

## INTRODUCTION

Two popular models of mutation in population-genetics are the infinite-sites and finite-sites models; outlined below. We here present a model of mutation which admits computations of likelihoods in a manner similar to the former model, while retaining the flexibility of the latter.

**Finite Sites Model (FSM)**  
 Mutations affect sites distributed uniformly on  $\{1, \dots, T\}$ . State of population of size  $n$  encoded by  $S \in \{A, G, C, T\}^{n \times T}$ .

- + Can encode any set of sequences; easy to communicate.
- Computing likelihoods scales abysmally; even for "simple" datasets.

**Infinite Sites Model (ISM)**  
 Mutations affect sites distributed uniformly on  $[0, 1]$ . For a countable number of mutations, we only need to track if a site is segregating or not. State of population encoded by  $S \in \{0, 1\}^{n \times L}$ , where  $L \leq T$  denotes the number of segregating sites

- + Compact representation of data when  $L \ll T$ .
- + Each step will almost surely lead to a lower-rank-state.<sup>a</sup>
- + Small state-space  $\Rightarrow$  traversing state-space easier  $\Rightarrow$  recursive computation of likelihoods more feasible.
- Restrictive modeling assumptions (violated in practice).
- Unable to account for sites with  $>2$  observed nucleotides.

Our **Almost Infinite Sites Model** bridges the gap between these two models of mutation, by applying FSM-analogues of the simplifying assumptions of the ISM—notably the indistinguishability of sites. By imposing an upper bound,  $b \geq L$ , on the total number of mutations in a genealogical history, the number of possible ancestral histories of a population may be kept finite, and likelihoods may be recursively computed.

<sup>a</sup> $\text{rank}(S)$  is the sum of the number of distinct haplotypes in  $S$  and the number of segregating sites in  $S$ .

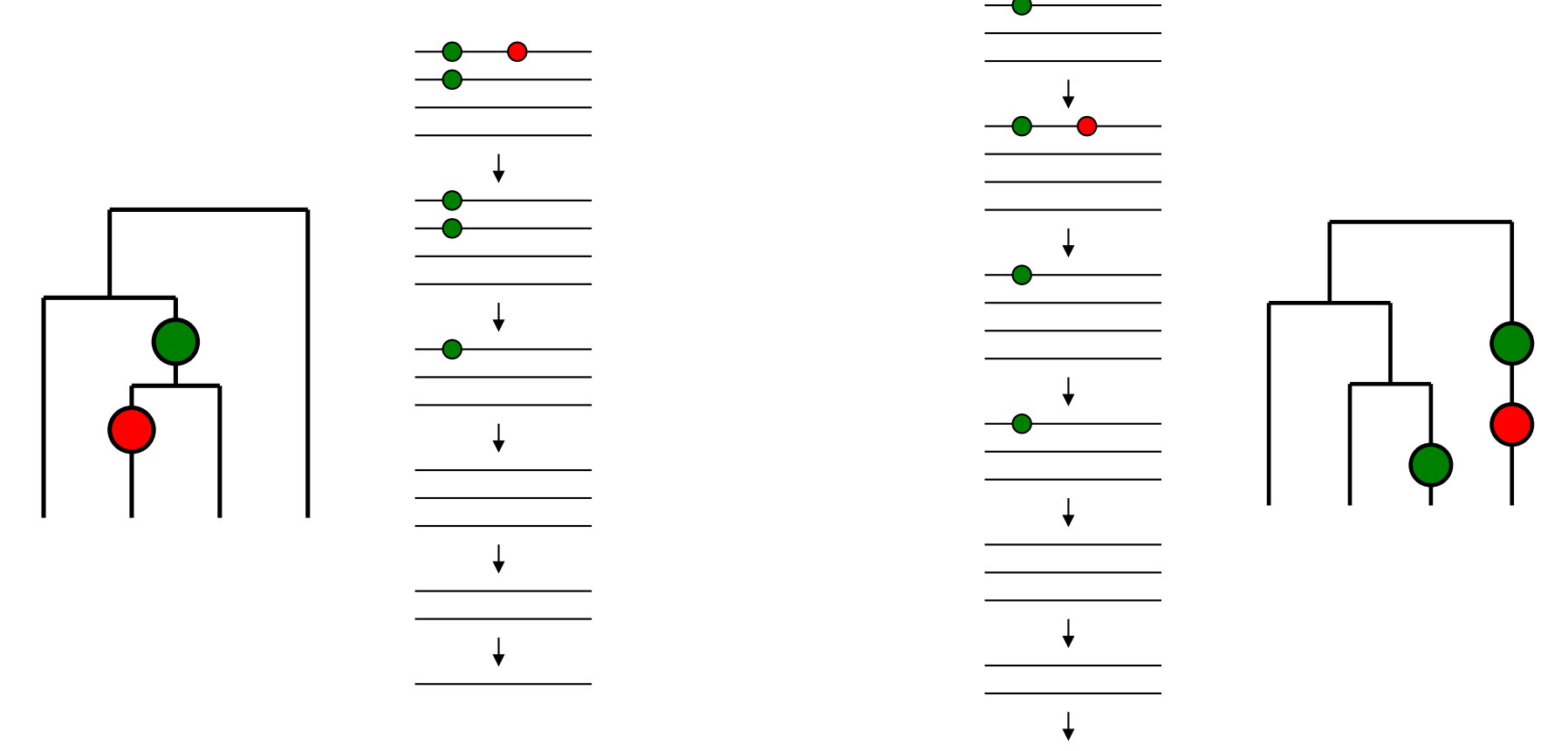
## ENCODING SEQUENCES

Since populations and mutations are presumed to be exchangeable, it is worth noting that  $S_1$  and  $S_2$  encode the same state (denoted  $S_1 \sim S_2$ ) if and only if permutation-matrices  $P_r, P_c$  exist satisfying  $S_1 = P_r S_2 P_c$ . Hence we distinguish between *encodings* of states  $S$  (not unique), and *states*  $\psi = [S]_{\sim}$  (unique).



A sequence of states is referred to as *admissible* if it can be constructed from a coalescent process undergoing mutation by removing events "from the bottom up".

An admissible sequence starting from  $\psi$  and terminating in a single sequence is referred to as a *genealogical history* of  $\psi$ .



Two distinct genealogical histories of the same state (along with coalescents demonstrating admissibility). Note that the right genealogical history is not admissible under the ISM.

## RECURSION – COMPUTING LIKELIHOODS

By summing the probabilities of all genealogical histories of a dataset, we obtain its associated likelihood. Using the Markovian structure of the coalescent, we may express the likelihood recursively.

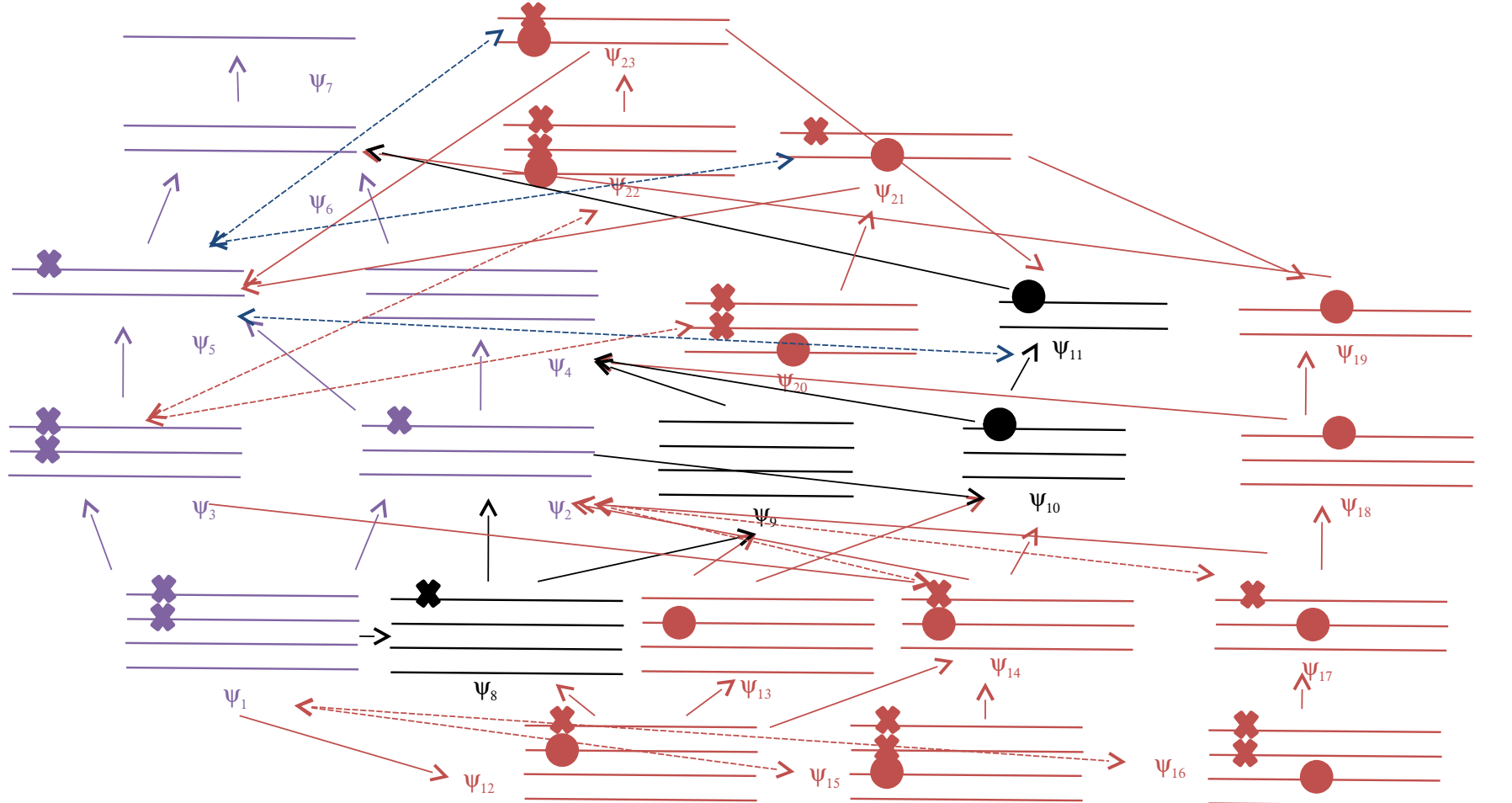
$$\mathbb{P}(\text{---}) = \sum_{\Psi \text{ potential history of ---}} \mathbb{P}(\Psi), \quad \mathbb{P}(\psi) = \sum_{\psi': \mathbb{P}(\psi' \rightarrow \psi) > 0} \mathbb{P}(\psi') \mathbb{P}(\psi' \rightarrow \psi)$$

We aim to determine  $\mathbb{P}(\psi, B \leq b)$ ; the joint probability of observing sample configuration  $\psi = [S]_{\sim} = (A, \mathbf{n}^r, \mathbf{n}^c)$  and the number of mutations in the genealogical history of  $\psi$  (denoted by  $B$ ) being bounded from above by  $b > 0$ . To this end we make the following assumptions:

- Mergers in the genealogical history of  $\psi$  are governed by an underlying Kingman-coalescent process.
- Mutation-events occur at rate  $\frac{\theta}{2}$ , and affect a site uniformly chosen from  $\{1 \dots T\}$

Taking a "forward-backward" approach similar to Griffiths (1989), we obtain the following recurrence:

$$\mathbb{P}((A, \mathbf{n}^r, \mathbf{n}^c), B \leq b) = \frac{N-1}{N-1+\theta} \sum_{i: n_i^r > 1} \frac{n_i^r - 1}{N-1} \mathbb{P}((A, \mathbf{n}^r - e_i, \mathbf{n}^c), B \leq b) + \frac{\theta}{N-1+\theta} \sum_{\substack{1 \leq i \leq L \\ 1 \leq j \leq L \\ x \in \{A, G, T, C\}}} \frac{N_{i,j,x}^r(A, \mathbf{n}^r) N_{i,j,x}^c(A, \mathbf{n}^c)}{N} P_{x, \alpha, i, j} \mathbb{P}(\mathcal{M}_{i,j,x}(A, \mathbf{n}^r, \mathbf{n}^c), B \leq b-1)$$



Ancestral Configurations with upper bounds  $b = 2, 3, 4$  (purple, black and red respectively) for a simple dataset.

## ALGORITHM – SOLVING THE RECURSION (1)

Although the recursion (1) can be used directly for computing  $\mathbb{P}(\psi, B \leq b)$  via naïve tail-recursion, this is highly inefficient. We explore the following approach based on storing all values already computed.

```

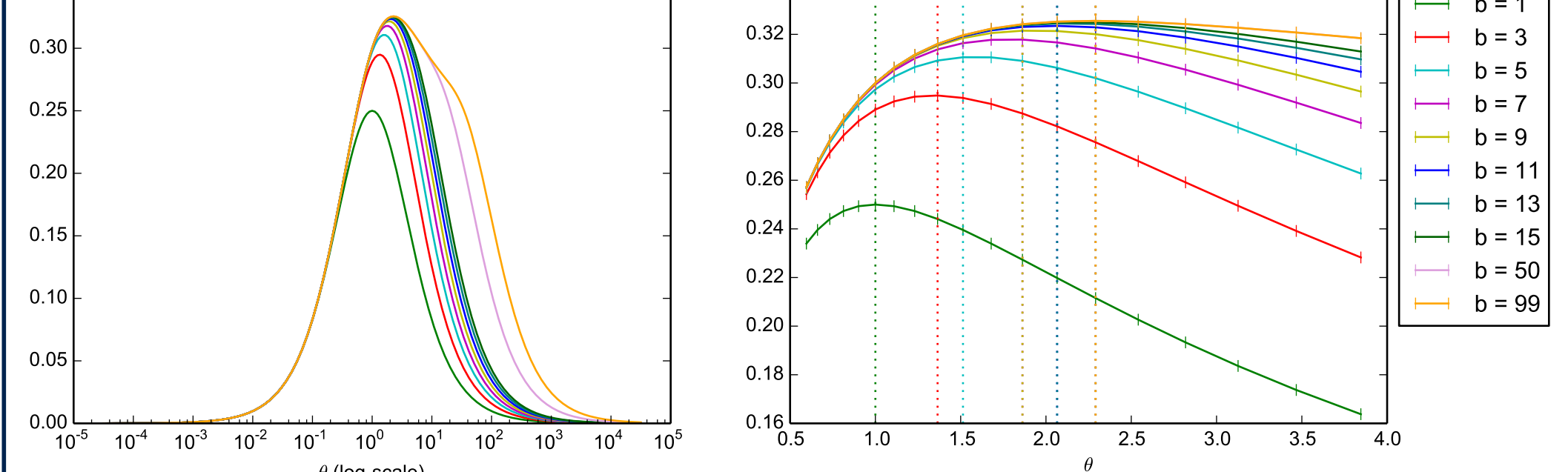
Input S, b, and model-specific constants for (1)
Output P([S]_{\sim}, B \le b)
Initialise empty hash table H
if S = [0] (No segregating sites; 1 active lineage)
  Return 1
Else if \exists S' \sim S : (S', b) \in Keys(H)
  Look up P([S]_{\sim}, B \le b) = H(S', b)
  Return P([S]_{\sim}, B \le b)
Else if b \ge b_{min}([S]_{\sim})
  Compute P([S]_{\sim}, B \le b) using recursion in formula (1)
  Add key-value pair ((S, b), P([S]_{\sim}, B \le b)) to H
Else
  Return 0
End
  
```

The algorithm has been implemented in Python; source, code is available at: <https://github.com/Cronjaeger/almost-infinite-sites-recursions>

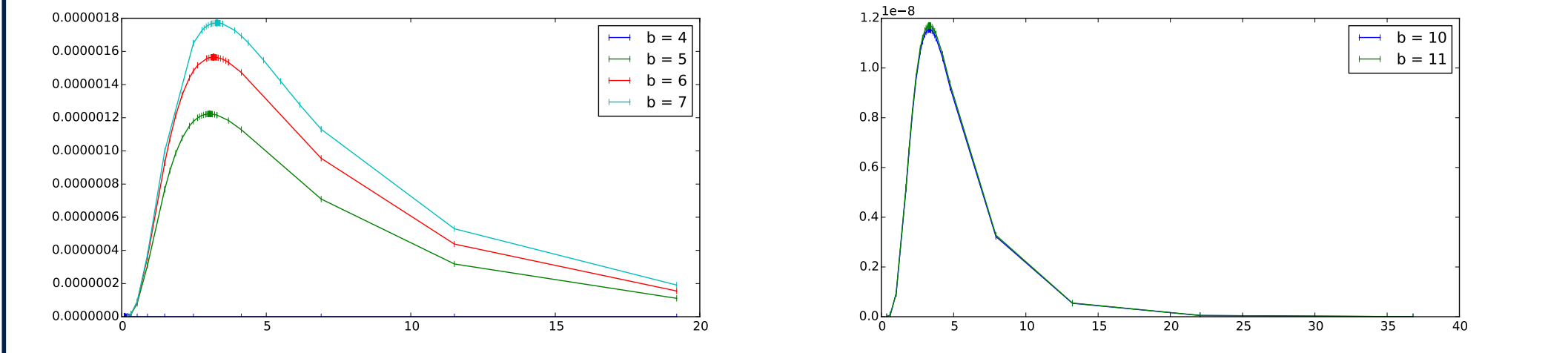
<sup>a</sup>Checking this condition quickly is non-trivial. In our implementation it is done by applying a  $\sim$ -invariant hash-function and resolving collisions by checking if suitable  $P_r, P_c$  can be constructed using a backtracking algorithm.

## RESULTS & BENCHMARKING

We may estimate  $\hat{\theta}_{ML}$  numerically using recursion (1).

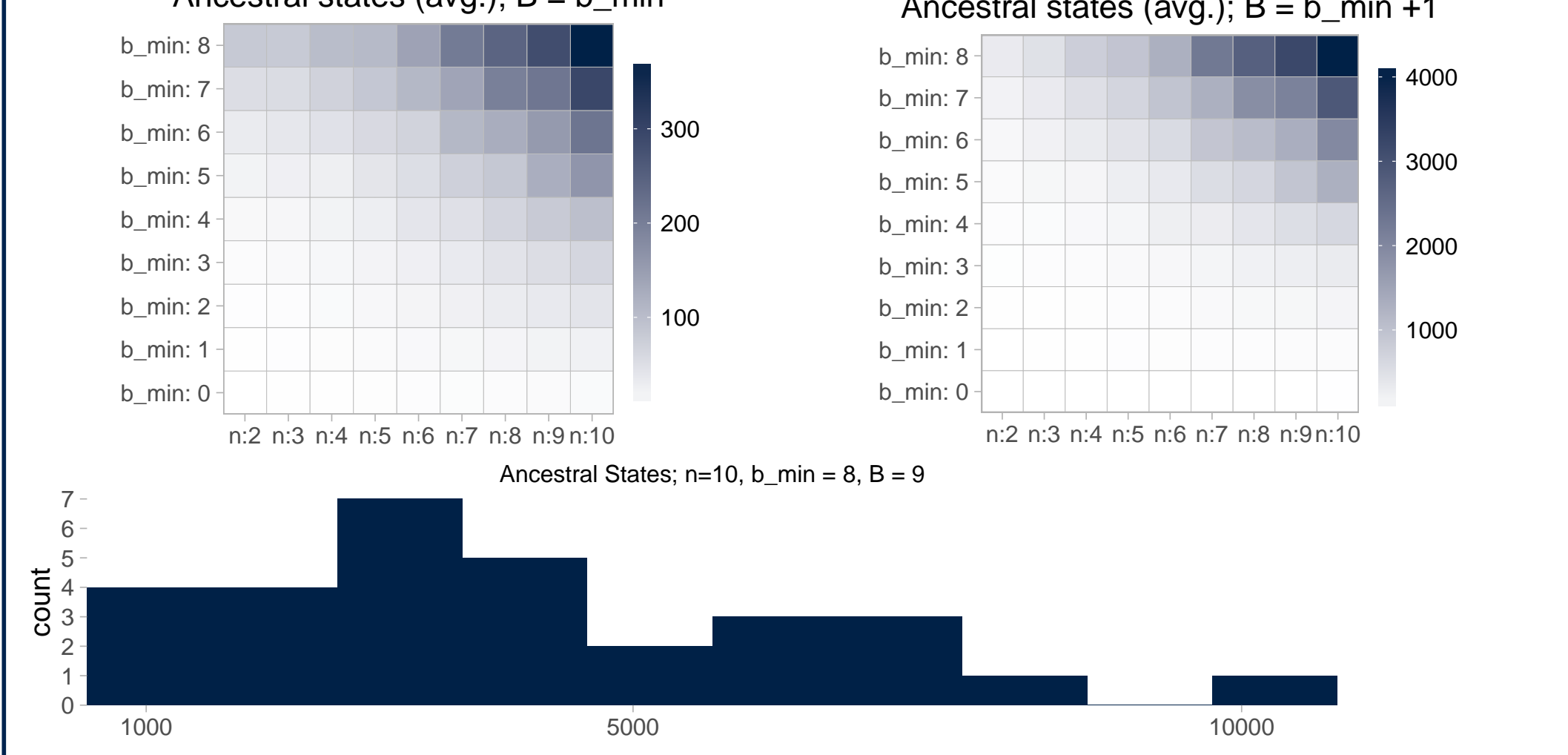


Likelihoods for a very simple dataset  $S = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$  and accompanying MLE-estimates for varying  $b$ .



The number of non-segregating sites impacts the mass gap significantly: on the left 50% of sites are segregating; on the right 1%.

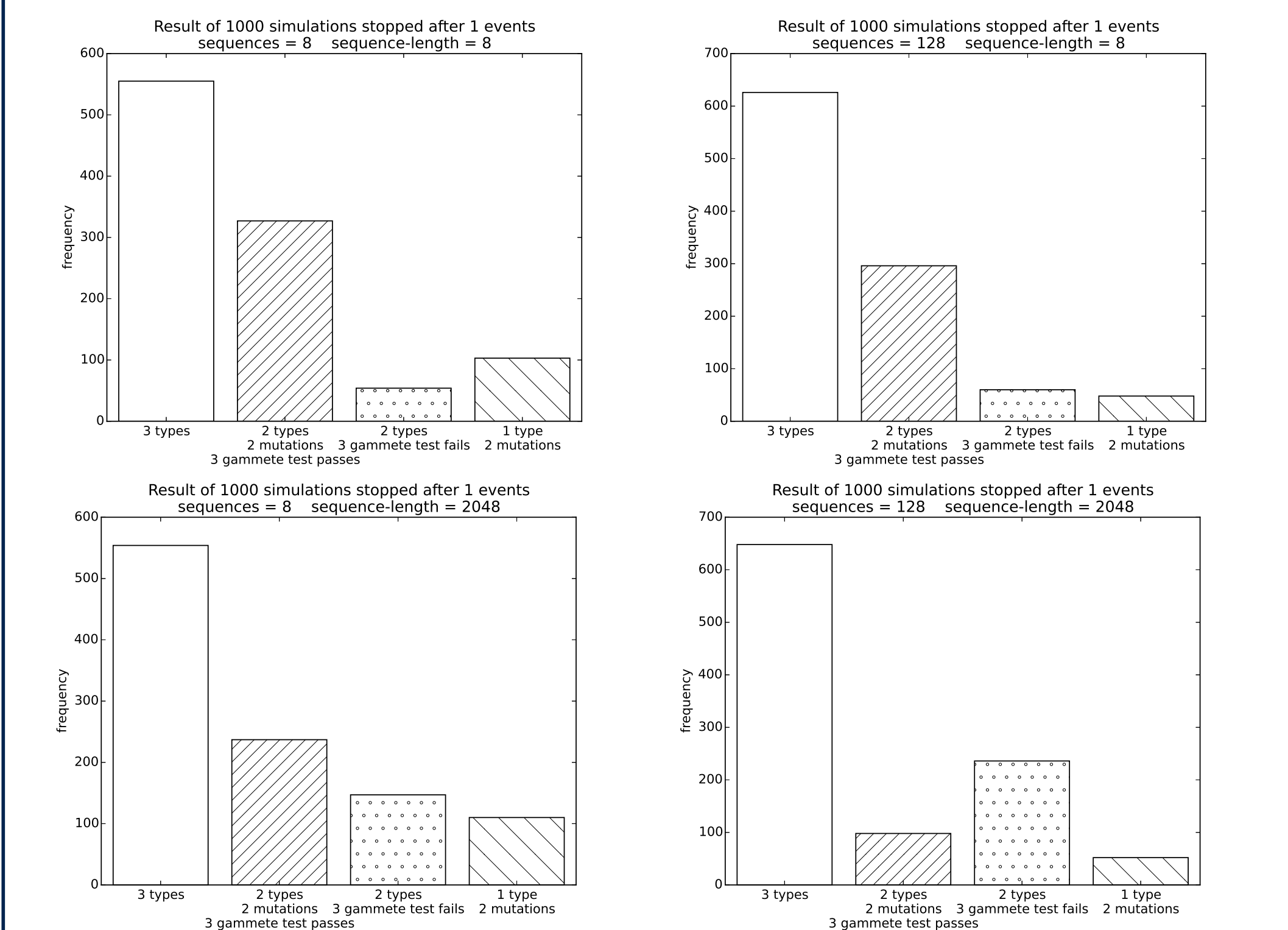
Some benchmarking-results for simulated data:



## SIMULATION – HOW IS THE INFINITE SITES ASSUMPTION VIOLATED?

We 1000 times simulate a Kingman-coalescent, and add mutations under the finite sites hypothesis, until one of the following events have occurred twice:

1. A site with  $> 2$  nucleotides occurs.
2. A site with 2 nucleotides has been affected by  $> 2$  mutations.
3. An *incompatibility* has occurred, eg.  $S = \begin{bmatrix} A & T & T & A & T \\ A & T & A & A & T \end{bmatrix}$
4. Two mutations "cancel out"



Code for simulation and plot-generation code available at: <https://github.com/Cronjaeger/coalescent-simulations>.

## FUTURE WORK AND INTERESTING OPEN QUESTIONS

- This work was mainly motivated by wanting to analyse viral sequence-data for sampled regions with sufficiently high mutation-rates that the ISA might be violated at a significant number of sites. This still remains to be done.
- Upper bounds the on run-time of the suggested algorithm remain unknown. Even for the case when  $b$  is equal to the number of segregating sites (recursion reduces to that of Griffiths '89), no literature could be found. The main factor of interest is the total number of states considered, given an initial input  $\psi$ .
- Being able to give bounds on  $\mathbb{P}(B \leq b | \psi)$  would allow a more informed way of picking  $b$  when applying our algorithm.
- Verifying if  $\exists S' \sim S : (S', b) \in \text{Keys}(\mathbf{H})$  holds is currently slow.
- A range of combinatorial and asymptotic questions regarding the set of all configurations remain open; solving them might bring us closer to understanding the nature and limitations of the algorithm.

1. Given  $\psi$ , how many distinct genealogical histories of  $\psi$  with up to  $b$  mutation-events exist?
2. Given  $\psi$ , how many distinct states  $\psi'$  exist, which occur in at least one such genealogical history?
3. What are the maxima and/or maximizers of the above counts taken over the set of all  $\psi$  with  $\text{rank}(\psi) < R$ ?
4. Can asymptotic bounds on the above quantities be obtained?

It is worth noting that answering questions 1. and 2. corresponds to counting paths and nodes respectively in a graph such as the one depicted on this poster.

## REFERENCES

Árnason, E., 2004. Mitochondrial Cytochrome b DNA Variation in the High-Fecundity Atlantic Cod: Trans-Atlantic Clines and Shallow Gene Genealogy. *Genetics*, 166(4), pp.1871–1885.  
 Dialdestoro, K. et al., 2016. Coalescent Inference Using Serially Sampled, High-Throughput Sequencing Data from Intra-host HIV Infection. *Genetics*, 202(4), pp.1449–72.  
 Ethier, S.N. & Griffiths, R.C., 1987. The Infinitely-Many-Sites Model as a Measure-Valued Diffusion. *The Annals of Probability*, 15(2), pp.515–545.  
 Griffiths, R.C., 1987. Counting genealogical trees. *Journal of Mathematical Biology*, 25(4), pp.423–431.  
 Griffiths, R.C., 1989. Genealogical-tree probabilities in the infinitely-many-site model. *Journal of Mathematical Biology*, (27), pp.667–680.  
 Hobolth, A. & Wiuf, C., 2009. The genealogy, site frequency spectrum and ages of two nested mutant alleles. *Theoretical population biology*, 75(4), pp.260–5.  
 Jenkins, P.A. & Song, Y.S., 2011. The effect of recurrent mutation on the frequency spectrum of a segregating site and the age of an allele. *Theoretical Population Biology*, 80(2), pp.158–173.

## ACKNOWLEDGEMENTS

The work outlined here grew out of a 10 Week joint mini-project by Mr. Cronjäger and Ms. Avalos-Pacheco, who are both students in the Oxford-Warwick Statistical Programme; a Centre for Doctoral Training in Next Generational Statistical Science supported by the Engineering and Physical Sciences Research Council and the Medical Research Council. Ms. Avalos Pacheco is furthermore supported by the Mexican National Council of Science and Technology (CONACYT). The project was supervised by Dr. Jotun Hein (Oxford) and Dr. Paul Jenkins (Warwick), who acting as supervisors both contributed substantially.