most other gene sets were consistently negative or near zero. We also found that the EBFs for the *shuffled null* gene sets were also consistently below the *perturbed gene set*, indicating that the EBFs are well-calibrated.

To further test the robustness of the gene set hypothesis tests, we ran the tests in two other simulation settings. First, we tested the robustness of the EBFs to the size of the gene sets. To do this, we generated a similar dataset as before — containing 500 genes, 250 of which belong to gene sets — but this time we varied the number of genes in each gene set to be in the set $\{1, 5, 10, 15, 20, 25\}$. As expected the EBFs gradually declined when the *perturbed gene set* contained fewer genes (Supplementary Fig 2b). However, the test remained robust even for gene sets containing as few as 5 genes.

Second, we ran the hypothesis test in a setting in which the gene sets were misspecified. In particular, we again constructed gene sets of size 25, but here only 12 of the genes in the *perturbed gene set* truly showed a difference between conditions. Even when the gene sets were misspecified as such, we found that the EBFs for the *perturbed gene set* remained substantially above those of the unperturbed gene sets (Supplementary Supplementary Fig 2a).

Together, these results imply that the CPLVM gene set hypothesis tests can detect targeted, pathway-specific changes between conditions.

# 5   Application to Perturb-seq data

Next, we applied our models to data from the Perturb-seq platform (Adamson et al., 2016; Dixit et al., 2016).

## 5.1   Data

Perturb-seq is a scRNA-seq platform designed to measure the RNA transcript levels in cells that have been exposed to a set of CRISPR lentivirus guides (Adamson et al., 2016; Dixit et al., 2016). Each guide targets a specific gene, deactivating it by "cutting" it out of the genome using a Cas9 nuclease.

In our experiments, for the foreground data, we leveraged Perturb-seq data that contains scRNA-seq measurements on pools of bone marrow-derived dendritic cells (BMDCs), each of which was infected with a unique CRISPR guide (Dixit et al., 2016). Each CRISPR guide in this study was designed to target one of 24 unique transcription factors. For the background data, we use control data from cells that did not receive any treatment. To

preprocess the data, for each targeted gene, we pooled data from all CRISPR guides that target that gene. We subsetted each experiment to the 500 most variable genes, according to the Poisson deviance (Supplementary material, Townes et al. 2019).

We fit the CPLVM separately to the datasets from each of the 24 experiments, using the transcript counts from the untreated cells as the background data $\mathbf{Y}$ and the counts from CRISPR-treated cells as the foreground data $\mathbf{X}$.

## 5.2    Identifying covariance shifts in Perturb-seq data

To evaluate the CPLVM's ability to capture shifts in variation in Perturb-seq data, we fit the model for each of the 24 experiments. For comparison, we also fit a Poisson GLM (4)-(5) that only identifies changes in the marginal distribution of each gene.

Examining the CPLVM's latent factors, we found that they identified several shifts in gene-gene covariation that were not picked up by the GLM. One such instance was observed in the *HIF1A*-perturbed experiment. Here, we found that two genes (*LYZ2* and *CCL4*) showed positive correlation across cells in the foreground data, but no correlation in the background data. The CPLVM captured this gene-gene relationship in one of its components (Figure 6c), while the GLM failed to detect this relationship. Instead, the GLM identified a shift in the marginal expression of *CCL4* alone.

This result suggests that the CPLVM is useful for identifying shifts in variation that occur across multiple genes, and that univariate linear models are not able to detect this type of change.

## 5.3    Perturb-seq global hypothesis tests

Next, we sought to more broadly explore the main sources of variation in each Perturb-seq experiment, and the extent to which each guide induced a substantial change in expression patterns. To do so, we first evaluated the magnitude of the overall change in variation by running global hypothesis tests for each experiment. We computed global EBFs for each (Figure 7a). To calibrate each test, we also computed EBFs for a second dataset in which cells were randomly reassigned to the foreground or background condition. This shuffled dataset is intended to remove any biologically meaningful patterns that are specific to the foreground data.

Examining the global EBFs, this analysis revealed that most of the experiments showed substantial change in gene expression variation between the untreated and treated conditions. This suggests that most of the guides used in this study had an effect on transcription levels globally across genes, which is expected for these transcription factors. The EBFs for the shuffled datasets were also mostly positive, which was expected from the simulation experiments. However, the EBFs from the shuffled data were consistently lower than their corresponding global EBFs.

Figure 6: **CPLVM applied to Perturb-seq data.** (a) Scatter plot of $\log(x+1)$ expression for *LYZ2* and *CCL4* in the*HIF1A* experiment. *CCL4* shows a positive shift in its marginal expression between conditions, but *LYZ2* does not. However, the correlation between these two genes changes between conditions (Pearson $\rho = 0.12$ in the background, and $\rho = 0.74$ in the foreground). (b) GLM coefficients for the *HIF1A* experiment. Only *CCL4* is identified as differentially expressed. (c) CPLVM loadings from one CPLVM component for the *HIF1A* experiment. Both *CCL4* and *LYZ2* are identified as having differential variation in this component.

## 5.4 Perturb-seq gene set hypothesis tests

To more narrowly characterize the variation in the Perturb-seq experiments, we performed a series of gene set hypothesis tests. To do this, we leveraged the MSigDB Hallmark gene sets, which categorize genes into a collection of established pathways (Liberzon et al., 2015). For each experiment, we computed the EBF for each Hallmark gene set.

Many gene sets emerged as perturbed from this analysis. For example, in the *HIF1A*-perturbed experiment, a number of coordinated gene sets appeared as top hits, including *TNF-α signaling* and *inflammatory response* (Figure 7b). Moreover, we found that the magnitude of the gene set EBFs were not correlated with the size of the gene sets (Pearson $\rho = 0.03$), suggesting that the tests were not biased by the sizes of the gene sets. These gene set hypothesis test results suggest that the CPLVM is able to identify coordinated changes in gene expression even among small sets of genes.

# 6 Application to small molecule perturbation data

As further investigation of the CPLVM's behavior on real data, we next applied our model to a scRNA-seq dataset from the MIX-seq platform (McFarland et al., 2020).

Figure 7: **Hypothesis testing with Perturb-seq data.** (a) Global hypothesis tests for Perturb-seq experiments. Blue bars represent EBFs for each experiment, and red bars are the EBFs for the shuffled data. Vertical ticks represent 95% confidence intervals. (b) Gene set EBFs for the *HIF1A* experiment.

## 6.1 Data

The MIX-seq platform provides scRNA-seq readouts of cancer cell lines' transcriptional responses after being treated with a panel of small molecule therapies (McFarland et al., 2020). We used a MIX-seq dataset that contains data for 24 cell lines, and we focused on an experiment in which the cells were exposed to idasanutlin, which inhibits the activity of *MDM2*. *MDM2* is known negatively regulate the tumor-suppressor gene *TP53* (Vassilev et al., 2004). Furthermore, idasanutlin has been shown to elicit a selective transcriptional and death response in cells that have wild-type *TP53*, while cells with a mutated copy of this gene do not respond (McFarland et al., 2020).

## 6.2 Application to idasanutlin data

We fit the CPLVM to the MIX-seq data and analyzed the fitted parameters. We used the transcript counts from idasanutlin-treated cells as the foreground matrix and a the counts from a pool of DMSO-treated cells as the background matrix. In the CPLVM model, we set $k_1 = k_2 = 2$ for visualization.

Visualizing the foreground-specific latent variables for each cell, we found that the CPLVM factors were able to partially separate cells with mutated *TP53* and cells with wild-type *TP53* (Figure 8b). Meanwhile, a PCA projection of the foreground cells did not clearly identify this subgroup structure (Figure 8a). A cluster analysis found that the CPLVM latent variables showed tighter clustering of these subgroups compared to PCA (Figure 8c).

Furthermore, we ran the CPLVM gene set hypothesis tests on the idasanutlin data, again using the MSigDB Hallmark gene sets. This analysis revealed that the *P53 pathway* gene set was among the top enriched pathways (Figure 8d). This observation coincides with the known mechanism of action of idasanutlin, namely, its direct effect on the *MSM2/TP53*

24

Figure 8: **Contrastive latent variable models applied to chemical perturbation data.** (a) PCA projection of the foreground cells from the idasanutlin experiment. Points (cells) are colored by their *TP53* mutation status. (b) Foreground cells projected into the foreground-specific latent space of the CPLVM. (c) Silhouette score for the clusters of TP53-mutated cells and wild-type cells in the PCA and CPLVM projections. (d) Top gene set EBFs for idasanutlin. The P53 pathway appears as the gene set with the second-highest EBF.

pathway (Vassilev et al., 2004; McFarland et al., 2020). These results suggests that the CPLVM and corresponding statistical test is able to accurately identify axes of heterogeneity in the response to chemical perturbations, and our model's representation can recover subgroup structure in the foreground data.

# 7 Application to GTEx data

Beyond perturbational data, the CPLVM can be used more generally for count datasets with two conditions. In this section, we demonstrate one such application using bulk RNA-seq data from the Genotype-Tissue Expression (GTEx) Consortium v8 study (Consortium et al., 2017, 2020).

The GTEx data contains gene expression measurements in a large number of tissues, collected from thousands of healthy donors. For this experiment, we focused on a subset of the data to answer a specific question: whether there are differences in gene expression variation in coronary artery tissue between donors with and without ischemic heart disease. To do this, we treated gene expression samples from donors with heart disease as the foreground matrix and samples from healthy donors as the background matrix. We subsetted to the 200 most variable genes, and fit the CPLVM on the RNA-seq counts for these samples.

Examining the foreground-specific components of the CPLVM, we found that one of the factors picked up on several genes related to oxygen intake (Figure 9b). In particular, the genes with the largest coefficients in this factor — *SFTPB*, *SFTPA2*, *SFTPC*, and *SFTPA1* — primarily belonged to the pulmonary surfactant protein complex. This complex is known to aid the lung and heart's oxygen-passing abilities.

Figure 9: **CPLVM applied to RNA-seq data from coronary artery tissue in patients with and without heart disease.** (a) Sorted loadings values for one component of the shared loadings matrix **S**. The top genes are related to typical heart function and heart muscle regulation. (b) Sorted loadings values for one component of the foreground-specific loadings matrix **W**. The top genes are related to oxygen delivery in the heart and lungs, a process that is dysregulated in ischemic artery disease.

Furthermore, we examined the parameters of the CPLVM that are shared between the foreground and background samples. We expect these factors to detect variation in expression that exists in both patients with and without heart disease. Indeed, we found that the genes with the highest loading values in one component — *MYH7*, *DES*, and *MYL2* — were related to basic heart functioning (Figure 9a). These results suggest that the CPLVM can be used for settings beyond perturbation experiments, such as for examining structural differences between biological conditions. It also implies that the CPLVM can be used to investigate the shared structure between conditions.

## 8   Discussion

In this study, we presented latent variable models, CPLVM and CGLVM, for case-control sequencing. These models capture the change in gene expression variation that is specific to the case condition, as well as the variation that exists in both the case and control conditions. Our modeling framework provides a set of low-dimensional latent factors that describe this variation. Furthermore, we provide a flexible hypothesis testing framework for characterizing transcriptional structure in case-control experiments.

Through a series of simulations and experiments with gene expression data, we showed that the CPLVM captures foreground-specific structure and structure that exists in both conditions. In simulations, we showed that that the CPLVM captures the transcriptional variation better than linear models, can identify the proper number of latent dimensions,

and enable reliable hypothesis testing of both global and pathway-specific shifts in gene expression. In the context of CRISPR- and drug-treated scRNA-seq data, we showed that CPLVMs can be used to generate biological insights and identify subgroup structure. These insights go beyond traditional differential expression measurements, enabling discovery of differential relationships between genes and cells, such as estimating changes in gene-gene correlations and identifying foreground-specific heterogeneity within a population of cells.

Several future directions remain to be explored. First, the modeling approach could be extended in several ways. Experiments with more than two conditions could be considered. For example, in gene expression datasets measured across several tissues, it could be useful to model the shared variation among the tissues, as well as the variation that is specific to each tissue (Consortium et al., 2020). Second, more complex inference schemes could be considered. While we used a mean-field variational approximation to the CPLVM, more flexible posterior approximations could be used, such as a variational autoencoder (VAE, Kingma and Ba 2014; Lopez et al. 2018). Finally, while our hypothesis testing procedure proved to be well-calibrated, several improvements could be made. Our method uses the ELBO to approximate Bayes factors, but better approximations to the marginal likelihood could be used. Furthermore, our procedure implicitly assigns equal prior weight to both hypotheses, $p(\mathcal{M}_0) = p(\mathcal{M}_1) = 0.5$. This choice proved to be robust in practice, but further investigation into the effect of choosing these priors is warranted.

## Acknowledgements

# References

Abid, A., Zhang, M. J., Bagaria, V. K., and Zou, J. (2017). Contrastive principal component analysis. *arXiv preprint arXiv:1709.06716*.

Abid, A., Zhang, M. J., Bagaria, V. K., and Zou, J. (2018). Exploring patterns enriched in a dataset with contrastive principal component analysis. *Nature Communications*, 9(1):1–7.

Adamson, B., Norman, T. M., Jost, M., Cho, M. Y., Nuñez, J. K., Chen, Y., Villalta, J. E., Gilbert, L. A., Horlbeck, M. A., Hein, M. Y., et al. (2016). A multiplexed single-cell crispr screening platform enables systematic dissection of the unfolded protein response. *Cell*, 167(7):1867–1882.

Anderson, T. W. (1958). *An introduction to multivariate statistical analysis*, volume 2. Wiley New York.

Aoshima, M. and Yata, K. (2018). Two-sample tests for high-dimension, strongly spiked eigenvalue models. *Statistica Sinica*, pages 43–62.

Becht, E., McInnes, L., Healy, J., Dutertre, C.-A., Kwok, I. W., Ng, L. G., Ginhoux, F., and Newell, E. W. (2019). Dimensionality reduction for visualizing single-cell data using umap. *Nature Biotechnology*, 37(1):38–44.

Boileau, P., Hejazi, N. S., and Dudoit, S. (2020). Exploring high-dimensional biological data with sparse contrastive principal component analysis. *Bioinformatics*, 36(11):3422–3430.

Cai, T., Liu, W., and Xia, Y. (2013). Two-sample covariance matrix testing and support recovery in high-dimensional and sparse settings. *Journal of the American Statistical Association*, 108(501):265–277.

Chandrasekaran, V., Sanghavi, S., Parrilo, P. A., and Willsky, A. S. (2009). Sparse and low-rank matrix decompositions. *IFAC Proceedings Volumes*, 42(10):1493–1498.

Consortium, G. et al. (2017). Genetic effects on gene expression across human tissues. *Nature*, 550(7675):204.

Consortium, G. et al. (2020). The gtex consortium atlas of genetic regulatory effects across human tissues. *Science*, 369(6509):1318–1330.

Delmans, M. and Hemberg, M. (2016). Discrete distributional differential expression ($d^3$e)- a tool for gene expression analysis of single-cell rna-seq data. *BMC Bioinformatics*, 17(1):110.

Dillon, J. V., Langmore, I., Tran, D., Brevdo, E., Vasudevan, S., Moore, D., Patton, B., Alemi, A., Hoffman, M., and Saurous, R. A. (2017). Tensorflow distributions. *arXiv preprint arXiv:1711.10604*.

Ding, J., Condon, A., and Shah, S. P. (2018). Interpretable dimensionality reduction of single cell transcriptome data with deep generative models. *Nature Communications*, 9(1):1–13.

Dixit, A., Parnas, O., Li, B., Chen, J., Fulco, C. P., Jerby-Arnon, L., Marjanovic, N. D., Dionne, D., Burks, T., Raychowdhury, R., et al. (2016). Perturb-seq: dissecting molecular circuits with scalable single-cell rna profiling of pooled genetic screens. *Cell*, 167(7):1853–1866.

Finak, G., McDavid, A., Yajima, M., Deng, J., Gersuk, V., Shalek, A. K., Slichter, C. K., Miller, H. W., McElrath, M. J., Prlic, M., et al. (2015). Mast: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell rna sequencing data. *Genome biology*, 16(1):1–13.

Glass, K., Huttenhower, C., Quackenbush, J., and Yuan, G.-C. (2013). Passing messages between biological networks to refine predicted interactions. *PloS One*, 8(5):e64832.

Goodman, S. N. (1999). Toward evidence-based medical statistics. 2: The bayes factor. *Annals of Internal Medicine*, 130(12):1005–1013.

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

Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. (2013). Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347.

Ishii, A., Yata, K., and Aoshima, M. (2019). Equality tests of high-dimensional covariance matrices under the strongly spiked eigenvalue model. *Journal of Statistical Planning and Inference*, 202:99–111.

Kass, R. E. and Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical Association*, 90(430):773–795.

Kharchenko, P. V., Silberstein, L., and Scadden, D. T. (2014). Bayesian approach to single-cell differential expression analysis. *Nature Methods*, 11(7):740–742.

Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Kinker, G. S., Greenwald, A. C., Tal, R., Orlova, Z., Cuoco, M. S., McFarland, J. M., Warren, A., Rodman, C., Roth, J. A., Bender, S. A., et al. (2020). Pan-cancer single-cell rna-seq identifies recurring programs of cellular heterogeneity. *Nature Genetics*, 52(11):1208–1218.

Korthauer, K. D., Chu, L.-F., Newton, M. A., Li, Y., Thomson, J., Stewart, R., and Kendziorski, C. (2016). A statistical approach for identifying differential distributions in single-cell rna-seq experiments. *Genome Biology*, 17(1):222.

Leek, J. T. and Storey, J. D. (2008). A general framework for multiple testing dependence. *Proceedings of the National Academy of Sciences*, 105(48):18718–18723.

Li, D., Jones, A., and Engelhardt, B. (2020). Probabilistic contrastive principal component analysis. *arXiv preprint arXiv:2012.07977*.

Li, J., Chen, S. X., et al. (2012). Two sample tests for high-dimensional covariance matrices. *The Annals of Statistics*, 40(2):908–940.

Liberzon, A., Birger, C., Thorvaldsdóttir, H., Ghandi, M., Mesirov, J. P., and Tamayo, P. (2015). The molecular signatures database hallmark gene set collection. *Cell Systems*, 1(6):417–425.

Lopez, R., Regier, J., Cole, M. B., Jordan, M. I., and Yosef, N. (2018). Deep generative modeling for single-cell transcriptomics. *Nature Methods*, 15(12):1053–1058.

Love, M. I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome Biology*, 15(12):1–21.

McFarland, J. M., Paolella, B. R., Warren, A., Geiger-Schuller, K., Shibue, T., Rothberg, M., Kuksenko, O., Colgan, W. N., Jones, A., Chambers, E., et al. (2020). Multiplexed single-cell transcriptional response profiling to define cancer vulnerabilities and therapeutic mechanism of action. *Nature Communications*, 11(1):1–15.

Miao, Z., Deng, K., Wang, X., and Zhang, X. (2018). Desingle for detecting three types of differential expression in single-cell rna-seq data. *Bioinformatics*, 34(18):3223–3224.

Nabavi, S., Schmolze, D., Maitituoheti, M., Malladi, S., and Beck, A. H. (2016). Emdomics: a robust and powerful method for the identification of genes differentially expressed between heterogeneous classes. *Bioinformatics*, 32(4):533–541.

O'Brien, P. C. (1992). Robust procedures for testing equality of covariance matrices. *Biometrics*, pages 819–827.

Qiu, X., Hill, A., Packer, J., Lin, D., Ma, Y.-A., and Trapnell, C. (2017). Single-cell mrna quantification and differential analysis with census. *Nature Methods*, 14(3):309–315.

Robinson, M. D., McCarthy, D. J., and Smyth, G. K. (2010). edger: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1):139–140.

Severson, K. A., Ghosh, S., and Ng, K. (2019). Unsupervised learning with contrastive latent variable models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4862–4869.

Srivastava, M. S. and Yanagihara, H. (2010). Testing the equality of several covariance matrices with fewer observations than the dimension. *Journal of Multivariate Analysis*, 101(6):1319–1329.

Stuart, J. M., Segal, E., Koller, D., and Kim, S. K. (2003). A gene-coexpression network for global discovery of conserved genetic modules. *Science*, 302(5643):249–255.

Townes, F. W., Hicks, S. C., Aryee, M. J., and Irizarry, R. A. (2019). Feature selection and dimension reduction for single-cell rna-seq based on a multinomial model. *Genome Biology*, 20(1):1–16.

Vassilev, L. T., Vu, B. T., Graves, B., Carvajal, D., Podlaski, F., Filipovic, Z., Kong, N., Kammlott, U., Lukacs, C., Klein, C., et al. (2004). In vivo activation of the p53 pathway by small-molecule antagonists of mdm2. *Science*, 303(5659):844–848.

Wang, T., Li, B., Nelson, C. E., and Nabavi, S. (2019). Comparative analysis of differential gene expression analysis tools for single-cell rna sequencing data. *BMC Bioinformatics*, 20(1):40.

Xia, Y., Cai, T., and Cai, T. T. (2015). Testing differential networks with applications to the detection of gene-gene interactions. *Biometrika*, 102(2):247–266.

Young, M. D., Mitchell, T. J., Braga, F. A. V., Tran, M. G., Stewart, B. J., Ferdinand, J. R., Collord, G., Botting, R. A., Popescu, D.-M., Loudon, K. W., et al. (2018). Single-cell transcriptomes from human kidneys reveal the cellular identity of renal tumors. *Science*, 361(6402):594–599.

Zou, H., Hastie, T., and Tibshirani, R. (2006). Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 15(2):265–286.

Zou, J. Y., Hsu, D. J., Parkes, D. C., and Adams, R. P. (2013). Contrastive learning using spectral methods. *Advances in Neural Information Processing Systems*, 26:2238–2246.

# 9 Supplementary material

## 9.1 Selecting variable genes

For all scRNA-seq experiments, we subsetted the data to the most variable 500 genes. We computed the closed-form Poisson deviance for each gene in each dataset using intercept only GLM-PCA as suggested by Townes et al. (2019), and took the genes with the highest deviance.

## 9.2 Perturb-seq data

The Perturb-seq data were downloaded from GEO: `https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE90063`.

## 9.3 Mix-seq data

The MIX-seq data were downloaded from Figshare: `https://figshare.com/articles/MIX-seq_data/10298696`.

## 9.4 Code

Code for the model and experiments is available at `https://github.com/andrewcharlesjones/cplvm`.

Supplementary Fig 1: **Global ELBO Bayes factors for the CGLVM.** Global hypothesis testing using the CGLVM with the same data as used in Figure 4a. Global tests were run on three datasets: data generated from a null model (left box), shuffled data approximating truly null data (middle box), and data generated from an alternative model (right box).



Supplementary Fig 2: **Gene set hypothesis tests are robust to gene set misspecification and gene set size.** (a) EBFs for 10 gene sets, each containing 25 genes. Set 1 was "perturbed", but only a fraction (12 of 25) of the genes had substantial differences in variation between conditions. Solid horizontal line shows the mean of the shuffled null, and dotted horizontal lines indicate the 95% confidence interval for the shuffled null. (b) EBFs for the "perturbed" gene set, but at varying gene set sizes.