

# Are profile mixture models over-parameterized?

HECTOR BAÑOS<sup>1,2,3,\*</sup>, EDWARD SUSKO<sup>2,3</sup>, AND ANDREW J. ROGER<sup>1,3</sup>

<sup>1</sup> *Department of Biochemistry and Molecular Biology, Dalhousie University, Halifax, Nova Scotia, B3H 4R2, Canada*

<sup>2</sup> *Department of Mathematics and Statistics, Dalhousie University, Halifax, Nova Scotia, B3H 4R2, Canada*

<sup>3</sup> *Institute for Comparative Genomics and Evolutionary Bioinformatics, Dalhousie University, Halifax, Nova Scotia, B3H 4R2, Canada*

\**hbanos@dal.ca*

## ABSTRACT

1 Site heterogeneity of the amino acid substitution process accounts for the biochemical  
2 constraints on the range of admissible amino acids at specific sites. Phylogenetic models of  
3 protein sequence evolution that do not account for site heterogeneity are more prone to  
4 long-branch attraction artifacts.

5 Profile mixture models are used to model site heterogeneity. Even though model,  
6 tree, and mixing parameters are statistically consistent, the performance of these models  
7 with short alignments is unclear. Here we explore the behavior of tree topology estimates  
8 and marginal cumulative distributions with short simulated alignments. We find that  
9 over-parameterization is not a problem for complex profile mixture models and that simple  
10 models behave poorly. Misspecification of the frequency distributions does not cause a  
11 problem if the estimated cumulative distribution function adequately approximates the  
12 true one. Also, we find that misspecification of the exchangeabilities can severely affect  
13 parameter estimation and that an increase in likelihood does not necessarily reflect better  
14 tree estimation. Although the inclusion of more taxa often helps, it can hurt estimation if  
15 the exchangeabilities are badly misspecified.

16 Finally, we explore the effects of including an ‘F-class’ with the overall amino acid

17 frequencies of the dataset as an additional class in the profile mixture model. Surprisingly,  
18 the F-class does not seem to help parameter estimation significantly, and it can decrease  
19 the probability of correct tree estimation, depending on the scenario, despite the fact that  
20 it tends to improve likelihood scores. We also investigate this with several empirical data  
21 sets.

22 *Key words:* Phylogenetics; Mixture model; Frequency profile mixtures; Long-branch  
23 attraction.

24

25 Phylogenetic methods have been used to resolve many deep phylogenetic problems  
26 in the tree of life (Brown et al. (2013); Daubin (2002); Pisani et al. (2015); Raymann et al.  
27 (2015); Wickett et al. (2014)). To decrease estimate variability, these models require a  
28 large number of orthologous genes. The alignments of multiple genes (proteins) are either  
29 concatenated into a ‘supermatrix’ from which trees are estimated or individual  
30 gene/protein trees are first estimated and then combined using supertree or ‘species tree’  
31 methods. In either case, as more genes or proteins are considered, systematic biases can  
32 arise (Philippe et al. (2011)), underscoring the importance of adequately modeling the  
33 nucleotide or amino acid substitution process.

34 The substitution process of amino acid sequences is usually modeled as a  
35 site-independent Markov process in a tree. The most common approach assumes constant  
36 stationary frequencies of the amino acids and a constant matrix of exchangeabilities  
37 throughout the tree. The amino acid frequencies are usually estimated from the observed  
38 frequencies in the entire alignment. The matrix of exchangeabilities is fixed *a priori*,  
39 chosen from a set of empirically defined matrices, see for example Jones et al. (1992); Le  
40 and Gascuel (2008a); Whelan and Goldman (2001). Also, it is customary to consider  
41 different rates across sites to accommodate faster or slower substitution processes at  
42 different sites (see for example Yang (1994)). These models with almost the same

43 substitution process at each site and where the only difference comes from distinct rates  
44 across sites, are known as *site frequency homogeneous models*. From now on, and as it is  
45 customary, we refer to these as *site-homogeneous*, although these are not strictly  
46 homogeneous since site-rates may differ per site.

47       However, there are different ranges of amino acids admissible at sites in proteins  
48 because of functional or structural restrictions (Franzosa and Xia (2009); Goldstein (2008);  
49 Pál et al. (2006)). These ranges can vary widely, from a few or just one, to essentially all  
50 possible amino acids at a site (Halpern and Bruno (1998); Lartillot et al. (2007); Lartillot  
51 and Philippe (2004); Wang et al. (2008)). Consequently, site-homogeneous models, which  
52 overlook this across-site amino acid frequency heterogeneity, are less biologically plausible  
53 and are prone to long-branch attraction (LBA) artifacts (Feuda et al. (2017); Lartillot  
54 et al. (2007); Simion et al. (2017); Wang et al. (2008); Williams et al. (2013)). LBA is a  
55 pervasive systematic bias in tree estimation whereby distantly-related groups with long  
56 branches are artefactually grouped together (Felsenstein (1978); Philippe and Laurent  
57 (1998)). Partition (Lanfear et al. (2016); Pupko et al. (2002); Yang (1996)) and mixture  
58 models (Lartillot and Philippe (2004); Le and Gascuel (2008a); Schrepf et al. (2020); Le  
59 and Gascuel (2008b); Wang et al. (2008)) have been used to model heterogeneity of the  
60 amino acid substitution process across sites.

61       The CAT model (Lartillot and Philippe (2004)), a popular Bayesian mixture model,  
62 was shown to be less prone to LBA artifacts and to fit data better than site-homogeneous  
63 models (Lartillot et al. (2007)). In this model, all frequency vectors are assumed to be  
64 independently and identically drawn from a Dirichlet process model which effectively  
65 allows non-parametric estimation of the mixing distribution. Unfortunately, in its current  
66 implementation (Lartillot et al. (2013)), convergence may not be achieved in practice for  
67 large data sets. In the maximum likelihood framework, models accounting for site  
68 heterogeneity include mixture of the substitution rate matrices predefined for sites coming  
69 from different secondary structural elements and surface accessibility classes (Goldman

70 et al. (1998); Le and Gascuel (2008a, 2010)), or for different site rates (Le et al. (2012)),  
71 and a mixture of amino acid site frequency profiles (Schrempf et al. (2020); Le and Gascuel  
72 (2008b); Wang et al. (2008, 2014)). The latter, known as *profile mixture models*, have  
73 become widely used for analyses of deep phylogenetic problems.

74 Frequency vectors and weights pre-estimated from data bases of alignments are  
75 frequently used to reduce the complexity and computational cost of estimation with profile  
76 mixture models (Schrempf et al. (2020); Le and Gascuel (2008b); Wang et al. (2008)). We  
77 refer to a set of frequency vectors with their corresponding weights as a *mixing*  
78 *distribution*, and we refer to the indices of the frequency vectors in the mixing distribution  
79 as the *classes*. Similar to the empirical estimates of rate matrices, such mixing distributions  
80 are estimated from large data sets such as those described in Dufayard et al. (2005)  
81 and Sander and Schneider (1994). The techniques used to obtain empirical estimates of  
82 these mixing distributions vary. For example, in Le and Gascuel (2008b) the authors  
83 introduced six mixing distributions having 10, 20, 30, 40, 50, and 60 classes that were  
84 estimated from large data sets by ML estimation. These are known as the C10-C60 (or  
85 generically CXX) mixing distributions. Schrempf and colleagues (Schrempf et al. (2020))  
86 used K-means and the CAT model to estimate empirical mixing distributions ranging from  
87 4, 8, 16, up to 4096 classes. These are known as the UDM mixing distributions.

88 Profile mixture models are less susceptible to LBA than site-homogeneous models  
89 (Wang et al. (2008)). Also, as carefully described in the next section, these models have  
90 desirable properties when inferring parameters. For example, identifiability of the tree and  
91 mixing distribution is known to hold for a large subclass of models (Yourdkhani et al.  
92 (2021)). As a consequence, the tree and the mixing parameters are statistically consistent  
93 even with a large number of profiles. Informally speaking, this means that, if the model is  
94 correctly specified, one can effectively estimate the true parameters as the number of sites  
95 increases.

96 Although the identifiability and consistency properties satisfied by profile mixture

97 models are desirable in any modeling context, there is no guarantee that such models will  
98 have good small sample properties. Nor it is known what problems may arise from model  
99 misspecification. Thus there is a need to explore the performance of profile mixture models  
100 through simulations of short alignments both with and without model misspecification.  
101 One of our main motivations is to determine whether the large numbers of parameters in  
102 profile mixtures create problems with small samples. Specifically, we want to determine if  
103 there is excessive variability in estimates from models with too many parameters relative  
104 to the sample size.

105 By varying several parameters of empirically-derived profile mixture models, we  
106 simulated distinct alignments of lengths 300, 600, and 1000 (the approximate lengths of  
107 true alignments of single proteins). A detailed explanation of the different simulation  
108 model settings is provided below. We fit distinct mixing distributions and matrices of  
109 exchangeabilities to each simulation. We assessed model performance using four criteria  
110 described in detail in the Materials and Methods. Two of these criteria concern the tree  
111 topology MLE accuracy and variability. The other two are a measure for comparing the  
112 marginal cumulative distribution functions (CDFs) inherited from the observed and  
113 expected mixing distributions. These CDFs, properly described in the following section,  
114 are an alternative way to re-parameterize profile mixture models.

115 All findings are presented later in detail, but the major highlights include the  
116 following:

- 117 • When the exchangeabilities and frequency vectors are correctly specified, there is no  
118 evidence of model over-parameterization (over-fitting), even when there are many  
119 classes with zero weight estimates. This relates to a concern articulated in several  
120 studies (e.g. Anderson and Lindgren (2021); Li et al. (2021)) that fitted complex  
121 mixture models with classes estimated to have zero weights are over-parameterized  
122 and should be avoided in favor of simpler models. Also, the inclusion of more taxa  
123 improves tree estimation.

- 124 • When there is misspecification of the frequency classes, we observe that tree  
125 estimation is not necessarily acutely affected. If the set of frequency vectors is  
126 sufficiently rich, the estimated CDF closely approximates the CDF of the generating  
127 model, and when that occurs, the frequency of correct tree estimations is large.
- 128 • Severe problems can arise from the misspecification of the exchangeabilities. We  
129 observe that this scenario can lead to a bias in the MLE mixture weights. This bias  
130 favors parameters that maximize the likelihood but decrease the similarity between  
131 observed and expected CDF. This produces a decline in tree estimation accuracy.  
132 Under these conditions, adding taxa does not necessarily improve tree estimation.

133 We also explore the effects of the “F-class,” a class that is defined from the  
134 empirical frequencies of amino acids from the overall alignment, that is often included as  
135 an additional class in models to account for remaining sites in the data that are not well  
136 modeled by the fixed empirically-derived frequency vectors. However, we find that the  
137 F-class does not significantly improve tree estimation, and, in some cases, may compromise  
138 accuracy. This exploration is complemented by looking at empirical data. Our analyses  
139 suggest that while the F-class increases the likelihood significantly, it may lead to  
140 erroneous tree estimation.

## 141 MATERIALS AND METHODS

### 142 *Mixture models and over-parameterization with large samples*

143 In this section, we define and elaborate on some theoretical properties of profile  
144 mixture models. Re-parameterizing the mixing distributions as CDFs provides insight into  
145 why over-parameterization is less of a problem for profile mixture models than it might be  
146 for models without stringent parameter constraints. Also, we briefly discuss known  
147 identifiability results for such models.

148 Roughly speaking, profile mixture models are mixtures of time-reversible models,

149 with a common exchangeability matrix  $R$ . The parameter space  $\Theta$  of a profile mixture  
150 model with  $C$  classes is defined by:

- 151 (i) A rooted metric tree  $T$  on  $N$  taxa.
- 152 (ii) A symmetric  $20 \times 20$  matrix of non-negative exchangeabilities  $R$ .
- 153 (iii) For  $c = 1, 2, \dots, C$ , a frequency distribution vector  $\boldsymbol{\pi}_c$ , and a weight  $w_c$ , with  $w_c > 0$   
154 and  $\sum_{c=1}^C w_c = 1$ .
- 155 (iv) A collection of  $K$  scalar rate parameters  $\{r_k\}$ , with  $r_k \geq 0$ , and rate weight  $d_k$ , with  
156  $d_k > 0$  and  $\sum_{k=1}^K d_k = 1$ .

157 The substitution process of a profile mixture model is as follows: for each site, a  
158 frequency vector  $\boldsymbol{\pi}_c$  is sampled with probability  $w_c$ ; and a rate parameter  $r_k$  is sampled  
159 with probability  $d_k$ . Evolution of a sequence at a site is then according to a continuous  
160 Markov substitution process over tree  $T$  with exchangeabilities  $R$ , root distribution  $\boldsymbol{\pi}_c$ , and  
161 rate  $r_k$ . For a given site pattern  $\boldsymbol{x}_i$ , the likelihood function is determined by a weighted  
162 average of partial site likelihoods conditional on each site-profile class and site-rate class:

$$L(\theta|\boldsymbol{x}_i) = \sum_{c=1}^C w_c \sum_{k=1}^K d_k P(\boldsymbol{x}_i|T, R, \boldsymbol{\pi}_c, r_k),$$

163 where  $\theta = (T, R, \{\boldsymbol{\pi}_c\}, \{w_c\}, \{r_k\}, \{d_k\}) \in \Theta$ .

164 For a model with fixed frequency variables, a natural way of parameterizing the  
165 mixture model is in terms of its weights,  $w_c$ ,  $c = 1, \dots, C$ . This leads to models of differing  
166 dimensions  $C$  that, as mentioned before, can get very large, raising concerns about  
167 over-parameterization.

168 An alternative way of parameterizing the mixture is in terms of its cumulative  
169 distribution function (CDF):

$$G(\boldsymbol{\pi}) = P(\mathbf{\Pi} \leq \boldsymbol{\pi}),$$

170 where  $\mathbf{\Pi}$  represents the random frequency vector for a site. This allows one to express  
171 mixing distributions with differing components ( $C = 20$  say or  $C = 60$ ) as being in the

172 same parameter space. However, now the frequency mixture parameter space, which is a  
173 space of distribution functions, is infinite-dimensional, and would appear an extreme case  
174 of over-parameterization.

175 Surprisingly, even for this infinite dimensional space of distribution functions,  
176 estimation of both the mixing distribution and structural parameters like the tree is  
177 frequently still consistent as was shown by Kiefer and Wolfowitz (1956) for a wide class of  
178 models under mild regularity conditions. Most of these conditions are expected to hold for  
179 the models considered here. A more detailed explanation of how the results in Kiefer and  
180 Wolfowitz (1956) apply to our context is given in the Appendix.

181 The implication of Kiefer and Wolfowitz (1956) is that class frequency mixture  
182 models are not overparameterized, at least with large samples. The reason for this is that  
183 the space of all distribution functions as a space of functions is relatively “small” in the  
184 mathematical sense of being a compact space (a closed and bounded space in our setting).  
185 An alternative way of seeing why this is the case is to note that the  $w_c$  are restricted to be  
186 non-negative and sum to one. By contrast in cases of true over-parameterization,  
187 parameters are unrestricted (for example, a regression model where there are more  
188 predictors than observations).

189 Another surprising result of estimation within the mixing distribution setting is  
190 that even if parameter estimation is unrestricted and any mixing distribution is allowed,  
191 the maximum likelihood estimator will be a finite mixing distribution: i.e. it will be  
192 describable in terms of a fixed set of weights  $w_1, \dots, w_C$  for some  $C$ . This is an implication  
193 of the results of Lindsay (1983) as detailed in the Appendix.

194 Adding to this, in Yourdkhani et al. (2021) identifiability results are reported for a  
195 large family of profile mixture models. These authors showed generic identifiability of the  
196 tree and mixing parameters for models with  $C \cdot K < 72$ , where  $C$  is the number of classes  
197 and  $K$  the number of rates, trees with more than 8 taxa, and where no parameters in  $\Theta$   
198 are assumed to be fixed. As discussed in the Appendix, these results also apply to some of



199 the models with fixed frequency vectors considered here. We conjecture that, generically,  
200 identifiability of the tree topology can be achieved for models with fixed frequency vectors  
201 with  $C$  classes,  $K$  rates and  $m$  taxa, for some  $C \cdot K > 72$ , and all  $m > 8$ . This will be  
202 explored in future work.

### 203 *Simulation setting*

204 In this section, we describe all the different parameters used to simulate alignments  
205 under profile mixture models. These are presented below in the following order: (1) the  
206 trees; (2) mixing distributions; (3) exchangeability matrices; (4) sequence lengths; and (5)  
207 rate parameters.

208 By choosing different combinations of parameters, we simulated a total of 108  
209 scenarios. For each scenario, 100 simulations were performed using `Alisim` (Ly-Trong et al.  
210 (2021)). We now proceed to describe all the choices of parameters.

211 *Trees* Nine different trees are considered for these simulations. All trees have the  
212 ‘structure’ of tree  $T$  shown in Figure 1. The features that vary per tree are the length of a  
213 single edge  $l$ , where  $l \in \{0.005, 0.02, 0.05\}$  and the number of taxa at each polytomy  $m$ ,  
214 where  $m \in \{1, 2, 3\}$ . We denote each of the trees by  $T_{6m}(l)$ .

215 The structure of  $T$  is chosen since it is often a tree susceptible to LBA artifacts. For  
216 fixed  $l$  and changing  $m$ , the simulations are, in effect, all from the same tree but with  
217 differing levels of taxonomic sampling from the 6 clades. By increasing the number of taxa,  
218 we obtain more information on the frequency vectors. By decreasing the edge length  $l$ , we  
219 make the tree more susceptible to LBA artifacts whereby the f-clade (i.e. clade including  
220 taxa  $f_1, f_2, \dots, f_m$ ) and the e-clade (i.e. clade with taxa  $e_1, e_2, \dots, e_m$ ) group together to the  
221 exclusion of the other clades.

222 *Mixing distributions* For each tree, we simulated data using ten different mixing  
223 distributions. One of these is the model C60 as defined in Le and Gascuel (2008b), which

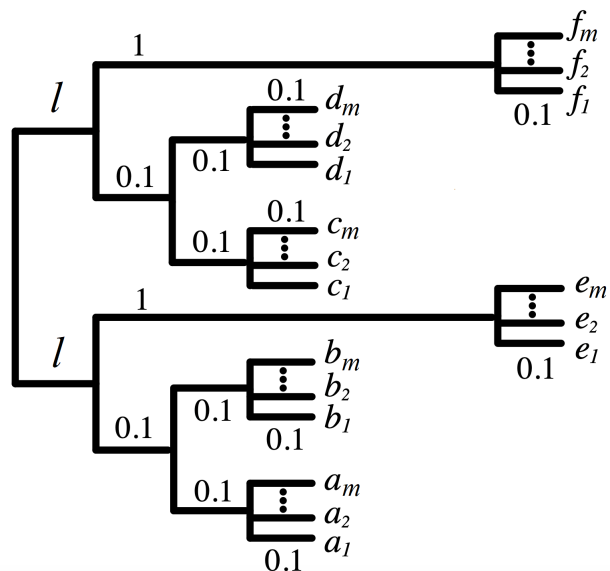


Fig. 1. The main structure of the tree where all the simulations were conducted. Only the edge lengths  $l$  and the number of taxa at each polytomy  $m$  are variable. This tree has  $6m$  taxa

224 has 60 classes. Another seven are built on the frequency vectors of C60. Specifically, for  
 225  $i \in \{10, 15, 20, 30, 40, 50, 60\}$ , we defined the mixing distribution  $C60[i]$  by choosing  $i$   
 226 frequency vectors from C60 at random and assigning them non-zero weights sampled from  
 227 a Dirichlet distribution with concentration parameters  $\alpha = \mathbf{1}$ . We denote by  $\mathfrak{C}60$  the set of  
 228 mixing distributions  $C60[i]$ , for all  $i$ , including C60. The two remaining mixing  
 229 distributions used to simulate data are known as UDM-0256, and UDM-4096 (Schrepf  
 230 et al. (2020)), that have 256 and 4096 classes respectively. Specifically, we used the  
 231 non-transformed mixing distributions denoted as UDM-0256-None and UDM-4096-None in  
 232 Schrepf et al. (2020).

233 In reality, every site in a protein has a unique physicochemical environment that  
 234 would likely be better fit by a site-specific frequency vector. Therefore profile mixture  
 235 models with fixed frequencies try to approximate commonly occurring patterns amongst  
 236 sites represented by site classes. By simulating under C60, the most complex of the CXX  
 237 mixing distributions, and the complex UDM mixing distributions, we try to emulate real  
 238 data. The  $C60[i]$  distributions reflect the scenario wherein, for a small sample, not all

239 relevant frequency vectors are represented, nor are distributed as in C60. We note that the  
240 UDM mixing distributions include additional components that are not in the CXX  
241 distributions but because the data-sets used to estimate both CXX and UDM mixing  
242 distributions overlap, so some similarities in their frequency profiles are expected.

243 *Exchangeability Matrices, Sequence Lengths, and Site Rate Variation* For all  
244 combinations of mixing distributions and trees, we simulated data using the  
245 exchangeability matrix from the LG model (Le and Gascuel (2008a)); we refer to this as  
246 the LG matrix. We also used a “POISSON” exchangeability matrix, a matrix with equal  
247 exchangeabilities, but only for simulations involving the two UDM mixing distributions.

248 For all combinations of mixing distributions, trees, and matrices, we used sequence  
249 lengths of 300, 600, and 1000 amino acids. In practice, alignments of length 300 are more  
250 typical for single protein data sets. The other choices of sequence lengths are meant to  
251 capture the effects in tree estimation with increasing sequence length.

252 Lastly, for all simulations we used four rate parameters coming from a discrete- $\Gamma(4)$   
253 distribution (Yang (1994)) with  $\alpha = 0.5$ .

#### 254 *Fitted models and Precision of parameter estimation*

255 In this section, we describe how different choices of fixed frequency vectors and  
256 exchangeabilities were fitted to the simulations. We also introduce four criteria used to  
257 evaluate model fitness.

258 When fitting a model to data, estimated parameters were obtained by maximum  
259 likelihood. IQ-tree2 (Minh et al. (2020)) was used to get the maximum log-likelihoods and  
260 the estimators (MLEs) for all simulations. The parameters that were optimized by  
261 IQ-tree2 are the tree topology, edge lengths, and weights of the frequency vectors. All  
262 other parameters including the frequency vectors, the matrix of exchangeabilities, and the  
263 rate parameters, were supplied as fixed values to IQ-tree2.

264 The LG matrix was the only matrix used to simulate data under the mixing  
265 distributions in  $\mathfrak{C}60$ . To these simulations, we fit the LG matrix, the F-class, and various  
266 frequency vectors. Specifically, the frequency vectors fitted included C60, C40, C30, C20,  
267 defined in Le and Gascuel (2008b), LG (Le and Gascuel (2008a)), LG4X (Le et al. (2012)),  
268 CK36 as defined below, as well as the frequency vectors in  $\mathfrak{C}60$  used to simulate the data.  
269 In these analyses, we explored mainly two things: (1) Possible model over-parameterization  
270 from fitting more general models to short alignments; and (2) misspecification of the  
271 frequency vectors by using models that have different frequency vectors than those of the  
272 generating model.

273 The frequency vectors CK36 were obtained from the cluster centers derived from a  
274 k-means algorithm (Hartigan and Wong (1979)) on all the classes of C20, C30, C40, and  
275 C60 with  $k = 36$  (the choice of  $k$  was determined by the elbow method (Thorndike  
276 (1953))).

277 For data generated under the UDM mixing distributions, we fitted the  
278 exchangeability matrices POISSON and LG, with and without the F-class, and frequency  
279 vectors of C60, C40, C30, C20, CK36, LG4X. We also fitted the POISSON model for data  
280 generated under POISSON exchangeabilities and the LG model for data generated under  
281 LG exchangeabilities. In these analyses we primarily explored the effect of: (1)  
282 Misspecification of the frequency vectors; (2) misspecification of the matrix of  
283 exchangeabilities; and (3) use of the F-class in estimation. Recall that only the mixing  
284 weights, the tree topology, and edge-lengths are optimized, everything else is fixed before  
285 the maximization of the likelihood.

286 To simplify the presentation of results, we only considered a subspace of tree space.  
287 In preliminary results, we computed the likelihoods of all 105 6-taxon unrooted topologies  
288 for data generated under  $T_6(0.005)$ . We noted that there were two tiers in terms of  
289 log-likelihood values; these occurred regardless of the generating and fitted models. One  
290 tier consists of the log-likelihoods of the 35 topologies shown in Table S1 in the

291 Supplementary Material, and the other consists of the remaining 70 trees. The first tier  
292 showed significantly larger log-likelihood values, see Figure S1 in the Supplementary  
293 Material. The 35 topologies in the first tier are all the topologies displaying the embedded  
294 quartet tree  $AB|CD$ . This shows, as expected, there are no problems with estimating the  
295 relationships amongst the taxa in this quartet. The main difficulty was instead determining  
296 the correct placement of the long branches because of the LBA-related artifacts.

297 Therefore, for simplicity, in all cases we restrict the tree space considered to just  
298 these 35 tree topologies by substituting  $X \in \{A, B, C, D, E, F\}$  by the adequate  $m$ -taxon  
299 polytomy. We also consider in this tree space the ‘star tree’ topology obtained from  $T$  in  
300 Figure 1 by setting  $l = 0$ .

301 We now proceed to introduce the criteria used to compare the overall model  
302 performance.

303 *Mean Integrated Squared Error and Maximal Difference* As mentioned earlier,  
304 profile mixture models can be parameterized as CDFs. Therefore, a reasonable way to  
305 measure the precision of parameter estimation is comparing how closely the observed CDF  
306 resembles the true one. Unfortunately, assessing the precision of parameter estimate *via*  
307 the expected and observed CDFs is burdensome due to the complexity of the  
308 20-dimensional space in which they reside. Here we use marginal CDFs instead to assess  
309 the precision of parameter estimation.

To compare marginal CDFs we used two measures. The first one is the *Mean Integrated Squared error* (MISE) (Scott (1992)), also known as  $L^2$  risk function, which is defined as follows

$$\text{MISE} = \frac{1}{20} \sum_{i=0}^{20} \int_0^1 (G_i(x) - \hat{G}_i(x))^2 dx,$$

310 where  $G_i$  is the true marginal CDF corresponding to amino acid  $i$ , and  $\hat{G}_i(x)$  is the  
311 marginal CDF obtained from the estimated mixing distribution. Here all the integrals were  
312 computed using the function `integrate` from the R package `pracma` with default settings.

The second measure consists in the maximum difference (MD) between estimated marginal CDFs and the true ones. This is borrowed from the Kolmogorov–Smirnov test (Massey (1951)), and it is defined as follows:

$$\text{MD} = \frac{1}{20} \sum_{i=0}^{20} \max_{x \in [0,1]} \{|G_i(x) - \hat{G}_i(x)|\},$$

313 where  $G_i$  and  $\hat{G}_i(x)$  are as defined for the MISE. Here the maximum of the absolute  
314 difference is computed using the function `optimize` from the R package `stats` with default  
315 settings.

316 For a given choice of parameters, we report the mean MISE and MD for all 100  
317 simulations. Note that in both cases, as these measures approach to zero, the observed  
318 CDF approaches the true CDF. While these two measures cannot be used in practice since  
319 the true CDF is then unknown, in this case and as detailed in the Results section, these  
320 measures help us assess model over-parameterization among other things.

321 In order to assess the precision of parameter estimation in a more standard way, we  
322 also looked at two more criteria based on the tree topology estimate.

323 *Overall Accuracy and Proportional Mode* The following criteria, *overall accuracy*  
324 denoted OA, is a standard way to evaluate the precision of parameter estimation. Given a  
325 set of simulations, OA is defined as the proportion of these where the true tree topology  
326 was the one maximizing the likelihood.

327 The next criterion, the *proportion of settings where the true tree was the mode of*  
328 *the distribution of estimated trees*, denoted PM, evaluates the precision of parameter  
329 estimation in a manner that can indicate whether there are biases in estimation. PM  
330 consists of the proportion of settings where the true topology was chosen the most. Given  
331 several sets of simulations, PM is the proportion of these sets where the true tree topology  
332 was the MLE more often than any other topology. For a given method, the tree estimated  
333 most frequently is the mode of the distribution of estimated trees. So alternatively, PM is  
334 the frequency with which the true tree is the mode of the distribution of estimated trees.

335 To account for some sampling variability in PM, in each scenario we tested if the  
336 proportion of times the true topology was chosen was significantly higher than the  
337 runner-up topology using a binomial test. PM is obtained from dividing the number of  
338 scenarios where the true topology is significantly more likely than the runner-up (rejecting  
339 the null hypothesis  $H_0 : p = 0.5$ ) by the total number of scenarios.

340 Note that OA and PM may not necessarily be strongly correlated. Suppose, for  
341 instance, that for a given model there is a bias in estimation towards certain trees under  
342 certain settings. That could result in a large OA because of the large frequency of  
343 estimations of the true tree in settings where the model is biased towards it. But if the  
344 true tree is not always favored, the true tree would not be most frequently estimated,  
345 leading to a small PM. On the other hand, one could have really low OA but no other tree  
346 is chosen more times, leading to a high PM. Ideally, one would like to see both high OA  
347 and PM, which would indicate that the true tree is being chosen the most and with little  
348 variability across scenarios.

## 349 RESULTS

350 We present the results according to the set of mixing distributions we used to  
351 simulate the data. We start with  $\mathfrak{C}60$ , then UDM, followed by the empirical data-sets. For  
352 the remainder of the text and for simplicity, we denote by C60L the model C60 as defined  
353 in Le and Gascuel (2008b) including both frequency classes and optimal weights from that  
354 study. Then, when we discuss fitting one of the CXX models we are referring strictly to  
355 fitting just its frequency classes.

### 356 *Simulated Data Under $\mathfrak{C}60$ Mixing Distributions*

357 Table 1 displays the mean MISE for data generated under 72 different simulation  
358 conditions total; i.e. under eight mixing distributions in  $\mathfrak{C}60$  and nine trees. The eight  
359 mixing distributions included C60L and C60[ $i$ ] for  $i$  in  $\{10, 15, 20, 30, 40, 50, 60\}$ . We see in

	<b>PF+F</b>	<b>C60+F</b>	<b>C40+F</b>	<b>CK36+F</b>	<b>C30+F</b>	<b>C20+F</b>
<b>300</b>	0.00037	0.00042	0.00154	0.00118	0.00199	0.00333
<b>600</b>	0.00020	0.00023	0.00137	0.00102	0.00179	0.00309
<b>1000</b>	0.00013	0.00015	0.00129	0.00095	0.00170	0.00299

Table 1. The mean MISE for data generated under distributions in  $\mathfrak{C}60$ . For each of the 72 scenarios (9 trees and 8 models) per sequence length, we compute the mean MISE for the 100 simulations. Then we take the mean of all these values, which are the entries in this table. Label PF+F represents the overall performance in MISE when fitting the generating classes per scenario with the F-class.

360 this table, as expected, that estimating models that included only the frequency classes  
 361 used to generate the data plus an F-class, denoted here as perfect fit (PF)+F, are those  
 362 with the lowest MISE. The second lowest value is achieved by C60+F model, followed by  
 363 CK36+F, C40+F, C30+F, and C20+F, respectively. In this case, we note that the MISE  
 364 of PF+F and C60+F are really close. This behavior is similar for the MD criterion, as seen  
 365 in Table S2 in the Supplementary Material. These observations suggest that even though  
 366 C60+F fits 61 classes to the data, it still has a much better fit than the CXX+F models  
 367 with fewer classes even though, for many of the simulation settings, there were many fewer  
 368 classes present. This demonstrates that over-parameterization is not a problem for  
 369 complex, correctly specified models, even for models having several classes with zero  
 370 weights. Table S3 in the Supplementary Material shows the mean normalized MISE for the  
 371 same data as in Table 1, where MISE for any given scenario was re-scaled to give a sum of  
 372 1 over all fitted models.

373 This shows that the mean MISE values in Table 1 adequately consolidate all  
 374 scenarios and no biases between classes are introduced by a scenario with considerably  
 375 larger MISE values.

376 Figure 2 shows the plots of average OA (A) and PM (B) over all data generated  
 377 under mixing distributions in  $\mathfrak{C}60$ . Each dot in OA represents the proportion of times the  
 378 tree was correctly inferred over 2400 simulations (3 values of  $l$  and 8 sets of classes in  $\mathfrak{C}60$ ,  
 379 with 100 repetitions each), and PM is a proportion over 24 distinct scenarios. For  
 380 alignments of length 300, fitting C60+F, C40+F, C30+F, C20+F, or CK36+F, produces,



381 on average, no significant difference in OA. However, C60+F and CK36+F have better PM  
382 values. On the other hand, models LG+F and LG4X+F perform poorly as reflected by  
383 both OA and PM.

384 For longer alignments, C20+F has a significantly lower OA than C60+F, C40+F,  
385 C30+F, and CK36+F. This shows how as sample size increases, more complex models that  
386 approximate the true CDF better have superior performance. In this case, fitting C60+F  
387 and CK36+F still yields the best PM values. Overall these results reinforce the inference  
388 that over-parameterization of C60+F does not cause problems and misspecification of the  
389 frequency vectors (e.g. for CK36+F) does not compromise tree estimation if the estimated  
390 CDF adequately approximates the true CDF.

391 One concern is that C60+F could be fitting better on average because it closely  
392 resembles two of the mixing distributions – C60L and C60[60] – used in the foregoing  
393 simulations. However, C60+F also behaves well for data generated with fewer classes.  
394 Table S4 in the Supplementary Material shows the OA for data generated under C60[ $i$ ],  
395 with  $i \in 10, 30, 60, T_6(0.005)$ , and different sequence lengths. This table shows that there is  
396 little evidence of over-fitting when the estimating model has many more classes than the  
397 generating model. This is also reflected in the MISE scores when considering different  
398 generating classes (Supplementary Material Table S5).

399 Figure 2 also shows that increasing numbers of taxa improves tree estimation.  
400 When the number of taxa increases, MISE and MD scores decrease/improve (Table 2) and  
401 this behavior is consistent across models. As expected, we also observe that tree estimation  
402 accuracy also improves as sequence length increases. Analogously, MISE and MD scores  
403 decrease as sequence length increases (Table 2).

#### 404 *Simulated Data Under UDM Mixing Distributions*

405 Table 3 displays the mean MISE values over data generated under UDM-0256 with  
406 POISSON exchangeabilities. The MISE values are considered separately when fitting CXX

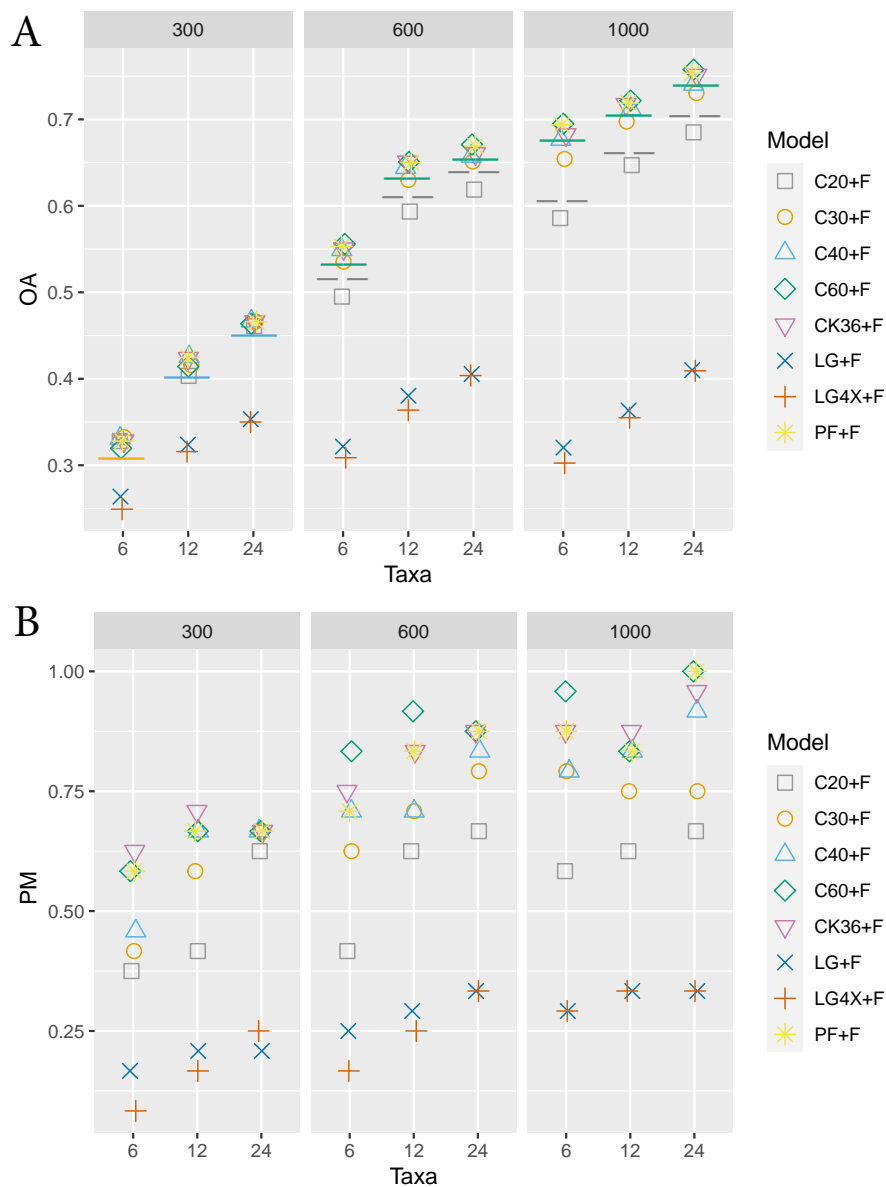


Fig. 2. (A) the plot of the OA values for various models fitting data generated under mixing distributions in  $\mathbb{C}60$  fitted to all frequency vectors. Label PF denotes the frequency vectors used to generate the data. The  $x$ -axis represents the number of taxa on the tree. The plot is divided by sequence lengths of 300, 600, and 1000. The lower bound of the 95% confidence interval (CI) of the model with the highest OA is depicted with a solid line whose color is in agreement with such model. Depicted with a dashed gray line is the higher bound of the CI of classes in C20 in the cases where such model was significantly worse than the best model. (B) A similar plot to that on top but for PM. For this case, no confidence interval can be computed.

Taxa	MISE			MD		
	300	600	1000	300	600	1000
<b>6</b>	0.00040	0.00020	0.00012	0.178	0.136	0.111
<b>12</b>	0.00031	0.00016	0.00008	0.155	0.117	0.085
<b>24</b>	0.00024	0.00013	0.00007	0.127	0.097	0.071

Table 2. The MISE and MD for fitted C60+F to data generated under  $T_{6m}(0.02)$  and C60[15] for all  $m$ . One can see a decrease in MISE and MD by either increasing the number of taxa, or the sequence length.

	LG matrix			POISSON matrix		
	300	600	1000	300	600	1000
<b>C60+F</b>	0.00724	0.00783	0.00822	0.00083	0.00064	0.00055
<b>C40+F</b>	0.00767	0.00822	0.00861	0.00091	0.00071	0.00062
<b>CK36+F</b>	0.00807	0.00856	0.00890	0.00089	0.00071	0.00063
<b>C30+F</b>	0.00852	0.00901	0.00931	0.00096	0.00076	0.00068
<b>C20+F</b>	0.00919	0.00964	0.00990	0.00111	0.00092	0.00085
<b>C60</b>	0.00357	0.00378	0.00388	0.00080	0.00062	0.00054
<b>C40</b>	0.00333	0.00337	0.00343	0.00087	0.00069	0.00061
<b>CK36</b>	0.00435	0.00448	0.00457	0.00086	0.00069	0.00062
<b>C30</b>	0.00422	0.00438	0.00444	0.00090	0.00073	0.00065
<b>C20</b>	0.00398	0.00399	0.00405	0.00103	0.00088	0.00082

Table 3. The mean MISE for data generated under UDM-0256 and POISSON exchangeabilities. For each of the 9 scenarios (9 trees) per sequence length, we compute the mean MISE for the 100 simulations. This is done separately when fitting the LG and POISSON matrices.

407 models and CK36 (with and without the F-class) to the POISSON and LG matrices. In  
 408 the former case, there is misspecification just of the classes whereas, for the latter, both  
 409 the classes and exchangeabilities are misspecified.

410 When fitting the LG matrix, i.e where there is misspecification of the  
 411 exchangeabilities, we see a significant difference in model fit between models including the  
 412 F-class and omitting it. The MISE scores are elevated for models including the F-class  
 413 relative to those without and this effect is independent of the sequence length (Table 3).  
 414 We believe this is not directly an artifact of the F-class *per se*, as it is discussed below.

415 When fitting the correctly specified POISSON matrix, for any given sequence  
 416 length and set of fitted frequency vectors, the mean MISE is comparable whether the  
 417 F-class is included or omitted. In most cases, we see a slightly better MISE when excluding  
 418 the F-class although the difference is minuscule and most likely not significant. We note

419 that the classes in C60 have the best MISE scores across all sequence lengths, but many of  
420 the other CXX models also performed well; C20 yielded the poorest scores overall. In this  
421 same table, we also see how MISE decreased as sequence length increased when there is no  
422 misspecification but this trend did not necessarily hold when there was. We believe this  
423 comes from the fact that as more data becomes available, more classes can be  
424 misestimated introducing more error.

425 Figures 3 and 4 show the plots of average OA (A) and PM (B) for data generated  
426 under UDM distributions and POISSON exchangeabilities. The former figure shows data  
427 fitted using models with POISSON exchangeabilities, and the latter, LG exchangeabilities.  
428 In these plots, each dot in the OA plot represents the proportion of times the tree was  
429 correctly inferred over 600 simulations (3 values of  $l$ , 2 sets of UDM classes, and 100  
430 repetitions), and for PM the proportion over 6 distinct scenarios.

431 Figure 3 shows the case when there was misspecification of classes only. Recall that  
432 the mixing distributions in UDM have 256 and 4096 frequency vectors, thus there is  
433 significant misspecification when fitting CXX and the CK36 models. Nevertheless, we saw  
434 that in all cases, fitting with any of the CXX mixtures led to reasonable performance. The  
435 OA estimates were even close to the case of no misspecification of the frequency classes in  
436 Figure 2.

437 In this case, there seems to be no significant difference in OA between fitting the  
438 F-class or discarding it. However, fitting without the F-class yielded, in some cases, better  
439 PM values. Similar to the C60 case, fitting with LG4X and POISSON with no site-profile  
440 mixture model also behaved both very poorly. In this case, we also see how the increase of  
441 taxa monotonically improved tree estimation. The exception was for sequence lengths of  
442 600 and 1000 when increasing taxa from 12 to 24. We consider these to reflect a tie in  
443 performance accounting for variability due to the finite number of simulations. Similar  
444 conclusions can be drawn for the case where the only difference was that the data was  
445 generated and fitted using the LG matrix instead of POISSON (Supplementary Material

446 Figure S2).

447 When there is misspecification of both classes and exchangeabilities as shown in  
448 Figure 4, we observe that for all cases, except one, OA was equal or significantly worse  
449 compared to Figure 3 and Figure S2 in the Supplementary material. We believe this is  
450 similar to the behavior shown in Table 3 described above. Furthermore, in some cases, an  
451 increase in numbers of taxa led to a decrease in OA (see Figure 4, 1000 sites).

452 The most striking, and perhaps surprising, feature in Figure 4 is that fitting with  
453 the F-class often led to a substantial drop in OA compared to when it is omitted. To  
454 investigate this further, we explored the impact of the F-class on the likelihood scores and  
455 mixing weights of the fitted models when exchangeabilities were misspecified (LG) versus  
456 correctly specified (POISSON). These analyses were based on 27000 observations (9 trees,  
457 3 sequence lengths, 2 UDM distributions, 5 fitted models, 100 simulations) and the results  
458 are shown in Figure 5. When there is misspecification of the exchangeabilities, the F-class  
459 frequently improves the likelihood values substantially, whereas when the exchangeabilities  
460 are correctly specified, only modest increases in likelihood are seen (see Figure 5 (A)). The  
461 frequently large increases in likelihood values with the F-class in the case of the LG  
462 exchangeabilities is surprising because, as indicated in Figure 4 and Table 3, better OA  
463 estimates and lower MISE scores are obtained in this case when there is no F-class.

464 Since models without the F-class are special cases of the comparable models that  
465 include F-classes, we should always expect an increase in likelihood when fitting the latter  
466 models. When there is no misspecification of the frequency classes, a crude approximation  
467 is that the null distribution has a mean of 5 and a standard deviation of 5.48 (from a  
468 mixture of a degenerate uniform[0] and a  $\chi^2$  with 20 degrees of freedom, see Self and Liang  
469 (1987)). When there is no misspecification, simulations have a mean of 3.8 and a standard  
470 deviation of 3.4, smaller than the natural increase in likelihood, alluded to above, that are  
471 expected with increases in the number of parameters estimated. By contrast, with model  
472 misspecification, the mean and standard deviation are 34.9 and 19.7, respectively. Thus

473 when there is no misspecification of the exchangeabilities, differences in log-likelihood are  
474 pretty small and would be judged small relative to crude chi-square approximations. By  
475 contrast, with misspecification, such differences are very large compared to expectations  
476 based on the number of parameters estimated.

477 We also investigated the impact of misspecification of exchangeabilities on the  
478 estimated weight of the F-class (Figure 5B). Figure 5 (B) contains two overlapping plots.  
479 When exchangeabilities are correctly specified, the F-class weights tend to be relatively  
480 small (e.g. mean weight of non-misspecified = 0.11), whereas when they are misspecified  
481 the weight distribution shifts dramatically to adopt larger values, often exceeding 0.5 (e.g.  
482 mean weight of misspecified = 0.62). Clearly, the weights of the F-class are far from zero in  
483 the latter case. We explore this bias further below in relation to the uniformity of  
484 frequencies at sites as measured by Shannon entropy.

The Shannon entropy, as defined in our context

$$H(\boldsymbol{\pi}) = - \sum_{j=1}^{20} \pi_j \ln(\pi_j),$$

485 is a common measure of the degree of uniformity of the amino acid frequencies at sites. We  
486 note that when there is misspecification of the exchangeabilities, there is a bias towards  
487 frequency classes with high entropy. In more than 80% of the 27000 data sets where we  
488 fitted a model with the F-class, this was the class with the highest entropy. When the  
489 F-class had the highest entropy, it was assigned, on average, more than half the total  
490 weight (average weight = 0.59). Moreover, when fitting without the F-class, we noted that,  
491 in general, the class with the highest entropy is assigned a really large weight. For  
492 example, for all 5400 simulations (9 trees, 3 sequence lengths, 2 UDM distributions, and  
493 100 repetitions per condition) when fitting the classes in either C20, C30, CK36, or C60,  
494 the class with the highest entropy had the largest weight, and on average, that weight was  
495 4.71 times more than the weight assigned to that class when there is no misspecification.  
496 When fitting the classes in C40, the class with the second-highest entropy is the one with  
497 the largest weight, and it was also, on average, 3.99 times more than its weight when there

498 is no misspecification. Therefore we observe that this bias may not directly related to  
499 entropy but may instead be some other factor that is correlated with entropy.

500 We also explored the effects of misspecification of the exchangeabilities for data  
501 generated using the LG matrix but fitted the POISSON matrix. In this case, the tree  
502 estimation accuracy is also affected (for eg, OA values are in the range of 0.64 to 0.75 for  
503 sequence length of 1000 when it is correctly specified vs. a range of 0.60 to 0.71 when  
504 misspecified). Figure S3 in the Supplementary Material shows the average OA (A) and PM  
505 (B) for this case. In contrast to the reverse misspecification scenario (e.g. Figure 4), the  
506 F-class does not seem to hinder tree estimation. Figure S4 in the Supplementary Material,  
507 the analog of Figure 5, shows how, in this case, the weight of the F-class is close to zero,  
508 and therefore there is no likelihood difference with or without it. Furthermore, we did not  
509 find this phenomenon to be as strongly correlated to the Shannon entropy, as in the  
510 previous case. Although we found a shift in the correlation of entropy and class weight for  
511 all models. The average correlation between entropy and class weight for C60, C40, C30,  
512 and CK36 is 0.21 when there is no misspecification and -0.56 when there is. For C20 the  
513 shift is in a different direction, i.e the correlation between entropy and class weight is -0.12  
514 when there is no misspecification and 0.11 when there is. Nonetheless, for this model and  
515 when there is misspecification, the class with the highest entropy is the one with the  
516 second lowest weight (average weight = 0.008). This same class has the highest weight  
517 when there is no misspecification (average weight = 0.142). This also suggests that entropy  
518 is somehow related to or affected by, model misspecification.

519 Finally, we note that it is formally possible that the generally good performance of  
520 CXX models in the foregoing analyses could be related to a tendency of these models to  
521 prefer topologies where long branches are apart (i.e. they could have a long-branch  
522 repulsion (LBR) bias). To test this, we simulated from a topology with long branches  
523 together, obtained from the tree in Figure 1 after swapping the clade composed of taxa  
524  $c_1, \dots, c_m, d_1, \dots, d_m$  together with the edge leading to it and the clade composed of taxa



Fig. 3. (A) The plot of the OA values per model of the data generated under the UDM distributions and POISSON exchangeabilities. Different classes are fitted but in all cases, we fit POISSON exchangeabilities. The  $x$ -axis represents the number of taxa on the tree. The plot is divided by the sequence length. The lower bound for the 95% confidence interval (CI) of the model with the highest OA is depicted with an arrow whose color represents such model. An arrow pointing to the left represents CXX+F, an arrow pointing to the right represents CXX without F. (B) A similar plot to that on top but for PM. For this case, no confidence interval can be computed.



ARE PROFILE MIXTURE MODELS OVER-PARAMETERIZED?

25

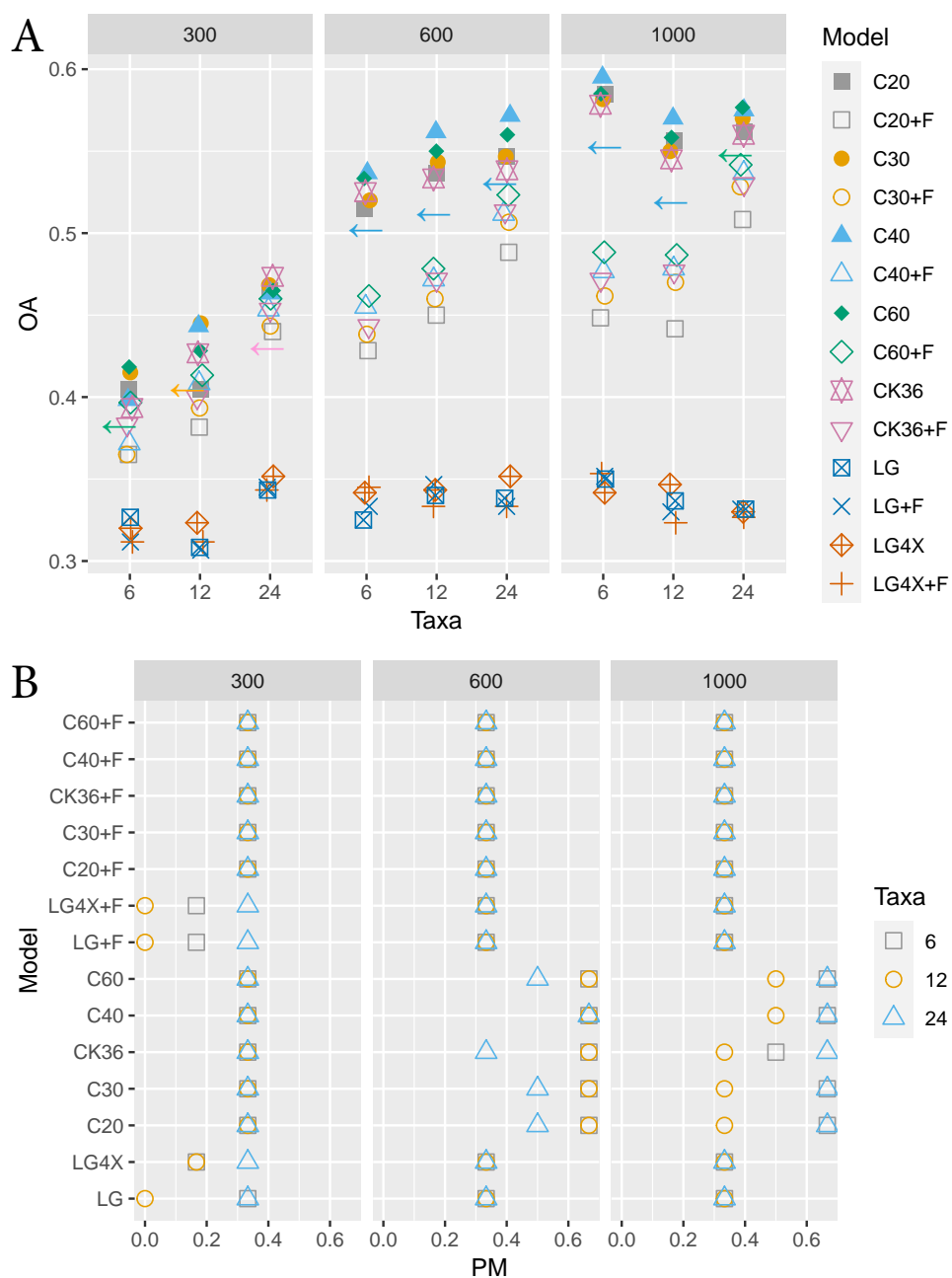


Fig. 4. (A) The plot of the OA values per model of data generated under the UDM distributions and POISSON exchangeabilities. Different classes are fitted but in all cases, we fit LG exchangeabilities. The  $x$ -axis represents the number of taxa on the tree. The plot is divided by the sequence length. The lower bound for the 95% confidence interval of the best model per scenario is depicted with a gray line. (B) A similar plot to that on top but for PM. For this case, no confidence interval can be computed.

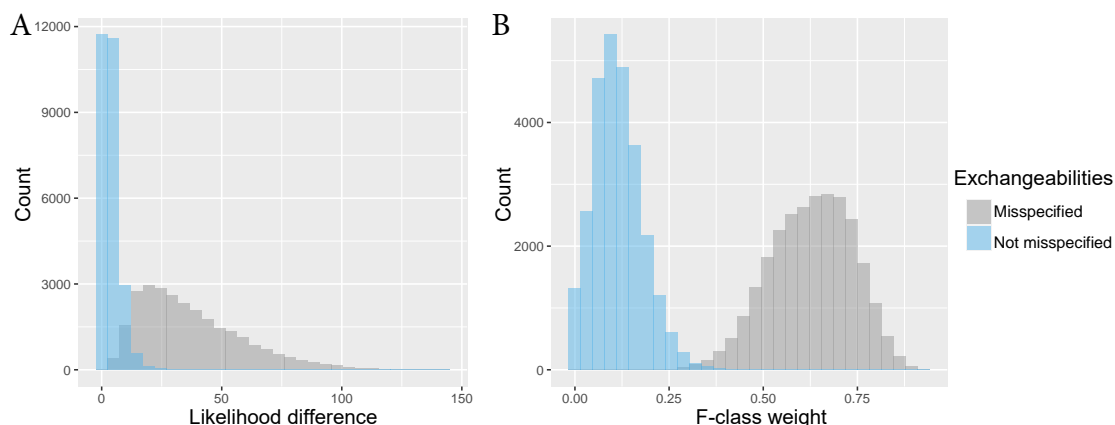


Fig. 5. (A) Histograms showing the difference between likelihoods values of models fitted with and without the F-class. No misspecification of the exchangeabilities is depicted in blue and in gray when there is. (B) Histograms showing the inferred F-class weight when there is no misspecification of the exchangeabilities (blue) and when there is (gray). The histograms consist of data generated under all trees, number of taxa, both UDM mixing distributions, and POISSON exchangeabilities. Misspecification of exchangeabilities refers to fitting using LG matrix instead of the POISSON matrix.

525  $e_1, \dots, e_m$  together with the edge leading to it. Consequently, the long branches group  
526 together to the exclusion of short branches. Figure S5 in the Supplementary Material, gives  
527 the results for the 12 taxon case for this simulating scenario. On average, the estimates in  
528 OA and PM for the CXX models agree with those found when exploring LBA. We note  
529 that models with fewer classes tend to have better performance in these cases, suggesting  
530 these show a slight LBR bias. In contrast, the LG and LG4X models are strongly affected  
531 by LBR; i.e., they have notably poor performance under the LBA conditions and  
532 extremely good performance under the LBR conditions. A bias towards either the LBR or  
533 LBA topologies is not desirable in general. In this sense the CXX models, especially  
534 C30-C60, show little bias and are clearly better choices under these simulation conditions.

535

### *Real Data*

536

To investigate the impact of mixture model choice on real data, we also analyzed  
537 three empirical data sets. These data sets are concatenated supermatrices:

538

I) a 133-protein dataset (24,291 sites  $\times$  40 taxa) assembled to assess the phylogenetic

539 position of the microsporidia in the tree of eukaryotes (Brinkmann et al. (2005)). The  
540 microsporidia are specifically related to Fungi but are sometimes recovered as  
541 branching outside of all eukaryotes because of an LBA artefact in which they are  
542 attracted to the outgroup archaeal sequences. We consider two trees: the correct tree  
543 recovered with the LG+C20+F+G model ( $T_I^{C20}$ ) (Susko et al. (2018)) and the LBA  
544 tree recovered with the LG+F+G model ( $T_I^{LG}$ )

545 II) a dataset of 146 proteins (35,371 sites  $\times$  37 taxa) assembled to assess the  
546 phylogenetic position of the nematodes in the animal tree of life (Lartillot et al.  
547 (2007)). In this case the two competing topologies are the correct topology (recovered  
548 with LG+C20+F+G:  $T_{II}^{C20}$ ) where nematodes branch as sister to arthropods (i.e. the  
549 Ecdysozoa group) versus the artefactual topology recovered with LG+F+G ( $T_{II}^{LG}$ ).

550 III) a dataset of 146 proteins (35,371 sites  $\times$  32 taxa) assembled to assess the  
551 phylogenetic position of the platyhelminths in the animal tree of life (Lartillot et al.  
552 (2007)). The correct position of platyhelminths within the Protostomia is reflected in  
553 the tree recovered by CAT+GTR ( $T_{III}^{CAT}$ ) instead of the artefactual Coelomata  
554 topology ( $T_{III}^{LG}$ ) recovered by LG+F+G and many mixture models (see Lartillot et al.  
555 (2007), Susko et al. (2018) and Wang et al. (2017))

556 For each of these trees, we computed the log-likelihoods when fitting classes in C20,  
557 C40, and C60, with and without the F-class, and with both LG and POISSON matrices.  
558 These likelihoods are shown in Table 4. In all cases, fitting with the LG matrix produces  
559 higher likelihood values.

560 For data set I, we observe that the correct tree is obtained with the largest  
561 log-likelihood differences over the incorrect tree when C60 and C20 are fit with LG  
562 exchangeabilities. For this data set the F-class only sometimes negatively affects  
563 topological estimation, but never improves it. For data set II, POISSON exchangeabilities  
564 strongly favor the correct topology over the incorrect one relative to LG. Here the F-class  
565 makes little difference but again never increases support for the correct tree. For data set

Model	Fitted	$L(T_I^{C20})$	$D(T_I^{C20})$	$L(T_{II}^{C20})$	$D(T_{II}^{C20})$	$L(T_{III}^{CAT})$	$D(T_{III}^{CAT})$
C60	LG+F	-715744	14	-712614	24	-626777	-2
C60	LG	-716579	20	-713137	27	-627241	-1
C60	POI+F	-722917	13	-718761	60	-631550	29
C60	POI	-722917	13	-718761	60	-631550	29
C40	LG+F	-716584	4	-713345	26	-627478	-6
C40	LG	-717555	9	-713882	30	-627977	-5
C40	POI+F	-724391	9	-720761	58	-633321	26
C40	POI	-724391	9	-720761	58	-633321	26
C20	LG+F	-718315	7	-714772	19	-628597	-10
C20	LG	-719775	19	-715659	23	-629393	-8
C20	POI+F	-727735	9	-723988	56	-635979	17
C20	POI	-727737	9	-723988	56	-635979	17

Table 4. The log-likelihoods of the trees estimated from the empirical data sets, where  $D(T_J)$  denotes the log-likelihood of the ‘correct tree’ (e.g. C20 or CAT superscripts) minus the ‘incorrect’ tree (e.g. LG superscripts). POI stands for POISSON matrix of exchangeabilities.

566 III, POISSON exchangeabilities favor the correct tree over the incorrect tree, with C60  
 567 showing the biggest log-likelihood difference. LG exchangeabilities seem to always favor the  
 568 incorrect tree.

569 Overall, the F-class never increases support for the correct tree and sometimes  
 570 decreases it. Whether LG improves estimation versus POISSON depends on the data set.  
 571 However, we note that the proteins and taxa in datasets II and III heavily overlap so the  
 572 outcomes of these analyses are not technically independent.

573 DISCUSSION

574 By extending earlier results (Kiefer and Wolfowitz (1956); Lindsay (1983);  
 575 Yourdkhani et al. (2021)), we confirm that, for profile mixture models, the tree and mixing  
 576 parameters of profile mixture models are statistically consistent even with a large number  
 577 of classes. However, since good performance is not guaranteed for short alignments, we  
 578 conducted an extensive simulation study of the performance and properties of profile  
 579 mixture models with smaller data sets with the goal of determining if  
 580 over-parameterization was a problem. We also investigated the effects of model

581 misspecification through both misspecification of the frequency classes and the  
582 exchangeabilities. Finally, the effects of the F-class was also investigated in all possible  
583 settings. These analyses provide useful theoretical and practical insights regarding model  
584 fit. Our main findings are the following:

585 (A) *Over-parameterization is not a problem for complex models:* For all  
586 alignment sizes explored here, we saw no evidence (in terms of MISE, MD, OA and  
587 PM) that use of more complex profile mixture models led to more variable or poorer  
588 estimation. This is true even for models having several classes with zero weights  
589 estimates. Consistent with the theoretical results for large numbers of sites,  
590 over-parameterization of these mixture models does not appear to be a problem for  
591 shorter alignments.

592 Since it is the mixture structure that is important in assessing whether models are  
593 overparameterized, large sample results likely extend to rates-across-sites mixtures  
594 (Yang (1994); Felsenstein and Churchill (1996); Mayrose et al. (2005); Susko et al.  
595 (2003)) and the types of mixtures used to infer selection pressure (Yang et al.  
596 (2000)). We also speculate that some of the small sample results found here may  
597 extend to those settings too but additional work is needed.

598 (B) *Misspecification of the frequency vectors does not necessarily imply bad*  
599 *fit:* Misspecification of the frequency vectors in profile mixture models does not  
600 cause problems if the estimated CDF can adequately approximate the true CDF.  
601 The more data available, the more classes are likely needed to closely approximate  
602 the true CDF.

603 (C) *Simple models behave poorly:* Likely as a consequence of (B), both the  
604 site-homogeneous POISSON and LG models, and the site-heterogeneous LG4X model  
605 perform very poorly in all scenarios. We believe this is because these have one (LG  
606 and POISSON) or very few (LG4X) classes. Although inference using simple models

607 can be much faster we do not recommend their use given their poor performance (i.e.  
608 susceptibility to LBA) under realistic site-heterogeneous simulation conditions.

609 (D) ***Misspecification of exchangeabilities and the presence of an F-class can***  
610 ***severely affect tree estimation:***

611 A severe decrease in accuracy of tree estimation is observed for data generated under  
612 the POISSON matrix and fitted using the LG matrix. In this scenario, it is clear that  
613 the F-class degrades performance. Such misspecification resulted in large weights of  
614 the F-class. For data generated under the LG matrix and fitted using the POISSON  
615 matrix, performance is affected but not as severely as in the previous case.

616 From this, we hypothesize that misspecification of exchangeabilities is more  
617 problematic when the matrix used to fit is less uniform than the true exchangeability  
618 matrix. In that case, the F-class tends to be accorded a large weight that leads to  
619 poorer tree estimation performance. Although the reverse misspecification scenario  
620 also degrades performance somewhat, the F-class has a little role in that case. We  
621 suspect that because the LG matrix was originally estimated as an “approximation”  
622 of a GTR matrix for many alignments in a site-homogeneous context, the LG matrix  
623 is less uniform than it would be if it was estimated in the presence of profile mixture  
624 models like the CXX set. Thus, we suspect the pathological behavior of the F-class  
625 and poor performance may apply to real estimation settings. We note that use of  
626 both the LG matrix and the F-class lead to higher likelihoods, so model selection  
627 criteria like AIC will frequently favor their use in real settings. Since the F-class  
628 never appeared to improve estimation in any of the simulations or real data analysis  
629 settings we examined, we discourage its use in site-profile mixture models.

630 (E) ***Better likelihood estimates do not imply better tree estimates:*** As a  
631 consequence of (D), and also observed in the data, better estimates in likelihood do  
632 not imply better tree estimates. This is also weakly observed even when there is no

633 misspecification of the exchangeabilities.

634 (F) *Adding more taxa can improve or hurt tree estimation accuracy*: Adding  
635 more taxa generally improves MISE and tree estimation. Surprisingly, when the  
636 model is misspecified (e.g. using UDM frequencies and misspecification of  
637 exchangeabilities) adding taxa does not always improve estimation; in one case it  
638 actually decreases performance (Fig. 4).

639 The poor performance of the methods when exchangeabilities are misspecified  
640 provides a strong motivation to develop software tools that allow ML estimation of a GTR  
641 matrix over all sites in the presence of a profile mixture model. In future work, we plan to  
642 construct mixing distributions that closely approximate the true CDF for data, hoping this  
643 would lead to more accurate tree estimation than current models.

644 To finalize, we give some practical recommendations for single gene phylogeny  
645 inference. First, we do not discourage the use of ‘rich models’ (those with many frequency  
646 classes), even when several classes have zero weight estimates. We suggest avoiding models  
647 with one or very few frequency classes. We also discourage the use of the F-class, unless  
648 both scenarios, with and without the F-class, can be explored.

#### 649 FUNDING

650 This work and H.B. were supported by the Moore-Simons Project on the Origin of  
651 the Eukaryotic Cell, Simons Foundation grant 735923LPI (DOI:  
652 <https://doi.org/10.46714/735923LPI>) and by NSERC Discovery Grants awarded to  
653 A.J.R. and E.S.

#### 654 REFERENCES

655 Anderson, F. E. and A. R. Lindgren. 2021. Phylogenomic analyses recover a clade of  
656 large-bodied decapodiform cephalopods. *Molecular Phylogenetics and Evolution*

- 657 156:107038.
- 658 Billera, L. J., S. P. Holmes, and K. Vogtmann. 2001. Geometry of the space of  
659 phylogenetic trees. *Advances in Applied Mathematics* 27:733–767.
- 660 Brinkmann, H., M. van der Giezen, Y. Zhou, G. P. de Raucourt, and H. Philippe. 2005. An  
661 Empirical Assessment of Long-Branch Attraction Artefacts in Deep Eukaryotic  
662 Phylogenomics. *Systematic Biology* 54:743–757.
- 663 Brown, M. W., S. C. Sharpe, J. D. Silberman, A. A. Heiss, B. F. Lang, A. G. B. Simpson,  
664 and A. J. Roger. 2013. Phylogenomics demonstrates that breviate flagellates are related  
665 to opisthokonts and apusomonads. *Proc R Soc B* 280.
- 666 Daubin, V. 2002. A phylogenomic approach to bacterial phylogeny: Evidence of a core of  
667 genes sharing a common history. *Genome Research* 12:1080–1090.
- 668 Dufayard, J.-F., L. Duret, S. Penel, M. Gouy, F. Rechenmann, and G. Perrière. 2005. Tree  
669 pattern matching in phylogenetic trees: automatic search for orthologs or paralogs in  
670 homologous gene sequence databases. *Bioinformatics* 21:2596–2603.
- 671 Felsenstein, J. 1978. Cases in which Parsimony or Compatibility Methods will be  
672 Positively Misleading. *Systematic Biology* 27:401–410.
- 673 Felsenstein, J. and G. A. Churchill. 1996. A Hidden Markov Model approach to variation  
674 among sites in rate of evolution. *Molecular Biology and Evolution* 13:93–104.
- 675 Feuda, R., M. Dohrmann, W. Pett, N. Lartillot, G. Wörheide, and D. Pisani. 2017.  
676 Improved Modeling of Compositional Heterogeneity Supports Sponges as Sister to All  
677 Other Animals. *Current Biology* 27:3864–3870.
- 678 Franzosa, E. A. and Y. Xia. 2009. Structural Determinants of Protein Evolution Are  
679 Context-Sensitive at the Residue Level. *Molecular Biology and Evolution* 26:2387–2395.
- 680 Gaston, D., E. Susko, and A. J. Roger. 2011. A phylogenetic mixture model for the  
681 identification of functionally divergent protein residues. *Bioinformatics* 27:2655–2663.



- 682 Goldman, N., J. L. Thorne, and D. T. Jones. 1998. Assessing the impact of secondary  
683 structure and solvent accessibility on protein evolution. *Genetics* 149:445–458.
- 684 Goldstein, R. A. 2008. The structure of protein evolution and the evolution of protein  
685 structure. *Current Opinion in Structural Biology* 18:170–177.
- 686 Halpern, A. L. and W. J. Bruno. 1998. Evolutionary distances for protein-coding sequences:  
687 modeling site-specific residue frequencies. *Molecular Biology and Evolution* 15:910–917.
- 688 Hartigan, J. A. and M. A. Wong. 1979. Algorithm as 136: A k-means clustering algorithm.  
689 *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 28:100–108.
- 690 Jones, D. T., W. R. Taylor, and J. M. Thornton. 1992. The rapid generation of mutation  
691 data matrices from protein sequences. *Bioinformatics* 8:275–282.
- 692 Kiefer, J. and J. Wolfowitz. 1956. Consistency of the maximum likelihood estimator in the  
693 presence of infinitely many incidental parameters. *The Annals of Mathematical*  
694 *Statistics* 27:887–906.
- 695 Lanfear, R., P. B. Frandsen, A. M. Wright, T. Senfeld, and B. Calcott. 2016.  
696 *PartitionFinder 2: New Methods for Selecting Partitioned Models of Evolution for*  
697 *Molecular and Morphological Phylogenetic Analyses. Molecular Biology and Evolution*  
698 *34:772–773.*
- 699 Lartillot, N., H. Brinkmann, and H. Philippe. 2007. Suppression of long-branch attraction  
700 artefacts in the animal phylogeny using a site-heterogeneous model. *BMC Evol Biol*  
701 *7:S4.*
- 702 Lartillot, N. and H. Philippe. 2004. A Bayesian Mixture Model for Across-Site  
703 Heterogeneities in the Amino-Acid Replacement Process. *Molecular Biology and*  
704 *Evolution* 21:1095–1109.
- 705 Lartillot, N., N. Rodrigue, D. Stubbs, and J. Richer. 2013. *PhyloBayes MPI: Phylogenetic*

- 706 Reconstruction with Infinite Mixtures of Profiles in a Parallel Environment. *Systematic*  
707 *Biology* 62:611–615.
- 708 Le, S. Q., C. C. Dang, and O. Gascuel. 2012. Modeling Protein Evolution with Several  
709 Amino Acid Replacement Matrices Depending on Site Rates. *Molecular Biology and*  
710 *Evolution* 29:2921–2936.
- 711 Le, S. Q. and O. Gascuel. 2008a. An Improved General Amino Acid Replacement Matrix.  
712 *Molecular Biology and Evolution* 25:1307–1320.
- 713 Le, S. Q. and O. Gascuel. 2008b. An Improved General Amino Acid Replacement Matrix.  
714 *Molecular Biology and Evolution* 25:1307–1320.
- 715 Le, S. Q. and O. Gascuel. 2010. Accounting for Solvent Accessibility and Secondary  
716 Structure in Protein Phylogenetics Is Clearly Beneficial. *Systematic Biology* 59:277–287.
- 717 Li, Y., X.-X. Shen, B. Evans, C. W. Dunn, and A. Rokas. 2021. Rooting the Animal Tree  
718 of Life. *Molecular Biology and Evolution* Msab170.
- 719 Lindsay, B. G. 1983. The Geometry of Mixture Likelihoods: A General Theory. *The*  
720 *Annals of Statistics* 11:86 – 94.
- 721 Ly-Trong, N., S. Naser-Khdour, R. Lanfear, and B. Q. Minh. 2021. Alisim: A fast and  
722 versatile phylogenetic sequence simulator for the genomic era. *bioRxiv* .
- 723 Massey, F. J. 1951. The Kolmogorov-Smirnov test for goodness of fit. *Journal of the*  
724 *American Statistical Association* 46:68–78.
- 725 Mayrose, I., N. Friedman, and T. Pupko. 2005. A Gamma mixture model better accounts  
726 for among site rate heterogeneity. *Bioinformatics* 21:ii151–ii158.
- 727 Minh, B. Q., H. A. Schmidt, O. Chernomor, D. Schrempf, M. D. Woodhams, A. von  
728 Haeseler, and R. Lanfear. 2020. IQ-TREE 2: New Models and Efficient Methods for  
729 Phylogenetic Inference in the Genomic Era. *Molecular Biology and Evolution*  
730 37:1530–1534.

- 731 Pál, C., B. Papp, and M. J. Lercher. 2006. An integrated view of protein evolution. *Nature*  
732 *Reviews Genetics* 7:337–348.
- 733 Philippe, H., H. Brinkmann, D. V. Lavrov, D. T. J. Littlewood, M. Manuel, G. Wörheide,  
734 and D. Baurain. 2011. Resolving difficult phylogenetic questions: Why more sequences  
735 are not enough. *PLOS Biology* 9:1–10.
- 736 Philippe, H. and J. Laurent. 1998. How good are deep phylogenetic trees? *Current Opinion*  
737 *in Genetics and Development* 8:616–623.
- 738 Pisani, D., W. Pett, M. Dohrmann, R. Feuda, O. Rota-Stabelli, H. Philippe, N. Lartillot,  
739 and G. Wörheide. 2015. Genomic data do not support comb jellies as the sister group to  
740 all other animals. *Proceedings of the National Academy of Sciences* 112:15402–15407.
- 741 Pupko, T., D. Huchon, Y. Cao, N. Okada, and M. Hasegawa. 2002. Combining Multiple  
742 Data Sets in a Likelihood Analysis: Which Models are the Best? *Molecular Biology and*  
743 *Evolution* 19:2294–2307.
- 744 Raymann, K., C. Brochier-Armanet, and S. Gribaldo. 2015. The two-domain tree of life is  
745 linked to a new root for the archaea. *Proceedings of the National Academy of Sciences*  
746 112:6670–6675.
- 747 Sander, C. and R. Schneider. 1994. The HSSP database of protein structure-sequence  
748 alignments. *Nucleic Acids Research* 22:3597–3599.
- 749 Schrempf, D., N. Lartillot, and G. Szöllösi. 2020. Scalable Empirical Mixture Models That  
750 Account for Across-Site Compositional Heterogeneity. *Molecular Biology and Evolution*  
751 37:3616–3631.
- 752 Scott, D. W. 1992. *Multivariate density estimation theory, practice, and visualization*. J.  
753 Wiley.
- 754 Self, S. G. and K.-Y. Liang. 1987. Asymptotic properties of maximum likelihood

- 755 estimators and likelihood ratio tests under nonstandard conditions. *Journal of the*  
756 *American Statistical Association* 82:605–610.
- 757 Simion, P., H. Philippe, D. Baurain, N. King, G. Wörheide, and M. Manuel. 2017. A Large  
758 and Consistent Phylogenomic Dataset Supports Sponges as the Sister Group to All  
759 Other Animals. *Current Biology* 27:958–967.
- 760 Susko, E., C. Field, C. Blouin, and A. J. Roger. 2003. Estimation of Rates-Across-Sites  
761 Distributions in Phylogenetic Substitution Models. *Systematic Biology* 52:594–603.
- 762 Susko, E., L. Lincker, and A. J. Roger. 2018. Accelerated Estimation of Frequency Classes  
763 in Site-Heterogeneous Profile Mixture Models. *Molecular Biology and Evolution*  
764 35:1266–1283.
- 765 Thorndike, R. L. 1953. Who belongs in the family? *Psychometrika* 18:267–276.
- 766 Wang, H.-C., K. Li, E. Susko, and A. J. Roger. 2008. A class frequency mixture model  
767 that adjusts for site-specific amino acid frequencies and improves inference of protein  
768 phylogeny. *BMC Evol Biol* 8:331.
- 769 Wang, H.-C., B. Q. Minh, E. Susko, and A. J. Roger. 2017. Modeling Site Heterogeneity  
770 with Posterior Mean Site Frequency Profiles Accelerates Accurate Phylogenomic  
771 Estimation. *Systematic Biology* 67:216–235.
- 772 Wang, H.-C., E. Susko, and A. J. Roger. 2014. An Amino Acid Substitution-Selection  
773 Model Adjusts Residue Fitness to Improve Phylogenetic Estimation. *Molecular Biology*  
774 *and Evolution* 31:779–792.
- 775 Whelan, S. and N. Goldman. 2001. A General Empirical Model of Protein Evolution  
776 Derived from Multiple Protein Families Using a Maximum-Likelihood Approach.  
777 *Molecular Biology and Evolution* 18:691–699.
- 778 Wickett, N. J., S. Mirarab, N. Nguyen, T. Warnow, E. Carpenter, N. Matasci,  
779 S. Ayyampalayam, M. S. Barker, J. G. Burleigh, M. A. Gitzendanner, B. R. Ruhfel,

- 780 E. Wafula, J. P. Der, S. W. Graham, S. Mathews, M. Melkonian, D. E. Soltis, P. S.  
781 Soltis, N. W. Miles, C. J. Rothfels, L. Pokorny, A. J. Shaw, L. DeGironimo, D. W.  
782 Stevenson, B. Surek, J. C. Villarreal, B. Roure, H. Philippe, C. W. dePamphilis,  
783 T. Chen, M. K. Deyholos, R. S. Baucom, T. M. Kutchan, M. M. Augustin, J. Wang,  
784 Y. Zhang, Z. Tian, Z. Yan, X. Wu, X. Sun, G. K.-S. Wong, and J. Leebens-Mack. 2014.  
785 Phylotranscriptomic analysis of the origin and early diversification of land plants.  
786 Proceedings of the National Academy of Sciences 111:E4859–E4868.
- 787 Williams, T. A., P. G. Foster, C. J. Cox, and T. M. Embley. 2013. An archaeal origin of  
788 eukaryotes supports only two primary domains of life. Nature 504:231–236.
- 789 Yang, Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with  
790 variable rates over sites: Approximate methods. Journal of Molecular Evolution 39:306 –  
791 314.
- 792 Yang, Z. 1996. Maximum-Likelihood Models for Combined Analyses of Multiple Sequence  
793 Data. J Mol Evol 42:587–596.
- 794 Yang, Z., R. Nielsen, N. Goldman, and A.-M. K. Pedersen. 2000. Codon-Substitution  
795 Models for Heterogeneous Selection Pressure at Amino Acid Sites. Genetics 155:431–449.  
796
- 797 Yourdkhani, S., E. S. Allman, and J. A. Rhodes. 2021. Parameter Identifiability for a  
798 Profile Mixture Model of Protein Evolution. Journal of Computational Biology  
799 28:570–586.

800 APPENDIX

801 *Statistical consistency of the MLE*

802 As mentioned in the section entitled “*Mixture models and over-parameterization*  
803 *with large samples*”, the results of Kiefer and Wolfowitz (1956) imply that the tree and

804 mixing parameters are consistent. In this section we elaborate on why this is true in the  
805 context of site-profile mixture models. To do this, we first introduce the concept of a *forest*.  
806 This allows us to then extend the parameter space of the model (to a compact one) so the  
807 regularity conditions in Kiefer and Wolfowitz (1956) are satisfied.

808 A phylogenetic *forest*  $F$  on  $X$  is a collection of phylogenetic trees  $F_q$  on  $X_q$ , known  
809 as *components*, where  $X = \cup X_q$  and  $X_q \cap X_p = \emptyset$  for any two components.

810 We extend the parameter space  $\Theta$  defined in Section “*Mixture models and*  
811 *over-parameterization with large samples*” by substituting point (i) of that section with:

812 (i') A metric forest  $F$  on  $N$  taxa obtained after removing a, possibly empty, set of edges  
813 from a rooted metric binary tree on  $N$  taxa and retaining only components that  
814 display taxa.

In this case, the substitution process of a profile mixture model is as follows: for each site,  
a class  $\pi_c$  is sampled with probability  $w_c$ , and a rate parameter  $r_k$  is sampled with  
probability  $d_k$ . On each component  $F_q$  of  $F$  an independent substitution process on  $F_q$   
with exchangeabilities  $R$ , root distribution  $\pi_c$ , and a rate parameter  $r_k$  is conducted. For a  
given site pattern  $\mathbf{x}_i$ , the likelihood function is determined by a weighted average of the  
product of partial site likelihoods conditional on each site-profile class and site-rate class  
per component:

$$L(\theta|\mathbf{x}_i) = \sum_{c=1}^C w_c \sum_{k=1}^K d_k \prod_{q=1}^Q P(\mathbf{x}_i|F_q, R, \pi_c, r_k),$$

815 where  $\theta = (F, R, \{\pi_c\}, \{w_c\}, \{r_k\}, \{d_k\}) \in \Theta$  and  $Q$  are the number of components  
816 of  $F$ . The reasoning behind forests is to account for infinite edge lengths, which are  
817 represented by the edges missing from the tree defining  $F$ . This not only allows us to  
818 consider this limiting case but also, it weakly depicts the effects of functional divergence  
819 (Gaston et al. (2011)). Note that when no edges are removed from the tree in (i'), the  
820 resulting forest has one component and the likelihood is the same as the one defined in the  
821 section “*Mixture models and over-parameterization with large samples.*”

822 We now show we “compactified” the parameter space  $\Theta$  of the profile mixture. We  
823 show this by proving (i’) above, and elements (iii) and (iv) of parameter space  $\Theta$  are  
824 compact. We do not consider (ii) (i.e. the matrix of exchangeabilities) in the parameter  
825 space because, for the models we consider here, it is fixed beforehand.

826 Clearly (iii) is compact since both, the root distribution vectors and the class  
827 weights are closed and bounded. To argue (i’) is compact, we need to recall that the space  
828 of rooted metric phylogenetic trees on  $n$  taxa can be viewed as a collection of  $(2n - 3)!!$   
829 open cubes corresponding to all different tree topologies (Billera et al. (2001)). The  
830 limiting cases in these cubes correspond to infinite edge lengths on the trees. We can  
831 ensure boundedness by re-parameterizing edge lengths via the logistic function  $p = \frac{e^t}{1+e^t}$ .  
832 Cases with an edge length  $p = 1$  corresponds to those edge lengths being infinite. The  
833 limiting likelihoods in those cases correspond to forests. Therefore (i’) can be viewed as a  
834 compact space.

835 For (iv), the rate weights are clearly compact (closed and bounded). Now, even if  
836 the rate parameters are unbounded from above, the limiting case, i.e. when  $r \rightarrow \infty$ , is  
837 equivalent to the process occurring in the forest where all components are just single taxa.  
838 Therefore  $\Theta$  is compact, and the results of Kiefer and Wolfowitz (1956) hold in our context.

839 Another important statement mentioned in the section “*Mixture models and*  
840 *over-parameterization with large samples*” is: even when the parameter estimation is  
841 unrestricted and any mixing distribution is allowed, the maximum likelihood estimator will  
842 be a finite mixing distribution. This is an implication from Theorem 3.1 in Lindsay (1983).  
843 For this result to hold, the trace of the likelihood curve over the mixing parameters must  
844 be compact. This follows immediately from the fact that: (1) the mixing parameter space  
845 is compact; and (2) the image of continuous functions, such as the trace of the likelihood  
846 curve, is compact whenever the domain is compact.

847

### *Identifiability*

848

849

850

In this section we argue why, as mentioned in the section “*Mixture models and over-parameterization with large samples,*” for many of the cases considered here the tree parameter is identifiable.

851

852

853

854

855

856

857

858

859

In Theorem 5.7 in Yourdkhani et al. (2021) it is shown that for profile mixture models with  $C \cdot K < 72$ , where  $C$  is the number of classes and  $K$  the number of rates, and more than 8 taxa, the tree and numerical parameters are generically identifiable, up to arbitrary re-scaling of the tree and the exchangeability matrix. Generically identifiable means identifiable except maybe in a set of measure zero; informally speaking this means identifiable except maybe in a tiny subset of parameters relative to the full parameter space. Although in such work there is no description of the generic setting of the parameter space, we argue that with just a small perturbation of the parameters one can always guarantee the result to hold.

860

861

862

863

864

865

866

Since in all models considered here fixed parameters are obtained empirically (that is these have no structure for eg. be solutions of a phylogenetic invariant) and the parametric function is continuous, there exists  $\epsilon$  such that a translation of the numerical parameters by  $\epsilon$  will make these generic. This is true for the frequency vectors, weights, rate parameters, edge lengths, and the exchangeabilities matrices. While the POISSON matrix is not generated from data, the proof of Theorem 5.7 in Yourdkhani et al. (2021) is built on this matrix up to a constant and therefore identifiability also holds in this case.