

*Research article***Starless bias and parameter-estimation bias in the likelihood-based phylogenetic method****Xuhua Xia**<sup>1,2,\*</sup><sup>1</sup> Department of Biology, University of Ottawa, Ottawa, Canada, K1N 6N5<sup>2</sup> Ottawa Institute of Systems Biology, Ottawa, Canada, K1H 8M5**\* Correspondence:** Email: [Xuhua.Xia@uottawa.ca](mailto:Xuhua.Xia@uottawa.ca); Tel: +6135625800.

**Abstract:** I analyzed various site pattern combinations in a 4-OTU case to identify sources of starless bias and parameter-estimation bias in likelihood-based phylogenetic methods, and reported three significant contributions. First, the likelihood method is counterintuitive in that it may not generate a star tree with sequences that are equidistant from each other. This behaviour, dubbed starless bias, happens in a 4-OTU tree when there is an excess (i.e., more than expected from a star tree and a substitution model) of conflicting phylogenetic signals supporting the three resolved topologies equally. Special site pattern combinations leading to rejection of a star tree, when sequences are equidistant from each other, were identified. Second, fitting gamma distribution to model rate heterogeneity over sites is strongly confounded with tree topology, especially in conjunction with the starless bias. I present examples to show dramatic differences in the estimated shape parameter  $\alpha$  between a star tree and a resolved tree. There may be no rate heterogeneity over sites (with the estimated  $\alpha > 10000$ ) when a star tree is imposed, but  $\alpha < 1$  (suggesting strong rate heterogeneity over sites) when an (incorrect) resolved tree is imposed. Thus, the dependence of “rate heterogeneity” on tree topology implies that “rate heterogeneity” is not a sequence-specific feature, cautioning against interpreting a small  $\alpha$  to mean that some sites are under strong purifying selection and others not. Thirdly, because there is no existing (and working) likelihood method for evaluating a star tree with continuous gamma-distributed rate, I have implemented the method for JC69 in a self-contained R script for a four-OTU tree (star or resolved), in addition to another R script assuming a constant rate over sites. These R scripts should be useful for teaching and exploring likelihood methods in phylogenetics.

**Keywords:** maximum likelihood; molecular phylogenetics; rate heterogeneity; starless; star-tree paradox

---

## 1. Introduction

I explore two phylogenetic issues here. The first is the starless bias. If a set of aligned sequences are equidistant from each other, i.e., the number of various types of substitutions between any two sequences is exactly the same, then we intuitively would expect a star tree. Distance-based methods such as neighbor-joining [1] or FastME [2] will indeed give us a star tree whenever pairwise distances are all equal. The starless bias refers to the inability of a phylogenetic method to generate a star tree with equidistant sequences. It was first alluded to in a study of potential bias in maximum likelihood method involving missing data and rate heterogeneity over sites [3], but its occurrence is more general than that. I will illustrate this bias here with four sequences, identify the source of the bias, and discuss alternative approaches relevant to the problem.

The second issue, related to the first, is the confounding effect of tree topology on phylogenetic parameter estimation, in particular the shape parameter  $\alpha$  of gamma distribution used to model rate heterogeneity over sites. It may seem obvious that  $\alpha$  depend on topology, but the issue needs to be studied for two reasons. First, how such dependence occurs is not well dissected. Second, one frequently encounters interpretation of  $\alpha$  as if it is a sequence-specific feature, with a small alpha interpreted as indicating some sites strongly constrained by purifying selection and other sites not. It is therefore relevant to caution against such interpretation with real examples. For simplicity, I will work on 4-OTU trees only so that all site patterns can be conveniently considered.

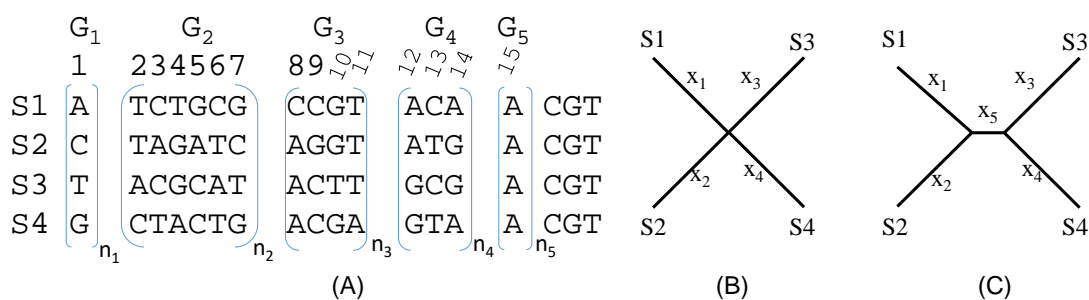
I provide two self-contained R script files New2.R and NewGamma3.R (s001.R and s002.R within the Supplementary genetics-05-04-212.zip) implementing the likelihood method with JC69 model [4] with four OTUs. New2.R (s001) assumes a constant rate over sites and estimates branch lengths for the star tree (with four branches) and a resolved tree (with five branches). NewGamma3.R (s002) does the same but with a continuous gamma-distributed rate over sites. They are plain text files that can be copied and pasted into an R window to generate results presented in the paper, and should also be useful for teaching and exploring likelihood methods in phylogenetics.

### 1.1. Site patterns classification and sequence representation

With four sequences, there are 256 possible site patterns. For sequences to evolve according to the JC69 substitution model [4], the 256 site patterns would become equally frequently after sequences experienced full substitution saturation. Different combinations of these 256 site patterns support different substitution models and different topologies. We will focus on the JC69 substitution model and classify these 256 possible site patterns into 15 classes (Figure 1A) so that we only need to calculate likelihood for these 15 site pattern classes. The transition probabilities needed for likelihood calculation for the JC69 model, together with other frequently used Markovian nucleotide substitution models such as F84 (used in DNAML since 1984) [5,6], {XE “substitution model: HKY85”} [7], TN93 {XE “substitution model: TN93”} [8], and GTR {XE “substitution model: GTR”} [9,10] have been numerically illustrated and re-derived with three different approaches [11,12]. These illustrations are sufficiently detailed for one to extend the JC69 model in the two R script files to other models.

The 15 site pattern classes (Figure 1A) can be further lumped into five groups ( $G_1$  to  $G_5$ , Figure 1A). Sites in  $G_1$  have all four OTUs (S1 to S4) with different nucleotides, with 24 unique site patterns represented as a single  $G_1$  site in Figure 1A (because JC69 sees these 24 site patterns as identical).

Sites in  $G_2$  each have three different nucleotides, with a total of 144 unique site patterns. Only six sites are used to represent them in Figure 1A because JC69 sees these 144 unique sites in this group to be identical to one of the six representative sites. Sites in  $G_3$  feature two nucleotides, with three OTUs having the same nucleotide, and a total of 48 unique sites represented by four sites in Figure 1A. Sites in  $G_4$  also feature two nucleotides, with two OTUs sharing one nucleotide and the other two OTUs sharing the other (i.e., they are the traditional informative sites in Fitch parsimony). There are 36 unique  $G_4$  sites represented by three sites in Figure 1A.  $G_5$  sites are monomorphic, with four unique site patterns represented by site 15 in Figure 1A.



**Figure 1.** Representative site patterns and alternative trees used. (A) A total of 256 possible site patterns with four sequences, classed into 15 site patterns relevant to the JC69 model and further boxed into five site pattern groups ( $G_1$  to  $G_5$ ). Sites within each group jointly do not support any of the three alternative resolved trees. (B) A star tree. (C) One of three resolved trees with  $x_5 > 0$ .

Given a JC69 model, the five groups of sites ( $G_1$  to  $G_5$  in Figure 1A) support the three unrooted topologies equally (i.e., they do not preferentially support any of the three). This is obvious for  $G_1$  and  $G_5$  sites. The six sites in  $G_2$ , shown in Figure 1A, jointly also support the three topologies equally, so do the four sites in  $G_3$  and three sites in  $G_4$ . Note that, with a star tree and JC69, we cannot have  $G_1$ ,  $G_2$  and  $G_4$  sites without having  $G_3$  sites first. With low sequence divergence, almost all site patterns from a star tree should be  $G_3$  sites.

Subscripts  $n_1$  to  $n_5$  in Figure 1A mean multiples of the enclosed sites, e.g., an  $n_2$  of 10 means that the six  $G_2$  sites in Figure 1A is repeated 10 times (for a total of 60 sites). A combination of ( $n_1$ ,  $n_2$ ,  $n_3$ ,  $n_4$ ,  $n_5$ ), where  $n_i$  corresponds to those in Figure 1A, means a set of aligned sequences containing  $n_1$   $G_1$  sites,  $n_2$   $G_2$  sites (i.e.,  $n_2 \cdot 6$  sites),  $n_3$   $G_3$  sites (for a total of  $n_3 \cdot 4$  sites),  $n_4$   $G_4$  sites (for a total of  $n_4 \cdot 3$  sites), and  $n_5$   $G_5$  sites (Figure 1A). A set of four sequences of length 256 containing all 256 possible site patterns is specified as (24, 24, 12, 12, 4). Such a set of sequences will naturally have equal nucleotide frequencies and twice as many transversions as transitions, i.e., the equilibrium ratio of substitution saturation. A set of sequences with any combinations of ( $n_1$ ,  $n_2$ ,  $n_3$ ,  $n_4$ ,  $n_5$ ) are equidistant from each other from a JC69 perspective, and we would desire to have a star tree (Figure 1B) instead of one of the three resolved trees (Figure 1C). Hereafter I specify a set of four aligned sequences equidistant from each other simply by ( $n_1$ ,  $n_2$ ,  $n_3$ ,  $n_4$ ,  $n_5$ ) which guarantee that the four sequences are equidistant from each other.

Not all ( $n_1$ ,  $n_2$ ,  $n_3$ ,  $n_4$ ,  $n_5$ ) combinations are equally likely under the JC69 model with a star tree. For example, the site pattern combination (24, 24, 12, 96, 32) has a near-zero chance to occur with a

star tree and a strict JC69 model, because, given a star tree with  $x_5 = 0$ ,  $G_4$  sites can only emerge through independent substitutions along each of the four branches (i.e., all  $G_4$  sites result from convergent substitutions). This means that we cannot have  $G_4$  sites in a star tree without first having  $G_3$  sites, so  $G_4$  sites should not be more frequent than  $G_3$  sites given a star tree and JC69. The ratio of  $G_3/G_4$  sites will decrease from  $\infty$  (when the first substitution occurs) towards  $4/3$  (i.e., 48 possible  $G_4$  site patterns and 36 possible site patterns) when sequences have gradually experienced full substitution saturation. The site pattern combination above with a ratio of  $G_3/G_4$  equal to  $(12*4)/(96*3)$  cannot happen with a star tree because the ratio is far too small.

However, when different genes evolving under JC69 models with different rates (and potentially with different evolutionary histories and conflicting phylogenetic signals) are concatenated, strange site pattern combinations such as (24, 24, 12, 96, 32) may occur and cannot be dismissed as unreal. In fact, this site pattern combination results from a concatenation of site patterns from simulations with four different trees, i.e., the star tree and the three resolved trees (all with JC69 with no rate heterogeneity over sites), with slight modifications to ensure 1) that the nucleotide frequencies are exactly equal to 0.25, 2) that the number of transversions is exactly trice as many as the number of transitions, and 3) that the number of transitional and transversional differences between each pairwise comparison is exactly the same.

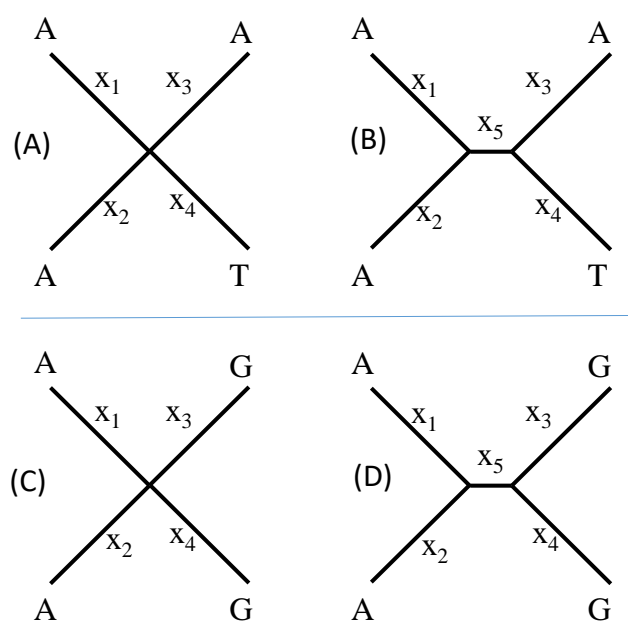
### 1.2. *The starless bias*

The site pattern combination (24, 24, 12, 96, 32) mentioned above represents a special deviation from any of the four topologies (the star tree plus the three resolved trees). While the four sequences will remain equidistant from each other with the site pattern combination of (24, 24, 12, 96, 32), the likelihood method will not favour a star tree in spite of equidistance among sequences. In general, as illustrated in Figure 2, we expect that increasing number of  $G_4$  sites relative to  $G_3$  sites would increase the likelihood of rejecting a star tree. As I mentioned before,  $G_3$  sites are the first site pattern to appear in sequence evolution along a star tree. In contrast, a  $G_4$  site requires a minimum of two substitutions when  $x_5 = 0$  (Figure 2C), but only a minimum of one with a resolved tree from a parsimony perspective (Figure 2D). Tree lnL would be greater with  $x_5 > 0$  (so that a single substitution can occur along the internal branch to be shared by two descendants) than with  $x_5 = 0$  (which would force two independent substitutions). Thus,  $G_4$  sites favor a resolved tree ( $x_5 > 0$ ). In short, increasing number of  $G_4$  sites relative to  $G_3$  sites increases the likelihood of rejecting the star tree. This is true even with  $G_4$  sites support all three resolved topologies equally because  $1/3$  of the  $G_4$  sites will support a resolved topology.

### 1.3. *Characterization of rate heterogeneity is confounded by topology*

I will use  $G_4$  sites to illustrate the effect of topology on rate heterogeneity over sites. With a star tree in Figure 2C, all  $G_4$  sites require exactly the same number of changes (a minimum of two substitutions per site from a parsimony perspective). In other words, there is no rate heterogeneity among  $G_4$  sites with a star tree. However, for a resolved tree in Figure 2D,  $1/3$  of the  $G_4$  sites (with identical nucleotides between sister groups as in Figure 2D) would require only one substitution from the parsimony perspective. The other  $2/3$  of the  $G_4$  sites (with different nucleotides between sister groups) would require two substitutions from the parsimony perspective. Thus, from a parsimony

perspective, 1/3 of the  $G_4$  sites evolve at a rate half as fast as the other 2/3 of the sites. Likewise with the likelihood perspective when multiple substitutions are corrected, 2/3 of the  $G_4$  sites will have a substitution rate at least twice as large as 1/3 of the  $G_4$  sites, giving rise to rate heterogeneity not present with a star tree. Thus, the relative numbers of  $G_4$  and  $G_3$  sites (denoted  $n_4$  and  $n_3$ , respectively, Figure 1A) affect not only the tendency to reject the star tree, but also the characterization of the rate heterogeneity over sites. In what follows, I assess the effect of  $(n_4 - n_3)$  on the tendency to reject the star tree and the confounding effect of topology on estimating the shape parameter  $\alpha$ .



**Figure 2.**  $G_3$  and  $G_4$  sites support star tree and resolved tree differently. (A) and (B) a  $G_3$  site mapped to a star tree and a resolved tree, respectively. (C) and (D) a  $G_4$  site mapped to a star tree and a resolved tree, respectively.

## 2. Materials and methods

### 2.1. Sets of sequences to make different points

I provide five additional sets of sequences in fasta format as representative of various site pattern combinations. Sample1.fas and Sample2.fas (s003.fas and s004.fas within the Supplementary genetics-05-04-212.zip) have sequences equivalent to site pattern combinations (0, 14, 43, 4, 40) and (1, 11, 41, 5, 42), respectively (Figure 1A). They are from sequence simulation under JC69 model with a star tree using Evolver in PAML [13] to show that likelihood methods will indeed recover a star tree if sequences indeed evolve according to a star tree and a specific substitution model.

Sample3.fas (s005.fas within the Supplementary genetics-05-04-212.zip) includes all 256 possible site patterns, i.e., all 24  $G_1$  sites, 144  $G_2$  sites, 48  $G_3$  sites, 36  $G_4$  sites and 4  $G_5$  sites. It is represented by the combination of (24, 24, 12, 12, 4). Such site patterns are expected from sequence evolution under JC69 for an infinitely long time so that all sites have experienced full substitution saturation. It is used to highlight the point that the likelihood method correctly recovers the star tree as it should even with full substitution saturation.

Supplementary sequence files Sample4.fas and Sample5.fas (s006.fas and s007.fas within the Supplementary genetics-05-04-212.zip) are used to make the main points in the paper. S006 features a site pattern combination of (24, 24, 12, 96, 32) from which likelihood method will not recover a star tree, although the sequences are equidistant from each other with distances computed with any substitution models. As mentioned earlier, this data set is a concatenation of site patterns from simulations with four different trees, i.e., the star tree and the three resolved trees based on JC69 with a constant rate over sites, with slight modifications to ensure that 1) nucleotide frequencies are exactly equal to 0.25, 2) all pairwise comparisons lead to exactly the same number of transitional and transversional differences, and 3) the ratio of transitional and transversional differences between each pairwise comparison is exactly 1/2. Because of the concatenation of sequences simulated from four different topologies, the sequence set represents a special deviation from any of the four topologies (the star tree plus the three resolved trees). While the four sequences will remain equidistant from each other with the site pattern combination of (24, 24, 12, 96, 32), the likelihood method will strongly reject the star tree in spite of equidistance among sequences. The sequence set is also used to investigate the confounding effect of tree topology on the estimation of shape parameter  $\alpha$  of gamma distribution.

Supplementary file Sample5.fas is used to assess the effect of increasing  $(n_4 - n_3)$  on decreasing support for the star tree. It contains 16 sets of sequences (with each set having four sequences) of equal sequence length of 2528 derived as follows. A simulation of 100 sets of sequences with a star tree, JC69 model without rate heterogeneity, a sequence length of 2528, and a tree length of 4.8, leads to site pattern combinations that average to (168, 216, 156, 120, 80). This site pattern combination, analyzed by either New2.R or NewGamma3.R, leads to the same tree lnL. The continuous gamma model with a star tree yields lnL equal to -13965.18, shape parameter  $\alpha$  equal to 5000 (maximum set in optimization), and branch lengths equal to 1.118709. Exactly the same results were obtained with a resolved tree with an additional internal branch length equal to 0. I changed  $n_3$  and  $n_4$  values to explore the effect of  $(n_4 - n_3)$  on the changing support for the star tree and on parameter estimation. In order to maintain the sequence length of 2528,  $n_1$ ,  $n_2$  and  $n_5$  were also adjusted.

## 2.2. Likelihood method implemented in R for star tree and gamma-distributed rate

For anyone to recreate results in this paper, I have implemented the likelihood method in R for the JC69 model and four OTUs to evaluate a star tree and a rooted tree, either with a constant rate over sites (Supplementary New2.R) or with a continuous gamma-distributed rate (NewGamma3.R). These are plain text files, self-contained, and well-annotated at the beginning of the file. One can copy and paste them into an R window to obtain results reported here. Computing likelihood with the pruning algorithm has been numerically illustrated in detail [14].

Because the continuous gamma requires integration over the rate, I used the R function 'integral' in the 'pracma' package which therefore needs to be installed before running the R scripts. The tree evaluation with a continuous gamma may take three minutes to complete on a PC with an i7-4770 CPU. I also used PhyML [15] to obtain lnL for the resolved tree. I set the tree improvement option '-s' to 'BEST' (best of NNI and SPR search), and the '-o' option to 'tlr' which optimizes the topology, the branch lengths and rate parameters. For JC69+ $\Gamma$  model, four categories of rates were used to estimate the shape parameter. PhyML does not generate lnL for a star tree, hence the need for the R scripts.

### 3. Results and discussion

#### 3.1. Likelihood method recovers the true star tree when sequences evolve under a substitution model

This part, while trivial, is needed for methodological validation. I have simulated sequence evolution of four OTUs under JC69 model and a star tree, by using the Evolver program in the PAML package [13], with different tree lengths and sequence lengths from 500 to 3000, each with 100 sets of sequences, with a constant rate over site. The sequences are then processed in DAMBE [16] so that the six site patterns in  $G_2$  (Figure 1A) occur exactly equally, so do the four site patterns in  $G_3$  and three site patterns in  $G_4$  (Figure 1A). This ensures that the resulting sequences do not support any one of the three alternative topologies. I analyzed these site patterns by both PhyML 3.1 [15] and the likelihood methods implemented in R in the Supplementary New2.R file and NewGamma3.R. All these methods recover the star tree, which has the same tree lnL as any of the resolved tree. Furthermore, the branch lengths are equal and nearly identical to the tree length in the input tree used for simulation.

Table 1 includes results for two such simulated sets of sequences without rate heterogeneity (A and B in the column headed by ‘Data’, Table 1). These two sets of sequences are in Supplementary Sample1.FAS and Sample2.Fas (s003 and s004). Their sequence patterns are represented by the combination of (0, 14, 43, 4, 40) and (1, 11, 41, 5, 42), respectively, for input to the R scripts. The resolved tree and the star tree consistently have the same tree lnL, whether evaluated by New2.R or PhyML (Table 1). Take Sample1.fas (Data A in Table 1) for example. The resolved tree and the star tree both have lnL of  $-1537.68$ , and the branch length for the internal branch of the resolved tree is  $b_5 = 0$  (Table 1). PhyML outputs  $b_5 = 0.000001$  which is likely the lower bound set for lnL maximization. Analyzing these two sets of sequences with JC69+ $\Gamma$  by PhyML or my NewGamma3.R script generates the same branch lengths and tree lnL, except for a very large shape parameter  $\alpha$  in the order of 10000 (which is expected because these two data sets were simulated with no rate heterogeneity over sites).

**Table 1.** Results of likelihood-based phylogenetic estimation of branch lengths ( $b_1$  to  $b_5$  corresponding to  $x_1$  to  $x_5$  in Figure 1B and Figure 1C) without using gamma-distributed rates to accommodate rate heterogeneity over sites, from running the attached R script (New2.R) and PhyML 3.1 [15].

Data <sup>(1)</sup>	Method	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$ <sup>(2)</sup>	lnL
A	New2.R	0.426784	0.426784	0.426784	0.426784	0.000000	-1537.68
		0.426784	0.426784	0.426784	0.426784	N/A	-1537.68
	PhyML	0.426785	0.426783	0.426793	0.426769	0.000001	-1537.68
B	New2.R	0.400581	0.400575	0.400584	0.400578	0.000000	-1416.09
		0.400584	0.400584	0.400584	0.400584	N/A	-1416.09
	PhyML	0.400586	0.400578	0.400570	0.400584	0.000001	-1416.09
C	New2.R	5.655316	2.585728	10	10	5.63553	-1419.57
		10	10	9.999938	10	N/A	-1419.57
	PhyML	10.000000	6.303193	10.000000	6.180340	6.303193	-1419.57

\*Note: (1) A to C corresponding to Supplementary file Sample1.fas to Sample3.fas (s003 to s005), respectively, with site pattern combination (0, 14, 43, 4, 40), (1, 11, 41, 5, 42), and (24, 24, 12, 12, 4), respectively. (2) N/A indicates  $b_5$  set to 0 for the star tree.

Data set C (Table 1) is in Sample3.fas with the site pattern combination of (24, 24, 12, 12, 4). It represents sequences having experienced full substitution saturation under JC69 model and should have infinitely long branch lengths. However, for practical computation one always set an upper limit for branch lengths, e.g., 10, partly because 1) it is rare for branch lengths to be longer than 10 in practice, and 2) lnL changes little when a branch length increases beyond 10. Minimum and maximum branch lengths can be specified in my two R scripts. Because PhyML 3.1 does not seem to allow branch lengths to go above 10, for comparison I have presented branch lengths and tree lnL with maximum branch lengths set to 10 (Table 1). Both the resolved tree and the star tree have the same lnL of -1419.57. One can change the maximum branch lengths to 100, 1000 or greater in my R scripts, and the star tree, which has one fewer branch to estimate, always has the same lnL as the resolved tree. Thus, the likelihood method does its job well as long as sequences evolve strictly according to substitution model, even with sequences having experienced full substitution saturation.

### 3.2. Likelihood method fails to recover a star tree with equidistant sequences

The likelihood method rejects a star tree for sequences in Supplementary file Sample4.fas (s006), with a site pattern combination of (24, 24, 12, 96, 32), in spite of equidistance among the sequences, either with or without gamma-distributed rate (Table 2). When a constant rate is assumed, the tree lnL is -2943.148 for a resolved tree in contrast to -2947.793 for a star tree (Table 2). A likelihood ratio test would lead to  $2\Delta\ln L = 9.29$ ,  $DF = 1$ ,  $p = 0.0023$  and a rejection of the star tree in favour of a resolved tree. The phylogenetic result with JC69+ $\Gamma$  rejects the star tree more strongly, with  $2\Delta\ln L = 58.03$ ,  $DF = 1$ ,  $p = 2.58 \times 10^{-14}$ . Note that the sequences length for this data set is only 536 nt. Longer alignment length would lead to even stronger rejection of the star tree.

**Table 2.** Evaluate two alternative topologies and their branch lengths ( $b_i$ ) in Figure 1B (with  $b_5$  set to 0 and not evaluated) and Figure 1C. Data in Sample4.fas file (s006), with site pattern combination (24, 24, 12, 96, 32).

Method <sup>(1)</sup>	$b_1$	$b_2$	$b_3$	$b_4$	$b_5$ <sup>(2)</sup>	$\alpha$ <sup>(3)</sup>	lnL
New2.R	0.849814	0.849814	1.693866	0	0.850	N/A	-2943.148
	1.020276	1.020276	1.020276	1.020276	N/A	N/A	-2947.793
PhyML	0.849942	0.849764	1.694148	0.000015	0.850	N/A	-2943.148
New $\Gamma$ 3	2.783640	1.133961	1.575337	2.343977	10	1	-2930.116
	3.359848	1.413517	1.886961	2.886549	30	0.843	-2921.743
	3.767304	1.850750	2.302868	3.315211	50	0.745	-2918.782
	1.020429	1.020429	1.020429	1.020429	N/A	10000	-2947.797
PhyML	2.213152	1.976554	2.206997	1.968304	10	0.814	-2927.806

\*Note: (1) Results from running Supplementary R scripts New2.R and NewGamma3.R (New $\Gamma$ 3), and PhyML. (2)  $b_5$  is  $x_5$  in Figure 1. N/A indicates  $b_5$  set to 0 for a star tree. (3) Shape parameter in gamma distribution, inapplicable (N/A) in phylogenetic reconstruction with a constant rate.

While a site pattern of (24, 24, 12, 96, 32) might be considered too extreme, one may take a milder site pattern combination such as (120, 192, 120, 204, 164), with a sequence length of 2528. The resulting lnL based on JC69+ $\Gamma$  is -13880.12 for a resolved tree, but -13889.4900 for a star tree, so  $2\Delta\ln L = 18.74$ . The star tree is rejected with  $p = 0.000015$ .



One might argue that there is nothing wrong with the likelihood method itself because the data set is pathological. As mentioned in the Materials and Methods section, this data set is from a concatenation of simulated sequences from four topologies (the star tree and the three resolved trees) with modifications so that nucleotide frequencies are equal and all pairwise comparisons lead to the same number of transition and transversional differences. Thus, although the sequences are equidistant from each other, the star tree is not an appropriate model for the concatenated sequences (neither is any of the three resolved topologies). However, the issue at hand is not on the validity of the likelihood approach, but on its robustness in phylogenetic reconstruction with conflicting phylogenetic signals. With sequences equidistant from each other, we do desire a star tree instead of having it conclusively rejected. Furthermore, we have the problem of inconsistent parameter estimation that I highlight below.

The PAML package [13] has a `baseml` program which implements a continuous gamma-distributed rate. However, it does not produce correct results. Take for example the sequences in `Sample4.fas` (s006). If I use a tree with  $b_1 = b_2 = b_3 = b_4 = 4$  and  $b_5 = 0.000001$ , `baseml` will generate an  $\ln L$  of  $-2957.037660$  and estimated  $b_1 = b_2 = b_3 = b_4 = 2.59636$ , together with a warning of “check convergence”. If I change the  $b_1$  to  $b_4$  values to other values such as 1, 2 or 3, the same output was generated. The output  $\ln L$  ( $= -2957.037660$ ), with JC+Gamma, is far too small. Even without gamma,  $\ln L$  should be  $-2947.797$  instead of  $-2957$ . I wrote `NewGamma3.R` main because of the inadequate performance of `baseml`.

### 3.3. Inconsistent parameter estimation of likelihood method

The estimated parameter values in Table 2 are disconcerting in two ways. First, when the star tree is imposed, there is no rate heterogeneity over sites, with the shape parameter  $\alpha$  in the order of 10000 (Table 2). However, the estimated  $\alpha$  becomes 0.745 when  $b_5$  ( $x_5$  in Figure 1) is allowed to be greater than 0, indicating strong rate heterogeneity over sites. I have given reasons in Figure 2, illustrated with  $G_4$  sites, that an excess of  $G_4$  sites will result in rate heterogeneity when a resolved tree is imposed but no rate heterogeneity when a star tree is imposed. That is, given a set of  $G_4$  sites supporting the three resolved topologies equally, 1/3 of the  $G_4$  sites will share a low rate of substitution whereas the other 2/3 will share a high rate of substitution when a resolved topology is imposed. In contrast, all  $G_4$  sites will have the same rate of substitution when a star tree is imposed. This highlights the point that we do need a star tree instead of having three equally supported resolved trees because we need the star tree to perform proper parameter estimation in this case. As I mentioned before, the sequence data is a concatenation of sequence simulation on four different topologies (star and three resolved trees), all with JC69 with a constant rate. However, the maximum likelihood criterion does not allow us to choose the star tree.

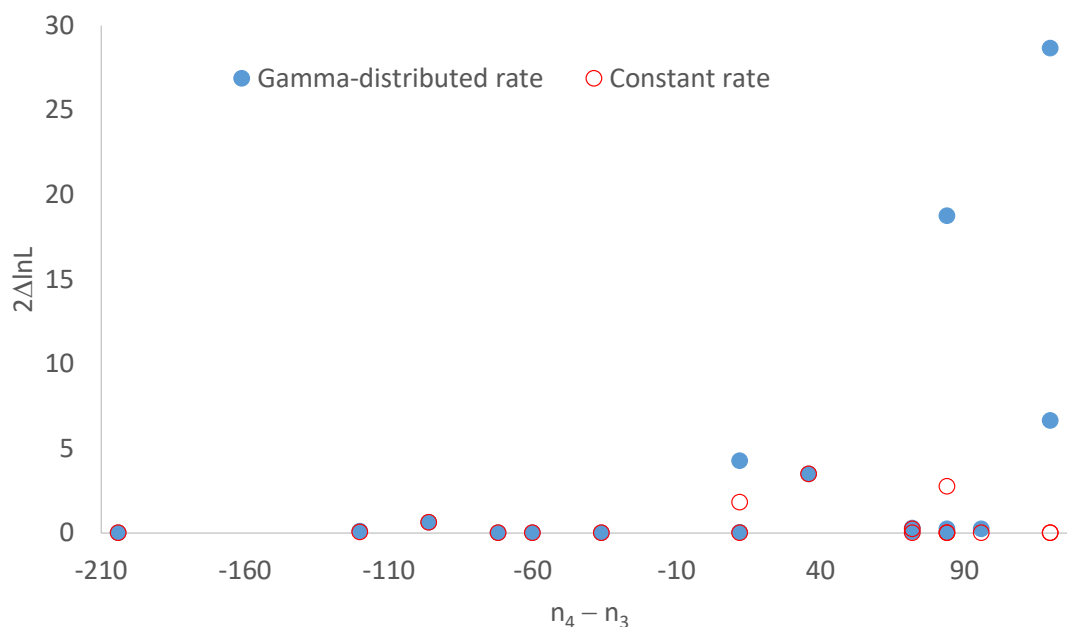
Second, the substantial change of  $\ln L$  with  $b_5$  under JC69+ $\Gamma$ , associated with a change in  $\alpha$ , is also unpleasant. For all practical consideration, the range of values for  $b_5$  from 10 to 50 all just indicates a branch experiencing substitution saturation, and one would expect  $\ln L$  to change little with  $b_5$  in this range of values. However,  $\ln L$  is  $-2930.116$  with  $b_5 = 10$  (and  $\alpha = 1$ , Table 2), but  $-0.2918.782$  with  $b_5 = 50$  (and  $\alpha = 0.745$ , Table 2). Given that the sequences are equidistant from each other, we do desire a tree with  $b_5 = 0$ , not a method that says that a tree with an unrealistically long  $b_5$  is the best tree. Note that PhyML (version 3.1) appears to have set the maximum branch length to 10, so its generated  $\ln L$  ( $= -2927.806$ , Table 2) has not yet reached its maximum.

While  $\alpha$  is known to depend on topology, phylogenetic researchers often interpret  $\alpha$  as if it is a sequence-specific property. A small  $\alpha$  is often interpreted to mean strong purifying selection constraining certain sites but not others. I hope that the illustration above will caution against such interpretation.

### 3.4. The tendency to reject the star tree with increasing $(n_4 - n_3)$

I used sequences in Sample5.fas (s007) to assess the effect of  $(n_4 - n_3)$ . The file contains 16 sets of sequences of the same length (= 2528) but with different site patterns  $(n_1, n_2, n_3, n_4, n_5)$ , in particular with different  $(n_4 - n_3)$  values. We may take likelihood chi-square statistic  $2\Delta\ln L$  as a measure of the tendency to reject the star tree, and expect it to increase with  $(n_4 - n_3)$  for reasons outlined in Figure 2. This expectation is consistent with the empirical evidence (Figure 3). The relationship is stronger between  $2\Delta\ln L$  and  $n_4/n_3$ . Also, if we change on  $n_4$ , but keep  $n_1, n_2, n_3$ , and  $n_5$  constant, then  $2\Delta\ln L$  increases smoothly with  $(n_4 - n_3)$  or with  $n_4/n_3$ .

Figure 3 indicates that the relationship between  $2\Delta\ln L$  and  $(n_4 - n_3)$  is weaker when rate heterogeneity over sites is modelled by a gamma distribution. However, it is not always so, and one can easily find a counter example in which a constant-rate model leads to rejection of the star tree and a gamma-distributed rate model does not. For example, for a site pattern combination of (400,0,0,0,400) we will have  $\ln L = -3984.64$  for a resolved tree with JC69,  $\ln L = -3987.998$  for a star tree. This leads to  $2\Delta\ln L = 6.716$  and a rejection of the star tree (DF = 1,  $p = 0.0096$ ). With JC69+ $\Gamma$ , we will have  $\ln L = -3546.648$  for a resolved tree,  $\ln L = -3547.94$  for a star tree. This leads to  $2\Delta\ln L = 2.584$  and the star tree is not rejected (DF = 1,  $p = 0.1079$ ).



**Figure 3.** Support for a resolved tree against the star tree, measured by  $2\Delta\ln L = 2(\ln L_{\text{resolved tree}} - \ln L_{\text{star tree}})$ , increases with increased  $G_4$  sites relative to  $G_3$  sites (measured by  $n_4 - n_3$ ). The open circles indicates the mean  $(n_4 - n_3)$  value from 100 simulated sequences with a star tree and  $x_1 = x_2 = x_3 = x_4 = 1.2$ .

I should finally mention that the starless bias and parameter-estimation bias illustrated by the site pattern combination of (24, 24, 12, 96, 32) in Sample4.fas (s006) is also caused by model misspecification, in the sense that no tree model is appropriate for the data (because it is concatenation of sequences simulated from the star tree and three resolved trees). Existing model-testing methods do not address this type of model misspecification. The conventional selection of the best substitution model will choose JC69 because lnL is the same from JC69 to GTR. This would give us false confidence that we are using an appropriate substitution model for data analysis, without realizing that no tree model is appropriate (i.e., neither the star tree nor the resolved trees) for this set of concatenated sequences.

#### 4. Conclusions

Two approaches are relevant to this problem of different topologies contributing conflicting phylogenetic signals caused by concatenating sequences of potentially different evolutionary history. The first is to use mixture model with different tree models each contributing a subset of sites. This is different from mixture models with the same topology but different nucleotide or amino acid substitution models. The second approach is as follows. For each phylogeny (T) obtained from a set of sequence data (S) under a substitution model (M), we should obtain the empirical site patterns from S and compare these empirical site patterns against their expectation from sequence simulation based on T and M. If the empirical site patterns deviate significantly from the expectation, then we conclude that the tree model is not appropriate. Applying this approach to sequences in Sample4.fas (s006), we will find the empirical site patterns from the sequence data differ highly significantly from the expectation, and we can conclude that no tree model is appropriate for this set of concatenated sequences or, more specifically, that the sequences are from a mixture of incompatible tree models.

#### Acknowledgements

This study is supported by Discovery Grant from Natural Science and Engineering Research Council (NSERC, RGPIN/ 2018-03878) of Canada. I thank Guy Baele, Sudhir Kumar, Laura Kubatko and Arindam RoyChoudhury for discussion. Z. Yang and two anonymous reviewers provided comments that improved clarity of the manuscript.

#### Conflict of interest

The authors declare no conflict of interest.

#### References

1. Saitou N, Nei M (1987) The neighbor-joining method: A new method for reconstructing phylogenetic trees. *Mol Biol Evol* 4: 406–425.
2. Desper R, Gascuel O (2004) Theoretical foundation of the balanced minimum evolution method of phylogenetic inference and its relationship to weighted least-squares tree fitting. *Mol Biol Evol* 21: 587–598.

3. Xia X (2014) Phylogenetic bias in the likelihood method caused by missing data coupled with Among-Site rate variation: An analytical approach. In: Basu M, Pan Y, Wang J, editors. *Bioinformatics Research and Applications: Springer*, 12–23.
4. Jukes TH, Cantor CR (1969) Evolution of protein molecules. In: Munro HN, editor. *Mammalian Protein Metabolism*. New York: Academic Press, 21–123.
5. Hasegawa M, Kishino H (1989) Heterogeneity of tempo and mode of mitochondrial DNA evolution among mammalian orders. *Jpn J Genet* 64: 243–258.
6. Kishino H, Hasegawa M (1989) Evaluation of the maximum likelihood estimate of the evolutionary tree topologies from DNA sequence data, and the branching order in Hominoidea. *J Mol Evol* 29: 170–179.
7. Hasegawa M, Kishino H, Yano T (1985) Dating of the human-ape splitting by a molecular clock of mitochondrial DNA. *J Mol Evol* 22: 160–174.
8. Tamura K, Nei M (1993) Estimation of the number of nucleotide substitutions in the control region of mitochondrial DNA in humans and chimpanzees. *Mol Biol Evol* 10: 512–526.
9. Lanave C, Preparata G, Saccone C, et al. (1984) A new method for calculating evolutionary substitution rates. *J Mol Evol* 20: 86–93.
10. Tavaré S (1986) Some probabilistic and statistical problems in the analysis of DNA sequences. In: Miura RM, editor. *Lectures on Mathematics in the Life Sciences*. Providence, RI: Amer Math Soc: 57–86.
11. Xia X (2017) Deriving transition probabilities and evolutionary distances from substitution rate matrix by probability reasoning. *J Genet Genome Res* 3: 031.
12. Xia X (2018) Nucleotide substitution models and evolutionary distances. *Bioinf Cell*: 269–314.
13. Yang Z (2007) PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol* 24: 1586–1591.
14. Xia X (2018) Maximum likelihood in molecular phylogenetics. *Bioinf Cell*: 381–395.
15. Guindon S, Dufayard JF, Lefort V, et al. (2010) New algorithms and methods to estimate maximum-likelihood phylogenies: Assessing the performance of PhyML 3.0. *Syst Biol* 59: 307–321.
16. Xia X (2018) DAMBE7: New and improved tools for data analysis in molecular biology and evolution. *Mol Biol Evol* 35: 1550–1552.



AIMS Press

© 2018 the Author(s), licensee AIMS Press. This is an open access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0>)