
A Variational Inference Approach to Single-Cell Gene Regulatory Network Inference using Probabilistic Matrix Factorization

Omar Mahmood*
Center for Data Science
New York University

Claudia Skok Gibbs*
Center for Data Science
New York University

Richard Bonneau
Genomics & Systems Biology,
New York University
and
Prescient Design, Genentech

Kyunghyun Cho
Center for Data Science,
New York University
and
Prescient Design, Genentech
kyunghyun.cho@nyu.edu

Abstract

1 Inferring gene regulatory networks (GRNs) from single-cell gene expression
2 datasets is a challenging task. Existing methods are often designed heuristically
3 for specific datasets and lack the flexibility to incorporate additional information
4 or compare against other algorithms. Further, current GRN inference methods do
5 not provide uncertainty estimates with respect to the interactions that they predict,
6 making inferred networks challenging to interpret. To overcome these challenges,
7 we introduce Probabilistic Matrix Factorization for Gene Regulatory Network infer-
8 ence (PMF-GRN). PMF-GRN uses single-cell gene expression data to learn latent
9 factors representing transcription factor activity as well as regulatory relationships
10 between transcription factors and their target genes. This approach incorporates
11 available experimental evidence into prior distributions over latent factors and
12 scales well to single-cell gene expression datasets. By utilizing variational infer-
13 ence, we facilitate hyperparameter search for principled model selection and direct
14 comparison to other generative models. To assess the accuracy of our method,
15 we evaluate PMF-GRN using the model organisms *Saccharomyces cerevisiae* and
16 *Bacillus subtilis*, benchmarking against database-derived gold standard interactions.
17 We discover that, on average, PMF-GRN infers GRNs more accurately than current
18 state-of-the-art single-cell GRN inference methods. Moreover, our PMF-GRN ap-
19 proach offers well-calibrated uncertainty estimates, as it performs gene regulatory
20 network (GRN) inference in a probabilistic setting. These estimates are valuable
21 for validation purposes, particularly when validated interactions are limited or a
22 gold standard is incomplete.

23 **Keywords**— Probabilistic Matrix Factorization, Variational Inference, Gene Regulatory Network Inference,
24 Single Cell, Gene Expression.

*Equal contribution

25 1 Background

26 An essential problem in systems biology is to extract information from genome wide sequencing data to unravel
27 the mechanisms controlling cellular processes within heterogeneous populations (1). Gene regulatory networks
28 (GRNs) that annotate regulatory relationships between transcription factors (TFs) and their target genes (2) have
29 proven to be useful models for stratifying functional differences between cells (3; 4; 5; 6) that can arise during
30 normal development (7), responses to environmental signals (8) and dysregulation in the context of disease
31 (9; 10; 11).

32 GRNs cannot be directly measured with current sequencing technology. Instead, methods must be developed
33 to piece together snapshots of transcriptional processes in order to reconstruct a cell's regulatory landscape
34 (12). Initial approaches to GRN inference relied on Microarray technology (13; 14; 15), a hybridization-based
35 method to measure the expression of thousands of genes simultaneously (16). This technology was biased
36 as it was limited to only those genes that were annotated at the time, which in turn presented challenges for
37 inferring the complete regulatory landscape (1). Subsequently, the high-throughput sequencing method RNA-seq
38 provided a genome wide readout of transcriptional output, allowing for the detection of novel transcripts (17)
39 and thus improving GRN inference potential. More recently, single-cell RNA-seq technology has enabled the
40 characterization of gene expression profiles within heterogeneous populations (18), vastly increasing the potential
41 for GRN inference algorithms (19; 20). In contrast to bulk RNA experiments (Microarray and RNA-seq) that
42 average measurements of gene expression across heterogeneous cell populations, GRNs inferred from single-cell
43 data have the advantage of unmasking biological signal in distinct cells (21).

44 Several matrix factorization approaches have been proposed to overcome the limitations of reconstructing GRNs
45 from Microarray data (22). These include use of statistical techniques such as Singular Value Decomposition
46 and Principal Component Analysis (23), Bayesian Decomposition (24), and Non-negative Matrix Factorization
47 (25; 26; 27). More recently, matrix factorization approaches have been applied to integrative analysis of DNA
48 methylation and miRNA expression data (28), as well as single-cell RNA-seq and single-cell ATAC-seq data
49 (29). However, to the best of our knowledge, these matrix factorization approaches have not yet been used to
50 infer GRNs from single-cell gene expression data. Meanwhile, several regression-based methods have been
51 proposed to learn GRNs from single-cell RNA-seq and single-cell ATAC-seq to capture regulatory relationships
52 at single-cell resolution (30). So far, these integrative approaches to GRN inference have been successfully
53 implemented using regularized regression (31), self-organizing maps (32), tree-based regression (33), and
54 Bayesian Ridge regression (34).

55 Although regression-based methods for inferring GRNs from single-cell data are available, they still suffer
56 from significant limitations (35). Firstly, these methods heavily rely on the input data used to learn the GRN,
57 causing issues when new data becomes available or new assumptions are required in the model. This can result
58 in inaccurate predictions if the new data or assumptions are not well integrated into the existing model, leading
59 to the need for a complete re-design of the algorithm, which can be costly and time-consuming. Additionally,
60 these methods typically focus on inferring a single GRN that explains the available data, without performing
61 hyperparameter search to determine the optimal model. This can lead to heuristic model selection, with no
62 justification for the approach taken or evidence that the best possible model has been selected. Conversely,
63 hyperparameter search ensures the accuracy of the GRN inference algorithm by finding the optimal model
64 that fits the data well while avoiding overfitting. Regression-based GRN inference algorithms that do not
65 perform hyperparameter search may miss important data features or overemphasize irrelevant ones, leading to
66 inaccurate or incomplete models. Moreover, these methods do not provide an indication of their uncertainty
67 about the predictions that they make. Finally, several regression-based GRN inference algorithms struggle to
68 scale optimally to the size of typical single-cell datasets, limiting inference to small subsets of data or requiring
69 enormous amounts of computational time."

70 In this study, we introduce PMF-GRN, a novel approach that uses probabilistic matrix factorization (36) to
71 infer gene regulatory networks from single-cell gene expression and chromatin accessibility information. This
72 approach extends previous methods that applied matrix factorization for GRN inference with Microarray data,
73 to address the current limitations in regression-based single-cell GRN inference. We implement our approach
74 in a probabilistic setting with variational inference, which provides a flexible framework to incorporate new
75 assumptions or biological data as required, without changing the way the GRN is inferred. We also use a
76 principled hyperparameter selection process with the Evidence Lower Bound (ELBO) objective function, which
77 optimizes the parameters of our probabilistic model for automatic model selection. In this way, we replace
78 heuristic model selection by comparing a variety of generative models and hyperparameter configurations before
79 selecting the optimal parameters with which to infer a final GRN. Our probabilistic approach provides uncertainty
80 estimates for each predicted regulatory interaction, serving as a proxy for the model confidence in each predicted
81 interaction. Uncertainty estimates can be useful in the situation where there are limited validated interactions
82 or a gold standard is incomplete. By using stochastic gradient descent (SGD), we perform GRN inference on
83 a GPU, allowing us to easily scale to a large number of observations in a typical single-cell gene expression
84 dataset. Unlike many existing methods, PMF-GRN is not limited by pre-defined organism restrictions, making it
85 widely applicable for GRN inference.

86 To demonstrate the novelty and advantages of PMF-GRN, we apply our method to two single-cell gene expression
 87 datasets for the model organism *Saccharomyces cerevisiae*. We evaluate our model’s performance in a normal
 88 inference setting, as well as with cross-validation and noisy data. To assess the accuracy of predicted regulatory
 89 interactions, we evaluate all regulatory predictions using Area Under the Precision Recall Curve (AUPRC)
 90 against database derived gold standards. Our findings show that the uncertainty estimates are well-calibrated for
 91 inferred TF-target gene interactions, as the accuracy of predictions increases when the associated uncertainty
 92 decreases. Here, in comparison to three state-of-the-art regression-based methods for inferring single-cell GRNs,
 93 namely the Inferelator (31), Scenic (33), and Cell Oracle (34), our method demonstrates an overall improved
 94 performance in recovering the true underlying GRN. We also include GRNs inferred using two microarray
 95 datasets for *Bacillus subtilis* by converting expression values to integers to simulate a single-cell-like experiment,
 96 demonstrating our method’s performance on a second dataset.

97 2 Results

98 2.1 The PMF-GRN Model

99 The goal of our probabilistic matrix factorization approach is to decompose observed gene expression into latent
 100 factors, representing TF activity (TFA) and regulatory interactions between TFs and their target genes. These
 101 latent factors, which represent the underlying GRN, cannot be measured experimentally, unlike gene expression.
 102 We model an observed gene expression matrix $W \in \mathbb{R}^{N \times M}$ using a TFA matrix $U \in \mathbb{R}_{>0}^{N \times K}$, a TF-target gene
 103 interaction matrix $V \in \mathbb{R}^{M \times K}$, observation noise $\sigma_{obs} \in (0, \infty)$ and sequencing depth $d \in (0, 1)^N$, where N
 104 is the number of cells, M is the number of genes and K is the number of TFs. We rewrite V as the product
 105 of a matrix $A \in (0, 1)^{M \times K}$, representing the degree of existence of an interaction, and a matrix $B \in \mathbb{R}^{M \times K}$
 106 representing the interaction strength and its direction:

$$V = A \odot B,$$

107 where \odot denotes element-wise multiplication. An overview of the graphical model is shown in Figure 1A.

108 These latent variables are mutually independent *a priori*, i.e., $p(U, A, B, \sigma_{obs}, d) =$
 109 $p(U)p(A)p(B)p(\sigma_{obs})p(d)$. For the matrix A , prior hyperparameters represent an initial guess of the
 110 interaction between each TF and target gene which need to be provided by a user. These can be derived from
 111 genomic databases or obtained by analyzing other data types, such as the measurement of chromosomal
 112 accessibility, TF motif databases, and direct measurement of TF-binding along the chromosome, as shown in
 113 Figure 1B (see Methods section for details).

114 The observations W result from a matrix product UV^T . We assume noisy observations by defining a distribution
 115 over the observations with the level of noise σ_{obs} , i.e., $p(W|U, V = A \odot B, \sigma_{obs}, d)$.

116 Given this generative model, we perform posterior inference over all the unobserved latent variables; U, A, B, d
 117 and σ_{obs} , and use the posterior over A to investigate TF-gene interactions. Exact posterior inference with an
 118 arbitrary choice of prior and observation probability distributions is, however, intractable. We address this issue
 119 by using variational inference (37; 38), where we approximate the true posterior distributions with tractable,
 120 approximate (variational) posterior distributions.

121 We minimize the KL-divergence $D_{KL}(q||p)$ between the two distributions with respect to the parameters of
 122 the variational distribution q , where p is the true posterior distribution. This allows us to find an approximate
 123 posterior distribution q that closely resembles p . This is equivalent to maximizing the evidence lower bound
 124 (ELBO) i.e. a lower bound to the marginal log likelihood of the observations W :

$$\begin{aligned} \log p(W) \geq \mathbb{E}_{U, A, B, \sigma_{obs}, d \sim q(U, A, B, \sigma_{obs}, d)} [& \log p(W|U, V = A \odot B, \sigma_{obs}, d) \\ & + \log p(U, A, B, \sigma_{obs}, d) \\ & - \log q(U, A, B, \sigma_{obs}, d)] \end{aligned}$$

125 The mean and variance of the approximate posterior over each entry of A from maximizing the ELBO are then
 126 used as the degree of existence of an interaction between a TF and a target gene and its uncertainty, respectively.

127 It is important to note that matrix factorization based GRN inference is only identifiable up to a latent factor
 128 (column) permutation. In the absence of prior information, the probability that the user assigns TF names to the
 129 columns of U and V in the same order as the order in which the inference algorithm implicitly assigns TFs to
 130 these columns is $\frac{1}{K!}$, which is essentially 0 for any reasonable value of K . Incorporating prior-knowledge of
 131 TF-target gene interactions into the prior distribution over A is therefore essential to give the inference algorithm
 132 information about which column corresponds to which TF.

133 With this identifiability issue in mind, we design an inference procedure that can be used on any dataset, described
 134 in Figure 1C. The first step is to randomly hold out prior information for some percentage of the genes in $p(A)$

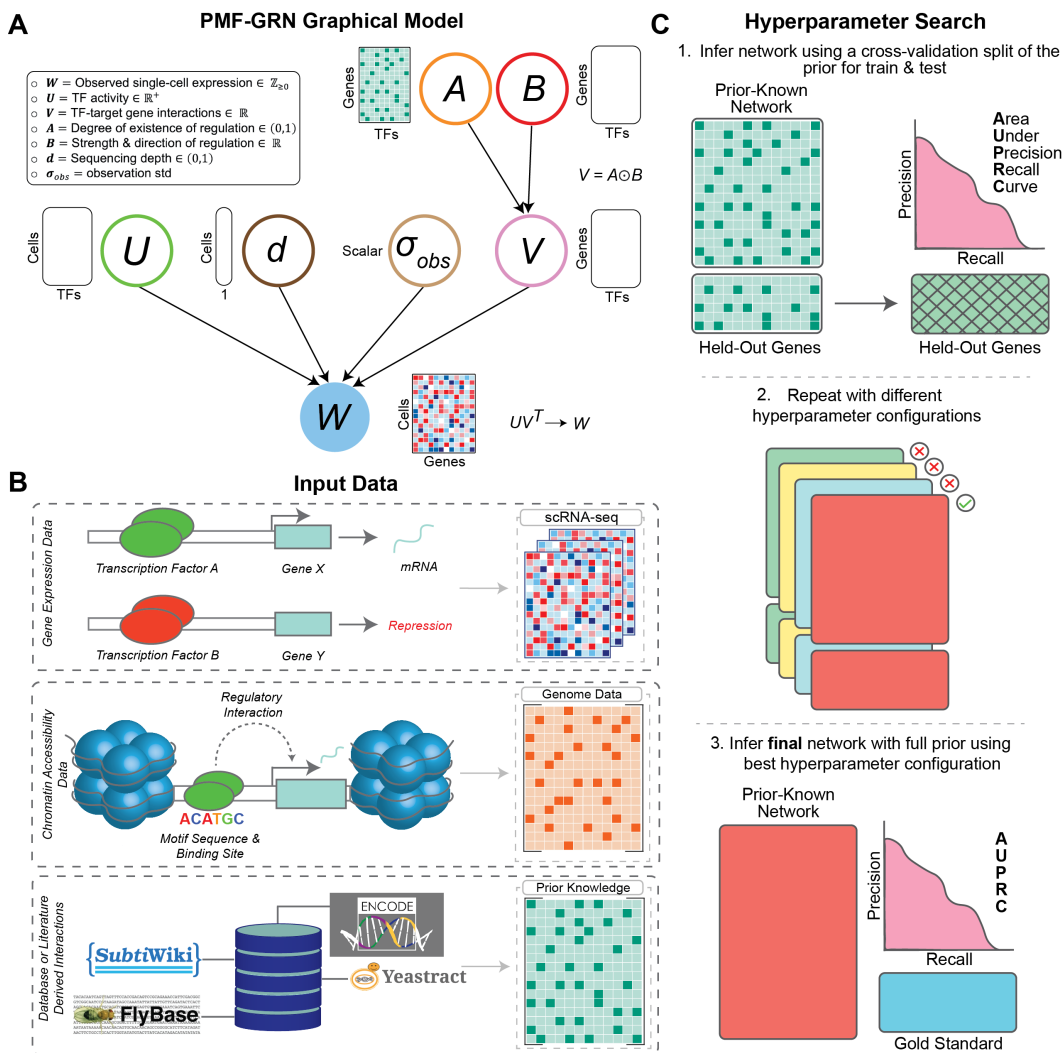


Figure 1: (A) PMF-GRN graphical model overview. Input single-cell gene expression W is decomposed into latent factors U and V , representing TF activity and TF-gene interactions respectively. V is further decomposed into A and B , representing the degree of existence of interaction, and the strength and direction of an interaction, respectively. Information obtained from chromatin accessibility data or genomics databases is incorporated into the prior distribution for A . Additional latent variables are included to model observation noise σ_{obs} and sequencing depth d , in order to better model our observed single-cell gene expression input data. (B) Input experimental data for PMF-GRN includes single-cell RNA-seq gene expression data. ATAC-seq is used to determine chromatin accessibility through peak calling. Motif enrichment within these accessible regions can be used to create a prior-known network to better inform the prior distribution. When experimental information is unavailable, databases can be used instead to construct a known-prior network. (C) Hyperparameter selection process is performed for model selection. The provided prior-known network is split into a train and validation dataset. 80% of the prior-known information is used to infer a GRN, while the remaining 20% is used for validation by computing AUPRC. This process is repeated multiple times, using different hyperparameter configurations in order to determine the optimal hyperparameters for the GRN inference task at hand. Finally, using the optimal hyperparameters, as determined by the highest achieved AUPRC, a final network is inferred using the full prior and evaluated using an independent gold standard.

135 (we choose 20%) by leaving the rows corresponding to these genes in A but setting the prior logistic normal
136 means for all entries in these rows to be the same low number.

137 The second step is to carry out a hyperparameter search using this modified prior-knowledge matrix. The early
138 stopping and model selection criteria are both the ‘validation’ AUPRC of the posterior point estimates of A
139 corresponding to the held out genes against the entries for these genes in the full prior hyperparameter matrix.
140 This step is motivated by the idea that inference using the selected hyperparameter configuration should yield a
141 GRN whose columns correspond to the TF names that the user has assigned to these columns.

142 The third step is to choose the hyperparameter configuration corresponding to the highest validation AUPRC and
143 perform inference using this configuration with the full prior. An importance weighted estimate of the marginal
144 log likelihood is used as the early stopping criterion for this step. The resulting approximate posterior provides
145 the final posterior estimate of A .

146 **2.2 Advantages of PMF-GRN**

147 Existing methods almost always couple the description of the data generating process with the inference
148 procedure used to obtain the final estimated GRN (31; 34; 33). Designing a new model thus requires designing a
149 new inference procedure specifically for that model, which makes it difficult to compare results across different
150 models due to the discrepancies in their associated inference algorithms. Furthermore, this *ad hoc* nature of
151 model building and inference algorithm design often leads to the lack of a coherent objective function that can be
152 used for proper hyperparameter search as well as model selection and comparison, as evident in (31). Heuristic
153 model selection in available GRN inference methods presents the challenge of determining and selecting the
154 optimal model in a given setting.

155 The proposed PMF-GRN framework decouples the generative model from the inference procedure. Instead of
156 requiring a new inference procedure for each generative model, it enables a single inference procedure through
157 (stochastic) gradient descent with the ELBO objective function above, across a diverse set of generative models.
158 Inference can easily be performed in the same way for each model. Through this framework, it is possible to
159 define the prior and likelihood distributions as desired with the following mild restrictions: we must be able to
160 evaluate the joint distribution of the observations and the latent variables, the variational distribution and the
161 gradient of the log of the variational distribution.

162 The use of stochastic gradient descent in variational inference comes with a significant computational advantage.
163 As each step of inference can be done with a small subset of observations, we can run GRN inference on a very
164 large dataset without any constraint on the number of observations. This procedure is further sped up by using
165 modern hardware, such as GPUs.

166 Under this probabilistic framework, we carry out model selection, such as choosing distributions and their
167 corresponding hyperparameters, in a principled and unified way. Hyperparameters can be tuned with regard
168 to a predefined objective, such as the marginal likelihood of the data or the posterior predictive probability of
169 held out parts of the observations. We can further compare and choose the best generative model using the same
170 procedure.

171 This framework allows us to encode any prior knowledge via the prior distributions of latent variables. For
172 instance, we incorporate prior domain knowledge about TF-gene interactions as hyperparameters that govern the
173 prior distribution over the matrix A . If prior knowledge about TFA is available, this can be similarly incorporated
174 into the model via the hyperparameters of the prior distribution over U .

175 Because our approach is probabilistic by construction, inference also estimates uncertainty without any separate
176 external mechanism. These uncertainty estimates can be used to assess the reliability of the predictions, i.e.,
177 more trust can be placed in interactions that are associated with less uncertainty. We verify this correlation
178 between the degree of uncertainty and the accuracy of interactions in the experiments.

179 Overall, the proposed approach of probabilistic matrix factorization for GRN inference is scalable, generalizable
180 and aware of uncertainty, which makes its use much more advantageous compared to most existing methods.

181 **2.3 PMF-GRN Recovers True Interactions in Simple Eukaryotes**

182 To demonstrate PMF-GRNs ability to infer informative and robust GRNs, we use two single-cell RNA-seq
183 datasets from the model organism *Saccharomyces cerevisiae*. *S.cerevisiae* is a relatively simple and well studied
184 eukaryote with an available and reliable gold standard, which allows us to test and evaluate our models
185 performance.

186 We perform three experiments using two independently collected single-cell RNA-seq *S. cerevisiae* datasets
187 (8; 39) to test PMF-GRN and compare our performance against three state-of-the-art GRN inference methods,
188 the Inferelator (AMuSR, BBSR, StARS) (31), Scenic (33), and CellOracle (34). In the first experiment, we infer
189 a GRN for each of the two single-cell datasets and average the posterior means of A to simulate a "multi-task"

190 GRN inference approach for building the final combined network. Using AUPRC, we show that PMF-GRN
191 outperforms AMuSR, StARS, and Scenic, while performing competitively with BBSR and CellOracle (Figure
192 2A). To provide a baseline for each method in the scenario where data cannot be cleanly separated into tasks, we
193 combine the two expression datasets into one observation before inferring a GRN. This baseline demonstrates a
194 large performance decrease for BBSR, indicating that the method may only be useful when gene expression
195 is organized into tasks. This could present challenges when attempting to infer GRNs in more complicated
196 organisms where cell-types or conditions are less easily defined. Here, we show the effectiveness of PMF-GRN
197 in recovering the true underlying GRN for both scenarios, as performance remains relatively stable whether a
198 network was inferred using the individual or combined data. We also provide an example as to how we can
199 use PMF-GRN on a single observation or multiple observation matrices to infer a consensus GRN by simple
200 averaging.

201 In the second experiment, we implement a 5 fold cross-validation approach to establish a baseline for each
202 model. Cross-validation is an essential technique for evaluating the performance of machine learning models like
203 PMF-GRN as it allows us to test our method's ability to generalize to new data. Cross-validation further allows
204 us to simulate the process of training and testing PMF-GRN on multiple subsets of the available data, providing
205 a more robust and reliable estimate of model accuracy. In the context of GRN inference, cross-validation is
206 particularly important because it helps us assess the performance of PMF-GRN in predicting TF-target gene
207 interactions based on limited data, which is often the case in experimental settings.

208 We first combine the two *S. cerevisiae* single-cell RNA-seq datasets into one observation matrix for simplicity. To
209 perform cross-validation, the gold standard is divided into an 80% – 20% split, where a network is inferred using
210 80% of the gold standard as "prior-known information", and evaluated using the remaining 20%. We repeat this
211 cross-validation process five times using different random splits of the gold standard to obtain meaningful results.
212 We observe that PMF-GRN outperforms Scenic and CellOracle, while achieving competitive performance to
213 BBSR and StARS (Figure 2B). We note that for this experiment, we are unable to implement the AMuSR
214 algorithm as it is a multi-task inference approach that requires more than one task (dataset).

215 Finally, in the third experiment, we demonstrate the robustness of each GRN inference method against noisy
216 prior information. To do so, we infer GRNs where increasing amounts of noise have been added to the input
217 prior-known information. Here, we show that as noise increases, PMF-GRN's AUPRC decreases similarly to
218 CellOracle, while on average, performing better than BBSR, StARS and CellOracle, demonstrating that it is one
219 of the most robust approaches to inferring accurate GRNs from noisy priors (Figure 2C).

220 From the results of our experiments on the *S. cerevisiae* data, we have the following observations. The first main
221 observation is that on average the proposed PMF-GRN performs better than the Inferelator in recovering the true
222 GRN, regardless of whether we pick the mean or median Inferelator algorithm in terms of AUPRC. Specifically,
223 we see that PMF-GRN performs markedly better than two Inferelator algorithms (AMuSR and StARS), and
224 similarly to the remaining algorithm (BBSR). However, when the expression data is not separated into tasks,
225 PMF-GRN outperforms BBSR. In comparison to CellOracle, we observe that PMF-GRN infers competitive
226 GRNs during normal inference. However, PMF-GRN greatly outperforms CellOracle when performing cross-
227 validation. Finally, we observe that PMF-GRN consistently outperforms Scenic in all experiments considered.

228 The second main observation is that our approach eliminates the high variance associated with choosing between
229 different inference algorithms. Implementing the Inferelator on the *S. cerevisiae* datasets in a normal setting
230 yields AUPRCs approximately in the range 0.2 to 0.4, without any a priori information on which of these
231 algorithms to use. The resulting inferred GRN could be arbitrarily accurate or inaccurate depending on which
232 algorithm is chosen. In contrast, our method is reliable as it provides one set of results, chosen using a
233 principled objective function, performing competitively with the best performing Inferelator algorithm (BBSR)
234 and CellOracle.

235 Finally, in order to highlight the identifiability issue and ensure that the prior-known information provided is
236 useful, we demonstrate the performance of PMF-GRN where prior information is not used (e.g. all prior logistic
237 normal means of A are set to the same low number). We use the same process for all other GRN inference
238 algorithms by providing an empty prior. Additionally, we demonstrate PMF-GRN's performance when we
239 randomly shuffle the prior-known TF-target gene interaction hyperparameters before using them to build the
240 prior distribution for A . We repeat this process for all other GRN inference algorithms by providing them with
241 prior-known information in which the gene names have been shuffled randomly. As anticipated, the resulting
242 AUPRC scores are close to 0, implying that our approach, as well as the Inferelator and CellOracle are capable
243 of taking into account such prior information well and that the prior information we provided is useful and
244 reliable (see Methods section for details). The results for GRNs inferred without prior-known information are
245 demonstrated by the black dots, while the gray dots demonstrate GRNs inferred with shuffled prior-known
246 information, shown in Figure 2A.

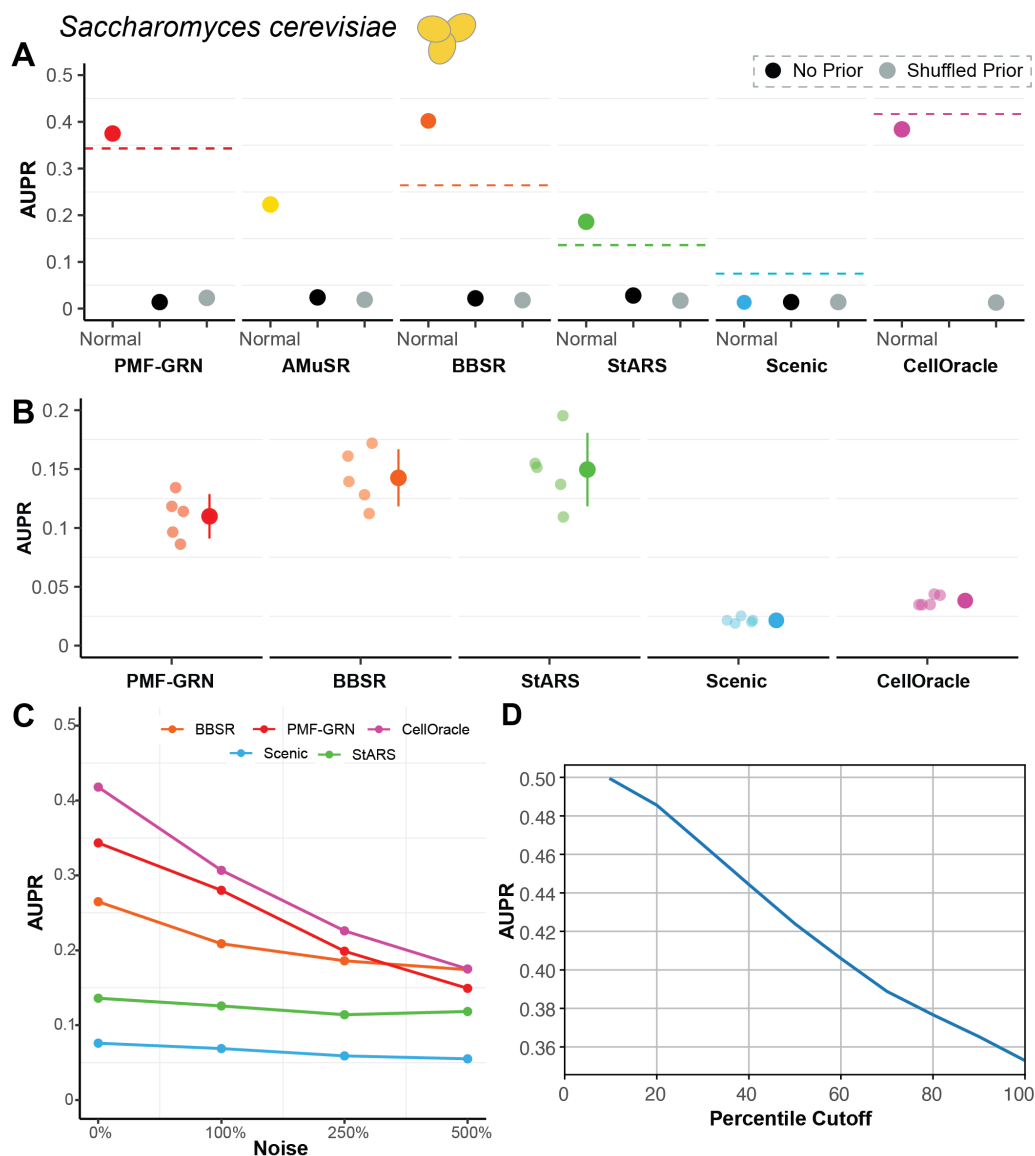


Figure 2: GRN inference in *Saccharomyces cerevisiae*. **(A)** Consensus network AUPR using gold standard network. Performance of PMF-GRN (red) is compared to three Inferelator algorithms, AMuSR (yellow), BBSR (orange), StARS (green), as well as to Scenic (blue), and CellOracle (purple). A baseline for each method (dashed line), demonstrates performance if the expression data is combined into one task. Two negative controls, no prior information (black) and shuffled prior information (gray), are inferred to ensure reliable results. **(B)** 5 fold cross-validation establishes a baseline for each model. Low-opacity dots represent each of the five cross-validation experiments. The mean AUPR \pm standard deviation for each GRN inference method is depicted by colored dot and line. **(C)** GRNs inferred with increasing amounts of noise added to the prior. **(D)** Calibration results on the *S.cerevisiae* (GSE144820 (8) only) dataset. Posterior means are cumulatively placed in bins based on their posterior variances. The x-coordinate x of each point in the plot represents all posterior means that correspond to the bottom $x\%$ of posterior variances. The y-coordinate is the ‘overlap’ AUPRC (see Methods section for details) calculated on these posterior means against the gold standard.

2.4 PMF-GRN Recovers True Interactions in Prokaryotes as Evaluated by Cross-Validation

To demonstrate GRN inference on a second dataset, we carry out experiments using two microarray datasets for the prokaryote *Bacillus Subtilis* (B1 - GSE27219 (40) and B2 - GSE67023 (41)). Although PMF-GRN is not primarily designed to learn GRNs from microarray data, we show that it is still possible to learn informative GRNs with this data. For our *B. subtilis* experiments, we have access to prior-knowledge derived from the subtiwiki database (42; 43; 44). Here, we implement a 5 fold cross-validation approach by using five random splits of the subtiwiki database-derived information, where 80% is used as prior knowledge and 20% is used as the gold standard for evaluation.

The two *B. subtilis* datasets were previously normalized after data collection as part of standard microarray processing. However, each dataset was normalized using different approaches (described in 4). For B1, the expression data underwent no further normalization and was simply converted to integers to simulate single-cell-like data. For B2, the expression data was re-scaled and then converted to integers, in order to contain only positive integers resembling single-cell-like data. The results from our experiments are shown in Figure 3, and the numbers used to create this figure are given in Supplementary Table 4 and 5. Using five repeats of cross-validation, we show the performance of GRNs inferred for the two *B. subtilis* datasets (B1 and B2). We remark that the difference in performance between B1 and B2 is likely a result of the chosen microarray processing normalization. To further support this claim, we demonstrate GRN performance after re-scaling the data with min-max scaling (Supplemental Figure S1).

'No Prior' and 'Shuffled' results are also shown in Figure 3 by black and gray dots respectively. Here, we are able to demonstrate that for B1 and B2, each GRN yields a better performance as compared to negative controls.

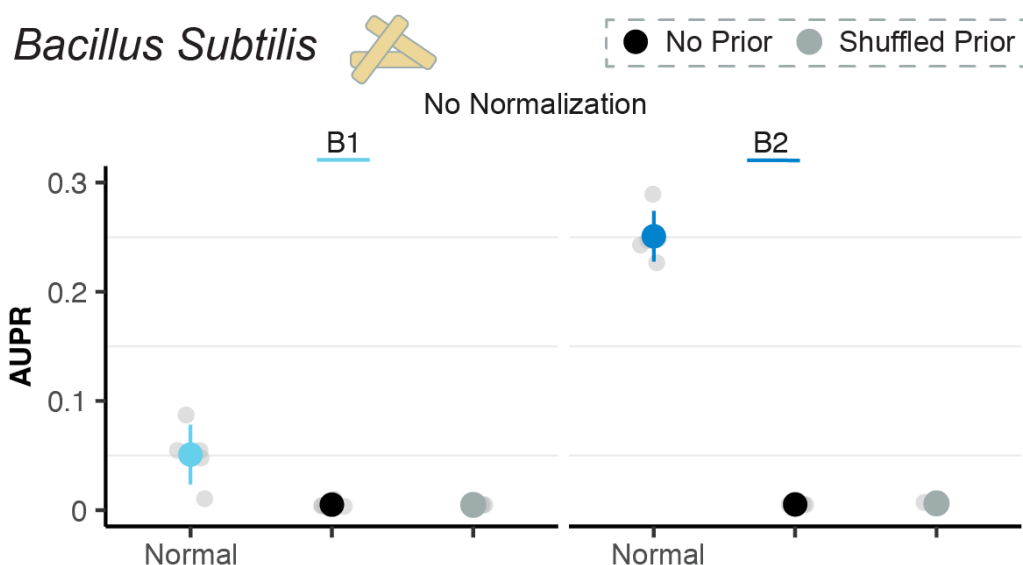


Figure 3: Results for GRNs learned in *B. subtilis* datasets B1 and B2 without data normalization. Light gray dots represent the results for each of the 5 cross-validation experiments. Colored dots represent the mean of the cross-validation experiments \pm standard deviation. Negative controls are demonstrated by black dots for "No Prior" and grey dots for "Shuffled Prior".

2.5 PMF-GRN Provides Well-Calibrated Uncertainty Estimates

Through our inference procedure, we obtain a posterior variance for each element of A , in addition to the posterior mean. We interpret each variance as a proxy for the uncertainty associated with the corresponding posterior point estimate of the relationship between a TF and a gene. Due to our use of variational inference as the inference procedure, our uncertainty estimates are likely to be underestimates. However, these uncertainty estimates still provide useful information as to the confidence the model places in its point estimate of each interaction. We expect posterior estimates associated with lower variances (uncertainties) to be more reliable than those with higher variances.

275 In order to determine whether this holds for our posterior estimates, we cumulatively bin the posterior means of
276 A according to their variances, from low to high. We then calculate the AUPRC for each bin as shown for the
277 GSE125162 (8) *S.cerevisiae* dataset in Figure 2 D. It is evident from the figures that the AUPRC decreases as
278 the posterior variance increases. Stated differently, inferred interactions associated with lower uncertainty are
279 more likely to be accurate than those associated with higher uncertainty. This is in line with our expectations.
280 The more certain the model is about the degree of existence of a regulatory interaction, the more accurate it is
281 likely to be, showing that our model is well-calibrated.

282 3 Conclusion

283 In this paper we present a framework for probabilistic matrix factorization, optimized using automatic variational
284 inference, for inferring GRNs from single cell gene expression data. In contrast with previous methods, our
285 framework decouples the model that defines the data generation process from the inference procedure. Concretely,
286 this means that we can modify the latent variables that constitute the model, along with their distributions,
287 without altering the inference procedure. This flexibility will allow for different sequencing data and modeling
288 assumptions to be readily incorporated into the model. Building new models no longer requires defining a new
289 inference procedure, which has previously been the case.

290 Additionally, PMF-GRN provides a principled way to carry out model selection and hyperparameter configuration
291 by using the same objective function and inference procedure across all models. This feature differs from previous
292 GRN methods, where it is often unclear which algorithm or hyperparameters to use for a given dataset. In
293 the PMF-GRN framework, we carry out hyperparameter searches across generative models and choose the
294 configuration that corresponds to the optimal value of the objective function. This greatly reduces the need for
295 heuristic model selection.

296 In order to demonstrate successful GRN inference, we infer a consensus GRN for *S.cerevisiae* using our
297 principled model selection method, and compare our results to GRNs inferred by the Inferelator, Scenic and
298 CellOracle with respect to a reliable gold standard. Whereas the Inferelator yields a set of highly varying
299 results across the variants, our approach results in a single inferred GRN. This GRN yields an AUPRC that
300 is higher than Scenic, as well as than the mean and median AUPRC achieved by the Inferelator's respective
301 algorithms (AMuSR and StARS), and is comparable to the AUPRC achieved by the best performing Inferelator
302 algorithm (BBSR), in addition to CellOracle. However, when the expression data is not separated into tasks,
303 we demonstrate that BBSR can no longer recover a competitive network. Our model hence yields a reliable
304 high-performing set of results without any need for heuristic model selection.

305 We further evaluate PMF-GRN by performing cross-validation and find that PMF-GRN, BBSR, and StARS
306 yield high performance, indicating that these methods are able to generalize well to new data and do not overfit
307 to training data. In contrast, Scenic and CellOracle do not perform well during cross-validation, indicating that
308 these methods may not be generalizable.

309 Finally, because prior-known networks are inherently noisy due to limited validated regulatory interactions,
310 we design an experiment to test each algorithm's robustness against increasing amounts of noise in the input
311 prior-known information. This experiment identifies PMF-GRN and CellOracle as the methods which are
312 overall most robust against noisy priors, indicating that inferred GRNs will remain reliable regardless of noisy
313 interactions introduced into the prior during inference.

314 Using two microarray *B. subtilis* datasets, we further demonstrate that PMF-GRN is capable of learning
315 informative GRNs. To include database information into both the prior knowledge and evaluation, we use an
316 approach motivated by cross-validation. Here, we find that scaling allow us to place these datasets on the same
317 scale, allowing them to be more comparable during inference. Although PMF-GRN is not primarily designed for
318 microarray data, we show that is still possible to learn informative networks by simply converting the expression
319 to integers to represent single-cell-like counts.

320 In order to determine the effect of incorporating our prior domain knowledge into the model, we compare results
321 obtained using shuffled and unshuffled hyperparameters for the matrix A . We observe that for both *S. cerevisiae*
322 and *B. subtilis*, not using prior information or shuffling the prior information results in very low AUPRCs,
323 whereas using the prior information as intended results in significantly better AUPRCs. This result holds for
324 PMF-GRN as well as for CellOracle and all Inferelator algorithms. However, for Scenic, we show that it is
325 challenging to obtain a GRN that performs better than these negative controls. This shows that prior information
326 is essential for addressing the latent factor identifiability issue and obtaining interpretable results from matrix
327 factorization, as well as regression based approaches.

328 In contrast to previous methods, our model provides well-defined uncertainty estimation in addition to point
329 estimates of GRNs. We evaluate these uncertainty estimates as provided by our model, by computing the AUPRC
330 for inferred TF-target gene interactions corresponding to different levels of posterior uncertainty. We find that the
331 AUPRC increases as the posterior variance decreases, demonstrating that when our model is more certain about
332 its estimates, it produces better rankings of TF-target gene interactions compared to when it is uncertain. This

333 indicates that our model is well-calibrated. For downstream experimental validation, biologists could therefore
 334 place more trust in model estimates that have a lower posterior variance. We also note that the computational
 335 cost of our model scales linearly with the number of cells in the dataset. This enables application of our method
 336 to single-cell RNA-seq datasets of any size.

337 We envision many possible directions for future work to design a better algorithm for inferring GRNs under
 338 our framework. This framework could be extended to explicitly model multiple expression matrices and their
 339 batch effects. We could probabilistically model prior information for A obtained from ATAC-seq and TF motif
 340 databases, and include this as part of the probabilistic model over which we carry out inference. Evaluating
 341 the posterior estimates of the direction of transcriptional regulation, provided by the matrix B , could provide
 342 a useful benchmark for the computational estimation of TF activation and repression. Research could also be
 343 carried out on improved self-supervised objectives for hyperparameter selection.

344 Future work could also focus on how to use results from our framework to guide experimental wet-lab work. For
 345 example, the uncertainty quantification provided by our model could open up new research directions in active
 346 learning for GRN inference. Highly ranked, uncertain interactions could be experimentally tested and the results
 347 fed back into the prior hyperparameter matrix for A . Inference with this updated matrix would ideally yield a
 348 better posterior GRN estimate. Posterior estimates of TFA provided by our model could be useful to wet lab
 349 scientists, as this quantity incorporates information on post-transcriptional modifications.

350 Most importantly, the study of GRN inference is far from complete. So far, this has required new computational
 351 models and assumptions in order to keep up with relevant sequencing technologies. It is thus essential to develop
 352 a model that can be easily adapted to new biological datasets as they become available, without having to
 353 completely re-build each model. We have therefore proposed PMF-GRN as a modular, principled, probabilistic
 354 approach that can be easily adapted to both new and different biological data without having to design a new
 355 GRN inference method.

356 4 Methods

357 4.1 Model Details

358 We index cells, genes and TFs using $n \in \{1, \dots, N\}$, $m \in \{1, \dots, M\}$ and $k \in \{1, \dots, K\}$, respectively.
 359 We treat each cell's expression profile W_n as a random variable, with local latent variables U_n and d_n , and
 360 global latent variables (that are shared among all cells) σ_{obs} and $V = A \odot B$. We use the following likelihood
 361 for each of our observations:

$$p(W_n|U, V, \sigma_{obs}, d) = \mathcal{N}(d_n * U_n V^\top, \sigma_{obs}^2).$$

362 We assume that U, V, σ_{obs} and d are independent i.e. $p(U, V, \sigma_{obs}, d) = p(U)p(V)p(\sigma_{obs})p(d)$. In addition to
 363 our iid assumption over the rows of U and d , We also assume that the entries of U_n are mutually independent,
 364 and that all entries of A and B are mutually independent. We choose a lognormal distribution for our prior over
 365 U and a logistic Normal distribution for our prior over d :

$$\begin{aligned} p(\log(U_{nk})) &= \mathcal{N}(\mu_u, \sigma_u^2), \\ p(\text{logit}(d_n)) &= \mathcal{N}(0, 9) \end{aligned}$$

366 where $\mu_u \in \mathbb{R}$ and $\sigma_u \in \mathbb{R}^+$.

367 We use a logistic Normal distribution for our prior over A , a Normal distribution for our prior over B and a
 368 logistic Normal distribution for our prior over σ_{obs} :

$$\begin{aligned} p(\text{logit}(A_{mk})) &= \mathcal{N}(\text{logit}(\text{clip}(\bar{A}_{mk}, a_{\max}, a_{\min})), \sigma_a^2), \\ p(B_{mk}) &= \mathcal{N}(0, \sigma_b^2). \\ p(\log(\sigma_{obs})) &= \mathcal{N}(0, 1), \end{aligned}$$

369 where $\bar{A}_{mk} \in \{0, 1\}$, $a_{\max} \in (0, 1)$, $a_{\min} \in (0, 1)$, $\sigma_a \in \mathbb{R}_{>0}$, $\text{clip}(\bar{A}_{mk}, a_{\max}, a_{\min}) =$
 370 $\max(\min(\bar{A}_{mk}, a_{\max}), a_{\min})$ and $\sigma_b \in \mathbb{R}_{>0}$. \bar{A}_{mk} is given by a pipeline that is used by other methods
 371 such as the Inferelator. The pipeline leverages ATAC-seq and TF binding motif data to provide binary initial
 372 guesses of gene-TF interactions. a_{\max} and a_{\min} are hyperparameters that determine how we clip these binary
 373 values before transforming them to the logit space.

374 For our approximate posterior distribution, we enforce independence as follows:

$$q(U, A, B, \sigma_{obs}, d) = q(U)q(A)q(B)q(\sigma_{obs})q(d).$$

375 We impose the same independence assumptions on each approximate posterior as we do for its corresponding
 376 prior. Specifically, we use the following distributions:

$$\begin{aligned} q(\log(U_{nk})) &= \mathcal{N}(\tilde{U}_{nk}, \tilde{\sigma}_{\tilde{U}_{nk}}^2) \\ q(\text{logit}(d_n)) &= \mathcal{N}(\tilde{d}_n, \tilde{\sigma}_{\tilde{d}_n}^2) \\ q(\text{logit}(A_{mk})) &= \mathcal{N}(\tilde{A}_{mk}, \tilde{\sigma}_{\tilde{A}_{mk}}^2) \\ q(B_{mk}) &= \mathcal{N}(\tilde{B}_{mk}, \tilde{\sigma}_{\tilde{B}_{mk}}^2) \\ q(\log(\sigma_{obs})) &= \mathcal{N}(\tilde{o}, \tilde{\sigma}_{\tilde{o}}^2), \end{aligned}$$

377 where the parameters on the right hand sides of the equations are called variational parameters; $\tilde{U}_{nk}, \tilde{d}_n, \tilde{A}_{mk},$
 378 $\tilde{B}_{mk}, \tilde{o} \in \mathbb{R}$ and $\tilde{\sigma}_{\tilde{U}_{nk}}, \tilde{\sigma}_{\tilde{d}_n}, \tilde{\sigma}_{\tilde{A}_{mk}}, \tilde{\sigma}_{\tilde{B}_{mk}}, \tilde{\sigma}_{\tilde{o}} \in \mathbb{R}^+$. To avoid numerical issues during optimization, we place
 379 constraints on several of these variational parameters.

380 4.2 Inference

381 We perform inference on our model by optimizing the variational parameters to maximize the ELBo. In doing
 382 so, we minimise the KL-divergence between the true posterior and the variational posterior. In practice, to help
 383 with addressing the latent factor identifiability issue, we use a modified version of the ELBo where the prior and
 384 posterior terms are weighted by a constant $\beta \geq 1$ (45):

$$\begin{aligned} &\mathbb{E}_{U,A,B,\sigma_{obs},d \sim q(U,A,B,\sigma_{obs},d)} [\log p(W|U, V = A \odot B, \sigma_{obs}, d) \\ &\quad + \beta(\log p(U, A, B, \sigma_{obs}, d) - \log q(U, A, B, \sigma_{obs}, d))] \end{aligned}$$

385 Inference is carried out using the Adam optimizer with learning rate 0.1 and beta values of 0.9 and 0.99. We clip
 386 gradient norms at a value of 0.0001. We set $a_{\min} = 0.005$, $a_{\max} = 0.995$, $\sigma_b^2 = 1$ and $\mu_u = 0$. We vary σ_a
 387 and σ_u as hyperparameters that control the strengths of the priors over A and U , respectively. We also vary β as
 388 a hyperparameter.

389 We choose a hyperparameter configuration using validation AUPRC as the objective function as well as the
 390 early stopping metric. We hold out hyperparameters for $p(A)$ for a fraction of the genes. We do this by setting
 391 $\tilde{A}_{mk} = 0$ for m corresponding to these genes for all k . During inference we regularly obtain posterior point
 392 estimates for these entries and measure the AUPRC against the original values of these entries as given in the
 393 full prior. This quantity is known as the validation AUPRC.

394 Once we have picked the hyperparameter configuration corresponding to the best validation AUPRC, we perform
 395 inference with this model using the full prior without holding out any information. We use an importance
 396 weighted estimate of the marginal log likelihood as our early stopping criterion:

$$\log p(W) = \log \left(\mathbb{E}_{U,A,B,\sigma_{obs},d \sim q(U,A,B,\sigma_{obs},d)} \left[\frac{p(W|U, A, B, \sigma_{obs}, d)p(U, A, B, \sigma_{obs}, d)}{q(U, A, B, \sigma_{obs}, d)} \right] \right),$$

397 where the expectation is computed using simple Monte Carlo and the \log - \sum -exp trick is used to avoid numerical
 398 issues.

399 4.3 Computing Summary Statistics for the Posterior

400 After training the model, we use \tilde{A} and $\tilde{\sigma}_A$, the variational parameters of $q(A)$, to obtain a mean and a variance
 401 for each entry of A . Since $q(A)$ is logistic normal, it admits no closed form solution for the mean and variance.
 402 We therefore use Simple Monte Carlo i.e. we sample each entry of A several times from its posterior distribution
 403 and then compute the sample mean and sample variance from these samples. We use each mean as a posterior
 404 point estimate of the probability of interaction between a TF and a gene, and its associated variance as a proxy
 405 for the uncertainty associated with this estimate.

406 4.4 Calculating AUPRC

407 The gold standards for the datasets used in this paper do not necessarily perfectly overlap with the genes and
 408 TFs that make up the rows and columns of A as defined by the prior hyperparameters i.e. there may be genes
 409 and TFs in the gold standard with a recorded interaction or lack of interaction, that do not appear in our model at
 410 all because they are not present in the prior. The reverse is also true: the prior may contain genes and TFs that

411 are not in the gold standard. For this reason, we compute the AUPRC using one of two methods: ‘keep all gold
412 standard’ or ‘overlap’, which correspond to evaluating only interactions that are present in the gold standard or
413 only interactions that are present in both the gold standard and the prior/posterior. We present results with ‘keep
414 all gold standard’ AUPRC as the evaluation metric when comparing our model to the Inferelator in Figures 2 and
415 3. For our evaluation of uncertainty calibration (Figure 2 D), we use the overlap AUPRC so that bins containing
416 a lower number of posterior means do not have artificially deflated AUPRCs (see the Evaluating Calibration of
417 Posterior Uncertainty part of the Methods Section for further information).

418 **4.5 Evaluating Calibration of Posterior Uncertainty**

419 We create 10 bins, corresponding to the lowest 10%, 20%, 30% and so on of posterior variances. We place
420 the posterior point estimates of TF-gene interactions associated with these variances into these bins and then
421 calculate the ‘overlap AUPRC’ for each bin using the corresponding gold standard. The AUPRC for each bin is
422 calculated using those interactions that are in the gold standard and also in the bin. We use such a cumulative
423 binning scheme because using a non-cumulative scheme could result in some bins having very small numbers of
424 posterior interactions that are present in the gold standard, which would lead to noisier estimates of the AUPRC.

425 **4.6 Inference and Evaluation on Multiple Observations of W**

426 The Inferelator method applies two scRNA-seq experiments separately on *S. cerevisiae*, with each resulting in a
427 distinct model. These models are used to infer TF-gene interaction matrices, which are then sparsified. The
428 final matrix is obtained by taking the intersection of the two matrices and retaining only the entries that are
429 non-zero in both matrices. In our approach, we also train a separate model on each expression matrix, and obtain
430 a posterior mean matrix for A for each of them. To obtain the final posterior mean matrix for A , we average
431 the posterior mean matrices from each model. While this approach works well, future research could focus on
432 explicitly modeling separate expression matrices within the model, as discussed in the Conclusion section.

433 **4.7 Measuring the Impact of Prior Hyperparameters**

434 We evaluate the utility of each of the prior hyperparameter matrices used in our experiments. In Figures
435 2A and 3A, we present with grey dots the AUPRCs achieved when performing inference using shuffled
436 prior hyperparameters for A . This corresponds to randomly assigning to each row (gene) of A , the prior
437 hyperparameters that correspond to a different row of A . Shuffling the hyperparameters should lead to worse
438 performance, as the posterior estimates should then also be shuffled, whereas the row/column labels for the
439 posterior will remain unshuffled. For the ‘no prior’ setting, shown with black dots in the figures, we set
440 $\bar{A}_{mk} = 0 \forall m, k$. The difference in AUPRC achieved using the unshuffled vs shuffled or no hyperparameters
441 measures the usefulness of the provided hyperparameters for the inference task on the dataset in question.

442 **4.8 Cross-Validation**

443 For each model organism, *S. cerevisiae* and *B. subtilis*, we perform a five-fold cross validation experiment.
444 Cross-validation is performed by partitioning the gold standard into an 80% - 20% split, where 80% of the data
445 represents prior-known information to be used as a prior for $p(A)$, and the remaining 20% is treated as the gold
446 standard for evaluation. This process is repeated five times to generate five random splits of the data in order to
447 robustly evaluate GRN inference. It is important to note that PMF-GRN performs hyperparameter search before
448 inferring a final GRN within each cross-validation split. For each of the five partitioned cross-validation folds
449 the 80%, or prior portion, is further split into 80% train and 20% test for hyperparameter search and evaluation.
450 Once the optimal hyperparameters have been determined, the initial 80% split is treated as the training data,
451 while the remaining 20%, which was not seen during hyperparameter selection, is used for evaluation.

452 **4.9 Datasets and Preprocessing**

453 We inferred each GRN using a single-cell RNA-seq expression matrix, a TF-target gene connectivity matrix, and
454 a gold standard for bench-marking purposes. We modeled the single-cell expression matrices based on the raw
455 UMI counts obtained from sequencing for the *S. cerevisiae* datasets, which were therefore not normalized for
456 the purpose of this work. For the two *B. subtilis* datasets used in this work, we demonstrate the effect of different
457 normalization and scaling techniques, and convert all data used to integers in order to create a single-cell-like
458 dataset. We further obtained binary TF-gene matrices representing prior-known interactions, which served as
459 prior hyperparameters over A , and were derived from the YEASTRACT and subtiwiki databases. We acquired a
460 gold standard for *S. cerevisiae* our datasets from independent work which is detailed below.

461 **Saccharomyces cerevisiae**

462 We used two raw UMI count expression matrices for the organism *S. cerevisiae* obtained from NCBI GEO
463 (GSE125162 (8) and GSE144820 (39)). For this well studied organism, we employed the YEASTRACT (46; 47)
464 literature derived network of TF-target gene interactions to be used as a prior over A in both *S. cerevisiae*
465 networks. A gold standard for *S. cerevisiae* was additionally obtained from a previously defined network (48)
466 and used for bench-marking our posterior network predictions. We note that the gold standard is roughly a reliable
467 subset of the YEASTRACT prior. Additional interactions in the prior can still be considered to be true but have
468 less supportive evidence than those in the gold standard.

469 **Bacillus subtilis**

470 We used two microarray datasets for *B. subtilis*, which we label as B1 (GSE27219) and B2 (GSE67023). Both
471 B1 and B2 underwent different normalization as part of standard microarray processing, described in detail in
472 (40) and (41). For the experiment "No Normalization", B1 was simply converted to integers, while B2 contained
473 negative numbers and had to be scaled and then converted to integers so that the data represented positive integers
474 similar to single-cell data.

475 For this well studied organism, we use the subtiwiki database (42; 43; 44) to obtain a network of prior-known
476 TF-target gene interactions to be used as a prior over A as well as a gold standard for benchmarking posterior
477 predictions.

478 **5 Data Availability**

479 The datasets used in this work are publicly available. They are referenced in the Methods section and are
480 available through <https://github.com/nyu-dl/pmf-grn>.

481 **6 Code Availability**

482 Code, inferred GRNs, and inference and evaluation scripts can be found at [https://github.com/nyu-dl/](https://github.com/nyu-dl/pmf-grn)
483 [pmf-grn](https://github.com/nyu-dl/pmf-grn).

484 **7 Author Contributions**

485 CSG and KC contributed to Conceptualization of the project. OM and KC designed the probabilistic model. OM
486 implemented PMF-GRN Software, Experiments and Validation. CSG implemented PMF-GRN Experiments,
487 Validation, and Inferelator Software. OM, CSG, and KC contributed to Methodology, Software, Validation,
488 Formal Analysis, Visualization, and Writing Original Draft Preparation. CSG contributed to Data Curation. KC
489 and RB contributed to Supervision, Project Administration and Funding Acquisition.

490 **8 Acknowledgements**

491 We thank members of the Bonneau lab for insightful discussions and feedback on this manuscript. We also thank
492 the the staff of the NYU IT High Performance Computing and Flatiron Institute Scientific Computing Core. This
493 work was supported by Samsung Advanced Institute of Technology (under the project *Next Generation Deep*
494 *Learning: From Pattern Recognition to AI*); NSF Award 1922658 NRT-HDR: FUTURE Foundations, Transla-
495 tion, and Responsibility for Data Science; the National Institutes of Health (RM1HG011014, R01NS116350,
496 R01NS118183, R01AI130945); and the Simons Foundation.

497 **References**

- 498 [1] Hecker M, Lambeck S, Toepfer S, Van Someren E, Guthke R. Gene regulatory network inference: data
499 integration in dynamic models—a review. *Biosystems*. 2009;96(1):86-103.
- 500 [2] Chai LE, Loh SK, Low ST, Mohamad MS, Deris S, Zakaria Z. A review on the computational approaches
501 for gene regulatory network construction. *Computers in biology and medicine*. 2014;48:55-65.
- 502 [3] Karlebach G, Shamir R. Modelling and analysis of gene regulatory networks. *Nature reviews Molecular*
503 *cell biology*. 2008;9(10):770-80.
- 504 [4] Äijö T, Lähdesmäki H. Learning gene regulatory networks from gene expression measurements using
505 non-parametric molecular kinetics. *Bioinformatics*. 2009;25(22):2937-44.

- 506 [5] Nachman I, Regev A, Friedman N. Inferring quantitative models of regulatory networks from expression
507 data. *Bioinformatics*. 2004;20(suppl_1):i248-56.
- 508 [6] Burdziak C, Azizi E, Prabhakaran S, Pe'er D. A nonparametric multi-view model for estimating cell
509 type-specific gene regulatory networks. *arXiv preprint arXiv:190208138*. 2019.
- 510 [7] Allaway KC, Gabitto MI, Wapinski O, Saldi G, Wang CY, Bandler RC, et al. Genetic and epigenetic
511 coordination of cortical interneuron development. *Nature*. 2021;597(7878):693-7.
- 512 [8] Jackson CA, Castro DM, Saldi GA, Bonneau R, Gresham D. Gene regulatory network reconstruction
513 using single-cell RNA sequencing of barcoded genotypes in diverse environments. *Elife*. 2020;9:e51254.
- 514 [9] Ciofani M, Madar A, Galan C, Sellars M, Mace K, Pauli F, et al. A validated regulatory network for Th17
515 cell specification. *Cell*. 2012;151(2):289-303.
- 516 [10] Ji Z, He L, Regev A, Struhl K. Inflammatory regulatory network mediated by the joint action of NF- κ B,
517 STAT3, and AP-1 factors is involved in many human cancers. *Proceedings of the National Academy of
518 Sciences*. 2019;116(19):9453-62.
- 519 [11] Yosef N, Shalek AK, Gaublot JM, Jin H, Lee Y, Awasthi A, et al. Dynamic regulatory network
520 controlling TH17 cell differentiation. *Nature*. 2013;496(7446):461-8.
- 521 [12] Mercatelli D, Scalambra L, Triboli L, Ray F, Giorgi FM. Gene regulatory network inference re-
522 sources: A practical overview. *Biochimica et Biophysica Acta (BBA)-Gene Regulatory Mechanisms*.
523 2020;1863(6):194430.
- 524 [13] Huynh-Thu VA, Irrthum A, Wehenkel L, Geurts P. Inferring regulatory networks from expression data
525 using tree-based methods. *PloS one*. 2010;5(9):e12776.
- 526 [14] Wang Y, Joshi T, Zhang XS, Xu D, Chen L. Inferring gene regulatory networks from multiple microarray
527 datasets. *Bioinformatics*. 2006;22(19):2413-20.
- 528 [15] Chang C, Ding Z, Hung YS, Fung PCW. Fast network component analysis (FastNCA) for gene regulatory
529 network reconstruction from microarray data. *Bioinformatics*. 2008;24(11):1349-58.
- 530 [16] Dufva M. Introduction to microarray technology. *DNA Microarrays for Biomedical Research: Methods
531 and Protocols*. 2009:1-22.
- 532 [17] Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. *Nature reviews
533 genetics*. 2009;10(1):57-63.
- 534 [18] Saliba AE, Westermann AJ, Gorski SA, Vogel J. Single-cell RNA-seq: advances and future challenges.
535 *Nucleic acids research*. 2014;42(14):8845-60.
- 536 [19] Akers K, Murali T. Gene regulatory network inference in single-cell biology. *Current Opinion in Systems
537 Biology*. 2021;26:87-97.
- 538 [20] Lähmann D, Köster J, Szczurek E, McCarthy DJ, Hicks SC, Robinson MD, et al. Eleven grand challenges
539 in single-cell data science. *Genome biology*. 2020;21(1):1-35.
- 540 [21] Chen G, Ning B, Shi T. Single-cell RNA-seq technologies and related computational data analysis.
541 *Frontiers in genetics*. 2019:317.
- 542 [22] Ochs MF, Fertig EJ. Matrix factorization for transcriptional regulatory network inference. In: 2012 IEEE
543 Symposium on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB). IEEE;
544 2012. p. 387-96.
- 545 [23] Alter O, Brown PO, Botstein D. Singular value decomposition for genome-wide expression data processing
546 and modeling. *Proceedings of the National Academy of Sciences*. 2000;97(18):10101-6.
- 547 [24] Moloshok TD, Klevecz R, Grant JD, Manion FJ, Speier IV W, Ochs MF. Application of Bayesian
548 decomposition for analysing microarray data. *Bioinformatics*. 2002;18(4):566-75.
- 549 [25] Kim PM, Tidor B. Subsystem identification through dimensionality reduction of large-scale gene expression
550 data. *Genome research*. 2003;13(7):1706-18.
- 551 [26] Brunet JP, Tamayo P, Golub TR, Mesirov JP. Metagenes and molecular pattern discovery using matrix
552 factorization. *Proceedings of the national academy of sciences*. 2004;101(12):4164-9.

- 553 [27] Gao Y, Church G. Improving molecular cancer class discovery through sparse non-negative matrix
554 factorization. *Bioinformatics*. 2005;21(21):3970-5.
- 555 [28] Yang Z, Michailidis G. A non-negative matrix factorization method for detecting modules in heterogeneous
556 omics multi-modal data. *Bioinformatics*. 2016;32(1):1-8.
- 557 [29] Duren Z, Chen X, Zamanighomi M, Zeng W, Satpathy AT, Chang HY, et al. Integrative analysis of
558 single-cell genomics data by coupled nonnegative matrix factorizations. *Proceedings of the National
559 Academy of Sciences*. 2018;115(30):7723-8.
- 560 [30] Hu X, Hu Y, Wu F, Leung RWT, Qin J. Integration of single-cell multi-omics for gene regulatory network
561 inference. *Computational and Structural Biotechnology Journal*. 2020;18:1925-38.
- 562 [31] Skok Gibbs C, Jackson CA, Saldi GA, Tjärnberg A, Shah A, Watters A, et al. High-performance single-cell
563 gene regulatory network inference at scale: the Inferelator 3.0. *Bioinformatics*. 2022.
- 564 [32] Jansen C, Ramirez RN, El-Ali NC, Gomez-Cabrero D, Tegner J, Merckenschlager M, et al. Building
565 gene regulatory networks from scATAC-seq and scRNA-seq using linked self organizing maps. *PLoS
566 computational biology*. 2019;15(11):e1006555.
- 567 [33] Van de Sande B, Flerin C, Davie K, De Waegeneer M, Hulselmans G, Aibar S, et al. A scalable SCENIC
568 workflow for single-cell gene regulatory network analysis. *Nature Protocols*. 2020;15(7):2247-76.
- 569 [34] Kamimoto K, Stringa B, Hoffmann CM, Jindal K, Solnica-Krezel L, Morris SA. Dissecting cell identity
570 via network inference and in silico gene perturbation. *Nature*. 2023:1-10.
- 571 [35] Äijö T, Bonneau R. Biophysically motivated regulatory network inference: progress and prospects. *Human
572 heredity*. 2016;81(2):62-77.
- 573 [36] Mnih A, Salakhutdinov RR. Probabilistic matrix factorization. *Advances in neural information processing
574 systems*. 2007;20.
- 575 [37] Blei DM, Kucukelbir A, McAuliffe JD. Variational inference: A review for statisticians. *Journal of the
576 American statistical Association*. 2017;112(518):859-77.
- 577 [38] Ranganath R, Gerrish S, Blei D. Black box variational inference. In: *Artificial intelligence and statistics*.
578 PMLR; 2014. p. 814-22.
- 579 [39] Jariani A, Vermeersch L, Cerulus B, Perez-Samper G, Voordeckers K, Van Brussel T, et al. A new protocol
580 for single-cell RNA-seq reveals stochastic gene expression during lag phase in budding yeast. *elife*.
581 2020;9:e55320.
- 582 [40] Nicolas P, Mäder U, Dervyn E, Rochat T, Leduc A, Pigeonneau N, et al. Condition-dependent transcriptome
583 reveals high-level regulatory architecture in *Bacillus subtilis*. *Science*. 2012;335(6072):1103-6.
- 584 [41] Arrieta-Ortiz ML, Hafemeister C, Bate AR, Chu T, Greenfield A, Shuster B, et al. An experimentally
585 supported model of the *Bacillus subtilis* global transcriptional regulatory network. *Molecular systems
586 biology*. 2015;11(11):839.
- 587 [42] Michna RH, Zhu B, Mäder U, Stülke J. Subti Wiki 2.0—an integrated database for the model organism
588 *Bacillus subtilis*. *Nucleic acids research*. 2016;44(D1):D654-62.
- 589 [43] Zhu B, Stülke J. Subti Wiki in 2018: from genes and proteins to functional network annotation of the
590 model organism *Bacillus subtilis*. *Nucleic acids research*. 2018;46(D1):D743-8.
- 591 [44] Pedreira T, Elfmann C, Stülke J. The current state of Subti Wiki, the database for the model organism
592 *Bacillus subtilis*. *Nucleic Acids Research*. 2022;50(D1):D875-82.
- 593 [45] Higgins I, Matthey L, Pal A, Burgess C, Glorot X, Botvinick M, et al. beta-VAE: Learning Basic
594 Visual Concepts with a Constrained Variational Framework. In: *International Conference on Learning
595 Representations*; 2017. Available from: <https://openreview.net/forum?id=Sy2fzU9gl>.
- 596 [46] Monteiro PT, Oliveira J, Pais P, Antunes M, Palma M, Cavalheiro M, et al. YEASTRACT+: a portal
597 for cross-species comparative genomics of transcription regulation in yeasts. *Nucleic acids research*.
598 2020;48(D1):D642-9.
- 599 [47] Teixeira MC, Monteiro PT, Palma M, Costa C, Godinho CP, Pais P, et al. YEASTRACT: an upgraded
600 database for the analysis of transcription regulatory networks in *Saccharomyces cerevisiae*. *Nucleic acids
601 research*. 2018;46(D1):D348-53.

- 602 [48] Tchourine K, Vogel C, Bonneau R. Condition-specific modeling of biophysical parameters advances
603 inference of regulatory networks. *Cell reports*. 2018;23(2):376-88.
- 604 [49] Faria JP, Overbeek R, Taylor RC, Conrad N, Vonstein V, Goelzer A, et al. Reconstruction of the regulatory
605 network for *Bacillus subtilis* and reconciliation with gene expression data. *Frontiers in Microbiology*.
606 2016;7:275.

607 **Supplementary Information: Probabilistic Matrix** 608 **Factorization for Gene Regulatory Network Inference**

609 **A Supplementary Tables**

Method	Prior Information		
	Regular	None	Shuffled
PMF-GRN	0.375	0.014	0.023
AmUSR	0.223	0.024	0.019
BBSR	0.402	0.022	0.018
StARS	0.186	0.028	0.017
Scenic	0.014	0.014	0.014
CellOracle	0.383	N/A	0.013

Supplementary Table 1: AUPRCs achieved by PMF-GRN, the Inferelator algorithms (AMuSR, BBSR, and StARS), Scenic and CellOracle on *S. cerevisiae* datasets.

Method	Cross Validation Split				
	Split 1	Split 2	Split 3	Split 4	Split 5
PMF-GRN	0.114	0.096	0.086	0.1342	0.118
BBSR	0.112	0.128	0.161	0.171	0.139
StARS	0.109	0.137	0.154	0.195	0.151
Scenic	0.020	0.021	0.018	0.025	0.021
CellOracle	0.034	0.042	0.034	0.043	0.034

Supplementary Table 2: AUPRCs achieved by PMF-GRN, the Inferelator algorithms (AMuSR, BBSR, and StARS), Scenic and CellOracle on *S. cerevisiae* datasets using the gold standard for 5-fold cross validation.

Method	Noise Added			
	No Noise	100% Noise	250% Noise	500% Noise
PMF-GRN	0.343	0.280	0.198	0.149
BBSR	0.264	0.208	0.186	0.174
StARS	0.136	0.125	0.114	0.118
Scenic	0.075	0.068	0.059	0.055
CellOracle	0.417	0.306	0.226	0.175

Supplementary Table 3: AUPRCs achieved by PMF-GRN, the Inferelator algorithms (AMuSR, BBSR, and StARS), Scenic and CellOracle on *S. cerevisiae* datasets using increasing amounts of noise added to the prior-knowledge data.

610 **B Supplementary Methods**

611 **B.1 TF Target Gene Connectivity Matrix Generation**

612 **B.1.1 *Saccharomyces cerevisiae***

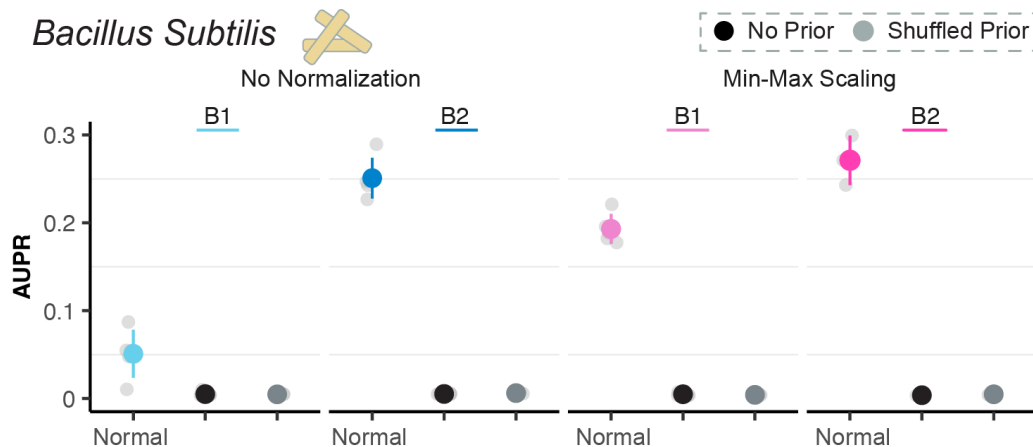
613 Datasets were obtained from (31) without further modification.

614 **B.1.2 *Bacillus subtilis***

615 A prior-known TF-target gene interactions matrix was obtained from the Subtiwiki database (49) from "regula-
616 tions" (downloaded 07/21/22). Using the columns "regulator locus" and "gene locus" a cross-tab integer matrix
617 was created, where 1 represents the existence of an interaction and 0 represents no interaction. This matrix
618 was randomly split 5 times in 80%-20% proportions along the gene axis to generate independent prior-known
619 information and gold standard matrices.

620 To demonstrate the importance of scaling microarray data to place independently collected datasets on the same
621 scale, we demonstrate how Min-Max Scaling improves inference in both *B. subtilis* datasets. For "Min-Max
622 Scaling", both B1 and the positive scaled B2 dataset were subsequently normalized using the following logic.

623 Using the observation axis, values were linearly transformed so that the minimum value was mapped to 0 and the
 624 maximum value was mapped to 1. Each value was then multiplied by 100 and converted to integers to produce
 625 the resulting expression matrix of scaled single-cell-like integers.



Supplemental Figure S1: Results for GRNs learned in *B. subtilis* datasets B1 and B2 comparing "No Normalization" to "Min-Max Scaling". Light gray dots represent the results for each of the 5 cross-validation experiments. Colored dots represent the mean of the cross-validation experiments \pm standard deviation. Negative controls are demonstrated by black dots for "No Prior" and grey dots for "Shuffled Prior".

Method	<i>B. subtilis</i> Cross-Validation Dataset B1		
	Regular	No Prior	Shuffled
No Normalization	0.0509 \pm 0.0273	0.0048 \pm 0.0003	0.0050 \pm 0.0028
Min-Max Scaling	0.1931 \pm 0.0171	0.0042 \pm 0.0003	0.0042 \pm 0.0018

Supplementary Table 4: AUPRCs achieved by PMF-GRN on *B. subtilis* B1 dataset. Results are reported as the mean AUPRC across five 'cross-validation' splits \pm standard deviation

Method	<i>B. subtilis</i> Cross-Validation Dataset B2		
	Regular	No Prior	Shuffled
No Normalization	0.2508 \pm 0.0232	0.0062 \pm 0.0006	0.0052 \pm 0.0004
Min-Max Scaling	0.2886 \pm 0.0312	0.0048 \pm 0.0008	0.0038 \pm 0.0005

Supplementary Table 5: AUPRCs achieved by PMF-GRN on *B. subtilis* B2 dataset. Results are reported as the mean AUPRC across five 'cross-validation' splits \pm standard deviation

626 B.2 Inferelator, Scenic, and CellOracle Networks

627 B.2.1 *Saccharomyces cerevisiae*

628 Networks were inferred using the "multitask" workflow setting of the Inferelator for the same single-cell
 629 *S. cerevisiae* datasets described in (31). For each algorithm, BBSR, StARS, and AMuSR, the following parameters
 630 were used: gold_standard_filter_method="keep_all_gold_standard", num_bootstraps=5. Aggregated multi-task
 631 networks were used for benchmarking, while single-task networks were disregarded for the purpose of this work.
 632 To make these networks directly comparable to PMF, we did not make use of normalization, count minimum, or
 633 meta-data options available within the Inferelator workflow.

634 Networks inferred with Scenic and CellOracle used the same input files, with no additional parameters specified.