# Generalizing Tree Probability Estimation via Bayesian Networks

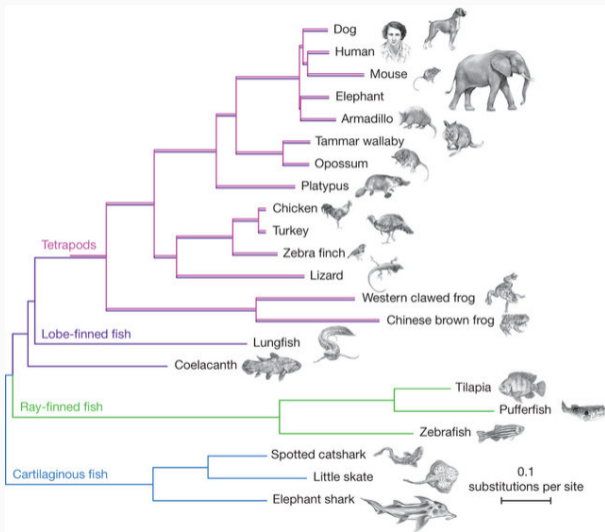Cheng Zhang and Frederick A. Matsen IV

December 1, 2018

Fred Hutchinson Cancer Research Center, Seattle, WA

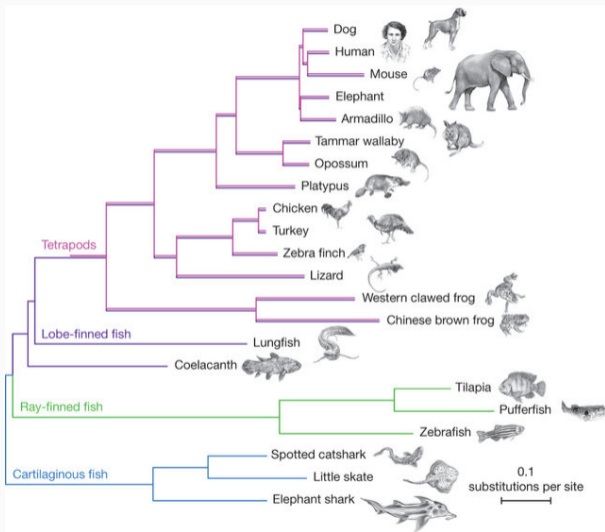**FRED HUTCH**
CURES START HERE™

Tree of life

In **Molecular Evolution**, phylogenetic trees are used to model the evolutionary relationship among various biological species or other entities.
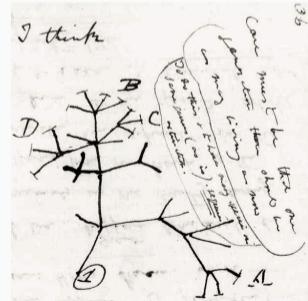
Tree of life

In **Molecular Evolution**, phylogenetic trees are used to model the evolutionary relationship among various biological species or other entities.



from Darwin's Notebook

Markov chain Monte Carlo

Markov chain Monte Carlo

Current approaches are unsatisfactory.

- Sample relative frequencies (**SRF**).
  - Do not generalize!
- Conditional clade distribution (**CCD**).
  - Not flexible enough for real data!

*Markov chain Monte Carlo*

Current approaches are unsatisfactory.

- Sample relative frequencies (**SRF**).
  - Do not generalize!
- Conditional clade distribution (**CCD**).
  - Not flexible enough for real data!

What is the best way to use MCMC samples?

*Markov chain Monte Carlo*

Current approaches are unsatisfactory.

- Sample relative frequencies (**SRF**).
  - Do not generalize!
- Conditional clade distribution (**CCD**).
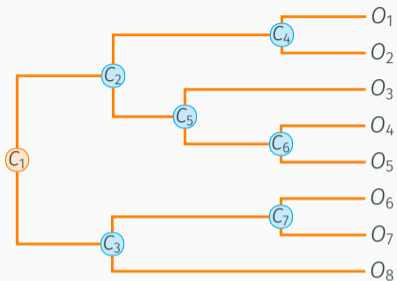  - Not flexible enough for real data!

What is the best way to use MCMC samples?

**Our Contribution**: Subsplit Bayesian Networks. A general probability estimation framework for phylogenetic trees based on MCMC samples that

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

Key: harness the similarity of trees properly.

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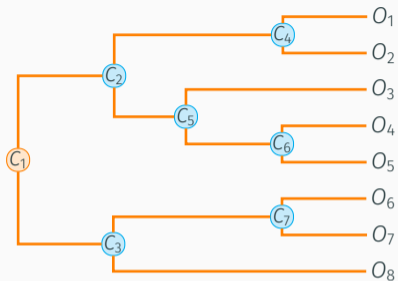- A *clade X* is a nonempty subset of $\mathcal{X}$.

$$C_5 = \{O_3, O_4, O_5\}, \ C_7 = \{O_6, O_7\}.$$

- *Clade Decomposition*

$$T_{\mathcal{C}} = \{C_2, C_3, C_4, C_5, C_6, C_7\}$$

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

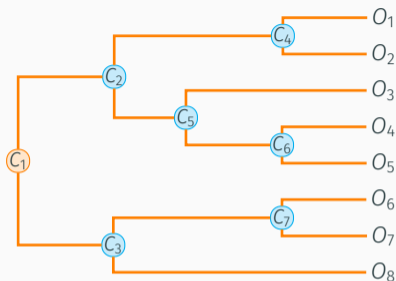$$C_5 = \{O_3, O_4, O_5\}, \ C_7 = \{O_6, O_7\}.$$

- *Clade Decomposition*

$$T_{\mathcal{C}} = \{C_2, C_3, C_4, C_5, C_6, C_7\}$$

A *subsplit* of a clade *X* is an ordered pair of disjoint subclades $(Y, Z)$ such that $Y \cup Z = X$, $Y \succ Z$. Examples: $C_1 \rightarrow (C_2, C_3)$, $C_2 \rightarrow (C_4, C_5)$.

## Subsplit Decomposition

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

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

$$C_5 = \{O_3, O_4, O_5\}, \ C_7 = \{O_6, O_7\}.$$

- *Clade Decomposition*
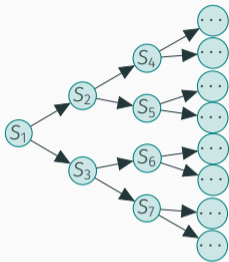
$$T_{\mathcal{C}} = \{C_2, C_3, C_4, C_5, C_6, C_7\}$$

A *subsplit* of a clade $X$ is an ordered pair of disjoint subclades $(Y, Z)$ such that $Y \cup Z = X$, $Y \succ Z$. Examples: $C_1 \rightarrow (C_2, C_3)$, $C_2 \rightarrow (C_4, C_5)$.

## Subsplit Decomposition

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

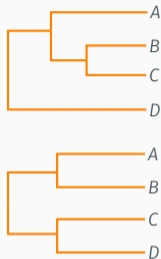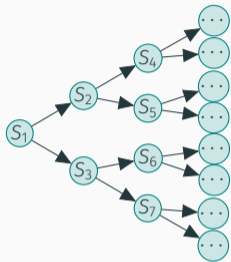$$p(T) = p(C_2, C_3)p(C_4, C_5 | C_2, C_3)p(C_6 | C_4, C_5)p(C_7 | C_2, C_3)$$

A **Subsplit Bayesian Network** on a leaf set $\mathcal{X}$ of size $N$ is a Bayesian network

- nodes take on subsplit / singleton clade values.
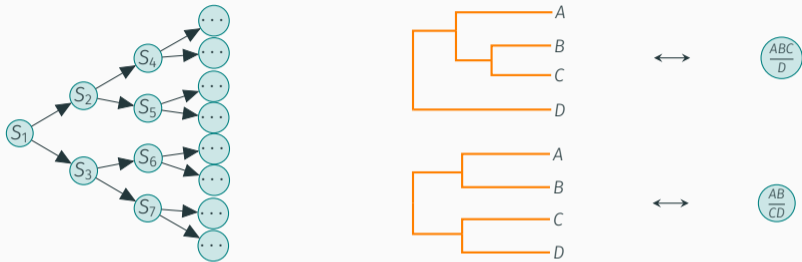- contains a full and complete binary tree.

A **Subsplit Bayesian Network** on a leaf set $\mathcal{X}$ of size $N$ is a Bayesian network

- nodes take on subsplit / singleton clade values.
- contains a full and complete binary tree.
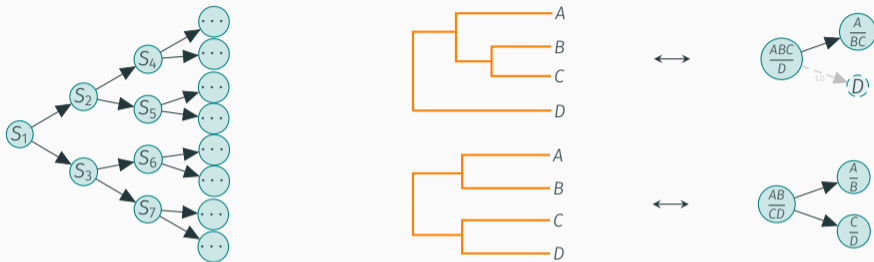
A **Subsplit Bayesian Network** on a leaf set $\mathcal{X}$ of size $N$ is a Bayesian network
- nodes take on subsplit / singleton clade values.
- contains a full and complete binary tree.
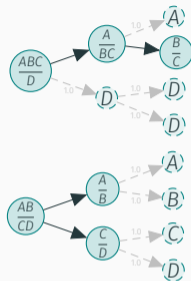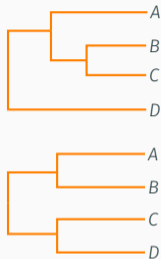
A **Subsplit Bayesian Network** on a leaf set $\mathcal{X}$ of size $N$ is a Bayesian network
- nodes take on subsplit / singleton clade values.
- contains a full and complete binary tree.
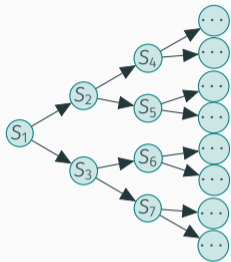
A **Subsplit Bayesian Network** on a leaf set $\mathcal{X}$ of size $N$ is a Bayesian network

- nodes take on subsplit / singleton clade values.
- contains a full and complete binary tree.
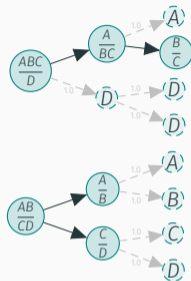
A **Subsplit Bayesian Network** on a leaf set $\mathcal{X}$ of size $N$ is a Bayesian network

- nodes take on subsplit / singleton clade values.
- contains a full and complete binary tree.
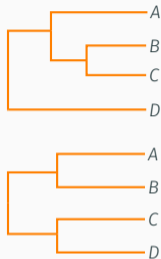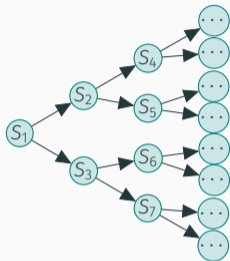
SBN probability for rooted trees

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

A **Subsplit Bayesian Network** on a leaf set $\mathcal{X}$ of size $N$ is a Bayesian network
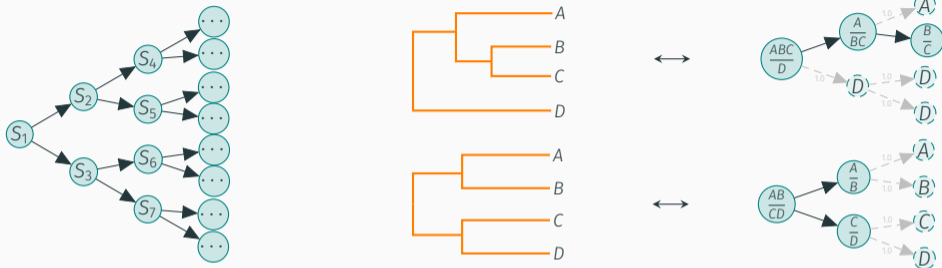
- nodes take on subsplit / singleton clade values.
- contains a full and complete binary tree.

SBN probability for rooted trees

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

SBNs provide valid probability distributions

A **Subsplit Bayesian Network** on a leaf set $\mathcal{X}$ of size $N$ is a Bayesian network
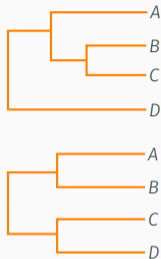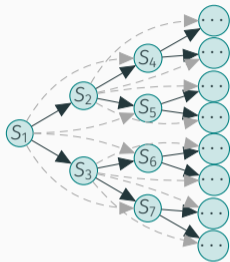
- nodes take on subsplit / singleton clade values.
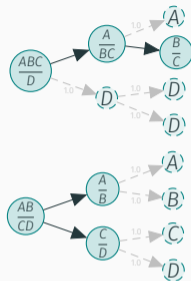- contains a full and complete binary tree.

SBN probability for rooted trees

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

SBNs provide valid probability distributions and are flexible.

### Rooted Trees
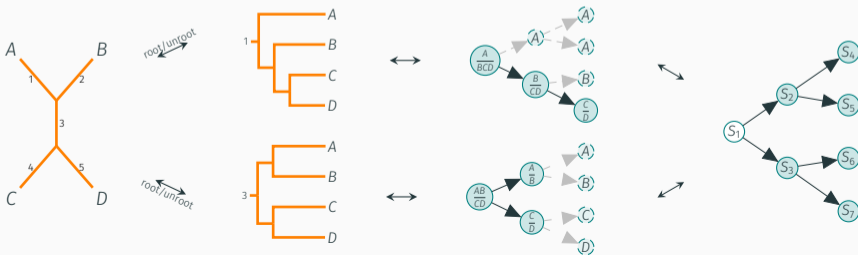
- maximum likelihood

## Rooted Trees

- maximum likelihood

## Unrooted Trees

## Rooted Trees
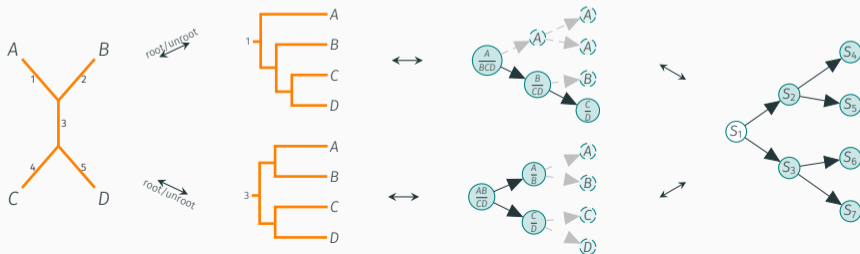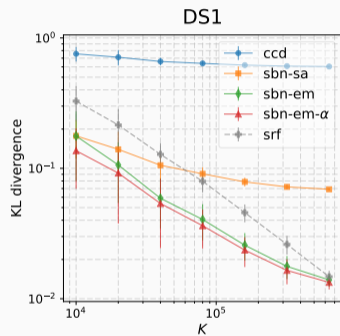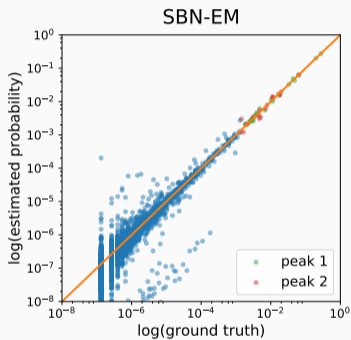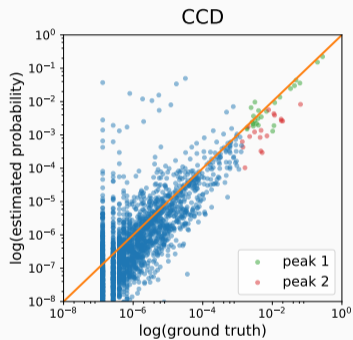
- maximum likelihood

## Unrooted Trees

- Expectation Maximization
- simple averaging lower bound maximization
- incorporate regularization when necessary

A real data set with multimodal posterior

# Experiments

| Data set | (#Taxa, #Sites) | Tree space size | Sampled trees | KL divergence to ground truth | | | | |
|----------|-----------------|-----------------|---------------|-------|-------|--------|--------|------------|
| | | | | SRF | CCD | SBN-SA | SBN-EM | SBN-EM-$\alpha$ |
| DS1 | (27, 1949) | $5.84 \times 10^{32}$ | 1228 | 0.0155 | 0.6027 | 0.0687 | 0.0136 | **0.0130** |
| DS2 | (29, 2520) | $1.58 \times 10^{35}$ | 7 | **0.0122** | 0.0218 | 0.0218 | 0.0199 | 0.0128 |
| DS3 | (36, 1812) | $4.89 \times 10^{47}$ | 43 | 0.3539 | 0.2074 | 0.1152 | 0.1243 | **0.0882** |
| DS4 | (41, 1137) | $1.01 \times 10^{57}$ | 828 | 0.5322 | 0.1952 | 0.1021 | 0.0763 | **0.0637** |
| DS5 | (50, 378) | $2.84 \times 10^{74}$ | 33752 | 11.5746 | 1.3272 | 0.8952 | 0.8599 | **0.8218** |
| DS6 | (50, 1133) | $2.84 \times 10^{74}$ | 35407 | 10.0159 | 0.4526 | **0.2613** | 0.3016 | 0.2786 |
| DS7 | (59, 1824) | $4.36 \times 10^{92}$ | 1125 | 1.2765 | 0.3292 | 0.2341 | 0.0483 | **0.0399** |
| DS8 | (64, 1008) | $1.04 \times 10^{103}$ | 3067 | 2.1653 | 0.4149 | 0.2212 | 0.1415 | **0.1236** |

# Poster # 123

- We proposed a general framework for tree probability estimation based on subsplit Bayesian networks.
- SBNs exploit the similarity among trees to provide flexible probability estimators that generalize to unsampled trees.
- Future work
  - extends to general trees
  - structure learning of SBNs
  - deeper investigation on the effect of parameter sharing
  - applications in other probabilistic learning problems in tree spaces (e.g., MCMC transition kernel design and variational inference)