

# Complex models of sequence evolution require accurate estimators as exemplified with the invariable site plus Gamma model

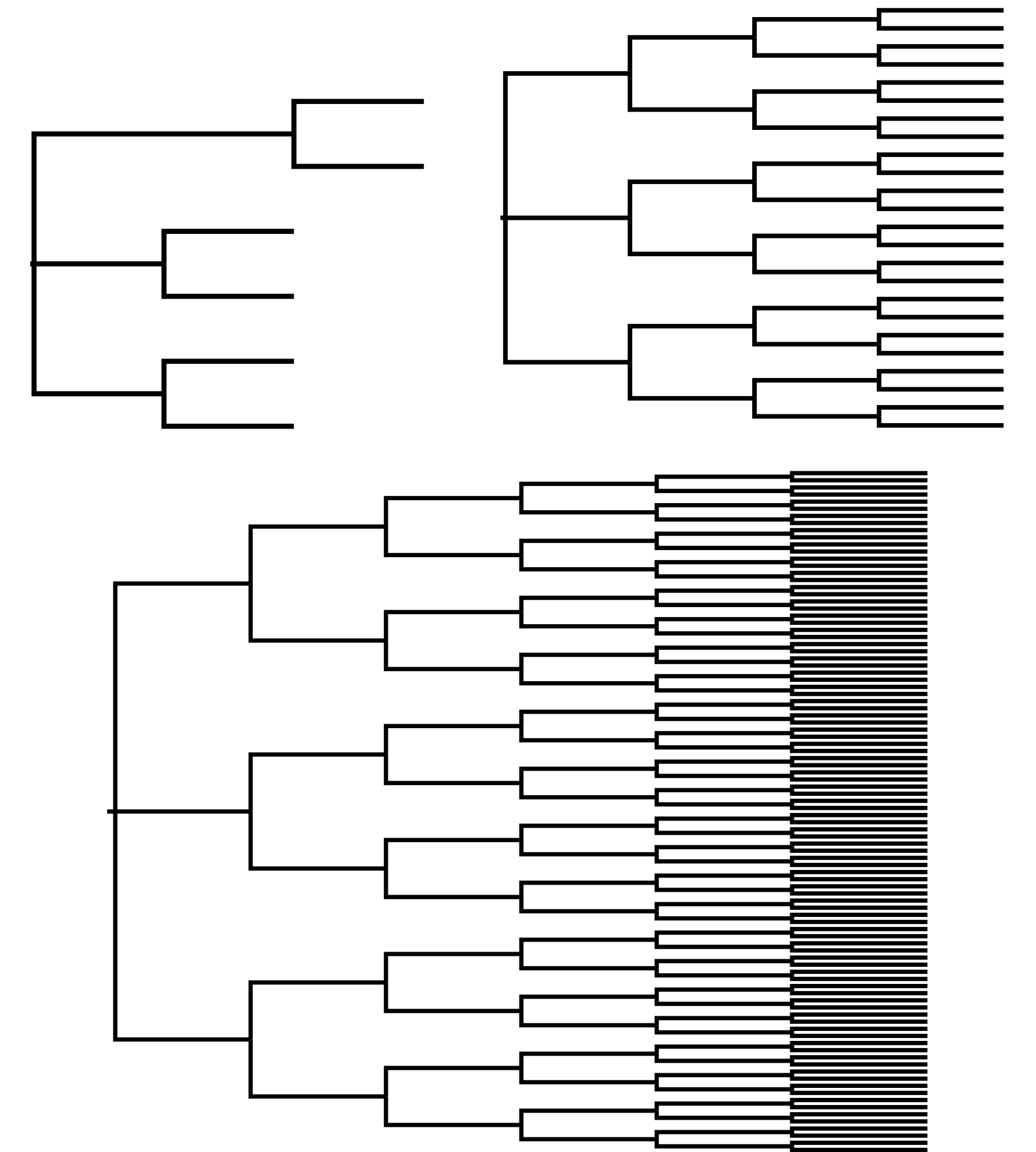
Lam-Tung Nguyen, Arndt von Haeseler and Bui Quang Minh

Center for Integrative Bioinformatics Vienna (CIBIV), Max F. Perutz Laboratories, Austria

## Introduction

- The I+ $\Gamma$  (Gu et al. 1995) model (invariable site and Gamma rate heterogeneity across sites) is widely-used in phylogenetics.
- I+ $\Gamma$  model is parameterized by two parameters:  $p_{inv}$  (proportion of invariable sites) and  $\alpha$  (shape of  $\Gamma$  distribution).
- Its use is controversial due to correlation between  $p_{inv}$  and  $\alpha$  for strong rate heterogeneity ( $\alpha < 1$ ) (Sullivan & Swofford 1999; Mayrose et al. 2005; Yang 2006):
  - The use of I+ $\Gamma$  is discouraged (Jia et al. 2014; Yang 2006).
  - However, the identifiability of I+*continuous*  $\Gamma$  model was proven (Rogers 2001; Allman & Rhodes 2008; Chai & Housworth 2011).
- We perform simulations to assess accuracy of I+*discrete*  $\Gamma$  estimators in maximum likelihood and Bayesian software.

## Simulations under K2P+I+ $\Gamma_4$ ( $\kappa = 4$ )



## Current phylogenetic software do not produce accurate estimates for I+ $\Gamma$ model!

PhyML										RAxML										MrBayes										IQ-TREE										
$\alpha$	$p_{inv}$									$\hat{p}_{inv}$	$\hat{\alpha}$	$\hat{i}$	$\hat{p}_{inv}$	$\hat{\alpha}$	$\hat{i}$	$\hat{p}_{inv}$	$\hat{\alpha}$	$\hat{i}$	$\hat{p}_{inv}$	$\hat{\alpha}$	$\hat{i}$	$\hat{p}_{inv}$	$\hat{\alpha}$	$\hat{i}$	$\hat{p}_{inv}$	$\hat{\alpha}$	$\hat{i}$													
6-taxon tree (tree length = 1.0)																																								
0.10	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.01	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.71	0.74	0.77	0.79	0.82	0.84	0.86	0.89	0.90	0.92
0.50	0.19	0.18	0.17	0.17	0.25	0.36	0.43	0.53	0.68	0.69	0.30	0.34	0.40	0.45	0.51	0.57	0.64	0.72	0.81	0.91	0.10	0.14	0.20	0.28	0.38	0.49	0.60	0.70	0.66	0.74	0.82	0.74	0.68	0.62	0.56	0.53	0.52	0.52	0.53	0.68
1.00	0.84	0.62	0.47	0.35	0.33	0.32	0.26	0.23	0.22	0.04	1.31	1.17	1.05	0.95	0.83	0.75	0.68	0.63	0.63	0.73	1.14	1.02	0.96	0.95	0.69	0.57	0.31	0.21	0.13	0.10	1.02	1.01	1.01	1.00	1.01	1.00	1.02	1.04	1.06	1.03
24-taxon tree (tree length = 4.5)																																								
0.10	0.41	0.28	0.35	0.30	0.43	0.51	0.69	0.70	0.80	0.93	0.64	0.67	0.70	0.73	0.76	0.73	0.77	0.82	0.88	0.94	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.64	0.67	0.70	0.73	0.76	0.79	0.76	0.81	0.87	0.95
0.50	0.01	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.01	0.11	0.21	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.87	5.18	5.87	6.74	7.69	8.91	16.64	20.97	28.88	8.72	3.38
1.00	1.01	1.01	1.00	1.00	0.99	0.99	1.00	0.99	1.00	0.87	1.02	1.02	1.01	1.01	1.00	1.00	1.00	1.00	1.00	1.00	4.50	4.50	4.50	4.50	4.50	4.50	4.50	4.50	4.51	7.44	4.50	4.50	4.50	4.50	4.50	4.50	4.50	4.51	4.51	4.59
96-taxon tree (tree length = 18.9)																																								
0.10	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.37	0.43	0.49	0.65	0.67	0.68	0.74	0.81	0.87	0.94	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.27	0.27	0.28	0.28	0.62	0.49	0.28	0.26	0.27	0.28
0.50	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	18.85	18.89	18.92	18.92	18.90	18.93	18.99	19.01	18.97	18.96
1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.03	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.89	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.99

## Cause, Solution and Recommendation

### CAUSE:

- Optimization heuristics in different software trapped in suboptimal estimates.

### SOLUTION:

- An expectation maximization (EM) algorithm to estimate  $p_{inv}$ .
- Restarting parameter estimation from 10 initial  $p_{inv}$  evenly spaced between 0 and fraction of constant sites.

### RECOMMENDATION:

- Other software developers to employ similar strategy.
- Critically scrutinize heuristics when implementing more complex models.

IQ-TREE-EM													
$\alpha$	$p_{inv}$									$\hat{p}_{inv}$	$\hat{\alpha}$	$\hat{i}$	
6-taxon tree (tree length = 1.0)													
0.10	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.00	0.10	0.10
0.50	0.06	0.10	0.19	0.27	0.38	0.48	0.59	0.70	0.80	0.90	0.06	0.10	0.10
1.00	1.05	1.00	0.98	0.98	1.02	1.03	1.03	1.04	1.00	1.02	1.00	1.00	1.00
24-taxon tree (tree length = 4.5)													
0.10	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.00	0.10	0.10
0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50
1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
96-taxon tree (tree length = 18.9)													
0.10	0.00	0.10	0.20	0.30	0.40	0.50	0.60	0.70	0.80	0.90	0.00	0.10	0.10
0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50
1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00

## References

- Allman ES, Rhodes JA (2008) *Math Biosci* 211:18-33.
- Chai J, Housworth EA (2011) *Syst Biol* 60:713-718.
- Gu X, Fu YX, Li WH (1995) *Mol Biol Evol* 12:546-557.
- Jia FZ, Lo N, Ho SYW (2014) *PLOS One* 9.
- Mayrose I, Friedman N, Pupko T (2005) *Bioinformatics* 21:ii151-ii158.
- Rogers JS (2001) *Syst Biol* 50:713-722.
- Sullivan J, Swofford DL, Naylor GJP (1999) *Mol Biol Evol* 16:1347-1356.
- Yang Z (2006) Oxford University Press.