

Introduction

Goal: Estimate the probability of phylogenetic (i.e. evolutionary) trees based on MCMC samples

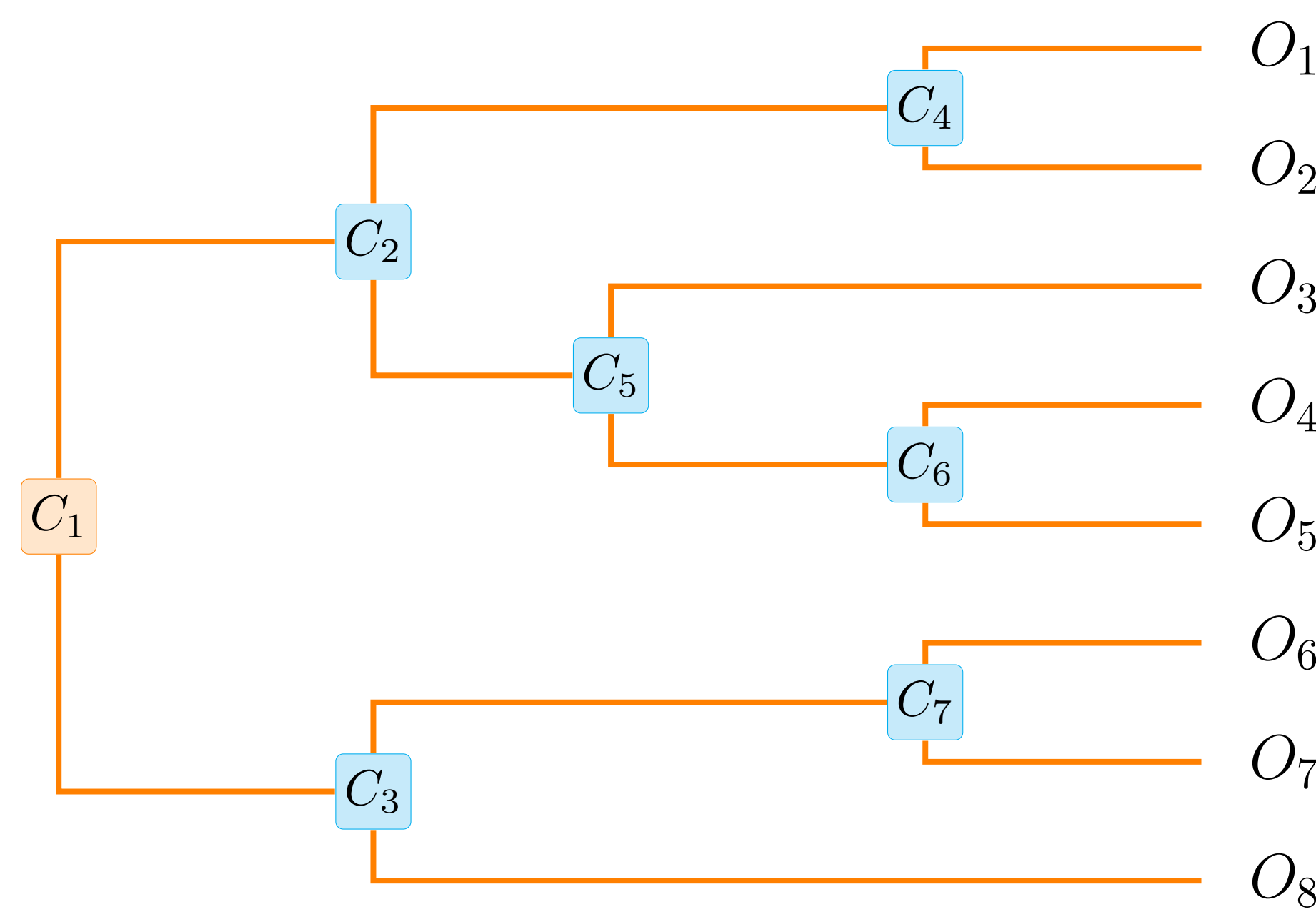
Motivation: Current methods are unsatisfactory

- The common practice of using simple sample relative frequencies (SRF) does not support unsampled trees, and is prone to large variance between different runs
- Previous efforts do extend to unsampled trees, but make too strong assumptions to provide accurate posterior estimation for real data.

By introducing a novel graphical model, **subsplit Bayesian networks**, we propose a general probability estimation framework for phylogenetic trees that

- generalizes to unsampled trees
- provides accurate approximation for real data posteriors

Problem Setup



A phylogenetic tree T is a binary tree with labeled leaves.

- label set $\mathcal{X} = \{O_1, \dots, O_N\}$, each label represents a species.
- A *clade* X is a nonempty subset of \mathcal{X}

Conditional Clade Distribution

- Clade Decomposition (follow the splitting process of the tree).

$$T_C = \{C_2, C_3, C_4, C_5, C_6, C_7\}$$

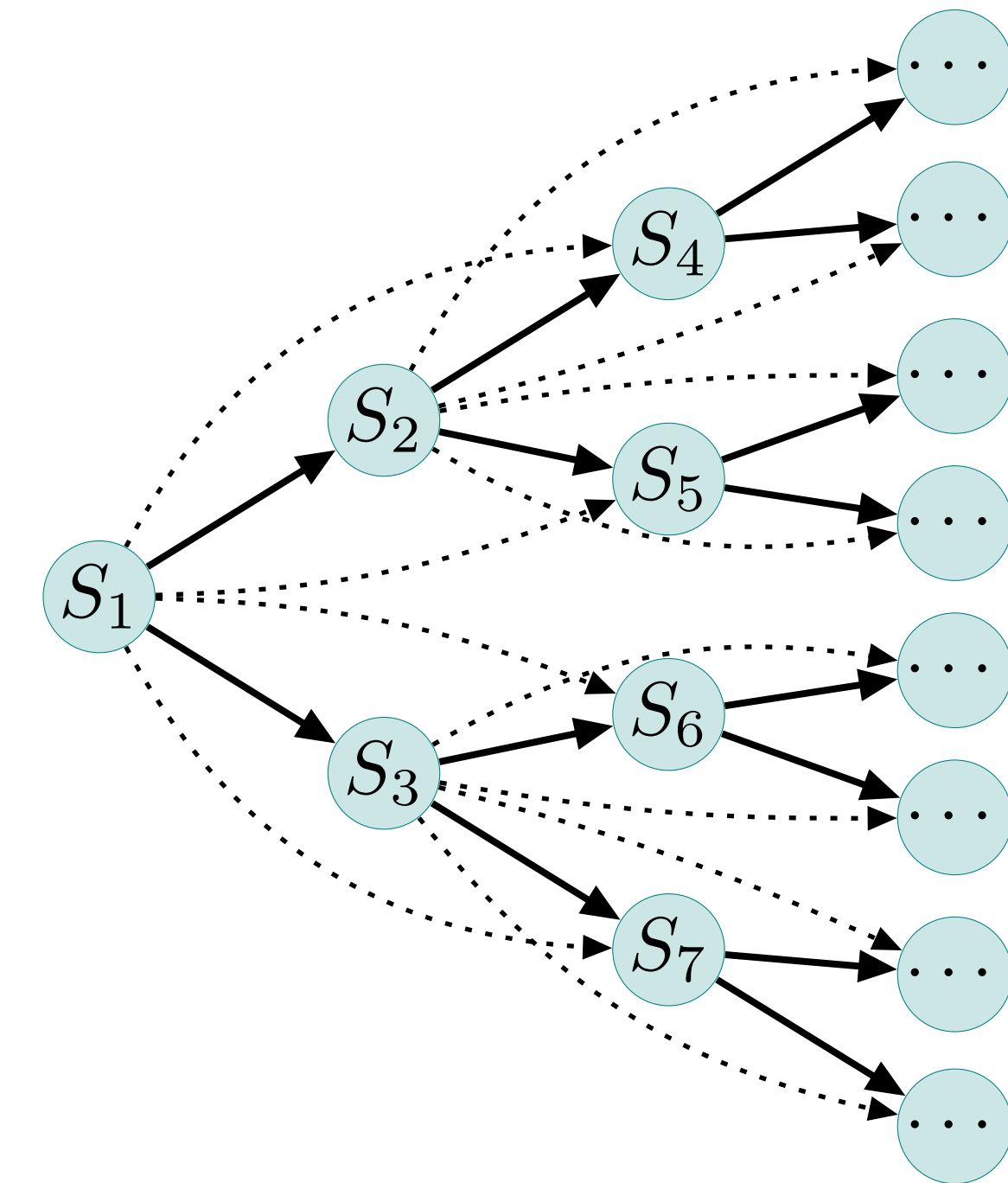
- Conditional Independent Approximation

$$\begin{aligned} p_{\text{cd}}(T) &= p(C_2, C_3, C_4, C_5, C_6, C_7) \\ &= p(C_2, C_3)p(C_4, C_5|C_2)p(C_6|C_5)p(C_7|C_3) \end{aligned}$$

Let \succ be a total order on clades. A **subsplit** (Y, Z) of a clade X is an ordered pair of disjoint subclades of X such that $Y \cup Z = X$, $Y \succ Z$.

- Subsplit Decomposition

$$T_S = \{(C_2, C_3), (C_4, C_5), (\{O_3\}, C_6), (C_7, \{O_8\})\}$$



Subsplit Bayesian Network

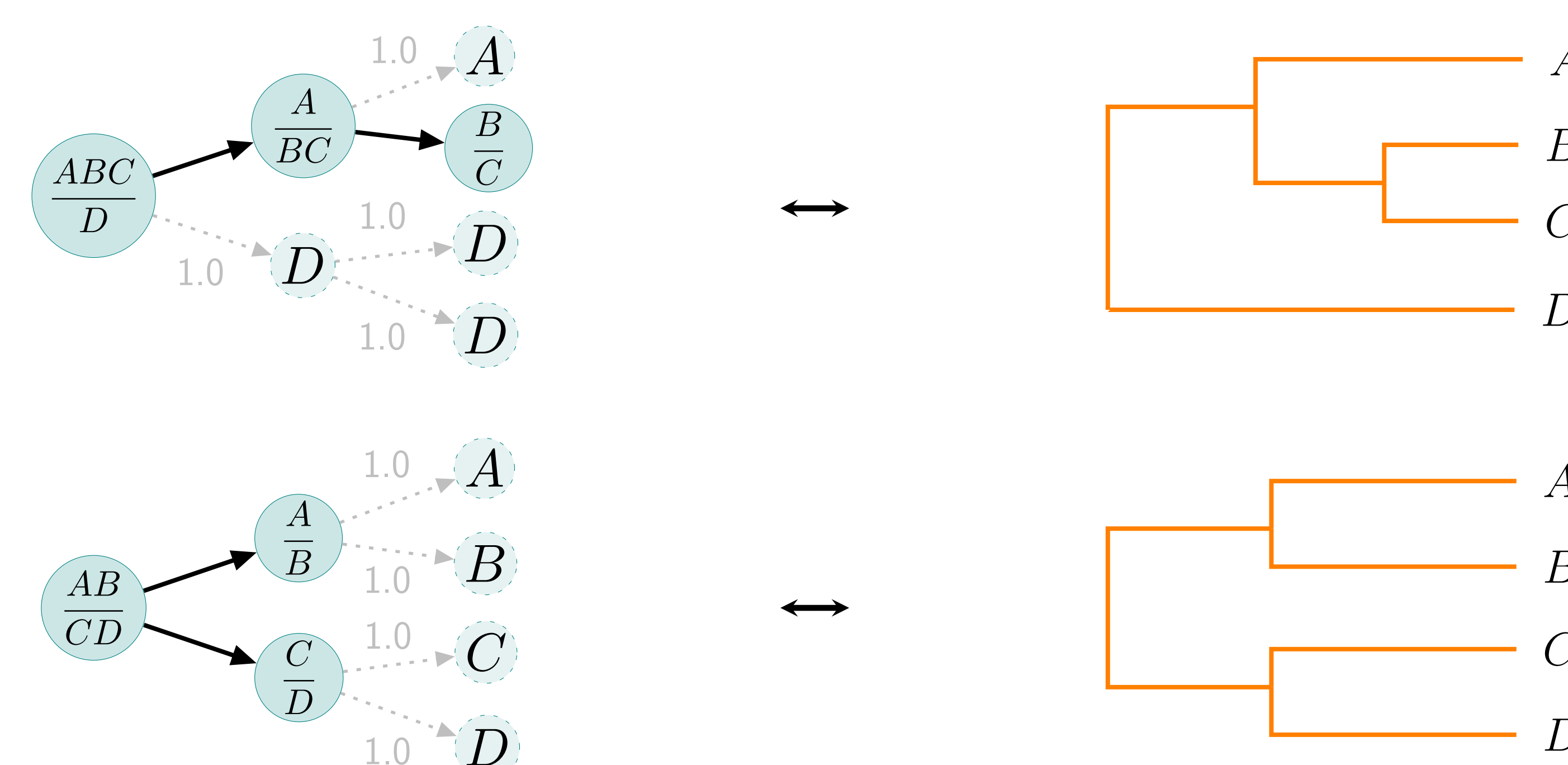
A **subsplit Bayesian network** (SBN) $\mathcal{B}_{\mathcal{X}}$ on a leaf set \mathcal{X} of size N is a Bayesian network of depth $N - 1$ whose nodes take on subsplit or singleton clade values of \mathcal{X} and

- the root node takes on subsplits of the entire leaf set \mathcal{X}
- contains a full and complete binary tree $\mathcal{B}_{\mathcal{X}}^*$

SBNs probability for rooted trees

$$p_{\text{sbn}}(T) = p(S_1) \prod_{i>1} p(S_i | S_{\pi_i})$$

SBNs provide **valid probability distributions** of the entire tree space and are **flexible** to capture complicated dependence structures.



ML for Rooted Trees

Given a sample of rooted trees $\mathcal{D} = \{T_k\}_{k=1}^K$, where $T_k = \{S_i = s_{i,k}, i \geq 1\}$, $k = 1, \dots, K$, the SBN log-likelihood function is

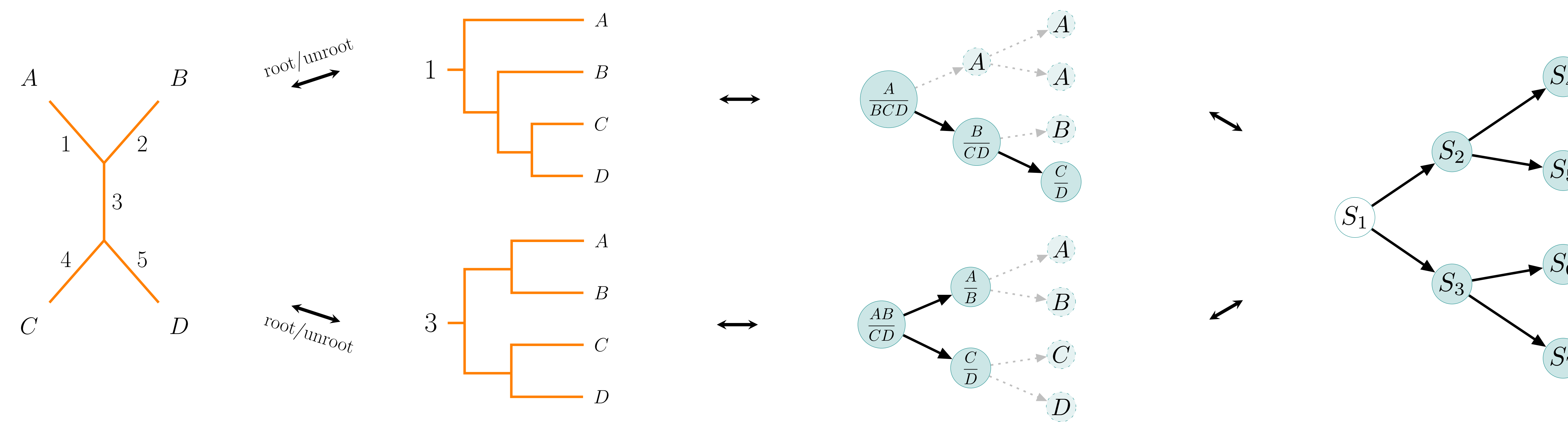
$$\log L(\mathcal{D}) = \sum_{k=1}^K \left(\log p(S_1 = s_{1,k}) + \sum_{i>1} \log p(S_i = s_{i,k} | S_{\pi_i} = s_{\pi_i,k}) \right)$$

Using **conditional probability sharing**, we have

$$\log L(\mathcal{D}) = \sum_{s_1 \in \mathbb{C}_r} m_{s_1} \log p(S_1 = s_1) + \sum_{s|t \in \mathbb{C}_{\text{ch|pa}}} m_{s,t} \log p(s|t)$$

where \mathbb{C}_r denotes the set of all observed root splits of S_1 , $\mathbb{C}_{\text{ch|pa}}$ denotes the set of all observed parent-child subsplit pairs, and $m_{s_1}, m_{s,t}$ denotes the corresponding frequency counts.

Learning SBNs for Unrooted Trees



Lower Bounds Maximization

- Simple Averaging**

$$q(S_1 = s_1) = \frac{1}{2N - 3}, \quad \forall s_1 \sim T^u$$

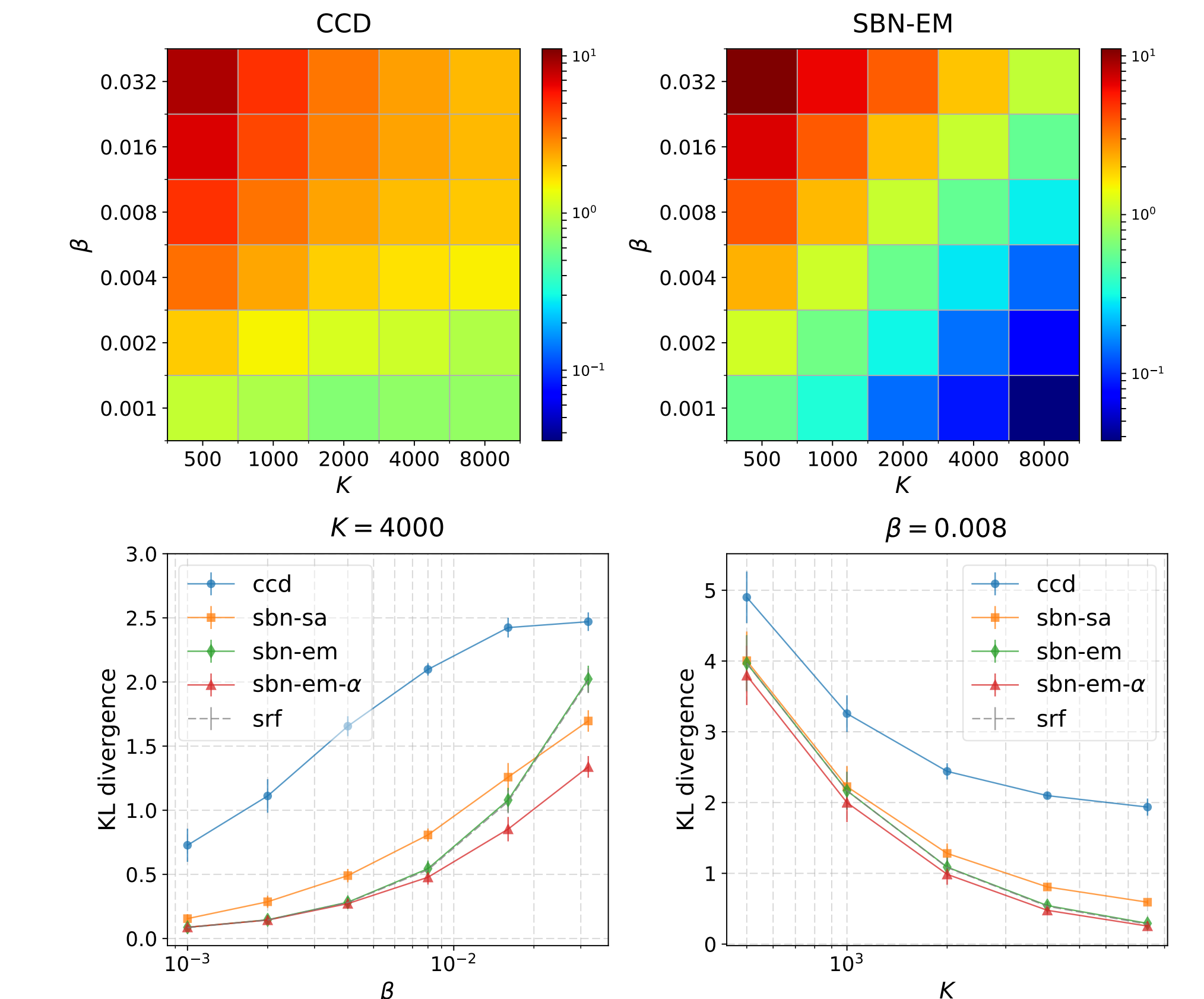
- Expectation Maximization**

$$q^{(n)}(S_1 = s_1) = p(S_1 = s_1 | T^u, \hat{p}^{\text{EM},(n)}), \quad \forall s_1 \sim T^u$$

Remark: can incorporate regularization when data is insufficient or the number of parameters is large.

Experiments

SBN algorithms perform consistently much better than CCD on a challenging tree probability estimation problem with simulated data.



SBN algorithms relax the conditional clade independence assumption and provides accurate approximation in multimodal distributions.

