
SBGGeHPF: Stochastic Block Graph Generative Hierarchical Poisson Factorization

Anonymous Authors¹

Abstract

Gene regulatory network (GRN) models provide insight into mechanisms underlying cellular function. While previous methods have attempted to infer GRNs from single cell RNA sequencing (scRNA-seq) data, they are limited in interpretability, and do not explain how GRNs impact cell state transitions and plasticity. We developed Stochastic Block Graph Generative Hierarchical Poisson Factorization, or SBGGeHPF, to identify regulons, defined as communities of genes, driving cell plasticity. SBGGeHPF combines tasks of community detection, graph structure learning, matrix factorization, and low rank estimation to learn an interpretable, joint cell and gene latent space from scRNA-seq and GRN data. We applied SBGGeHPF to a simulated dataset and real tumor scRNA-seq data. SBGGeHPF faithfully reconstructs expression matrices, refines GRNs as densely connected communities, and successfully associates them to heterogeneous cell populations.

1. Introduction

Single-cell RNA sequencing (scRNA-seq) has enabled the characterization of cell heterogeneity and plasticity at high resolution in complex diseases such as cancer. scRNA-seq also provides unique opportunities to study the role of gene regulatory circuitry at a fine scale, as alteration of gene regulatory relationships are often driving forces behind cell function and fate. Gene regulatory networks (GRNs) are popular models for disentangling interactions between transcription factors (TFs) and target genes in dynamic biological systems, and an active field of research is reverse-engineering GRNs from scRNA-seq (Mercatelli, 2020). Identification of regulons, or jointly co-regulated genes, driving cell plasticity would be a boon for uncovering mechanisms involved in therapeutic response, such as cancer and immune cell

plasticity during/after immunotherapy in melanoma (Madhamshettiwar, 2012; Marusyk, 2012). However, biological and technical noise and sparsity from high drop-out rates cause poor performance for most GRN inference methods, which were originally designed for bulk RNA sequencing data. More importantly, current methods are not capable of pinpointing regulatory mechanisms that drive cell plasticity (Mercatelli, 2020; Iglesias-Martinez, 2021; Osorio, 2020; Keyl, 2023). In particular, identifying topological structures of GRNs (e.g., hubs, communities) associated with heterogeneous and altered cell states in disease can reveal novel drug targets for reversing abnormal cell plasticity. We propose a novel framework with joint Bayesian modeling of scRNA-seq and GRN data, to identify GRN structures and regulons explaining cell state transitions.

We present **Stochastic Block Graph Generative Hierarchical Poisson Factorization**, or SBGGeHPF, a Bayesian hierarchical model to achieve both goals of generating refined GRN structure using scRNA-seq and associating GRN topologies to cell sub-populations to identify interpretable regulons driving cell plasticity. Novelly, SBGGeHPF combines tasks of community detection, graph structure learning, and matrix factorization and low rank estimation. Matrix factorization allows the model to learn joint gene and cell factors. A stochastic block graph is used to leverage the joint factors in defining community structures of genes, interpretable as regulons. Densely-connected communities are then assigned to sub-populations of cells, for which upregulation of key regulons may drive cell plasticity. Finally, inter-community edges are penalized and intra-community edges are rewarded, thus preserving only meaningful community interactions and refining expected network topology. We apply SBGGeHPF to both a simulated scRNA-seq dataset with pre-defined ground-truth co-expression patterns as well as a clinical dataset from melanoma patients treated with immunotherapy. We show SBGGeHPF holds promise in disentangling GRN and cell fate dynamics from scRNA-seq.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

055 2. Methods

056 2.1. The SBGGeHPF Model

057 Suppose we have a scRNA-seq dataset with C cells and
 058 G genes, unique molecular identifier (UMI) counts X is
 059 a $C \times G$ matrix with integer entries. We also consider an
 060 input GRN $N = \{E, V\}$ representing prior knowledge on
 061 possible regulatory links, where E represents the set of
 062 genes and V represents edges connecting pairs of genes
 063 that regulate each other. We define ‘factors’ as sets of co-
 064 expressed genes associated with clusters of cells. If we have
 065 k factors, our objective is to learn a cell loading $W \in \mathbf{R}^{C \times k}$
 066 and a gene loading $H \in \mathbf{R}^{G \times k}$ where the low-rank count
 067 matrix $\hat{X} = WH^T$ resembles the input matrix X . We also
 068 aim to jointly learn a refined GRN $A_G \in \mathbf{R}^{G \times G}$ where the
 069 factors correspond to community structures of the graph and
 070 the weights indicate the probability of two genes regulating
 071 each other.

072 We regularize the model to preserve the topological structure
 073 of the input graph. To achieve this, we designed the
 074 following generative model inspired by the ideas of hierarchical
 075 Poisson factorization (Gopalan, 2015) and stochastic
 076 block models (Lee, 2019) (Fig. 1):

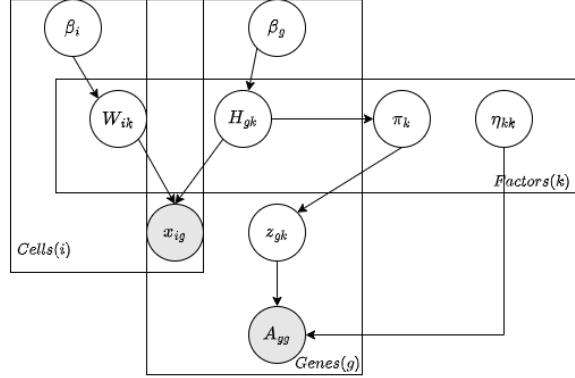
$$\begin{aligned}\beta_i &\sim \text{Gamma}(a', b'), \beta_g \sim \text{Gamma}(c', d') \\ W_{ik} &\sim \text{Gamma}(a, \beta_i), H_{gk} \sim \text{Gamma}(c, \beta_g) \\ x_{ig} | W_{ik}, H_{gk} &\sim \text{Poisson}(W_{ik}H_{gk}^T) \\ \pi_k &\sim \text{Dirichlet}(\mathbf{E}_g[H_{gk}]), z_{gk} \sim \text{Categorical}(\pi_k) \\ \eta_{kk} &\sim \text{Beta}(\alpha, \beta) \\ \Lambda &= z_{gk}\eta_{kk}z_{gk}^T, \Lambda \in \mathbf{R}^{G \times G} \\ A_G | \Lambda &\sim \text{Bernoulli}(\Lambda)\end{aligned}$$

077 $a', b', c', d', a, c, \alpha, \beta$ are hyperparameters. W_{ik} are cell
 078 weights across k factors based on cell budget β_i . H_{gk} are
 079 gene weights across k factors based on gene budget β_g . π
 080 is the factor membership assignment probability parameterized
 081 by a Dirichlet distribution with Gamma prior given by
 082 the mean gene weights in each factor. z is the per-gene factor
 083 membership sampled from a Categorical distribution. η is the
 084 block matrix sampled from Beta distributions where the
 085 diagonal entries represent the probability of within-factor
 086 connections, and the off-diagonal entries represent the
 087 probability of cross-factor connections, to capture cascades of
 088 TFs between communities. z and η jointly determine the
 089 edge probability between pairs of genes, and the generated
 090 GRN A_G is sampled from a Bernoulli distribution.

091 2.2. Hyperparameter Initialization and Regularization

092 We choose to initialize hyperparameters b' and d' to preserve
 093 the variance-to-mean ratio of total UMI counts per cell or gene in the sampling distributions of gene/cell budgets
 094 (Mendes Levitin, 2019). Specifically, we set b' and d' as:

$$b' = \frac{\text{Var}[\sum_g x_{ig}]}{\mathbf{E}[\sum_g x_{ig}]}, d' = \frac{\text{Var}[\sum_i x_{ig}]}{\mathbf{E}[\sum_i x_{ig}]}$$



095 **Figure 1.** Graphical representation of the model. Circles represent
 096 latent variables. Colored circles represent observations.

097 where x_{ig} are individual entries of count matrix X . To enforce sparsity, we initialize the Gamma shape parameters a', c' and a, c as 1.0 and 0.3, respectively. The Beta distribution parameters α, β for η are initialized to non-informative values 1.0 and 1.0. Since we expect the generated GRN to have densely connected communities and loose connections between communities, we learn a sparsity regularizer ρ and construct a mask $M \in \mathbf{R}^{k \times k}$. Assuming $\rho \sim \text{Beta}(\gamma, 5\gamma)$, the mask matrix is constructed as:

$$M_{ij} = \begin{cases} 1, & \text{if } i = j \\ \rho, & \text{if } i \neq j \end{cases}$$

098 Then, the regularized Bernoulli rate is given by $\Lambda = z(\eta \odot M)z^T$. This regularization method only allows strong signals of cross-community connection to be preserved, which prunes the input GRN and refines the graph’s topological structure according to observed phenotypic states.

100 2.3. Model Inference

101 Inference is conducted by Markov Chain Monte Carlo
 102 (MCMC) sampling. The traditional Metropolis-Hastings
 103 (MH) algorithm does not scale well for high dimensional
 104 distributions due to the random walk nature of its movement.
 105 Therefore, we perform parameter inference under the
 106 Hamiltonian Monte Carlo (HMC) framework, where
 107 gradient information of the target distribution is used to
 108 guide the sampler movement and make distant proposals
 109 with high acceptance probabilities. Our model contains both
 110 discrete and continuous latent variables, which renders traditional
 111 HMC-based algorithms such as No U-Turn Sampling
 112 (NUTS) ineffective. To overcome the mixed nature of the
 113 latent space, we apply Gibbs sampling at discrete sites and
 114 NUTS at continuous sites to sample from the target distribution.
 115 The models are trained with 4 chains with a warm-up
 116 distance of 100 and a sampling distance of 400. The sampling
 117 converged with an \hat{R} value less than 1.05, indicating
 118 that the chains are well-mixed.

110 2.4. Data Preparation

111 2.4.1. SIMULATED SCRNA-SEQ DATA

112 We simulated an scRNA-seq dataset with known co-
 113 expression patterns and a known set of gene factors as-
 114 sociated with cell grouping patterns, to test SBGGeHPF
 115 performance in learning and reconstructing interpretable
 116 joint cell and gene latent spaces. We used ESCO (Tian,
 117 2021) which simulates scRNA-seq by incorporating varia-
 118 tion in expression from cell heterogeneity (i.e., differentially
 119 expressed genes or DEGs), intrinsic variation in gene ex-
 120 pression between similar cell types, technical noise, and
 121 co-expression patterns using a Gaussian copula. We gen-
 122 erated data for 2000 cells and 200 genes with global DEG
 123 probability of 0.5 (for 100 DEGs total) and group-specific
 124 DEG probability of 0.3. Cells were divided into three groups
 125 with 60% of cells in group 1 and 20% in groups 2 and 3.
 126 To create SBGGeHPF input, we filtered the count matrix
 127 for the 100 DEGs and constructed a GRN by calculating
 128 empirical correlation between gene pairs for the 100 DEGs.

129 2.4.2. MELANOMA CLINICAL SCRNA-SEQ DATA

130 We also tested SBGGeHPF on a recently published clinical
 131 dataset consisting of scRNA-seq of tumor samples from
 132 melanoma patients treated with immunotherapy (Wang,
 133 2023). Data are from biopsies from a single patient before
 134 and while receiving immunotherapy (about 8,000 cells total).
 135 We performed feature selection on the union of the top 3,000
 136 highly variable genes (HVGs) across both samples and a list
 137 of about 1,900 known TFs in humans. The unnormalized
 138 count matrix was input into SBGGeHPF. After normalizing
 139 and log transforming counts, we attempted GRN inference
 140 with GRNBoost from the SCENIC pipeline (Aibar, 2017),
 141 but the result GRN was densely connected with no local
 142 structure. Subsequently, we constructed a filtered covariance
 143 matrix between each HVG pair (binarized by non-zero co-
 144 variance interactions and only keeping TF-involved interac-
 145 tions) as the initial GRN for SBGGeHPF. Alternative graph
 146 inputs to SBGGeHPF can be obtained using other GRN
 147 inference methods (Lachmann, 2016; Passemiers, 2022) or
 148 causal graphs (Squires, preprint; Lopez, preprint).

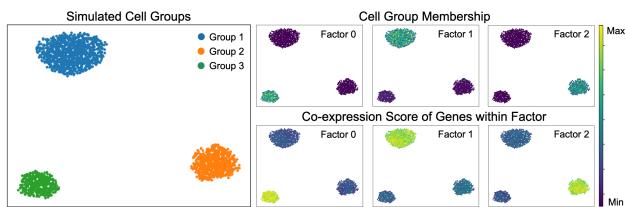
149 3. Results

150 3.1. Simulated scRNA-seq Data

151 SBGGeHPF correctly reconstructed input scRNA-seq
 152 counts and the structure of the ground-truth GRN, while
 153 recovering factors representing densely-connected com-
 154 munities. Reconstructed counts and GRN were sampled
 155 from SBGGeHPF following its generative process with
 156 parameters that are posterior means. Reconstructed and
 157 original counts had tight linear fit with slope close to 1
 158 ($R^2 = 0.990$, regression line $y = 0.984x + 0.066$), confirm-

159 ing that SBGGeHPF accurately reconstructs UMI counts.
 160 The ground-truth GRN contained three densely-connected
 161 modules representing three sets of DEGs across three dis-
 162 tinct cell groups, which SBGGeHPF recovered with high
 163 fidelity and refined structure (Supp Fig. 6).

164 Additionally, SBGGeHPF summarized groups of genes
 165 from densely-connected modules into factors describing
 166 cell heterogeneity, and it correctly assigned factors to dis-
 167 tinct subgroups of cells (Fig. 2). The simulated dataset
 168 was designed such that specific cell groups were defined by
 169 sets of DEGs. Hence, we computed co-expression of genes
 170 within factors and observed selective enrichment in their
 171 respective cell group (Fig. 2), indicating that learned factors
 172 estimate ground truth DEG sets defining cell heterogeneity.



173 **Figure 2.** Left: UMAP of ESCO simulated scRNA-seq data based
 174 on 100 DEGs, visualizing pre-defined cell groups. Top right:
 175 SBGGeHPF correctly associates learned factors to the ground
 176 truth cell groups. Cells are colored by the cell weights normalized
 177 by the cell capacity. Bottom right: SBGGeHPF learns factors
 178 that resemble ground truth DEG sets. Cells are colored by co-
 179 expression scores given by the sum of gene expression values of
 180 genes in each factor.

3.2. Melanoma Clinical scRNA-seq Data

181 To investigate the effect of SBGGeHPF on real scRNA-seq
 182 data, we tested the model on the melanoma clinical dataset
 183 and compared the performance to a naive Hierarchical Pois-
 184 son Factorization (HPF) model (Mendes Levitin, 2019).
 185 SBGGeHPF achieved better reconstruction performance
 186 ($R^2 = 0.775$, regression line $y = 0.677x + 0.113$) than
 187 HPF ($R^2 = 0.734$, regression line $y = 0.605x + 0.152$).
 188 SBGGeHPF also generated a GRN capturing both global
 189 and local topology of the input network (Fig. 3 left, middle).

190 The melanoma dataset contained cells with no clear prior re-
 191 lationship between gene modules and cell groupings across
 192 samples. Nonetheless, SBGGeHPF uncovered interpretable
 193 factors and associated them to distinct groups of cells, while
 194 HPF failed on these tasks (Fig. 4).

195 To further interpret factors learned by SBGGeHPF, we exam-
 196 ined tumor clonality predicted by InferCNV (Patel, 2014).
 197 The factors were associated with clones 1 and 3, which are
 198 two clones that exist across both samples. Factor 1 asso-
 199 ciated with on-treatment clone 1 cells, which retained its
 200 rough size. Factor 2 represented both clone 1 and clone 3

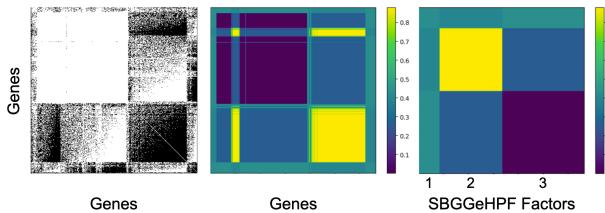


Figure 3. Left: Hierarchically sorted binary adjacency matrix of the original GRN from melanoma scRNA-seq data. Two genes can be connected if their empirical covariance is non-zero and at least one is a TF. Middle: Adjacency matrix of the generated GRN, with genes ordered as in the left figure and weights corresponding to learned edge probabilities. Right: Adjacency matrix of the generated GRN, with genes sorted by corresponding learned factors.

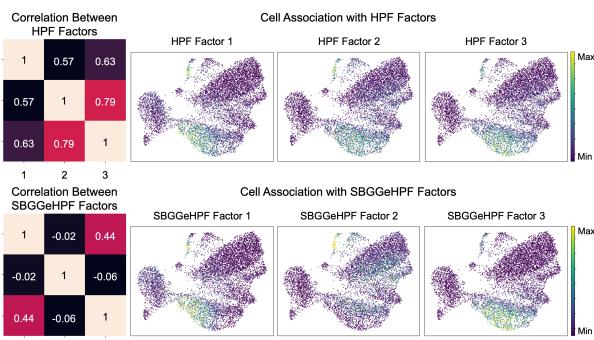


Figure 4. Top: Correlations between factors identified by HPF model on melanoma dataset, and enrichment of HPF factors across melanoma cells. Bottom: Correlations between factors identified by SBGGeHPF model on melanoma dataset, and enrichment of SBGGeHPF factors across melanoma cells.

pre-treatment cells. Factor 3 was enriched in on-treatment clone 3 cells, which expanded significantly between pre- and on-treatment sample timepoints (Fig. 5). Genes in the factors also displayed different regulatory patterns corresponding to respective clonal dynamics. Factor 1 represented a set of TFs with global regulation effect. Factor 2 was a self-regulating gene module but had weak regulatory effect on Factor 3 genes. Factor 3 genes were most likely targets with no interactions between each other (Fig. 3 right).

Finally, we performed gene set enrichment analysis (GSEA) on the factors to support that the learned modules characterizing tumor heterogeneity were also biologically relevant (Subramanian, 2005; Fang, 2023). Most factors had significant enrichment for gene sets with false discovery rate (FDR) below 0.250, including genes from both well-profiled (e.g., IL-6/JAK/STAT3, PI3K/AKT/mTOR signaling) and under-studied pathways (e.g., cholesterol homeostasis) in melanoma (Gu, 2022) (Supp Fig. 7). Taken together, when applied to the melanoma dataset, SBGGeHPF links gene regulation to tumor heterogeneity using interpretable factors, and these factors represent regulons driving cell plasticity

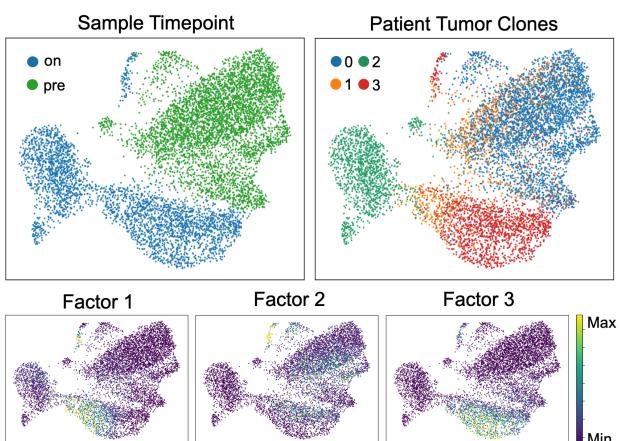


Figure 5. Top left: UMAP of melanoma scRNA-seq data after feature selection, visualizing patient samples before and on immunotherapy. Top right: Tumor clonality, inferred using InferCNV. Bottom: SBGGeHPF factor association to heterogenous cell groups. Each factor correspond to a row of the learned cell weight matrix, and cells are colored by the cell weights normalized by the learned cell capacity.

through various pathways in melanoma.

4. Conclusion and Applications

SBGGeHPF shows promise in disentangling gene regulatory elements and their role in transcriptional fate. We tested the model on simulated and real tumor sample scRNA-seq data. We showed that on both datasets, SBGGeHPF learned interpretable factors corresponding to GRN neighborhoods and mapped topological structures to heterogenous cell groups. Future work includes fine-tuning the graph learning process to work on directed causal graphs, and further interpreting network structures within learned communities. Additionally, we are working to expand this model to additional patient tumor and immune cell data. By expanding our focus to both tumor and immune cells, we hope to gain a comprehensive understanding of mechanisms associated with melanoma progression and effector cell heterogeneity driving patient response or resistance to cancer immunotherapy. Especially as clinical scRNA-seq datasets become increasingly available across cancer types and treatment conditions, SBGGeHPF and its ability to identify key cell state-specific regulons can bolster the importance of well-characterized pathways as well as uncover understudied biological mechanisms, towards improving cancer therapeutics and patient outcomes.

References

- Aibar, S. SCENIC: single-cell regulatory network inference and clustering. *Nature methods*, 2017. doi: 10.1038/nmeth.4463.

- 220 Fang, Z. GSEApY: a comprehensive package for performing
 221 gene set enrichment analysis in python. *Bioinformatics*,
 222 2023. doi: 10.1093/bioinformatics/btac757.
- 223 Gopalan, P. Scalable recommendation with hierarchical
 224 poisson factorization. *Association for Uncertainty in AI
 225 Proceedings*, 2015.
- 226 Gu, J. Cholesterol homeostasis and cancer: a new perspec-
 227 tive on the low-density lipoprotein receptor. *Cellular
 228 Oncology*, 2022. doi: 10.1007/s13402-022-00694-5.
- 229 Iglesias-Martinez, L. KBoost: a new method to infer gene
 230 regulatory networks from gene expression data. *Scientific
 231 Reports*, 2021. doi: 10.1038/s41598-021-94919-6.
- 232 Keyl, P. Single-cell gene regulatory network prediction
 233 by explainable AI. *Nucleic Acids Research*, 2023. doi:
 234 10.1093/nar/gkac1212.
- 235 Lachmann, A. ARACNe-AP: gene network reverse engi-
 236 neering through adaptive partitioning inference of mutual
 237 information. *Bioinformatics*, 2016. doi: doi.org/10.1093/
 238 bioinformatics/btw216.
- 239 Lee, C. A review of stochastic block models and extensions
 240 for graph clustering. *Applied Network Science*, 2019. doi:
 241 10.1007/s41109-019-0232-2.
- 242 Lopez, R. Large-scale differentiable causal discovery of
 243 factor graphs. preprint. doi: 10.48550/arXiv.2206.07824.
- 244 Madhamshettiwar, P. B. Gene regulatory network inference:
 245 evaluation and application to ovarian cancer allows the
 246 prioritization of drug targets. *Genome medicine*, 2012.
 247 doi: 10.1186/gm340.
- 248 Marusyk, A. Intra-tumour heterogeneity: a looking glass
 249 for cancer? *Nature reviews. Cancer*, 2012. doi: 10.1038/
 250 nrc3261.
- 251 Mendes Levitin, H. De novo gene signature identifica-
 252 tion from single-cell RNA-seq with hierarchical poisson
 253 factorization. *Molecular Systems Biology*, 2019. doi:
 254 10.15252/msb.20188557.
- 255 Mercatelli, D. Gene regulatory network inference resources:
 256 A practice overview. *BBA - Gene Regulatory Mechanisms*,
 257 2020. doi: 10.1016/j.bbagr.2019.194430.
- 258 Osorio, D. scTenifoldNet: A machine learning workflow
 259 for constructing and comparing transcriptome-wide gene
 260 regulatory networks from single-cell data. *Cell Press
 261 Patterns*, 2020. doi: 10.1016/j.patter.2020.100139.
- 262 Passemiers, A. Fast and accurate inference of gene regu-
 263 latory networks through robust precision matrix estima-
 264 tion. *Bioinformatics*, 2022. doi: 10.1093/bioinformatics/
 265 btac178.
- 266 Patel, A. Single-cell RNA-seq highlights intratumoral het-
 267 erogeneity in primary glioblastoma. *Science*, 2014. doi:
 268 10.1126/science.1254257.
- 269 Squires, C. Permutation-based causal structure learning
 270 with unknown intervention targets. preprint. doi: 10.
 271 48550/arXiv.1910.09007.
- 272 Subramanian, A. Gene set enrichment analysis: A
 273 knowledge-based approach for interpreting genome-wide
 274 expression profiles. *Proceedings of the National Academy
 275 of Sciences of the United States of America*, 2005. doi:
 276 10.1073/pnas.0506580102.
- 277 Tian, J. ESCO: single cell expression simulation incorpo-
 278 rating gene co-expression. *Bioinformatics*, 2021. doi:
 279 10.1093/bioinformatics/btab116.
- 280 Wang, Y. Multimodal single-cell and whole-genome se-
 281 quencing of small, frozen clinical specimens. *Nature
 282 Genetics*, 2023. doi: 10.1038/s41588-022-01268-9.

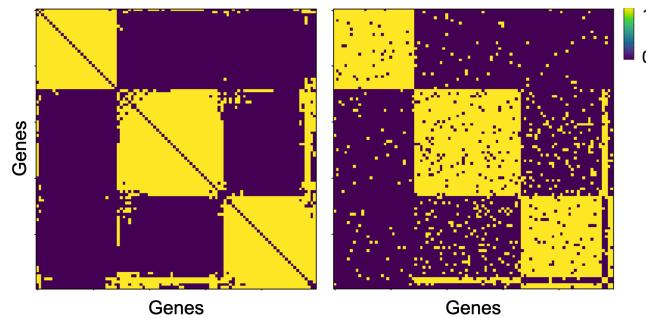
A. Supplementary Figures

Figure 6. Left: Binary adjacency matrix of original GRN from simulated scRNA-seq data, hierarchically sorted. Right: Binary adjacency matrix of reconstructed GRN, with unchanged sorting.

330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384

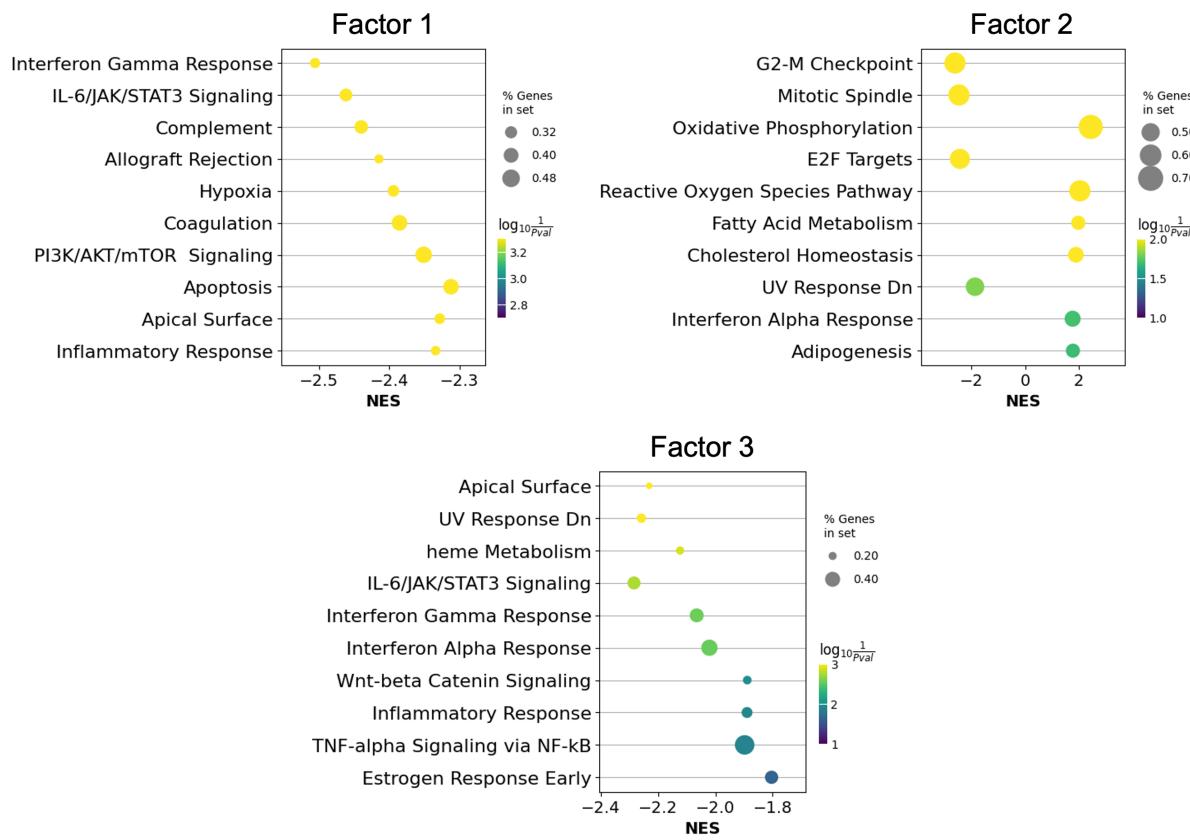


Figure 7. Top enriched pathways from gene set enrichment analysis (GSEA) performed for each SBGGeHPF factor. Cells were assigned to Factor 0, 1, or 2 based on learned cell weights, and signal-to-noise ratio was used to rank genes from cells, using factor membership as the condition. Gene enrichment was identified using hallmark gene sets, and false discovery rate cutoff of 0.25 was used to select for enriched gene sets.