



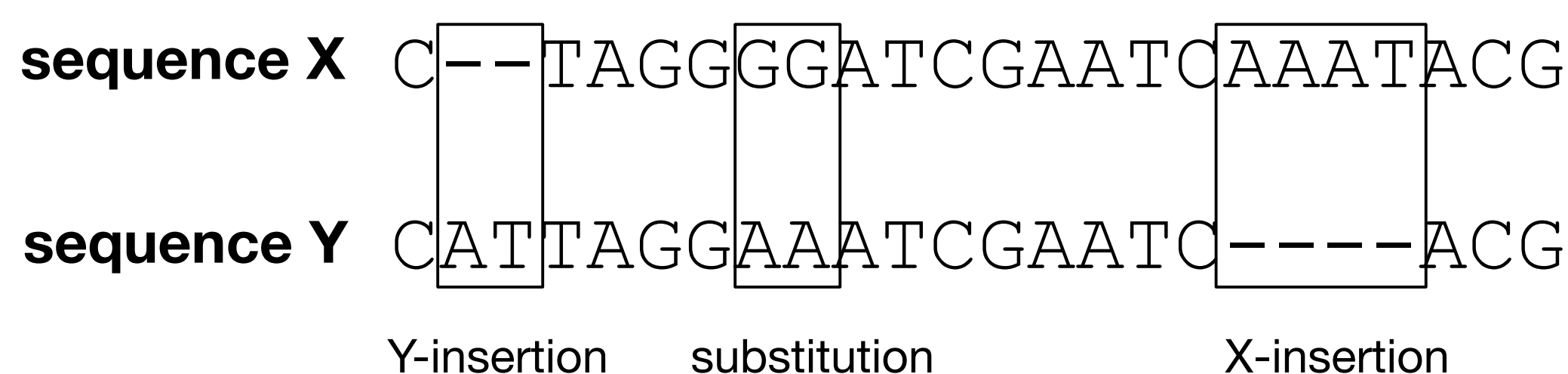
Abstract

Although extensive research has been done on biological sequence alignment problem, model selection problem remains untouched. We assume the model selection enables us to extract biological knowledge by interpreting resulting models, e.g. by comparing resulting models of pairwise DNA sequence alignment between some different pairs of species. As a model selection method, we introduce Factorized Asymptotic Bayesian Pairwise Hidden Markov Models (FAB-PHMM), based on asymptotic consistent information criterion with model evidence. We conducted an experiment on a synthetic dataset to illustrate model selection capability of proposed method. On a real DNA sequence data experiment, we observed that much more complex models are selected than previously utilized models, and the result is consistent with evolutionary distance.

Background

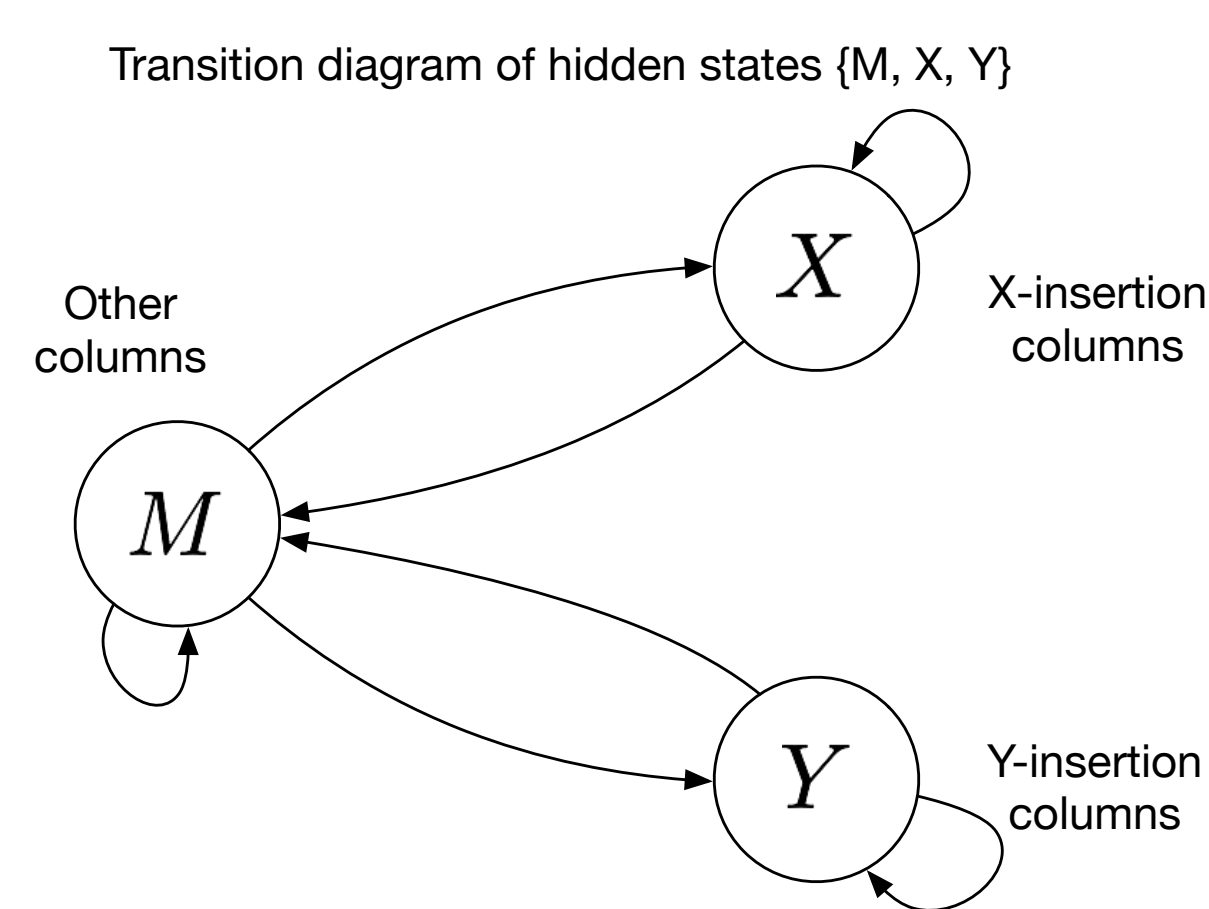
Pairwise sequence alignment

Pairwise sequence alignment aims to assess similarity of two sequences by introducing “gap”. Gaps are represented as “-” and the gap columns have the nucleotides only in single side. In score-based alignment [Smith+1981], a similarity score is calculated as the sum of column scores; match columns make positive contribution and insertion/substitution columns affect negatively. The optimal alignment is given by Dynamic Programming, which seek the optimal gap insertions that maximize the similarity scores. Similarity of the sequence pair can be assessed by the score of optimal alignment.



Pairwise Hidden Markov Models (PHMM)

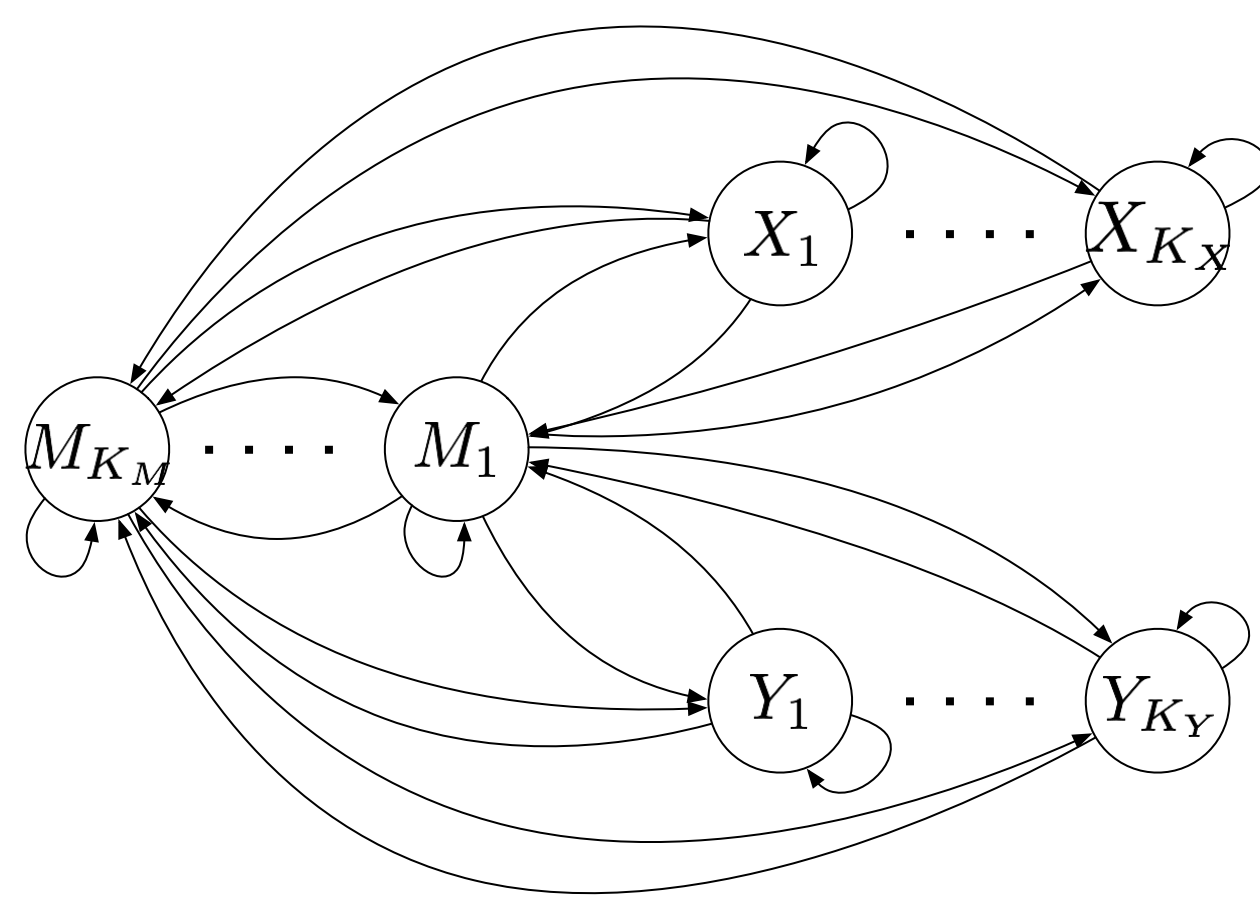
Pairwise Hidden Markov Models are probabilistic models [Durbin+, 1999] of pairwise sequence alignment. Using this model, alignment can be obtained via MAP decoding. In addition, parameters can be learned using EM algorithm.



Problem Setting

Question: Which is the “best” model of PHMM?

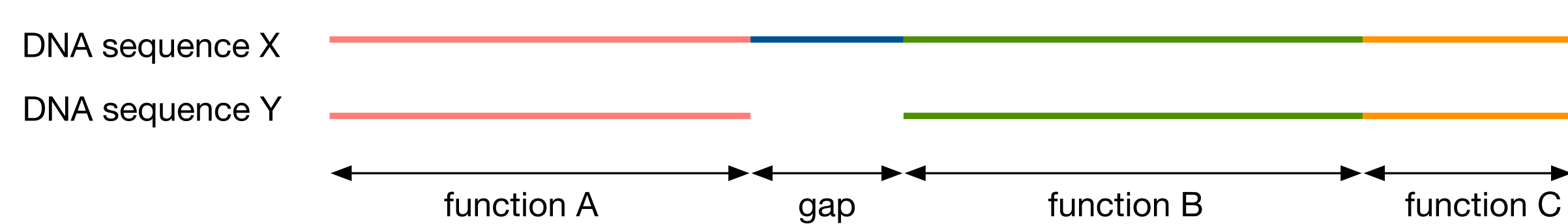
The model of PHMM is parameterized by the number of hidden states for each type $M = (K_M, K_X, K_Y)$. Our goal is to select the best model M that maximize a marginal likelihood given a observed data.



Motivations

1) Biological function annotation

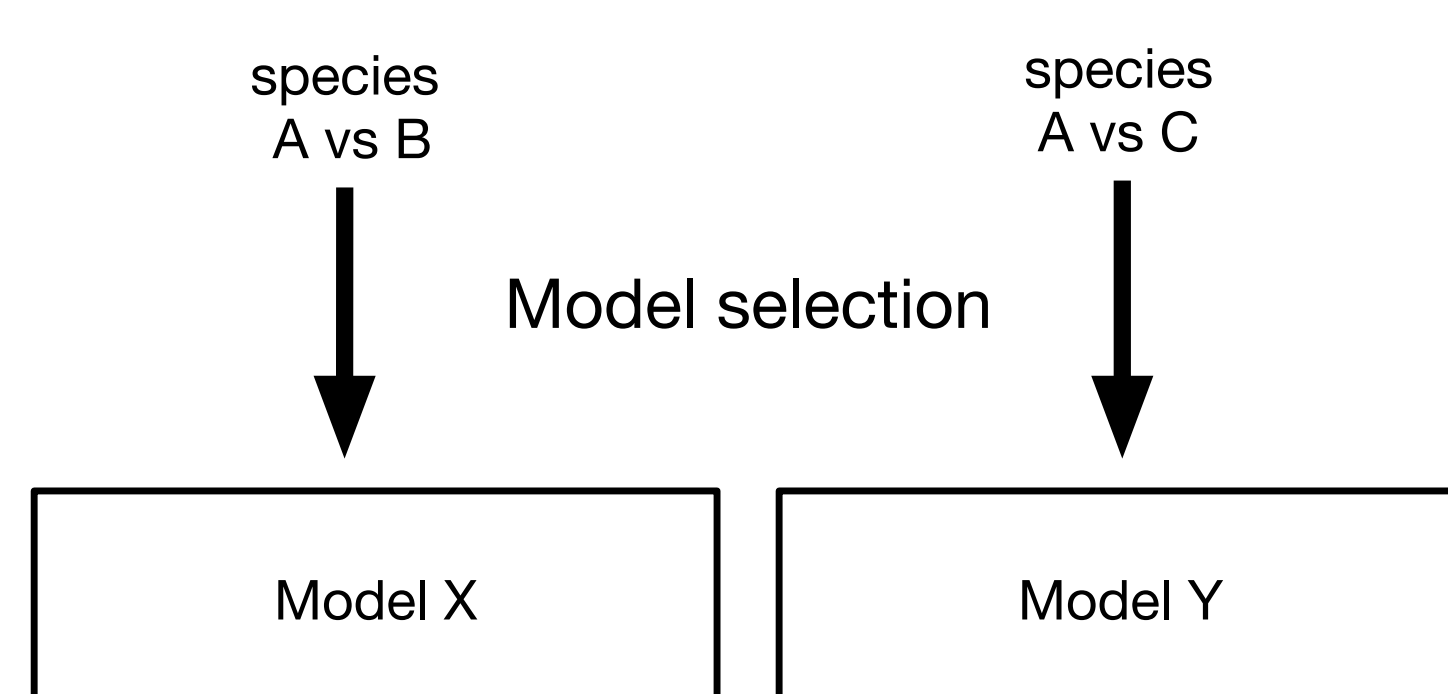
Each DNA region has different biological functions. Assuming that each of those regions are generated from different probabilistic distributions corresponding to a specific hidden states, we can annotate the regions with function type by inspecting posterior of hidden states.



2) Evolutionary distance

We assume resulting model reflect evolutionary distance; selected model of evolutionary further species should have complex model. The selected model might summarize evolutionary behavior of them.

Pairwise sequence comparison among some species



Difference between model X and Y might reflect difference between B and C

Methods

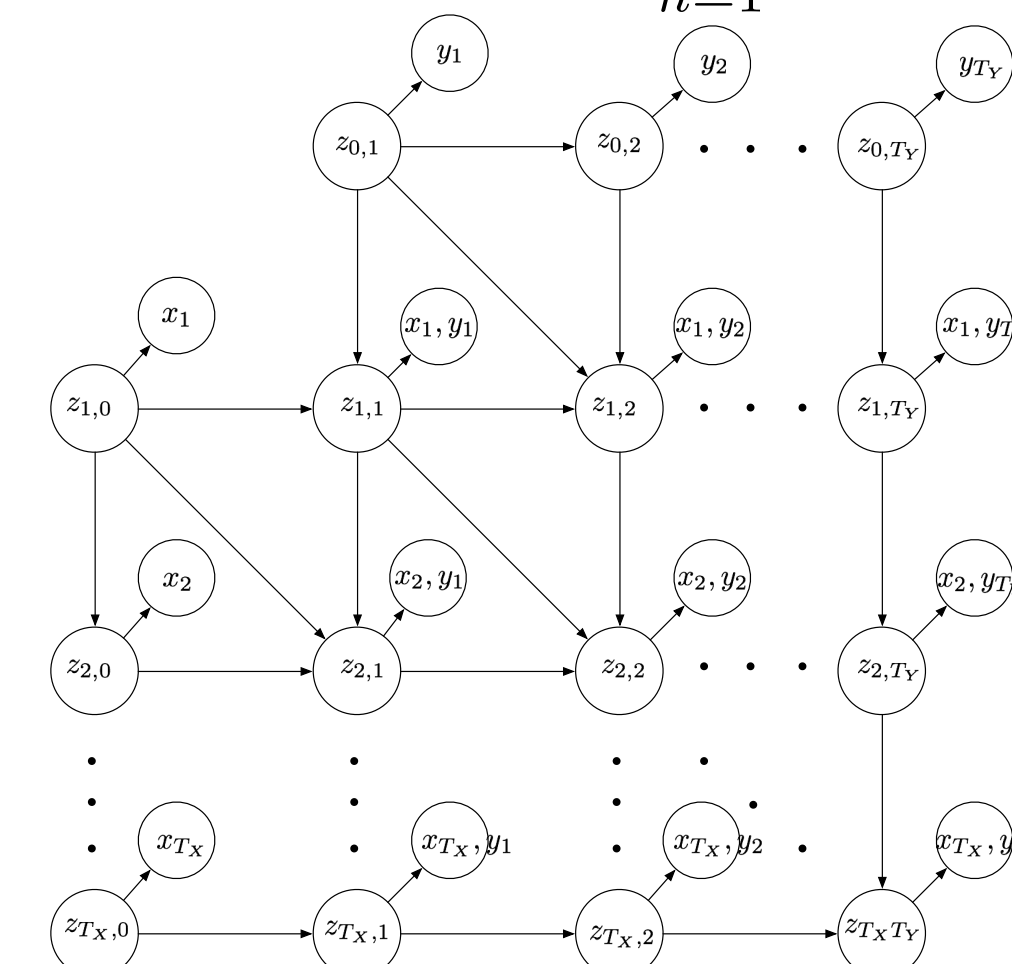
Formulation of HMM and PHMM

HMM

$$\log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) = \sum_{n=1}^N \left[\log p(\mathbf{z}_1^n | \boldsymbol{\alpha}) + \sum_{t=1}^T \left(\log p(\mathbf{z}_t^n | \mathbf{z}_{t-1}^n, \boldsymbol{\beta}) + \log p(\mathbf{x}_t^n | \mathbf{z}_t^n, \boldsymbol{\phi}) \right) \right]$$

PHMM

$$\log p(\mathbf{x}, \mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) = \sum_{n=1}^N \left[\log p(\mathbf{z}_{\text{in}}^n | \boldsymbol{\alpha}) + \sum_{t=0}^{T_X} \sum_{u=0}^{T_Y} \left(\log p(\mathbf{z}_{tu}^n | \mathbf{p}(\mathbf{z}_{tu}^n), \boldsymbol{\beta}) + \log p(\mathbf{x}_t^n, \mathbf{y}_u^n | \mathbf{z}_{tu}^n, \boldsymbol{\phi}) \right) \right]$$



$$\mathbf{x}^n = [GG]$$

$$\mathbf{y}^n = [AAGG]$$

$$\mathbf{z}^n = \begin{bmatrix} 0 & \mathbf{z}_Y & \mathbf{z}_Y & 0 & 0 \\ 0 & 0 & 0 & \mathbf{z}_M & 0 \\ 0 & 0 & 0 & 0 & \mathbf{z}_M \end{bmatrix} \begin{matrix} \text{G} \\ \text{G} \\ \text{A} \end{matrix} \mathbf{x}^n$$

Factorized Asymptotic Bayesian Pairwise Hidden Markov Models (FAB-PHMM)

Following the FAB-HMM [Fujimaki+2012], we apply Factorized Information Criterion (FIC) for PHMM. Our goal is to maximize FIC, which is an asymptotically consistent approximation of marginal likelihood.

Maximization of marginal likelihood

$$\mathcal{M}^* = \arg \max_{\mathcal{M}} p(\mathbf{x}, \mathbf{y} | \mathcal{M}) \quad p(\mathbf{x}, \mathbf{y} | \mathcal{M}) = \int \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{y}, \mathbf{z}, \boldsymbol{\theta} | \mathcal{M}) d\boldsymbol{\theta}$$

Variational Lower Bound

We introduce variational distribution $q(\mathbf{z})$ to take lower-bound of the marginal likelihood.

$$p(\mathbf{x}, \mathbf{y} | \mathcal{M}) \geq \sum_{\mathbf{z}} q(\mathbf{z}) \log \left(\frac{p(\mathbf{x}, \mathbf{y}, \mathbf{z} | \mathcal{M})}{q(\mathbf{z})} \right).$$

Factorized Information Criterion (FIC)

$$p(\mathbf{x}, \mathbf{y}, \mathbf{z} | \mathcal{M}) = \int \prod_{n=1}^N p(\mathbf{z}_{\text{in}}^n | \boldsymbol{\alpha}) \prod_{t=0}^{T_X} \prod_{u=0}^{T_Y} \prod_{k=1}^K p_k(\mathbf{z}_{tu}^n | \boldsymbol{\beta}_k)^{p(\mathbf{z}_{tu}^n)_k} p(\mathbf{x}_t^n, \mathbf{y}_u^n | \mathbf{z}_{tu}^n, \boldsymbol{\phi}_k)^{z_{tu}^n} p(\boldsymbol{\theta} | \mathcal{M}) d\boldsymbol{\theta}$$

Laplace approximate on each terms and ignore asymptotically small terms w.r.t. N

$$FIC(\mathbf{x}, \mathbf{y} | \mathcal{M}) = \max_{q(\mathbf{z})} \sum_{\mathbf{z}} q(\mathbf{z}) \left(\log p(\mathbf{x}, \mathbf{y}, \mathbf{z} | \bar{\boldsymbol{\theta}}) - \frac{D_{\alpha}}{2} \log N - \sum_{k=1}^K \frac{D_{\beta_k}}{2} \log \left(\sum_{n,t,u} z_{tuk}^n \right) - \sum_{k=1}^K \frac{D_{\phi_k}}{2} \log \left(\sum_{n,t,u} z_{tuk}^n \right) - \log q(\mathbf{z}) \right)$$

FIC Lower Bound

We optimize FICLB through EM-like iteration. This lower-bound is obtained by using 1) linear approximation of logarithm function $\log a \leq \log b - (a - b)/b$ and 2) optimality of ML estimator $\log p(\mathbf{x}, \mathbf{y}, \mathbf{z} | \bar{\boldsymbol{\theta}}) \geq \log p(\mathbf{x}, \mathbf{y}, \mathbf{z} | \boldsymbol{\theta})$

$$FIC \geq FICLB(\mathbf{x}, \mathbf{y}, q, \bar{q}, \boldsymbol{\theta}, \mathcal{M}) = \sum_n \sum_{\mathbf{z}^n} q(\mathbf{z}^n) \left(\log p(\mathbf{x}^n, \mathbf{y}^n, \mathbf{z}^n | \boldsymbol{\theta}) + \sum_{tuk} z_{tuk}^n \log \delta_{tuk} - \log q(\mathbf{z}) \right) - \frac{D_{\alpha}}{2} \log N - \sum_{k=1}^K \frac{D_{\beta_k}}{2} \log \left(\sum_{n,t,u} \bar{q}(z_{tuk}^n) - 1 \right) - \sum_{k=1}^K \frac{D_{\phi_k}}{2} \log \left(\sum_{n,t,u} \bar{q}(z_{tuk}^n) - 1 \right)$$

$$\delta_{tuk} = \begin{cases} \exp \left(- \frac{D_{\phi_k}}{2 \sum_{ntu} z_{tuk}^n} \right) & \text{if } t = T_X \text{ and } u = T_Y \\ \exp \left(- \frac{D_{\phi_k}}{2 \sum_{ntu} \bar{q}(z_{tuk}^n)} - \frac{D_{\beta_k}}{2 \sum_{ntu} \bar{q}(z_{tuk}^n)} \bar{q}(z_{tuk}^n) \right) & \text{otherwise} \end{cases}$$

Heuristic model Pruning

Algorithm 1 The model pruning algorithm

Input: data (\mathbf{x}, \mathbf{y}) , initial model $\mathcal{M} = (K_M, K_X, K_Y)$, initial variational distribution q , initial parameter set $\boldsymbol{\theta}$ and pruning threshold ϵ

while not converged **do**

$\bar{q} \leftarrow q$

$q \leftarrow \arg \max_q FICLB(\mathbf{x}, \mathbf{y}, q, \bar{q}, \boldsymbol{\theta}, \mathcal{M})$

for all k that suffice $\sum_{n,t,u} q(z_{tuk}^n) \leq \epsilon$ **do**

 delete the k -th hidden state of the model \mathcal{M}

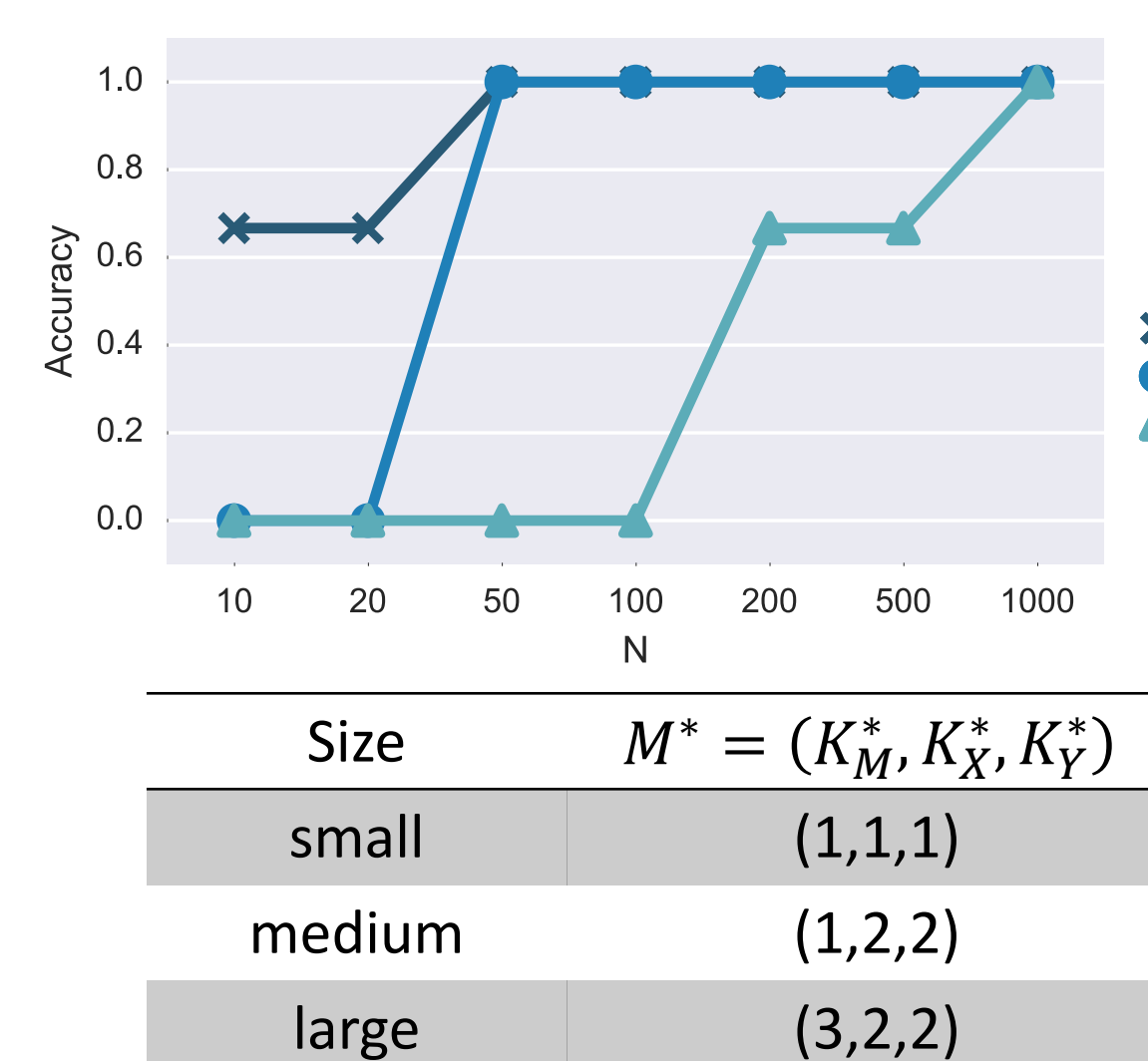
end for

$\boldsymbol{\theta} \leftarrow \arg \max_{\boldsymbol{\theta}} FICLB(\mathbf{x}, \mathbf{y}, q, \bar{q}, \boldsymbol{\theta}, \mathcal{M})$

end while

Experiments

Model selection on synthetic data



Real DNA sequence

