

Isolating salient variations of interest in single-cell transcriptomic data with contrastiveVI

Ethan Weinberger^{1,*}, Chris Lin^{1,*}, and Su-In Lee^{1,✉}

¹Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle

* denotes equal contribution.

✉ corresponding: suinlee@cs.washington.edu

Abstract

Single-cell RNA sequencing (scRNA-seq) technologies enable a better understanding of previously unexplored biological diversity. Oftentimes, researchers are specifically interested in modeling the latent structures and variations enriched in one *target* scRNA-seq dataset as compared to another *background* dataset generated from sources of variation irrelevant to the task at hand. For example, we may wish to isolate factors of variation only present in measurements from patients with a given disease as opposed to those shared with data from healthy control subjects. Here we introduce Contrastive Variational Inference (contrastiveVI; <https://github.com/suinleelab/contrastiveVI>), a framework for end-to-end analysis of target scRNA-seq datasets that decomposes the variations into shared and target-specific factors of variation. On three target-background dataset pairs we demonstrate that contrastiveVI learns latent representations that recover known subgroups of target data points better than previous methods and finds differentially expressed genes that agree with known ground truths.

Main

Single-cell RNA sequencing (scRNA-seq) technologies have emerged as powerful tools for understanding previously unexplored biological diversity. Such technologies have enabled advances in our understanding of biological processes such as those underlying cancer [38],

27 Alzheimer’s disease [13, 28], and COVID-19 [36]. In many settings, scRNA-seq data ana-
28 lysts are specifically interested in patterns that are enriched in one dataset, referred to as the
29 *target*, as compared to a second related dataset, referred to as the *background*. Such target
30 and background dataset pairs arise naturally in many biological research contexts. For ex-
31 ample, data from healthy controls versus a diseased population or from pre-intervention and
32 post-intervention groups form intuitive background and target pairs. Moreover, with the de-
33 velopment of new technologies for measuring the effects of large numbers of perturbations in
34 parallel, such as Perturb-Seq [9] and MIX-Seq [29], tools for better understanding variations
35 unique to such perturbed cell lines compared to control populations will be critical.

36 Isolating salient variations present only in a target dataset is the subject of *contrastive*
37 *analysis* (CA) [40, 3, 17, 22, 32, 2]. While many recent studies have modeled scRNA-seq data
38 by fitting probabilistic models and representing the data in a lower dimension [23, 30, 16, 26,
39 24, 25], few were designed for CA. Such methods are thus unlikely to capture the enriched
40 variations in a target dataset, which are often subtle compared to the overall variations in
41 the data [3]. One recent study [17] designed a probabilistic model for analyzing scRNA-seq
42 data in the CA setting. However, this method assumes that a generalized linear model is
43 sufficiently expressive to model the variations in scRNA-seq data, even though previous work
44 has demonstrated substantial improvements by using more expressive nonlinear methods [23].

45 To address these limitations, we developed *contrastiveVI*, a deep generative model that
46 enables analysis of scRNA-seq data in the CA setting. *contrastiveVI* learns a probabilistic
47 representation of the data that accounts for the specific technical biases and noise characteris-
48 tics of scRNA-seq data as well as batch effects. Moreover, to handle CA tasks, *contrastiveVI*
49 models the variations underlying scRNA-seq data using two sets of latent variables: the
50 first, called the *background variables*, are shared across background and target cells while
51 the second, called the *salient variables*, are used to model variations specific to target data.
52 *contrastiveVI* can be used for a number of analysis tasks, including dimensionality reduction,
53 target dataset subgroup discovery, and differential gene expression testing. To highlight this
54 functionality, we applied *contrastiveVI* to three publicly available background and target
55 scRNA-seq dataset pairs, and demonstrated strong performance on all of them.

56 Results

57 The *contrastiveVI* Model

58 *contrastiveVI* is a probabilistic latent variable model that represents the uncertainty in ob-
59 served RNA counts as a combination of biological and technical factors. The input to

60 the contrastiveVI model consists of an RNA unique molecular identifier (UMI) count matrix
 61 along with labels denoting each cell as belonging to the background or target dataset (**Figure**
 62 **1a**). Additional categorical covariates such as anonymized donor ID or experimental batch
 63 are optional inputs to the model that can be used to integrate datasets.

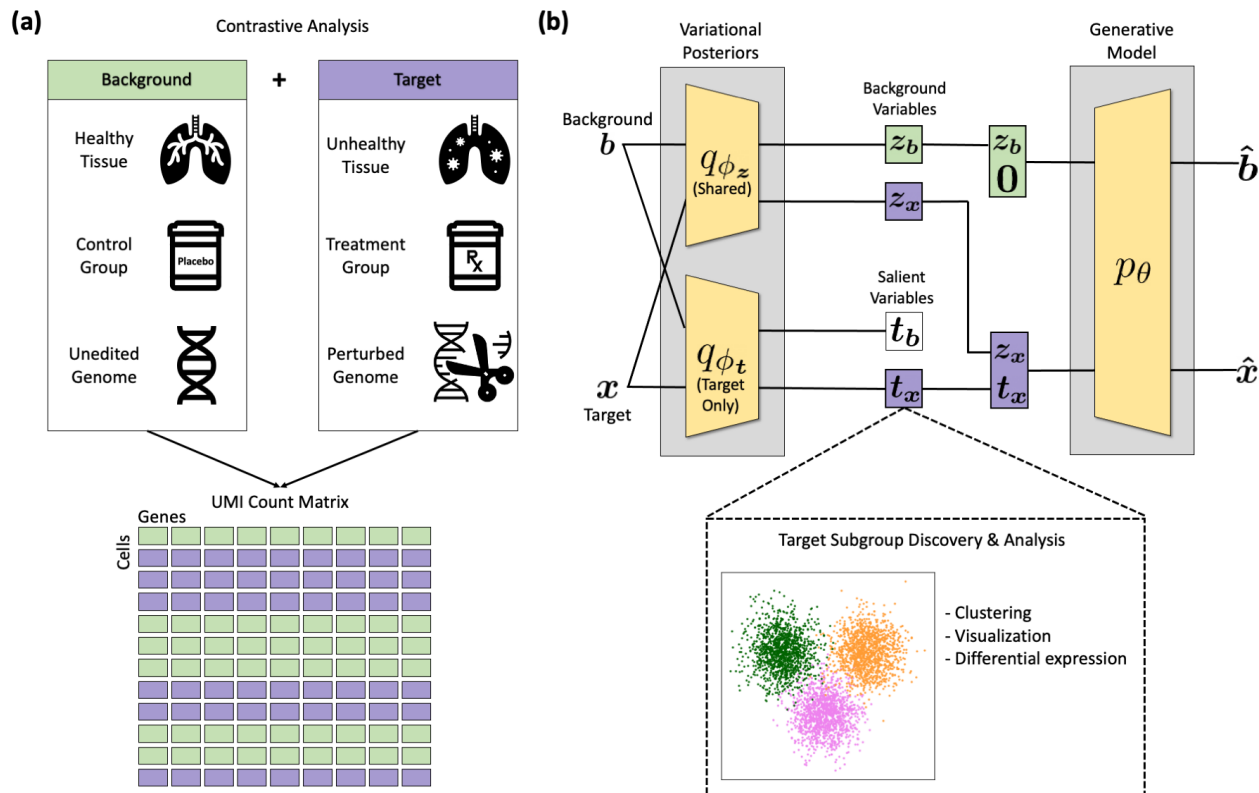


Figure 1: Overview of contrastiveVI. Given a reference background dataset and a second target dataset of interest, contrastiveVI separates the variations shared between the two datasets and the variations enriched in the target dataset. **a**, Example background and target data pairs. Samples from both conditions produce an RNA count matrix with each cell labeled as background or target. **b**, Schematic of the contrastiveVI model. A shared encoder network q_{ϕ_z} transforms a cell into the parameters of the posterior distribution for z , a low-dimensional set of latent factors shared across target and background data. For target data points only, a second encoder q_{ϕ_t} encodes target data points into the parameters of the posterior distribution for t , a second set of latent factors encoding variations enriched in the target dataset and not present in the background.

64 contrastiveVI encodes each cell as the parameters of a distribution in a low-dimensional
 65 latent space. This latent space is divided into two parts, each with its own encoding function.
 66 The first set of latent variables, called the background variables, capture factors of variation
 67 that are shared among background and target data. The second set of variables, denoted
 68 as the salient variables, capture variations unique to the target dataset. Only target data
 69 points are given salient latent variable values; background data points are instead assigned

70 a zero vector for these variables to represent their absence. As with scVI [23], contrastiveVI
71 also provides a way to estimate the parameters of the distributions underlying the observed
72 RNA measurements given a cell’s latent representation. Such distributions explicitly account
73 for technical factors in the observed data such as sequencing depth and batch effects. All
74 distributions are parameterized by neural networks.

75 The contrastiveVI model is based on the variational autoencoder (VAE) framework [21].
76 As such, its parameters can be learned using efficient stochastic optimization techniques,
77 easily scaling to large scRNA-seq datasets consisting of measurements from tens or hundreds
78 of thousands of cells. Following optimization, we can make use of the different components
79 of the contrastiveVI model for downstream analyses. For example, the salient latent repre-
80 sentations of target data can be used as inputs to clustering or visualization algorithms to
81 discover subgroups of target points. Moreover, the distributional parameters can be used for
82 additional tasks such as imputation or differential gene expression analysis. A more detailed
83 description of the contrastiveVI model can be found in **Methods**.

84 **contrastiveVI isolates subtle variations in target cells**

85 To evaluate the performance of contrastiveVI and other methods, we rely on datasets with
86 known biological variations in the target condition that are not present in the background
87 condition. One such dataset consists of expression data from bone marrow mononuclear
88 cells (BMMCs) from two patients with acute myeloid leukemia (AML) and two healthy
89 controls. The two patients underwent allogeneic stem-cell transplants, and BMMC samples
90 were collected before and after the transplant. It is known that gene expression profiles of
91 BMMCs differ pre- and post-transplant [39]. Therefore, the known biological variations in
92 this target dataset (AML patient BMMCs) correspond to pre- vs. post-transplant cellular
93 states. A performant model should learn a salient latent space separating pre- vs. post-
94 transplant status, while the latent space from a non-performant model does not make this
95 distinction.

96 Qualitatively, pre- and post-transplant cells are well separated in the salient latent space
97 learned by contrastiveVI (**Figure 2a**). We also quantified how well contrastiveVI’s salient
98 latent space separates the two groups of target cells using three metrics—the average silhou-
99 ette width, adjusted Rand Index (ARI), and adjusted mutual information (AMI; **Methods**).
100 We find that contrastiveVI performs well on all of these metrics (**Figure 2b**), indicating that
101 it successfully recovers the variations enriched in the target dataset. Furthermore, we exper-
102 imented with a workflow for using contrastiveVI for end-to-end biological discovery. After
103 embedding the AML patient samples into the contrastiveVI salient latent space, we used

104 k-means clustering to divide the samples into two groups. Highly differentially expressed
105 genes across the two clusters were then obtained by Monte Carlo sampling of denoised,
106 library size-normalized expressions from the contrastiveVI decoder (**Methods**). Finally,
107 pathway enrichment analysis (**Methods**) was performed with these differentially expressed
108 genes using the Kyoto Encyclopedia of Genes and Genomes (KEGG) 2016 pathway database
109 [18]. Based on our quantitative results, our two clusters exhibited strong agreement with
110 the two ground-truth groups (ARI: 0.77 ± 0.01). Moreover, the pathways enriched by the
111 differentially expressed genes between the two clusters are related to immune response and
112 graft rejection (**Figure 2c**). We provide a full list of enriched pathways in **Supplementary**
113 **Table 1**. These results align with known cellular state transitions of BMDCs before and
114 after a transplant.

115 **contrastiveVI outperforms other modeling approaches**

116 To illustrate the advantages of contrastiveVI, we benchmarked its performance against that
117 of three previously proposed methods for analyzing raw scRNA-seq count data. First, to
118 demonstrate that our contrastive approach is necessary for capturing enriched variations in
119 target datasets, we compared against scVI [23]. scVI has achieved state-of-the-art results
120 on many tasks; however, it was not specifically designed for the CA setting and thus may
121 struggle to isolate salient variations of interest. We also compared against two contrastive
122 methods designed for analyzing scRNA-seq count data: contrastive Poisson latent variable
123 model (CPLVM) and contrastive generalized latent variable model (CGLVM) [17]. While
124 these methods are designed for the contrastive setting, they both make the strong assumption
125 that linear models can accurately capture the complex variations in scRNA-seq data. To our
126 knowledge, the CPLVM and CGLVM methods are the only existing contrastive methods for
127 analyzing scRNA-seq count data.

128 Qualitatively (**Figure 2a**), we find that none of these baseline models are able to sep-
129 arate pre- and post-transplant cells as well as contrastiveVI can. This finding is further
130 confirmed by quantitative results (**Figure 2b**). Across all of our metrics we find that con-
131 trastiveVI significantly outperforms baseline models, with especially large gains in the ARI
132 and AMI. These results indicate that contrastiveVI recovered the variations enriched in the
133 AML patient data far better than baseline models.

134 **contrastiveVI separates intestinal epithelial cells by infection type**

135 We next applied contrastiveVI to data collected in Haber et al. [15]. This data consists of
136 gene expression measurements of intestinal epithelial cells from mice infected with either

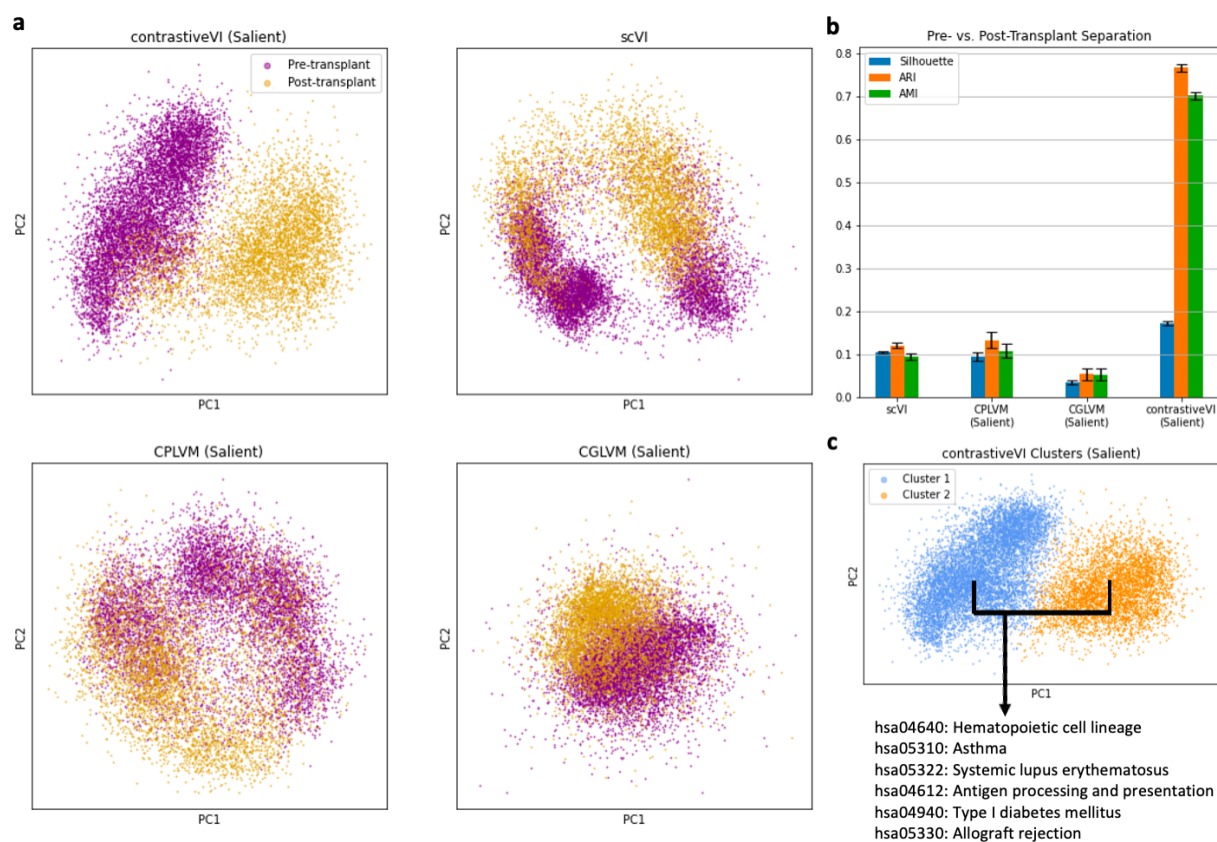


Figure 2: contrastiveVI successfully captures enriched variations in scRNA-seq data. **a**, Principal component (PC) plots of contrastiveVI and baseline models' latent representations. For scVI, the first two PCs of the model's single latent representations are plotted, while for contrastive methods the PCs from their salient latent representations are plotted. **b**, Quantitative measures of separation between pre- and post-transplant cells. Silhouette is the average silhouette width of pre-annotated subpopulations, ARI is the adjusted Rand index, and AMI is the adjusted mutual information. Higher values indicate better performance for all metrics. For each method, the mean and standard error across five random trials are plotted. **c**, contrastiveVI's salient latent representations of the target dataset were clustered into two groups. Pathway enrichment analysis was then performed on the differentially expressed genes between the two clusters.

137 *Salmonella* or *Heligmosomoides polygyrus* (*H. poly*). As a background dataset we used
 138 measurements collected from healthy cells released by the same authors. Here our goal
 139 is to separate cells by infection type in the salient latent space. On the other hand, any
 140 separations in the background latent space should reflect variations shared between healthy
 141 and infected cells, such as those due to cell type differences. We present our results in Figure
 142 3.

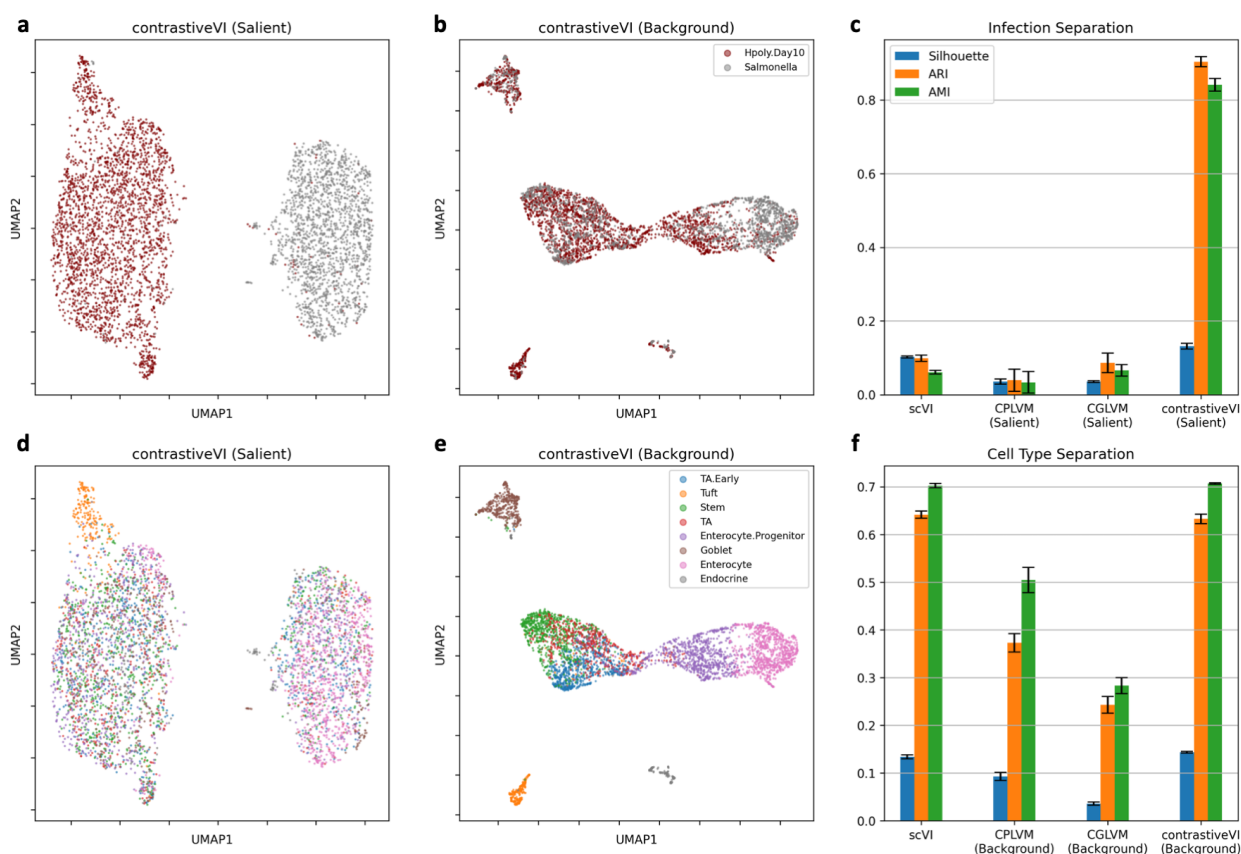


Figure 3: contrastiveVI isolates responses to different infections in intestinal epithelial cells. **a,b**, UMAP plots of contrastiveVI's salient and background representations colored by infection type. Cells are correctly separated by infection type in the salient space, while they mix across infection types in the background space. **c**, Clustering metrics quantify how well cells separate by infection type for scVI's single latent space and contrastive models' salient latent spaces, with means and standard errors across five random trials plotted. **d,e**, UMAP plots of contrastiveVI's salient and background representations colored by cell type. Cells separate well by cell type in the background space, while they mix across cell types in the salient space. **f**, Quantifying how well cells separate by cell type in scVI's single latent space and contrastive models' background latent spaces, with means and standard errors across five random trials for each method.

143 We find that contrastiveVI successfully separates the cells by infection type in its salient

144 latent space (**Figure 3a**). Moreover we find that cells mix across infection types in our
145 background latent space as expected (**Figure 3b**). These results indicate that enriched
146 variations due to infection response are correctly being relegated to the salient latent space.
147 Once again we find that previously proposed methods fail to stratify the two classes of
148 target samples in their salient latent spaces as demonstrated by a set of quantitative metrics
149 (**Figure 3c**). For this dataset we were able to further validate contrastiveVI's separation
150 of target and background variations using ground truth cell type labels provided by the
151 authors (**Supplementary Table 2**). In particular, we found strong mixing across cell types
152 in contrastiveVI's salient latent space (**Figure 3d**), while cell types separated clearly in the
153 background latent space (**Figure 3e**). Our quantitative metrics indicate that contrastiveVI's
154 background latent space is competitive with if not outright superior to other methods' at
155 capturing variations between cell types (**Figure 3f**). Taken together, these results further
156 indicate that contrastiveVI successfully disentangles variations enriched in target data from
157 those shared across the target and background, even when other methods struggle.

158 **contrastiveVI stratifies cells by response to molecular perturbations**

159 In addition to studying transplant outcome and infection response, contrastiveVI can be
160 applied to examine drug treatment response. We demonstrate this capability using cancer
161 cell lines treated with vehicle control dimethyl sulfoxide (DMSO) or idasanutlin collected by
162 McFarland et al. [29]. The small molecule idasanutlin is an antagonist of MDM2, a negative
163 regulator of the tumor suppressor protein p53, hence offering cancer therapeutic opportunity
164 [35]. In the CA context, DMSO-treated samples are considered the background dataset,
165 and idasanutlin-treated samples the target dataset. Based on the mechanism of action of
166 idasanutlin, activation of the p53 pathway is observed in cell lines with wildtype *TP53* (gene
167 of p53) and not in transcriptionally inactive mutant *TP53* cell lines [35]. Therefore, unique
168 variations in the target dataset should be related to *TP53* mutation status. This stratifica-
169 tion of cell response based on *TP53* mutation is readily identified by the salient latent space
170 of all methods (**Figure 4a** and **Figure 4b**). Notably, contrastiveVI outperforms other meth-
171 ods based on ARI and AMI, providing better separated clusters for downstream analyses.
172 (**Figure 4b**). Particularly, the two clusters identified using the contrastiveVI salient latent
173 space have differentially expressed genes enriched for the p53 signaling pathway (**Figure**
174 **4c**). It is worth noting that the p53 signaling pathway is the only statistically significant
175 (under 0.05 false discovery rate) pathway identified by contrastiveVI. All these results show
176 that contrastiveVI finds salient variations in the target samples treated with idasanutlin that
177 specifically relate to the biological ground truth effect of idasanutlin perturbation.

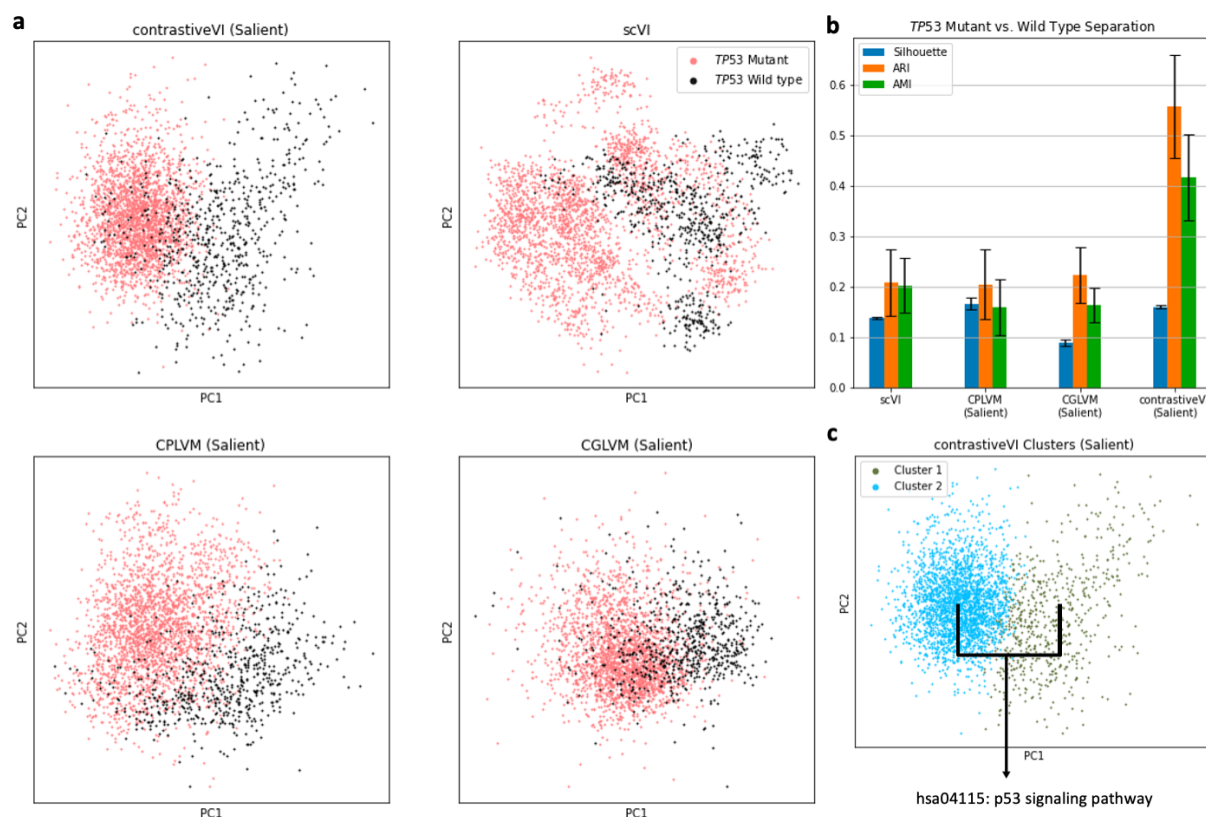


Figure 4: contrastiveVI stratifies cancer cell lines by response to idasanutlin. **a**, PC plots of target data latent representations from contrastiveVI and baseline models. The first two PCs of scVI's single latent space are plotted. For contrastive methods, the first two PCs of their salient latent space are plotted. **b**, The average silhouette (silhouette), adjusted Rand Index (ARI) and adjusted mutual information (AMI), with mean and standard error across five random trials plotted for each method. **c**, Two clusters identified by k-means clustering with contrastiveVI's salient latent representations of the target dataset. Highly differentially expressed genes were identified from the two clusters, and these genes were used to perform pathway enrichment analysis.

178 Discussion

179 In this work we introduced contrastiveVI, a scalable probabilistic framework for isolating
180 enriched variations in a target scRNA-seq dataset as compared to a related background dataset.
181 contrastiveVI is the first method designed to analyze scRNA-seq data in the contrastive anal-
182 ysis setting that both explicitly models the technical factors of variation in scRNA-seq data
183 and takes advantage of the expressive power of deep generative modeling techniques. More-
184 over, contrastiveVI includes a number of other capabilities relevant to scRNA-seq analysis
185 out of the box, such as batch effect correction and differential expression testing.

186 In three different contexts—response to cancer treatment, infection by different pathogens,
187 and exposure to small-molecule drug perturbations—we demonstrated that contrastiveVI iso-
188 lated enriched variations in target cells while other methods struggled. With the recent de-
189 velopment of new sequencing technologies for efficiently measuring transcriptomic responses
190 to various perturbations, such as Perturb-Seq and MIX-Seq, we expect contrastiveVI to be
191 of immediate interest to the scRNA-seq research community. Moreover, contrastiveVI was
192 implemented using the scvi-tools [11] Python library, thereby enabling interoperability with
193 Scanpy [37] and Seurat [33] analysis pipelines.

194 The ideas behind contrastiveVI admit multiple potential directions for future work. Sim-
195 ilar contrastive disentanglement techniques could be used to extend models that make use
196 of multimodal data, such as totalVI [12], to better understand how variations enriched in
197 target datasets are expressed across different modalities of single-cell data. Moreover, recent
198 work [10, 14, 31, 27, 34] in learning biologically meaningful representations of gene expression
199 data could be incorporated to better understand the different sources of variation learned
200 by the model. For example, using a constrained architecture such that latent variables cor-
201 respond to gene pathways could shed more light on the biological phenomena captured in
202 the different latent spaces.

203 Methods

204 The contrastiveVI model

205 Here we present the contrastiveVI model in more detail. We begin by describing the model’s
206 generative process and then the model’s inference procedure.

207 The contrastiveVI generative process

208 For a target data point x_n we assume that each expression value x_{ng} for sample n and gene
 209 g is generated through the following process:

$$\begin{aligned}
 210 \quad & z_n \sim \text{Normal}(0, I) \\
 211 \quad & t_n \sim \text{Normal}(0, I) \\
 212 \quad & \ell_n \sim \text{log normal}(\ell_\mu, \ell_\sigma^2) \\
 213 \quad & \rho_n = f_w(z_n, t_n, s_n) \\
 214 \quad & w_{ng} \sim \text{Gamma}(\rho_{ng}, \theta_g) \\
 215 \quad & y_{ng} \sim \text{Poisson}(\ell_n w_{ng}) \\
 216 \quad & h_{ng} \sim \text{Bernoulli}(f_h^g(z_n, t_n, s_n)) \\
 217 \quad & x_{ng} = \begin{cases} y_{ng} & \text{if } h_{ng} = 1 \\ 0 & \text{otherwise} \end{cases}
 \end{aligned}$$

218 In this process z_n and t_n both refer to sets of latent variables underlying variations in
 219 scRNA-seq expression data. Here z_n represents variables that are shared across background
 220 and target cells, while t_n represents variations unique to the target cells. We place a stan-
 221 dard multivariate Gaussian prior on both sets of latent factors, as such a specification is
 222 computationally convenient for inference in the VAE framework [21]. To encourage the
 223 disentanglement of latent factors, for background data points b_n we assume the same gener-
 224 ative process but instead set $t_n = \mathbf{0}$ to represent the absence of salient latent factors in the
 225 generative process. Categorical covariates such as experimental batches are represented by
 226 s_n .

227 ℓ_μ and $\ell_\sigma \in \mathbb{R}_+^B$, where B denotes the cardinality of the categorical covariate, parameterize
 228 the prior for latent RNA library size scaling factor on a log scale. For each category
 229 (e.g. experimental batch), ℓ_μ and ℓ_σ^2 are set to the empirical mean and variance of the
 230 log library size. The gamma distribution is parameterized by the mean $\rho_{ng} \in \mathbb{R}_+$ and
 231 shape $\theta_g \in \mathbb{R}_+$. Furthermore, following the generative process, θ_g is equivalent to a gene-
 232 specific inverse dispersion parameter for a negative binomial distribution, and $\theta \in \mathbb{R}_+^G$ is
 233 estimated via variational Bayesian inference. f_w and f_g in the generative process are neural
 234 networks that transform the latent space and batch annotations to the original gene space,
 235 i.e.: $\mathbb{R}^d \times \{0, 1\}^B \rightarrow \mathbb{R}^G$, where d is the latent dimension. The network f_w is constrained
 236 during inference to encode the mean proportion of transcripts expressed across all genes by

237 using a softmax activation function in the last layer. That is, letting $f_w^g(z_n, t_n, s_n)$ denote
238 the entry in the output of f_w corresponding to gene g , we have $\sum_g f_w^g(z_n, t_n, s_n) = 1$. The
239 neural network f_h encodes whether a particular gene’s expression has dropped out in a cell
240 due to technical factors.

241 Our generative process closely follows that of scVI [23], with the addition of the salient
242 latent factors t_n . While scVI’s modeling approach has been shown to excel at many scRNA-
243 seq analysis tasks, our empirical results demonstrate that it is not suited for contrastive
244 analysis (CA). By dividing the RNA latent factors into shared factors z_n and target-specific
245 factors t_n , contrastiveVI successfully isolates variations enriched in target datasets missed by
246 previous methods. We depict the full contrastiveVI generative process as a graphical model
247 in **Supplementary Figure 1**.

248 Inference with contrastiveVI

249 We cannot compute the contrastiveVI posterior distribution using Bayes’ rule as the integrals
250 required to compute the model evidence $p(x_n|s_n)$ are analytically intractable. As such, we
251 instead approximate our posterior distribution using variational inference [5]. For target
252 data points we approximate our posterior with a distribution factorized as follows:

$$253 \quad q_{\phi_x}(z_n, t_n, \ell_n|x_n, s_n) = q_{\phi_z}(z_n|x_n, s_n)q_{\phi_t}(t_n|x_n, s_n)q_{\phi_\ell}(\ell_n|x_n, s_n). \quad (1)$$

254 Here ϕ_x denotes a set of learned weights used to infer the parameters of our approximate
255 posterior. Based on our factorization, we can divide ϕ_x into three disjoint sets ϕ_z , ϕ_t and ϕ_ℓ
256 for inferring the parameters of the distributions of z , t and ℓ respectively. Following the VAE
257 framework [21], we then approximate the posterior for each factor as a deep neural network
258 that takes in expression levels as input and outputs the parameters of its corresponding
259 approximate posterior distribution (e.g. mean and variance). Moreover, we note that each
260 factor in the posterior approximation shares the same family as its respective prior distri-
261 bution (e.g. $q(z_n|x_n, s_n)$ follows a normal distribution). We can simplify our likelihood by
262 integrating out w_{ng} , h_{ng} , and y_{ng} , yielding $p_\nu(x_{ng}|z_n, t_n, s_n, \ell_n)$, which follows a zero-inflated
263 negative binomial (ZINB) distribution (**Supplementary Note 1**) and where ν denotes the
264 parameters of our generative model. As with our approximate posteriors, we realize our
265 generative model as a deep neural network. For Equation 1 we can derive (**Supplementary**
266 **Note 2**) a corresponding variational lower bound:

$$p(x|s) \geq \mathbb{E}_{q(z,t,\ell|x,s)} \log p(x|z, t, \ell, s) - D_{KL}(q(z|x, s)||p(z)) - D_{KL}(q(t|x, s)||p(t)) - D_{KL}(q(\ell|x, s)||p(\ell|s)). \quad (2)$$

Next, for background data points we approximate the posterior using the factorization:

$$q_{\phi_b}(z_n, \ell_n|b_n, s_n) = q_{\phi_z}(z_n|b_n, s_n)q_{\phi_\ell}(\ell_n|b_n, s_n), \quad (3)$$

where ϕ_b denotes a set of learned parameters use to infer the values of z_n and ℓ_n for background samples. Following our factorization, we divide ϕ_b into the disjoint sets ϕ_z and ϕ_ℓ . We note that ϕ_z and ϕ_ℓ are shared across target and background samples; this encourages the posterior distributions q_{ϕ_z} and q_{ϕ_ℓ} to capture variations shared across the datasets, while q_{ϕ_t} captures variations unique to the target data. Once again we can simplify our likelihood by integrating out w_{ng} , h_{ng} , and y_{ng} to obtain $p_\nu(x_{ng}|z_n, \mathbf{0}, s_n, \ell_n)$, which follows a ZINB distribution. We similarly note that the parameters of our generative model ν are shared across target and background points to encourag z to capture shared variations across target and background points while t captures target-specific variations. We then have the following variational lower bound for our background data points:

$$p(b|s) \geq \mathbb{E}_{q(z,\ell|x,s)} \log p(b|z, \ell, s) - D_{KL}(q(z|b, s)||p(z)) - D_{KL}(q(\ell|b, s)||p(\ell|s)). \quad (4)$$

We then jointly optimize the parameters of our generative model and inference networks using stochastic gradient descent to maximize the sum of these two bounds over our background and target data points. All neural networks used to implement the variational and generative distributions were feedforward and used standard activation functions. We used the same network architecture and hyperparameter values for all experiments, and we refer the reader to **Supplementary Note 3** for more details.

Differential gene expression analysis with contrastiveVI

For two cell groups $A = (a_1, a_2, \dots, a_n)$ and $B = (b_1, b_2, \dots, b_m)$ in the target dataset, the posterior probability of gene g being differentially expressed in the two groups can be obtained as proposed by Boyeau et al. [6]. For any arbitrary cell pair a_i, b_j , we have two mutually exclusive models

$$\mathcal{M}_1^g : |r_{a_i, b_j}^g| > \delta \text{ and } \mathcal{M}_0^g : |r_{a_i, b_j}^g| \leq \delta$$

293 where $r_{a_i, b_j}^g := \log_2(\rho_{a_i}^g) - \log_2(\rho_{b_j}^g)$ is the log fold change of the denoised, library size-
294 normalized expression of gene g , and δ is a pre-defined threshold for log fold change mag-
295 nitude to be considered biologically meaningful. The posterior probability of differential
296 expression is therefore expressed as $p(\mathcal{M}_1^g | x_{a_i}, x_{b_j})$, which can be obtained via marginaliza-
297 tion of the latent variables and categorical covariates:

$$298 \quad p(\mathcal{M}_1^g | x_{a_i}, x_{b_j}) = \sum_s \int_{z_{a_i}, t_{a_i}, z_{b_j}, t_{b_j}} p(\mathcal{M}_1^g | z_{a_i}, t_{a_i}, z_{b_j}, t_{b_j}) p(s) dp(z_{a_i}, t_{a_i} | x_{a_i}, s) dp(z_{b_j}, t_{b_j} | x_{b_j}, s),$$

299 where $p(s)$ is the relative abundance of target cells in category s , and the integral can be
300 computed via Monte Carlo sampling using the variational posteriors q_{ϕ_z}, q_{ϕ_t} . Finally, the
301 group-level posterior probability of differential expression is

$$302 \quad \int_{a, b} p(\mathcal{M}_1^g | x_a, x_b) dp(a) dp(b),$$

303 where, assuming that the cells are independent, $a \sim \mathcal{U}(a_1, \dots, a_m)$ and $b \sim \mathcal{U}(b_1, \dots, b_m)$.
304 Computationally, this quantity can be estimated by a large random samples of pairs from
305 the cell group A and B . In our experiments, 10,000 cell pairs were sampled, 100 Monte Carlo
306 samples were obtained from the variational posteriors for each cell, and the δ threshold was
307 set to 0.25, which is the default value recommended by the scvi-tools Python library [11].
308 Genes with group-level posterior probability of differential expression greater than 0.95 were
309 considered for downstream pathway enrichment analysis.

310 Pathway enrichment analysis

311 Pathway enrichment analysis refers to a computational procedure for determining whether
312 a predefined set of genes (i.e., a gene pathway) have statistically significant differences in
313 expression between two biological states. Many tools exist for performing pathway enrich-
314 ment analysis (see [19] for a review). In our analyses we use Enrichr [8], a pathway analysis
315 tool for non-ranked gene lists based on Fisher's exact test, to find enriched pathways from
316 the KEGG 2016 pathways database [18]. Specifically, the Enrichr wrapper implemented in
317 the open-source GSEAPy¹ Python library was used for our analyses. Pathways enriched at
318 false discovery rate smaller than 0.05 (adjusted by the Benjamini-Hochberg procedure [4])
319 are reported in this study.

¹<https://gseapy.readthedocs.io/en/latest/>

320 Baseline models

321 Because the choice of library size normalization method tends to drastically impact dimen-
322 sion reduction and subsequent clustering results of methods not designed for modeling library
323 sizes [30], we consider CA methods specifically tailored for scRNA-seq count data as base-
324 lines in this study. To our knowledge, CPLVM (contrastive Poisson latent variable model)
325 and CGLVM (contrastive generalized latent variable model) are the only CA methods that
326 explicitly model count-based scRNA-seq normalization [17]. We present a summary of pre-
327 vious work in CA in **Supplementary Table 4**. We also consider scVI, a deep generative
328 model for UMI count data that takes batch effect, technical dropout, and varying library
329 size into modeling considerations [23], to illustrate the need for models specifically designed
330 for CA. Below we describe the CA methods CPLVM and CGLVM in more detail.

331 In CPLVM, variations shared between background and target condition are assumed to
332 be captured by two sets of latent variables $\{z_i^b\}_{i=1}^n$ and $\{z_j^t\}_{j=1}^m$, and target condition-specific
333 variations are described by latent variables $\{t_j\}_{j=1}^m$, where n, m are the number of background
334 and target cells, respectively. Library size differences between the two conditions are modeled
335 by $\{\alpha_i^b\}_{i=1}^n$ and $\{\alpha_j^t\}_{j=1}^m$, whereas gene-specific library sizes are parameterized by $\delta \in \mathbb{R}_+^G$,
336 where G is the number of genes. Each data point is considered Poisson distributed, with rate
337 parameter determined by $\alpha_i^b \delta \odot (S^\top z_i^b)$ for a background cell i and by $\alpha_j^t \delta \odot (S^\top z_j^t + W^\top t_j)$ for
338 a target cell j , where S, W are model weights that linearly combine the latent variables, and
339 \odot represents an element-wise product. The model weights and latent variables are assumed
340 to have Gamma priors, δ has a standard log-normal prior, and α_i^b, α_j^t have log-normal priors
341 with parameters given by the empirical mean and variance of log total counts in each dataset.
342 Posterior distributions are fitted using variational inference with mean-field approximation
343 and log-normal variational distributions.

344 The CA modeling approaches of CGLVM and CPLVM are similar. In CGLVM, however,
345 the relationships of latent factors are considered additive and relate to the Poisson rate
346 parameter via an exponential link function (similar to a generalized linear modeling scheme).
347 All the priors and variational distributions are Gaussian in CGLVM.

348 Model optimization details

349 For all datasets, contrastiveVI models were trained with 80% of the background and target
350 data, and with 20% of the data reserved as a validation set for early stopping to determine
351 the number of training epochs needed. Training was early stopped when the validation
352 variational lower bound showed no improvement for 45 epochs, typically resulting in 127 to
353 500 epochs of training. All contrastiveVI models were trained with the Adam optimizer [20]

354 with $\varepsilon = 0.01$, learning rate at 0.001, and weight decay at 10^{-6} . The same hyperparameters
355 and training scheme were used to optimize the scVI models using only target data, usually
356 with 274 to 500 epochs of training based on the early stopping criterion. As in Jones et al.,
357 the CPLVMs were trained via variational inference using all background and target data for
358 2,000 epochs with the Adam optimizer with $\varepsilon = 10^{-8}$ and learning rate at 0.05, and the
359 CGLVMs were similarly trained for 1,000 epochs and learning rate at 0.01 [17]. All models
360 were trained with 10 salient and 10 background latent variables for five times with different
361 random weight initializations. We also trained models with varying salient latent dimension
362 sizes and obtained overall consistent results (**Supplementary Figure 2**).

363 **Datasets and preprocessing**

364 Here we briefly describe all datasets used in this work along with any corresponding pre-
365 processing steps. All preprocessing steps were performed using the Scanpy Python package
366 [37]. All our code for downloading and preprocessing these datasets is publicly available
367 at <https://github.com/suinleelab/contrastiveVI>. For all experiments we retained the
368 top 2,000 most highly variable genes returned from the Scanpy `highly_variable_genes`
369 function with the `flavor` parameter set to `seurat_v3`. For all datasets, the number of cells
370 in the background vs. target can be found in **Supplementary Table 3**.

371 **Zheng et al., 2017**

372 This dataset consists of single-cell RNA expression levels of a mixture of bone marrow
373 mononuclear cells (BMMCs) from 10x Genomics [1]. For our target dataset, we use samples
374 taken from patients with acute myeloid leukemia (AML) before and after a stem cell trans-
375 plant. For our background dataset, we use measurements taken from two healthy control
376 patients released as part of the same study. All data is publicly available: files containing
377 measurements from first patient pre- and post-transplant can be found here and here, re-
378 spectively; from the second patient pre- and post-transplant here and here, respectively; and
379 from the two healthy control patients here and here.

380 **Haber et al., 2017**

381 This dataset (Gene Expression Omnibus accession number GSE92332) used scRNA-seq
382 measurements to investigate the responses of intestinal epithelial cells in mice to differ-
383 ent pathogens. In particular, in this dataset responses to *Salmonella* and the parasite *H.*
384 *polygyrus* were investigated. Here our target dataset consisted measurements of cells in-
385 fected with *Salmonella* and from cells 10 days after being infected with *H. polygyrus*, while

386 our background consisted of measurements from healthy control cells released as part of the
387 same study.

388 **McFarland et al., 2020**

389 This dataset measured cancer cell lines' transcriptional responses after being treated with
390 various small-molecule therapies. For our target dataset, we used data from cells that were
391 exposed to idasanutlin, and for our background we used data from cells that were exposed to
392 a control solution of dimethyl sulfoxide (DMSO). *TP53* mutation status was determined by
393 cross-referencing with a list of cell lines with mutations provided by the authors in the code
394 repository accompanying the paper. The data was downloaded from the authors' Figshare
395 repository.

396 **Evaluation metrics**

397 Here we describe the quantitative metrics used in this study. All metrics were computed
398 using their corresponding implementations in the scikit-learn Python package [7].

399 **Silhouette width**

400 We calculate silhouette width using the latent representations returned by each method. For
401 a given sample i , the silhouette width $s(i)$ is defined as follows. Let $a(i)$ be the average
402 distance between i and the other samples with the same ground truth label, and let $b(i)$ be
403 the smallest average distance between i and all other samples with a different label. The
404 silhouette score $s(i)$ is then

$$405 \quad s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}.$$

406 A silhouette width close to one indicates that i is tightly clustered with cells with the
407 same ground truth label, while a score close to -1 indicates that a cell has been grouped with
408 cells with a different label.

409 **Adjusted Rand index**

410 The adjusted Rand index (ARI) measures agreement between reference clustering labels and
411 labels assigned by a clustering algorithm. Given a set of n samples and two sets of clustering
412 labels describing those cells, the overlap between clustering labels can be described using a
413 contingency table, where each entry indicates the number of cells in common between the
414 two sets of labels. Mathematically, the ARI is calculated as

$$\text{ARI} = \frac{\sum_{ij} \binom{n_{ij}}{2} - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}},$$

where n_{ij} is the number of cells assigned to cluster i based on the reference labels and cluster j based on a clustering algorithm, a_i is the number of cells assigned to cluster i in the reference set, and b_j is the number of cells assigned to cluster j by the clustering algorithm. ARI values close to 1 indicate agreement between the reference labels and labels assigned by a clustering algorithm.

Adjusted mutual information

The adjusted mutual information (AMI) is a corrected-for-chance version of the normalized mutual information, and it is another measure of the agreement between reference clustering labels and labels assigned by a clustering algorithm. For two clusterings U and V , we have

$$\text{AMI}(U, V) = \frac{I(U; V) - \mathbb{E}[I(U; V)]}{(H(U) + H(V))/2 - \mathbb{E}[I(U; V)]}$$

where I represents mutual information, and H represents entropy. AMI values closer to 1 indicate greater agreement between U and V .

References

- [1] 10x Genomics. 10x genomics. support: single cell gene expression datasets. <https://support.10xgenomics.com/single-cell-gene-expression/datasets>, 2021.
- [2] A. Abid and J. Zou. Contrastive variational autoencoder enhances salient features. *arXiv preprint arXiv:1902.04601*, 2019.
- [3] A. Abid, M. J. Zhang, V. K. Bagaria, and J. Zou. Exploring patterns enriched in a dataset with contrastive principal component analysis. *Nature Communications*, 9(1): 1–7, 2018.
- [4] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.
- [5] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.

- 441 [6] P. Boyeau, R. Lopez, J. Regier, A. Gayoso, M. I. Jordan, and N. Yosef. Deep gener-
442 ative models for detecting differential expression in single cells. *Machine Learning in*
443 *Computational Biology (MLCB)*, October 2019.
- 444 [7] L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae,
445 P. Prettenhofer, A. Gramfort, J. Grobler, et al. Api design for machine learning software:
446 experiences from the scikit-learn project. *arXiv preprint arXiv:1309.0238*, 2013.
- 447 [8] E. Y. Chen, C. M. Tan, Y. Kou, Q. Duan, Z. Wang, G. V. Meirelles, N. R. Clark, and
448 A. Ma’ayan. Enrichr: interactive and collaborative html5 gene list enrichment analysis
449 tool. *BMC Bioinformatics*, 14(1):1–14, 2013.
- 450 [9] A. Dixit, O. Parnas, B. Li, J. Chen, C. P. Fulco, L. Jerby-Arnon, N. D. Marjanovic,
451 D. Dionne, T. Burks, R. Raychowdhury, et al. Perturb-seq: dissecting molecular circuits
452 with scalable single-cell rna profiling of pooled genetic screens. *Cell*, 167(7):1853–1866,
453 2016.
- 454 [10] N. Fortelny and C. Bock. Knowledge-primed neural networks enable biologically in-
455 terpretable deep learning on single-cell sequencing data. *Genome biology*, 21(1):1–36,
456 2020.
- 457 [11] A. Gayoso, R. Lopez, G. Xing, P. Boyeau, K. Wu, M. Jayasuriya, E. Melhman,
458 M. Langevin, Y. Liu, J. Samaran, G. Misrachi, A. Nazaret, O. Clivio, C. Xu, T. Ashuach,
459 M. agha Lotfollahi, V. Svensson, E. da Veiga Beltrame, C. Talavera-López, L. Pachter,
460 F. J. Theis, A. M. Streets, M. I. Jordan, J. Regier, and N. Yosef. scvi-tools: a library
461 for deep probabilistic analysis of single-cell omics data. *bioRxiv*, 2021.
- 462 [12] A. Gayoso, Z. Steier, R. Lopez, J. Regier, K. L. Nazor, A. Streets, and N. Yosef. Joint
463 probabilistic modeling of single-cell multi-omic data with totalvi. *Nature Methods*, 18
464 (3):272–282, 2021.
- 465 [13] A. Grubman, G. Chew, J. F. Ouyang, G. Sun, X. Y. Choo, C. McLean, R. K. Simmons,
466 S. Buckberry, D. B. Vargas-Landin, D. Poppe, et al. A single-cell atlas of entorhinal
467 cortex from individuals with alzheimer’s disease reveals cell-type-specific gene expression
468 regulation. *Nature neuroscience*, 22(12):2087–2097, 2019.
- 469 [14] G. Gut, S. G. Stark, G. Rätsch, and N. R. Davidson. Pmvae: Learning interpretable
470 single-cell representations with pathway modules. *bioRxiv*, 2021.

- 471 [15] A. L. Haber, M. Biton, N. Rogel, R. H. Herbst, K. Shekhar, C. Smillie, G. Burgin,
472 T. M. Delorey, M. R. Howitt, Y. Katz, et al. A single-cell survey of the small intestinal
473 epithelium. *Nature*, 551(7680):333–339, 2017.
- 474 [16] Y. Hao, S. Hao, E. Andersen-Nissen, W. M. Mauck III, S. Zheng, A. Butler, M. J. Lee,
475 A. J. Wilk, C. Darby, M. Zager, et al. Integrated analysis of multimodal single-cell data.
476 *Cell*, 2021.
- 477 [17] A. Jones, F. W. Townes, D. Li, and B. E. Engelhardt. Contrastive latent variable
478 modeling with application to case-control sequencing experiments. *arXiv preprint*
479 *arXiv:2102.06731*, 2021.
- 480 [18] M. Kanehisa and S. Goto. Kegg: kyoto encyclopedia of genes and genomes. *Nucleic*
481 *Acids Research*, 28(1):27–30, 2000.
- 482 [19] P. Khatrı, M. Sirota, and A. J. Butte. Ten years of pathway analysis: current approaches
483 and outstanding challenges. *PLoS Computational Biology*, 8(2):e1002375, 2012.
- 484 [20] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint*
485 *arXiv:1412.6980*, 2014.
- 486 [21] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint*
487 *arXiv:1312.6114*, 2013.
- 488 [22] D. Li, A. Jones, and B. Engelhardt. Probabilistic contrastive principal component
489 analysis. *arXiv preprint arXiv:2012.07977*, 2020.
- 490 [23] R. Lopez, J. Regier, M. B. Cole, M. I. Jordan, and N. Yosef. Deep generative modeling
491 for single-cell transcriptomics. *Nature Methods*, 15(12):1053–1058, 2018.
- 492 [24] M. Lotfollahi, F. A. Wolf, and F. J. Theis. scgen predicts single-cell perturbation
493 responses. *Nature methods*, 16(8):715–721, 2019.
- 494 [25] M. Lotfollahi, A. Klimovskaia, C. De Donno, Y. Ji, I. L. Ibarra, F. A. Wolf, N. Yakubova,
495 F. J. Theis, and D. Lopez-Paz. Learning interpretable cellular responses to complex
496 perturbations in high-throughput screens. *bioRxiv*, 2021.
- 497 [26] M. Lotfollahi, M. Naghipourfar, M. D. Luecken, M. Khajavi, M. Büttner, M. Wagen-
498 stetter, Ž. Avsec, A. Gayoso, N. Yosef, M. Interlandi, et al. Mapping single-cell data to
499 reference atlases by transfer learning. *Nature Biotechnology*, pages 1–10, 2021.

- 500 [27] W. Mao, E. Zaslavsky, B. M. Hartmann, S. C. Sealfon, and M. Chikina. Pathway-level
501 information extractor (plier) for gene expression data. *Nature methods*, 16(7):607–610,
502 2019.
- 503 [28] H. Mathys, J. Davila-Velderrain, Z. Peng, F. Gao, S. Mohammadi, J. Z. Young,
504 M. Menon, L. He, F. Abdurrob, X. Jiang, et al. Single-cell transcriptomic analysis
505 of alzheimer’s disease. *Nature*, 570(7761):332–337, 2019.
- 506 [29] J. M. McFarland, B. R. Paoella, A. Warren, K. Geiger-Schuller, T. Shibue, M. Roth-
507 berg, O. Kuksenko, W. N. Colgan, A. Jones, E. Chambers, et al. Multiplexed single-cell
508 transcriptional response profiling to define cancer vulnerabilities and therapeutic mech-
509 anism of action. *Nature Communications*, 11(1):1–15, 2020.
- 510 [30] D. Risso, F. Perraudeau, S. Gribkova, S. Dudoit, and J.-P. Vert. A general and flexible
511 method for signal extraction from single-cell rna-seq data. *Nature Communications*, 9
512 (1):1–17, 2018.
- 513 [31] S. Rybakov, M. Lotfollahi, F. J. Theis, and F. A. Wolf. Learning interpretable latent
514 autoencoder representations with annotations of feature sets. *bioRxiv*, 2020.
- 515 [32] K. A. Severson, S. Ghosh, and K. Ng. Unsupervised learning with contrastive latent
516 variable models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol-
517 ume 33, pages 4862–4869, 2019.
- 518 [33] T. Stuart, A. Butler, P. Hoffman, C. Hafemeister, E. Papalexi, W. M. Mauck III, Y. Hao,
519 M. Stoeckius, P. Smibert, and R. Satija. Comprehensive integration of single-cell data.
520 *Cell*, 177(7):1888–1902, 2019.
- 521 [34] V. Svensson, A. Gayoso, N. Yosef, and L. Pachter. Interpretable factor models of single-
522 cell rna-seq via variational autoencoders. *Bioinformatics*, 36(11):3418–3421, 2020.
- 523 [35] L. T. Vassilev, B. T. Vu, B. Graves, D. Carvajal, F. Podlaski, Z. Filipovic, N. Kong,
524 U. Kammlott, C. Lukacs, C. Klein, N. Fotouhi, and E. A. Liu. In vivo activation of
525 the p53 pathway by small-molecule antagonists of mdm2. *Science*, 303(5659):844–848,
526 2004. doi: 10.1126/science.1092472.
- 527 [36] A. J. Wilk, A. Rustagi, N. Q. Zhao, J. Roque, G. J. Martínez-Colón, J. L. McKechnie,
528 G. T. Ivison, T. Ranganath, R. Vergara, T. Hollis, et al. A single-cell atlas of the
529 peripheral immune response in patients with severe covid-19. *Nature Medicine*, 26(7):
530 1070–1076, 2020.

- 531 [37] F. A. Wolf, P. Angerer, and F. J. Theis. Scanpy: large-scale single-cell gene expression
532 data analysis. *Genome Biology*, 19(1):1–5, 2018.
- 533 [38] F. Wu, J. Fan, Y. He, A. Xiong, J. Yu, Y. Li, Y. Zhang, W. Zhao, F. Zhou, W. Li,
534 et al. Single-cell profiling of tumor heterogeneity and the microenvironment in advanced
535 non-small cell lung cancer. *Nature Communications*, 12(1):1–11, 2021.
- 536 [39] G. X. Y. Zheng, J. M. Terry, P. Belgrader, P. Ryvkin, Z. W. Bent, R. Wilson, S. B.
537 Ziraldo, T. D. Wheeler, G. P. McDermott, J. Zhu, M. T. Gregory, J. Shuga, L. Mon-
538 tesclaros, J. G. Underwood, D. A. Masquelier, S. Y. Nishimura, M. Schnall-Levin, P. W.
539 Wyatt, C. M. Hindson, R. Bharadwaj, A. Wong, K. D. Ness, L. W. Beppu, H. J. Deeg,
540 C. McFarland, K. R. Loeb, W. J. Valente, N. G. Ericson, E. A. Stevens, J. P. Radich,
541 T. S. Mikkelsen, B. J. Hindson, and J. H. Bielas. Massively parallel digital transcrip-
542 tional profiling of single cells. *Nature Communications*, 8(1):14049, 2017.
- 543 [40] J. Y. Zou, D. J. Hsu, D. C. Parkes, and R. P. Adams. Contrastive learning using spectral
544 methods. *Advances in Neural Information Processing Systems*, 26:2238–2246, 2013.