

## Complex Models of Sequence Evolution Require Accurate Estimators as Exemplified with the Invariable Site Plus Gamma Model

LAM-TUNG NGUYEN<sup>1</sup>, ARNDT VON HAESLER<sup>1,2</sup>, AND BUI QUANG MINH<sup>1,\*</sup>

<sup>1</sup>Center for Integrative Bioinformatics Vienna, Department of Structural and Computational Biology, Max F. Perutz Laboratories, University of Vienna, Medical University Vienna, Campus Vienna Biocenter 5, A-1030, Vienna, Austria; and <sup>2</sup>Bioinformatics and Computational Biology, Faculty of Computer Science, University of Vienna, Waehringer Strasse 29, A-1090 Vienna, Austria

\*Center for Integrative Bioinformatics Vienna (CIBIV), Max F. Perutz Laboratories, Campus Vienna Biocenter 5 (VBC5), A-1030 Vienna, Austria; E-mail: [minh.bui@univie.ac.at](mailto:minh.bui@univie.ac.at)

Received 8 March 2017; reviews returned 22 November 2017; accepted 23 November 2017

Associate Editor: Edward Susko

**Abstract.**—The invariable site plus  $\Gamma$  model (I+ $\Gamma$ ) is widely used to model rate heterogeneity among alignment sites in maximum likelihood and Bayesian phylogenetic analyses. The proof that the I+ continuous  $\Gamma$  model is identifiable (model parameters can be inferred correctly given enough data) has increased the credibility of its application to phylogeny reconstruction. However, most phylogenetic software implement the I+ discrete  $\Gamma$  model, whose identifiability is likely but unproven. How well the parameters of the I+ discrete  $\Gamma$  model are estimated is still disputed. Especially the correlation between the fraction of invariable sites and the fractions of sites with a slow evolutionary rate is discussed as being problematic. We show that optimization heuristics as implemented in frequently used phylogenetic software (PhyML, RAxML, IQ-TREE, and MrBayes) cannot always reliably estimate the shape parameter, the proportion of invariable sites, and the tree length. Here, we propose an improved optimization heuristic that accurately estimates the three parameters. While research efforts mainly focus on tree search methods, our results signify the equal importance of verifying and developing effective estimation methods for complex models of sequence evolution. [Gamma model; invariable sites; maximum likelihood; phylogenetic inference; rate heterogeneity among sites.]

In model based phylogenetic analysis, the invariable site plus  $\Gamma$  model (Yang 1994; Gu et al. 1995), hereafter referred to as I+ $\Gamma$ , is widely used to model rate heterogeneity among sites, because it often fits the data better than the  $\Gamma$  model or the invariable-sites model alone (Sullivan and Swofford 1997). Thus, the I+ $\Gamma$  model is frequently selected by MODELTEST (Posada and Crandall 1998). The I+ $\Gamma$  model has two parameters: the proportion of invariable sites  $p_{\text{inv}}$  ( $0 \leq p_{\text{inv}} < 1$ ) and the shape parameter  $\alpha$  ( $> 0$ ) of the  $\Gamma$  distribution. A small  $\alpha$  ( $< 1$ ) indicates strong rate heterogeneity, whereas a large  $\alpha$  ( $> 1$ ) corresponds to weak rate heterogeneity. Under certain conditions  $p_{\text{inv}}$  and  $\alpha$  compete with each other for the same phylogenetic signal. For example,  $\alpha \leq 1$  already accounts for sites with low rates; that interferes with  $p_{\text{inv}}$  and causes a correlation between the parameters making reliable estimation of those parameters difficult (Sullivan et al. 1999; Mayrose et al. 2005). Despite this interference, it has been shown that the I+ continuous  $\Gamma$  model is identifiable for “all but members of the F81 family of rate matrices on any phylogeny with more than two distinct interspecies distances” (Rogers 2001; Allman and Rhodes 2008; Chai and Housworth 2011). Since the I+ continuous  $\Gamma$  model is identifiable, reliable parameter estimation for this model should be possible for sufficiently long multiple sequence alignments.

However, most phylogenetic software only implement the I+ discrete  $\Gamma$  (Yang 1994) model as an approximation to the continuous  $\Gamma$  model because of its computational efficiency. The discussed competition between  $p_{\text{inv}}$  and

$\alpha$  is based on the analysis of the discrete  $\Gamma$ -distribution. The results have led to the suggestion to discourage the use of the I+ discrete  $\Gamma$  model (Yang 2006; Jia et al. 2014; Stamatakis 2014).

On the other hand, the identifiability of the I+ discrete  $\Gamma$  model is likely, but unproven (Chai and Housworth 2011), and it is unclear how accurately popular phylogenetic software estimate parameters of the I+ discrete  $\Gamma$  model.

Thus, we used simulations to assess the accuracy of the I+ discrete  $\Gamma$  estimators implemented in three maximum likelihood (ML) phylogenetic software: RAxML (Stamatakis 2014), PhyML (Guindon et al. 2010), IQ-TREE (Nguyen et al. 2015), and one Bayesian inference program MrBayes (Ronquist et al. 2012). More precisely, we simulated 100,000-bp long alignments along three balanced trees of 6, 24, and 96 taxa. The lengths of the alignments ensure the recovery of the correct tree topology. The three trees have uniform branch lengths of 0.1 substitutions per site except for one internal branch on the 6-taxon tree whose length equals 0.2 to allow for three distinct distances between the sequences as required for identifiability in the continuous case (Chai and Housworth 2011). We assumed the K2P model (Kimura 1980) with a transition/transversion ratio of 2.0 and the rate heterogeneity model I+ discrete  $\Gamma$  with four rate categories. For each tree and each pair  $(p_{\text{inv}}, \alpha) \in \{0.0, 0.1, \dots, 0.9\} \times \{0.1, 0.5, 1.0\}$ , we simulated 100 alignments using Seq-Gen (Rambaut and Grassly 1997). We used RAxML version 8.2.2, PhyML

version 20141029, IQ-TREE version 1.3.7, and MrBayes version 3.2.6 compiled with the BEAGLE library (Ayres et al. 2012) to infer the invariable proportion, the shape parameter, and the tree length from the simulated alignments. For RAxML, PhyML, and IQ-TREE, we used the default options.

For MrBayes we used the default priors, that is, uniform distribution within interval [0,1] for  $p_{\text{inv}}$ , exponential distribution with mean 1.0 for  $\alpha$ , nonclocklike uniform Dirichlet distribution for branch lengths and  $\Gamma$  distribution with mean of 10 for tree lengths (Unconstrained:GammaDir(1.0,0.1,1.0,1.0)). The sequential version of MrBayes was run with four chains (one hot and three cold chains) and one million MCMC generations. One thousand four hundred and eighty-nine (16.5%) nonconvergent MrBayes runs, where the effective sample sizes (ESS) on  $p_{\text{inv}}$ ,  $\alpha$ , or tree lengths are smaller than 100, were repeated with five million generations. However, 52 of the extended reruns were stopped after 4 weeks without completing all five million generations. We note that 207 of the 1489 reruns still did not converge. MrBayes estimates are then summarized as the mean of the posterior distribution with a default burn-in of 25%.

#### CURRENT PHYLOGENETIC PROGRAMS DO NOT PRODUCE ACCURATE ESTIMATES FOR THE I+ DISCRETE $\Gamma$ MODEL

Figure 1 displays the averages  $\bar{\alpha}$  of the estimated shape parameter  $\alpha$ , the averages  $\bar{p}_{\text{inv}}$  of the estimated invariable fraction  $p_{\text{inv}}$  and the average  $\bar{l}$  of the estimated tree length  $l$  produced by PhyML, RAxML, IQ-TREE, and MrBayes for the 100 alignments simulated from each parameter combinations. A program is called *accurate* if the estimated averages  $\bar{\alpha}$ ,  $\bar{p}_{\text{inv}}$ ,  $\bar{l}$  deviate no more than 10% from the true values.

None of the tested programs estimated all parameter combinations accurately. The problem is especially pronounced for the 6-taxon alignments. For extreme rate heterogeneity ( $\alpha=0.1$ ) MrBayes and PhyML recovered the true  $\alpha$ ,  $p_{\text{inv}}$ , and  $l$  for 9/10 and 5/10 parameter combinations respectively, whereas the average estimates from IQ-TREE and RAxML were inaccurate. For strong rate heterogeneity ( $\alpha=0.5$ ), the degrees of inaccuracy observed among all programs differ unsystematically. On the one hand, IQ-TREE and MrBayes accurately estimated the parameters in four and six settings. On the other hand, RAxML and PhyML could not estimate accurately the three parameters for any of the ten parameter-combinations. For medium rate variation ( $\alpha=1.0$ ), only IQ-TREE produced the accurate estimates for all settings. All other programs exhibited varying degrees of inaccuracy.

For the 24- and 96-taxon alignments we observed an increase in the number of accurate estimates for all programs. These results corroborate a previous study (Sullivan et al. 1999) showing that increased taxon sampling leads to more reliable estimates. However, under extreme rate heterogeneity ( $\alpha=0.1$ ),

only MrBayes estimated all parameter sets accurately. We note that our measure of accuracy correlates well with the Bayesian coverage probabilities, the frequency with which true parameter values are included in the 95% credible interval of the estimates (Supplementary Fig. S1 available on Dryad at <https://doi.org/10.5061/dryad.4j5c7>). Two hundred and seven (2.3%) nonconvergent MrBayes runs (effective sample size of  $\alpha$  or  $p_{\text{inv}}$  are smaller than 100) partly overlap with cases where MrBayes was not accurate for 6-taxon simulations ( $\alpha=0.1$  and  $p_{\text{inv}}=0.9$ ;  $\alpha=0.5$  and  $p_{\text{inv}}\leq 0.5$ ;  $\alpha=1.0$  and  $0.2\leq p_{\text{inv}}\leq 0.5$ ). Hence, nonconvergence is a predictor of difficult settings but does not fully explain the inaccuracy of MrBayes (Fig. 1).

We also observed that inaccurate estimates of  $\alpha$  and  $p_{\text{inv}}$  could sometimes lead to tree lengths that substantially deviate from the simulated lengths. For instance, for the 96-taxon alignments simulated with  $\alpha=0.1$  and  $p_{\text{inv}}=0.8$  (expected tree length = 18.9) IQ-TREE estimated an average tree length of 177.0 that is nine times longer than the simulated tree length. The other programs also sometimes produced tree lengths that were considerably longer than the simulated ones.

In terms of computing times PhyML, RAxML, and IQ-TREE needed for all analyses 62,441, 12,563, and 7,675 CPU hours, respectively. MrBayes needed 740,681 CPU hours to complete one million MCMC generations, thus it is 96.5 times slower than the fastest ML program. We note that this is only a lower-bound for the effective time one needs to wait for MrBayes results because 16.5% of MrBayes runs did not converge after one million MCMC generations. These runs were repeated with five times more generations, that led to significantly more computations.

#### MULTIPLE LOCAL OPTIMA ON THE LIKELIHOOD SURFACE CAUSE INACCURACY

Because the tested programs performed quite differently with respect to the accuracy of parameter estimation, the number of taxa cannot be the only explanation. We suspected that the optimization heuristics as implemented in these programs drive the accuracy. Examining the likelihood surfaces for many simulated alignments revealed a common feature that the parameter space has two distinct peaks of high log-likelihoods (Fig. 2): one global close to the true parameters and one suboptimal peak with slightly lower log-likelihood ( $\Delta\text{LnL}=-27$  in this example) separated by a flat valley from the true parameters. In this particular instance, MrBayes and PhyML found the true parameters whereas RAxML and IQ-TREE were trapped in the local maximum (not necessarily the case for other instances). In fact, whether the global or local optimum is detected depends on the starting values of the numerical optimization routines.

To summarize, we compared for each simulated alignment the log-likelihoods of the estimates with the log-likelihoods obtained for the true parameters.

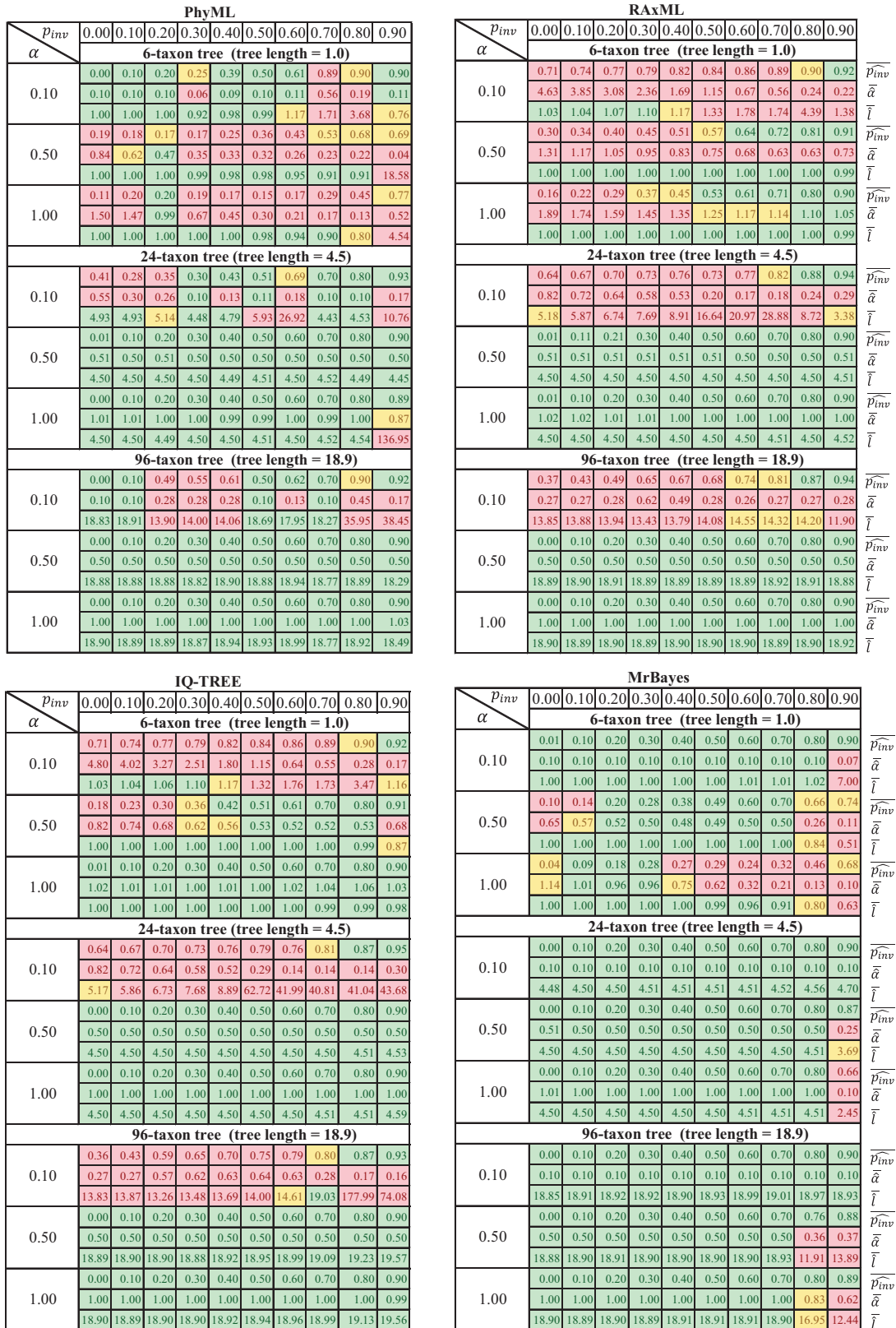


FIGURE 1. The averages  $\overline{\alpha}$  of the estimated shape parameter  $\alpha$ , the averages  $\overline{p_{inv}}$  of the estimated invariable fraction  $p_{inv}$  and the average  $\overline{l}$  of the estimated tree length  $l$  produced by PhyML, RAxML, IQ-TREE, and MrBayes for the 100 alignments simulated from each parameter combinations. The averages are highlighted according to their differences from the true values: inaccurate (more than 25% deviation, red in online version), moderately inaccurate (10% to 25% deviation yellow in online version), and accurate (less than 10% deviation green in online version). For  $p_{inv} = 0.0$  the estimated  $\overline{p_{inv}}$  is accurate if  $0 \leq \overline{p_{inv}} \leq 0.01$ , moderately inaccurate if  $0.01 < \overline{p_{inv}} \leq 0.05$ , and inaccurate if  $0.05 < \overline{p_{inv}}$ .

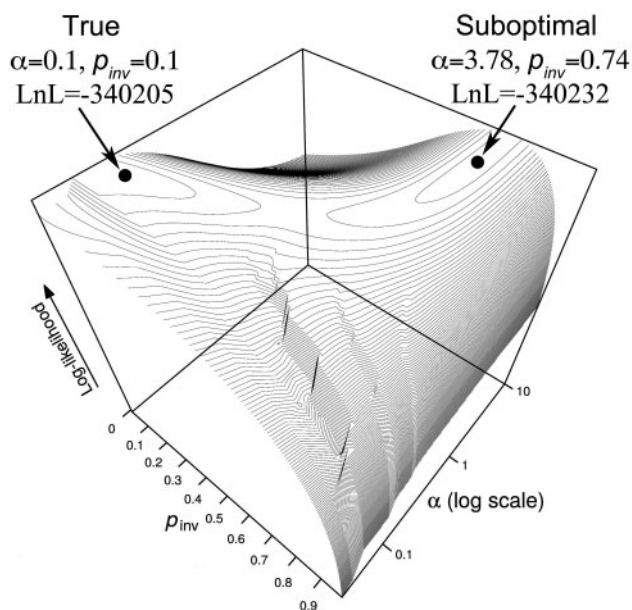


FIGURE 2. The likelihood surface for one simulated alignment as a function of  $\alpha$  and  $p_{\text{inv}}$ .

TABLE 1. Percentage of alignments where the true simulation parameters result in higher log-likelihoods than the inferred parameters from the programs for three simulation scenarios (6-, 24-, and 96-taxon trees)

Program	6-taxon tree (%)	24-taxon tree (%)	96-taxon tree (%)
MrBayes	34.7	6.6	5.6
PhyML	60.2	21.9	45.0
RAxML	89.2	37.7	49.5
IQ-TREE	36.0	34.1	44.0

Table 1 show how often the true parameter combination produced a higher likelihood than the inferred parameters from MrBayes, PhyML, RAxML, and IQ-TREE. These fractions are particularly high for the 6-taxon tree and for the ML inference programs. Most ML phylogenetic programs use general-purpose numerical methods to find  $\alpha$  and  $p_{\text{inv}}$  (e.g., Brent 1973). These methods are obviously not well adapted to the complex likelihood surface (Fig. 2) and explain the poor overall performance of the ML programs (Fig. 1).

#### EFFECTIVE OPTIMIZATION HEURISTIC PRODUCES ACCURATE ESTIMATES

As remedy, we propose an alternative optimization heuristic which employs the Expectation-Maximization (EM) algorithm (Dempster et al. 1977) to estimate  $p_{\text{inv}}$ . We assume a discrete  $\Gamma$  distribution with  $k$  rate categories. Under the I+ discrete  $\Gamma$  model, the site rates follow a discrete mixture model consisting of  $k+1$  categories with rates  $r_0, \dots, r_k$ , where  $r_0=0$  represents invariable sites and  $r_i > 0$  ( $i=1, \dots, k$ ) are the  $k$  rates determined from the shape parameter  $\alpha$  of the discrete  $\Gamma$  distribution (Yang 1994). Given a tree topology, the

optimization heuristic does the following:

1. Choose initial values for  $\alpha$  and  $p_{\text{inv}}$ .
2. Optimize branch lengths by the Newton–Raphson method.
3. Optimize substitution model parameters by the Broyden–Fletcher–Goldfarb–Shanno algorithm.
4. For each alignment site  $D_i$  compute its posterior probability of being invariable ( $1 \leq i \leq n$ , where  $n$  is the number of alignment sites):

$$P(r_0|D_i) = \frac{P(D_i|r_0)w_0}{\sum_{j=0}^k P(D_i|r_j)w_j}$$

where  $P(D_i|r_j)$  is the likelihood of site  $D_i$  having rate  $r_j$  and  $w_0 = p_{\text{inv}}$ ,  $w_j = \frac{1-p_{\text{inv}}}{k}$  ( $1 \leq j \leq k$ ).

5. Update  $p_{\text{inv}} = \frac{1}{n} \sum_{i=1}^n P(r_0|D_i)$ .
6. Optimize  $\alpha$  by the Brent method.
7. If the log-likelihood improvement is greater than a predefined  $\epsilon$  value, go back to Step 2. Otherwise, stop the parameter optimization.

Steps 4 and 5 correspond to the E- and M-step of the EM algorithm, respectively. To avoid being stuck in local optima, we repeat this optimization procedure from ten starting values of  $p_{\text{inv}}$  evenly spaced between 0 and the fraction of constant sites observed in the alignment. The initial value of  $\alpha$  is always set to 1.0.

We implemented the new optimization heuristic in IQ-TREE now called IQ-TREE-EM (IQ-TREE version 1.4.3) and repeated the previous simulations. Figure 3 shows that IQ-TREE-EM successfully recovered the true parameters for all but one parameter combination (6-taxon,  $\alpha=0.5$  and  $p_{\text{inv}}=0.0$ ) where the average estimates ( $\hat{\alpha}=0.59$  and  $\hat{p}_{\text{inv}}=0.06$ ) slightly deviated from the true values.

Also, the percentage of instances where the estimated log-likelihoods were lower than the log-likelihood for the true parameters dropped considerably (0.06% 6-taxon tree, 0.0% 24-taxon tree, and 0.03% 96-taxon tree; compare also with Table 1).

This increase in accuracy comes at the cost of an increased total computing time by a factor of 1.3 compared to IQ-TREE.

Thus, we conclude that the inaccurate parameter estimation of the I+ discrete  $\Gamma$  shown for the tested phylogenetic programs is caused by ineffective optimization methods.

#### IMPACT ON REAL DATA

To investigate the impact of accuracy on real data for ML estimates, we analyzed 70 DNA and 45 protein

IQ-TREE-EM											
$\alpha \backslash p_{inv}$	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	
<b>6-taxon tree (tree length = 1.0)</b>											
0.10	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	$\bar{p}_{inv}$
	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	$\bar{\alpha}$
	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.01	0.93
0.50	0.06	0.10	0.19	0.27	0.38	0.48	0.59	0.70	0.80	0.90	$\bar{p}_{inv}$
	0.59	0.52	0.50	0.48	0.49	0.48	0.48	0.52	0.52	0.54	$\bar{\alpha}$
	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.99	0.98	$\bar{l}$
1.00	0.01	0.09	0.19	0.29	0.40	0.50	0.60	0.70	0.80	0.90	$\bar{p}_{inv}$
	1.05	1.00	0.98	0.98	1.02	1.03	1.03	1.04	1.00	1.02	$\bar{\alpha}$
	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.99	1.00	0.99	$\bar{l}$
<b>24-taxon tree (tree length = 4.5)</b>											
0.10	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	$\bar{p}_{inv}$
	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	$\bar{\alpha}$
	4.50	4.49	4.49	4.50	4.51	4.51	4.51	4.52	4.52	4.19	$\bar{l}$
0.50	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	$\bar{p}_{inv}$
	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	$\bar{\alpha}$
	4.50	4.50	4.50	4.50	4.50	4.50	4.50	4.50	4.50	4.50	$\bar{l}$
1.00	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	$\bar{p}_{inv}$
	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	$\bar{\alpha}$
	4.50	4.50	4.50	4.50	4.50	4.50	4.50	4.51	4.50	4.50	$\bar{l}$
<b>96-taxon tree (tree length = 18.9)</b>											
0.10	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	$\bar{p}_{inv}$
	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	$\bar{\alpha}$
	18.90	18.89	18.90	18.91	18.89	18.92	18.97	18.98	19.02	18.80	$\bar{l}$
0.50	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	$\bar{p}_{inv}$
	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	$\bar{\alpha}$
	18.90	18.90	18.91	18.89	18.90	18.90	18.90	18.92	18.91	18.86	$\bar{l}$
1.00	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	$\bar{p}_{inv}$
	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	$\bar{\alpha}$
	18.90	18.89	18.90	18.89	18.90	18.90	18.90	18.89	18.90	18.87	$\bar{l}$

FIGURE 3. The averages  $\bar{\alpha}$  of the estimated shape parameter  $\alpha$ , the averages  $\bar{p}_{inv}$  of the estimated invariable fraction  $p_{inv}$  and the average  $\bar{l}$  of the estimated tree length  $l$  produced by IQ-TREE-EM for the 100 alignments simulated from each parameter combinations. The highlighting is explained in Fig. 1.

TreeBase alignments (Nguyen et al. 2015). We applied the GTR+I+ $\Gamma$ 4 and LG+I+ $\Gamma$ 4 models for DNA and protein data, respectively. Among 115 alignments, we detected 15 (5 DNA and 10 protein) alignments where the estimated  $\alpha$  and  $p_{inv}$  by PhyML, RAxML, or IQTREE deviated more than 10% from those by IQ-TREE-EM (Fig. 4; Supplementary Table S1 available on Dryad). The estimates by PhyML and IQ-TREE deviated from those by IQ-TREE-EM only for one and two alignments, respectively. However, RAxML estimated  $\alpha$  and  $p_{inv}$  dramatically different from IQ-TREE-EM, PhyML, and IQ-TREE for all 15 alignments. Interestingly, RAxML systematically overestimated  $\alpha$  and  $p_{inv}$  for all 5 DNA and underestimated them for all 10 protein alignments ( $p_{inv}$  sometimes very close to zero).

## DISCUSSION

Our simulations revealed a major issue for parameter estimation of the I+ discrete  $\Gamma$  model as implemented in phylogenetic software. Despite using very long alignments, none of the tested programs recovered the true  $\alpha$ ,  $p_{inv}$ , and tree length for all parameter combinations. Often, the estimates deviated heavily from the true values and different programs estimated different values for the same evolutionary parameters, although all programs inferred the true tree. Our further analysis of 115 TreeBase alignments showed that PhyML, IQ-TREE, and IQ-TREE-EM estimates generally agree with each other except for two alignments. However, we identified 15 (13%) alignments where RAxML systematically overestimated  $\alpha$  and  $p_{inv}$  for

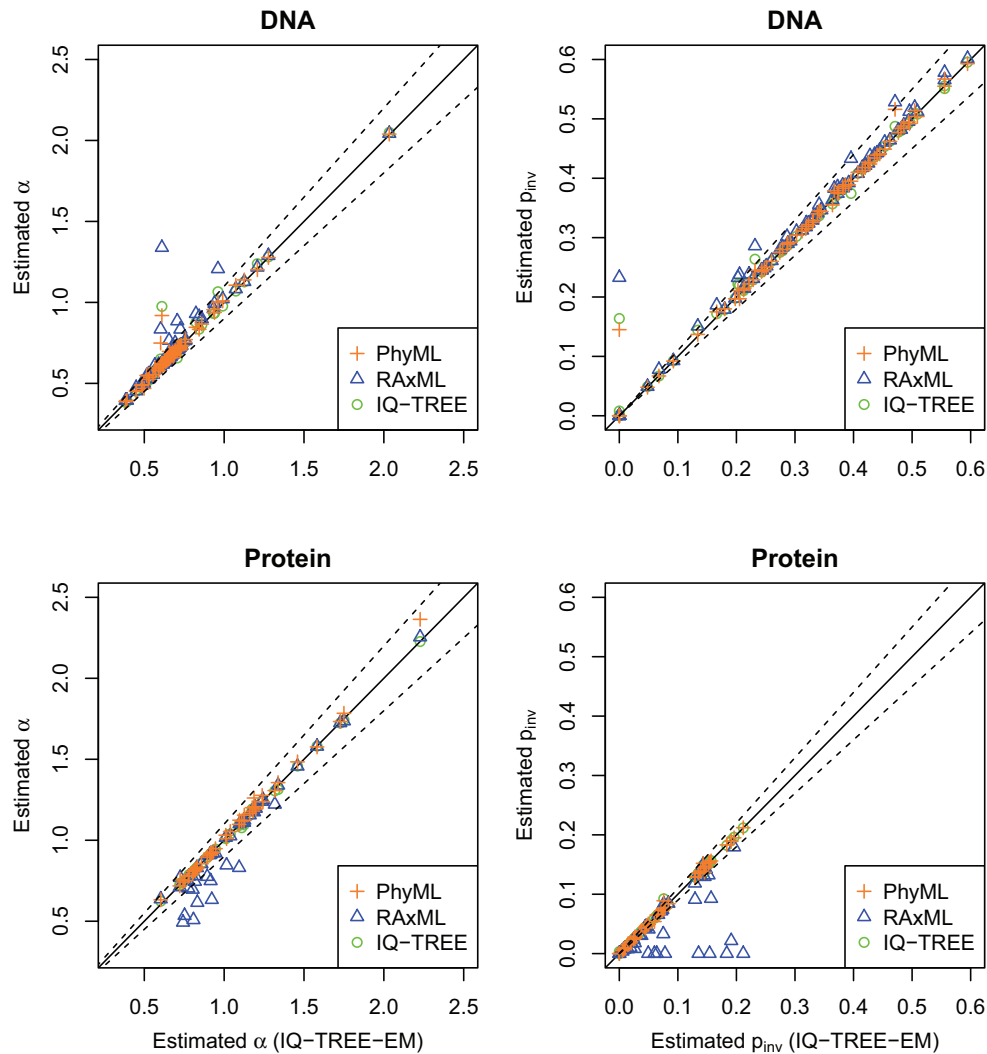


FIGURE 4. Estimation of  $\alpha$  (left) and  $p_{\text{inv}}$  (right) for TreeBase alignments using IQ-TREE-EM ( $x$ -axis) and IQ-TREE (circle), PhyML (cross) and RAxML (triangle). Dashed lines show the boundaries of 10% deviation from the IQ-TREE-EM estimates. Points above the upper dashed lines indicate overestimation compared with IQ-TREE-EM, whereas points under the lower dashed lines indicate underestimation.

DNA and underestimated for protein, compared with other programs. The reasons for that behavior are unclear and deserve further analyses. While this result may not be extrapolated to other data sets, phylogenetic software should benefit from the more robust optimization described for IQ-TREE-EM.

We showed that the estimation heuristics implemented in popular phylogenetic programs causes such inaccurate estimates and the  $I+\Gamma$  model *per se* is not problematic. The relatively good performance of MrBayes is likely attributed to the Bayesian sampling of the parameter space but comes at the cost of excessive computing time.

With IQ-TREE-EM, we provided an alternative optimization heuristic for ML methods that allows accurate estimation of the parameters for the  $I+$  discrete  $\Gamma$  model. IQ-TREE-EM combines two optimization techniques: the multiple starting point strategy and the EM algorithm. We note that the EM algorithm alone

will not achieve this accuracy (Supplementary Fig. S2 available on Dryad). Therefore, while the former allows to escape local optima, the latter helps to speed-up the optimization using analytical formula for  $p_{\text{inv}}$ . This new approach effectively infers the true evolutionary parameters for long alignments. Thus, it is tempting to speculate that the GTR+ $I$ + discrete  $\Gamma$  model is also identifiable as shown for the GTR+ $I$ + continuous  $\Gamma$  model (Chai and Housworth 2011).

Our observations show that as models of sequence evolution become more and more complex (e.g., Dirichlet rate and other mixture models), tailored numerical optimization methods are necessary to achieve accurate estimates of evolutionary parameters. It is not enough to recover the true tree, if one wants to understand how evolutionary forces shaped contemporary genomes. The effect of wrong parameter estimates for the substitution model on the total tree length is sometimes dramatic (see Fig. 1). This may

in turn bias downstream analysis such as divergence time dating, inference of site-specific evolutionary rates, and ancestral sequence reconstruction, which are sensitive to the parameter estimates. Thus, one should critically scrutinize the heuristics implemented in popular programs. A more thorough evaluation of phylogenetic inference programs allowing for very complicated models of sequence evolution is necessary, but beyond the scope of this article.

Finally, we would like to point out that we only addressed the accurate computation of  $p_{\text{inv}}$  and  $\alpha$  for the widely used I+ discrete  $\Gamma$  model. We do not discuss the biological interpretation of  $p_{\text{inv}}$ . The estimate of  $p_{\text{inv}}$  depends very much on the multiple sequence alignment at hand.  $p_{\text{inv}}$  may change if we enlarge the alignment. Thus, drawing an absolute conclusion from  $p_{\text{inv}}$  is in any case questionable.

#### SUPPLEMENTARY MATERIAL

Data available from the Dryad Digital Repository: <https://doi.org/10.5061/dryad.4j5c7>.

#### FUNDING

This work was supported by the Austrian Science Fund – FWF (Grant Nos I-2805-B29 and I-1824-B22).

#### ACKNOWLEDGEMENTS

The authors would like to thank Heiko A. Schmidt for fruitful discussions, two anonymous reviewers, Fredrik Ronquist, and Edward Susko for constructive and helpful comments on an earlier version of the manuscript. The computational results presented have been achieved using the Vienna Scientific Cluster 3 (VSC-3).

#### REFERENCES

Allman E.S, Rhodes J.A. 2008. Identifying evolutionary trees and substitution parameters for the general Markov model with invariable sites. *Math. Biosci.* 211:18–33.  
 Ayres D.L., Darling A., Zwickl D.J., Beerli P., Holder M.T., Lewis P.O., Huelsenbeck J.P., Ronquist F., Swofford D.L., Cummings

M.P., Rambaut A., Suchard M.A. 2012. BEAGLE: an application programming interface and high-performance computing library for statistical phylogenetics. *Syst. Biol.* 61:170–173.  
 Brent R.P. 1973. Algorithms for minimization without derivatives. Englewood Cliffs, New Jersey: Prentice-Hall. p. 1–195.  
 Chai J., Housworth E.A. 2011. On Rogers' proof of identifiability for the GTR+ Gamma+ I model. *Syst. Biol.* 60:713–718.  
 Dempster A.P., Laird N.M., Rubin D.B. 1977. Maximum likelihood from incomplete data via the EM algorithm. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 75:1–38.  
 Gu X., Fu Y.X., Li W.H. 1995. Maximum-likelihood-estimation of the heterogeneity of substitution rate among nucleotide sites. *Mol. Biol. Evol.* 12:546–557.  
 Guindon S., Dufayard J.F., Lefort V., Anisimova M., Hordijk W., Gascuel O. 2010. New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst. Biol.* 59:307–321.  
 Jia F.Z., Lo N., Ho S.Y.W. 2014. The impact of modelling rate heterogeneity among sites on phylogenetic estimates of intraspecific evolutionary rates and timescales. *PLoS One* 9:e95722.  
 Kimura M. 1980. A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide-sequences. *J. Mol. Evol.* 16:111–120.  
 Mayrose I., Friedman N., Pupko T. 2005. A Gamma mixture model better accounts for among site rate heterogeneity. *Bioinformatics* 21:ii151–ii158.  
 Nguyen L.T., Schmidt H.A., von Haeseler A., Minh B.Q. 2015. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* 32:268–274.  
 Posada D., Crandall K.A. 1998. MODELTEST: testing the model of DNA substitution. *Bioinformatics* 14:817–818.  
 Rambaut A., Grassly N.C. 1997. Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees. *Comp. Appl. Biosci.* 13:235–238.  
 Rogers J.S. 2001. Maximum likelihood estimation of phylogenetic trees is consistent when substitution rates vary according to the invariable sites plus gamma distribution. *Syst. Biol.* 50:713–722.  
 Ronquist F., Teslenko M., van der Mark P., Ayres D.L., Darling A., Höhna S., Larget B., Liu L., Suchard M.A., Huelsenbeck J.P. 2012. MrBayes 3.2: efficient Bayesian phylogenetic inference and model choice across a large model space. *Syst. Biol.* 61:539–542.  
 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30:1312–1313.  
 Sullivan J., Swofford D.L. 1997. Are guinea pigs rodents? The importance of adequate models in molecular phylogenetics. *J. Mamm. Evol.* 4:77–86.  
 Sullivan J., Swofford D.L., Naylor G.J.P. 1999. The effect of taxon sampling on estimating rate heterogeneity parameters of maximum-likelihood models. *Mol. Biol. Evol.* 16:1347–1356.  
 Yang Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. *J. Mol. Evol.* 39:306–314.  
 Yang Z. 2006. Computational molecular evolution. New York: Oxford University Press. p. 113–114.