

# Aggregating statistically correlated metabolic pathways into groups to improve prediction performance

Abdur Rahman M. A. Basher<sup>1</sup> <sup>a</sup> and Steven J. Hallam<sup>1,2</sup> <sup>b</sup>

<sup>1</sup>Graduate Program in Bioinformatics, University of British Columbia, Vancouver, BC V5Z 4S6, Canada

<sup>2</sup>Department of Microbiology & Immunology, University of British Columbia, Vancouver, BC V6T 1Z3, Canada  
arbasher@student.ubc.ca, shallam@mail.ubc.ca

**Keywords:** Pathway group, Relabeling, Data augmentation, Correlated models, Metabolic pathway prediction, MetaCyc

**Abstract:** Metabolic pathway prediction from genomic sequence information is an essential step in determining the capacity of living things to transform matter and energy at different levels of biological organization. A detailed and accurate pathway map enables researchers to interpret and engineer the flow of biological information from genotype to phenotype in both organismal and multi-organismal contexts. In this paper, we propose two novel hierarchical mixture models, SOAP (sparse correlated pathway group) and SPREAT (distributed sparse correlated pathway group), to improve pathway prediction outcomes. Both models leverage pathway abundance to represent an organismal genome as a mixed distribution of groups, and each group, in turn, is a mixture of pathways. Moreover, both models deal with missing potential pathways in the training set by provisioning supplementary pathways into the learning framework as part of noise reduction efforts. Because the introduction of supplementary pathways may lead to overestimation of some pathways, dual sparseness is applied. The resulting pathway group dataset is then used to train multi-label learning algorithms. Model effectiveness was evaluated on metabolic pathway prediction where correlated models, in particular, SOAP was able to equal or exceed the performance of previous pathway prediction algorithms on organismal genomes.


## 1 INTRODUCTION


Rapid advances in high-throughput sequencing and mass spectrometry over the past two decades have produced a veritable tidal wave of multi-omic information spanning the central dogma of biology (DNA, RNA, protein and metabolites) at the individual, population and community levels of organization (Wang et al., 2015) (Hassa et al., 2018) (Aguilar-Pulido et al., 2016) (Loh et al., 2012). As the ubiquity and abundance of these datasets increases there is a concomitant need to develop bioinformatics applications that scale with data volume and complexity. In particular, methods for predicting metabolic pathways have become essential to interpret and engineer the flow of biological information from genotype to phenotype (Lawson et al., 2019) (Hahn et al., 2016).

A metabolic pathway can be defined as a series of linked chemical reactions occurring within or between cells, often catalyzed and coordinated by a group of enzymes, resulting in metabolic flux from

substrate to product and so on to completion. A variety of rule-based and machine learning prediction methods have been developed to model these pathways in both organismal and multi-organismal contexts (Mascher et al., 2019) (Baranwal et al., 2020) (Yamanishi et al., 2015) (Tabei et al., 2016) (Ye and Doak, 2009) (Dale et al., 2010) (Karp et al., 2016) (M. A. Basher et al., 2020) (M. A. Basher et al., 2021b). While these methods rely on reference metabolic pathway databases (e.g., MetaCyc (Caspi et al., 2019) and KEGG (Kanehisa et al., 2017)) to reconstruct pathways, other computational methods ignore the use of reference database or follow an agnostic approach by ignoring pathway boundaries in the reconstruction process (Zhao et al., 2012) (Qi et al., 2014) (Shafiei et al., 2014) (Jiao et al., 2013).

Among recently developed pathway prediction methods is triUMPF (M. A. Basher et al., 2021b) which uses several layers of interactions among pathways and enzymes within a network to improve the precision of pathway predictions in terms of communities represented by a cluster of nodes (pathways and enzymes). Despite triUMPF's predictive gains, its performance remained error prone because its pre-

<sup>a</sup>  <https://orcid.org/0000-0002-3407-1187>

<sup>b</sup>  <https://orcid.org/0000-0002-4889-6876>

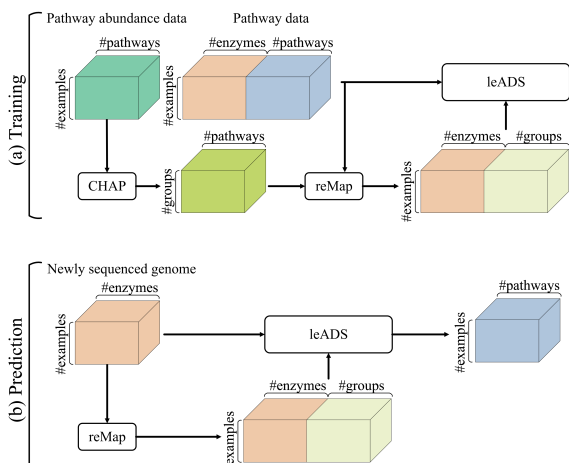


Figure 1: Group-based pathway prediction workflow. The training phase (a) takes pathway abundance data to discover groups using any correlated models in the CHAP package. Then, groups are used to map examples in the pathway abundance data to groups using reMap (Hallam Lab, 2021b). Then, the results of this mapping are used in leADS (Hallam Lab, 2021a), along with pathway data, to learn the model. After training, pathways can be predicted for a newly sequenced genome (b), by first inferring groups using reMap, and then apply the pretrained leADS to predict pathways from groups.

diction process depends on the quality of communities detected from both pathway and enzyme networks that are learned from pathway datasets which contain missing pathway information (M. A. Basher et al., 2020).

Previously, Shafiei and colleagues developed BiomeNet (Shafiei et al., 2014), an extension of MetaNetSim (Jiao et al., 2013), which is a hierarchical Bayesian network to reconstruct metabolic networks in a purely data-driven manner by leveraging enzyme abundances present in multi-organismal datasets. Instead of relying on defined pathway boundaries (Khatri et al., 2012), BiomeNet discovers functions that are referred to as subnetworks, where a subnetwork constitutes a group of connected reactions. Applications of BiomeNet to the human gut microbiome revealed subnetworks that are common among healthy and inflammatory bowel disease (IBD) microbiome patients as well as distinct subnetworks associated with IBD patients.

Inspired by BiomeNet, we developed CHAP (correlated pathway-group) a software package comprised of three correlated mixed-membership hierarchical Bayesian models, CTM (Blei and Lafferty, 2006), SOAP, and SPREAT, to capture mixed components given pathway abundance data. The component is referred to as a “pathway group”, which is comprised of a set of correlated pathways, while path-

ways are permitted to be inter-mixed across groups with different proportions, resulting in overlapping pathways on groups. Modeling explicitly correlations among pathways, using a Gaussian covariance matrix, is fundamental as functions of similar organisms or communities are shared. Moreover, due to noise or missing pathways information, the two novel models: SOAP and SPREAT provision supplementary pathways into the learning framework as part of noise reduction efforts. Because the introduction of supplementary pathways may lead to overestimation of some pathways, dual sparseness is applied where each example in the pathway abundance data is represented by a few focused mixing groups and each pathway group consists of a few relevant pathways. These last two properties were not included in CTM. By modeling examples as mixing groups, one may use results from correlated models for downstream group-based pathway prediction (Fig. 1).

Using CTM, SOAP, and SPREAT, we evaluated groups on metabolic pathway prediction. Resulting pathway group datasets were used to train reMap (Hallam Lab, 2021b) to map examples to groups. For pathway prediction using groups, we applied leADS software (Hallam Lab, 2021a) using the recommended settings discussed in (M. A. Basher and Hallam, 2021). The results were then compared to two heuristic or rule-based pathway prediction algorithms: MinPath (Ye and Doak, 2009) and Pathologic (Karp et al., 2016), and to two machine learning algorithms: mILGPR (M. A. Basher et al., 2020) and triUMPF (M. A. Basher et al., 2021a) on a set of Tier 1 (T1) pathway genome databases (PGDBs) and genomes used in the Critical Assessment of Metagenome Interpretation (CAMI) initiative (Sczyrba et al., 2017) following established benchmarks (M. A. Basher et al., 2020).

## 2 CORRELATED MODELS

In this section, we present three correlated pathway models: i)-CTM (correlated topic model) (Blei and Lafferty, 2006), ii)- SOAP (sparse correlated pathway group) and iii)- SPREAT (distributed sparse correlated pathway group). These models incorporate pathway abundance information to encode each example as a mixture distribution of groups, and each pathway group, in turn, is a mixture of pathways with different mixing proportions. The pathway abundance information can be obtained by mapping enzyme –with abundances– onto a reference pathway database (e.g. MetaCyc (Caspi et al., 2019)). Before we discuss these three models, let us first provide

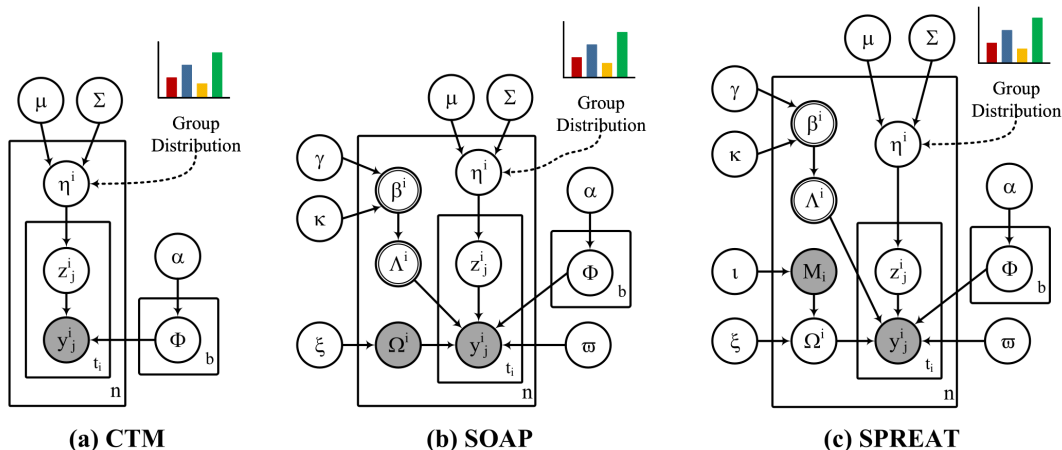


Figure 2: Graphical model representation of the correlated group models. The boxes are “plates” representing replicates. The outer plate represents examples, while the inner plate represents the repeated choice of pathways for an example. The logistic normal distribution, used to model the latent group proportions for an example, captures correlations among groups that are impossible to capture using a single Dirichlet. The observed data for each example  $i$  are a set of annotated pathways  $\mathbf{y}^{(i)}$  and a set of hypothetical pathways  $\mathbf{M}_i$ . The hidden variables are: per-example group proportions  $\eta^{(i)}$ , per-example group selection parameters  $\Lambda^{(i)}$ , per-example hypothetical pathway distributions  $\Omega^{(i)}$ , per-pathway group assignment parameter  $z_j^{(i)}$ , and per-group distribution over pathways  $\Phi_a$ .

some definitions and notations.

**Pathway Abundance Data.** Let  $\mathcal{P} = \{\mathbf{y}^{(i)} : 1 < i \leq n\}$  be a collection of  $n$  examples corresponding organismal or multi-organismal genomes (e.g. *Escherichia coli* K-12), where each example  $\mathbf{y}^{(i)} = (y_1^{(i)}, y_1^{(i)}, \dots, y_t^{(i)})$  is a vector encoding the unnormalized abundance information of pathways and  $t$  is the pathway size. Let  $\mathcal{Y} = \{h_1, h_2, \dots, h_t\}$  be a set of all known metabolic pathways obtained from a reference database (e.g., MetaCyc (Caspi et al., 2019)), and  $\mathcal{Y}_i \subseteq \mathcal{Y}$  corresponds to a subset of true pathways associated with the  $i$ th example.

**Group Modeling.** Given  $\mathcal{P}$ , a pathway group distribution for the  $i$ th example is a multinomial distribution vector, denoted by  $\eta^{(i)}$  of size  $b$  groups, i.e.,  $\{p(\Phi_a | \eta^{(i)})\}_{a=1}^b$ , where  $\Phi_j$  is a multinomial pathway distribution over the group  $j$ , i.e.,  $\{p(y_k | \Phi_j)\}_{k=1}^t$ . The overall goal of group modeling is to discover  $b$  hidden groups for each example.

The definition states that a pathway is distributed over groups, implying group correlation, i.e., if  $l \in \Phi_j \geq 0$  and  $l \in \Phi_k > 0$  then  $\Sigma_{j,k} \neq 0$ , where  $\Sigma \in \mathbb{R}^{b \times b}$  is a group-correlation matrix.

**Group Correlation.** Given  $\mathcal{P}$ , the pairwise group-correlation is defined by a Gaussian covariance matrix, denoted by  $\Sigma \in \mathbb{R}^{b \times b}$ . Each entry  $s_{i,j}$  in  $\Sigma$  characterizes the magnitude of correlation between  $i$  and  $j$  pathway groups, where a larger score indicates both pathway groups are highly correlated.

Missing pathway information in  $\mathcal{P}$  is common in both organismal and multi-organismal contexts due

to errors in open reading frame prediction or annotation as well as unknown protein function. Previously, Hanson and colleagues (Hanson et al., 2014) reported missing a set of potential pathways for the Hawaii Ocean Time-series data (Stewart et al., 2011), such as *tricarboxylic acid cycle (TCA)*. These missing pathways have negative implications in group modeling as  $\mathcal{P}$ , in this case, would be exposed to extreme noise. Although manually incorporating missing pathways to  $\mathcal{P}$  may provide a solution to model groups, this solution has the potential to increase false discovery pending experimental validation. A good compromise would be to record missing pathways in a separate list while keeping the original pathway abundance data intact. Let us denote  $\mathbf{M} \in \mathbb{Z}_{\geq 0}^{n \times t}$  a matrix holding a set of missing pathways where each entry is an integer value indicating the abundance of a pathway for an example. This matrix is called the *background* or the *supplementary* matrix. Now, with these definitions, let us describe the research problem.

**Problem Statement.** Given  $\mathcal{P}$  and  $\mathbf{M}$ , the objective is to recover the group distribution  $\eta$  for each example such that applying group based metabolic pathway prediction would recover more accurate pathways for an organismal or multi-organismal genome.

## 2.1 Correlated Topic Model

The correlated topic model (CTM) is a probabilistic graphical model that extends the generative story of latent Dirichlet allocation (LDA) (Blei et al., 2003) to

```

1 for  $a \in \{1, \dots, b\}$  do
2   Sample a distribution over pathways
    $\Phi_a \sim \text{Dir}(\cdot|\alpha)$ ;
3 for  $i \in \{1, \dots, n\}$  do
4   Draw per example group weight
    $\eta^{(i)} \sim \mathcal{N}(\cdot|\mu, \Sigma)$ ;
5   Draw group proportions  $\theta^{(i)} = \text{softmax}(\eta^{(i)})$ ;
6   for  $j \in \{1, \dots, t^{(i)}\}$  do
7     Sample a group assignment
      $z_j^{(i)} \sim \text{Mult}(\cdot|\theta^{(i)})$ ;
8     Sample a pathway  $y_j^{(i)} \sim \text{Mult}(\cdot|\Phi_{z_j^{(i)}})$ ;

```

**Algorithm 1:** The generative process for CTM given a collection of examples

incorporate correlation among groups (or topics in the original paper). Fig. 2a shows the Bayesian graphical model for CTM using plate notation. Like LDA, CTM is composed of a hierarchical Bayesian mixture model, where features (words in the original paper) are mixed to constitute groups that are assumed to be correlated, modeled by a Gaussian covariance matrix. Note that in this paper, we use the terms *feature* and *pathway* interchangeably.

Formally, let  $n$  be the total number of examples in  $\mathcal{P}$ , where each example  $i$  consists of features, i.e.,  $\mathbf{y}^{(i)}$ . Then, the generative process for CTM is described as follows. First, we draw a multinomial feature distribution  $\Phi_a$  from a Dirichlet prior  $\alpha > \mathbb{R}_{>0}$  for each group  $a \in \{1, \dots, b\}$ . Then, for each example  $i$ , a Gaussian random variable is drawn  $\eta^{(i)} \sim \mathcal{N}(\mu, \Sigma)$ , where  $\mu$  is a  $b$  dimensional mean and  $\Sigma \in \mathbb{R}^{b \times b}$  is the covariance matrix. The random variable  $\eta^{(i)}$  is projected onto the probability simplex to obtain the group distributions  $\theta^{(i)} = \text{softmax}(\eta^{(i)})$ , corresponding the logistic-normal distribution, from which a group indicator  $z_j^{(i)} \in \{1, \dots, b\}$  is sampled. Finally, each observed feature  $j \in \{1, \dots, t^{(i)}\}$  is drawn from the associated feature distribution, indicated by its group assignment  $z_j$ , i.e.,  $y_j^{(i)} \sim \Phi_{z_j^{(i)}}$ . This generative process (Algorithm 1) is identical to LDA except that the group distributions is sampled from the logistic normal instead from a Dirichlet prior as in LDA.

## 2.2 Correlated Pathway-Group Model

Correlated pathway group models are extensions to CTM: i)- SOAP (Fig. 2b) and ii)- SPREAT (Fig. 2c). Both models incorporate dual sparseness and supplementary pathways in modeling group proportions. The two properties were not implemented in CTM. Let us discuss these two models.

```

1 for  $a \in \{1, \dots, b\}$  do
2   Sample a distribution over pathways
    $\Phi_a \sim \text{Dir}(\cdot|\alpha)$ ;
3 for  $i \in \{1, \dots, n\}$  do
4   Draw per example group weight
    $\eta^{(i)} \sim \mathcal{N}(\cdot|\mu, \Sigma)$ ;
5   Draw group proportions  $\theta^{(i)} = \text{softmax}(\eta^{(i)})$ ;
6   Draw beta distribution  $\beta^{(i)} \sim \text{Beta}(\cdot|\gamma, \kappa)$ ;
7   Draw a sparsity indicator vector
    $\Lambda^{(i)} \sim \text{Bernoulli}(\cdot|\beta^{(i)})$ ;
8   if SPREAT then
9     Sample a vector  $\mathbf{M}_i \sim \text{Prior}(\cdot|1)$ ;
10    Sample background distribution
     $\Omega^{(i)}|\mathbf{M}_i \sim \text{Dir}(\cdot|\xi)$ ;
11  else
12    Draw background feature proportions
     $\Omega^{(i)} \sim \text{Dir}(\cdot|\xi)$ ;
13  for  $j \in \{1, \dots, t^{(i)}\}$  do
14    Sample a group assignment
     $z_j^{(i)} \sim \text{Mult}(\cdot|\Lambda^{(i)} \odot \theta^{(i)})$ ;
15    Sample a pathway
     $y_j^{(i)} \sim \text{Mult}(\cdot|(1 - \Omega_{z_j^{(i)}}^{(i)}) \odot \Phi_{z_j^{(i)}})$ ;

```

**Algorithm 2:** The generative process for SOAP and SPREAT

Analogous to CTM, given  $n$  number of examples in  $\mathcal{P}$  and a matrix encoding missing pathways  $\mathbf{M}$ , the generative process for SOAP and SPREAT can be described as follows. First, we draw a multinomial pathway distribution  $\Phi_a$  from asymmetric Dirichlet prior  $\alpha \in \mathbb{R}_{>0}$  for each group  $a \in \{1, \dots, b\}$ , where  $b$  is assumed to be known and fixed in advance. The symmetric assumption is appropriate, in such a scenario, because our prior knowledge, associated with these pathways, is inaccessible. For each example  $i$ , a group proportion is drawn  $\theta^{(i)} = \text{softmax}(\eta^{(i)})$ , where  $\eta^{(i)}$  is a Gaussian random variable with mean and covariance are denoted by  $\mu$  and  $\Sigma$ , respectively.

To sample a group, it is reasonable to expect that: i)- each example is usually explained with a handful set of a mixed proportion of groups and ii)- a group should consist only of a few related pathways. Therefore, we apply dual sparsity (Lin et al., 2014) (Airoldi et al., 2008) (Bien and Tibshirani, 2011) (He et al., 2017) to retain those relevant focused groups and pathways by: i)- introducing an auxiliary Bernoulli variable  $\Lambda^{(i)}$  of size  $b$  to determine whether a group is selected for the  $i$ th example or ignored and ii)- applying a cutoff threshold to keep top  $k \ll t$  pathways, based on their probabilities, for each group. Instead of sampling each entry in  $\Lambda^{(i)}$  directly from a Bernoulli coin toss, we assume that each entry is sampled from

Table 1: Correspondence between variational and original parameters.

Original parameter	$\Phi$	$\mu$	$\Sigma$	$\Lambda$	$\Omega$	$z$
Variational parameter	$\phi$	$\mathbf{v}$	$\zeta^2$	$\lambda$	$\omega$	$\zeta$

a Beta distribution  $\beta^{(i)}$ , parameterized by two hyperparameters  $\gamma \in \mathbb{R}_{>0}$  and  $\kappa \in \mathbb{R}_{>0}$ . Applying this dual sparsity, we aim to enhance the interpretability of the learned pathway groups while reducing the negative correlation among groups on  $\Sigma$ .

Next, a group indicator  $z_j^{(i)} \in \{1, \dots, b\}$  is drawn according to the example-specific mixture proportion  $\Lambda^{(i)} \odot \theta^{(i)}$ , where  $\odot$  represents the Hadamard product. Now each pathway  $y_j^{(i)}$  for the  $i$ th example is generated from a weighted distribution  $\Omega_j^{(i)} \odot \Phi_{z_j^{(i)}}$  using a smoothing prior  $\omega \in \mathbb{R}_{>0}$ . The parameter  $\Omega^{(i)} \in \mathbb{R}^t$ , derived from  $\mathbf{M}_i$ , represents a normalized supplementary pathway of size  $t$ , which is assumed to be drawn from a symmetric Dirichlet prior  $\xi \in \mathbb{R}_{>0}$ . For SPREAT, this parameter encodes distribution, where each element of  $\Omega_j^{(i)}$  corresponds to the pathway probability  $y_j^{(i)} \in \mathbf{M}_i$  for  $i$ th example. Here, the background pathway is assumed to be drawn from a sparse binary vector prior  $\mathbf{1} \in \mathbb{R}_{>0}$  that is included for completeness because pathways in  $\mathbf{M}$  for each example are known.

Representing SOAP and SPREAT as layer-wise mixing components supports the hierarchical modularity of metabolic pathway generation, where the components of one level (e.g., pathways) permit to contribute to groups with different degrees of granularity. The generative process of SOAP and SPREAT models is summarized in Algorithm 2. Notice that by setting all entries in  $\Omega$ ,  $\Lambda$ , and  $\omega$  to 1, SOAP and SPREAT are reduced to CTM (“collapse2ctm” or c2m), reflecting the flexibility of these models.

### 3 INFERENCE AND PARAMETER ESTIMATION FOR SPREAT

Here, we discuss the inference for the SPREAT model. A similar expression is derived for SOAP. Given  $\mathcal{P}$ , the goal of inference is to compute the posterior distribution of per-example group proportions ( $\eta$ ), per-example group selection parameters ( $\Lambda$ ) and the associated beta distributions ( $\beta$ ), per-example background pathway distributions ( $\Omega$ ), per-pathway group assignment ( $z$ ), and per-group distribution over pathways ( $\Phi$ ). By denoting all parameters as  $\Theta$  and variables as  $\mathbf{V}$  while omitting hyperparameters, we apply the Jensen’s inequality on a variational distri-

bution over hidden variables  $q(\Theta, \mathbf{V})$  to obtain the evidence lower bound (ELBO) as:

$$\mathcal{L}(q) = \mathbb{E}_q[\log p(\mathbf{Y}, \mathbf{M}, \Theta, \mathbf{V})] + \mathbb{H}(q) \quad (3.1)$$

where  $p(\mathbf{Y}, \mathbf{M}, \Theta, \mathbf{V})$  represents the joint distribution of all observed and latent variables of the model. The ELBO contains two terms. The first term,  $\mathbb{E}_q[\log p(\mathbf{Y}, \mathbf{M}, \Theta, \mathbf{V})]$ , captures how well  $q(\Theta, \mathbf{V})$  describes a distribution of the model. The second term is the entropy of the variational distribution,  $\mathbb{E}_q[-\log q(\Theta, \mathbf{V})]$ , which protects the variational distribution from “overfitting”. The two terms depends on  $q(\Theta, \mathbf{V})$  which is defined as:

$$q(\Theta, \mathbf{V}) = \prod_{a=1}^b q(\Phi_a | \phi_a) \left[ \prod_{i=1}^n q(\eta^{(i)} | \mathbf{v}, \zeta^2) \times q(\Lambda^{(i)} | \lambda^{(i)}) q(\Omega^{(i)} | \omega^{(i)}) \prod_{j=1}^{j=t_i} q(z_j^{(i)} | \zeta_j^{(i)}) \right] \quad (3.2)$$

where  $\phi, \mathbf{v}, \zeta^2, \lambda, \omega$  and  $\zeta$  are variational free parameters. Table 1 shows the correspondence between variational and the original parameters. Now, the first term in Eq. 3.1 is decomposed into:

$$\begin{aligned} \mathbb{E}_q[\log p(\mathbf{Y}, \mathbf{M}, \Theta, \mathbf{V})] &= \sum_{a=1}^{a=b} \mathbb{E}_q[\log p(\Phi_a | \alpha)] \\ &+ \sum_{i=1}^{i=n} \left( \mathbb{E}_q[\log p(\eta | \mu, \Sigma)] + \mathbb{E}_q[\log p(\Lambda^{(i)} | \beta^{(i)})] \right. \\ &+ \mathbb{E}_q[\log p(\beta^i | \gamma, \kappa)] + \mathbb{E}_q[\log p(\Omega^{(i)} | \mathbf{M}^{(i)}, \xi)] \\ &+ \sum_{j=1}^{j=t_i} \left( \mathbb{E}_q[\log p(y_j^{(i)} | z_j^{(i)}, \Omega_j^{(i)}, \Lambda^{(i)}, \Phi, \omega)] \right. \\ &\left. \left. + \mathbb{E}_q[p(z_j^{(i)} | \eta)] \right) \right) \end{aligned} \quad (3.3)$$

where,

$$\begin{aligned} \mathbb{E}_q[\log p(\Phi_a | \alpha)] &= \log \Gamma\left(\sum_{j=1}^{j=t} \alpha_j\right) - \sum_{j=1}^{j=t} \log \Gamma(\alpha_j) \\ &+ \sum_{j=1}^{j=t} (\alpha_j - 1) \mathbb{E}_q[\log \Phi_{a,j}] \end{aligned}$$

$$\begin{aligned} \mathbb{E}_q[\log p(\eta | \mu, \Sigma)] &= \frac{1}{2} \log |\Sigma^{-1}| - \frac{b}{2} \log 2\pi \\ &- \frac{1}{2} \left( \text{tr}(\text{diag}(\zeta^2) \Sigma^{-1}) \right. \\ &\left. + (\mathbf{v} - \mu)^\top \Sigma^{-1} (\mathbf{v} - \mu) \right) \end{aligned}$$

$$\begin{aligned} \mathbb{E}_q[\log p(\Lambda^{(i)} | \beta^{(i)})] &= \sum_{a=1}^{a=b} \left( \lambda_a^{(i)} \log \beta_a^{(i)} + (1 - \lambda_a^{(i)}) \right. \\ &\left. \times \log(1 - \beta_a^{(i)}) \right) \end{aligned}$$

$$\begin{aligned}
\mathbb{E}_q[\log p(\beta^{(i)}|\gamma, \kappa)] &= \sum_{a=1}^{a=b} \left( (\gamma-1) \log(\beta_a^{(i)}) \right. \\
&\quad \left. + (\kappa-1) \log(1-\beta_a^{(i)}) - \log(B(\gamma, \kappa)) \right) \\
\mathbb{E}_q[\log p(\Omega_i|\mathbf{M}^{(i)}, \xi)] &= \log \Gamma \left( \sum_{j=1}^{j=t} \xi_j + \mathbf{M}_j^{(i)} \right) \\
&\quad - \sum_{j=1}^{j=t} \log \Gamma(\xi_j + \mathbf{M}_j^{(i)}) \\
&\quad + \sum_{j=1}^{j=t} (\xi_j + \mathbf{M}_j^{(i)} - 1) \mathbb{E}_q[\log \Omega_j^{(i)}] \\
\mathbb{E}_q[\log p(y_j^{(i)}|z_j^{(i)}, \Omega_j^{(i)}, \Lambda^{(i)}, \Phi, \omega)] &= \log \omega \\
&\quad + \sum_{c=1}^{c=t} \sum_{a=1}^{a=b} \left( y_{j,c}^{(i)} \zeta_{a,j}^{(i)} \lambda_a^{(i)} \mathbb{E}_q[(1-\Omega_c^{(i)})] \mathbb{E}_q[\log \Phi_{a,j}] \right) \\
\mathbb{E}_q[\log p(z_j^{(i)}|\eta)] &\approx 1 - \log \rho + \sum_{a=1}^{a=b} \nu_a \zeta_{a,j}^{(i)} \\
&\quad - \left( \sum_{k=1}^{k=b} \mathbb{E}_q[\exp(\eta_k)] \right) \rho^{-1}
\end{aligned}$$

The second term  $\mathbb{H}(q)$  in Eq. 3.1 has the following parametric forms (see Eq. 3.2):

$$\begin{aligned}
\mathbb{H}(q) &= - \sum_{a=1}^{a=b} \mathbb{E}_q[\log q(\Phi_a|\phi_a)] - \sum_{i=1}^{i=n} \left( \mathbb{E}_q[\log q(\eta^{(i)}|\nu, \zeta^2)] \right. \\
&\quad + \mathbb{E}_q[\log q(\Lambda^{(i)}|\lambda^{(i)})] + \mathbb{E}_q[\log q(\Omega^{(i)}|\omega^{(i)})] \\
&\quad \left. + \sum_{j=1}^{j=t} \mathbb{E}_q[\log q(z_j^{(i)}|\zeta_j^{(i)})] \right)
\end{aligned} \tag{3.4}$$

where,

$$\begin{aligned}
\mathbb{E}_q[\log q(\Phi_a|\phi_a)] &= \log \Gamma \left( \sum_{j=1}^{j=t} \phi_{a,j} \right) - \sum_{j=1}^{j=t} \log \Gamma(\phi_{a,j}) \\
&\quad + \sum_{j=1}^{j=t} (\phi_{a,j} - 1) \mathbb{E}_q[\log \Phi_{a,j}] \\
\mathbb{E}_q[\log q(\eta^{(i)}|\nu, \zeta^2)] &= - \sum_{a=1}^{a=b} \frac{1}{2} \left( \log \zeta_a^2 + \log(2\pi) + 1 \right) \\
\mathbb{E}_q[\log q(\Lambda^{(i)}|\lambda^{(i)})] &= \sum_{a=1}^{a=b} \left( \lambda_a^{(i)} \log \lambda_a^{(i)} \right. \\
&\quad \left. + (1-\lambda_a^{(i)}) \log(1-\lambda_a^{(i)}) \right) \\
\mathbb{E}_q[\log q(\Omega^{(i)}|\omega^{(i)})] &= \log \Gamma \left( \sum_{j=1}^{j=t} \omega_j^{(i)} \right) - \sum_{j=1}^{j=t} \log \Gamma(\omega_j^{(i)}) \\
&\quad + \sum_{j=1}^{j=t} (\omega_j^{(i)} - 1) \mathbb{E}_q[\log \Omega_j^{(i)}] \\
\mathbb{E}_q[\log q(z_j^{(i)}|\zeta_j^{(i)})] &= \mathbb{E}_q \left[ \log \prod_{a=1}^{a=b} (\zeta_{a,j}^{(i)})^{z_{a,j}^{(i)}} \right] = \sum_{a=1}^{a=b} \zeta_{a,j}^{(i)} \log \zeta_{a,j}^{(i)}
\end{aligned}$$

The exceptions that correspond to the above equations are:

$$\begin{aligned}
\mathbb{E}_q[\log \Phi_{a,j}] &= \left( \Psi(\phi_{a,j}) - \Psi \left( \sum_{k=1}^{k=t} \phi_{a,k} \right) \right) \\
\mathbb{E}_q[\log \Omega_j^{(i)}] &= \left( \Psi(\omega_j^{(i)}) - \Psi \left( \sum_{k=1}^{k=t} \omega_k^{(i)} \right) \right) \\
\mathbb{E}_q[(1-\Omega_c^{(i)})] &= \frac{1-\omega_c^{(i)}}{\sum_{k=1}^{k=t} (1-\omega_k^{(i)})} \\
\mathbb{E}_q[\exp(\eta_k)] &= \exp(\nu_a + \frac{1}{2} \zeta_a^2) \\
B(\gamma, \kappa) &= \frac{\Gamma(\gamma)\Gamma(\kappa)}{\Gamma(\gamma+\kappa)}
\end{aligned}$$

where  $\Gamma$  denotes the Gamma function while  $\Psi$  is the logarithmic derivative of the Gamma function.

After expanding both terms in Eq. 3.1, we can now maximize the bound in Eq. 3.1 with respect to each variational parameters using mini-batch coordinate ascent updates (Hoffman et al., 2013) as:

**Optimize  $\zeta$ .** The analytical expression of the variational group assignment  $q(\zeta)$  for each pathway  $j$  and group  $a$  for the  $i$ th example is not amenable due to the non-conjugacy of logistic-normal with latent variables. Instead, we approximate the solution as:

$$\begin{aligned}
\zeta_{a,j}^{(i)} &\propto \exp \left( \sum_{c=1}^{c=t} y_{j,c}^{(i)} \lambda_a^{(i)} \frac{1-\omega_c^{(i)}}{\sum_{k=1}^{k=t} (1-\omega_k^{(i)})} \left( \Psi(\phi_{a,j}) \right. \right. \\
&\quad \left. \left. - \Psi \left( \sum_{k=1}^{k=t} \phi_{a,k} \right) + \nu_a - 1 \right) \right)
\end{aligned} \tag{3.5}$$

where  $\Psi(\cdot)$  is the digamma function. Notice that the variational parameter  $\omega_*^{(i)}$  acts as a smoothing parameter to selecting groups for each pathways, either from  $\mathbf{M}_j$  or from  $\mathcal{L}$ .

**Optimize  $\nu$ .** Collecting terms in the ELBO bound that contain only  $\nu$  and taking derivatives w.r.t.  $\nu_a$  for each group  $a$ , we obtain:

$$\begin{aligned}
\frac{\partial \mathcal{L}(q)|_{\nu}}{\partial \nu_a} &= -\Sigma^{-1}(\nu - \mu) + \sum_{j=1}^{j=t} \zeta_{a,j}^{(i)} \\
&\quad - \left( \exp(\nu_a + \frac{1}{2} \zeta_a^2) \right) t_i \rho^{-1}
\end{aligned} \tag{3.6}$$

where  $\rho$  is another variational parameter, as in CTM (Blei and Lafferty, 2006). The above equation is hard to optimize, instead, we use a conjugate gradient algorithm.

**Optimize  $\zeta^2$ .** By symmetry, we gather all the terms that has  $\zeta^2$  from Eq. 3.1, and take derivatives w.r.t.  $\zeta_a^2$  for each group  $a$  to obtain:

$$\frac{\partial \mathcal{L}(q)|_{\zeta^2}}{\partial \zeta_a^2} = -\frac{1}{2} \left( \Sigma_{a,a}^{-1} + t_i \rho^{-1} \exp \left( \nu_a + \frac{1}{2} \zeta_a^2 \right) - \frac{1}{\zeta_a^2} \right) \tag{3.7}$$

```

1 Initialize  $\phi, v, \zeta^2, \lambda, \omega, \zeta, \gamma, \kappa, \xi, \alpha, \Theta, \mathbf{t}, s = 0,$ 
    $l \geq 0, g \in (0.5, 1]$ 
2 repeat
3    $s = s + 1;$ 
4   example a minibatch randomly  $\mathcal{B} \subset \mathcal{P};$ 
5   for  $i \in \mathcal{B}$  do
6     repeat
7       Update  $\zeta^{(i)}$  with Eq. 3.5;
8       Update  $v^{(i)}$  with Eq. 3.6 using
       conjugate gradient algorithm;
9       Update  $\zeta^{2,(i)}$  with Eq. 3.7 using
       Newton's method;
10      Update  $\rho^{(i)}$  with Eq. 3.8;
11      Update  $\omega^{(i)}$  with Eq. 3.9;
12      Update  $\lambda^{(i)}$  with Eq. 3.10;
13    until local variational parameters
       converge;
14    Compute optimal values  $\mu = \frac{v}{|\mathcal{B}|},$ 
        $\Sigma = \text{diag}(\frac{\zeta^2}{|\mathcal{B}|}) + \mu\mu^\top;$ 
15    Compute global optimal values  $\phi$  with Eq.
       3.11;
16    Update the current estimate of the global
       variational parameters,
        $x = (1 - \tau)x + \tau x,$  where  $x \in \{\phi, \mu, \Sigma\};$ 
17    Update the learning rate  $\tau = (s + l)^{-g};$ 
18 until global convergence criterion is satisfied;

```

**Algorithm 3:** Minibatch variational inference for SPREAT

Again, there is no analytical solution to the above formula. Instead, we use Newton's method for each coordinate such that  $\zeta_a \in \mathbb{R}_{>0}$ .

**Optimize  $\rho$ .** We extract terms involved with the variational parameter  $\rho$ , and equating it's derivative to zero, we get:

$$\rho = \sum_{k=1}^{k=b} \exp(v_k + \frac{1}{2} \zeta_k^2) \quad (3.8)$$

**Optimize  $\omega$ .** We next isolate only the terms in the bound that contain variational background pathway distributions  $q(\omega)$ . However, setting it's derivatives to zero does not lead to a closed-form solution, instead, we approximate  $\omega_c^{(i)}$  for each example  $i$  according to:

$$\omega_c^{(i)} \propto \xi_c + \mathbf{M}_c^{(i)} - \left( \frac{1 - \omega_c^{(i)} - \sum_{k=1}^{k=t} (1 - \omega_k^{(i)})}{(\sum_{k=1}^{k=t} (1 - \omega_k^{(i)}))^2} \right) \times \sum_{j=1}^{j=t} \sum_{a=1}^{a=b} y_{j,c}^{(i)} \zeta_{a,j}^{(i)} \lambda_a^{(i)} \left( \Psi(\phi_{a,j}) - \Psi(\sum_{k=1}^{k=t} \phi_{a,k}) \right) \quad (3.9)$$

**Optimize  $\lambda$ .** To optimize  $\lambda$ , we use the canonical parameterisation of the Bernoulli distribution to get the following updates for each group  $a$  for each example:

$$\lambda_a^{(i)} = \frac{1}{1 + \exp(-(\log(\beta_a^{(i)}) - \log(1 - \beta_a^{(i)}))} \quad (3.10)$$

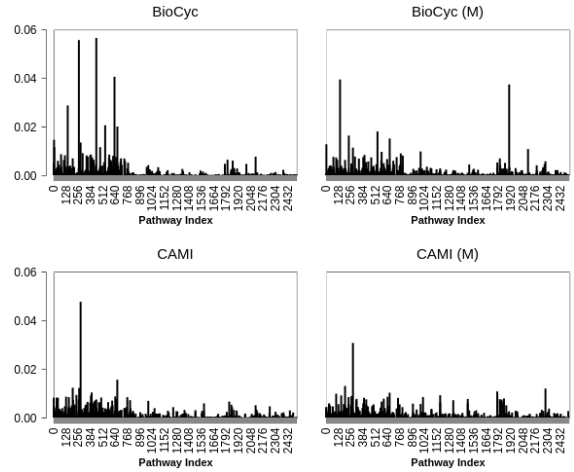


Figure 3: Pathway frequency (averaged on all examples) in BioCyc (v20.5 T2 &3) and CAMI data, and their background pathways, indicated by  $\mathbf{M}$ .

**Optimize  $\phi$ .** Finally, the optimal solution of the variational pathway distribution  $q(\Phi_a | \phi_a)$  for each group  $a$  is obtained by isolating terms involved in the ELBO bound in Eq. 3.1 and setting it's gradient to zero:

$$\phi_{a,c} = \alpha_c + \sum_{i=1}^{i=n} \sum_{j=1}^{j=t_i} y_{j,c}^{(i)} \zeta_{a,j}^{(i)} \lambda_a^{(i)} \frac{1 - \omega_c^{(i)}}{\sum_{k=1}^{k=t} (1 - \omega_k^{(i)})} \quad (3.11)$$

The variational inference algorithm samples a mini-batch from a collection, and uses it to compute the local latent parameters in Eqs 3.5, 3.6, 3.7, 3.8, 3.9, and 3.10 until the evidence lower bound in Eq. 3.1 converges. Then, the global variational parameter  $\phi$  is updated in Eq. 3.11 using the posteriors ( $\beta, \Lambda, \eta, z, \Omega$ ) collected from the previous step after being scaled according to the learning rate  $\tau = (s + l)^{-g}$ , where  $s$  is the current step,  $l \geq 0$  is the delay factor, and  $g \in (0.5, 1]$  is the forgetting rate. This process for SPREAT is summarized in Algorithm 3.

### 3.1 Posterior Predictive Distribution

The posterior predictive distribution estimates the distribution of an unobserved value ( $\tilde{\mathbf{y}}$ ) given the observed values ( $\mathbf{Y}_{obs}$ ) and parameters ( $\Theta$  and  $\mathbf{V}$ ) that are trained on a held-out training set (Hoffman et al., 2013). The predictive distribution for SPREAT is:

$$p(\tilde{\mathbf{y}} | \mathbf{Y}_{obs}, \tilde{\mathbf{M}}, \mathbf{M}_{obs}) = \int p(\tilde{\mathbf{y}} | \Theta, \tilde{\mathbf{M}}) p(\Theta | \mathbf{Y}_{obs}, \mathbf{M}_{obs}) d\Theta \approx \sum_{a=1}^{a=b} \left( \eta_a^{(i)} \times \sum_{j=1}^{j=t} (\Phi_{a,j} \times \tilde{\mathbf{y}}_j^{(i)}) \right) \times q(\Theta, \mathbf{V}) \quad (3.12)$$

where  $\tilde{\mathbf{M}}$  is  $\tilde{\mathbf{y}}$ 's background pathways and  $q(\Theta, \mathbf{V})$  corresponds to Eq. 3.2, trained on  $\mathbf{Y}_{obs}$  and  $\mathbf{M}_{obs}$ .

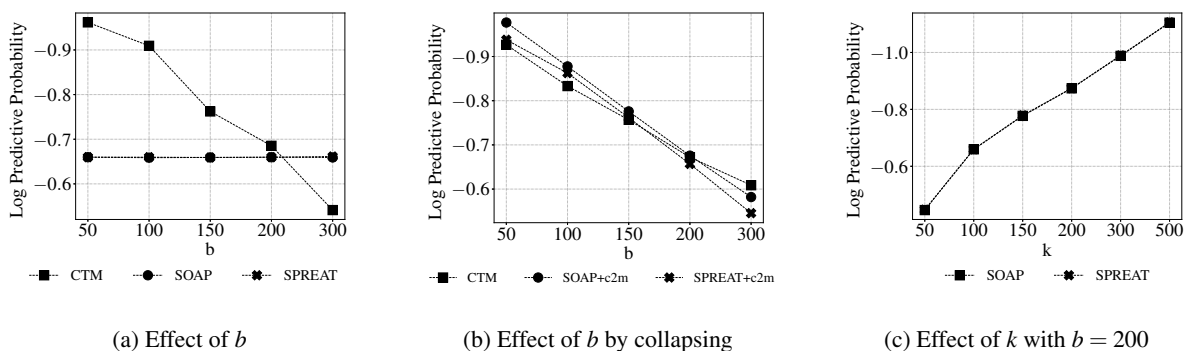


Figure 4: Log predictive distribution on CAMI data. Figs 4a and 4b show the effect of group size  $b$  to the performance of CTM, SOAP, and SPREAT and the collapsed models, respectively. Fig. 4c demonstrates the effect of retaining top  $k$  pathways to the performance SOAP and SPREAT. The performance is measures according to the log predictive probability where higher values indicate better performances.

## 4 EXPERIMENTAL SETTINGS

In this section, we describe the experimental datasets and settings used to validate the performance of the three correlated models. The CHAP package was written in Python v3 and is available under the GNU license at [github.com/hallamlab/chap](https://github.com/hallamlab/chap). All tests were conducted on a Linux server using 10 cores of Intel Xeon CPU E5-2650.

### 4.1 Description of Datasets

The three models were evaluated on diverse pathway datasets traversing the genomic information hierarchy (M. A. Basher et al., 2020): i)- T1 golden consisting of EcoCyc, HumanCyc, AraCyc, YeastCyc, LeishCyc, and TrypanoCyc; ii)- BioCyc (v20.5 T2 & 3) (Caspi et al., 2016); iii)- Critical Assessment of Metagenome Interpretation (CAMI) dataset composed of 40 genomes (Sczyrba et al., 2017); and iv)- Synset-2, a noisy training dataset, introduced in (M. A. Basher et al., 2020). For training, we applied BioCyc (v20.5 T2 & 3) data while for evaluating and testing we used T1 golden and CAMI data. The Synset-2 data was used to obtain supplementary pathways (see Section 4.2). The preprocessed experimental datasets can be obtained from [zenodo.org/record/5630322#.YYXur2DMK3B](https://zenodo.org/record/5630322#.YYXur2DMK3B) while information about these data is provided in (M. A. Basher et al., 2020).

### 4.2 Parameter Settings

Three experiments were conducted: i)- parameter sensitivity analysis, ii)- groups visualization, and iii)- metabolic pathway prediction. Unless otherwise mentioned, we applied the following default configura-

tions: the pathway distribution over groups  $\Phi$  were initialized using gamma distribution (with shape and scale parameters were fixed to 100 and  $1/100$ , respectively), the forgetting rate was  $g = 0.9$ , the delay rate was  $l = 1$ , the batch size was 100, the number of epochs was 3, the number of groups was  $b = 200$ , top  $k$  pathways was 100 (only for SOAP and SPREAT), the Dirichlet hyperparameters  $\alpha$  and  $\xi$  were 0.0001, and the beta hyperparameters  $\gamma$  and  $\kappa$  were 2 and 3, respectively. The supplementary pathways  $\mathbf{M}$  for BioCyc, CAMI, and golden T1 datasets were obtained using mILGPR (elastic-net with enzymatic reaction and pathway evidence features)(M. A. Basher et al., 2020) trained on Synset-2. A schematic view of pathway frequency for BioCyc T2 & 3 and CAMI data with their background pathways is depicted in Fig. 3.

After obtaining groups, we followed the pathway prediction pipeline in Fig. 1 by first mapping examples to groups using reMap software (Hallam Lab, 2021b) and, then, the pathway prediction is achieved using leADS software (Hallam Lab, 2021a). All hyperparameters in reMap, leADS, and mILGPR, were fixed to their default values.

## 5 EXPERIMENTAL RESULTS

This section analyzes the three models using the settings explained in the previous section.

### 5.1 Sensitivity Analysis

**Experimental setup.** Following the common practice, here we study the effect of hyperparameters on the performance of correlated models. First, we compare the sensitivity of SOAP and SPREAT against CTM by incorporating the background pathways  $\mathbf{M}$



while varying the number of groups according to  $b \in \{50, 100, 150, 200, 300\}$ . Next, we examine the SOAP and SPREAT with collapsed option (or c2m) to compare their performances to CTM, where the former models should exhibit similar performances as CTM. Finally, we conduct sparsity analysis of group distribution by varying the cutoff threshold value according to  $k \in \{50, 100, 150, 200, 300, 500\}$  (Section 2.2). For comparative analysis, we apply CAMI as test data to report the log predictive distribution (Section 3.1), where a lower score entails higher generalization capability for the corresponding model.

**Experimental results.** While the log predictive scores for SOAP and SPREAT in Fig. 4a appear to be flat across various group sizes, the CTM model projects a more realistic view where its performances are seen to be gaining by including more groups. Both SOAP and SPREAT incorporate supplementary pathways in modeling the pathway distribution, therefore, it is expected to learn additional pathways from  $\mathbf{M}$  that has an average of  $\sim 500$  pathways in relation to BioCyc v20.5 T2 & 3 which has  $\sim 195$  pathways on average. By excluding  $\mathbf{M}$  (“c2m” in Section 2.2) in the SOAP and SPREAT training, the log predictive distribution of these models exhibit similar performance as CTM (Fig. 4b), asserting our previous discussion. From Figs 4a and 4b, it is evident that  $b = 200$  represents an optimum group size. To find an optimum  $k$  value, we fixed  $b = 200$  and re-trained all models. From Fig. 4c, the performances for SOAP and SPREAT are seen to decline ( $< -0.6$ ) when  $k > 100$ .

Results from this experiment suggest that the settings  $b \in \mathbb{Z}_{[150,300]}$  and  $k \in \mathbb{Z}_{[50,100]}$  are optimum for discovering pathway groups in  $\mathcal{P}$ .

## 5.2 Groups Visualization

**Experimental setup.** Recall that groups constitute overlapping pathways. In this experiment, we visually explore the recovered groups from the three correlated models trained on BioCyc (v20.5 T2 & 3) data using configurations discussed in Section 4.2. We investigate group correlations, reflected in  $\Sigma$ , for SOAP, SPREAT, CTM, SOAP+c2m, and SPREAT+c2m models, to analyze the influence of dual sparseness (Section 2.2) and background pathways on  $\Sigma$ .

**Experimental results.** Fig. 5 demonstrates 50 randomly picked groups and their correlations as represented by  $\Sigma$  for all models. The width of edges indicates the strength of correlations. Essentially for every group in these models, there are approximately 12 to 19 closely related groups. This indicates that metabolic pathways are distributed over

multiple groups, therefore, forming overlapping pathways. With regard to  $\mathbf{M}$ , as explained in Section 5.1, background pathways in  $\mathbf{M}$  consist of  $\sim 500$  pathways for organismal or multi-organismal genomes in comparison to BioCyc (v20.5 T2 & 3) data that has an average of  $\sim 195$  pathways. These additive pathways have influenced the construction of group correlation for both SOAP and SPREAT. Pathway groups in SOAP consist of more associated groups ( $\sim 19$  groups) than the remaining models. This has an important implication for pathway prediction outcomes, discussed in Section 5.3. Sparse models share a similar group structure as CTM (also they have similar log predictive scores in Section 5.1), therefore, they may exhibit similar effects on pathway prediction performance. Results from this experiment show that SOAP and SPREAT are better contenders than CTM. Specifically, both models incorporate supplementary pathways and apply dual sparseness to reduce both the group size and the statistically irrelevant pathways.

## 5.3 Metabolic Pathway Prediction

**Experimental setup.** Pathway groups obtained from correlated models are used for pathway prediction. We consider five models: CTM, two models with background pathways (SOAP and SPREAT), and two collapsed models (SOAP+c2m and SPREAT+c2m). After obtaining groups, we trained reMap using the configuration discussed in Section 4.2. The results are reported on golden T1 data using four evaluation metrics: *Hamming loss*, *average precision*, *average recall*, and *average F1 score*. For comparative analysis, four pathway prediction algorithms are used: i)- MinPath v1.2 (Ye and Doak, 2009), ii)- PathoLogic v21 (Karp et al., 2016), iii)- mLGPR (elastic net with enzymatic reaction and pathway evidence features) (M. A. Basher et al., 2020), and iv)- triUMPF (M. A. Basher et al., 2021a).

**Experimental results.** Table 2 shows that groups from SOAP results in competitive performance against the other methods in terms of average F1 score with optimal performance on EcoCyc (0.8336). However, it seems to be underperforming on AraCyc, YeastCyc, and LeishCyc, yielding average F1 scores of 0.4764, 0.4914, and 0.4144, respectively. This is attributed to incorrect background pathways in  $\mathbf{M}$  (see Section 5.1), hence, impacting the training process. Interestingly, SPREAT’s performances are shown to be inferior to SOAP. As alluded in Section 5.2, the average number of correlated groups for SOAP is significantly larger than SPREAT (Section 5.2), enforcing to revisit a true positive pathway for an organism multiple times across groups in SOAP to signal its pres-

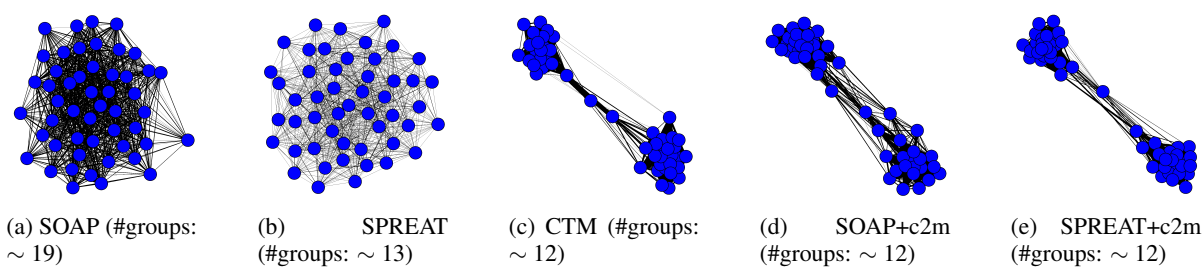


Figure 5: 50 randomly picked pathway groups, represented by blue circles, and their correlations, indicated by black links, for each model. The average number of related groups to each pathway group is indicated by #groups. CTM, SOAP+c2m, and SPREAT+c2m form two distinct clusters of groups, indicating pathways are less shared among groups while SOAP and SPREAT have more shared pathways in their groups.

ence in contrast to groups from SPREAT. With respect to the sensitivity score, correlated models, in general, resulted in higher scores than triUMPF, therefore, at-testing the novelty of modeling groups to improve predictions.

Results from this experiment demonstrate that the group-based approach, in particular SOAP, improves metabolic pathway prediction outcomes. We suggest applying SOAP for pathway predictions using the default configurations discussed in Section 4.2.

## 6 CONCLUSIONS

In this paper, we presented two novel statistical hierarchical mixture models, SOAP and SPREAT, to uncover correlated pathway groups given pathway abundance data. The work is motivated by the problem of missing pathways, which is very common in pathway prediction from organismal and multi-organismal datasets. We empirically evaluated correlated models for pathway prediction using golden T1 data and compared results to other prediction methods including PathoLogic, MinPath, mLGP, and triUMPF. Overall, correlated models showed promising results in boosting prediction performance over ML-based algorithms, such as triUMPF. There are several directions for future study. Foremost, we intend to build a model that combines both graph-based (M. A. Basher et al., 2021a) and group-based methods to improve metabolic pathway prediction with emphasis on multi-organismal genomes. Additional attention should be paid to sparseness induction in the covariance matrix for better interpretability (Fan et al., 2016).

## ACKNOWLEDGEMENTS

This work was performed under the auspices of Genome Canada, Genome British Columbia, the Natural Sciences and Engineering Research Council (NSERC) of Canada, and Compute/Calcul Canada. ARMA was supported by a UBC four-year doctoral fellowship (4YF) administered through the UBC Graduate Program in Bioinformatics.

## REFERENCES

- Aguiar-Pulido, V., Huang, W., Suarez-Ulloa, V., Cickovski, T., Mathee, K., and Narasimhan, G. (2016). Metagenomics, metatranscriptomics, and metabolomics approaches for microbiome analysis: supplementary issue: bioinformatics methods and applications for big metagenomics data. *Evolutionary Bioinformatics*, 12:EBO-S36436.
- Airoldi, E. M., Blei, D. M., Fienberg, S. E., and Xing, E. P. (2008). Mixed membership stochastic blockmodels. *Journal of machine learning research*, 9(Sep):1981–2014.
- Baranwal, M., Magner, A., Elvati, P., Saldinger, J., Violi, A., and Hero, A. O. (2020). A deep learning architecture for metabolic pathway prediction. *Bioinformatics*, 36(8):2547–2553.
- Bien, J. and Tibshirani, R. J. (2011). Sparse estimation of a covariance matrix. *Biometrika*, 98(4):807–820.
- Blei, D. and Lafferty, J. (2006). Correlated topic models. *Advances in neural information processing systems*, 18:147.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022.
- Caspi, R., Billington, R., Foerster, H., and et al. (2016). Biocyc: Online resource for genome and metabolic pathway analysis. *The FASEB Journal*, 30(1 Supplement):lb192–lb192.
- Caspi, R., Billington, R., Keseler, I. M., Kothari, A., Krumnacker, M., Midford, P. E., Ong, W. K., Paley, S., Subhraveti, P., and Karp, P. D. (2019). The metacyc

Table 2: Predictive performance of each comparing algorithm on 6 benchmark datasets. For each performance metric, ‘↓’ indicates the smaller score is better while ‘↑’ indicates the higher score is better. Bold text suggests the best performance in each column.

Methods	Hamming Loss ↓					
	EcoCyc	HumanCyc	AraCyc	YeastCyc	LeishCyc	TrypanoCyc
PathoLogic	0.0610	0.0633	0.1188	<b>0.0424</b>	<b>0.0368</b>	<b>0.0424</b>
MinPath	0.2257	0.2530	0.3266	0.2482	0.1615	0.2561
mLGPR	0.0804	0.0633	<b>0.1069</b>	0.0550	0.0380	0.0590
triUMPF	0.0435	0.0954	0.1560	0.0649	0.0443	0.0776
SOAP	<b>0.0392</b>	<b>0.0400</b>	0.1714	0.0934	0.0772	0.0479
SPREAT	0.0519	0.0827	0.1489	<u>0.0748</u>	0.0629	0.0503
CTM	0.0558	0.0835	<u>0.1425</u>	0.0804	0.0622	0.0503
SOAP+c2m	0.0590	0.0780	0.1457	0.0772	0.0614	0.0534
SPREAT+c2m	0.0542	0.0796	0.1520	0.0772	<u>0.0598</u>	0.0558
Methods	Average Precision Score ↑					
	EcoCyc	HumanCyc	AraCyc	YeastCyc	LeishCyc	TrypanoCyc
PathoLogic	0.7230	0.6695	0.7011	0.7194	<b>0.4803</b>	0.5480
MinPath	0.3490	0.3004	0.3806	0.2675	0.1758	0.2129
mLGPR	0.6187	0.6686	0.7372	0.6480	0.4731	0.5455
triUMPF	0.8662	0.6080	0.7377	<b>0.7273</b>	0.4161	0.4561
SOAP	0.8611	<b>0.7871</b>	0.6215	0.4851	0.2805	0.5985
SPREAT	<b>0.9400</b>	0.6750	0.8350	<u>0.6000</u>	0.3200	<b>0.6200</b>
CTM	0.9150	0.6700	<b>0.8750</b>	<u>0.5650</u>	0.3250	<b>0.6200</b>
SOAP+c2m	0.8950	0.7050	0.8550	0.5850	0.3300	0.6000
SPREAT+c2m	0.9250	0.6950	0.8150	0.5850	<u>0.3400</u>	0.5850
Methods	Average Recall Score ↑					
	EcoCyc	HumanCyc	AraCyc	YeastCyc	LeishCyc	TrypanoCyc
PathoLogic	0.8078	0.8423	0.7176	0.8734	0.8391	0.7829
MinPath	<b>0.9902</b>	<b>0.9713</b>	<b>0.9843</b>	<b>1.0000</b>	<b>1.0000</b>	<b>1.0000</b>
mLGPR	0.8827	0.8459	0.7314	0.8603	0.9080	0.8914
triUMPF	0.7590	0.3835	0.3529	0.3319	0.7126	0.6229
SOAP	<u>0.8078</u>	<u>0.8746</u>	<u>0.3863</u>	0.4978	<u>0.7931</u>	<u>0.9371</u>
SPREAT	0.6124	0.4839	0.3275	<u>0.5240</u>	0.7356	0.7086
CTM	0.5961	0.4803	0.3431	0.4934	0.7471	0.7086
SOAP+c2m	0.5831	0.5054	0.3353	0.5109	0.7586	0.6857
SPREAT+c2m	0.6026	0.4982	0.3196	0.5109	0.7816	0.6686
Methods	Average F1 Score ↑					
	EcoCyc	HumanCyc	AraCyc	YeastCyc	LeishCyc	TrypanoCyc
PathoLogic	0.7631	0.7460	0.7093	<b>0.7890</b>	0.6109	0.6447
MinPath	0.5161	0.4589	0.5489	0.4221	0.2990	0.3511
mLGPR	0.7275	0.7468	<b>0.7343</b>	0.7392	<b>0.6220</b>	0.6768
triUMPF	0.8090	0.4703	0.4775	0.4735	0.5254	0.5266
SOAP	<b>0.8336</b>	<b>0.8285</b>	0.4764	0.4914	0.4144	<b>0.7305</b>
SPREAT	0.7416	0.5637	0.4704	<u>0.5594</u>	0.4460	0.6613
CTM	0.7219	0.5595	<u>0.4930</u>	<u>0.5268</u>	0.4530	0.6613
SOAP+c2m	0.7061	0.5887	<u>0.4817</u>	0.5455	0.4599	0.6400
SPREAT+c2m	0.7298	0.5804	0.4592	0.5455	<u>0.4739</u>	0.6240

database of metabolic pathways and enzymes-a 2019 update. *Nucleic acids research*.

Dale, J. M., Popescu, L., and Karp, P. D. (2010). Machine learning methods for metabolic pathway prediction. *BMC bioinformatics*, 11(1):1.

Fan, J., Liao, Y., and Liu, H. (2016). An overview of the estimation of large covariance and precision matrices. *The Econometrics Journal*, 19(1):C1–C32.

Hahn, A. S., Konwar, K. M., Louca, S., and et al. (2016). The information science of microbial ecology. *Current opinion in microbiology*, 31:209–216.

- Hanson, N. W., Konwar, K. M., Hawley, A. K., and et al. (2014). Metabolic pathways for the whole community. *BMC genomics*, 15(1):1.
- Hassa, J., Maus, I., Off, S., Pühler, A., Scherer, P., Klocke, M., and Schlüter, A. (2018). Metagenome, metatranscriptome, and metaproteome approaches unraveled compositions and functional relationships of microbial communities residing in biogas plants. *Applied microbiology and biotechnology*, 102(12):5045–5063.
- He, J., Hu, Z., Berg-Kirkpatrick, T., and et al. (2017). Efficient correlated topic modeling with topic embedding. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 225–233. ACM.
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. W. (2013). Stochastic variational inference. *Journal of Machine Learning Research*, 14(1):1303–1347.
- Jiao, D., Ye, Y., and Tang, H. (2013). Probabilistic inference of biochemical reactions in microbial communities from metagenomic sequences. *PLoS Comput Biol*, 9(3):e1002981.
- Kanehisa, M., Furumichi, M., Tanabe, M., and et al. (2017). Kegg: new perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Research*, 45(D1):D353–D361.
- Karp, P. D., Latendresse, M., Paley, S. M., and et al. (2016). Pathway tools version 19.0 update: software for pathway/genome informatics and systems biology. *Briefings in bioinformatics*, 17(5):877–890.
- Khatri, P., Sirota, M., and Butte, A. J. (2012). Ten years of pathway analysis: current approaches and outstanding challenges. *PLoS computational biology*, 8(2):e1002375.
- Lawson, C. E., Harcombe, W. R., Hatzenpichler, R., and et al. (2019). Common principles and best practices for engineering microbiomes. *Nature Reviews Microbiology*, pages 1–17.
- Lin, T., Tian, W., Mei, Q., and Cheng, H. (2014). The dual-sparse topic model: mining focused topics and focused terms in short text. In *Proceedings of the 23rd international conference on World wide web*, pages 539–550. ACM.
- Loh, P.-R., Baym, M., and Berger, B. (2012). Compressive genomics. *Nature biotechnology*, 30(7):627.
- M. A. Basher, A. R. and Hallam, S. J. (2021). Relabeling metabolic pathway data with groups to improve prediction outcomes. *BioRxiv*.
- M. A. Basher, A. R., McLaughlin, R. J., and Hallam, S. J. (2020). Metabolic pathway inference using multi-label classification with rich pathway features. *PLOS Computational Biology*, 16(10):1–22.
- M. A. Basher, A. R., McLaughlin, R. J., and Hallam, S. J. (2021a). Metabolic pathway prediction using non-negative matrix factorization with improved precision. *Journal of Computational Biology*.
- M. A. Basher, A. R., McLaughlin, R. J., and Hallam, S. J. (2021b). Metabolic pathway prediction using non-negative matrix factorization with improved precision. In *Computational Advances in Bio and Medical Sciences*, pages 33–44, Cham. Springer International Publishing.
- Mascher, M., Schreiber, M., Scholz, U., Graner, A., Reif, J. C., and Stein, N. (2019). Genebank genomics bridges the gap between the conservation of crop diversity and plant breeding. *Nature genetics*, 51(7):1076–1081.
- Qi, Q., Li, J., and Cheng, J. (2014). Reconstruction of metabolic pathways by combining probabilistic graphical model-based and knowledge-based methods. In *BMC proceedings*, volume 8, pages 1–10. Springer.
- Sczyrba, A., Hofmann, P., Belmann, P., and et al. (2017). Critical assessment of metagenome interpretation—a benchmark of metagenomics software. *Nature methods*, 14(11):1063.
- Shafiei, M., Dunn, K. A., Chipman, H., Gu, H., and Bielawski, J. P. (2014). Biomenet: A bayesian model for inference of metabolic divergence among microbial communities. *PLoS Comput Biol*, 10(11):e1003918.
- Stewart, F. J., Sharma, A. K., Bryant, J. A., and et al. (2011). Community transcriptomics reveals universal patterns of protein sequence conservation in natural microbial communities. *Genome biology*, 12(3):R26.
- Tabei, Y., Yamanishi, Y., and Kotera, M. (2016). Simultaneous prediction of enzyme orthologs from chemical transformation patterns for de novo metabolic pathway reconstruction. *Bioinformatics*, 32(12):i278–i287.
- Hallam Lab (2021a). leADS: <https://github.com/hallamlab/leADS>.
- Hallam Lab (2021b). reMap: <https://github.com/hallamlab/reMap>.
- Wang, W.-L., Xu, S.-Y., Ren, Z.-G., Tao, L., Jiang, J.-W., and Zheng, S.-S. (2015). Application of metagenomics in the human gut microbiome. *World journal of gastroenterology: WJG*, 21(3):803.
- Yamanishi, Y., Tabei, Y., and Kotera, M. (2015). Metabolome-scale de novo pathway reconstruction using regioisomer-sensitive graph alignments. *Bioinformatics*, 31(12):i161–i170.
- Ye, Y. and Doak, T. G. (2009). A parsimony approach to biological pathway reconstruction/inference for genomes and metagenomes. *PLoS Comput Biol*, 5(8):e1000465.
- Zhao, Y., Chen, M.-H., Pei, B., Rowe, D., Shin, D.-G., Xie, W., Yu, F., and Kuo, L. (2012). A bayesian approach to pathway analysis by integrating gene–gene functional directions and microarray data. *Statistics in biosciences*, 4(1):105–131.