# A Nonparametric Multi-view Model for Estimating Cell Type-Specific Gene Regulatory Networks

**Cassandra Burdziak** [* 1] **Elham Azizi** [* 1] **Sandhya Prabhakaran** [1] **Dana Pe'er** [1]

## Abstract

We present a Bayesian hierarchical multi-view mixture model termed *Symphony* that simultaneously learns clusters of cells representing cell types and their underlying gene regulatory networks by integrating data from two views: single-cell gene expression data and paired epigenetic data, which is informative of gene-gene interactions. This model improves interpretation of clusters as cell types with similar expression patterns as well as regulatory networks driving expression, by explaining gene-gene covariances with the biological machinery regulating gene expression. We show the theoretical advantages of the multi-view learning approach and present a Variational EM inference procedure. We demonstrate superior performance on both synthetic data and real genomic data with subtypes of peripheral blood cells compared to other methods.

## 1. Introduction

Joint analysis of different types of data that are associated with the same underlying phenomenon is more informative than analysis of individual data types, and increases signal to noise ratio (Xu et al., 2013; Wang et al., 2013; Rey & Roth, 2012). This approach, known as multi-view learning or learning with multiple distinct features, has been successfully used in various settings (Li et al., 2002; Jones & Viola, 2003; Hardoon et al., 2004; Pan et al., 2007).

Here, we apply such a multi-view learning approach to address an important biological problem by integrating two views of gene regulation. Our goal is to infer cell clusters (characterizing cell types) as well as their underlying gene regulatory networks (GRNs), which are directed weighted networks between genes depicting the extent to which a regulatory gene influences the expression of each of its downstream *target* genes. Understanding differences between regulatory mechanisms across different cell types provides valuable insight in normal development of cell types (Davidson, 2010), and mechanisms disrupted in cancer cells (Pe'er & Hacohen, 2011; Kreeger & Lauffenburger, 2009).

Recent advances in single-cell genomic technologies (Hashimshony et al., 2012; Jaitin et al., 2014; Shalek et al., 2013) which measure gene expression at the resolution of individual cells, present remarkable opportunities to characterize different cell types by clustering cells based on heterogeneity of gene expression (as observed features) (Satija et al., 2015; Macosko et al., 2015). Learning GRNs from gene expression data alone, however, leads to detection of spurious network links based on correlated genes, while integrative learning from multiple data sources has been shown to improve overall joint inference (Zhu et al., 2008; Hecker et al., 2009; Azizi et al., 2014). Therefore, we aim to identify GRNs driving heterogeneous cell types through integrating single-cell expression data with other genomic data types. In particular, epigenetic technologies such as ATAC-seq (Buenrostro et al., 2015a), scan the genome for accessible DNA regions, identifying potential interaction between a gene and regulator proteins translated from other genes. In other words, epigenetic data contains information about direct regulatory links between genes and incorporation of epigenetic data is a promising direction for improved inference of GRNs (Guo et al., 2017; Rotem et al., 2015).

We present a novel integrative model, which we refer to as *Symphony*, as a Dirichlet process mixture model that jointly learns clusters of cells and GRNs specific to each cluster. *Symphony* is an extension of the BISCUIT model (Prabhakaran et al., 2016; Azizi et al., 2018) which clusters cells while simultaneously distinguishing biological heterogeneity from technical noise in single-cell gene expression data. This is done through incorporating cell-specific parameters scaling the cluster means and covariances for a multivariate Gaussian mixture model.

We extend the BISCUIT model and replace the hyperparameters with a generative process exclusively driven by the paired epigenetic data, which captures the biological mech-

---

*Equal contribution [1]Computational & Systems Biology Program, Memorial Sloan Kettering Cancer Center, New York, NY, USA. Correspondence to: Cassandra Burdziak <cnb3001@med.cornell.edu>, Elham Azizi <mail@elhamazizi.com>.
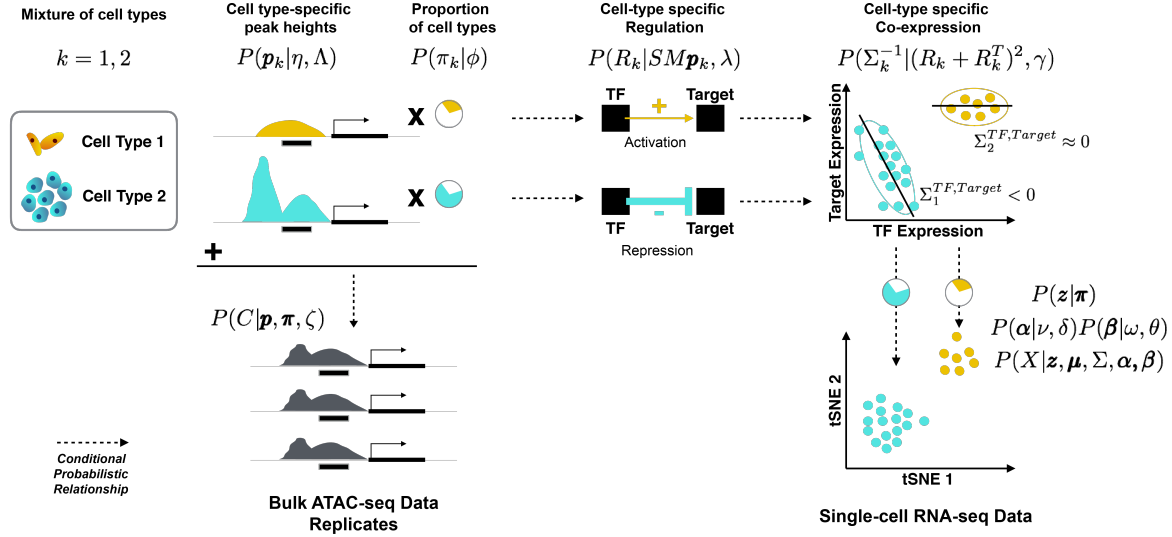
*Figure 1.* The generative process in *Symphony*: We aim to infer gene regulatory networks (GRNs) denoted as $R_k$ specific to each cluster of cells from integration of two views: Epigenetic data (indicative of network edges) from the bulk of cells ($C$) and expression of genes (network nodes) at resolution of single cells ($X$). GRNs are directed weighted networks with edge weights denoting regulatory impact of one gene on another (e.g. activation or repression); the impact of regulation is reflected in covariance between genes ($\Sigma_k$).

anism responsible for observed gene covariances per cell type. Briefly, the epigenetic profiles which denote accessible DNA in the bulk samples are deconvolved into cell-type specific accessible regions (Figure 1). Within these regions, the binding of regulatory proteins translated from genes impacts the expression of nearby genes, such that accessible regions may be mapped to gene-gene interactions. This mapping is based on prior knowledge of recurring DNA sequences (known as motifs) associated with these regulatory proteins which occur in regions of accessible DNA. Most importantly, the covariance in observed gene expression is related to a graph power of the regulatory network, capturing the propagated impact of regulation in the network (indirect regulation)(Figure 2).

This multi-view framework can also be applied to other settings, such as text characterization. For example, to learn the context of queries (vector of words), the bag-of-words simplification may not be sufficient (Biemann, 2005). However, the order of words, which can be represented as a latent directed network can imply the context, and incorporating observations from this network such as the frequency of one word following another (as a second view) can enhance extraction of context and clustering of queries (observed as the first view) (Landauer et al., 1997; Recchia & Jones, 2009).

**Related technologies and methods.** The problem of inferring GRNs specific to cell types involves identifying differences in gene-gene interactions across cell types. However, in most cases the cell types are not well-characterized, hence gene markers are not known to enable sorting of cell types

prior to measuring epigenetic data (gene-gene interactions). Therefore, epigenetic data measured on the bulk of cells represents a mixture of cell type-specific epigenetic profiles. One solution is measuring epigenetic profiles at the resolution of single cells. These technologies have only recently emerged (Buenrostro et al., 2015b); we therefore constructed a model to allow integration of bulk ATAC-seq data. *Symphony* can be easily adapted for inferring GRNs from single-cell ATAC-seq data as well.

Other works have attempted to apply computational deconvolution algorithms intended for bulk expression data, such as those using source separation techniques (Houseman et al., 2016), NMF-based methods (Repsilber et al., 2010) or Bayesian models (Erkkilä et al., 2010), to instead infer cell type-specific epigenetic profiles. Recent methods such as SCENIC (Aibar et al., 2017) infer GRNs from single-cell expression data alone and do not incorporate epigenetic or other types of data.

**Contributions.** In this paper, we show that a multi-view learning framework would improve the deconvolution of epigenetic data. Furthermore, using an integrative model, we improve the clustering of cells, and hence characterization of cell types. Most importantly, our model presents the advantage of inferring cell type-specific GRNs that give insight into heteroegeneity of underlying mechanisms across cell types. We present a Variational EM inference procedure and show that the integration guarantees model identifiability, while learning from the epigenetic view alone does not. While other works have attempted to integrate bulk multiomics data (Lake et al., 2017; Brown et al., 2013; Ritchie
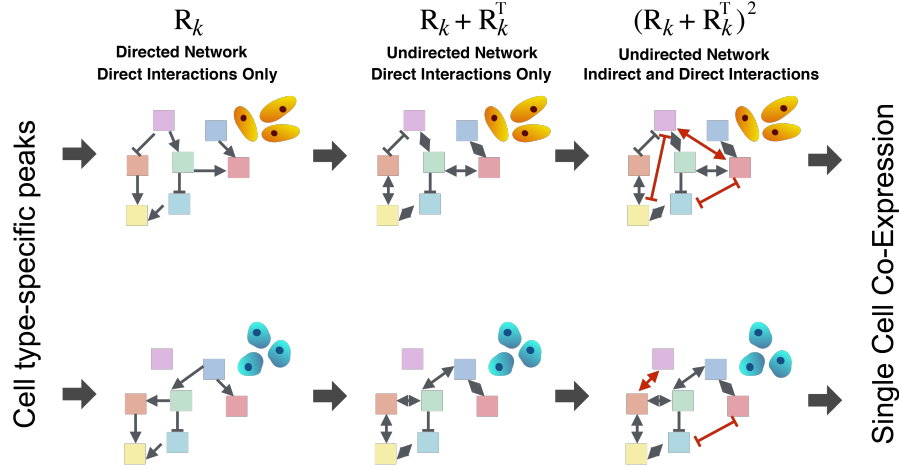
*Figure 2. Symphony* captures direct and indirect regulation. The impact of regulation ($R_k$) is propagated through the network up to path length of two and is reflected in covariance between indirectly connected genes ($\Sigma_k$).
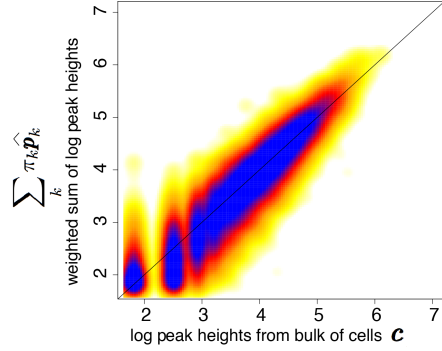


*Figure 3.* Validation of model assumption. Weighted sum of actual peak heights $\hat{p}_k$s measured from sorted clusters of CD34+ hematopoeitic cells using ATAC-seq, with weights proportional to proportions of cell types, compared to measured peak height from ATAC-seq on the bulk of cells $c$; heatmap shows density with yellow (low) to blue (high).

et al., 2015), there are no methods to our knowledge that infer heterogeneous GRNs through integrating epigenetic and single-cell resolution gene expression data.

## 2. Model

The observed data is considered as two views from the biological system (Figure 1):

**View 1.** Single-cell gene expression data from scRNA-seq technologies (Klein et al., 2015; Macosko et al., 2015) denoted as $X^{d \times n} = [\boldsymbol{x}_1, \cdots, \boldsymbol{x}_j, \cdots, \boldsymbol{x}_n]$ where each observation $\boldsymbol{x}_j \in \mathbb{R}^d$ for cell $j \in \{1, \cdots, n\}$ corresponds to $d$ genes (as features). Each entry $x_{ij}$ for $i = [1, \cdots, d]$ contains the expression of gene $i$ in cell $j$ (more precisely, the log of counts of mRNA molecules per gene $i$ from cell $j$ plus a pseudo-count).

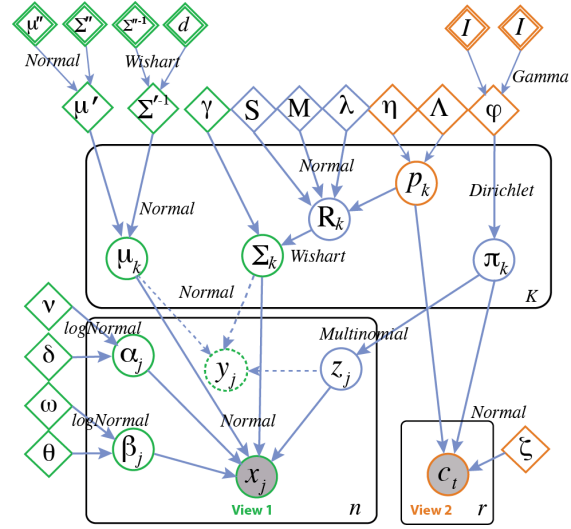**View 2.** Epigenetic data, for example measured with ATAC-



*Figure 4.* Plate model for *Symphony*. white circles denote latent variables of interest, diamonds are hyperparameters and double diamonds are hyperpriors calculated empirically.

seq technology (Buenrostro et al., 2015a), denoted as $C^{l \times r} = [\boldsymbol{c}_1, \cdots, \boldsymbol{c}_t, \cdots, \boldsymbol{c}_r]$ where each observation $\boldsymbol{c}_t \in \mathbb{R}^l$ for $t \in \{1, \cdots, r\}$ corresponds to $l$ genomic regions (as features). Specifically, $\boldsymbol{c}_t$ is an experimental replicate measuring accessibility of genomic regions $m = [1, \cdots, l]$.

**Prior knowledge.** The genomic regions in $C$ can be mapped to genes in $X$ with a pre-defined mapping function $g(i, i') = m$ that relates each genomic region $m \in \{1, \cdots, l\}$ to a gene-gene interaction $i' \to i$ for $i, i' \in \{1, \cdots, d\}$. We also define $M^{d \times d}$ based on prior knowledge containing binary values $M_{i,i'} = 1$ if the motif sequence for gene $i'$ exists in the genomic region $m$ in the vicinity of gene $i$, meaning a potential interaction can exist

from gene $i'$ to gene $i$.

## 2.1. Epigenetic Model (View 2)

The epigenetic data is informative of network structure, i.e. existence of edges between genes (features) and gene expression data contains information on network nodes. We aim to infer this directed weighted network (GRN) for each cluster $k \in \{1, \cdots, K\}$ of cells, denoted by the asymmetric matrix $R_k^{d \times d}$ in which entry $R_{k_{i,i'}} \neq 0$ if $i' \rightarrow i$, meaning gene $i'$ directly regulates gene $i$. $R_{k_{i,i'}}$ is the regulatory function of gene $i'$ on gene $i$ in cluster $k$ such that $R_{k_{i,i'}} > 0$ or $R_{k_{i,i'}} < 0$ represent activation or repression of expression respectively, with $|R_{k_{i,i'}}|$ being the strength of regulation.

We do not aim to distinguish all layers of the regulatory process (such as protein phosphorylation) and rather interpret GRNs as an approximation for the overall impact of TFs on target genes at the transcriptional level.

We model regulation of gene expression as follows: Genome accessibility in cluster $k$ is represented with latent variable $\boldsymbol{p}_k = [p_k^1, \cdots, p_k^l]^T \in \mathbb{R}^{+l}$ containing $l$ genomic regions (features). This represents log of peak heights plus 1 (to ensure positive domain) in all genomic regions, for each cell type $k$. We set a truncated multivariate Normal prior to capture the structure between genomic regions encompassing co-regulated genes (i.e. genes sharing regulators) with mean $\boldsymbol{\eta}$ and covariance $\Lambda$: $\boldsymbol{p}_k \sim truncN(\boldsymbol{\eta}, \Lambda, \boldsymbol{0}, +\infty)$. In this paper, we assume a setting where we do not observe $\boldsymbol{p}_k$s, and only observe epigenetic data from the bulk of cells which can be represented as a weighted sum of cluster-specific epigenetic profiles where $\pi_k$s are weights. Thus, our epigenetic model is:

$$\{\boldsymbol{c}\}_t^{(1,\cdots,l)}|\boldsymbol{p}_k, \pi_k \overset{\text{ind}}{\sim} \mathcal{N}(\sum_k \pi_k \boldsymbol{p}_k, \zeta I) \quad (1)$$

We validated the above assumption of weighted sum using ATAC-seq data from hematopoeitic progenitor cells from Corces et al. (2016) by computing the weighted sum of measurements from sorted cell types (actual $\boldsymbol{p}_k$s denoted with $\hat{\boldsymbol{p}}_k$) to measurements from the bulk of cells as $\boldsymbol{c}$ (Figure 3).

We have a $K$-order Dirichlet prior over $\pi_k$: $\pi_k|\varphi, K \sim Dir(\pi_k|\frac{\varphi}{K}, \cdots, \frac{\varphi}{K})$, where $\varphi^{-1} \sim Gamma(1,1)$. Then, in cell type $k$, if a genomic region $m$ in the vicinity of gene $i$ is accessible with log peak height $p_k^m$, and the motif sequence associated with one or more Transcription Factor (TF) proteins translated from genes $i', i'', \ldots$ exists in the region ($M_{i,i'} = M_{i,i''} = 1$), then the TF(s) can bind to the region and hence regulate the expression of gene $i$. Furthermore, we assume the peak height (strength of TF binding to genome) is informative of $i' \rightarrow i$ edge weight $|R_k|$ (strength of regulation). Thus, we model $R_k$ as follows:

$$R_k^{i,i'} \sim \mathcal{N}(S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}, \lambda) \quad (2)$$

The function $g(\cdot)$ maps gene pair $i, i'$ to genomic region $l$. $S$ denotes a sign indicator variable representing repression or activation function. We set $S$ according to the sign of the empirical covariance: $S^{i,i'} = sign(\Sigma''^{i,i'})$.

## 2.2. Single-cell Gene Expression Model (View 1)

We use the above epigenetic model to drive gene expression data based on the following key ideas: First, if a direct regulatory link exists from a TF associated with gene $i'$ to a target gene $i$ ($i' \rightarrow i$), we assume there is strong covariance between their expressions. Second, due to van der Corput's inequality (Montgomery, 2001), covariances can partly reflect the propagated impact of indirect regulation in cases where genes are not directly connected but exist on the same path in the network (Figure 2). For example if $i'' \rightarrow i'$ and $i' \rightarrow i$, we might also observe covariance between $i, i''$ even though they are not directly connected in the network (e.g. $R_k^{i,i''} = 0, \Sigma_k^{i,i''} \neq 0$). Here, we consider indirect effects with path length up to two using the square of the indirected network $(R_k + R_k^T)^2$ (Walker, 1992) such that:

$$\Sigma_k^{-1}|R_k \sim Wish((R_k + R_k^T)^{-2}, \gamma) \quad (3)$$

$(R_k + R_k^T)^2$ is positive semi-definite according to Lemma 3 in the following section making the above modeling assumption feasible. Additionally, this model can capture combinatorial regulation in the inferred covariances. In particular, a gene pair $i, i'$ will always have the same directionality of regulation (i.e. activation or repression relationship), but $\Sigma_k^{i,i'}$ can be positive in one cluster and negative in another cluster depending on the relative regulatory strength of activators and regulators in its path. An example of this variability in sign is shown in Supplmentary Figure 12.

Gene expression data for each cell $j$ denoted as $\boldsymbol{x}_j$ is then modeled similar to the multivariate Gaussian mixture model in BISCUIT:

$$\begin{aligned} \{\boldsymbol{x}\}_j^{(1,\cdots,d)}|z_j = k &\overset{\text{ind}}{\sim} \mathcal{N}(\alpha_j \boldsymbol{\mu}_k, \beta_j \Sigma_k) \\ \boldsymbol{\mu}_k \sim \mathcal{N}(\boldsymbol{\mu}', \Sigma'), \quad \boldsymbol{\mu}' &\sim \mathcal{N}(\boldsymbol{\mu}'', \Sigma'') \\ \Sigma'^{-1} \sim Wish(d, \frac{1}{d\Sigma''}), \quad z_j|\pi_k &\overset{\text{iid}}{\sim} Mult(z_j|\pi_k) \\ \alpha_j \sim log\mathcal{N}(\nu, \delta^2), \quad \beta_j &\sim log\mathcal{N}(\omega, \theta) \end{aligned} \quad (4)$$

where $z_j$ denotes assignment of cell $j$ to cluster $k \in \{1, \cdots, K\}$. With integrating two data modalities (gene expression and epigenetic data), we improve inference of clusters (cell types). The scaling parameters $\alpha_j, \beta_j$ specific to cell $j$ are used to normalize the data $\boldsymbol{x}_j$ in downstream analysis by transforming to $\boldsymbol{y}_j \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)$ according to the cluster it is assigned to $z_j = k$, similar to Prabhakaran et al. (2016). The plate model for *Symphony* is summarized in Figure 4.

## 3. Theory

We show theoretical advantages of integration of the two data types using *Symphony* as follows: We define $f(\boldsymbol{x}) := \mathcal{N}(\alpha\boldsymbol{\mu}_k, \beta\Sigma_k) \in \mathbb{R}^d$ as the multivariate Gaussian density of $\boldsymbol{x}$ and $f(\boldsymbol{c}) := \mathcal{N}(\sum_k \pi_k \boldsymbol{p}_k, \zeta I) \in \mathbb{R}^l$ as the multivariate Gaussian density of $\boldsymbol{c}$. First, we emphasize that the epigenetic model alone $f(\boldsymbol{c}|\boldsymbol{p}_k, \pi_k, \zeta)$ is not identifiable and therefore precise inference of deconvolved epigenetic profiles ($\boldsymbol{p}_k$s) is not possible:

**Lemma 1** *The epigenetic model $f(\boldsymbol{c}|\boldsymbol{p}_k, \pi_k)$ is non-identifiable (Proof in Supplementary section B)*

This motivated us to build an integrative model. Identifiability of the single-cell expression model $f(\boldsymbol{x}|\boldsymbol{\mu}_k, \Sigma_k, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{z})$ has been shown under certain enforced constraints on both $\boldsymbol{\alpha}, \boldsymbol{\beta}$:

**Lemma 2** *(Prabhakaran et al., 2016) Defining $\Phi = \{\forall j, k : (\alpha_j, \boldsymbol{\mu}_j, \beta_j, \Sigma_k)\} \cup \{\boldsymbol{\pi}\}\Phi = \Phi^\star$ if the following conditions hold: $\forall j : \boldsymbol{\mu}_k \geq \boldsymbol{\mu}' + diag(\Sigma')(\alpha_j - \nu)/\delta$ and $\forall j : \beta_j \geq \frac{\theta}{\omega+1}$*

While the above conditions guarantee identifiability, they are not inferred from data or biologically motivated and hence interpretation of parameters may not provide the best characterization of cell types. Here, we show that in the integrative model, the constraints for $\beta_j$s are no longer required and identifiability of the expression model $f(\boldsymbol{x}|\boldsymbol{\mu}_k, \Sigma_k, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{z}, \boldsymbol{p}_k, R_k)$ is guaranteed through the extension of the model that captures regulation, from which we observe additional epigenetic data $\boldsymbol{c}$. Hence the full model $f(X, C|\{\boldsymbol{\mu}_k, \Sigma_k, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{z}, \boldsymbol{p}_k, R_k, \pi_k, \zeta\})$ is identifiable. We will use the following lemma to show $(R_k + R_k^T)^2$ is positive semi-definite.

**Lemma 3** *Square of a symmetric matrix $H$ gives a symmetric positive semi-definite matrix $L$ (Proof in Supplementary section B).*

Since single-cell expression data $X$ usually contains expression of genes that do not have observations in their mapped genomic regions in $\boldsymbol{c}$ such that $d > l$, we first define a reduced version of the model where all genes in $\boldsymbol{x}$ do have mapped genomic regions: $f(X^0, C)$ where $X^0$ is a $l \times n$ subset of $X$. Then, $f(\boldsymbol{x}_{l \times 1}^0) = \mathcal{N}(\alpha\boldsymbol{\mu}_k^0, \beta\Sigma_k^0)$ with $\boldsymbol{\mu}_k^0, \Sigma_k^0$ being subsets of $\boldsymbol{\mu}$ and $\Sigma$. We next show the identifiability of the reduced model. Then, we use this result to extend the identifiability to the full model given $\boldsymbol{\beta}$ which is the parameter scaling $\Sigma_k$. Finally, we show the identifiability of the full model.

**Lemma 4** *In the reduced model: $f(X^0, C|\boldsymbol{\beta}, \boldsymbol{\mu}_k^0, \Sigma_k^0, \boldsymbol{\alpha}, \boldsymbol{z}, \boldsymbol{p}_k, R_k^0, \pi_k, \zeta)$, $\beta$s are identifiable under the conditions of: $\forall j : \boldsymbol{\mu}_k \geq \boldsymbol{\mu}' + diag(\Sigma')(\alpha_j - \nu)/\delta$ without the need for condition on $\beta$s (Proof in Supplementary section B).*

**Lemma 5** *For a given $\boldsymbol{\beta} = \beta^*$, identifiability of: $f(X, C|\boldsymbol{\mu}_k, \Sigma_k, \boldsymbol{\alpha}, \boldsymbol{\beta} = \boldsymbol{\beta}^*, \boldsymbol{z}, \boldsymbol{p}_k, R_k, \pi_k, \zeta)$ is guaranteed if $\forall j : \boldsymbol{\mu}_k \geq \boldsymbol{\mu}' + diag(\Sigma')(\alpha_j - \nu)/\delta$ (Proof in Supplementary section B).*

**Theorem 6** *The full model $f(X, C|\{\boldsymbol{\mu}_k, \Sigma_k, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{z}, \boldsymbol{p}_k, R_k, \pi_k, \zeta\}))$ is identifiable if $\forall j : \boldsymbol{\mu}_k \geq \boldsymbol{\mu}' + diag(\Sigma')(\alpha_j - \nu)/\delta$ (Proof in Supplementary section B).*

## 4. Inference

We applied the co-ordinate ascent mean field variational inference (CAVI) (Blei et al., 2017; Ghahramani & Beal, 2001) which assigns independent factors to the latent variables. The blueprint for corresponding CAVI updates are below and full derivations are presented in Supplementary section C.

**Variational E-step**

**a.** $q^*(z_j) = \prod_k r_{jk}^{z_{jk}}$ where

$$r_{jk} = \mathbb{E}_{\boldsymbol{z}} z_{jk} \propto |\widetilde{\beta_j \Sigma_k}|^{-1} \exp(-S_2)\widetilde{\pi_k}, \quad \sum_k r_{nk} = 1$$

$$S_2 = \frac{1}{2}\Big(tr(\Sigma_k^{-1}\beta_j^{-1}\Sigma''^{-1}) +$$
$$(\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)^T(\beta_j\Sigma_k)^{-1}(\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)\Big)$$

$$|\widetilde{\beta_j \Sigma_k}|^{-1} := \frac{1}{2}(-d\mathbb{E}_{\beta_j}\ln\beta_j +$$
$$d\ln(2) + \ln|R_k^*| + \sum_{i=1}^d \psi\Big(\frac{\gamma + 1 - i}{2}\Big))$$

$$(5)$$

where $\ln\widetilde{\pi_k} := \psi(\varphi_0) - \psi(\sum_k \varphi_k)$ and $\psi$ is the *digamma* function and $R_k^* = (R_k + R_k^T)^2$.

**Variational M-step**

**b.** $q^*(\pi_k) = \pi_k \sim$ Stick-breaking Beta$(1, \varphi)$

**c.** $q^*(\boldsymbol{\mu}_k) = \exp\Big(-\frac{1}{2}\sum_j r_{jk}\Big(tr(\Sigma_k^{-1}(\frac{\beta_j}{\alpha_j^2})^{-1}\Sigma''^{-1})$
$$+ (\bar{\boldsymbol{\mu}_k} - \frac{x_j}{\alpha_j})^T(\frac{\beta_j}{\alpha_j^2}\Sigma_k)^{-1}(\bar{\boldsymbol{\mu}_k} - \frac{x_j}{\alpha_j})$$
$$+ (\boldsymbol{\mu}_k - \boldsymbol{\mu}')^T\Sigma'^{-1}(\boldsymbol{\mu}_k - \boldsymbol{\mu}')\Big) + c\Big)$$
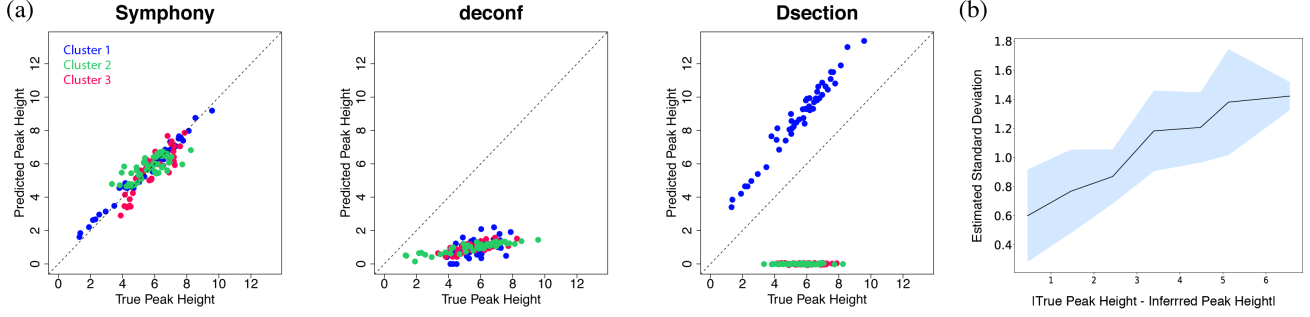
$$(6)$$

*Figure 5.* Performance in deconvolving epigenetic data. **(a)** Estimated peak heights $\boldsymbol{p}_k$ using *Symphony* for $K = 3$ synthetic clusters versus true peak heights, compared to two other deconvolution methods: deconf (Repsilber et al., 2010) and Dsection (Erkkilä et al., 2010); each dot represents a genomic region. **(b)** Moving average of estimated standard deviation for $\boldsymbol{p}_k$s using *Symphony* vs. estimate residual; shaded area shows 1 standard deviation in each window of length 1.
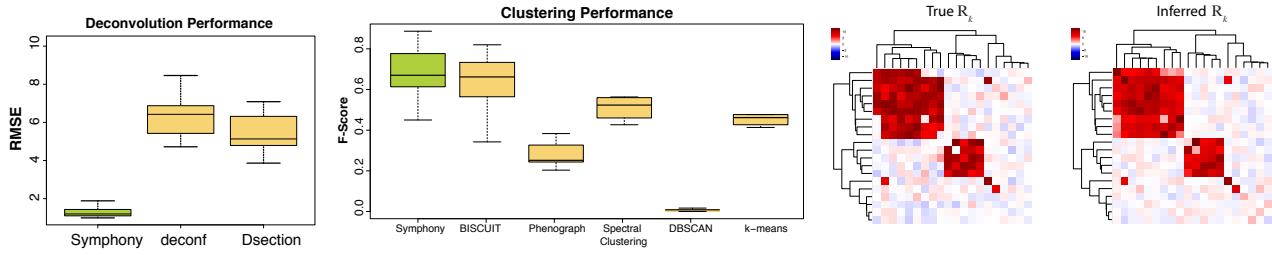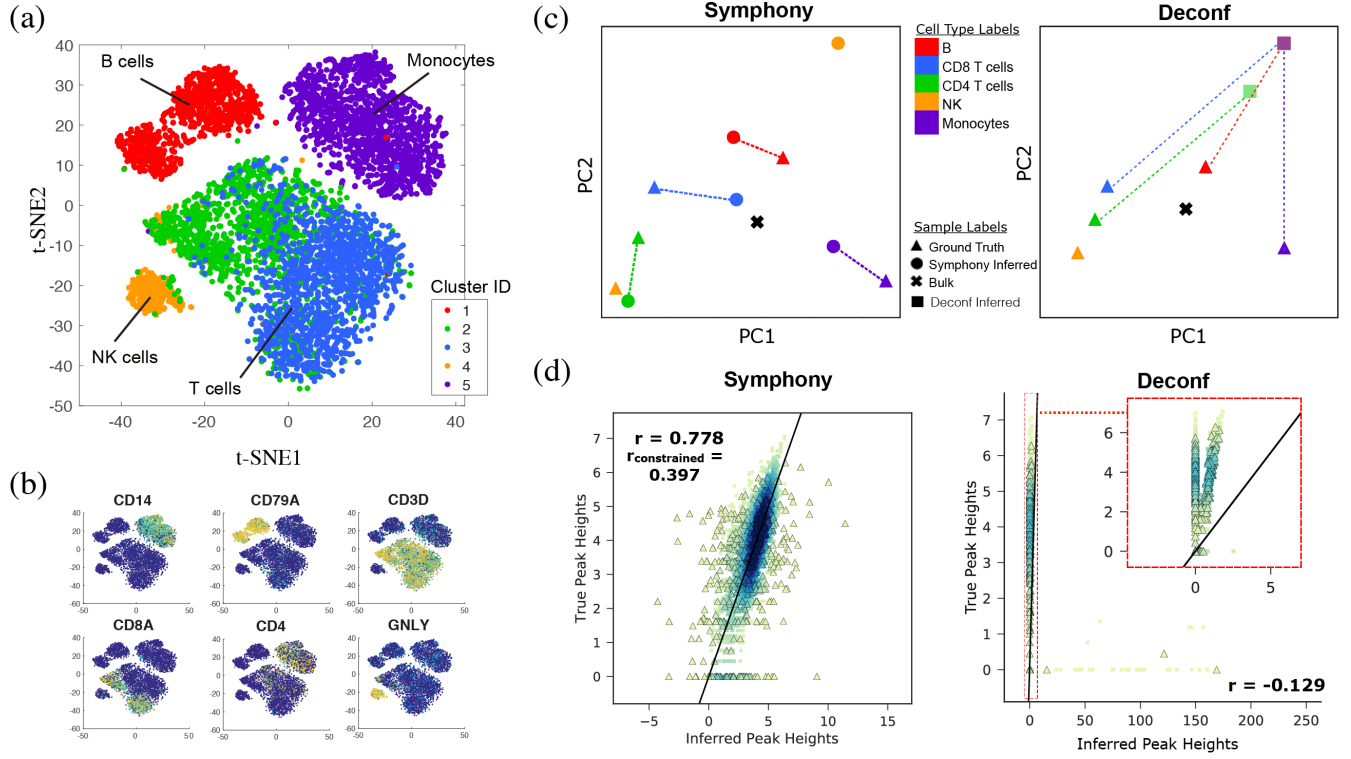


*Figure 6.* **Left:** Root-mean-square error (RMSE) in inferring $\boldsymbol{p}_k$s across 10 experiments compared to other deconvolution methods used for genomic data. **Middle:** Performance in clustering cells across 10 experiments compared to other clustering methods commonly used for single-cell gene expression data. **Right:** Heatmap depicting true versus inferred $R_k$ in synthetic data showing *Symphony's* superior capabilities in recovering underlying $R_k$.

**d.** $q^*(\Sigma_k^{-1}) = \exp\Big(-\frac{1}{2}\sum_j r_{jk}\Big(\ln\beta_j^d + \ln|\Sigma_k|$

$\qquad + \frac{1}{\beta_j}\Big(tr(\Sigma_k^{-1}\Sigma''^{-1})$

$\qquad + (\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)^T(\Sigma_k)^{-1}(\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)\Big)\Big)$

$\qquad + \ln|\Sigma_k^{-1}|^{\frac{\lambda-d-1}{2}} + \{-\frac{tr(R_k^{*-1}\Sigma_k^{-1})}{2}\}$

$\qquad - \ln(2^{\lambda d/2}|R_k^{*-1}|^{-\lambda/2}\Gamma_d(\lambda/2)) + c\Big)$

**e.** $q^*(\alpha_j) = \exp\Big(-\sum_k r_{jk}S_2$

$\qquad + \frac{1}{\sqrt{2\delta^2\pi}}\exp\Big(-\frac{(\ln\alpha_j - \nu)\mathbb{I}_{1\times 1}(\ln\alpha_j - \nu)}{2\delta^2}\Big) + c\Big)$

**f.** $q^*(\beta_j) = \exp\Big(-\sum_k r_{jk}S_2$

$\qquad + \frac{1}{\sqrt{2\theta^2\pi}}\exp\Big(-\frac{(\ln\beta_j - \omega)\mathbb{I}_{1\times 1}(\ln\beta_j - \omega)}{2\theta^2}\Big) + c\Big)$

**g.** $q^*(\boldsymbol{\mu}') \sim \mathcal{N}(\boldsymbol{\mu}_{\mu'}, \Sigma_{\mu'})$

$$(7)$$

**h.** $q^*(\Sigma^{-1\prime}) \sim \mathcal{W}(V_{\Sigma'^{-1}}, d_{\Sigma'^{-1}})$

**i.** $q^*(R_k) = \exp\Big(\{-\frac{tr(R_k^{*-1}\Sigma_k^{-1})}{2}\}$

$\qquad - \ln(2^{\lambda d/2}|R_k^{*-1}|^{-\lambda/2}\Gamma_d(\lambda/2))$

$\qquad -\frac{1}{2\delta^2}\sum_i\sum_{i'}\Big((R_k^{i,i'} - S^{i,i'}M^{i,i'}\boldsymbol{p}_k^{g(i,i')})\mathbb{I}_{1\times 1}(R_k^{i,i'} -$

$\qquad\qquad S^{i,i'}M^{i,i'}\boldsymbol{p}_k^{g(i,i')})\Big) + c\Big)$

**j.** $q^*(\boldsymbol{p}_k) = \exp\Big((\bar{c}_t - \sum_k(\pi_k\boldsymbol{p}_k))^T(\zeta\mathbb{I})^{-1}(\bar{c}_t - \sum_k(\pi_k\boldsymbol{p}_k))]$

$\qquad + \ln\Big(\frac{1}{\sqrt{2\pi\Lambda}}\exp(-\frac{1}{2}(\frac{\boldsymbol{p}_k-\eta}{\Lambda})^2))\Big)$

$\qquad -\frac{1}{2\delta^2}\sum_i\sum_{i'}\Big((R_k^{i,i'} - S^{i,i'}M^{i,i'}\boldsymbol{p}_k^{g(i,i')})\mathbb{I}_{1\times 1}$

$\qquad\qquad (R_k^{i,i'} - S^{i,i'}M^{i,i'}\boldsymbol{p}_k^{g(i,i')})\Big) + c\Big)$

$$(8)$$

where $c$ is the integration constant (details are included in Supplementary section C). Since the E-step takes $\approx \mathrm{O}(d^3)$ due to three matrix inversions, we substitute this with $z_{MAP}(\boldsymbol{x}) = \arg\max_z p(\boldsymbol{x}|\boldsymbol{z})p(\boldsymbol{z}|\boldsymbol{\pi})$ which takes $\approx \mathrm{O}(d^2)$

Figure 7. Performance on genomic data from PBMCs. **(a)** t-SNE map of 6825 single-cell gene expression data from 5 cell types after normalization, cells colored by clusters **(b)** expression of cell-type markers **(c)** PCA on expression-constrained peaks and those which show at least some accessibility in all cell types (to show effective peak *magnitude* estimation) showing global performance of Symphony in deconvolving epigenetic data; shown with projection of inferred peak heights on principal components of ground truth peak heights (from ATAC-seq on sorted cell types in Corces et al. (2016)) (left) compared to deconvolution using Deconf (right) which fails to deconvolve the majority of peaks shown with overlapping squares **(d)** Scatterplot of inferred peak heights for all clusters vs ground truth peak height using Symphony (left) compared to Deconf (right); peaks are colored by density; $r$ values show Pearson correlation; peaks constrained by expression data and bulk epigenetic data are triangular points with black outline; NK cells were not included in this plot due to the small cell proportion ($< 5\%$), making deconvolution impossible.

when $\Sigma_k^{-1}$s and $\Sigma''$ are apriori Cholesky decomposed.

Given the complexity of the model (refer to Supplementary section C), we implemented *Symphony* using the probabilistic programming language Stan (Carpenter et al., 2016). Furthermore, for a scalable implementation applicable to real genomic data containing thousands of cells, we used the probabilistic programming language, Edward (Tran et al., 2016; 2017) for variational EM (details and approximations presented in Supplementary section E).

## 5. Results

### 5.1. Synthetic Data

We first evaluated the performance in deconvolving epigenetic data, clustering cells, and inferring GRNs using data simulated from the *Symphony* model. We simulated data for $n = 100$ cells in $K = 3$ clusters with $d$ ranging from 5 to

20 genes and $l = 50$ using the *Symphony* model.

**Inference of $p_k, R_k$.** Figure 5(a) shows scatterplots of deconvolved peak heights ($p_k$) compared to actual data. We compared the performance of *Symphony* to two other deconvolution methods: *deconf* (Repsilber et al., 2010) which uses NMF, and *Dsection* (Erkkilä et al., 2010) which is based on a Bayesian model. While Dsection captures only the largest cluster, deconf underestimates the cluster-specific peak heights. The behavior of Dsection was reproducible across simulations, and is likely due to the lack of identifiability in the model for epigenetic data alone, as discussed above.

Figure 6 summarizes the error in estimating $p_k$s across 10 synthetic datasets with the same size as above. This shows the value of incorporating expression data (view 1) in deconvolution of epigenetic data (view 2). Figure 6 also shows a heatmap of inferred $R_k$ in one of the synthetic experiments
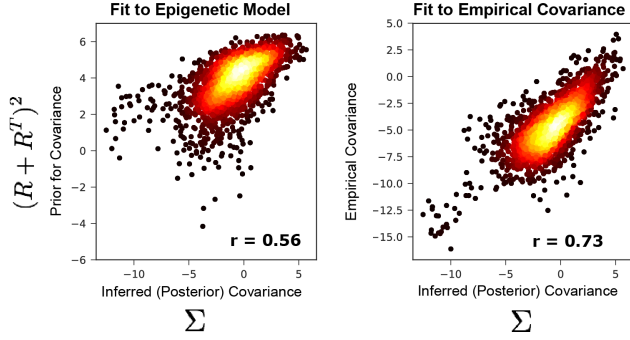
*Figure 8. Symphony* model fit. **Left:** Inferred covariance vs its prior capturing direct and indirect regulation $(R + R^T)^2$; **Right:** Empirical gene covariance compared to inferred covariance across all cell types. All axes are log-transformed.

as an example, compared to the actual $R_k$, confirming the ability of *Symphony* in inferring GRNs.

**Clustering performance.** We then show the performance in clustering with integrating both views as compared to only using gene expression data (view 1) by computing F-scores across 10 experiments with the same size as above. We compared the performance to BISCUIT (Prabhakaran et al., 2016), as well as other methods commonly used for clustering cells in single-cell gene expression data including DBscan (Satija et al., 2015), Phenograph (Levine et al., 2015), Spectral clustering (Ng et al., 2002) and k-means (with $K = 3$) (Figure 6). These results show improvement over BISCUIT due the epigenetic extension of the model and significant improvement over other methods. DBscan was unable to cluster the majority of cells, likely due to the small dimensionality of the feature space used in simulations. This shows the value of incorporating epigenetic data (view 2) in improving clustering performance, as compared to using expression data (view 1) alone. This has further value in biological interpretation of clusters as cell types that have both similar expression and similar underlying mechanisms driving expression.

### 5.2. Genomic Data

We also evaluated the performance of *Symphony* on real genomic data. We used previously published single-cell expression data for peripheral blood mononuclear cells (PBMCs) from Zheng et al. (2017) combined with ATAC-seq data for PBMCs from Corces et al. (2016). For single-cell expression data, we chose a subset of PBMCs from (Zheng et al., 2017) as $X$ containing $n = 6825$ cells which express known gene markers for either monocytes, B cells or T cells and NK cells. We chose to focus on $d = 28$ transcription factors which showed high standard deviation in expression data and are known to be lineage-defining factors.
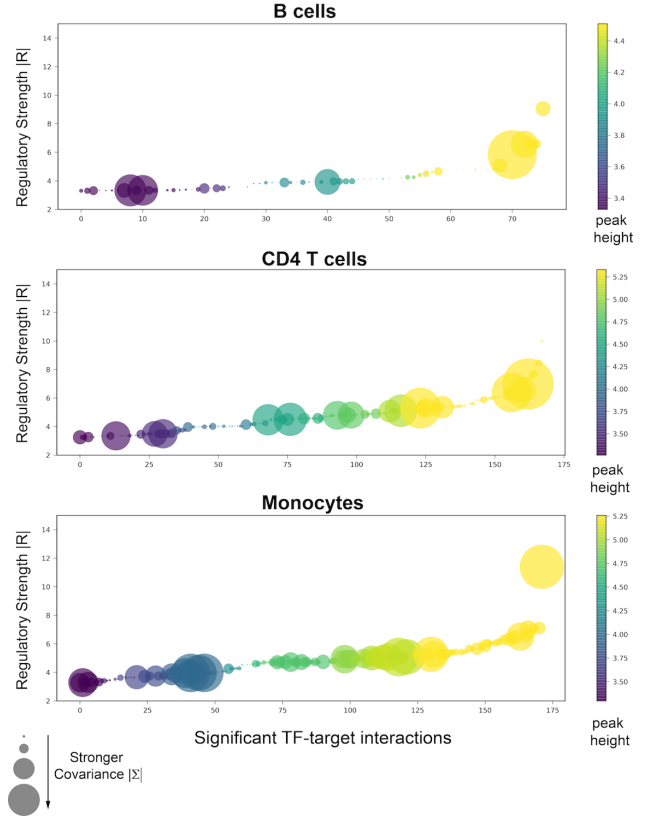


*Figure 9.* GRN Interpretation. Inferred regulations sorted by strength ($|R|$) show association with either peak heights or TF-target covariance or both; circle sizes are proportional to covariance strength ($|\Sigma|$) and they are colored by inferred peak height per cell type. Covariances are z-normalized for scale.

For epigenetic data matching the above cell types, we generated $r = 1$ mixture of epigenetic measurements with $l = 1053$ peaks from real ATAC-seq data collected from the above sorted cell types in Corces et al. (2016), with weights $\pi_k$ proportional to frequency of cell types in blood, and used this as observed epigenetic data $C$. We fixed the clustering in this experiment using Phenograph-derived assignments which we mapped onto nearest cell types to match with epigenetic data (Levine et al., 2015).

We determined non-zero entries in $M$ from ATAC-seq using the FIMO algorithm (Grant et al., 2011), which scans the sequence under the ATAC-seq peak for the occurrence of a motif. We associated a peak with the target gene closest in genomic distance to the peak in these experiments. This assignment is independent of the model structure and can be manually defined by the user.

In the following tests, we used a scalable implementation with Edward (Tran et al., 2016) detailed in Supplementary section E. Figure 14 shows the performance of this imple-
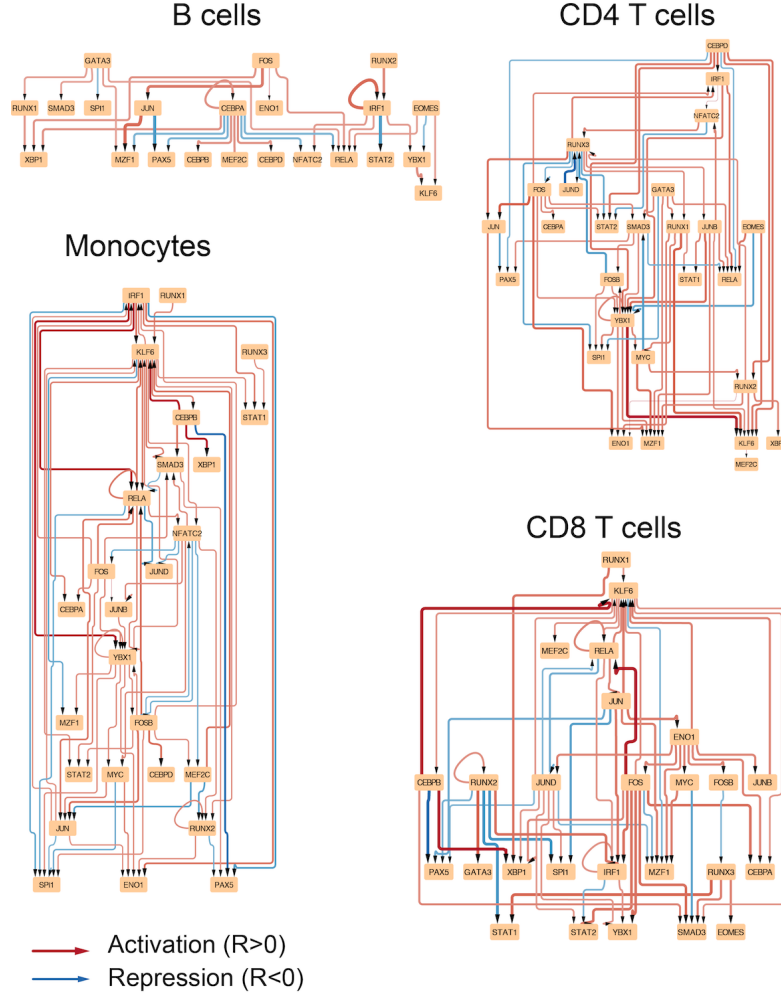
*Figure 10.* Subnetworks of cell type-specific GRNs between TFs with strongest regulations ($|R_k| > 4.5$); global structures show differences in connectivity of TFs across clusters (cell types); red and blue edges indicate activation and repression; edge widths are proportional to strength of regulation ($|R|$).

mentation on a larger number of cells and genes.

**Cell type characterization.** In this test, we pre-imputed and normalized data, and fix BISCUIT-derived normalization parameters ($\alpha, \beta$). Specifically, we normalized and imputed $\boldsymbol{x}_j$ for each cell $j$ by transforming it to $\boldsymbol{y}_j$ with $\boldsymbol{y}_j = A\boldsymbol{x}_j + b$, and setting $A = I/\beta_j$ and $b = (I - \alpha_j A)\boldsymbol{\mu}_k$ with BISCUIT -inferred parameters. This transformation corrects for cell-specific technical effects captured by $\alpha_j.\beta_j$, as $\boldsymbol{y}_j \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)$ while $\boldsymbol{x}_j \sim \mathcal{N}(\alpha_j\boldsymbol{\mu}_k, \beta_j\Sigma_k)$ (Figure 4) (Prabhakaran et al., 2016). Figure 7 (a) shows t-SNE projections (Maaten & Hinton, 2008) of $\boldsymbol{y}_j$s after normalizing expression data based on inferred parameters $\boldsymbol{\mu}_k, \Sigma_k, \boldsymbol{\alpha}, \boldsymbol{\beta}$ where cells are colored by cluster assignments $z_j$s for $K = 5$ distinct cell populations. Figure 7 (b) shows normalized expression of known marker genes used to characterize the clusters as monocyte, B, T and NK cell types.

**Deconvolving epigenetic data.** Figure 7 (c,d) show inferred peak heights $\boldsymbol{p}_k$ for all clusters using *Symphony* compared to ground truth cell type specific peak heights (measured with ATAC-seq from sorted cell types). Supplementary Figure 13 shows an example genomic region with differential peaks for three of these cell types showing distinct epigenetic profiles. Projection of peak heights to principal components of ground truth peaks (excluding peaks un-constrained by expression data and peaks which show 0 accessibility in some cell types, to show performance at deciphering magnitudes) shows superior performance in deconvolving all subsets of cells except for the smallest population (NK cells, $< 5\%$ of cells). The small error between estimated $\boldsymbol{p}_k$s and ground truth cell type-specific peak heights confirms that our model is a good fit for the biological mechanism of regulation. We also evaluated the deconvolution of epigenetic data using *deconf* (Figure 7 (c))
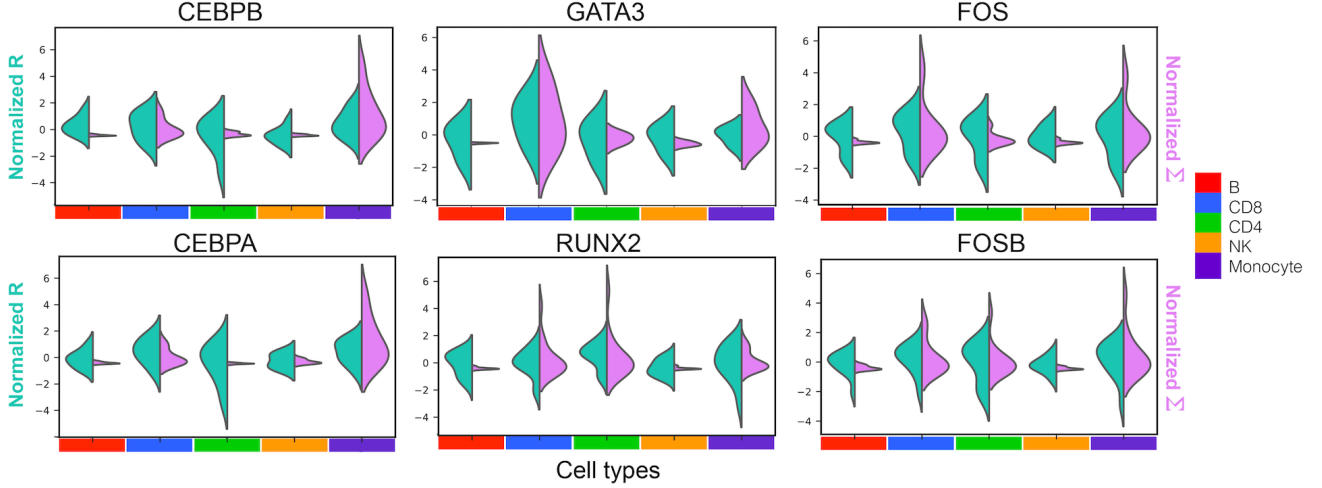
*Figure 11.* Identification of cell type-specific regulators. Violinplots of inferred cell type-specific regulatory function $R$ (green) for TFs that show variability across cell types, compared to inferred covariance $\Sigma$ (pink) between example TFs (regulators) and all of their target genes; positive and negative signs in $R$ denote activation and repression of expression, respectively; values have been z-score normalized to display both variables on the same y-axis.

which shows under-estimation or inaccurate estimation of peak heights. We could not test *Dsection* as the number of clusters exceeds the number of replicates.

**Inferring GRNs.** The main advantage of *Symphony* is the inference of GRNs. Figure 8 shows *Symphony* can successfully learn cell type-specific covariances comparable to empirical covariances. The gene-gene covariance matrix is then explained by direct regulation as well as the propagated impact of regulation through $(R + R^T)^2$, which we previously validated is capable of capturing epigenetic information through the accuracy of its prior $p$. Furthermore, Figure 9 reveals that inferred regulatory functions are explained by either TF binding strength (peak height) or TF-gene covariance or both.

Figure 10 shows the inferred GRNs between TFs with strongest inferred links ($|R_k| > 4.5$) in each cluster. The differences in the structure of the networks suggests different mechanisms driving cell type-specific expression.

We observe numerous regulatory interactions that are variable across clusters. Figure 11 shows examples of TF-gene interactions ($R$) that are also supported by known literature. It can be seen that regulatory functions are partially supported by gene-gene covariances ($\Sigma$). We observe $CEBPA, CEBPB$ differentially regulating target genes in monocytes (cluster 5). A recent study (Jaitin et al., 2016) has shown that knock-outs of $CEBPB$ block monocyte differentiation.

We observe regulatory edges in the GRN for T cells between Transcription Factors (TFs) $GATA3, RUNX2$ and

their target genes, with minimal interaction in B cells and NK cells (Figure 11), and indeed these TFs are known to be associated with activation of cytotoxic T cells (Pearce et al., 2003) and CD8 T cell development (Woolf et al., 2003). We also observe cases such as $FOS, FOSB$ with different regulatory functions despite belonging to the same TF family, showing an example of how expression-derived information can further distinguish genetic interactions which cannot be immediately deciphered from epigenetic data.

## 6. Conclusion

We present a hierarchical Bayesian mixture model named *Symphony* that infers clusters of cells representative of cell types and gene regulatory networks (GRNs) specific to cell types. This is done by modeling the regulatory mechanism driving gene expression in each cell type, and assuming two observations as two views from the system: epigenetic measurements, which are informative of network edges and single-cell gene expression data, informative of network node activity. To the best of our knowledge, this is the first computational method that integrates single-cell expression data with epigenetic data. We provide theoretical justifications for the model and an EM-VI procedure. *Symphony* shows great performance in clustering cells, deconvolving epigenetic profiles and inferring GRNs in both synthetic and real data from peripheral blood cells and shows superiority to other methods that only address one of these problems. Further, *Symphony* was able to deconvolve epigenetic data when only one replicate was available through integration of expression data, a potentially common task which would

be challenging for any source separation technique. Future iterations of the experiments will allow a *Symphony*-derived clustering of PBMCs to improve the mapping of cells to cell types, particularly for more similar cell types such as CD4+ and CD8+ T cells and when few genes are considered. Applied to the growing single-cell datasets, Symphony can reveal cell type-specific regulation in normal cells as well as disrupted regulation in cancerous cells.

# References

Aibar, S., González-Blas, C. B., Moerman, T., Imrichova, H., Hulselmans, G., Rambow, F., Marine, J.-C., Geurts, P., Aerts, J., van den Oord, J., et al. Scenic: single-cell regulatory network inference and clustering. *Nature methods*, 14(11):1083, 2017.

Aitchison, J. and Brown, J. A. The lognormal distribution with special reference to its uses in economics. 1957.

Azizi, E., Airoldi, E., and Galagan, J. Learning modular structures from network data and node variables. In *International conference on machine learning*, pp. 1440–1448, 2014.

Azizi, E., Carr, A. J., Plitas, G., Cornish, A. E., Konopacki, C., Prabhakaran, S., Nainys, J., Wu, K., Kiseliovas, V., Setty, M., et al. Single-cell map of diverse immune phenotypes in the breast tumor microenvironment. *bioRxiv*, pp. 221994, 2018.

Biemann, C. Ontology learning from text: A survey of methods. In *LDV forum*, volume 20, pp. 75–93, 2005.

Bishop, C. M. *Pattern Recognition and Machine Learning*. 2006.

Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.

Brown, C. D., Mangravite, L. M., and Engelhardt, B. E. Integrative modeling of eqtls and cis-regulatory elements suggests mechanisms underlying cell type specificity of eqtls. *PLoS genetics*, 9(8):e1003649, 2013.

Buenrostro, J. D., Wu, B., Chang, H. Y., and Greenleaf, W. J. Atac-seq: A method for assaying chromatin accessibility genome-wide. *Current protocols in molecular biology*, pp. 21–29, 2015a.

Buenrostro, J. D., Wu, B., Litzenburger, U. M., Ruff, D., Gonzales, M. L., Snyder, M. P., Chang, H. Y., and Greenleaf, W. J. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature*, 523(7561): 486–490, 2015b.

Carpenter, B., Gelman, A., Hoffman, M., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M. A., Guo, J., Li, P., and Riddell, A. Stan: A probabilistic programming language. *Journal of Statistical Software*, 20:1–37, 2016.

Corces, M. R., Buenrostro, J. D., Wu, B., Greenside, P. G., Chan, S. M., Koenig, J. L., Snyder, M. P., Pritchard, J. K., Kundaje, A., Greenleaf, W. J., et al. Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. *Nature genetics*, 48(10):1193, 2016.

Davidson, E. H. *The regulatory genome: gene regulatory networks in development and evolution*. Elsevier, 2010.

Erkkilä, T., Lehmusvaara, S., Ruusuvuori, P., Visakorpi, T., Shmulevich, I., and Lähdesmäki, H. Probabilistic analysis of gene expression measurements from heterogeneous tissues. *Bioinformatics*, 26(20):2571–2577, 2010.

Ghahramani, Z. and Beal, M. J. Propagation algorithms for variational bayesian learning. In *Advances in neural information processing systems*, pp. 507–513, 2001.

Grant, C. E., Bailey, T. L., and Noble, W. S. Fimo: scanning for occurrences of a given motif. *Bioinformatics*, 27(7): 1017–1018, 2011.

Guo, J., Grow, E. J., Yi, C., Mlcochova, H., Maher, G. J., Lindskog, C., Murphy, P. J., Wike, C. L., Carrell, D. T., Goriely, A., et al. Chromatin and single-cell rna-seq profiling reveal dynamic signaling and metabolic transitions during human spermatogonial stem cell development. *Cell stem cell*, 21(4):533–546, 2017.

Hardoon, D. R., Szedmak, S., and Shawe-Taylor, J. Canonical correlation analysis: An overview with application to learning methods. *Neural computation*, 16(12):2639–2664, 2004.

Hashimshony, T., Wagner, F., Sher, N., and Yanai, I. Cel-seq: single-cell rna-seq by multiplexed linear amplification. *Cell reports*, 2(3):666–673, 2012.

Hecker, M., Lambeck, S., Toepfer, S., Van Someren, E., and Guthke, R. Gene regulatory network inference: data integration in dynamic modelsa review. *Biosystems*, 96 (1):86–103, 2009.

Houseman, E. A., Kile, M. L., Christiani, D. C., Ince, T. A., Kelsey, K. T., and Marsit, C. J. Reference-free deconvolution of dna methylation data and mediation by cell composition effects. *BMC bioinformatics*, 17(1):259, 2016.

Jaitin, D. A., Kenigsberg, E., Keren-Shaul, H., Elefant, N., Paul, F., Zaretsky, I., Mildner, A., Cohen, N., Jung, S., Tanay, A., et al. Massively parallel single-cell rna-seq

for marker-free decomposition of tissues into cell types. *Science*, 343(6172):776–779, 2014.

Jaitin, D. A., Weiner, A., Yofe, I., Lara-Astiaso, D., Keren-Shaul, H., David, E., Salame, T. M., Tanay, A., van Oudenaarden, A., and Amit, I. Dissecting immune circuits by linking crispr-pooled screens with single-cell rna-seq. *Cell*, 167(7):1883–1896, 2016.

Jones, M. and Viola, P. Fast multi-view face detection. *Mitsubishi Electric Research Lab TR-20003-96*, 3:14, 2003.

Klein, A. M., Mazutis, L., Akartuna, I., Tallapragada, N., Veres, A., Li, V., Peshkin, L., Weitz, D. A., and Kirschner, M. W. Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell*, 161(5):1187–1201, 2015.

Kreeger, P. K. and Lauffenburger, D. A. Cancer systems biology: a network modeling perspective. *Carcinogenesis*, 31(1):2–8, 2009.

Kshirsagar, A. M. Bartlett decomposition and wishart distribution. *Ann. Math. Statist.*, 30(1):239–241, 03 1959. doi: 10.1214/aoms/1177706379. URL https://doi.org/10.1214/aoms/1177706379.

Lake, B., Cheng, S., Sos, B., Fan, J., Yung, Y., Kaeser, G., Duong, T., Gao, D., Chun, J., Kharchenko, P., et al. Integrative single-cell analysis by transcriptional and epigenetic states in human adult brain. *bioRxiv*, pp. 128520, 2017.

Landauer, T. K., Laham, D., Rehder, B., and Schreiner, M. E. How well can passage meaning be derived without using word order? a comparison of latent semantic analysis and humans. In *Proceedings of the 19th annual meeting of the Cognitive Science Society*, pp. 412–417, 1997.

Levine, J. H., Simonds, E. F., Bendall, S. C., Davis, K. L., El-ad, D. A., Tadmor, M. D., Litvin, O., Fienberg, H. G., Jager, A., Zunder, E. R., et al. Data-driven phenotypic dissection of aml reveals progenitor-like cells that correlate with prognosis. *Cell*, 162(1):184–197, 2015.

Li, S. Z., Zhu, L., Zhang, Z., Blake, A., Zhang, H., and Shum, H. Statistical learning of multi-view face detection. In *European Conference on Computer Vision*, pp. 67–81. Springer, 2002.

Maaten, L. v. d. and Hinton, G. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.

Macosko, E. Z., Basu, A., Satija, R., Nemesh, J., Shekhar, K., Goldman, M., Tirosh, I., Bialas, A. R., Kamitaki, N., Martersteck, E. M., et al. Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell*, 161(5):1202–1214, 2015.

Montgomery, H. L. Harmonic analysis as found in analytic number theory. In *Twentieth Century Harmonic AnalysisA Celebration*, pp. 271–293. Springer, 2001.

Ng, A. Y. et al. On spectral clustering: Analysis and an algorithm. 2002.

Pan, S. J., Kwok, J. T., Yang, Q., and Pan, J. J. Adaptive localization in a dynamic wifi environment through multi-view learning. In *AAAI*, pp. 1108–1113, 2007.

Pearce, E. L., Mullen, A. C., Martins, G. A., Krawczyk, C. M., Hutchins, A. S., Zediak, V. P., Banica, M., DiCioccio, C. B., Gross, D. A., Mao, C.-a., et al. Control of effector cd8+ t cell function by the transcription factor eomesodermin. *Science*, 302(5647):1041–1043, 2003.

Pe'er, D. and Hacohen, N. Principles and strategies for developing network models in cancer. *Cell*, 144(6):864–873, 2011.

Petersen, K. B. et al. The matrix cookbook.

Prabhakaran, S., Azizi, E., Carr, A., and Peer, D. Dirichlet process mixture model for correcting technical variation in single-cell gene expression data. In *International Conference on Machine Learning*, pp. 1070–1079, 2016.

Recchia, G. and Jones, M. N. More data trumps smarter algorithms: Comparing pointwise mutual information with latent semantic analysis. *Behavior research methods*, 41(3):647–656, 2009.

Repsilber, D., Kern, S., Telaar, A., Walzl, G., Black, G. F., Selbig, J., Parida, S. K., Kaufmann, S. H., and Jacobsen, M. Biomarker discovery in heterogeneous tissue samples-taking the in-silico deconfounding approach. *BMC bioinformatics*, 11(1):27, 2010.

Rey, M. and Roth, V. Copula mixture model for dependency-seeking clustering. *arXiv preprint arXiv:1206.6433*, 2012.

Ritchie, M. D., Holzinger, E. R., Li, R., Pendergrass, S. A., and Kim, D. Methods of integrating data to uncover genotype–phenotype interactions. *Nature Reviews Genetics*, 16(2):85, 2015.

Rotem, A., Ram, O., Shoresh, N., Sperling, R. A., Goren, A., Weitz, D. A., and Bernstein, B. E. Single-cell chip-seq reveals cell subpopulations defined by chromatin state. *Nature biotechnology*, 33(11):1165–1172, 2015.

Satija, R., Farrell, J. A., Gennert, D., Schier, A. F., and Regev, A. Spatial reconstruction of single-cell gene expression data. *Nature biotechnology*, 33(5):495–502, 2015.

Shalek, A. K., Satija, R., Adiconis, X., Gertner, R. S., Gaublomme, J. T., Raychowdhury, R., Schwartz, S., Yosef, N., Malboeuf, C., Lu, D., et al. Single-cell transcriptomics reveals bimodality in expression and splicing in immune cells. *Nature*, 498(7453):236–240, 2013.

Tran, D., Kucukelbir, A., Dieng, A. B., Rudolph, M., Liang, D., and Blei, D. M. Edward: A library for probabilistic modeling, inference, and criticism. *arXiv preprint arXiv:1610.09787*, 2016.

Tran, D., Hoffman, M. D., Saurous, R. A., Brevdo, E., Murphy, K., and Blei, D. M. Deep probabilistic programming. In *International Conference on Learning Representations*, 2017.

Walker, R. Implementing discrete mathematics: Combinatorics and graph theory with mathematica, 1992.

Wang, H., Nie, F., and Huang, H. Multi-view clustering and feature learning via structured sparsity. In *International conference on machine learning*, pp. 352–360, 2013.

Woolf, E., Xiao, C., Fainaru, O., Lotem, J., Rosen, D., Negreanu, V., Bernstein, Y., Goldenberg, D., Brenner, O., Berke, G., et al. Runx3 and runx1 are required for cd8 t cell development during thymopoiesis. *Proceedings of the National Academy of Sciences*, 100(13):7731–7736, 2003.

Xu, C., Tao, D., and Xu, C. A survey on multi-view learning. *arXiv preprint arXiv:1304.5634*, 2013.

Zheng, G. X., Terry, J. M., Belgrader, P., Ryvkin, P., Bent, Z. W., Wilson, R., Ziraldo, S. B., Wheeler, T. D., McDermott, G. P., Zhu, J., et al. Massively parallel digital transcriptional profiling of single cells. *Nature communications*, 8:14049, 2017.

Zhu, J., Zhang, B., Smith, E. N., Drees, B., Brem, R. B., Kruglyak, L., Bumgarner, R. E., and Schadt, E. E. Integrating large-scale functional genomic data to dissect the complexity of yeast regulatory networks. *Nature genetics*, 40(7):854, 2008.

# Supplementary Materials for A Nonparametric Multi-view Model for Estimating Cell Type-Specific Gene Regulatory Networks
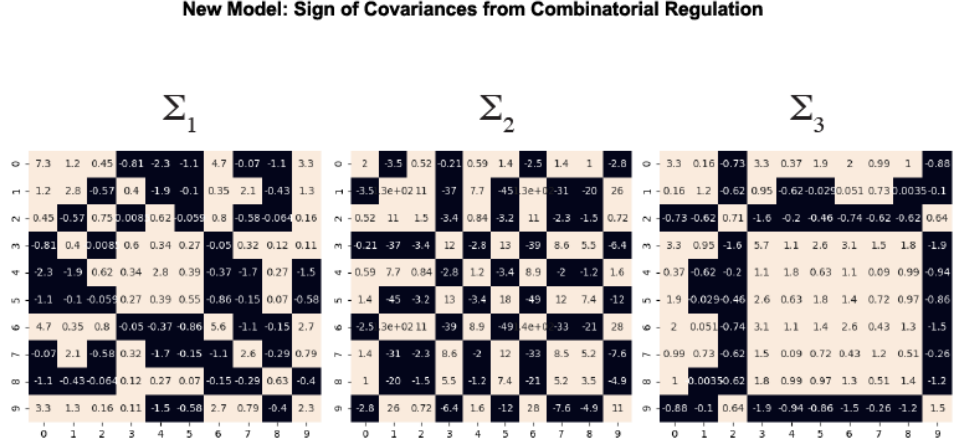
## A. Supplementary Figures



*Figure 12.* Examples of synthetic $\Sigma_k$ simulated from the same $S$ matrix for clusters $k = 1, 2, 3$, showing variability in sign of $\Sigma_k$ that can capture the impact of combinatorial regulations.

## B. Extended Theory

**Lemma 1** *The epigenetic model $f(\boldsymbol{c}|\boldsymbol{p}_k, \pi_k)$ is non-identifiable*
*Proof sketch.* Due to having $Kl + K$ unknowns in the mean parameters while having an $l$ dimensional Normal distribution, where $K$ is the maximum number of allowed clusters, we have an under-determined problem. Thus, we can provide multiple parameter sets $\pi_k$ and $\boldsymbol{p}_k$ leading to the same Normal distribution for $\boldsymbol{c}$.

**Lemma 3** *Square of a symmetric matrix $H$ gives a symmetric positive semi-definite matrix $L$.*
*Proof.* We show that there exists a symmetric positive semi-definite matrix $L$ iff there exists a symmetrix matrix $H$ that satisfies $H^2 = L$. $H, L \in \mathrm{R}^{d \times d}$. Orthogonal diagonalization of $L$ gives $QLQ^{-1} = D^2$ where $Q$ is the orthogonal matrix and $D$ is the square root diagonal matrix $diag(l_1^{\frac{1}{2}}, \cdots, l_d^{\frac{1}{2}})$. If there exists a matrix $H$ where $H = Q * diag(l_1^{\frac{1}{2}}, \cdots, l_d^{\frac{1}{2}}) * Q^{-1} = QDQ^{-1}$, then $QHQ^{-1} = D$. We now write:

$$QLQ^{-1} = D^2$$
$$= diag(l_1^{\frac{1}{2}}, \cdots, l_d^{\frac{1}{2}}) * diag(l_1^{\frac{1}{2}}, \cdots, l_d^{\frac{1}{2}})$$
$$= QHQ^{-1}QHQ^{-1} = QH^2Q^{-1}$$

showing $H^2 = L$. Next assume $H$ is symmetric and therefore all its eigenvalues are real. For some eigenvalue $h$ of $H$, $l$ is an eigenvalue of $L$ iff $l = h^2$ implying all eigenvalues of $L \in \mathrm{R}^*$ where $\mathrm{R}^* = \{0\} \cup \mathrm{R}^+$. Further, $L$ is symmetric given $H$ is symmetric. We now have $L$ as symmetric and has non-negative eigenvalues proving that $L$ is positive semi-definite. ■

**Lemma 4** *In the reduced model: $f(X^0, C|\boldsymbol{\beta}, \boldsymbol{\mu}_k^0, \Sigma_k^0, \boldsymbol{\alpha}, \boldsymbol{z}, \boldsymbol{p}_k, R_k^0, \pi_k, \zeta)$, $\beta s$ are identifiable under the conditions of: $\forall j : \boldsymbol{\mu}_k \geq \boldsymbol{\mu}' + diag(\Sigma')(\alpha_j - \nu)/\delta$ without the need for condition on $\beta s$.*

*Proof sketch.* Using Lemma 2 and focusing on the marginal distribution $f(X^0|\boldsymbol{\beta}, \boldsymbol{\mu}_k^0, \Sigma_k^0, \boldsymbol{\alpha}, \boldsymbol{z}, \pi_k)$ we know that $\beta_j \Sigma_k^0$ and $\pi_k$ are identified for all $j$ and $k$. Hence, $\beta_j^{1/2} \Sigma \pi_k (\Sigma_k^0)^{1/2}$ is identified. Using Lemma 3, $(R_k + R_k^T)^2$ is non-negative semi definite which allows us to define $(\Sigma_k^0)^{1/2}$ from the Wishart distribution.
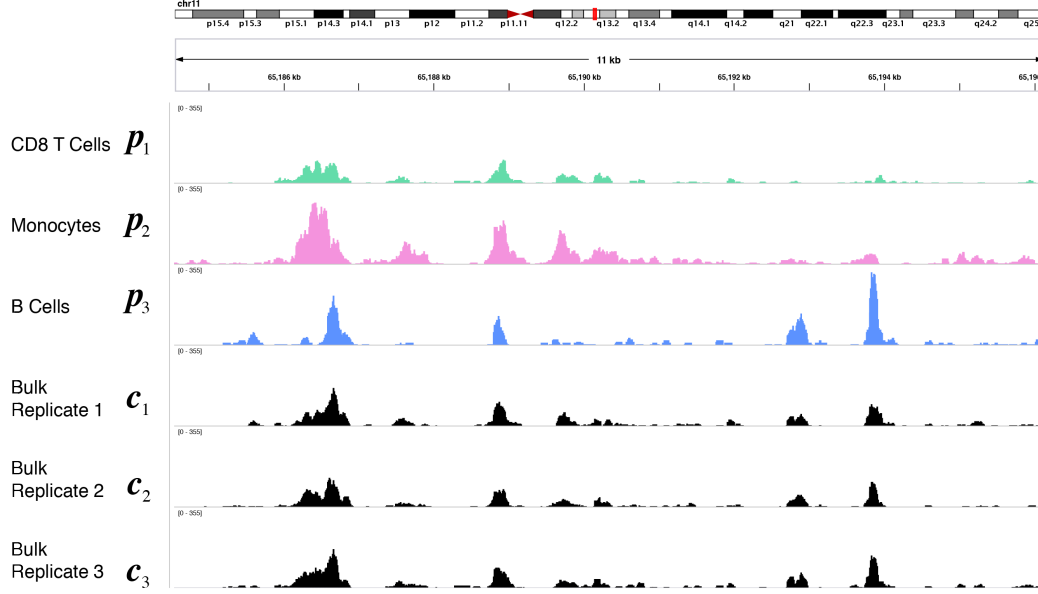
*Figure 13.* Epigenetic (ATAC-seq) data visualized for the three cell types used in section 7 showing an example region with differential peak heights $\boldsymbol{p}_k$s across the cell types and examples of simulated bulk data $\boldsymbol{c}_t$ for $t = 1, 2, 3$ from the weighted sum.

Therefore, $\beta_j^{1/2} \Sigma \pi_k (R_k + R_k^T)$ is identifiable. We also know that $\Sigma \pi_k \boldsymbol{p}_k$ is identified through the marginal distribution for $C$. Thus, putting the above together, it is not possible to have multiple values for $\beta_j$ since that would require the sum $\Sigma \pi_k (R_k + R_k^T)$ to be different for two sets of parameters while each element of this sum can be written as elements of $\Sigma \pi_k \boldsymbol{p}_k$ which is identified. ∎

**Lemma 5** *For a given $\boldsymbol{\beta} = \beta^*$, identifiability of: $f(X, C|\boldsymbol{\mu}_k, \Sigma_k, \boldsymbol{\alpha}, \boldsymbol{\beta} = \beta^*, \boldsymbol{z}, \boldsymbol{p}_k, R_k, \pi_k, \zeta)$ is guaranteed if $\forall j : \boldsymbol{\mu}_k \geq \boldsymbol{\mu}' + diag(\Sigma')(\alpha_j - \nu)/\delta$.*

*Proof sketch.* Parameters related to BISCUIT are identifiable using Lemma 2 without any of conditions on $\beta$s used in BISCUIT's proof, since $\beta$s are all given.

It remains to show identifiability for parameters $\boldsymbol{p}_k, R_k, \zeta$ from integrated model on $C$. $\zeta$ is identified due to Normality of $C$. Identifiability of $\boldsymbol{p}_k$ will lead to identifiability of $R_k$. We focus on proving identifiability of $\boldsymbol{p}_k$.

Using the identifiability of $\Sigma_k$ we have $d(d-1)/2$ equations based on generation of $\Sigma_k$ based on $\boldsymbol{p}_k$. Considering that $\boldsymbol{p}_k$ has $l$ unknowns and the relationships are at most polynomials of degree 2, as long as $d(d-1)/2 > l^2$ we have an overdetermined system of equations to identify $\boldsymbol{p}_k$ from $\Sigma_k$. ∎

**Theorem 6** *The full model $f(X, C|\Theta)$ where $\Theta := \{\boldsymbol{\mu}_k, \Sigma_k, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{z}, \boldsymbol{p}_k, R_k, \pi_k, \zeta\}$ is identifiable if $\forall j : \boldsymbol{\mu}_k \geq \boldsymbol{\mu}' + diag(\Sigma')(\alpha_j - \nu)/\delta$.*
*Proof sketch.* To prove identifiability of $\beta$, we use Lemma 5 on the following reduced distribution of $f(X, C|\Theta)$: $f(X^0, C|\boldsymbol{\beta}, \boldsymbol{\mu}_k^0, \Sigma_k^0, \boldsymbol{\alpha}, \boldsymbol{z}, \boldsymbol{p}_k, R_k^0, \pi_k, \zeta)$. Given the identified $\beta$s from Lemma 5, we use Lemma 6 to conclude identifiability of the rest of the parameters of the full model $f(X, C|\Theta)$, as desired. ∎

## C. Variational Inference update equation derivations

### C.1. Joint distribution for X and C

We use the graphical model in Figure 4 to write down the variational inference equations for Symphony. Note that this is constructed based on conditionally-conjugate priors for $\boldsymbol{\mu}$ and $\Sigma$ and on space discretised by $\boldsymbol{z} = \{z_1, \cdots, z_n\}$ where
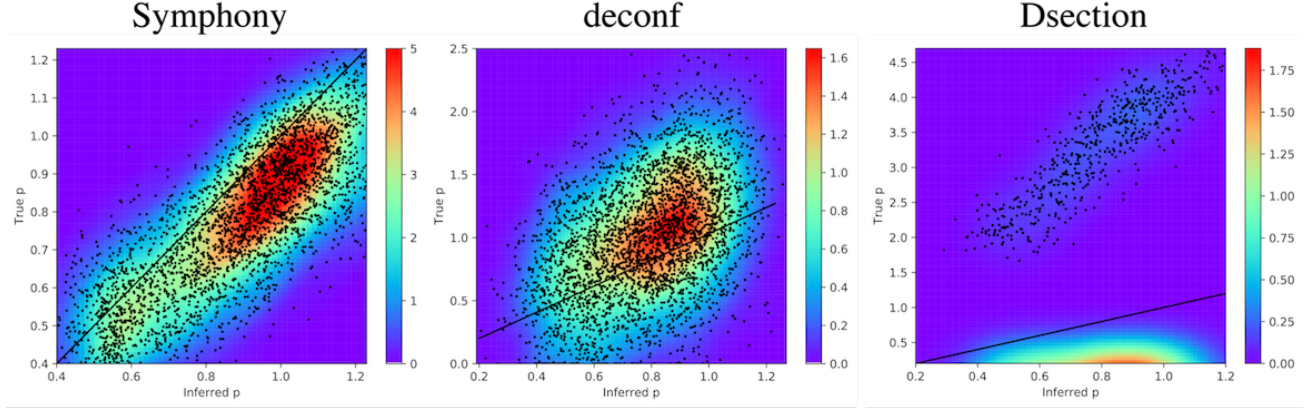
*Figure 14.* Performance in deconvolving epigenetic data with Edward implementation. Estimated peak heights $\boldsymbol{p}_k$ using *Symphony* for $n = 4000$ simulated cells in $K = 3$ clusters with $d = 100$ genes and $l = 550$ versus true peak heights, compared to two other deconvolution methods: deconf (Repsilber et al., 2010) and Dsection (Erkkilä et al., 2010); each dot represents a genomic region; heatmap shows density.

$z_j = \{j\}, j = [1, \cdots, k]$. The joint is written based on the Markov blanket for each parameter.

$$
\begin{aligned}
p(X, C, \boldsymbol{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda, \zeta) = {} & p(X|\boldsymbol{z}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\alpha}, \boldsymbol{\beta}) p(C|\boldsymbol{p}, \boldsymbol{\pi}, \zeta) \\
& p(\boldsymbol{z}|\boldsymbol{\pi}) p(\boldsymbol{\pi}|\boldsymbol{\varphi}) \\
& p(\boldsymbol{\mu}'|\boldsymbol{\mu}'', \Sigma'') p(\Sigma^{-1'}|\Sigma^{-1''}, d) \\
& p(\boldsymbol{\alpha}|\nu, \delta) p(\boldsymbol{\beta}|\omega, \theta) \\
& \prod_k p(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma') p(\Sigma_k^{-1}|(R_k + R_k^T)^2, \gamma) p(R_k|SM\boldsymbol{p_k}, \lambda) p(\boldsymbol{p}_k|\eta, \Lambda)
\end{aligned}
\tag{9}
$$

Next we expand each term in the RHS of Equation 9.

$$
\begin{aligned}
p(X|\boldsymbol{z}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\alpha}, \boldsymbol{\beta}) &= \prod_j p(\boldsymbol{x}_j|z_j, \boldsymbol{\mu}, \Sigma, \alpha_j, \beta_j) \\
&= \prod_j \prod_k p(\boldsymbol{x}_j|z_{jk}, \boldsymbol{\mu}_k, \Sigma_k, \alpha_j, \beta_j)^{z_{jk}} \\
&= \prod_j \prod_k \mathcal{N}(\boldsymbol{x}_j|\alpha_j \boldsymbol{\mu}_k, \beta_j \Sigma_k)^{z_{jk}}
\end{aligned}
\tag{10}
$$

$$
\begin{aligned}
p(\boldsymbol{z}|\boldsymbol{\pi}) &= \prod_j p(z_j|\boldsymbol{\pi}) \\
&= \prod_j \prod_k p(z_{jk}|\pi_k)^{z_{jk}} \\
&= \prod_j \prod_k \mathrm{Mult}(z_{jk}|\pi_k)^{z_{jk}} \\
&= \prod_j \prod_k \pi_k^{z_{jk}}
\end{aligned}
\tag{11}
$$

$$
\begin{aligned}
p(\boldsymbol{\pi}|\boldsymbol{\varphi}) &= \mathrm{Dir}(\boldsymbol{\pi}|\varphi_1, \cdots, \varphi_K) \\
&= \mathrm{Dir}(\boldsymbol{\pi}|\varphi_0) \\
&= \frac{1}{B(\varphi_0)} \prod_k \pi_k^{\varphi_0 - 1}
\end{aligned}
\tag{12}
$$

where $\sum_k \pi_k = 1$ and $\varphi_1 = \cdots = \varphi_K$ for symmetric prior. .

In the experiments, we used:

$$
\begin{aligned}
p(\boldsymbol{\pi}|\varphi) &= \mathrm{Stick}(\boldsymbol{\pi}|\varphi) \\
p(\pi'_k|\varphi) &= \mathrm{Beta}(\pi'_k|1, \varphi) \\
\pi_k &= \pi'_k \prod_i^{K-1} (1 - \pi'_i)
\end{aligned}
\tag{13}
$$

where $\sum_k \pi_k = 1$ and $\pi_k$ is the length of the *k-th* stick / proportion of the *k-th* cluster in stick breaking.

$$
\begin{aligned}
p(C|\boldsymbol{p}, \boldsymbol{\pi}, \zeta) &= \prod_t^r p(c_t|\boldsymbol{p}, \boldsymbol{\pi}, \zeta) \\
&= \prod_t^r \sum_k p(c_t|\boldsymbol{p}, \boldsymbol{\pi}, \zeta) \\
&= \prod_t^r \sum_k \mathcal{N}(c_t|\pi_k p_k, \zeta \mathrm{I})
\end{aligned}
\tag{14}
$$

$$
p(\boldsymbol{\mu}'|\boldsymbol{\mu}'', \Sigma'') = \mathcal{N}(\boldsymbol{\mu}'|\boldsymbol{\mu}'', \Sigma'')
\tag{15}
$$

$$
p(\Sigma'^{-1}|\Sigma''^{-1}, d) = \mathcal{W}(\Sigma'^{-1}|\Sigma''^{-1}, d)
\tag{16}
$$

$$
\begin{aligned}
p(\boldsymbol{\alpha}|\nu, \delta) &= \prod_j p(\alpha_j|\nu, \delta) \\
&= \prod_j \mathrm{logNormal}(\alpha_j|\nu, \delta)
\end{aligned}
\tag{17}
$$

$$
\begin{aligned}
p(\boldsymbol{\beta}|\omega, \theta) &= \prod_j p(\beta_j|\omega, \theta) \\
&= \prod_j \mathrm{logNormal}(\beta_j|\omega, \theta)
\end{aligned}
\tag{18}
$$

$$
\begin{aligned}
&\prod_k p(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma') p(\Sigma_k^{-1}|R_k, \gamma) p(R_k|SM\boldsymbol{p}_k, \lambda) p(\boldsymbol{p}_k|\eta, \Lambda) = \prod_k \mathcal{N}(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma') \mathcal{W}(\Sigma_k^{-1}|(R_k + R_k^T)^2, \gamma) \\
&\prod_i \prod_{i'} \mathcal{N}(R_k^{i,i'}|S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}, \lambda) \mathrm{trunc}\mathcal{N}(\boldsymbol{p_k}|\eta, \Lambda, \boldsymbol{0}, +\infty)
\end{aligned}
\tag{19}
$$

## C.2. Variational distributions

We now write the factorized distribution $q$ which will approximate the joint distribution in Equation 9. We follow Chapter 10 of Bishop (Bishop, 2006) and the Matrix cookbook Section 8.2 (Petersen et al.).

$$q(\boldsymbol{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda, \zeta | X, C) = \underbrace{q(\boldsymbol{z}|X,C)}_{Variational\ E-step}$$
$$\underbrace{q(\boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda, \zeta | X, C)}_{Variational\ M-step} \tag{20}$$

The sequential update equations can be written in terms of the E-step and M-step as follows:

**Variational E-step.** Take the expectation of the log of the joint distribution with respect to $\Theta := \{\boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda, \zeta\}$ i.e. all other parameters except $\boldsymbol{z}$. We have:

$$\begin{aligned}
\ln q^*(\boldsymbol{z}|X,C) &= \mathbb{E}_{\Theta}[\ln p(X, C, \boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda, \zeta )] + Const \\
&= \mathbb{E}_{\Theta}[\ln p(X|\boldsymbol{z}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\alpha}, \boldsymbol{\beta}) + p(C|\boldsymbol{p}, \boldsymbol{\pi}, \zeta) + \ln p(\boldsymbol{z}|\boldsymbol{\pi}) + \ln p(\boldsymbol{\pi}|\boldsymbol{\varphi}) + \\
&\quad \ln p(\boldsymbol{\mu}'|\boldsymbol{\mu}'', \Sigma'') + \ln p(\Sigma^{-1'}|\Sigma^{-1''}, d) + \\
&\quad \ln p(\boldsymbol{\alpha}|\nu, \delta) + \ln p(\boldsymbol{\beta}|\omega, \theta) + \\
&\quad \sum_k \left( \ln p(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma') + \ln p(\Sigma_k^{-1}|(R_k + R_k^T)^2, \gamma) + \ln p(R_k|SM\boldsymbol{p_k}, \lambda) + \ln p(\boldsymbol{p_k}|\eta, \Lambda) \right)] + Const
\end{aligned} \tag{21}$$

Taking terms in $\boldsymbol{z}$ alone:

$$\begin{aligned}
\ln q^*(\boldsymbol{z}|X,C) &= \mathbb{E}_{\Theta}[\ln p(X|\boldsymbol{z}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\alpha}, \boldsymbol{\beta}) + \ln p(\boldsymbol{z}|\boldsymbol{\pi})] + Const \\
&= \mathbb{E}_{\Theta}[\ln \prod_j \prod_k \mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k, \beta_j\Sigma_k)^{z_{jk}} + \ln \prod_j \prod_k \pi_k^{z_{jk}}] + Const \\
&= \mathbb{E}_{\Theta}[\sum_j \sum_k z_{jk} \ln \mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k, \beta_j\Sigma_k) + \sum_j \sum_k z_{jk} \ln \pi_k] + Const \\
&= \mathbb{E}_{\Theta}[\sum_j \sum_k z_{jk}[\ln \mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k, \beta_j\Sigma_k) + \ln \pi_k]] + Const \\
&= \mathbb{E}_{\Theta}[\sum_j \sum_k z_{jk}[(-\frac{d}{2}\ln 2\pi + \frac{1}{2}\ln |\beta_j\Sigma_k|^{-1} - \frac{1}{2}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)^T(\beta_j\Sigma_k)^{-1}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)) + \ln \pi_k]] + Const \\
&= \sum_j \sum_k z_{jk}[\mathbb{E}_{\Theta}(-\frac{d}{2}\ln 2\pi) + \mathbb{E}_{\Theta}(\frac{1}{2}\ln |\beta_j\Sigma_k|^{-1}) \\
&\quad - \mathbb{E}_{\Theta}(\frac{1}{2}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)^T(\beta_j\Sigma_k)^{-1}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)) + \mathbb{E}_{\Theta}(\ln \pi_k)] + Const \\
&= \sum_j \sum_k z_{jk} \ln \Delta_{jk} + Const \\
&\propto \sum_j \sum_k z_{jk} \ln \Delta_{jk}
\end{aligned} \tag{22}$$

Taking exponentials on both sides of Equation 22:

$$q^*(\boldsymbol{z}|X,C) \propto \prod_j \prod_k \Delta_{jk}^{z_{jk}} \tag{23}$$

where

$$\ln \Delta_{jk} := -\frac{d}{2}\ln 2\pi + \underbrace{\mathbb{E}_{\Theta}(\frac{1}{2}\ln|\beta_j\Sigma_k|^{-1})}_{S1} - \underbrace{\mathbb{E}_{\Theta}(\frac{1}{2}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)^T(\beta_j\Sigma_k)^{-1}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k))}_{S2} + \underbrace{\mathbb{E}_{\Theta}(\ln\pi_k)}_{S3} \tag{24}$$

Let us now expand the expectations given as S1, S2 and S3 in Equation 24.

1.

$$\begin{aligned}
S1 &:= \mathbb{E}_{\Theta}(\frac{1}{2}\ln|\beta_j\Sigma_k|^{-1}) \\
&= \frac{1}{2}\mathbb{E}_{(\beta_j,\Sigma_k)}(\ln|\beta_j\Sigma_k|^{-1}) \\
&= \frac{1}{2}\mathbb{E}_{(\beta_j,\Sigma_k)}(\ln\beta_j^{-d}|\Sigma_k^{-1}|) \\
&= \frac{1}{2}\mathbb{E}_{(\beta_j,\Sigma_k)}(\ln\beta_j^{-d} + \ln|\Sigma_k^{-1}|) \\
&= \frac{1}{2}(\mathbb{E}_{(\beta_j,\Sigma_k)}\ln\beta_j^{-d} + \mathbb{E}_{(\beta_j,\Sigma_k)}\ln|\Sigma_k^{-1}|) \\
&= \frac{1}{2}(-d\mathbb{E}_{\beta_j}\ln\beta_j + \mathbb{E}_{\Sigma_k}\ln|\Sigma_k|) \\
&= \frac{1}{2}\Big(-d\mathbb{E}_{\beta_j}\ln\beta_j + d\ln(2) + \ln|(\boldsymbol{R}_k + \boldsymbol{R}_k^t)^2| + \sum_{i=1}^{d}\psi\big(\frac{\gamma_k + 1 - i}{2}\big)\Big) \\
&\equiv \ln|\widetilde{\beta_j\Sigma_k}|^{-1}
\end{aligned} \tag{25}$$

where $\psi(.)$ is the digamma function and equals $\frac{d}{dx}\log\Gamma(x)$ and $\gamma_k := \gamma + n_k$. **Runtime complexity**$\sim \mathrm{O}(d^3)$

2.

$$\begin{aligned}
S_2 &:= \mathbb{E}_{\Theta}[\frac{1}{2}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)^T(\beta_j\Sigma_k)^{-1}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)] \\
&= \mathbb{E}_{(\Sigma_k,\beta_j,\alpha_j,\mu_k)}[\frac{1}{2}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)^T(\beta_j\Sigma_k)^{-1}(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)] \\
&= \frac{1}{2}\sum_i\sum_j(\beta\Sigma_{k_{ij}})^{-1}\mathbb{E}[(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)^T(\boldsymbol{x}_j - \alpha_j\boldsymbol{\mu}_k)] \quad \text{(linearity of expectation property)} \\
&= \frac{1}{2}\sum_i\sum_j(\beta\Sigma_{k_{ij}})^{-1}(\Sigma''^{-1} + (\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)^T(\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)) \quad \text{(via covariance formula)} \\
&= \frac{1}{2}\sum_i\sum_j(\beta\Sigma_{k_{ij}})^{-1}(\Sigma''^{-1}) + \frac{1}{2}\sum_i\sum_j(\beta\Sigma_{k_{ij}})^{-1}(\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)^T(\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k) \quad (\Sigma_k \text{ symmetricity}) \\
&= \frac{1}{2}\Big(trace(\Sigma_k^{-1}\beta^{-1}\Sigma''^{-1}) + (\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)^T(\beta\Sigma_k)^{-1}(\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)\Big)
\end{aligned} \tag{26}$$

**Runtime complexity**$\sim \mathrm{O}(d^3)$

3.

$$\begin{aligned}
S3 &:= \mathbb{E}_{\Theta}[\ln\pi_k] \\
&= \mathbb{E}_{\boldsymbol{\pi}}[\ln\pi_k] \\
&= \psi(\varphi_0) - \psi(\sum_k\varphi_k) \\
&\equiv \ln\widetilde{\pi}_k
\end{aligned} \tag{27}$$

where $\psi$ is the *digamma* function and $\psi(t) = \frac{d}{d_0}\ln(\Gamma(t)) = \frac{\Gamma'(t)}{\Gamma(t)}$. **Runtime complexity**$\sim \mathrm{O}(1)$

Therefore,

$$\ln \Delta_{jk} := -\frac{d}{2}\ln 2\pi + \ln|\widetilde{\beta_j \Sigma_k}|^{-1} - \mathbb{E}_{\boldsymbol{\Theta}}[\frac{1}{2}(\boldsymbol{x}_j - \alpha_j \boldsymbol{\mu}_k)^T (\beta_j \Sigma_k)^{-1}(\boldsymbol{x}_j - \alpha_j \boldsymbol{\mu}_k)] + \ln \widetilde{\pi}_k \qquad (28)$$

We require that $q^*(\boldsymbol{z}|X)$ is normalised and that for every observation $j$, there is only one non-zero $z_{jk} \forall k \in \{1, \cdots, K\}$. Therefore it is sufficient to normalise each $\Delta_{jk}$ as

$$r_{jk} = \frac{\Delta_{jk}}{\sum_{h=1}^K \Delta_{jh}} \qquad (29)$$

$$q^*(\boldsymbol{z}|X) = \prod_j \prod_k r_{jk}^{z_{jk}}$$
$$= \prod_j q^*(z_j) \qquad (30)$$

The expectation for the discrete distribution $q^*(z_{jk})$ gives the responsibilities $r_{jk}$ for point $x_j$ with the current $k^{th}$ cluster's parameters :

$$\mathbb{E}_{q^*(z_{jk})}[z_{jk}] := \frac{\pi_k \mathcal{N}(\boldsymbol{x}_j | \alpha_j \boldsymbol{\mu}_k, \beta_j \Sigma_k)}{\sum_{m=1}^K \pi_m \mathcal{N}(\boldsymbol{x}_j | \alpha_j \boldsymbol{\mu}_m, \beta_j \Sigma_m)} = r_{jk} \qquad (31)$$

with the overall runtime complexity for calculating $r_{jk}$ is $\mathrm{O}(d^3)$.

**MAP estimate for z**   Since the E-step has a runtime $\sim \mathrm{O}(d^3)$ due to three matrix inversions, we substitute this as $z_{MAP}(\boldsymbol{x}) = \arg\max_z p(\boldsymbol{x}|\boldsymbol{z})p(\boldsymbol{z}|\boldsymbol{\pi})$ which has a runtime of $\sim \mathrm{O}(d^2)$ when $\Sigma_k^{-1}$s and $\Sigma''$ are apriori Cholesky decomposed.

**Variational M-step.**   Take the expectation of the log of the joint distribution with respect to $\boldsymbol{z}$. We have:

$$\begin{aligned}
&\ln q^*(\boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda, \zeta | X, C) \\
&\quad = \mathbb{E}_{\boldsymbol{z}}[\ln p(X, C, \boldsymbol{z}, \boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda, \zeta)] + C \\
&\quad = \mathbb{E}_{\boldsymbol{z}}[\ln p(X|\boldsymbol{z}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\alpha}, \boldsymbol{\beta}) + \ln p(C|\boldsymbol{p}, \boldsymbol{\pi}, \zeta) + \ln p(\boldsymbol{z}|\boldsymbol{\pi}) + \ln p(\boldsymbol{\pi}|\boldsymbol{\varphi}) + \\
&\qquad \ln p(\boldsymbol{\mu}'|\boldsymbol{\mu}'', \Sigma'') + \ln p(\Sigma^{-1'}|\Sigma^{-1''}, d) + \\
&\qquad \ln p(\boldsymbol{\alpha}|\nu, \delta) + \ln p(\boldsymbol{\beta}|\omega, \theta)] + \\
&\qquad \sum_k \Big(\ln p(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma') + \ln p(\Sigma_k^{-1}|\boldsymbol{R}_k, \gamma) + \ln p(\boldsymbol{R}_k|SM\boldsymbol{p_k}, \lambda) + \ln p(\boldsymbol{p}_k|\eta, \Lambda)\Big)] + Const
\end{aligned} \qquad (32)$$

Taking terms in $(\boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda, \zeta)$, we get:

$$
\begin{aligned}
\ln q^*&(\boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda | X, C) = \\
& \mathbb{E}_{\boldsymbol{z}}[\ln p(X|\boldsymbol{z}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\alpha}, \boldsymbol{\beta}) + \\
& \ln p(\boldsymbol{z}|\boldsymbol{\pi}) + \ln p(C|\boldsymbol{p}, \boldsymbol{\pi}, \zeta) + \ln p(\boldsymbol{\pi}|\boldsymbol{\varphi}) + \\
& \ln p(\boldsymbol{\mu}'|\boldsymbol{\mu}'', \Sigma'') + \ln p(\Sigma^{-1'}|\Sigma^{-1'''}, d) + \ln p(\boldsymbol{\alpha}|\nu, \delta) + \ln p(\boldsymbol{\beta}|\omega, \theta)] + \\
& \sum_k \Big( \ln p(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma') + \ln p(\Sigma_k^{-1}|(R_k + R_k^T)^2, \gamma) + \ln p(R_k|SM\boldsymbol{p_k}, \lambda) + \ln p(\boldsymbol{p}_k|\eta, \Lambda) \Big)] \\
& + C \\
=& \mathbb{E}_{\boldsymbol{z}}\Big[ \ln \prod_j \prod_k \mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k, \beta_j\Sigma_k)^{z_{jk}} + \ln \prod_j \prod_k \pi_k^{z_{jk}} \Big] + \ln \mathrm{Dir}(\boldsymbol{\pi}|\varphi_0) + \\
& \ln \prod_t^r \prod_k \mathcal{N}(c_t | \sum_k \pi_k\boldsymbol{p}_k, \zeta\mathrm{I}) + \ln \mathcal{N}(\boldsymbol{\mu}'|\boldsymbol{\mu}'', \Sigma'') + \ln \mathcal{W}(\Sigma'^{-1}|\Sigma''^{-1}, d) + \\
& \ln \prod_j \log\mathrm{Normal}(\alpha_j|\nu, \delta) + \ln \prod_j \log\mathrm{Normal}(\beta_j|\omega, \theta) + \\
& \ln \Big( \prod_k \mathcal{N}(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma')\mathrm{Wishart}(\Sigma_k^{-1}|R_k, \gamma)\mathcal{N}((R_k + R_k^T)^2|SM\boldsymbol{p}_k, \lambda)\mathrm{trunc}\mathcal{N}(\boldsymbol{p}_k|\eta, \Lambda) \Big) \\
& + Const \\
=& \mathbb{E}_{\boldsymbol{z}}\Big[ \sum_j \sum_k z_{jk} \ln \mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k, \beta_j\Sigma_k) + \sum_j \sum_k z_{jk} \ln \pi_k \Big] + \ln \mathrm{Dir}(\boldsymbol{\pi}|\varphi_0) + \\
& \sum_t^r \sum_k \ln \mathcal{N}(c_t | \sum_k \pi_k\boldsymbol{p}_k, \zeta\mathrm{I}) + \ln \mathcal{N}(\boldsymbol{\mu}'|\boldsymbol{\mu}'', \Sigma'') + \ln \mathcal{W}(\Sigma'^{-1}|\Sigma''^{-1}, d) + \\
& \sum_j \ln \log\mathrm{Normal}(\alpha_j|\nu, \delta) + \sum_j \ln \log\mathrm{Normal}(\beta_j|\omega, \theta) + \\
& \sum_k \ln \mathcal{N}(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma') + \sum_k \ln \mathrm{Wishart}(\Sigma_k^{-1}|(R_k + R_k^T)^2, \gamma) + \\
& \sum_k \ln \mathcal{N}(R_k|SM\boldsymbol{p}_k, \lambda) + \sum_k \ln \mathrm{trunc}\mathcal{N}(\boldsymbol{p}_k|\eta, \Lambda) + Const
\end{aligned}
\tag{33}
$$

We assume that the set of latent variables is independent of the rest of the latent variables given $X$ and $C$. This independence assumption reduces the problem complexity and allows us to get closed-form solutions in the M-step. This is called the *mean-field assumption.* We use this assumption and proceed to factor the latent variables into conditionally-independent components, to perform co-ordinate ascent mean field VI (CAVI) on each variational component:

$$
\begin{aligned}
\ln q^*&(\boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma, \boldsymbol{\mu}', \Sigma', R, S, M, \lambda, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{p}, \eta, \gamma, \Lambda, \zeta | X, C) \\
=& \ln q^*(\boldsymbol{\pi}) + \ln q^*(\boldsymbol{\mu}) + \ln q^*(\Sigma) + \ln q^*(\boldsymbol{\alpha}) + \ln q^*(\boldsymbol{\beta}) \\
& + \ln q^*(\boldsymbol{\mu}') + + \ln q^*(\Sigma') + \ln q^*(R) + \ln q^*(\boldsymbol{p}) \\
=& \ln q^*(\boldsymbol{\pi}) + \ln \prod_k q^*(\mu_k) + \ln \prod_k q^*(\Sigma_k^{-1}) + \ln \prod_j q^*(\alpha_j) + \ln \prod_j q^*(\beta_j) \\
& + \ln q^*(\boldsymbol{\mu}') + + \ln q^*(\Sigma') + \ln \prod_k q^*(R) + \ln \prod_k q^*(\boldsymbol{p}) \\
=& \ln q^*(\boldsymbol{\pi}) + \sum_k \ln q^*(\mu_k) + \sum_k \ln q^*(\Sigma_k^{-1}) + \sum_j \ln q^*(\alpha_j) + \sum_j \ln q^*(\beta_j) \\
& + \ln q^*(\boldsymbol{\mu}') + + \ln q^*(\Sigma') + \sum_k \ln q^*(R) + \sum_k \ln q^*(\boldsymbol{p})
\end{aligned}
\tag{34}
$$

Let us find the approximate distributions $q^*(\cdot)$ for every parameter in the RHS of Equation 34 by comparing to RHS of

Equation 33.

1.

$$q^*(\boldsymbol{\pi}) = p(\boldsymbol{\pi}|\varphi) = \text{Stick}(\boldsymbol{\pi}|\varphi)$$
$$p(\pi_k'|\varphi) = \text{Beta}(\pi_k'|1, \varphi) \tag{35}$$
$$\pi_k = \pi_k' \prod_i^{K-1} (1 - \pi_i')$$

2.

$$\sum_k \ln q^*(\mu_k) = \mathbb{E}_{\boldsymbol{z}}\Big[\sum_j \sum_k z_{jk} \ln \mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k, \beta_j\Sigma_k)\Big] + \sum_k \ln \mathcal{N}(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma')$$

$$= \sum_k \Big(\mathbb{E}_{\boldsymbol{z}}\Big[\underbrace{\sum_j z_{jk} \ln \mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k, \beta_j\Sigma_k)\Big] + \ln \mathcal{N}(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma')\Big)}_{expanded\ below}$$

$$= \sum_j \mathbb{E}_{\boldsymbol{z}} z_{jk} \ln \mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k, \beta_j\Sigma_k) + \ln \mathcal{N}(\boldsymbol{\mu}_k|\boldsymbol{\mu}', \Sigma')$$

$$= -\frac{1}{2} \sum_j \mathbb{E}_{\boldsymbol{z}} z_{jk} (\boldsymbol{\mu}_k - \frac{\boldsymbol{x}_j}{\alpha_j})^T \Big(\frac{\beta_j}{\alpha_j^2}\Sigma_k\Big)^{-1} (\boldsymbol{\mu}_k - \frac{x_j}{\alpha_j})$$

$$-\frac{1}{2} \sum_j \mathbb{E}_{\boldsymbol{z}} z_{jk} \ln |\beta_j\Sigma_k| - \frac{d}{2} \sum_j \mathbb{E}_{\boldsymbol{z}} z_{jk} \ln(2\pi)$$

$$-\frac{1}{2} \ln |\Sigma'| - \frac{d}{2} \ln(2\pi) - \frac{1}{2} (\boldsymbol{\mu}_k - \boldsymbol{\mu}')^T \Sigma'^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}')$$

$$\propto -\frac{1}{2} \sum_j \mathbb{E}_{\boldsymbol{z}} [z_{jk} (\boldsymbol{\mu}_k - \frac{\boldsymbol{x}_j}{\alpha_j})^T \Big(\frac{\beta_j}{\alpha_j^2}\Sigma_k\Big)^{-1} (\boldsymbol{\mu}_k - \frac{\boldsymbol{x}_j}{\alpha_j})] - \frac{1}{2} (\boldsymbol{\mu}_k - \boldsymbol{\mu}')^T \Sigma'^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}') \tag{36}$$

(by taking terms in $\mu_k$ and $\mu'$)

$$\propto -\frac{1}{2} \sum_j \mathbb{E}_{\boldsymbol{z}} z_{jk} \mathbb{E}_{\boldsymbol{z}}[(\boldsymbol{\mu}_k - \frac{\boldsymbol{x}_j}{\alpha_j})^T \Big(\frac{\beta_j}{\alpha_j^2}\Sigma_k\Big)^{-1} (\boldsymbol{\mu}_k - \frac{\boldsymbol{x}_j}{\alpha_j})] - \frac{1}{2} (\boldsymbol{\mu}_k - \boldsymbol{\mu}')^T \Sigma'^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}')$$

$$\propto -\frac{1}{2} \sum_j \mathbb{E}_{\boldsymbol{z}} z_{jk} \Big(trace(\Sigma_k^{-1}(\frac{\beta_j}{\alpha_j^2})^{-1}\Sigma''^{-1}) + (\bar{\boldsymbol{\mu}}_k - \frac{\boldsymbol{x}_j}{\alpha_j})^T \Big(\frac{\beta_j}{\alpha_j^2}\Sigma_k\Big)^{-1} (\bar{\boldsymbol{\mu}}_k - \frac{\boldsymbol{x}_j}{\alpha_j})\Big)$$

$$-\frac{1}{2} (\boldsymbol{\mu}_k - \boldsymbol{\mu}')^T \Sigma'^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}')$$

$$\propto -\frac{1}{2} \sum_k \Big(\sum_j r_{jk} \Big(trace(\Sigma_k^{-1}(\frac{\beta_j}{\alpha_j^2})^{-1}\Sigma''^{-1}) + (\bar{\boldsymbol{\mu}}_k - \frac{\boldsymbol{x}_j}{\alpha_j})^T \Big(\frac{\beta_j}{\alpha_j^2}\Sigma_k\Big)^{-1} (\bar{\boldsymbol{\mu}}_k - \frac{\boldsymbol{x}_j}{\alpha_j})\Big)$$

$$+ (\boldsymbol{\mu}_k - \boldsymbol{\mu}')^T \Sigma'^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}')\Big)$$

3. Let us denote $R_k^* = (R_k + R_k^T)^2$. By construction of $R_k$, $R_k^*$ is a Gram matrix and a valid scale matrix for Wishart

parametrisation.

$$
\begin{aligned}
\sum_k \ln q^*(\Sigma_k^{-1}) &= \mathbb{E}_{\boldsymbol{z}}\Big[\sum_j\sum_k z_{jk}\ln\mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k,\beta_j\Sigma_k)\Big] + \sum_k\ln\mathcal{W}(\Sigma_k^{-1}|R_k^*,\lambda)\\
&= \sum_k\Big[\underbrace{\sum_j\mathbb{E}_{\boldsymbol{z}}z_{jk}\ln\mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k,\beta_j\Sigma_k) + \ln\mathcal{W}(\Sigma_k^{-1}|R_k^*,\lambda)}_{\text{expanded below}}\Big]\\
&= -\sum_j\mathbb{E}_{\boldsymbol{z}}z_{jk}\frac{d}{2}\ln 2\pi - \sum_j\mathbb{E}_{\boldsymbol{z}}z_{jk}\frac{1}{2}\ln\beta_j^d -\\
&\quad \sum_j\mathbb{E}_{\boldsymbol{z}}z_{jk}\frac{1}{2}\ln|\Sigma_k| - \sum_j\mathbb{E}_{\boldsymbol{z}}z_{jk}\frac{1}{2\beta_j}(\boldsymbol{x}_j-\alpha_j\boldsymbol{\mu}_k)^T(\Sigma_k)^{-1}(\boldsymbol{x}_j-\alpha_j\boldsymbol{\mu}_k)+\\
&\quad \ln|\Sigma_k^{-1}|^{\frac{\lambda-d-1}{2}} + \{-\frac{tr(R_k^{*-1}\Sigma_k^{-1})}{2}\} - \ln(2^{\lambda d/2}|R_k^{*-1}|^{-\lambda/2}\Gamma_d(\lambda/2))\\
&\propto -\sum_j\mathbb{E}_z z_{jk}(\tfrac{1}{2}\ln\beta_j^d + \tfrac{1}{2}\ln|\Sigma_k|)\\
&\quad -\sum_j\mathbb{E}_z z_{jk}\frac{1}{2\beta_j}\Big(trace(\Sigma_k^{-1}\Sigma''^{-1}) + (\boldsymbol{\mu}''-\alpha_j\boldsymbol{\mu}_k)^T(\Sigma_k)^{-1}(\boldsymbol{\mu}''-\alpha_j\boldsymbol{\mu}_k)\Big)\\
&\quad +\ln|\Sigma_k^{-1}|^{\frac{\lambda-d-1}{2}} + \{-\frac{tr(R_k^{*-1}\Sigma_k^{-1})}{2}\} - \ln(2^{\lambda d/2}|R_k^{*-1}|^{-\lambda/2}\Gamma_d(\lambda/2))
\end{aligned}
\tag{37}
$$

4. If $\alpha_j \sim \log\mathcal{N}(\nu,\delta)$, then $\alpha_j^* = \ln\alpha_j \sim \mathcal{N}(\nu,\delta)$ by properties of log Normal distribution (Aitchison & Brown, 1957).

$$
\begin{aligned}
\sum_j\ln q^*(\alpha_j) &= \mathbb{E}_{\boldsymbol{z}}\Big[\sum_j\sum_k z_{jk}\ln\mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k,\beta_j\Sigma_k)\Big] + \sum_j\ln\mathcal{N}(\alpha_j|\nu,\delta)\\
&= \mathbb{E}_{\boldsymbol{z}}\Big[\sum_j\sum_k z_{jk}\ln\mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k,\beta_j\Sigma_k)\Big] + \sum_j\mathcal{N}(\alpha_j^*|\nu,\delta)\\
&= \sum_j\Big(\underbrace{\Big[\sum_k\mathbb{E}_{\boldsymbol{z}}z_{jk}\ln\mathcal{N}(\boldsymbol{x}_j|\alpha_j\boldsymbol{\mu}_k,\beta_j\Sigma_k)\Big] + \mathcal{N}(\alpha_j^*|\nu,\delta)}_{\text{expanded below}}\Big)\\
&= -\frac{1}{2}\sum_k\mathbb{E}_{\boldsymbol{z}}z_{jk}\ln|\beta_j\Sigma_k| - \frac{d}{2}\sum_k\mathbb{E}_{\boldsymbol{z}}z_{jk}\ln(2\pi)-\\
&\quad \sum_k\mathbb{E}_{\boldsymbol{z}}z_{jk}\frac{1}{2}(\boldsymbol{x}_j-\alpha_j\boldsymbol{\mu}_k)^T(\beta_j\Sigma_k)^{-1}(\boldsymbol{x}_j-\alpha_j\boldsymbol{\mu}_k)+\\
&\quad \frac{1}{\sqrt{2\delta^2\pi}}\exp\Big(-\frac{(\alpha_j^*-\nu)\mathbb{I}_{1\times1}(\alpha_j^*-\nu)}{2\delta^2}\Big)\\
&\propto -\sum_k\mathbb{E}_{\boldsymbol{z}}z_{jk}\frac{1}{2}(\boldsymbol{x}_j-\alpha_j\boldsymbol{\mu}_k)^T(\beta_j\Sigma_k)^{-1}(\boldsymbol{x}_j-\alpha_j\boldsymbol{\mu}_k)+\\
&\quad \frac{1}{\sqrt{2\delta^2\pi}}\exp\Big(-\frac{(\alpha_j^*-\nu)\mathbb{I}_{1\times1}(\alpha_j^*-\nu)}{2\delta^2}\Big) \text{ (taking terms in } \alpha_j)\\
&\propto -\sum_k r_{jk}\frac{1}{2}(trace(\Sigma_k^{-1}\beta^{-1}\Sigma''^{-1}) + (\boldsymbol{\mu}''-\alpha_j\boldsymbol{\mu}_k)^T(\beta\Sigma_k)^{-1}(\boldsymbol{\mu}''-\alpha_j\boldsymbol{\mu}_k))+\\
&\quad \frac{\alpha_j}{\alpha_j\sqrt{2\delta^2\pi}}\exp\Big(-\frac{(\ln\alpha_j-\nu)\mathbb{I}_{1\times1}(\ln\alpha_j-\nu)}{2\delta^2}\Big)
\end{aligned}
\tag{38}
$$

(by replacing $\alpha_j^*$ with $\ln\alpha_j$ and making this a logNormal pdf )

5. Derivation is similar to that of Equation 38.

$$
\begin{aligned}
\sum_j \ln q^*(\beta_j) &= \mathbb{E}_{\mathbf{z}}\left[\sum_j \sum_k z_{jk} \ln \mathcal{N}(\mathbf{x}_j | \alpha_j \boldsymbol{\mu}_k, \beta_j \Sigma_k)\right] + \sum_j \ln \log \mathcal{N}(\beta_j | \omega, \theta) \\
&= \mathbb{E}_{\mathbf{z}}\left[\sum_j \sum_k z_{jk} \ln \mathcal{N}(\mathbf{x}_j | \alpha_j \boldsymbol{\mu}_k, \beta_j \Sigma_k)\right] + \sum_j \mathcal{N}(\beta_j^* | \omega, \theta) \\
&\propto -\sum_k r_{jk} \frac{1}{2}(trace(\Sigma_k^{-1} \beta^{-1} \Sigma''^{-1}) + (\boldsymbol{\mu}'' - \alpha_j \boldsymbol{\mu}_k)^T (\beta \Sigma_k)^{-1} (\boldsymbol{\mu}'' - \alpha_j \boldsymbol{\mu}_k)) + \\
&\quad \frac{\beta_j}{\beta_j \sqrt{2\theta^2 \pi}} \exp\left(-\frac{(\ln \beta_j - \omega)\mathbb{I}_{1\times 1}(\ln \beta_j - \omega)}{2\theta^2}\right)
\end{aligned}
\tag{39}
$$

(by replacing $\beta_j^*$ with $\ln \beta_j$ and making this a logNormal pdf )

6.

$$
\begin{aligned}
\ln q^*(\boldsymbol{\mu}') &= \ln \mathcal{N}(\boldsymbol{\mu}' | \boldsymbol{\mu}'', \Sigma'') + \sum_k \ln \mathcal{N}(\boldsymbol{\mu}_k | \boldsymbol{\mu}', \Sigma') \\
&\sim \ln \mathcal{N}(\boldsymbol{\mu}_{\mu'}, \Sigma_{\mu'}) \\
\boldsymbol{\mu}_{\mu'} &= \Sigma_{\mu'}(\Sigma''^{-1}\boldsymbol{\mu}'' + K^2 \Sigma'^{-1}\bar{\boldsymbol{\mu}}') \\
\Sigma_{\mu'} &= (\Sigma''^{-1} + K\Sigma'^{-1})^{-1}
\end{aligned}
\tag{40}
$$

7.

$$
\begin{aligned}
\ln q^*(\Sigma^{-1'}) &= \ln \mathcal{W}(\Sigma'^{-1} | \Sigma''^{-1}, d) + \sum_k \ln \mathcal{N}(\boldsymbol{\mu}_k | \boldsymbol{\mu}', \Sigma') \\
&\sim \ln \mathcal{W}(V_{\Sigma'^{-1}}, d_{\Sigma'^{-1}}) \\
V_{\Sigma'^{-1}} &= (d\Sigma'' + 2\Sigma_{rss})^{-1} \\
d_{\Sigma'^{-1}} &= d + K
\end{aligned}
\tag{41}
$$

8.

$$\sum_k \ln q^*(R_k) = \mathbb{E}_{(\pi_k, \boldsymbol{p}_k)}[\sum_t \sum_k \ln \mathcal{N}(c_t | \sum_k (\pi_k p_k), \zeta \mathbb{I})] + \sum_k \ln \mathcal{W}(\Sigma_k^{-1} | R_k^*, \gamma) +$$

$$\sum_k \sum_i \sum_{i'} \ln \mathcal{N}(R_k^{i,i'} | S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}, \lambda)$$

$$\ln q^*(R_k) = \mathbb{E}_{(\pi_k, \boldsymbol{p}_k)}[\sum_t^r \ln \mathcal{N}(c_t | \sum_k (\pi_k \boldsymbol{p}_k), \zeta \mathbb{I})] + \ln \mathcal{W}(\Sigma_k^{-1} | R_k^*, \gamma) +$$

$$\sum_i \sum_{i'} \ln \mathcal{N}(R_k^{i,i'} | S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}, \lambda)$$

$$= \sum_t^r \mathbb{E}_{(\pi_k, \boldsymbol{p}_k)}[(c_t - \sum_k (\pi_k \boldsymbol{p}_k))^T (\zeta \mathbb{I})^{-1} (c_t - \sum_k (\pi_k \boldsymbol{p}_k))] + \ln \mathcal{W}(\Sigma_k^{-1} | R_k^*, \gamma) +$$

$$\sum_i \sum_{i'} \ln \mathcal{N}(R_k^{i,i'} | S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}, \lambda) \tag{42}$$

$$\propto \sum_t^r trace(\zeta \mathbb{I})^{-1} + (\bar{c}_t - \sum_k (\pi_k \boldsymbol{p}_k))^T (\zeta \mathbb{I})^{-1} (\bar{c}_t - \sum_k (\pi_k \boldsymbol{p}_k))] +$$

$$\{-\frac{tr(R_k^{*-1} \Sigma_k^{-1})}{2}\} - \ln(2^{\lambda d/2} | R_k^{*-1} |^{-\lambda/2} \Gamma_d(\lambda/2))$$

$$- \frac{1}{2\delta^2} \sum_i \sum_{i'} \left( (R_k^{i,i'} - S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}) \mathbb{I}_{1\times 1} (R_k^{i,i'} - S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}) \right)$$

$$\propto \{-\frac{tr(R_k^{*-1} \Sigma_k^{-1})}{2}\} - \ln(2^{\lambda d/2} | R_k^{*-1} |^{-\lambda/2} \Gamma_d(\lambda/2))$$

$$- \frac{1}{2\delta^2} \sum_i \sum_{i'} \left( (R_k^{i,i'} - S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}) \mathbb{I}_{1\times 1} (R_k^{i,i'} - S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}) \right)$$

(taking terms only in $R_k$)

9.

$$\sum_k \ln q^*(\boldsymbol{p}_k) = \mathbb{E}_{(\pi_k, \boldsymbol{p}_k)}[\sum_t \sum_k \ln \mathcal{N}(c_t | \sum_k (\pi_k \boldsymbol{p}_k), \zeta \mathbb{I})] + \sum_k \ln \text{trunc}\mathcal{N}(\boldsymbol{p}_k | \eta, \Lambda) +$$

$$\sum_k \sum_i \sum_{i'} \ln \mathcal{N}(R_k^{i,i'} | S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}, \lambda)$$

$$\ln q^*(\boldsymbol{p}_k) = \mathbb{E}_{(\pi_k, \boldsymbol{p}_k)}[\sum_t^r \ln \mathcal{N}(c_t | \sum_k (\pi_k \boldsymbol{p}_k), \zeta \mathbb{I})] + \ln \text{trunc}\mathcal{N}(\boldsymbol{p}_k | \eta, \Lambda) +$$

$$\sum_i \sum_{i'} \ln \mathcal{N}(R_k^{i,i'} | S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}, \lambda)$$

$$= \sum_t^r \mathbb{E}_{(\pi_k, \boldsymbol{p}_k)}[(c_t - \sum_k (\pi_k \boldsymbol{p}_k))^T (\zeta \mathbb{I})^{-1} (c_t - \sum_k (\pi_k \boldsymbol{p}_k))] + \ln \left( \frac{1}{\sqrt{2\pi}\Lambda} \exp(-\frac{1}{2}(\frac{\boldsymbol{p}_k - \eta}{\Lambda})^2)) \right)$$

$$- \frac{1}{2} \ln \left( (1 + erf(\frac{\infty - \eta}{\Lambda \sqrt{2}})) - (1 + erf(\frac{-\eta}{\Lambda \sqrt{2}})) \right) + \sum_i \sum_{i'} \ln \mathcal{N}(R_k^{i,i'} | S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}, \lambda)$$

$$\propto \sum_t^r trace(\zeta \mathbb{I})^{-1} + (\bar{c}_t - \sum_k (\pi_k \boldsymbol{p}_k))^T (\zeta \mathbb{I})^{-1} (\bar{c}_t - \sum_k (\pi_k p_k))] + \ln \left( \frac{1}{\sqrt{2\pi}\Lambda} \exp(-\frac{1}{2}(\frac{\boldsymbol{p}_k - \eta}{\Lambda})^2)) \right)$$

$$- \frac{1}{2\delta^2} \sum_i \sum_{i'} \left( (R_k^{i,i'} - S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}) \mathbb{I}_{1\times 1} (R_k^{i,i'} - S^{i,i'} M^{i,i'} \boldsymbol{p}_k^{g(i,i')}) \right)$$

(taking terms only in $\boldsymbol{p}_k$)

$$\tag{43}$$

## D. Blueprint for Variational algorithm

1. **Perform** Variational E-step

   a. Compute $q^*(z_n) = \prod_k r_{jk}^{z_{jk}}$ where

$$r_{jk} \propto |\widetilde{\beta_j \Sigma_k}|^{-1} \exp(-S_2)\widetilde{\pi}_k, \tag{44}$$

$\sum_k r_{nk} = 1$ and $S_2 = \frac{1}{2}\Big(trace(\Sigma_k^{-1}\beta_j^{-1}\Sigma''^{-1}) + (\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)^T(\beta_j\Sigma_k)^{-1}(\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)\Big)$

2. **Perform** Variational M-step

   b. Compute $q^*(\pi_k) = \pi_k \sim$ Stick-breaking $\text{Beta}(1, \varphi)$

   c. Compute $q^*(\mu_k) = \exp\Big( -\frac{1}{2}\sum_j r_{jk}\Big(trace(\Sigma_k^{-1}(\frac{\beta_j}{\alpha_j^2})^{-1}\Sigma''^{-1}) + (\bar{\boldsymbol{\mu}}_{\boldsymbol{k}} - \frac{\boldsymbol{x}_j}{\alpha_j})^T(\frac{\beta_j}{\alpha_j^2}\Sigma_k)^{-1}(\bar{\boldsymbol{\mu}}_{\boldsymbol{k}} - \frac{\boldsymbol{x}_j}{\alpha_j}) + (\boldsymbol{\mu}_k - \boldsymbol{\mu}')^T\Sigma'^{-1}(\boldsymbol{\mu}_k - \boldsymbol{\mu}')\Big) + const\Big)$ from Equation 36.

   d. Compute $q^*(\Sigma_k^{-1}) = \exp\Big( -\frac{1}{2}\sum_j r_{jk}\Big( \ln \beta_j^d + \ln|\Sigma_k| + \frac{1}{\beta_j}\Big(trace(\Sigma_k^{-1}\Sigma''^{-1}) + (\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)^T(\Sigma_k)^{-1}(\boldsymbol{\mu}'' - \alpha_j\boldsymbol{\mu}_k)\Big)\Big) + \ln|\Sigma_k^{-1}|^{\frac{\lambda-d-1}{2}} + \{-\frac{tr(R_k^{*-1}\Sigma_k^{-1})}{2}\} - \ln(2^{\lambda d/2}|R_k^{*-1}|^{-\lambda/2}\Gamma_d(\lambda/2)) + const\Big)$ from Equation 37.

   e. Compute $q^*(\alpha_j) = \exp\Big( -\sum_k r_{jk}S_2 + \frac{1}{\sqrt{2\delta^2\pi}}\exp\Big( -\frac{(\ln\alpha_j - \nu)\mathbb{I}_{1\times 1}(\ln\alpha_j - \nu)}{2\delta^2}\Big) + const\Big)$ from Equation 38.

   f. Compute $q^*(\beta_j) = \exp\Big( -\sum_k r_{jk}S_2 + \frac{1}{\sqrt{2\theta^2\pi}}\exp\Big( -\frac{(\ln\beta_j - \omega)\mathbb{I}_{1\times 1}(\ln\beta_j - \omega)}{2\theta^2}\Big) + const\Big)$ from Equation 39.

   g. Compute $q^*(\boldsymbol{\mu}') \sim \mathcal{N}(\mu_{\mu'}, \Sigma_{\mu'})$ from Equation 40.

   h. Compute $q^*(\Sigma^{-1'}) \sim \mathcal{W}(V_{\Sigma'^{-1}}, d_{\Sigma'^{-1}})$ from Equation 41.

   i. Compute $q^*(R_k) = \exp\Big( \{-\frac{tr(R_k^{*-1}\Sigma_k^{-1})}{2}\} - \ln(2^{\lambda d/2}|R_k^{*-1}|^{-\lambda/2}\Gamma_d(\lambda/2)) - \frac{1}{2\delta^2}\sum_i\sum_{i'}\Big((R_k^{i,i'} - S^{i,i'}M^{i,i'}\boldsymbol{p}_k^{g(i,i')})\mathbb{I}_{1\times 1}(R_k^{i,i'} - S^{i,i'}M^{i,i'}\boldsymbol{p}_k^{g(i,i')})\Big) + const\Big)$

   j. Compute $q^*(\boldsymbol{p}_k) = \exp\Big( (\bar{c}_t - \sum_k(\pi_k\boldsymbol{p}_k))^T(\zeta\mathbb{I})^{-1}(\bar{c}_t - \sum_k(\pi_k\boldsymbol{p}_k))] + \ln\Big(\frac{1}{\sqrt{2\pi\Lambda}}\exp(-\frac{1}{2}(\frac{\boldsymbol{p}_k - \eta}{\Lambda})^2))\Big) - \frac{1}{2\delta^2}\sum_i\sum_{i'}\Big((R_k^{i,i'} - S^{i,i'}M^{i,i'}\boldsymbol{p}_k^{g(i,i')})\mathbb{I}_{1\times 1}(R_k^{i,i'} - S^{i,i'}M^{i,i'}\boldsymbol{p}_k^{g(i,i')})\Big) + const\Big)$

## E. Scalable Implementation

Given the complexity of the model (refer to Supplementary section C), for a scalable implementation applicable to biological data containing thousands of cells, we used probabilistic programming languages with several useful approximations and implementation tricks.

### E.1. Stan

A first implementation of *Symphony* intended for smaller-scale datasets utilizes probabilistic programming language Stan. As Stan uses the No-U-Turn Sampler (NUTS) MCMC algorithm, it produces more accurate results asymptotically and is therefore preferred if computational resources allow. However, since Stan does not support inference for an infinite mixture model, we simply use a finite mixture model in the experiments presented herein with the Stan implementation. In addition, we implemented a version with an optional asymmetric, user-defined Dirichlet prior for fair comparison against methods which require prior knowledge of mixture proportions.

### E.2. Edward

To scale *Symphony* to larger datasets, we implemented the model in probabilistic programming language Edward (Tran et al., 2016). As Edward is built on a tensorflow back-end, it allows GPU acceleration for faster matrix computation. In addition, the use of variational algorithms allows for faster approximations of the posterior than those which can be obtained with MCMC, albeit with a trade-off in accuracy observed in our case. In order to improve the fit of the model to real data with the Edward implementation, we made a number of approximations below which improved the empirical performance on PBMC and other datasets.

### E.3. Approximations

Inference of gene-gene GRNs and covariance matrices are the main goals of *Symphony*, yet accurate inference of such large matrices involves a number of computational challenges. In particular, constraints on covariance matrices of a multivariate normal distribution are difficult to enforce in the optimization setting of variational inference. For instance, large sparse matrices may very easily become non-singular during optimization, leading to un-defined loss functions.

We use several techniques to solve this problem. For one, we define the Wishart distribution in Edward using the Bartlett Decomposition, rather than the built-in Wishart function of tensorflow, which allows us to more easily define variational parameters. Specifically, we replace the sampling of covariance matrices $\Sigma_k \sim Wishart$ with a generative model constructing only univariate chi-squared distributions $c$ and normal distributions $n$, which can be shown to produce a valid sample from the Wishart distribution (Kshirsagar, 1959). In this setting, we define variational distributions corresponding to the dummy variables $n$ and $c$, as opposed to defining a matrix variate distribution which, during the course of optimization, must fit all the constraints of valid covariance matrices. Initialization of these parameters to large values additionally avoids problems with singularity in most cases.

In addition to the use of the Bartlett Decomposition, the Edward version of *Symphony* replaces the standard Wishart with a scaled Wishart for added flexibility of the model in the variational inference case. The scaled Wishart necessitates addition of a latent parameter $\delta$, such that

$$\Sigma_k' \sim Wishart$$

$$\delta_i \stackrel{\text{iid}}{\sim} Normal$$

$$\Sigma_k = \Delta\Sigma_k'\Delta, \text{ where } diag(\Delta) = \delta$$

Addition of the normal distribution above to the generative process infuses flexibility to the Wishart, whose variance is usually defined by a single degrees of freedom parameter. In addition, we allow separate inference of the diagonal and off-diagonal of the covariance matrices. This is a desirable property for *Symphony*, in that the model of gene regulation does not necessarily capture the diagonal of the covariance matrices representing variances of gene expression. Likewise, we solve additional issues caused by matrix inversion by simply replacing the prior on $\Sigma_k$ with a Wishart instead of Inverse-Wishart distribution. We note that, while this choice is not conjugate, this is valid as both distributions satisfy the requirements for priors on the covariance matrix.

We require a variational EM procedure with Edward, such that the cluster assignments $z$ are updated every several iterations with a maximization step. In particular, we choose $z_i$ for each cell based on the maximum likelihood of cluster assignment in the Gaussian mixture. This prevents the need for discrete optimization over categorical variational parameters. The performance of the variational EM algorithm is maximized when a good initial value for clustering is chosen. In this work, we initialized clusters based on cell-cell kNN graphs with Phenograph (Levine et al., 2015).

Finally, we made some other small distributional changes which seemed to produce better results in our experiments with this particular implementation. We replace the multivariate prior on $p$ with a univariate prior centered at a constant. In addition, we treated binary $M$ as a latent variable with a very tight variance. This was to add additional flexibility to the model, and to also assist with singularity issues by providing a dense matrix mean to $R$. We note that all fitted values of $M$ were similar within a small $\epsilon$ to their previously fixed values of either 0 or 1.