

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

<sup>1</sup>Graduate Program in Bioinformatics, <sup>2</sup>Department of Microbiology and Immunology, University of British Columbia

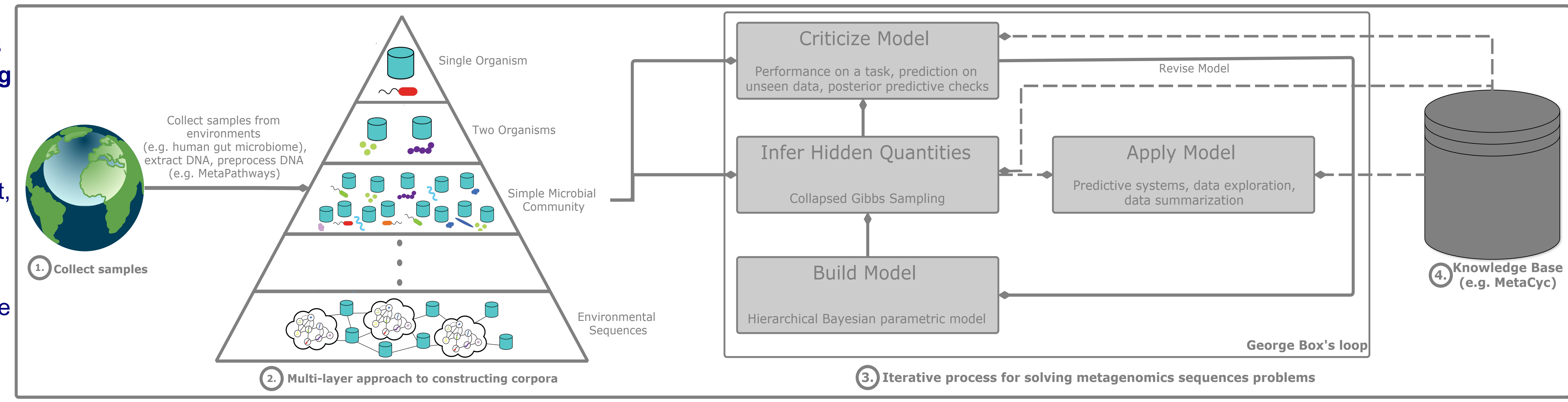
## 1 INTRODUCTION AND MOTIVATION

Microbial communities facilitate the majority of the biochemical activity on Earth, playing integral roles in energy and matter transformations in natural and engineered ecosystems. Metagenomics is used to analyze the genetic material of microbial communities directly from an environmental sample.

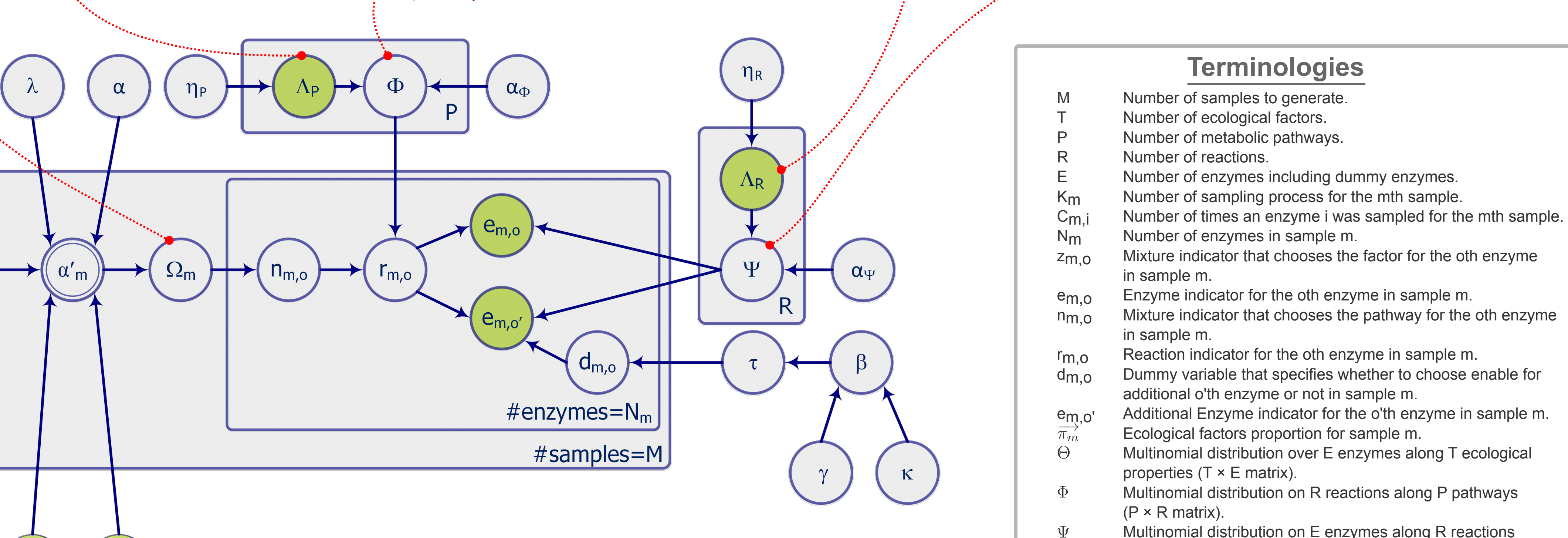
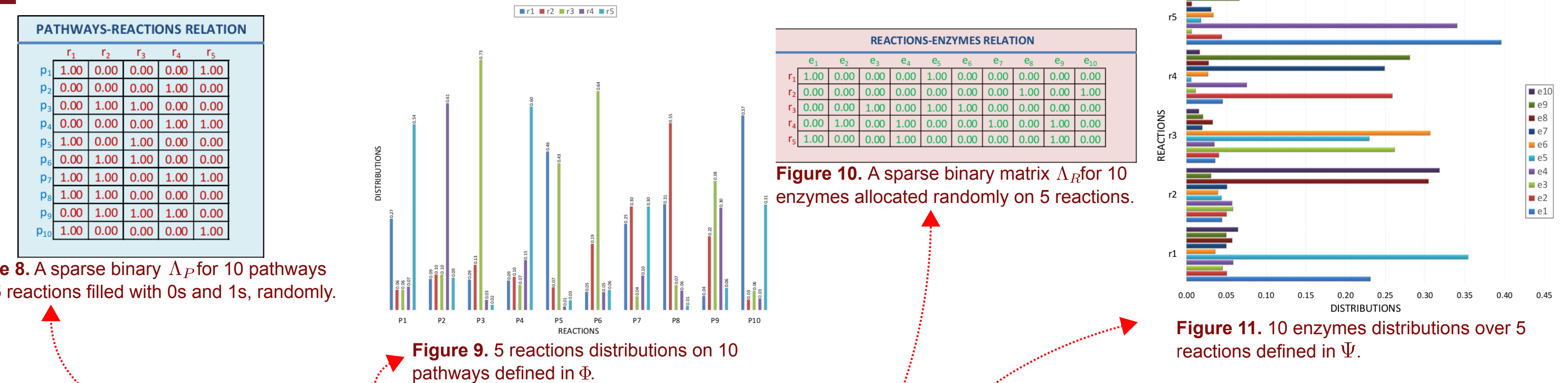
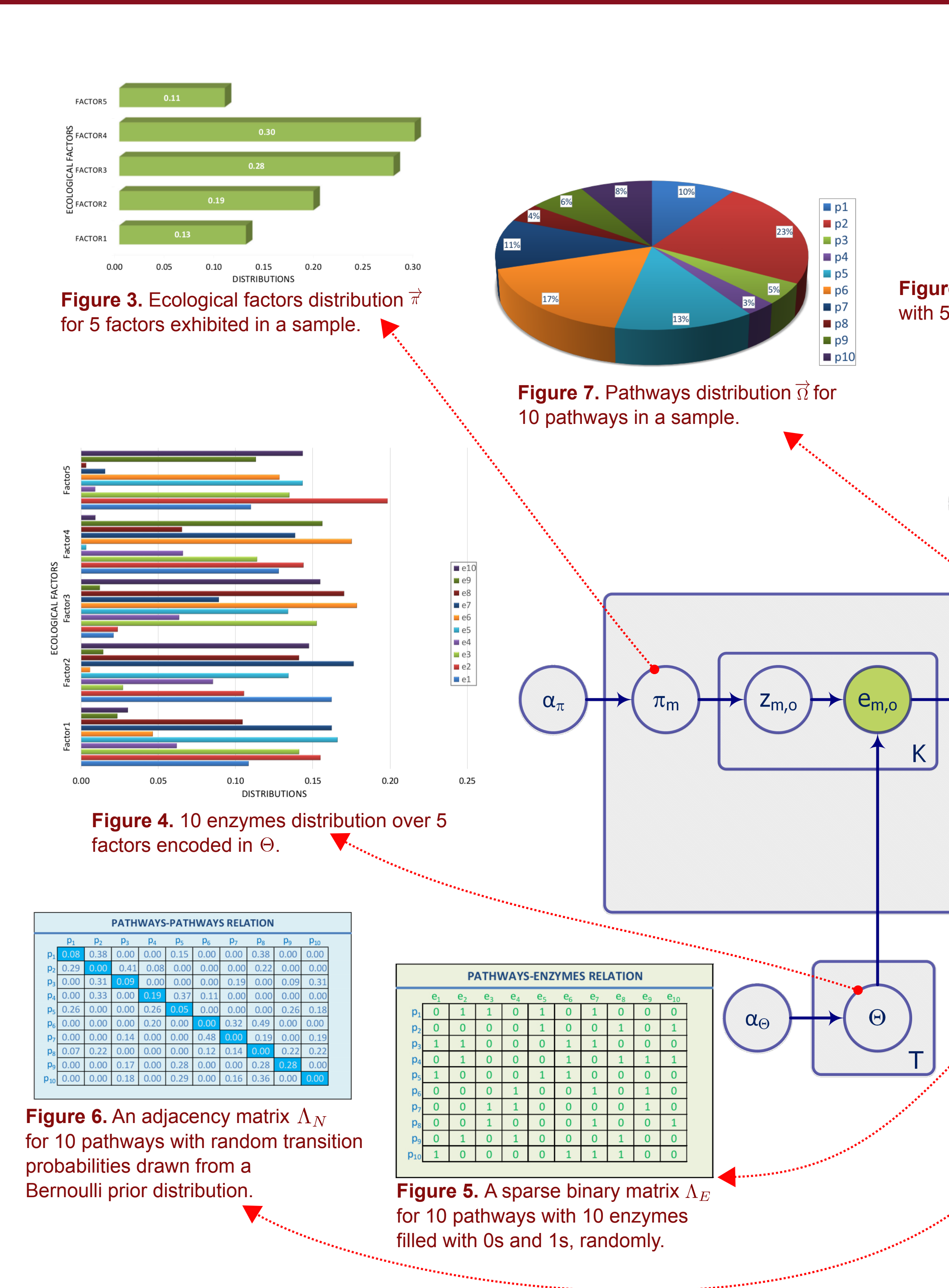
To estimate the metabolic potential of a metagenomic sample we devise a novel approach to reconstructing biological pathways from enzyme annotations and environmental parameters. Our approach enjoys a modular, flexible strategy based on statistical hierarchical Bayesian deep framework that encodes emergent information represented in MetaCyc, a highly curated database of enzyme sequences, reactions, and pathways. The model is based on graphical modeling techniques to infer latent pathways represented as mixture components in a sample. For the training, we adopt a collapsed Gibbs sampling technique to examine the genetic content of metagenomic datasets. Further, the model is well defined mathematically and aligns with the biological interpretations. Based on our preliminary analysis, we anticipate that our model can outperform the PathoLogic algorithm on a single organism task.

## 2 APPROACH

**Figure 1. A modern interpretation of the perspective of George Box's loop to iterative process for solving metagenomics data analysis problems.** Metagenomic sequence information is assembled and ORFs are predicted and annotated (1). Next, the enzymes and their associated pathways are curated based on a multilayer approach where the top layer dataset is synthesized using one organism (positive) mixed with non-overlapping pathways-enzymes (negative) (2). The next layer comprises of two organisms with partially overlapping pathways-enzymes (positive) and non-overlapping pathways-enzymes (negative). Continuing to add pathway information from more organisms approaches the metabolic potential of a microbial community. Afterward, Box's loop comes into play in which an iterative cycle of experimental design, model formulation, model criticism, and application refines the model (3). In the first step of the loop, a probability model is built with a well defined mathematical object. Observed data enter the picture in the second step of Box's loop. Here, the computational aspects are applied to infer the pathways using an inference algorithm to compute the posterior distribution (e.g. collapsed Gibbs sampling) and a knowledge base (e.g. MetaCyc) (4). Finally, the trained model is tested against real data, identifying the important ways that it succeeds and fails in extracting pathways.



## 3 HIERARCHICAL BAYESIAN PARAMETRIC MODEL



**Figure 2. Graphical model representation of an unsupervised Bayesian model for metabolic pathway inference.** The boxes are "plates" representing replicates. The outer plate represents metagenomic samples, while the inner left plate represents the repeated choice of factors and enzymes, and the inner right plate represents the repeated choice of reactions and enzymes within a sample. The model comprises of a hierarchical Bayesian mixture model, where enzymes constitute reactions, reactions are mixed to form pathways, pathways are determined by reactions and environmental parameters, and each sample is treated as a vector of environmental parameter distributions.

## 4 COLLAPSED GIBBS SAMPLING ALGORITHM

The overall generative process is summarized below:

- Sample a distribution  $\tau \sim \text{Beta}(\cdot, \gamma, \kappa)$
- For each ecological factors  $t \in \{1, \dots, T\}$ :
  - Sample a distribution over enzymes  $\theta_t = (\theta_{t,1}, \dots, \theta_{t,P})^T \sim \text{Dirichlet}(\cdot | \vec{\alpha}_t)$
- For each pathway  $i \in \{1, \dots, P\}$ :
  - Sample labels from  $\Lambda_P^i \in \{0, 1\}_1^R \sim \text{Bernoulli}(\cdot | \vec{\eta}_R^i) : \forall j (j \in i) = 1$  where  $j \in \{1, \dots, R\}$
  - Sample a distribution over reactions  $\Phi_i = (\phi_{i,1}, \dots, \phi_{i,R})^T \sim \text{Dirichlet}(\cdot | \vec{\alpha}_i^R), \Lambda_P^i$
- For each reaction  $j \in \{1, \dots, R\}$ :
  - Sample labels from  $\Lambda_R^j \in \{0, 1\}_1^E \sim \text{Bernoulli}(\cdot | \vec{\eta}_E^j) : \forall e (e \in j) = 1$  where  $e \in \{1, \dots, E\}$
  - Sample a distribution over enzymes  $\Psi_j = (\psi_{j,1}, \dots, \psi_{j,K})^T \sim \text{Dirichlet}(\cdot | \vec{\alpha}_j^K), \Lambda_R^j$
- For each sample  $m \in \{1, \dots, M\}$ :
  - Sample a distribution over ecological factors  $\vec{\pi}_m = (\pi_{m,1}, \dots, \pi_{m,T})^T \sim \text{Dirichlet}(\cdot | \vec{\alpha}_T)$
  - For each enzyme  $o \in \{1, \dots, K\}$ :
    - Sample a factor label  $z_{m,o} \sim \text{Multinomial}(\cdot | \vec{\pi}_m)$
    - Sample an enzyme label  $e_{m,o} \sim \text{Multinomial}(\cdot | \vec{\Theta}_{z_{m,o}})$
  - Compute the Dirichlet prior  $\vec{\alpha}_m = \vec{\lambda} \cdot \rho_m + \alpha$   
 $\rho_m = (\vec{s} \cdot \mathbf{A}_E^T) \cdot \mathbf{A}_N$  and  $\vec{s} = (\frac{C_{m,E}}{K_m}, \dots, \frac{C_{m,E}}{K_m})$
  - Sample a distribution over pathways  $\vec{\Omega}_m = (\Omega_{m,1}, \dots, \Omega_{m,P})^T \sim \text{Dirichlet}(\cdot | \vec{\alpha}_m^P)$
  - For each observed enzyme  $o \in \{1, \dots, N\}$ :
    - Sample a pathway label  $n_{m,o} \sim \text{Multinomial}(\cdot | \vec{\Omega}_m)$
    - Sample a reaction label  $r_{m,o} \sim \text{Multinomial}(\cdot | \vec{\Phi}_{n_{m,o}})$
    - Sample an enzyme  $e_{m,o} \sim \text{Multinomial}(\cdot | \vec{\Psi}_{r_{m,o}})$
    - Sample a binary label  $d_{m,o} \sim \text{Bernoulli}(\cdot | \tau)$
    - Sample an additional enzyme  $e_{m,o'}$  according to:  
 $(e_{m,o} \sim \text{Multinomial}(\cdot | \vec{\Psi}_{r_{m,o}}))$  if  $d_{m,o} = 0$  or  $(e_{m,o}, e_{m,o'}) \sim \{\text{Multinomial}(\cdot | \vec{\Psi}_{r_{m,o}})\}_1^2$  if  $d_{m,o} = 1$

## 5 CONCLUSION AND DISCUSSION

Inspired by Box's loop, an unsupervised hierarchical deep Bayesian architecture is developed to detect pathways that are present in ecological sequences, which are constructed in a multi-layer approach. Further work includes adopting a supervised strategy to recovering pathways, studying the correlated pathways and enzymes to better understanding the microbial interactions, examining pathways abundances to assist in capturing the global metabolic network in samples, constructing a well-defined set of ecological factors contributing to pathways inference, learning mixtures of taxa in microbiomes, suggesting ways to capture super-pathways in MetaCyc, and proposing ways to minimize computational intricacies in a sparse metagenomics sequences.

## 6 REFERENCES

- Karpe, Peter D., et al. (2011). "The pathway tools pathway prediction algorithm." Standards in genomic sciences 5:3: 424.
- Abubucker, Sahar, et al. (2012). "Metabolic reconstruction for metagenomic data and its application to the human microbiome." PLoS Comput Biol 8:6.
- Shafiei, Mahdi, et al. (2014). "BiomeNet: A Bayesian model for inference of metabolic divergence among microbial communities." PLoS Comput Biol 10:11.
- Hanson, Niels W., et al. (2014). "Metabolic pathways for the whole community." BMC genomics 15:1.
- Blei, David M. (2014). "Build, compute, critique, repeat: Data analysis with latent variable models." Annual Review of Statistics and Its Application 1: 203-232.
- Konwar, Kishori M., et al. (2015). "MetaPathways v2. 5: quantitative functional, taxonomic and usability improvements." Bioinformatics 31:20: 3345-3347.

Terminologies	
M	Number of samples to generate.
T	Number of ecological factors.
P	Number of metabolic pathways.
R	Number of reactions.
E	Number of enzymes including dummy enzymes.
K <sub>m</sub>	Number of sampling process for the mth sample.
C <sub>m,i</sub>	Number of times an enzyme i was sampled for the mth sample.
N <sub>m</sub>	Number of enzymes in sample m.
Z <sub>m,o</sub>	Mixture indicator that chooses the factor for the oth enzyme in sample m.
e <sub>m,o</sub>	Enzyme indicator for the oth enzyme in sample m.
n <sub>m,o</sub>	Mixture indicator that chooses the pathway for the oth enzyme in sample m.
r <sub>m,o</sub>	Reaction indicator for the oth enzyme in sample m.
d <sub>m,o</sub>	Dummy variable that specifies whether to choose enable for additional oth enzyme or not in sample m.
e <sub>m,o'</sub>	Additional Enzyme indicator for the oth enzyme in sample m.
$\vec{\pi}_m$	Ecological factors proportion for sample m.
$\vec{\Theta}$	Multinomial distribution over E enzymes along T ecological properties (T × E matrix).
$\Phi$	Multinomial distribution on R reactions along P pathways (P × R matrix).
$\Psi$	Multinomial distribution on E enzymes along R reactions (R × E matrix).
$\vec{\Omega}_m$	Pathways proportion for sample m.
$\vec{\alpha}_m$	Scaled, smoothed, normalized parameter for pathways in sample m (P-vector).
$\Lambda_E, \Lambda_P, \Lambda_R$	Sparse binary indicator matrices of sizes (P × E), (P × R), and (R × E), respectively.
$\Lambda_N$	Stochastic matrix of size P × P.
$\vec{\lambda}$	Hyperparameter that specifies the weight contributed to pathways (P-vector).
$\vec{\alpha}_\pi, \vec{\alpha}_\theta, \vec{\alpha}_\Phi, \vec{\alpha}_\Psi$	Hyperparameters on $\pi, \theta, \Phi$ , and $\Psi$ , respectively.
$\alpha$	Symmetric hyperparameter on $\vec{\alpha}_m$ .
$\vec{\eta}_E, \vec{\eta}_R, \vec{\eta}_P, \vec{\eta}_N$	Bernoulli prior distribution on $\Lambda_E, \Lambda_N, \Lambda_P$ , and $\Lambda_R$ respectively.
$\tau$	Bernoulli prior distribution on $d_{m,o}$ .
$\beta$	Beta distribution.
$\gamma, \kappa$	Hyperparameters for $\beta$ .