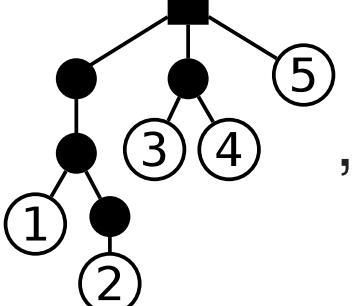


MOTIVATION: HOW COMPLEX ARE PROBABILITY RECURSIONS?

Given a set of aligned sequences, e.g.

Sequence 1: ... T ... G ... A ... A ...
Sequence 2: ... T ... G ... A ... G ...
Sequence 3: ... A ... A ... T ... A ...
Sequence 4: ... A ... A ... T ... A ...
Sequence 5: ... A ... A ... A ... A ...



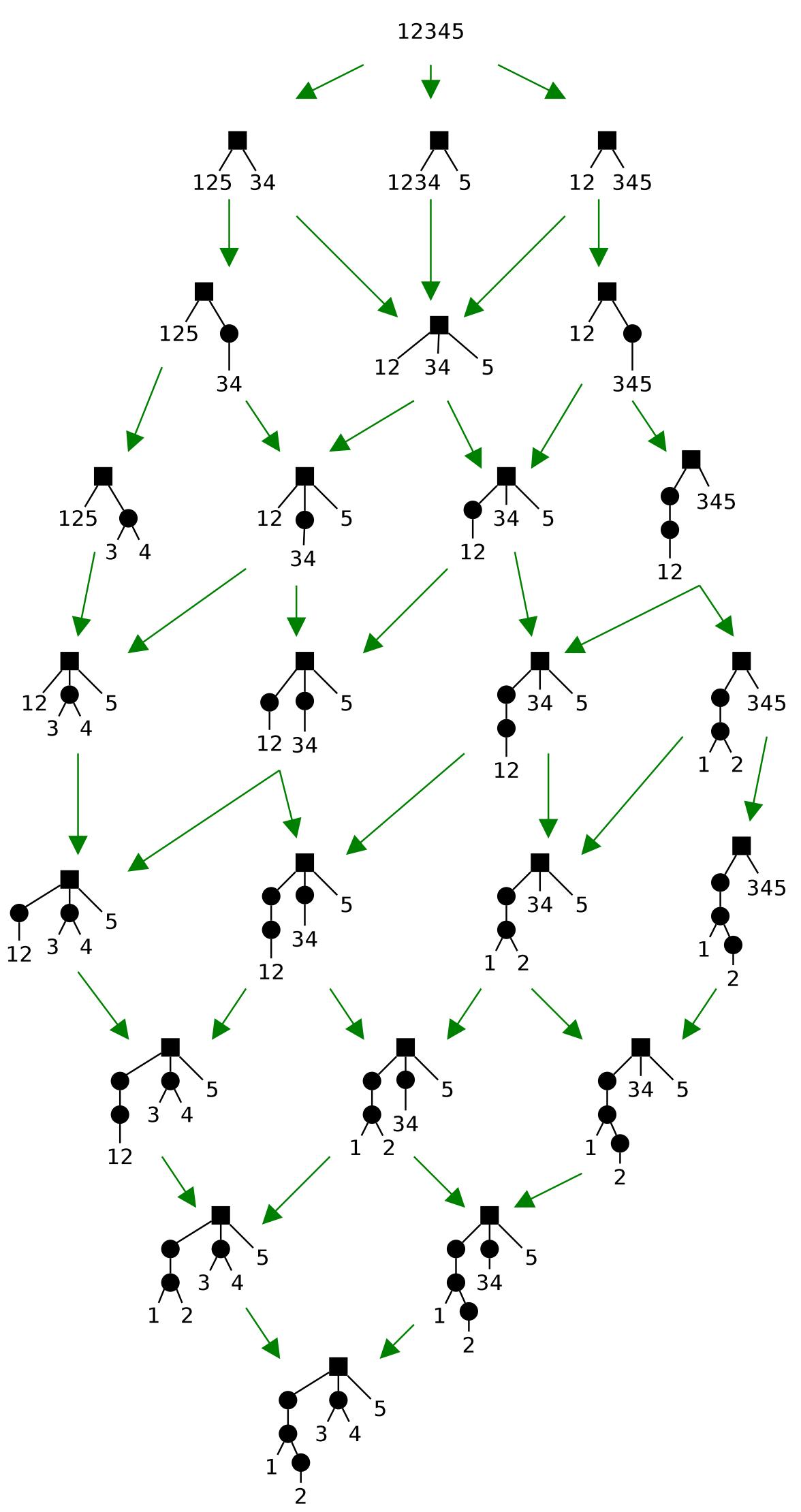
the probability of having evolved from some initial sequence—here ... A ... A ... A ... A ...—may be expressed:

$$\mathbb{P}\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array}\right) = \sum_{x \in \text{Histories}} \mathbb{P}(x)$$

which in turn may be expressed by conditioning on the most recent event:

$$\mathbb{P}\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array}\right) = \frac{1}{\theta+4} \mathbb{P}\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \end{array}\right) + \frac{\theta}{\theta+4} \frac{1}{5} \mathbb{P}\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array}\right).$$

If we apply this conditioning-trick recursively, computing $\mathbb{P}\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array}\right)$ reduces to computing a weighted sum over all paths from “12345” to “ $\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \end{array}$ ” in the graph below:



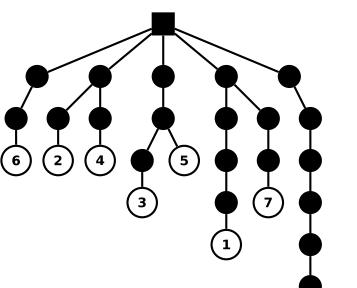
We may note that

- number of nodes = number of distinct terms in \mathbb{P} -recursion,
- number of paths “12345 → ... → $\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \end{array}$ ” = number of execution paths when evaluating $\mathbb{P}\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array}\right)$ via tail-recursion (without memoization).

We refer to the nodes in the above graph as **ancestral states**, and the paths “12345 → ... → $\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \end{array}$ ” as **ancestral histories**.

CENTRAL PROBLEM: COUNTING QUICKLY

The number of ancestral states and histories grows very quickly. Whereas $\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array}$ has 71 histories and 25 ancestral states (including itself), the dataset



for example already has 89,325,103,544,240,200 histories and 12,356 ancestral states. Since we do not know the graph a-priori our main aim is to **count the number of ancestral states and/or histories just by considering the initial data** (i.e. without having to determine the entire ancestral graph).

HOW TO COUNT (ABSTRACTLY)

We encode non-plane trees as nested systems of sets, e.g.

$$\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array} = \{\{\{1, \{2\}\}, \{3, 4\}\}, 5\}.$$

Given a tree $T = \{T_1, \dots, T_r\}$, we may determine the number of ancestral histories and ancestral states—denoted $h(T)$ and $a(T)$ respectively—as follows:

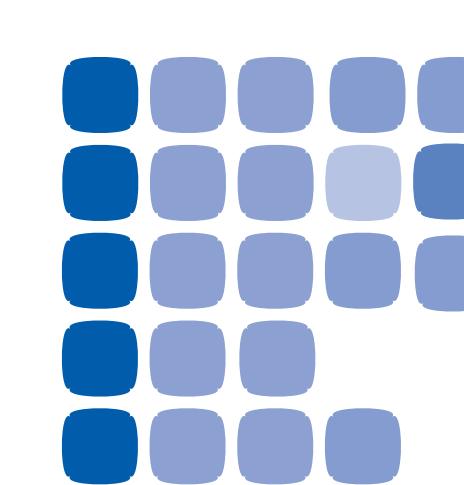
- If T is a leaf (i.e. has no children), $a(T) = h(T) = 1$ holds.
- If $T = \{T'\}$ (i.e. T has root degree 1), $a(T) = 1 + a(T')$ and $h(T) = h(T')$ holds.
- If $T = \{T_1, \dots, T_r\}$ holds for $r > 1$, then:

$$a'(T) = 1 + \sum_{S \subseteq [k]} \left(\sum_{\lambda \in S} \frac{|S|!}{\prod_i \lambda_i! \prod_{i \notin S} \alpha_i!} \right) \prod_{i \notin S} a'(T_i)$$

$$h(T) = \sum_{S \subseteq [r]} h(\{T_i \mid i \in S\}) h(\{T_i \mid i \notin S\}) \binom{\sum_{i=1}^r k_i - 2}{\sum_{i \in S} k_i - 1}$$

Whereby $a'(T') = a(T') + \mathbb{1}(\deg(\text{root}(T')) > 1)$ and k_i denotes the number of nodes in sub-tree T_i .

How many ways can we explain the same data?



HOW TO COUNT – A SIMPLE EXAMPLE

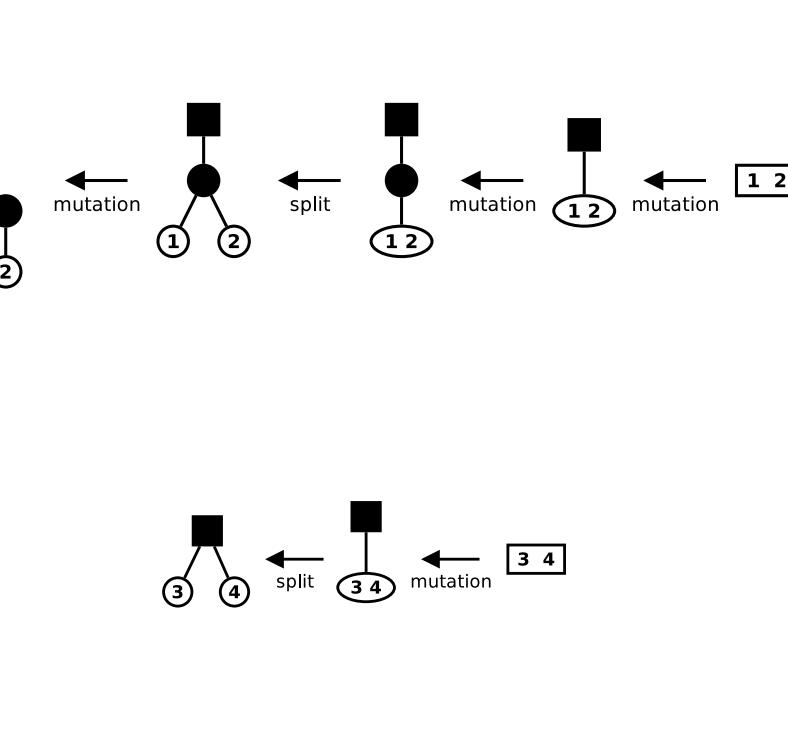
1. We decompose our gene-tree into sub-trees:

$$\begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array} = \left\{ \begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array}, \begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \\ \textcircled{5} \end{array} \right\}$$

2. We consider the problem for each sub-tree separately:

T_{Sub}	$h(T_{\text{Sub}})$	$a'(T_{\text{Sub}})$
$\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}$	1	5
$\begin{array}{c} \textcircled{3} \\ \textcircled{4} \end{array}$	1	3
$\begin{array}{c} \textcircled{5} \end{array}$	1	1

Ancestors of T_{Sub} (verification)



Computing $h\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right)$: we consider all three ways of ordering mergers immediately below the root:

$$h\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right) = h\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right) + h\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right) + h\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right)$$

Each of the right hand terms can then be computed with relative ease:

$$h(T_1) \left(h(T_2) h(T_3) \binom{3+1-2}{3-1} \right) \binom{5+3+1-2}{5-1} = 35$$

$$h(T_2) \left(h(T_1) h(T_3) \binom{5+1-2}{5-1} \right) \binom{5+3+1-2}{3-1} = 21$$

$$h(T_3) \left(h(T_1) h(T_2) \binom{5+3-2}{5-1} \right) \binom{5+3+1-2}{1-1} = 15$$

Which implies that $h\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right) = 35 + 21 + 15 = 71$ holds.

Computing $a\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right)$: Here, we first compute $a'\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right)$ using our recursion:

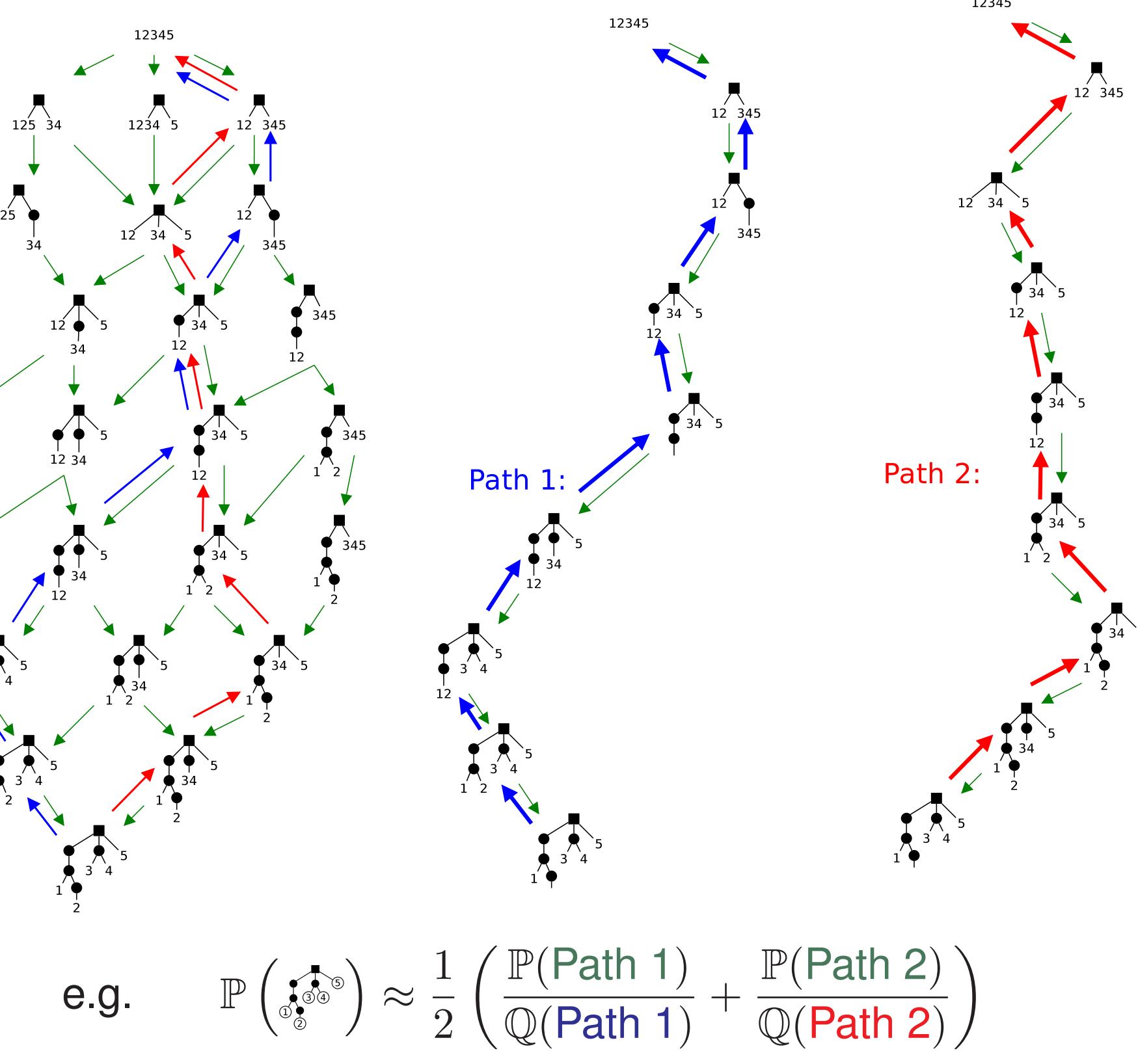
$$a'\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right) = 1 + \underbrace{5 \cdot 3 \cdot 1}_{S=\emptyset} + \underbrace{5}_{S=\{2,3\}} + \underbrace{3}_{S=\{1,3\}} + \underbrace{1}_{S=\{1,2\}} + \underbrace{1}_{S=\{1,2,3\}}$$

It therefore follows that $a\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right) = a'\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right) - 1 = 25$

APPLICATION: COMBINATORIAL IMPORTANCE SAMPLER

We can approximate probabilities of aligned sequences—e.g. $\mathbb{P}\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right)$ —by sampling ancestral histories $X_1, \dots, X_N \stackrel{iid}{\sim} \mathbb{Q} \ll \mathbb{P}$ and relying on the following approximation:

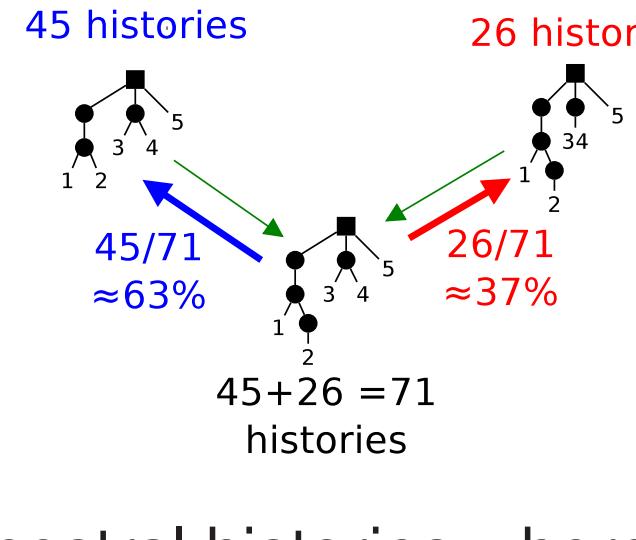
$$\mathbb{P}\left(\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}\right) = \sum_{x \in H} \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} \mathbb{Q}(x) = \mathbb{E}_{X \sim \mathbb{Q}} \left[\frac{\mathbb{P}(X)}{\mathbb{Q}(X)} \right] \approx \frac{1}{N} \sum_{i=1}^N \frac{\mathbb{P}(X_i)}{\mathbb{Q}(X_i)}$$



For this approach to work effectively, \mathbb{Q} should satisfy:

1. \mathbb{Q} must approximate \mathbb{P} well on the space of histories;
2. sampling $X_i \sim \mathbb{Q}$ should be fast;
3. computing $\mathbb{Q}(X_i)$ should be fast.

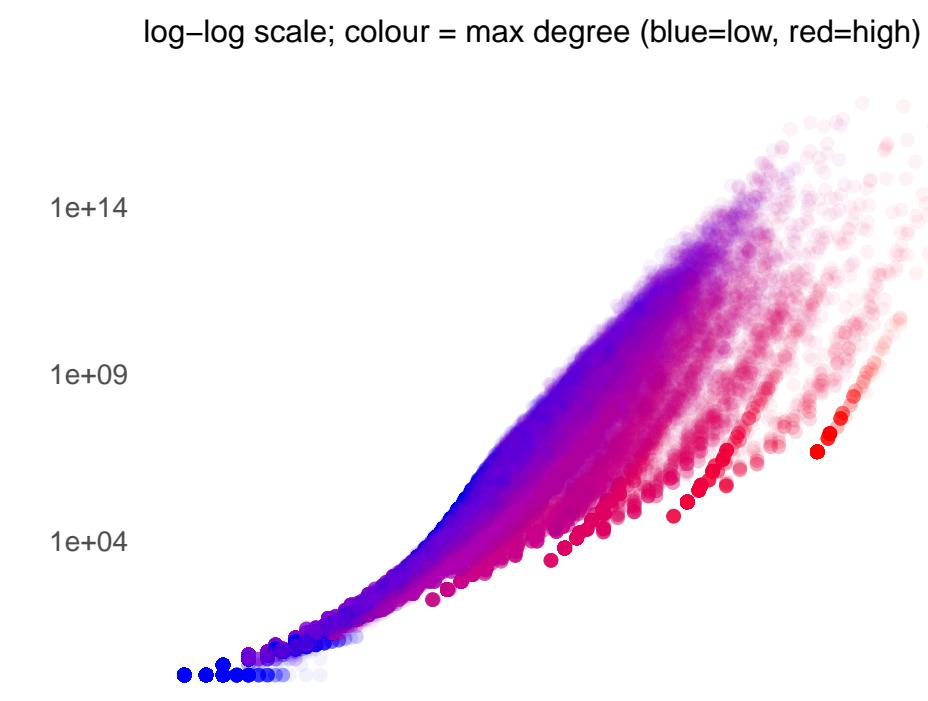
Our scheme *uniformly samples* elements from the space of ancestral histories by generating them “from the bottom up”. To sample an ancestral history of $\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}$, there are for example two “first steps” to pick from:



Since $\begin{array}{c} \textcircled{1} \\ \textcircled{2} \end{array}$ has 45 ancestral histories whereas $\begin{array}{c} \textcircled{1} \end{array}$ has only 26. Our sampler will therefore generate paths that “go left” as a first step with probability 45/71 ≈ 63%.

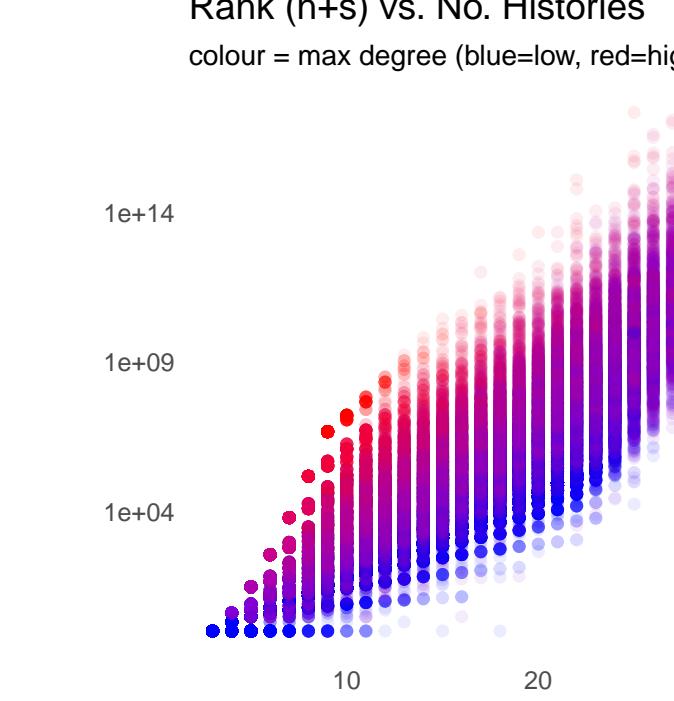
COMPUTATIONAL EXPERIMENTS

No. Ancestral States vs. No. Histories



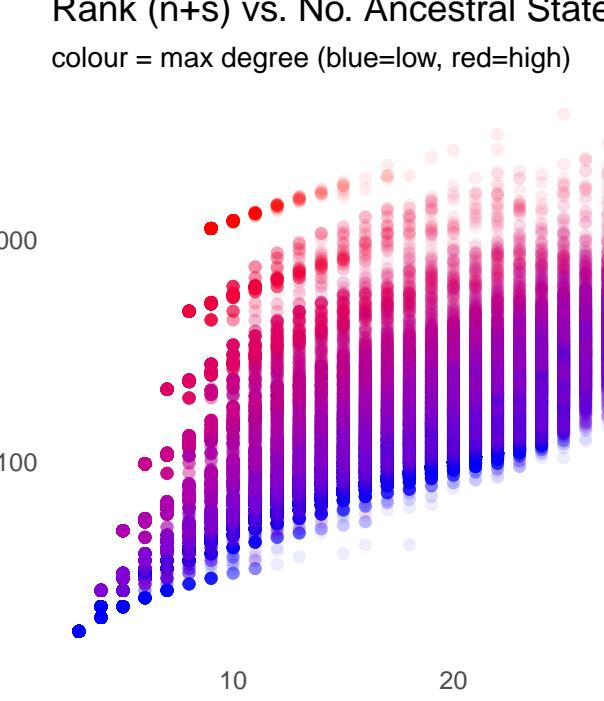
Counts for 70,000 simulated datasets (generated using ms); 500 per choice of $n = 2, \dots, 8$ sequences and $s = 1, \dots, 20$ segregating sites

Rank (n+s) vs. No. Histories

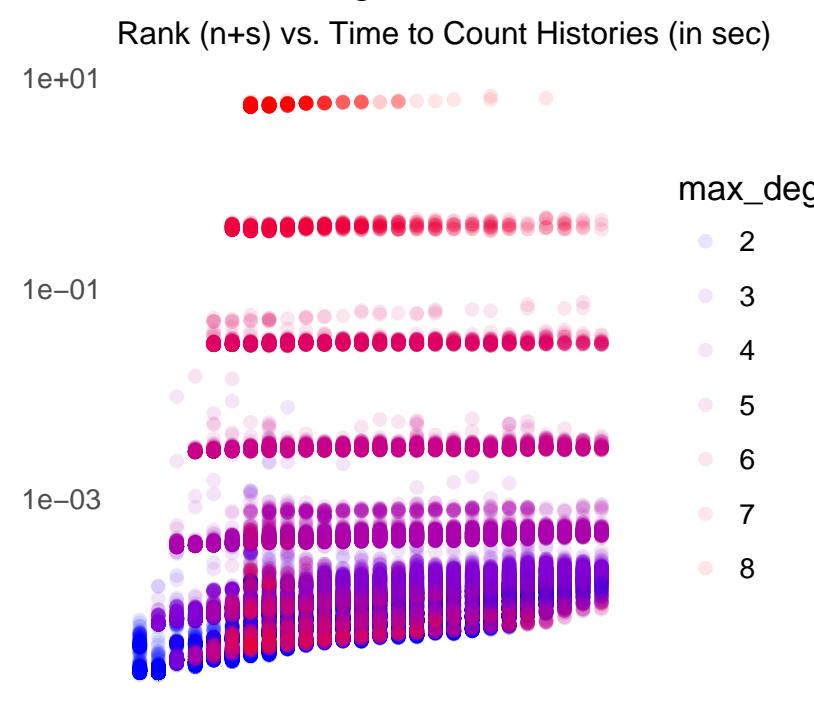


A natural measure of how complex a dataset is is the number of evolutionary events ($n + s - 1$). We see that some degree of proportionality exists (on a log-scale) and note that the maximal degree of the gene tree seems very important as well.

Rank (n+s) vs. No. Ancestral States



How Counting Scales



The number of events in the histories of a dataset is of little importance relative to the maximal degree of the gene tree

n	m	$n=2$	$n=3$	$n=4$	$n=5$	$n=6$	$n=7$	$n=8$
s = 1	1	4	13	30	630	12,600	340,200	7,200
s = 2	5	5	20	108	1,190	10,692	151,800	3,800
s = 3	6	6	21	119	14,540	226,800	3,200	700
s = 4	10	45	215	1,840	21,426	328,320	4,000	900
s = 5	18	70	460	3,024	35,736	452,790	5,000	1,000
s = 6	35	161	770	5,213	56,240	870,396	6,000	1,200
s = 7	56	252	1,270	13,860	124,514	1,804,374	7,000	1,400
s = 8	126	588	3,927	30,800				