

**1 Reconstruction of ancestral protein sequences using autoregressive  
2 generative models**

**3** Pierre Barrat-Charlaix,\* Matteo De Leonardis, and Andrea Pagnani

**4** *DISAT, Politecnico di Torino, Italy*

**5** (Dated:)

6

## Abstract

7 abstract

8 **I. INTRODUCTION**

9 Homologous proteins have a common evolutionary origin that can span billions of years.  
10 Throughout their evolution, they diversify through mutations while selection preserves their  
11 function. Consequently, many protein families consist of thousands of sequences that are  
12 highly variable and yet maintain similar structures and functions. On the other hand, even a  
13 few mutations can destabilize a protein and destroy its function. A quantitative description  
14 of sequence changes in proteins is thus a challenging problem, with important consequences  
15 for our understanding of the evolution of life.

- 16 • Proteins are important. Evolved for billions of years. Mutation + selection – $\downarrow$  varied  
17 sequences with relatively conserved function/structure. Challenge to describe this  
18 evolution
- 19 • How typical evolution models work, since JC69 and on. Applications to phylogenetic  
20 and ASR – $\downarrow$  very important bioinformatic tools. However, do not take epistasis into  
21 account.
- 22 • On the other hand, generative models. List all cool things
- 23 • Challenge to extend generative models to dynamics and to apply them to dynamical  
24 problems such as ASR. This is especially important for ASR, since conclusions depend  
25 on the quality of the reconstructed sequence
- 26 • Here, we proceed in two parts. First we develop an analytically tractable dynamical  
27 generative model based on ArDCA
- 28 • Second, we test its capacity to perform ASR on simulated and directed evolution data

---

\* Correspondance to: PBC, DISAT pierre.barratcharlaix@polito.it

29 **II. RESULTS**

30 **A. Autoregressive model of sequence evolution**

31 Models of evolution commonly used in phylogenetics rely on the assumptions that sequence  
 32 positions evolve independently and that evolution at each position  $i$  follows a continuous  
 33 time Markov chain (CTMC) parametrized by a substitution rate matrix  $\mathbf{Q}^i$ . Matrix  $\mathbf{Q}^i$  is  
 34 of dimensions  $q \times q$  where  $q = 4$  for DNA, 20 for amino acids or 64 for codon models. The  
 35 probability of observing a change from state  $a$  to state  $b$  during evolutionary time  $t$  is then  
 36 given by  $P(b|a, t) = \left( e^{t\mathbf{Q}^i} \right)_{ab}$ .

37 If the model is time-reversible, it is a general property of CTMCs that the substitution  
 38 rate matrix can be written as

$$\mathbf{Q} = \mathbf{H} \cdot \boldsymbol{\Pi} = \mathbf{H} \cdot \begin{pmatrix} \pi_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \pi_q \end{pmatrix}, \quad (1)$$

39 where  $\mathbf{H}$  is symmetric and  $\boldsymbol{\Pi}$  is diagonal with entries that sum to 1 [1]. The two matrices  
 40 have simple interpretations. On the first hand,  $\boldsymbol{\Pi}$  fixes the long-term equilibrium frequencies,  
 41 that is  $P(b|a, t) \xrightarrow[t \rightarrow \infty]{} \pi_b$ . On the other,  $\mathbf{H}$  influences the dynamics of the Markov chain  
 42 but does not change the equilibrium distribution. Most commonly,  $\boldsymbol{\Pi}$  is considered to be  
 43 independent of the sequence position  $i$ , while  $\mathbf{H}$  can be multiplied by position-dependent  
 44 rates in order to model the different variability of different sites [2–4].

45

46 In order to incorporate constraints coming from a protein’s structure and function into  
 47 the evolutionary model, we develop a family-specific model of protein sequence evolution  
 48 based on the the autoregressive generative model ArDCA [5]. ArDCA models the diversity  
 49 of sequences in a protein family using a set of learned conditional probabilities. In practice,  
 50 the model assigns a probability to any sequence  $\mathbf{a} = \{a_1, \dots, a_L\}$  of  $L$  amino acids:

$$P^{AR}(\mathbf{a}) = \prod_{i=1}^L p_i(a_i | a_{<i}), \quad (2)$$

51 where the product runs over positions in the sequence and  $a_{<i} = a_1, \dots, a_{i-1}$  represents  
 52 all amino acids before position  $i$ . Functions  $p_i$  represent the probability according to the

53 model to observe state  $a_i$  in position  $i$ , given that the previous amino acids were  $a_1, \dots, a_{i-1}$ .  
 54 Their precise functional form is given in the methods section. They are learned using the  
 55 aligned sequences of members of the family. In actual implementations, the order in which  
 56 the product in Eq. 2 is performed is not the natural  $(1, \dots, L)$  but rather an order where  
 57 positions are sorted by increasing variability. This does not significantly effect the model we  
 58 present below, and we keep the notation of Eq. 2 for simplicity.

59 It has been shown in [5] that the generative capacities of ArDCA are comparable to that  
 60 of state of the art models such as bmDCA [6]. This means that a set of sequences sampled  
 61 from the probability in Eq. 2 is statistically hard to distinguish from the natural sequences  
 62 used in training or, in other words, that the model can be used to sample new artificial  
 63 homologs of a protein family. Generative capacities of a protein model come from its ability  
 64 to represent epistasis, that is the relation between the effect of a mutation and sequence  
 65 context in which it occurs. Here, epistasis is modeled through the conditional probabilities  $p_i$ :  
 66 the distribution of amino acids at position  $i$  depends on the states of the previous positions  
 67  $1, \dots, i - 1$ .

68 We take advantage of the autoregressive architecture to define the following evolution  
 69 model. Given two amino acid sequences **a** and **b**, we propose

$$P(\mathbf{b}|\mathbf{a}, t) = \prod_{i=1}^L q_i(b_i|a_i, b_{<i}, t), \quad (3)$$

70 where the conditional propagator  $q_i$  is defined as

$$q_i(b_i|a_i, b_{<i}, t) = \left( e^{t \cdot Q^i(b_{<i})} \right)_{a_i, b_i}, \quad Q^i(b_{<i}) = \mathbf{H} \cdot \begin{pmatrix} p_i(1|b_{<i}) & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & p_i(q|b_{<i}) \end{pmatrix}. \quad (4)$$

71 According to these equations, evolution for each position  $i$  follows a standard CTMC.  
 72 However, we use the decomposition of Eq. 1 to set the equilibrium frequency at  $i$  to  $p_i(b|b_{<i})$ .  
 73 In other words, we consider that position  $i$  evolves in the context of  $b_1, \dots, b_{i-1}$ , and that its  
 74 dynamics are constrained by its long term frequency given by the autoregressive model. An  
 75 important consequence of this choice is that our evolution model will converge at long times  
 76 to the generative distribution  $P^{AR}$ .

$$q_i(b_i|a_i, b_{<i}, t) \xrightarrow{t \rightarrow \infty} p_i(b_i|b_{<i}), \quad P(\mathbf{b}|\mathbf{a}, t) \xrightarrow{t \rightarrow \infty} P^{AR}(\mathbf{b}). \quad (5)$$

77 We argue here that such a property is essential to build a realistic protein sequence  
 78 evolution model, particularly when considering evolution over a relatively long time frame.  
 79 Note that to converge to a generative distribution, accurate modeling of epistasis is required.  
 80 Using site-specific frequencies would not be sufficient, as the effect of mutations in a protein  
 81 sequence typically depends on the context [7]. The technique proposed here allows us to  
 82 represent epistasis through the context dependent probabilities  $p_i$ , while still considering  
 83 each sequence position one at a time.

84 Interestingly, we note that the model in Eq. 3 is not time reversible, although context  
 85 dependent site propagators in Eq. 4 are reversible. We show in the SM that this is mainly  
 86 an artifact of the autoregressive nature of the model coupled with epistasis. Using non-time  
 87 reversible evolutionary models is uncommon in the field, but this is mainly due to practical  
 88 considerations and there are no fundamental reasons for evolution itself to be reversible [8].  
 89 In practice, this means that algorithms using this model have to be adapted accordingly.

90

91 We underline that this approach has important differences with standard models of  
 92 evolution used in phylogenetics. In phylogenetic reconstruction, the tree and the sequence  
 93 evolution model are usually inferred at the same time and from the same data. The number  
 94 of parameters of the evolution model is then kept low to reduce the risk of overfitting, for  
 95 instance by using site specific rates to account for variable and conserved sites. Methods that  
 96 introduce more complex models such as site specific frequencies do so by jointly inferring the  
 97 parameters and the tree, leading to relatively complex algorithms [9, 10].

98 Here instead, parameters of the generative model in Eq. 2 are learned from a protein family,  
 99 *i.e.* a set of diverged homologous protein sequences. While it is true that these sequences  
 100 share a common evolutionary history and cannot be considered as independent samples,  
 101 common learning procedures only account for this in a very crude way [5, 11]. Despite this,  
 102 it appears that the generative properties of such models are not strongly affected by ignoring  
 103 the phylogeny [12, 13]. This allows us to proceed in two steps: first construct the model  
 104 from data while ignoring phylogeny, and then only use it for phylogenetic inference tasks.

105 An advantage of this approach is that once the model of Eq. 2 is inferred, the propagator  
 106 in Eq. 3 comes “for free” as no additional parameters are required. Importantly, our model

107 does not use site specific substitution rates. Indeed, it has been shown that these can be seen  
108 as emergent properties when using more complex evolution models such as the one presented  
109 here [14]. However, a disadvantage is that the technique is only applicable to a given protein  
110 family at a time, and requires the existence of an appropriate training set for the model.

111 **B. Ancestral sequence reconstruction**

112 We apply our evolutionary model to the task of ancestral sequence reconstruction (ASR).  
113 The goal of ASR is the following: given a set of extant sequences with a shared evolutionary  
114 history and the corresponding phylogenetic tree, is it possible to reconstruct the sequences of  
115 extinct ancestors at the internal nodes of the tree? Along with the autoregressive evolutionary  
116 model described above, we thus need two inputs to perform ASR: a known phylogenetic  
117 tree, and the multiple sequence alignment of the leaf sequences. The length of the aligned  
118 sequences has to exactly correspond to that of the autoregressive model.

119 To reconstruct ancestral sequences using the autoregressive model, we proceed as follows:

120 *i* for sequence position  $i = 1$ , use the evolution model defined by the equilibrium  
121 frequencies  $p_1$  to reconstruct a state  $a_1^n$  at each internal node  $n$  of the tree;

122 *ii* iterating through subsequent positions  $i > 1$ : reconstruct state  $a_i^n$  at each internal  
123 node  $n$  using the model defined in Eq. 4, with the context  $a_{<i}^n$  having been already  
124 reconstructed in the previous iterations.

125 It is important to note that when any position  $i > 1$  is reconstructed, the context at different  
126 internal nodes of the tree may differ. For a branch joining two nodes  $(n, m)$  of the tree, the  
127 evolution model will thus differ if we go down or up the branch: in one case the context at  
128 node  $n$  must be used, in the other case the context at node  $m$ . This is a consequence of the  
129 time-irreversibility of the model. For this reason, we use a variant of Felsenstein’s pruning  
130 algorithm that is adapted to irreversible models [15]. This comes at no computational cost.

131 Note that this method is adapted to both maximum likelihood and Bayesian inference. In  
132 the ML case, each iteration reconstructs the most probable residue is at a position  $i$  given  
133 the already reconstructed context. In the Bayesian case, residue  $a_i^n$  is instead sampled from  
134 the posterior distribution.

135 In any realistic application, the phylogenetic tree has to be reconstructed from the aligned  
136 sequences. In principle, a consistent approach would use the same evolutionary model for  
137 tree inference and ASR. However, our model does not allow us to reconstruct the tree.  
138 Therefore, in any realistic application, the tree is reconstructed using an evolutionary model  
139 that typically will differ from ours.

140 To reduce issues related to this evolutionary model discrepancy, we adopt the following  
141 strategy: our ASR method blindly trusts the topology of the input tree, but recomputes the  
142 branch length using the sequences. As explained in the Methods, there is no direct way to  
143 optimize branch length with the autoregressive model. For simplicity, we use profile model  
144 with position-specific amino acid frequencies for this task. This provides a relatively accurate  
145 estimate of the branch lengths, as shown in Figure S1.

### 146 C. Results on simulated data

147 There are two difficulties when evaluating the capacity of a model to perform ASR. The  
148 first is that in the case of biological data, the real phylogeny and ancestral sequences are  
149 usually not known. As a consequence, one must rely on simulated data to measure the  
150 quality of reconstruction. The second is that the reconstruction of an ancestral sequence is  
151 always uncertain, as evolutionary models are typically stochastic. The uncertainty becomes  
152 higher for nodes that are remote from the leaves. This means that it is only possible to make  
153 a statistical assessment about the quality of a reconstruction.

154 To test our approach, we adopt the following setup. We first generate phylogenetic trees  
155 by sampling from a coalescent process. We decide to use Yule’s coalescent instead of the  
156 more common Kingman. The latter tends to produce a large majority of internal nodes  
157 in close vicinity to the leaves with the others separated by very long branches, resulting  
158 in a trivial reconstruction for most nodes and a very hard one for the deep nodes. Yule’s  
159 coalescent generates a more even distribution of node depths, allowing us to better evaluate  
160 reconstruction quality, see Supplementary Material and Figure S2. For each tree, we simulate  
161 the evolution of sequences using a model that we refer to as “evolver” to obtain two multiple  
162 sequence alignments, one for the leaves and one for the internal nodes of the tree. We then  
163 reconstruct internal nodes using the desired approach by using the leaf alignment and the  
164 tree topology as input data.

165 We will consider two kinds of evolver models: (*i*) the same autoregressive model that we  
166 will then use for reconstruction, which is an ideal case and (*ii*) an evolutionary model based  
167 on a Metropolis sampling of a Potts model. These two evolvers come from models trained  
168 on actual protein families: we use evolvers based on the PF00072 response regulator family  
169 for results of the main text, and show results for three other families (PF00014, PF00076  
170 and PF00595) in the Supplementary Material . It is important to note that the approach  
171 that we propose only makes sense when considering the evolution of a precise protein family,  
172 on which the model in Eq. 2 is trained. Hence, any evolver model used in our simulations  
173 should reproduce at long times the statistics of the considered protein family, *i.e.* it should  
174 satisfy Eq. 5. For this reason, we only consider the two evolvers above and do not use more  
175 traditional evolutionary models such as an arbitrary GTR on amino-acids [16].

176 For reconstruction, we compare our autoregressive approach to the commonly used  
177 IQ-TREE program [17]. When supplied with a protein sequence alignment and a tree,  
178 IQ-TREE infers a joint substitution rate matrix for all sequence positions, with rates that  
179 can differ accross positions. Both methods run on a fixed tree topology, with branch lengths  
180 being re-inferred using maximum likelihood (see Methods).

181

182 *Autoregressive evolver.* We first investigate the case of the autoregressive evolver. This  
183 setting is of course ideal for our method, as there is perfect coincidence between the model  
184 used to generate the data and to perform ASR. We first evaluate the quality of reconstruction  
185 by computing the Hamming distance of the real and inferred sequences for each internal node  
186 of the simulated phylogenies. The left and central panels of Figure 1 show this Hamming  
187 distance as a function of the node depth, that is the distance separating the node from  
188 the leaves, and for a maximum likelihood reconstruction. Hamming distance is computed  
189 including gap characters in the aligned sequences on the right panel, while they are ignored  
190 on the central one. We see that the autoregressive reconstruction clearly outperforms the  
191 state of the art method: the improvement in Hamming distance increases with node depths,  
192 and the distance to the real ancestor drops from  $\sim 0.4$  to  $\sim 0.3$  when using the autoregressive  
193 approach. The increase in reconstruction quality with node depths is consistent with recent  
194 findings that epistasis only becomes important at relatively large sequence divergences  
195 [18, 19].

196 Interestingly, the performance of IQ-TREE degrades if Hamming distance is computed

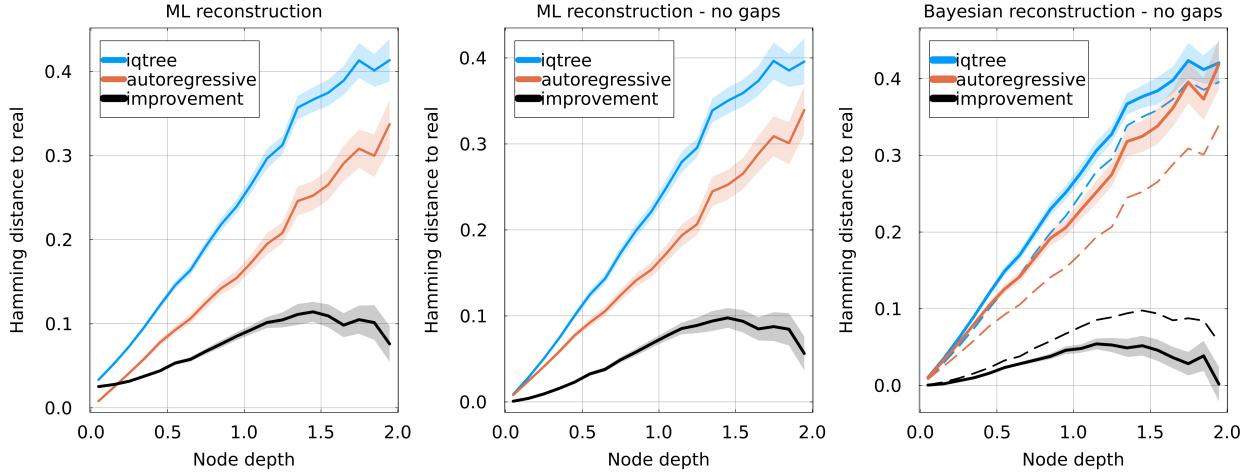
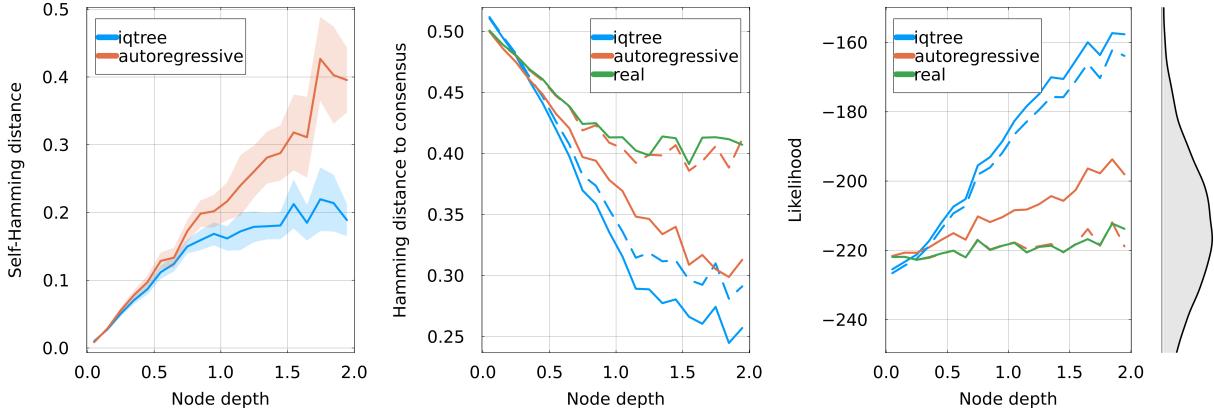


FIG. 1. Hamming distance between reconstructed and real sequences as a function of node depth, using IQ-TREE and our autoregressive approach. The difference between the two methods (“improvement”) is shown as a black curve. Estimation of the incertitude is shown as a ribbon. The evolver and reconstruction autoregressive models are learned on the PF00072 family. **Left:** Hamming distance between the full aligned sequences, gaps included, using maximum likelihood reconstruction. **Center:** Hamming distance ignoring gapped positions, using maximum likelihood reconstruction. **Right:** comparison of Bayesian (solid lines) and maximum likelihood (dashed lines) reconstructions, ignoring gaps.

including gaps, as in the left panel. This is because like other popular methods, IQ-TREE treats gaps in input sequences as unknown amino acids, and reconstructs an ancestral amino acid for gapped positions [17, 20]. On the contrary, our autoregressive approach treats gaps as if they were an additional amino acid and will reconstruct ancestral sequences that can contain gaps. This benefits the autoregressive approach as aligned ancestral sequences can in fact contain gaps. This effect is particularly visible at low node depths. However, ignoring the effects of gaps in the Hamming distance also leads to a clear improvement when using the autoregressive approach as shown in the central panel.

The right panel of Figure 1 shows the quality of the reconstruction for Bayesian reconstruction. In this case, a ensemble of sequences is reconstructed for each internal node, and the metric is the average Hamming distance between this ensemble and the real ancestor. Gaps are again ignored when computing the Hamming distance. We again observe an improvement when using the autoregressive method, of slightly lesser magnitude than in the



**FIG. 2. Left:** for Bayesian reconstruction, average pairwise Hamming distance among sequences reconstructed for each internal node. This quantifies the diversity of sequences obtained using Bayesian reconstruction. **Center:** Hamming distance between reconstructed sequences and the consensus sequence of the alignment. Solid lines represent maximum likelihood reconstruction or the real internal sequences, and dashed lines Bayesian reconstruction. IQ-TREE appears more biased towards the consensus sequence. **Right:** Log-likelihood of reconstructed and real sequences in the autoregressive model, *i.e.* using the logarithm of Eq. 2. Maximum likelihood methods (orange and blue solid lines) are biased towards more probable sequences. Bayesian autoregressive reconstruction gives sequences that are at the same likelihood level than the real ancestors. The equilibrium distribution of likelihood of sequences generated by Eq. 2 is shown on the right.

210 maximum likelihood case.

211

212 *Reconstructed sequences.* To further analyze the reconstructed sequences, we first look at  
 213 the diversity of generated ancestors in Bayesian reconstruction. The left panel of Figure 2  
 214 shows the average Hamming distance between sequences reconstructed at the same internal  
 215 node, as a function of depth. For deeper nodes (depth  $\gtrsim 1$ ), the Bayesian autoregressive  
 216 approach reconstructs a significantly more diverse set of sequences than IQ-TREE: Hamming  
 217 distance between reconstructions saturates at 0.2 for the latter, while it steadily increases  
 218 for the former. Higher diversity can be interpreted as a greater uncertainty concerning the  
 219 ancestral sequence. However, this must be put in the context of Figure 1: sequences obtained  
 220 by Bayesian autoregressive reconstruction are more varied but also on average closer to the  
 221 real ancestor.

222 The difference in sequence diversity for the two methods is in part explained by the central  
223 panel of Figure 2, which shows the Hamming distance between reconstructed ancestors and  
224 the consensus sequence of the multiple sequence alignment at the leaves. It appears there that  
225 for deep nodes, IQ-TREE reconstructs sequences that are relatively similar to the consensus,  
226 with an average distance between the Bayesian reconstruction and the consensus of about  
227 0.3. Contrasting with that, results of the autoregressive method shows less bias towards the  
228 consensus with an average distance of 0.4 for deep nodes, in line with the real ancestors.  
229 We also note that maximum likelihood sequences for both method are always closer to the  
230 consensus than Bayesian ones, a bias of maximum likelihood already observed in the past  
231 [21].

232 The bias induced by ignoring the equilibrium distribution of the sequences is also visible  
233 in the right panel of Figure 2: it shows the log-likelihood of reconstructed and real ancestral  
234 sequences according to the generative model. Note that the log-likelihood here comes from  
235 the log-probability of Eq. 2 and can be interpreted as the “quality” of a sequence according  
236 to the generative model. It is unrelated to the likelihood computed in Felsenstein’s pruning  
237 algorithm. Reconstructions with IQ-TREE increase in likelihood when going deeper in  
238 the tree, eventually resulting in sequences that are very uncharacteristic of the equilibrium  
239 generative distribution as can be seen from the histogram on the right. This effect also  
240 happens with the maximum likelihood reconstruction of the autoregressive model, although  
241 to a lesser extent. The Bayesian autoregressive reconstruction does not suffer from this bias  
242 and reconstructs sequences with a log-likelihood that is similar to that of the real ancestors.

243

244 *Potts evolver.* We then assess the performance of our reconstruction method in the case  
245 where the evolver is a Potts model. Potts models are a simple type of generative model and  
246 have been used extensively to model protein sequences. They can be used to predict contact  
247 in three dimensional structures, effects of mutations, protein-protein interaction partners [11].  
248 They can be sampled to generate novel sequences which are statistically similar to natural  
249 ones and often functional [6, 22]. Additionally, it has recently been shown that they can be  
250 used to describe the evolution of protein sequences both qualitatively and quantitatively [23].

251 Potts and autoregressive models both accurately reproduce the statistical properties of  
252 protein families. In this sense, it can be considered that they correspond to similar long  
253 term generative distributions in the sense of Eq. 5. However, the dynamics of a Potts

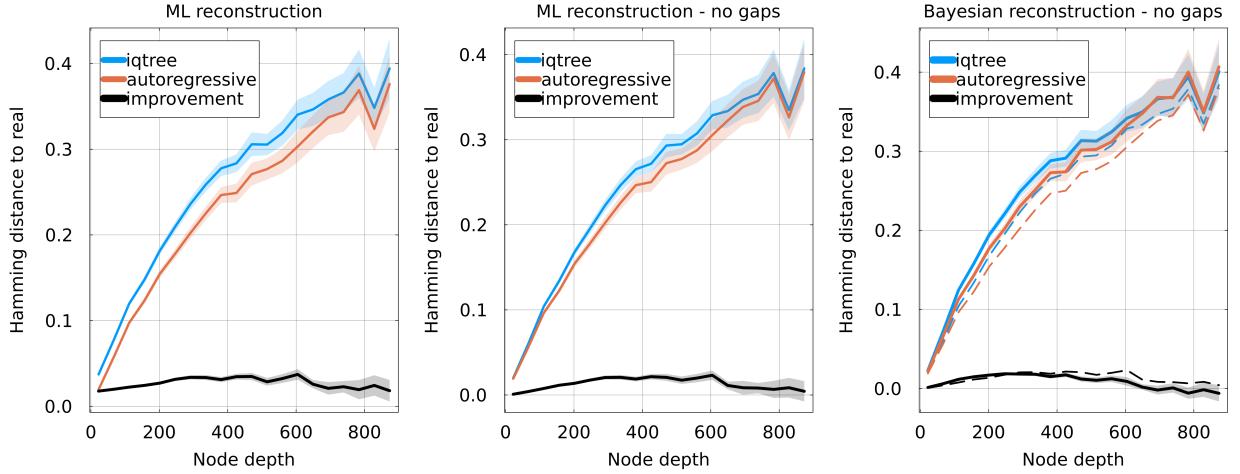


FIG. 3. Analogous to Figure 1, but using a Potts model as the evolver. Hamming distance between reconstructed and real sequences as a function of node depth, using IQ-TREE and our autoregressive approach. The difference between the two methods is shown as a black curve. The evolver and reconstruction autoregressive models are learned on the PF00072 family. **Left:** Hamming distance between the full aligned sequences, gaps included, using maximum likelihood reconstruction. **Center:** Hamming distance ignoring gapped positions, using maximum likelihood reconstruction. **Right:** comparison of Bayesian (solid lines) and maximum likelihood (dashed lines) reconstructions, ignoring gaps.

model are fundamentally different from the ones of usual evolutionary models, including our autoregressive one. First, they are described by a *discrete* time Markov chain, instead of the continuous time used in models based on substitution rate matrices such as in Eq. 1 [24]. The discrete time in these dynamics corresponds to attempts at mutation, which can be either accepted or rejected depending on the effect of the mutation according to the model. Secondly, these dynamics naturally give rise to different evolutionary timescales for various sequence positions, as well as interesting qualitative behavior such as the entrenchment of mutations.

To see how this change in dynamics affects our results, we (*i*) sample a large and varied ensemble of sequences from the Potts model and use it to train an autoregressive model, in a way to guarantee consistent long term distributions between the Potts and autoregressive, and (*ii*) evolve the Potts model along random phylogenies, generating alignments for the leaves and the internal nodes in the same way as above. We then attempt reconstruction

267 of internal nodes using the inferred autoregressive model and IQ-TREE. Figure 3 shows  
268 the results of reconstruction, with panels directly comparable to Figure 1. We again see a  
269 consistent improvement when using the autoregressive model over IQ-TREE, although of a  
270 much smaller amplitude, with an absolute improvement gain in Hamming distance of about  
271 2% for deep internal nodes.

272 **D. Results on experimental evolution data.**

273 We take advantage of recent developments in directed evolution experiments to test our  
274 method in a controlled setting. We use the data published in [25]: in this work, authors  
275 evolved the antibiotic resistant proteins  $\beta$ -lactamase PSE-1 and acetyltransferase AAC6 by  
276 submitting them to cycles of mutagenesis and functional selection. Starting from a wild-type  
277 protein, they obtained thousands of diverse functional sequences after the directed evolution.  
278 An interesting result of this work is that it is possible to recover structural information about  
279 the wild-type from the set of evolved sequences.

280 Here, we use this data as a test setting for ASR: the sequences obtained after directed  
281 evolution all derive from a common ancestor, the wild-type, for which we know the amino  
282 acid sequence. We can thus reconstruct the wild-type sequence using different ASR methods  
283 and compare it to the ground truth. The phylogeny is not known, but given the large  
284 population size during the experiment and the relatively low number of selection rounds, it  
285 is reasonable to approximate it using a star-tree, *i.e.* a tree with a single coalescent event  
286 taking place at the root (see Methods). Since the reconstruction task is most interesting  
287 when using relatively varied sequences, we decide to use data for the PSE-1 wild-type where  
288 20 cycles of mutagenesis & selection have been performed, resulting in a mean Hamming  
289 distance of 12% to the wild-type.

290 Our ASR procedure is as follows. We randomly pick  $M$  amino acid sequences among the  
291 proteins evolved from PSE-1 and with 20 cycles of mutagenesis & selection, with  $3 \leq M \leq 640$ .  
292 The total number of sequences at round 20 of directed evolution is much larger, making  
293 it computationally hard to use all of them. We then construct a star-like phylogeny and  
294 place the  $M$  selected sequences at the leaves, and perform ASR using either IQ-TREE or  
295 our autoregressive method which we have trained on an alignment of PSE-1 homologs. We  
296 obtain the reconstructed amino acid sequence of the root, which we can then compare to the

297 actual wild-type. As a comparison, and because our approximation of the phylogeny is very  
298 simple, we also attempt to reconstruct the root by taking the consensus sequence of the  $M$   
299 leaves. We repeat this procedure 100 times for each value of  $M$  for a statistical assessment  
300 of the different methods.

301 The results are shown in Figure 4. The left panel shows the average non-normalized  
302 Hamming distance to the wild-type as a function of the number of leaves used  $M$ . For a  
303 low  $M$ , all methods understandably make a large number of errors, with a mean Hamming  
304 distance larger than 10 for  $M = 3$ . For a higher  $M$ , IQ-TREE and the autoregressive method  
305 stabilize to a fixed number of errors: we find a Hamming distance of  $\sim 4.3$  for IQ-TREE and  
306  $\sim 2.9$  for the autoregressive. The consensus curiously reaches a minimum at intermediate  $M$ ,  
307 a fact commented in the Supplementary Material , and saturates at a Hamming distance of 6  
308 when considering all sequences of the round 20. The reconstruction errors are overwhelmingly  
309 located at six sequence positions. In the central panel, the fraction of mistakes made at  
310 these six positions over the 100 repetitions of  $M = 640$  leaves is shown for each method.  
311 We observe that there are two positions (129 and 152) where IQ-TREE systematically fails  
312 at recovering the wild-type state while the autoregressive model’s reconstruction is correct.  
313 Inversely, IQ-TREE recovers the wild-type state more often at position 68. The right panel  
314 shows the logo of the set of reconstructed sequences at these 6 positions and for each method.

315 Overall, we see that the reconstruction of the autoregressive model is more accurate. This  
316 gain in accuracy comes from the representation of the functional constraints acting on the  
317 PSE-1 protein by the generative model, which are inferred separately using an alignment  
318 of homologs. The improvement in reconstruction errors is modest, going from an average  
319 Hamming distance of 4.3 to 2.9. However, the gain is intrinsically limited by the data itself:  
320 the evolved sequences have an average Hamming distance of about 12% to the ancestor,  
321 which is experimentally challenging but remains small compared to the divergence found  
322 in the homologs of PSE-1. For instance, the root-to-tip distance estimated by IQ-TREE  
323 and the autoregressive model are respectively 0.13 and 0.15, corresponding to the regime of  
324 shallow trees when comparing with Figure 1.

325 In order

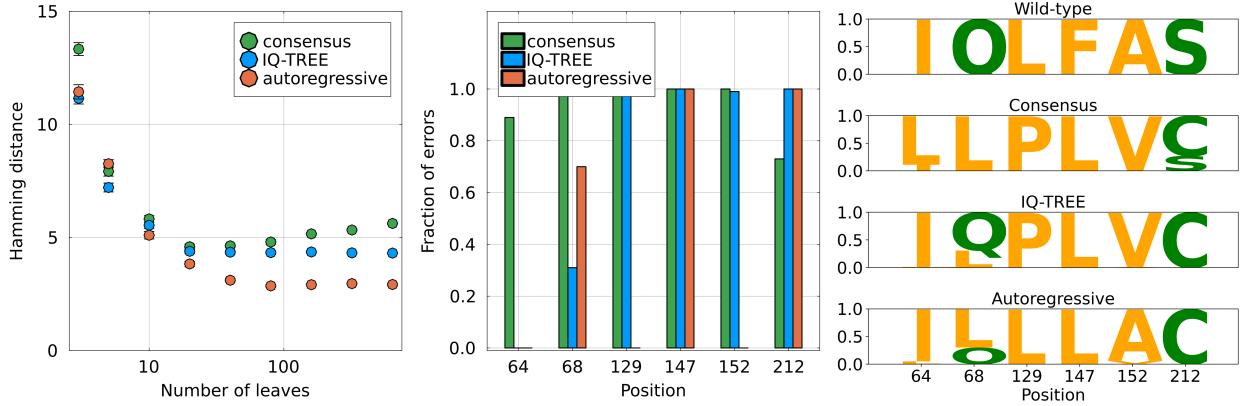


FIG. 4. Reconstruction of the wild-type PSE1 sequence used in [25] using sequences from round 20 of the directed evolution. **Left.** Non normalized Hamming distance to the wild-type PSE1 sequence as a function of the number of sequences  $M$  used for reconstruction. The fact that the consensus method has a local minimum is discussed in the Supplementary Material . For comparison, the average distance between a leaf sequence and the wild-type is 25. The error bars are computed using the standard deviation obtained from the 100 choices of sequences. **Middle.** For the six sequence positions where most of the reconstruction errors are located, fraction of errors of each method out of 100 independent reconstructions using different sets of  $M = 640$  leaves. **Right.** Sequence logo of the reconstructed sequence for the three methods, obtained using 100 independent reconstructions with different sets of  $M = 640$  leaves. The logo is only shown for the six positions where most errors are located. For example, all three methods fail 100 times at position 147, reconstructing a leucine  $L$  instead of a phenylalanine  $F$ .

326 **III. DISCUSSION**

327 **IV. METHODS**

328 **A. ArDCA**

329 The ArDCA model assigns a probability to any sequence of amino acids of length  $L$  given  
330 by

$$P^{AR}(\mathbf{a}) = \prod_{i \in \sigma(L)} p_i(a_i | a_{<i}), \quad (6)$$

331 where  $\sigma(L)$  is a permutation of the  $L$  first integers and  $a_{<i}$  stands for  $a_1, \dots, a_{i-1}$ . This

means that the order in which the conditional probabilities  $p_i$  are applied is not necessarily the sequence order. The permutation  $\sigma$  is fixed at model inference.

Conditional probabilities  $p_i$  are defined as

$$p_i(b|a_{<i}) = \frac{1}{Z_i} \exp \left( \sum_{j < i} J_{ij}(b, a_j) + h_i(b) \right), \quad (7)$$

with the  $i$   $q$ -dimensional vectors  $J_i$ . and  $h_i$  are learned parameters. The model is normally trained using a multiple sequence alignment of homologous proteins, *i.e.* a protein family, by finding the parameters  $J$  and  $h$  that maximize the likelihood of the sequences. It was shown in [5] that this specific parametrization captures essential features of the variability of members of a protein family.

By definition, homologous proteins share a joint evolutionary history and cannot be considered as statistically independent. To avoid biases, a reweighting is applied to sequences based on their vicinity to other sequences. This scheme has been showed to substantially increase the performance of such models [11].

## 344 B. Branch length inference

345 • how it works generally

346 • why I can't use AR directly

347 • how we do it with the profile model

348 • refer to the SI figure for results

## 349 C. Simulations

350 A simulation is performed as follows. First, a random tree of  $n = 100$  leaves is generated  
351 from Yule's coalescent. We then normalize its height to a fixed value  $H$  that depends on the  
352 evolver model used: for the autoregressive model we use  $H = 2.0$ , while for the Potts model  
353 combined with Metropolis steps, we use  $H = 8$  sweeps, *i.e.*  $H = 8L$  Metropolis steps where  
354  $L$  is the length of the sequences.

355 A root sequence is sampled from the evolver model's equilibrium distribution, and evolution  
356 is simulated along each branch independently starting from the root. In the case of the

357 autoregressive evolve, the dynamics is the one of Eq. 3. In the case of the Potts model,  
358 we use a Markov chain with the Metropolis update rule. In this way, we obtain for each  
359 repetition a tree and the alignments for internal and leaf nodes. Results presented in this  
360 work are obtained by averaging over  $M = 100$  such simulations for each protein family.

361 **D. Experimental evolution data**

362 We use the data coming from the experimental evolution of the PSE-1 protein, published  
363 in [25].

- 364 • A word on the wild-type and what it does  
365 • A brief summary of the experiment. What diversity is obtained after 20 rounds? What  
366 is the population size?  
367 • what the corresponding pfam is, where we got the sequences (ref.)  
368 • the alignment methodology

369 It has been noticed in [23] that taking into account the transition possibilities between  
370 amino acids allowed by the genetic code is important when describing short term evolutionary  
371 dynamics with generative models. In our framework, a natural way to include these is by using  
372 the symmetric matrix  $\mathbf{H}$  in the decomposition of Eq. 1. Terms of the  $\mathbf{H}$  matrix do not affect  
373 the equilibrium distribution of the model, which thus remains generative, but influences the  
374 short term dynamics. Here, we simply counted the number of possibilities to transition from  
375 any amino acid to any other based on the genetic code, and we constructed the corresponding  
376  $\mathbf{H}$  matrix. The diagonal matrix remains given by the equilibrium probabilities of amino  
377 acids in the context of the sequence, as given by Eq. 4. We found that this substantially  
378 improves the results of the autoregressive reconstruction for the experimental evolution data.

379 **E. iqtree**

380 We run IQ-TREE using the **-asr**

- 381 • how I run iqtree + results of model finder

- 
- 382 [1] Ziheng Yang. *Computational Molecular Evolution*. Oxford Series in Ecology and Evolution.  
383 Oxford University Press, Oxford, New York, October 2006. ISBN 978-0-19-856702-8.
- 384 [2] Ziheng Yang. Maximum likelihood phylogenetic estimation from DNA sequences with variable  
385 rates over sites: Approximate methods. *Journal of Molecular Evolution*, 39(3):306–314,  
386 September 1994. ISSN 1432-1432. doi:10.1007/BF00160154.
- 387 [3] Alexandros Stamatakis. RAxML version 8: A tool for phylogenetic analysis and post-analysis  
388 of large phylogenies. *Bioinformatics (Oxford, England)*, 30(9):1312–1313, May 2014. ISSN  
389 1367-4811. doi:10.1093/bioinformatics/btu033.
- 390 [4] Lam-Tung Nguyen, Heiko A. Schmidt, Arndt von Haeseler, and Bui Quang Minh. IQ-TREE:  
391 A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phyloge-  
392 nies. *Molecular Biology and Evolution*, 32(1):268–274, January 2015. ISSN 0737-4038. doi:  
393 10.1093/molbev/msu300.
- 394 [5] Jeanne Trinquier, Guido Uguzzoni, Andrea Pagnani, Francesco Zamponi, and Martin Weigt.  
395 Efficient generative modeling of protein sequences using simple autoregressive models. *Nature  
396 Communications*, 12(1):5800, October 2021. ISSN 2041-1723. doi:10.1038/s41467-021-25756-4.
- 397 [6] Francisco McGee, Sandro Hauri, Quentin Novinger, Slobodan Vucetic, Ronald M. Levy,  
398 Vincenzo Carnevale, and Allan Haldane. The generative capacity of probabilistic protein  
399 sequence models. *Nature Communications*, 12(1):6302, November 2021. ISSN 2041-1723.  
400 doi:10.1038/s41467-021-26529-9.
- 401 [7] Michael Socolich, Steve W. Lockless, William P. Russ, Heather Lee, Kevin H. Gardner, and  
402 Rama Ranganathan. Evolutionary information for specifying a protein fold. *Nature*, 437(7058):  
403 512–518, September 2005. ISSN 1476-4687. doi:10.1038/nature03991.
- 404 [8] Felsenstein, Joseph. *Inferring Phylogenies*. Oxford university press edition, September 2003.  
405 ISBN 978-0-87893-177-4.
- 406 [9] A L Halpern and W J Bruno. Evolutionary distances for protein-coding sequences: Modeling  
407 site-specific residue frequencies. *Molecular Biology and Evolution*, 15(7):910–917, July 1998.  
408 ISSN 0737-4038. doi:10.1093/oxfordjournals.molbev.a025995.
- 409 [10] Vadim Puller, Pavel Sagulenko, and Richard A Neher. Efficient inference, potential, and  
410 limitations of site-specific substitution models. *Virus Evolution*, 6(2), August 2020. ISSN

- 411 2057-1577. doi:10.1093/veaa066.
- 412 [11] Simona Cocco, Christoph Feinauer, Matteo Figliuzzi, Remi Monasson, and Martin Weigt.  
413 Inverse Statistical Physics of Protein Sequences: A Key Issues Review. *Reports on Progress in*  
414 *Physics*, 81(3):032601, March 2018. ISSN 0034-4885, 1361-6633. doi:10.1088/1361-6633/aa9965.
- 415 [12] Adam J. Hockenberry and Claus O. Wilke. Phylogenetic Weighting Does Little to Improve  
416 the Accuracy of Evolutionary Coupling Analyses. *Entropy*, 21(10):1000, October 2019. ISSN  
417 1099-4300. doi:10.3390/e21101000.
- 418 [13] Edwin Rodriguez Horta and Martin Weigt. On the effect of phylogenetic correlations in  
419 coevolution-based contact prediction in proteins. *PLoS computational biology*, 17(5):e1008957,  
420 May 2021. ISSN 1553-7358. doi:10.1371/journal.pcbi.1008957.
- 421 [14] Jose Alberto de la Paz, Charisse M. Nartey, Monisha Yuvaraj, and Faruck Morcos. Epistatic  
422 contributions promote the unification of incompatible models of neutral molecular evolution.  
423 *Proceedings of the National Academy of Sciences*, page 201913071, March 2020. ISSN 0027-8424,  
424 1091-6490. doi:10.1073/pnas.1913071117.
- 425 [15] Bastien Boussau and Manolo Gouy. Efficient Likelihood Computations with Nonreversible  
426 Models of Evolution. *Systematic Biology*, 55(5):756–768, October 2006. ISSN 1063-5157.  
427 doi:10.1080/10635150600975218.
- 428 [16] Andrew Rambaut and Nicholas C. Grass. Seq-Gen: An application for the Monte Carlo  
429 simulation of DNA sequence evolution along phylogenetic trees. *Bioinformatics*, 13(3):235–238,  
430 June 1997. ISSN 1367-4803. doi:10.1093/bioinformatics/13.3.235.
- 431 [17] Bui Quang Minh, Heiko A Schmidt, Olga Chernomor, Dominik Schremppf, Michael D Woodhams,  
432 Arndt von Haeseler, and Robert Lanfear. IQ-TREE 2: New Models and Efficient Methods for  
433 Phylogenetic Inference in the Genomic Era. *Molecular Biology and Evolution*, 37(5):1530–1534,  
434 May 2020. ISSN 0737-4038. doi:10.1093/molbev/msaa015.
- 435 [18] Yeonwoo Park, Brian P. H. Metzger, and Joseph W. Thornton. Epistatic drift causes gradual  
436 decay of predictability in protein evolution. *Science*, 376(6595):823–830, May 2022. doi:  
437 10.1126/science.abn6895.
- 438 [19] Leonardo Di Bari, Matteo Bisardi, Sabrina Cotogno, Martin Weigt, and Francesco Zamponi.  
439 Emergent time scales of epistasis in protein evolution, March 2024.
- 440 [20] Ziheng Yang. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Molecular Biology*  
441 *and Evolution*, 24(8):1586–1591, August 2007. ISSN 0737-4038. doi:10.1093/molbev/msm088.

- 442 [21] Paul D. Williams, David D. Pollock, Benjamin P. Blackburne, and Richard A. Goldstein.  
443 Assessing the Accuracy of Ancestral Protein Reconstruction Methods. *PLOS Computational*  
444 *Biology*, 2(6):e69, June 2006. ISSN 1553-7358. doi:10.1371/journal.pcbi.0020069.
- 445 [22] William P. Russ, Matteo Figliuzzi, Christian Stocker, Pierre Barrat-Charlaix, Michael Socol-  
446 ich, Peter Kast, Donald Hilvert, Remi Monasson, Simona Cocco, Martin Weigt, and Rama  
447 Ranganathan. An evolution-based model for designing chorismate mutase enzymes. *Science*,  
448 369(6502):440–445, July 2020. ISSN 0036-8075, 1095-9203. doi:10.1126/science.aba3304.
- 449 [23] Matteo Bisardi, Juan Rodriguez-Rivas, Francesco Zamponi, and Martin Weigt. Modeling  
450 Sequence-Space Exploration and Emergence of Epistatic Signals in Protein Evolution.  
451 *Molecular Biology and Evolution*, page msab321, November 2021. ISSN 1537-1719. doi:  
452 10.1093/molbev/msab321.
- 453 [24] Sophia Alvarez, Charisse M. Nartey, Nicholas Mercado, Jose Alberto de la Paz, Tea Huseinbe-  
454 govic, and Faruck Morcos. In vivo functional phenotypes from a computational epistatic model  
455 of evolution. *Proceedings of the National Academy of Sciences*, 121(6):e2308895121, February  
456 2024. doi:10.1073/pnas.2308895121.
- 457 [25] Michael A. Stiffler, Frank J. Poelwijk, Kelly P. Brock, Richard R. Stein, Adam Riesselman,  
458 Joan Teyra, Sachdev S. Sidhu, Debora S. Marks, Nicholas P. Gauthier, and Chris Sander.  
459 Protein Structure from Experimental Evolution. *Cell Systems*, 10(1):15–24.e5, January 2020.  
460 ISSN 24054712. doi:10.1016/j.cels.2019.11.008.

461      **Supplementary Material: Reconstruction of ancestral protein**  
462      **sequences using autoregressive generative models**

463      **Appendix A: Autoregressive evolution model**

- 464      • Show the simplified expression for H=1  
465      • discuss irreversibility, with example

466      **Appendix B: Supplementary figures**

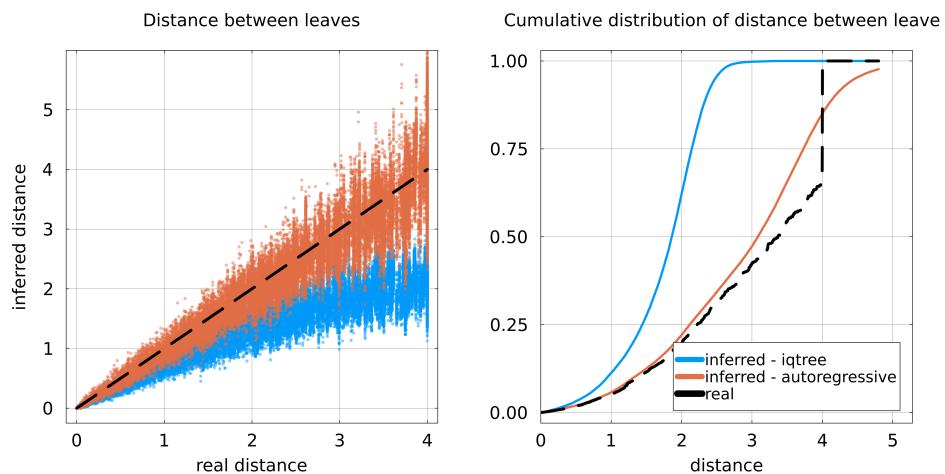


Figure S 1. Quality of branch length inference using data simulated with the autoregressive evolver. Inference is performed using the topology of a tree and leaf sequences generated using the autoregressive evolver. Two techniques are compared: IQ-TREE and the profile model corresponding to the autoregressive evolver. **Left:** inferred distance vs distance in the real trees for every pair of leaves. **Right:** Cumulative distribution of pairwise distance between leaves for the two inference methods and for the real trees. The discontinuity in the curve for the real tree is caused by the ultrametricity and fixed total height of the generated trees.

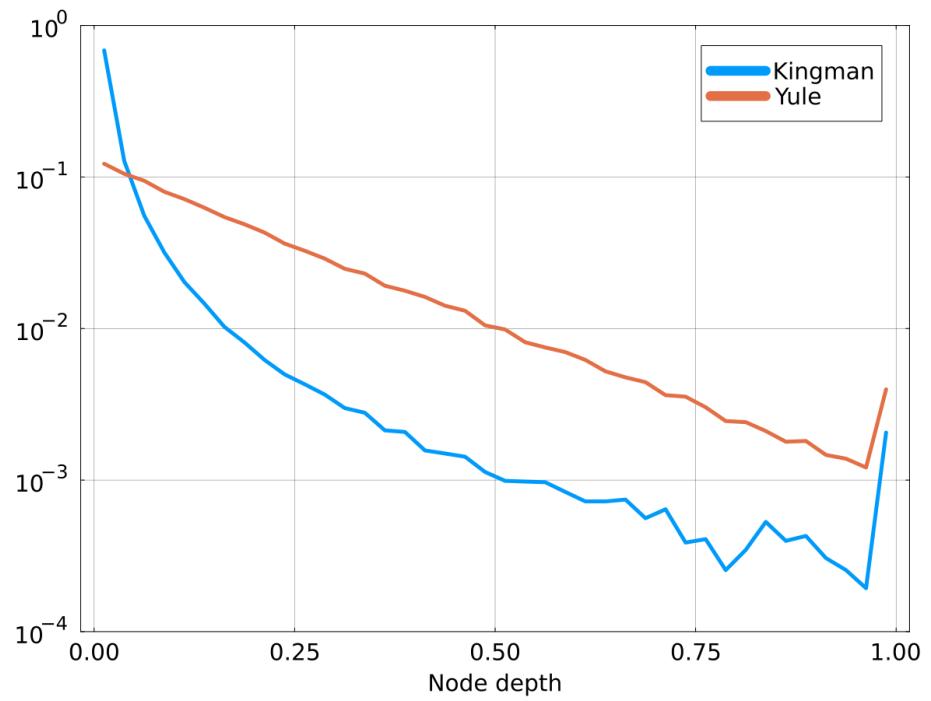


Figure S 2. Distribution of node depth for trees coming from the Kingman and Yule coalescents. Node depth is defined as the distance from a node to the closest leaf. Data is obtained by sampling several trees from each coalescent. Heights of trees are normalized to one. The Kingman process concentrates most of the nodes in close vicinity to the leaves, while the Yule process spreads them more evenly.