# Improved Importance Sampling of Phylogenies

Mathias C. Cronjäger

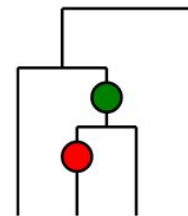Department of Statistics
University of Oxford

# Outline

# Setup

The data:
Single Nucleotide Polymorphisms (SNPs) in $n$ aligned genetic sequences sampled in the present.

The model:
Individuals are leaves on an an unobserved random tree with the most recent common ancestor at the root (a coalescent). New segregating sites occur along branches at rate $\theta$.

The challenge:
Computing likelihoods does not scale as $n$ increases (there are too many possible trees to consider)



$$\mathbb{P}\left(\begin{array}{c}\includegraphics\end{array}\right) = \sum_{\Psi \text{ genealogical history of}} \mathbb{P}(\Psi)$$

# Importance sampling of ancestral histories

$$\mathbb{P}\left(\vcenter{\hbox{\includegraphics{tree}}}\right) = \sum_{x \in H(\vcenter{\hbox{\tiny tree}})} \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} \mathbb{Q}(x) = \mathbb{E}_{X \sim \mathbb{Q}}\left[\frac{\mathbb{P}(X)}{\mathbb{Q}(X)}\right] \approx \frac{1}{N} \sum_{i=1}^{N} \frac{\mathbb{P}(X_i)}{\mathbb{Q}(X_i)}$$
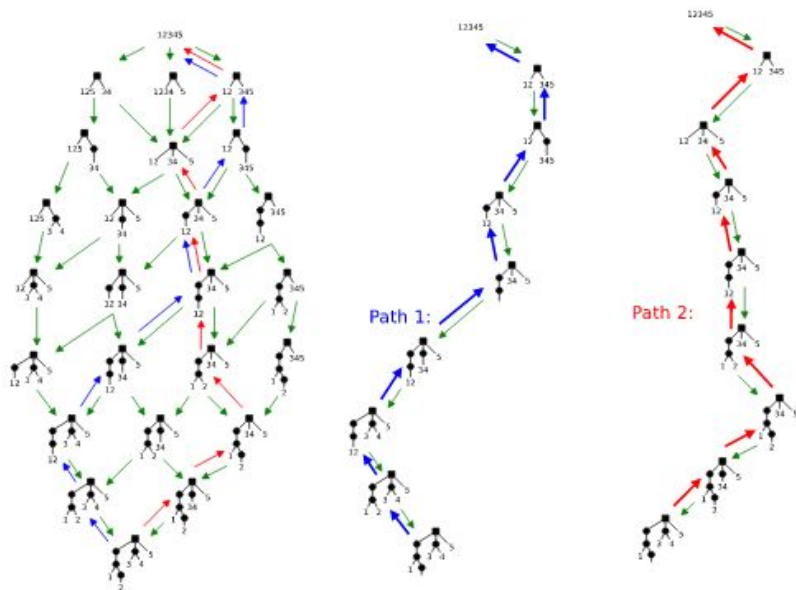


Path 1:

Path 2:

Ideally, the proposal distribution **Q**, should

1. Approximate the target distribution **P** well
2. Sampling $X \sim$ **Q** should be fast
3. Computing weights $f_\mathbf{Q}(X)$ should be fast

Conditions 2 & 3 allow us to pick large $N$ (number of particles).
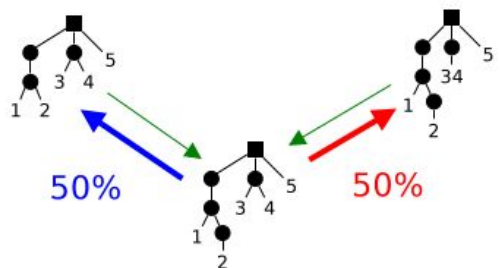Condition 1 gives us better convergence in $N$.

<u>For this problem</u>: Known (feasible) proposals are **sequential**. Sampling from **Q** is done by constructing paths step-by-step, **from the bottom up**.
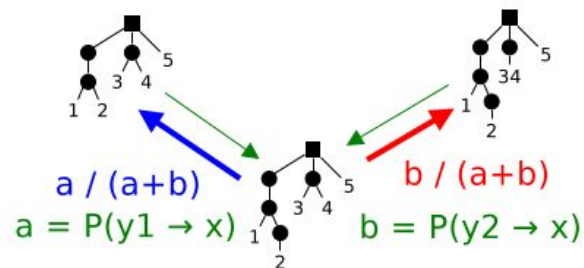
e.g.  $\mathbb{P}\left(\vcenter{\hbox{\tiny tree}}\right) \approx \frac{1}{2}\left(\frac{\mathbb{P}(\text{Path 1})}{\mathbb{Q}(\text{Path 1})} + \frac{\mathbb{P}(\text{Path 2})}{\mathbb{Q}(\text{Path 2})}\right)$

4

# Existing proposals (a single step)
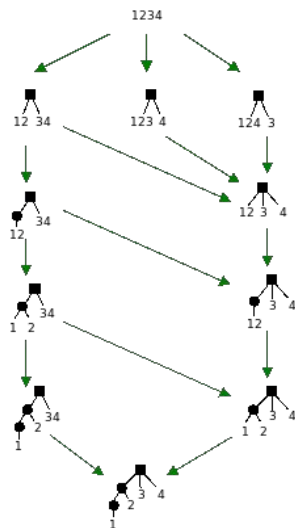


Stephens and Donnelly proposal

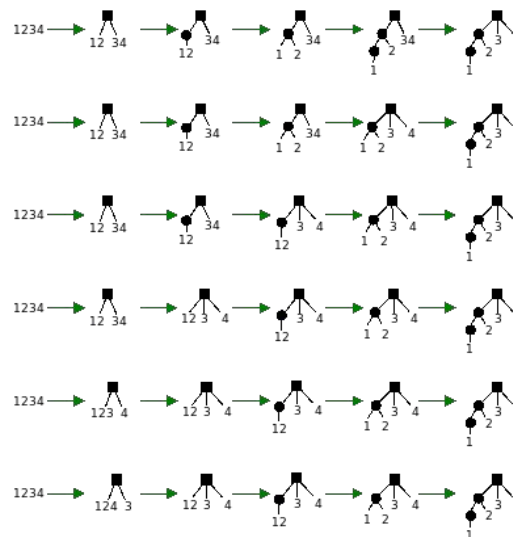Griffiths and Taveré (G&T)

# Sequential proposals and path density bias

Consider for example $\mathbf{Q}_{SD}$ which samples ancestral histories as follows:

1. Start with state $x_0$ = Observed data
2. Given a partial path $p = [x_0, \dots , x_i]$, pick a state $x_{i+1}$ uniformly at random from the predecessors of $x_i$, and update the partial path to $p' = [x_0, \dots , x_i , x_{i+1}]$
3. Terminate when hitting a state with no ancestors.

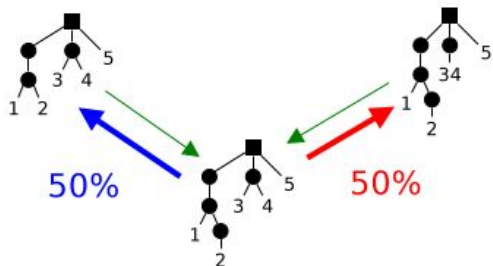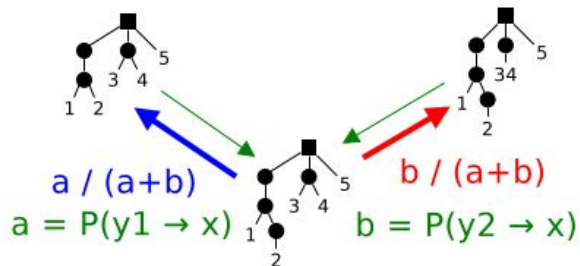Major Issue: Oversampling of paths containing nodes of low in-degree.



| Path | Weight |
|---|---|
| | $1/2$ |
| | $1/4$ |
| | $1/8$ |
| | $1/24$ |
| | $1/24$ |
| | $1/24$ |

# Combinatorial importance sampling



7

# Likelihoods (toy example 1: 4 seq, 2 sites)

# Likelihoods (toy example 2: 5 seq; 4 sites)



9

# Thank you.
# Any Questions?

I also have a poster →

Collaborators (and supervisors)

Jotun Hein
(Oxford, Statistics)

Paul Jenkins
(Warwick, Statistics & CS)

Code: https://github.com/Cronjaeger/combIS

Get in touch:
Email: cronjager@stats.ox.ac.uk
Twitter: @mcCronjaeger