

Phylogenetics

State aggregation for fast likelihood computations in molecular evolution

Iakov I. Davydov^{1,2}, Marc Robinson-Rechavi^{1,2} and Nicolas Salamin^{1,2,*}

¹Department of Ecology and Evolution, Biophore, University of Lausanne, 1015 Lausanne, Switzerland and ²Swiss Institute of Bioinformatics, Genopode, Quartier Sorge, 1015 Lausanne, Switzerland

*To whom correspondence should be addressed.

Associate Editor: Alfonso Valencia

Received on May 13, 2016; revised on September 7, 2016; accepted on September 23, 2016

Abstract

Motivation: Codon models are widely used to identify the signature of selection at the molecular level and to test for changes in selective pressure during the evolution of genes encoding proteins. The large size of the state space of the Markov processes used to model codon evolution makes it difficult to use these models with large biological datasets. We propose here to use state aggregation to reduce the state space of codon models and, thus, improve the computational performance of likelihood estimation on these models.

Results: We show that this heuristic speeds up the computations of the M0 and branch-site models up to 6.8 times. We also show through simulations that state aggregation does not introduce a detectable bias. We analyzed a real dataset and show that aggregation provides highly correlated predictions compared to the full likelihood computations. Finally, state aggregation is a very general approach and can be applied to any continuous-time Markov process-based model with large state space, such as amino acid and coevolution models. We therefore discuss different ways to apply state aggregation to Markov models used in phylogenetics.

Availability and Implementation: The heuristic is implemented in the godon package (<https://bitbucket.org/Davydov/godon>) and in a version of FastCodeML (<https://gitlab.isb-sib.ch/phylo/fastcodeml>).

Contact: nicolas.salamin@unil.ch

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

Evolutionary models are necessary to study the processes governing the evolution of genes, genomes and organisms. While relatively simple models are often sufficient to provide a good estimation of species or gene trees, inferring the specific processes that govern the evolution of molecular data (e.g. selection or co-evolution) requires more complex models. The ability to apply these complex models to large datasets involving many genes and/or species offers the promise to better understand evolution in a more general context. This approach has, however, an important computational cost because of the large numbers of parameters and/or the large size of the state space involved in these complex models.

The computational performance of phylogenetic methods has always been an important issue in molecular evolution. Likelihood-based methods in phylogenetics would not be possible without the use of Felsenstein's tree pruning algorithm (Felsenstein, 1981) coupled with the growth of computer performance. However, these methods only became commonly used with the heuristics implemented in software such as PhyML and RaXML (Guindon *et al.*, 2010; Stamatakis, 2014). Recent years have thus seen tremendous decreases in computing times, to the extent that datasets with thousands of sequences can now be analyzed. However, most progress has been made on simple models of DNA or amino-acid evolution. More complex models, such as codon models used to detect

selection, are still computationally too costly to be applied on large genomic datasets (e.g. all Ensembl Compara; Vilella *et al.*, 2009).

The complexity of codon models comes from the large state-space that is necessary to represent the 61 codons (excluding the three stop codons). The simplest codon model, which is called M0 (Goldman and Yang, 1994), assumes a single parameter ω to model a constant selective pressure occurring on all sites and branches of a phylogenetic tree. The M0 model is probably not realistic enough and more complex models that involve multiple transition matrices have been developed to detect episodic positive selection on a subset of sites and of phylogenetic branches (Murrell *et al.*, 2012; Smith *et al.*, 2015; Zhang *et al.*, 2005). One of the most commonly used complex models is the branch-site model (Zhang *et al.*, 2005), which assumes three classes of selection (parameters $\omega_0, \omega_1, \omega_2$ with ω_2 allowing positive selection) on sites along specific branches of the tree (called foreground branches) and two classes (parameters ω_0 and ω_1) on the other branches.

Since an accurate phylogenetic tree is critical to evolutionary and comparative studies, most developments to speedup the parameter estimation of evolutionary models have focused first on the optimization of search strategies to find the tree topology and branch lengths. Examples include the choice of the starting tree topology (Guindon and Gascuel, 2003; Huelsenbeck *et al.*, 2001; Nguyen *et al.*, 2014; Stamatakis *et al.*, 2004; Stamatakis, 2014), improved tree rearrangements strategies (Guindon and Gascuel, 2003; Hordijk and Gascuel, 2005; Nguyen *et al.*, 2014; Stamatakis *et al.*, 2005; Swofford and Olsen, 1990), computation economy (Gladstein, 1997; Goloboff, 1993; Ronquist, 1998) and independent branch-length estimation (Guindon and Gascuel, 2003).

However, an important part of the computational cost is spent calculating the likelihood function itself. Although this part is not the most limiting step for tree searching methods using simple models, it becomes a major bottleneck for the evaluation of more complex evolutionary scenarios such as codon models. In this case, the reuse of the eigenvectors and eigenvalues for a set of branches can improve computational performance (Schabauer *et al.*, 2012; Valle *et al.*, 2014). Other optimization techniques that involve, for example, transforming the problem of exponentiating an asymmetric matrix into a symmetric problem, or performing matrix-matrix multiplication rather than matrix-vectors for the estimation of conditional vectors, have also been shown to speedup the calculations of the likelihood (Schabauer *et al.*, 2012). There has also been some progress on Bayesian computation, e.g. using data augmentation (de Koning *et al.*, 2012; Lartillot, 2006; Rodrigue *et al.*, 2008). Despite these improvements, likelihood calculations still remain computationally intensive.

The size of the state-space of the continuous-time Markov chain directly impacts the most computationally intensive steps of this likelihood computation, since it affects the size of the rate and probability matrices (Q and P , see below), as well as of the conditional probability vector. A method allowing a reduction of the number of states while affecting minimally the precision of the likelihood estimation is therefore a potentially interesting avenue to further reduce the computational burden of these methods.

We propose here a heuristic method to speedup matrix exponentiation and partial likelihood calculations by reducing the number of states in a continuous-time Markov chain without losing the complexity of the model. We use state aggregation techniques to selectively combine states of the instantaneous rate matrix. We illustrate this technique with a simple and a complex codon model, since their state-space is relatively large (61 states). We show using simulations and the analysis of an empirical dataset that aggregation can provide significant speedup for codon models, with a very low cost in terms of accuracy. We further discuss the potential biological applications

that could benefit from this approach to illustrate the wide applicability of state aggregation.

1.1 Key steps of likelihood computation in phylogenetics

The performance of the likelihood calculations are governed by two computationally intensive steps: matrix exponentiation and matrix-vector multiplication.

Matrix exponentiation is at the heart of models based on continuous-time Markov chains. The rate of change from one state to any other in an infinitesimally small time interval is given by the instantaneous rate matrix Q . The probability of changing between the states of the process in a time interval t is then given by the probability matrix P : $P(t) = e^{Qt}$. For computational purposes, the rate matrix is first diagonalized such that $Q = \Lambda U^{-1}$, where U is the matrix of eigenvectors and Λ is a matrix whose diagonal elements correspond to the eigenvalues of the instantaneous matrix Q . This matrix decomposition allows the probability matrix P to be quickly computed for any time interval t as $P(t) = e^{Qt} = Ue^{\Lambda t}U^{-1}$.

Branches of a phylogenetic tree represent the evolutionary path between an ancestral sequence and its descendants. We therefore need to compute the matrix P for every branch of a tree. The instantaneous rate matrix Q needs thus to be exponentiated for every branch length. The probabilities of observing the states in the ancestral sequence are then calculated by multiplying the conditional probability vectors for each descendant branch. These probability vectors are obtained by multiplying the P matrix for branch i with the conditional vector of the corresponding descendant. This procedure, known as Felsenstein's tree pruning algorithm, is repeated for every node of the phylogenetic tree until we reach the root of the tree (Felsenstein, 1973).

2 Algorithm

2.1 State aggregation

The computational cost of the two steps described above highly depends on the state-space of the continuous-time Markov chain used. Any reduction in the state-space can therefore increase the efficiency of the likelihood calculations. We investigate here the use of state aggregation to combine states of a Markov chain into several groups and therefore reduce the complexity of matrix exponentiation and matrix-vector multiplication.

Let us consider a Markov chain taking values in a finite set $S = \{A_1, A_2, \dots, A_n\}$ with transition matrix P and stationary frequencies $\pi_1, \pi_2, \dots, \pi_n$. Let $S^c = \{A_1, A_2, \dots, A_m\}$ be a set of states to be aggregated, where $m < n$.

The aggregated chain will have a space of

$$\tilde{S} = \{A^C, A_{m+1}, A_{m+2}, \dots, A_n\},$$

where A^C is the aggregated state. The new aggregated state A^C changes the entries of the probability matrix P in the following way:

$$\begin{aligned} \tilde{p}_{A_i, A_j} &= p_{A_i, A_j}, \\ \tilde{p}_{A_i, A^C} &= \sum_{A_k \in S^c} p_{A_i, A_k}, \\ \tilde{p}_{A^C, A_j} &= \frac{1}{\tilde{\pi}_{A^C}} \cdot \sum_{k \in S^c} \pi_k p_{k, j}, \\ \tilde{p}_{A^C, A^C} &= \frac{1}{\tilde{\pi}_{A^C}} \sum_{k \in S^c} \sum_{l \in S^c} \pi_k p_{k, l}, \end{aligned}$$

where $A_i, A_j \notin S^c$.

The stationary frequencies are estimated as $\tilde{\pi}_i = \pi_i$ for $A_i \notin S^C$. These stationary frequencies are consistent with frequencies of the original Markov chain. The frequencies of the aggregated state is estimated as $\tilde{\pi}_{AC} = \sum_{k \in S^C} \pi_k$.

The same method can be applied at the level of the instantaneous rate matrix Q . The diagonal elements of the matrix must however be set to $-\sum_{j \neq i} q_{i,j}$ to ensure that the sum of every row is equal to zero (Supplementary Fig. S1B, C) (Aldous and Fill, 2002, chapter2).

2.2 Aggregation for codon models

An obvious question in performing aggregation is the definition of ‘similar states’ to aggregate. We define all non-observed states for a position of the alignment to be ‘similar’ in the context of that position. The rationale is that the codons that are not observed at this site in any of the sequences at the tips of the tree have low probability to occur as ancestral states. The lack of some possible codons could be due to chance, but in many cases we expect a subset of codons to occur at a site because of natural selection or mutational bias. For example, a protein site which is constrained to be negatively charged will only use codons encoding such amino acids. It is thus justified to call all other codons ‘similar’ relative to this site. We therefore aggregated all states unobserved at a position (i.e. triplet of columns of the DNA alignment) into a single state (Fig. 1). The approach that we use here to aggregate states in codon models resembles the models of amino acid and nucleotide substitutions proposed in Yang *et al.* (1998); Susko and Roger (2007); Phillips *et al.* (2004); Vera-Ruiz *et al.* (2014). However, we propose to select the new aggregated state-space independently for each position of the alignment, which was not done in the amino acid and nucleotide contexts. Note that we performed in this study only the aggregation on the probability matrix P . We discuss the advantages of aggregation on the P or Q matrices in the Discussion section.

For the aggregated process to have a Markovian property, it has to satisfy the lumpability condition, i.e. $p_{A_i, A_k} = p_{A_j, A_k}$ should be true for any $i, j \in A^C$ and $k \notin A^C$ (Kemeny and Snell, 1983), or equivalently for the instantaneous rate matrix (Hillston, 1995). This condition is generally not satisfied with respect to an arbitrary aggregation scheme, as this would require all the transition rates or substitution probabilities to have the same value. Moreover, in the widely used codon substitution models, double substitutions are not allowed and their respective transition rates are set to zero, which makes lumpability condition unsatisfiable. Thus the proposed technique should be viewed as a heuristic.

The intensity of state aggregation can be modified and we tested three different approaches by implementing them for the codon model M0 (Goldman and Yang, 1994). The first and least aggressive approach aggregates only the positions that were absolutely conserved in any sites of the alignment. The state-space for these sites is thus reduced to two states: the observed (conserved) codon, and the

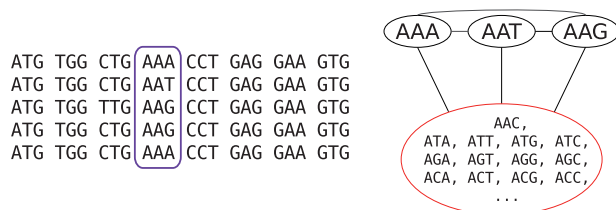


Fig. 1. Example of state aggregation for one position (highlighted in purple) in a codon alignment

‘meta-state’ of the 60 other non-stop codons. In the second approach, all positions were aggregated and the ‘meta-state’ included all codons not present in the position subjected to aggregation. In the third approach, all the positions were aggregated, but the ‘meta-state’ included only codons corresponding to the amino-acids not present at the current position. This can be viewed as a less-aggressive version of the second approach utilizing properties of the genetic code. The first two approaches represent extreme cases of the application of aggregation, while the third one is more moderate.

Given the small speedup of the first and third approaches on M0 (see Results), only the second approach was employed for the more complex branch-site model.

Additionally, two random aggregation strategies were evaluated. These strategies were used as a control to determine if our choice of state partitioning is better than random. In the first strategy, ‘meta-states’ of full aggregation were shuffled between the alignment positions. This should give a speedup similar to the full aggregation, while not relying on codons present at each position. In the second random strategy, the state-space was randomly split into ‘meta-states’, while keeping the total number of states per position. The number of states stays the same in this case, but the computations are expected to take more time since multiple ‘meta-states’ are present.

3 Materials and methods

3.1 Software

State aggregation for the M0 model was implemented in the godon package (<https://bitbucket.org/Davydov/godon>). We selected an optimization algorithm with a large but fixed number of iterations (10 000 iterations in this case) to reduce the influence of random factors associated with the optimization trajectory on the total computation time.

State aggregation for the branch-site model was implemented in a version of FastCodeML (<https://gitlab.isb-sib.ch/phylo/fastcodeml>, branch agg), which is a software that has been optimized for computational efficiency of the calculation of the matrix exponentiation and the matrix-vector multiplication (Valle *et al.*, 2014).

All sequence simulations were performed using the evolver program from the PAML package (Yang, 2007).

3.2 Dataset

Six datasets were simulated for the M0 model (see Supplementary Table S1). We varied one parameter at a time, based on the following settings: 300 codons, 18 sequences, $\omega_0 = 0.3$, $\kappa = 2$, equal codon frequencies ($\pi_i = 1/61$), default tree length (4). We used the the ENSGT00680000099620 gene tree from the Ensembl database (Cunningham *et al.*, 2015) for topology and relative branch lengths.

For the branch-site model, 2000 alignments were simulated with stochastic birth-death trees and κ , ω_0 , ω_2 , p_0 , p_1 , alignment length and number of tips sequences chosen randomly (Supplementary Table S2A, Fig. S2). One thousand of the alignments were simulated under the branch-site model null hypothesis with $\omega_2 = 1$, while the other 1000 alignments represented the alternative hypothesis with $\omega_2 > 1$. In these simulations every parameter was drawn randomly from a specific distributions (Supplementary Table S2) to obtain more biologically realistic datasets. We chose at random a single foreground branch to perform the simulations and the same

foreground branch was used foreground the inference. We used evolver from PAML 4.8 for simulation (Yang, 2007).

We also simulated datasets using an *extended branch-site model*. In this model, ω_0 and ω_2 (when $\omega_2 > 1$) were replaced by a set of discrete categories created from *Beta* and *Gamma* distributions respectively (five discrete categories were used). We also incorporated *Gamma* distributed site-rate variations (Rubinstein et al., 2011). The parameters for these distributions are described in [Supplementary Table S2B](#) and Figure S3. We used the *cosim* package for the simulations (<https://bitbucket.org/Davydov/cosim>).

The likelihood ratio test (LRT) was used for model selection, with a significance level of $\alpha = 0.05$.

We performed multiple hypothesis testing correction using the *qvalue* R package, π_0 was estimated using the bootstrap method (Storey et al., 2004).

Finally, a Primates dataset from the Selectome database (Moretti et al., 2014; Proux et al., 2009) release 6 was used to study the behaviour of the method on a real dataset. The dataset consists of 15 669 gene trees and alignments (<http://selectome.unil.ch/cgi-bin/download.cgi>). We tested the inference of selection on every non-terminal branch of the Primates trees.

4 Results

4.1 M0 model

For the simple M0 model, we first compare the performance of likelihood maximization in three different modes: full likelihood (no aggregation), aggregation for conserved positions and full aggregation. Here we kept the branch lengths fixed and optimized ω and κ .

The parameter values obtained for all datasets using both aggregation modes are highly correlated with values estimated by the full likelihood ([Supplementary Figs. S4, S5, S6, S7, S8, S9](#)). The error in estimation of the parameters is small and is not dependent on the simulation parameters ([Figs. 2, S10, S11](#)), with the exception of tree length ([Fig. 3](#)). The bias in parameter estimation associated with the long trees is smaller for less aggressive aggregation strategies ([Supplementary Fig. S12](#)). Comparisons with the two random aggregation strategies show noticeably better accuracies in parameter estimation with the observation-based aggregation ([Supplementary Fig. S13](#)).

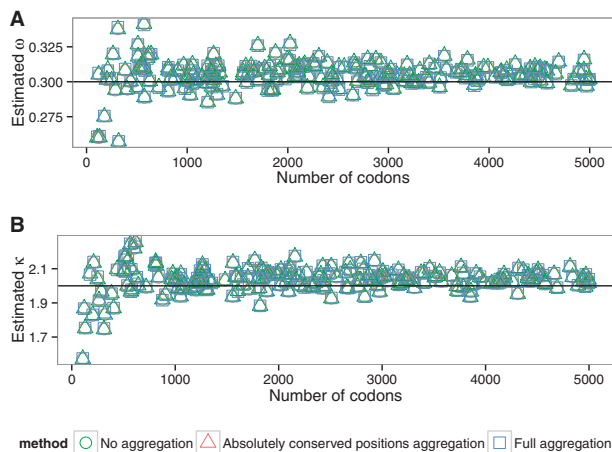


Fig. 2. Estimated ω (A) and κ (B) values depending on the alignment length (alen dataset, M0 model). Lines correspond to the simulation parameter values

The mean computational speedup is approximately 1.7 for aggregation on all positions ([Fig. 4](#)), but only 1.2 and 1.02 ([Supplementary Fig. S14](#)) for genetic code based aggregation and aggregation limited to fixed positions respectively. We thus only analyzed in details the behavior of the full aggregation mode.

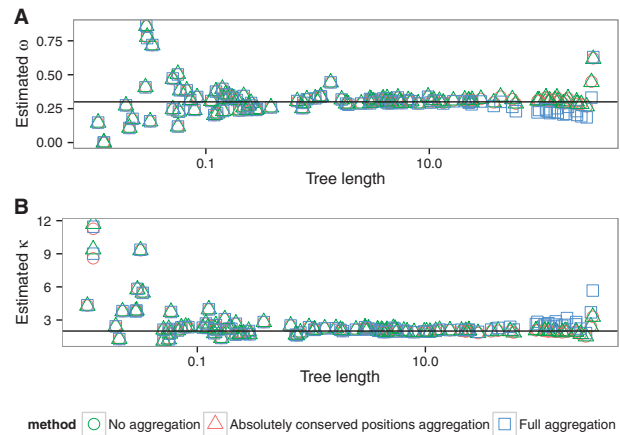


Fig. 3. Estimated ω (A) and κ (B) values depending on the tree length (tlen dataset, M0 model). Lines correspond to the simulation parameter values. Tree length limited to the range [0.01; 300], see text

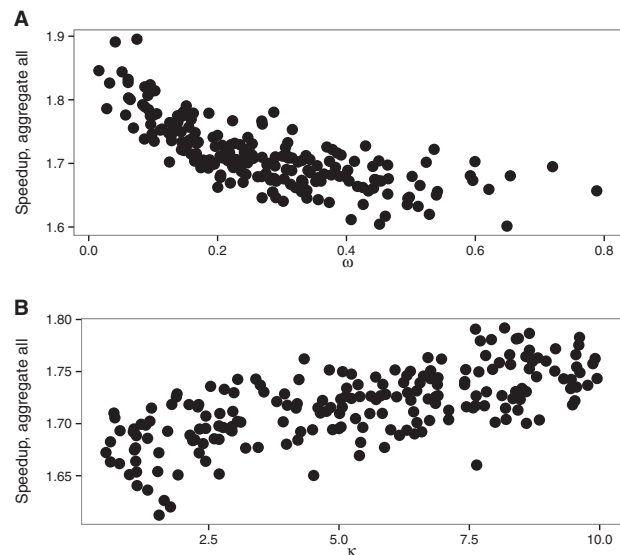


Fig. 4. Speedup of aggregation on all alignment positions with M0 model. (A) vvar dataset with variable ω value, (B) kvar dataset with variable κ

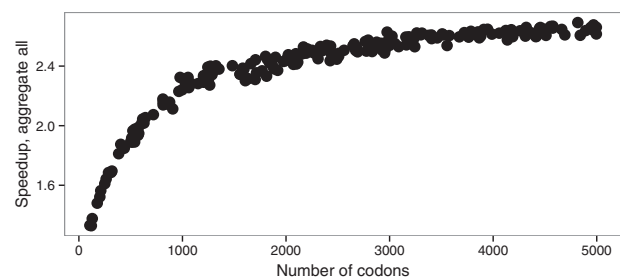


Fig. 5. Speedup depending on a number of the codons (alen dataset, M0 model)

First, we see a strong effect of the alignment length on the speedups obtained (Fig. 5). Matrix eigendecomposition is performed only once per likelihood evaluation and the decrease in the state-space between the full likelihood and the aggregation does not have any impact on the eigendecomposition performance. However, a longer alignment will increase the number of times the tree pruning step is performed (i.e. once per site), which becomes more important in the overall computational cost. For instance, the speedup obtained with an alignment of 500 codons is 1.8 with full aggregation. The maximum speedup of 6.8 fold was achieved on extremely long alignments (above 10 000 codons) and short trees (total length < 0.05).

While there is a larger error on the estimation of model parameters (κ and ω) with shorter alignments, this effect is identical with or without aggregation (Fig. 2). The heuristic that we propose does therefore not increase error on a simple model even with short alignments. Interplay between eigendecomposition and pruning times explains the direct effect of the relationship between the number of sequences and the speedup (Fig. 6). A large number of sequences decreases the proportion of time spent in the eigendecomposition phase and subsequently increases the speedup.

Changes in the other parameters impact the speedup of the aggregation mode insofar as they change the number of codon states per alignment site. The latter has then a direct effect on the number of non-aggregated states. Indeed, we aggregate into one state all codons which are not observed in a given position. Thus any processes that reduces the number of different codons per position also increases the efficiency of aggregation. Hence, the speedup is slightly higher for smaller ω values because, as ω approaches 0, more and more codons at a particular site are only part of a synonymous codon set. The number of possible codons is thus greatly reduced and there is a higher chance that the aggregation will lead to very few states. In contrast, increasing ω values will lead to an increasing number of states observed. Similarly, extremely short branches limit state variety at each site, which in turn increase the level of aggregation possible and thus increase speedup (Fig. 7). Biased codon frequencies can also reduce diversity of states and thus increase aggregation speedup. In our simulations, codon frequencies were drawn from a Dirichlet distribution and we varied the concentration

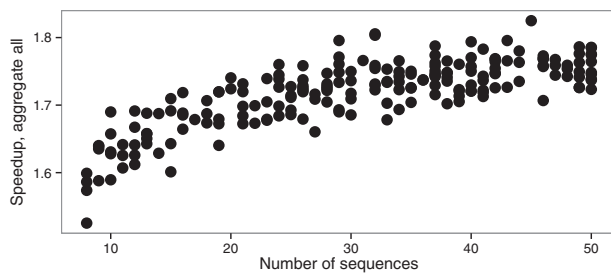


Fig. 6. Speedup versus number of sequences (nseq dataset, M0 model)

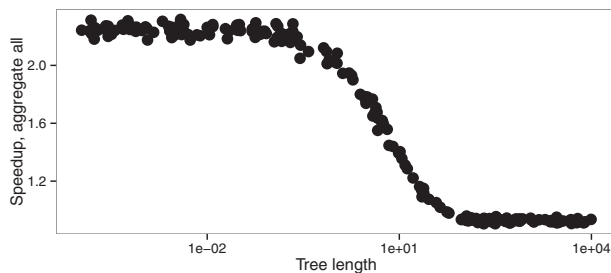


Fig. 7. Speedup versus tree length (tlen dataset, M0 model)

parameter α to estimate the effect of codon frequencies on codon aggregation. We see a better speedup associated with smaller values of the α parameter, which leads to a higher variance between codon frequencies (Supplementary Fig. S15).

Total tree length is the only parameter in our simulations that also affects the accuracy of the estimation of model parameters (Figs. 3, S16). Longer trees tend to improve the accuracy of the estimation of the parameters ω and κ . However, extremely long trees lead to an increase in error both in aggregated and in full likelihood mode, probably because of saturation. It appears that under reasonable conditions of applicability of the M0 model (i.e. total tree length < 20 substitutions per codon), aggregation does not lead to any detectable bias, while for extremely long trees aggregation can introduce a slight bias.

We also estimated the branch lengths during the optimization of the M0 model. There was no systematic bias in branch lengths estimation for short trees (Supplementary Fig. S17), while we observed an increased error in branch lengths estimation on extremely long trees (Supplementary Fig. S18). The error on branch lengths was accompanied by increased errors on ω and κ .

Thus, overall speedup on the simple M0 model can be explained by average observed codons count and by alignment length (Fig. 8). The relationship between speedup and ω , κ , tree length and codon frequencies is effectively explained by a reduced size of the state space of the continuous-time Markov chain. Aggregation is thus all the more effective when sequence data are biased or when analyses contain closely related species, which is probably the case for many real multiple sequence alignments.

4.2 Branch-site model

Given the small speedup that we obtained for the aggregation on fixed positions, we implemented only the full aggregation mode for the branch-site model in FastCodeML. We then compared this new implementation with the standard FastCodeML. We see a slight increase in both false positives and true positives with aggregation (Tables 1, 2). Overall, ROC curves show that the performance of

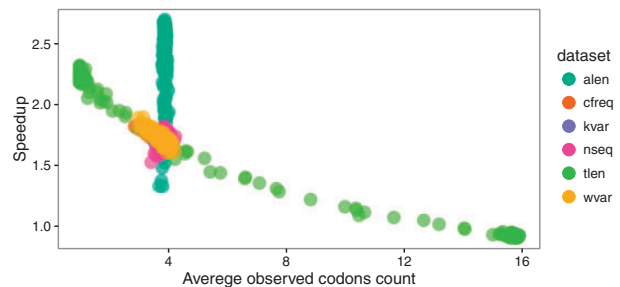


Fig. 8. Speedup versus average codon count for M0 model. Each point represents one simulated alignment, dataset indicated by color. For the dataset with changing alignment length (alen), variation in speedup does not depend on the observed codon count (which does not vary significantly), but longer alignments lead to higher speedup, see Figure 5

Table 1. Statistical performance of FastCodeML in normal and aggregated modes on simulated data

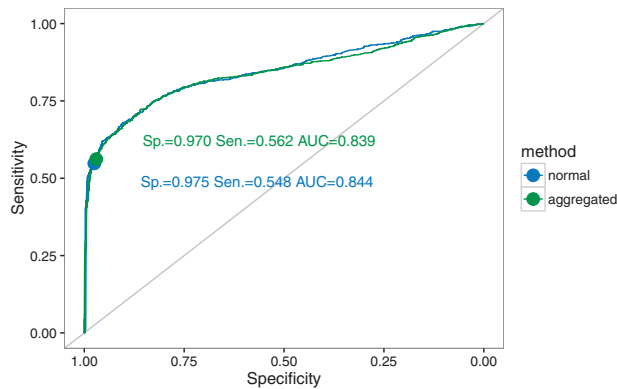
Mode	True positives	True negatives	False positives	False negatives
normal	551	973	27	449
aggregated	562	970	30	438

Numbers in the cells correspond to the number of performed tests. A single branch was tested per tree.

Table 2. Statistical performance of FastCodeML on the simulated dataset

A	Selection detected (aggregated)		B	Selection detected (aggregated)			
	-	+		-	+		
Selection detected (normal)	-	963	12	Selection detected (normal)	-	429	23
	+	7	18		+	9	539

Numbers in the cells correspond to the number of performed tests. A single branch was tested per tree. A) Without positive selection; B) With positive selection.

**Fig. 9.** ROC (receiver operating characteristic) curves for FastCodeML in full likelihood and aggregated likelihood modes for the branch-site model simulations. Specificity, sensitivity and area under curve (AUC) indicated

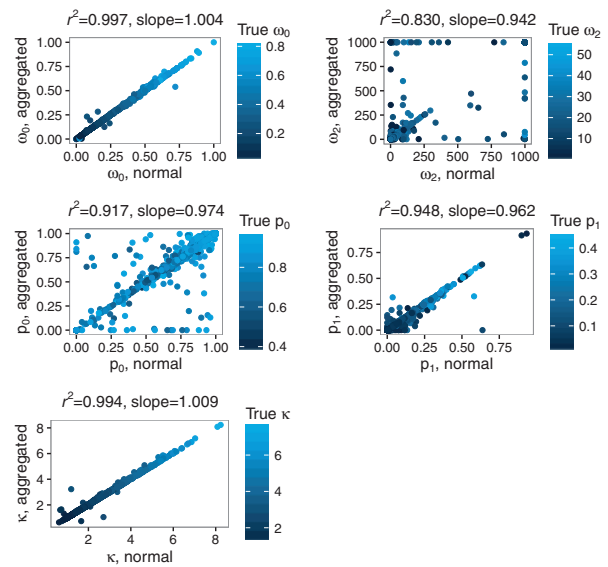
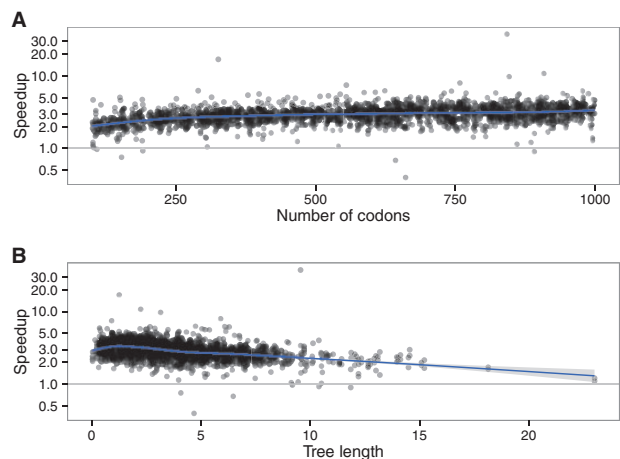
FastCodeML in aggregation mode is similar to the full likelihood mode (Fig. 9). Thus any errors in estimation under aggregation seem to have very little impact on the Likelihood ratio test (LRT) used to test for the presence of positive selection with the branch-site model.

For ω_0 , κ and p_1 , Pearson's correlation coefficients between aggregated and full likelihood estimates are 0.9986, 0.9969 and 0.9735, respectively (Fig. 10). A lower correlation is observed for p_0 and ω_2 (0.9578 and 0.9109, respectively). Yet, these correlations are much higher than those obtained between the full likelihood estimate and simulated values: 0.35 for p_0 and 0.20 for ω_2 .

As with the M0 model, speedup is mostly affected by sequence length and tree length (Fig. 11) through their effects on observed codon counts (Supplementary Fig. S19). We reached a maximum speedup of 4.4 fold per likelihood computation for the branch-site model.

The extended branch-site model violates several assumptions of the branch-site model (homogeneous synonymous rate, fixed ω_0 and ω_2 values). In those cases, where data is more complex than the model, the performance of the aggregated mode becomes slightly worse compared to the full likelihood mode S23, although it remains very close (AUC 0.812 versus 0.818).

Finally, we used FastCodeML in normal and aggregated modes on a real dataset from Primates (Tables 3, S3). After correction for multiple testing (false discovery rate cutoff 0.05), 20 branches were identified to be under positive selection using full likelihood computations and 18 using aggregation, with 13 branches in common. We did not encounter multiple branches detected for an individual tree. The predictions are consistent between the two methods in 99.97% of the cases, which is higher than the consistency of 97.45% for the

**Fig. 10.** Correlation between aggregated and non-aggregated parameter estimates for the branch-site model**Fig. 11.** Effect of (A) alignment length and (B) tree length on the speedup, branch-site model**Table 3.** Statistical performance of FastCodeML on the Primates dataset

	Selection detected (aggregated)		
	-	+	
Selection detected (normal)	-	77576	280
	+	190	1114

Detected selection in normal and aggregated modes of FastCodeML. Numbers in the cells correspond to the number of performed tests. Every non-terminal branch was tested.

simulated data (Table 2). Aggregation gives a median speedup of 2.7 on this real dataset, confirming that real data can be sufficiently biased to make aggregation quite efficient.

5 Discussion

We propose state aggregation as a technique for speeding up the computation of likelihood in a phylogenetic context. By reducing the size of the state space of the Markov process, aggregation accelerates the phase of tree pruning during the likelihood computation and, in some cases, the eigendecomposition of the transition rate matrix. We show that aggregation can be applied to the likelihood calculation of two of the most commonly used codon models. It can also be used for other types of models (see below), in both maximum likelihood and Bayesian frameworks.

The speedup for codon models depends on the alignment length and the observed codon counts, the latter being mostly affected by the tree length (Figs. 8, S19).

These effects are especially strong with the M0 model, because the likelihood optimizer uses a fixed number of iterations. A similar trend is observed with a variable number of iterations, but with increased stochasticity (Supplementary Fig. S20). In general, state aggregation does not appear to have a systematic influence on the total number of iterations (Supplementary Fig. S21). The total runtime is therefore proportional to the likelihood computation time.

Alignment length and observed codon counts have a similar effect on a speedup for the branch-site model. We see more explicitly the dependency if we normalize for the number of likelihood function computations (Supplementary Fig. S22).

The most time consuming stages of the likelihood computation are matrix exponentiation and tree pruning. FastCodeML uses highly optimized algorithms to do matrix exponentiation (Schabauer *et al.*, 2012) and state aggregation improves the time to perform the tree pruning steps of the likelihood calculations (Supplementary Fig. S1A, B).

While the dependency of the speedup on the alignment length and the codon counts make intuitive sense, we can understand it in more details by considering the steps of the likelihood computation. Let us consider the computation time of the total likelihood (Supplementary Fig. S1A):

$$T_{\text{full}} = t_{\text{eigen}} + Kt_{\text{exp}} + NK61t_{\text{prun}},$$

where t_{eigen} is the time to decompose the instantaneous rate matrix, t_{exp} is the time to exponentiate the rate matrix for each internal node, t_{prun} is the time to compute the partial likelihood vector per internal node per position, K is the number of internal nodes and N is the number of positions in the alignment. The number of states is 61 for Markov chains modeling codon sequences.

Similarly for the state aggregation (Supplementary Fig. S1B):

$$T_{\text{aggr}} = t_{\text{eigen}} + Kt_{\text{exp}} + NKt_{\text{agg}} + NKMt_{\text{prun}},$$

where t_{agg} is the matrix aggregation time per internal node per position, and 61 is replaced by M , the number of states after aggregation. For a given branch and site combination, the aggregation time is comparable to the time spent computing a single element of the partial likelihood vector. In the full mode, 61 elements of the vector should be computed. The gain of computing time observed with the aggregation methods comes from the need to do a single aggregation step, which is fast, followed by the computation of M ($M < 61$) vector elements.

Aggregation speedup is thus:

$$\text{Speedup} = \frac{T_{\text{full}}}{T_{\text{aggr}}} = \frac{t_{\text{eigen}} + Kt_{\text{exp}} + NK61t_{\text{prun}}}{t_{\text{eigen}} + Kt_{\text{exp}} + NKt_{\text{agg}} + NKMt_{\text{prun}}}.$$

Generally performance is limited by eigendecomposition and pruning, so we can approximate speedup as:

$$\text{Speedup} \approx \frac{t_{\text{eigen}} + NK61t_{\text{prun}}}{t_{\text{eigen}} + NKt_{\text{agg}} + NKMt_{\text{prun}}}.$$

This representation gives a clear explanation for the dependency of the speedup on the alignment length and the observed codon counts. Increasing the alignment length causes a weaker effect on the non-accelerated eigendecomposition phase, which results in a more efficient acceleration. In contrast, a higher codon diversity in each alignment position increases the number of states in the aggregated Markov process (M), thus reducing the advantage of the aggregated process relative to the full one.

Not only does aggregation provide diminishing speedup for longer trees (more observed states, larger M), it also introduces a bias in the parameter estimation for extremely long trees. Consequently, for trees longer than 100 expected substitutions per position it is not practical to use state aggregation: biased results would be obtained without any significant speedup. In practice, however, extremely long trees are rare, for example in the Selectome database 99% of the trees has total length below 18 expected substitutions per position.

State aggregation can be applied either to the probability matrix P or to the instantaneous matrix Q (Supplementary Fig. S1B, C). In this work we were focused on applying aggregation to the probability matrix P . In this case (Supplementary Fig. S1B), aggregation is applied after exponentiation and must be performed for every position independently. The performance improvement is therefore achieved during the tree pruning phase. In the case of the matrix Q (Supplementary Fig. S1C), aggregation is applied prior to the exponentiation. This leads to smaller dimensions of P matrices, but eigendecomposition and exponentiation have to be performed for every position independently, since those positions will differ in the states aggregated. Moreover, aggregation of the matrix Q is expected to introduce more bias that will accumulate along the branches. Aggregation performed on the Q matrix will discard differences in substitution trajectories passing through unobserved states. There will thus be an accumulation of the error during both exponentiation and pruning phases. Aggregation done after the exponentiation phase only introduces error during the tree pruning phase. Preliminary results do not show an advantage of aggregating the matrix Q for codon models (not shown). A solution might be to perform a ‘softer’ aggregation on clusters of sites with similar patterns of codons. This would be done by first clustering alignment positions and then producing aggregated instantaneous rate matrices for each cluster. This should diminish the bias and allow to exponentiate a smaller number of Q matrices than for the aggregation per site, while still computing on smaller Q matrices than in non-aggregated mode. It is also possible that aggregation of the Q matrix could be more useful for other types of models, especially those with large instantaneous rate matrices, such as coevolution models (Dib *et al.*, 2014). Finally, a second round of aggregation might be performed after the exponentiation in order to speedup the tree pruning stage (Supplementary Fig. S1D). The computational and statistical performance of such approaches has yet to be investigated.

It is also possible to implement aggregation on a subset of the data only. In our case, we chose an extreme situation and aggregated only the most conserved positions. The result was a large loss in speedup relative to aggregation on all positions without any gain in accuracy. But there might be other cases where aggregation on a subset of data only makes most sense in terms of the cost (accuracy)—benefit (speedup) trade-off. Moreover, there are multiple ways to perform the aggregation itself. Here, we collapsed all of the codons or amino-acids which are not observed at the position of the

alignment. It is also possible to use other approaches to aggregation, e.g. aggregate all the codons reachable by more than a single mutation (or follow the amino-acids similarity properties). Models of amino acid substitutions have been derived from codon-based Markov models by aggregating codons separated by only synonymous substitutions. These models were, however, not built nor evaluated for computational efficiency (Kosiol and Goldman, 2011; Ren et al., 2005; Susko and Roger, 2007; Yang et al., 1998). Less aggressive aggregation shows increase in the accuracy at the price of reduced speedup, although, in our tests, accuracy was already good with the aggressive aggregation.

The combined use of both aggregated and non-aggregated modes in the same analysis could be efficient in several scenarios. First, aggregation could be used during likelihood maximization, but the final likelihood value computed without aggregation, providing a more accurate value. Second, aggregation could be used to obtain a starting point for non-aggregated likelihood maximization. Third, aggregation could be used in a preprocessing step to detect datasets of interest (e.g. gene families with a signal of positive selection). These datasets could then be analyzed with full likelihood to get an accurate estimation of the parameters and model comparison. Finally, aggregation could be used during the burn-in period in a Bayesian approach (e.g. MCMC). There are probably other scenarios where aggregation can provide a faster estimation of likelihood within a more complex analysis.

For the specific case of the branch-site model, we have tested the second scenario of using aggregation as a starting point and we do not obtain a significant speedup (Supplementary Fig. S24).

Aggregation can also have an impact on memory usage. Aggregation on the probability matrix P will reduce the size of the partial likelihood vectors (by a factor of $r = 61/N_{\text{states}}$). Additionally the Q matrix aggregation reduces the size of the P -matrices (by a factor of r^2). On the other hand, the actual improvement strongly depends on the details of the implementation, as partial likelihood vectors and probability matrices can be reused in a number of ways. Our implementation did not focus on the reduction of memory footprint and we thus do not discuss this aspect further.

Obviously, state aggregation in phylogeny and evolution is not limited to the branch-site and M0 codon models. First, it is universally applicable to Markov process-based codon models, such as the commonly used M1a/M2a, M8a/M8 (Wong et al., 2004), aBSREL (Smith et al., 2015), RELAX (Wertheim et al., 2014), or any other GY94 (Goldman and Yang, 1994) or MG94-based (Muse and Gaut, 1994) model. Second, it is not limited to codon models. Given a trade-off between per-position matrix aggregation slowdown and tree pruning speedup, aggregation is unlikely to give a significant performance improvement for models with a small number of states (e.g. nucleotide models). But even for amino acids models we can expect some degree of speedup. In contrast, we expect state aggregation to provide a significant performance improvement for the models with a large number of states, such as amino acid coevolution models that can include up to 400 states (Dib et al., 2014; Yeang and Haussler, 2007).

The aggregation of states in a Markov process is a powerful technique used in a variety of fields including computational biology, such as protein network interaction analysis (Petrov et al., 2012), reaction modeling (Ullah et al., 2012), single molecule photobleaching (Messina et al., 2006), or disease-progression models (Regnier and Shechter, 2013). Its application to phylogenetic models has not been systematically studied, although it has been implemented in some software (Lartillot and Philippe, 2004, e.g. PhyloBayes). This is, to

our knowledge, the first systematic study of state aggregation biases and computational efficiency for molecular evolution.

In conclusion, we demonstrate that state aggregation is a powerful method which improves computational performance of codon-based models, with little cost in accuracy. State aggregation is not limited to codon models, and we expect it to be useful for a large variety of phylogenetic models and methods.

Acknowledgments

We would like to thank Bastien Boussau, Christophe Dessimoz, Stephane Guindon and anonymous referees for useful comments. The computations were performed at the Vital-IT (<http://www.vital-it.ch>) centre for high-performance computing of the Swiss Institute of Bioinformatics.

Funding

This work was supported by the Swiss National Science Foundation [grant numbers CR32I3_143768, IZLRZ3_163872].

Conflict of Interest: none declared.

References

- Aldous, D. and Fill, J.A. (2002). Reversible markov chains and random walks on graphs. Unfinished monograph.
- Cunningham, F. et al. (2015) Ensembl 2015. *Nucleic Acids Res.*, **43**, D662–D669.
- de Koning, A.J. et al. (2012) Phylogenetics, likelihood, evolution and complexity. *Bioinformatics*, **28**, 2989–2990.
- Dib, L. et al. (2014) Evolutionary footprint of coevolving positions in genes. *Bioinformatics*, **30**, 1241–1249.
- Felsenstein, J. (1973) Maximum likelihood and minimum-steps methods for estimating evolutionary trees from data on discrete characters. *Syst. Zool.*, **22**, 240–249.
- Felsenstein, J. (1981) Evolutionary trees from DNA sequences: a maximum likelihood approach. *J. Mol. Evol.*, **17**, 368–376.
- Gladstein, D.S. (1997) Efficient incremental character optimization. *Cladistics*, **13**, 21–26.
- Goldman, N. and Yang, Z. (1994) A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Mol. Biol. Evol.*, **11**, 725–736.
- Goloboff, P.A. (1993) Character optimization and calculation of tree lengths. *Cladistics*, **9**, 433–436.
- Guindon, S. and Gascuel, O. (2003) A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Syst. Biol.*, **52**, 696–704.
- Guindon, S. et al. (2010) New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst. Biol.*, **59**, 307–321.
- Hillston, J. (1995) Compositional Markovian modelling using a process algebra. In: William, J.S. (eds.) *Computations with Markov Chains*, pp. 177–196. Springer, New York.
- Hordijk, W. and Gascuel, O. (2005) Improving the efficiency of SPR moves in phylogenetic tree search methods based on maximum likelihood. *Bioinformatics*, **21**, 4338–4347.
- Huelsenbeck, J.P. et al. (2001) MRBAYES: Bayesian inference of phylogenetic trees. *Bioinformatics*, **17**, 754–755.
- Kemeny, J. and Snell, J. (1983). *Finite Markov Chains: With a New Appendix "Generalization of a Fundamental Matrix"*. Undergraduate Texts in Mathematics. Springer, New York.
- Kosiol, C. and Goldman, N. (2011) Markovian and non-markovian protein sequence evolution: aggregated markov process models. *J. Mol. Biol.*, **411**, 910–923.
- Lartillot, N. (2006) Conjugate Gibbs sampling for Bayesian phylogenetic models. *J. Comput. Biol.*, **13**, 1701–1722.

- Lartillot,N. and Philippe,H. (2004) A Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Mol. Biol. Evol.*, **21**, 1095–1109.
- Messina,T.C. *et al.* (2006) Hidden Markov model analysis of multichromophore photobleaching. *J. Phys. Chem. B*, **110**, 16366–16376.
- Moretti,S. *et al.* (2014) Selectome update: quality control and computational improvements to a database of positive selection. *Nucleic Acids Res.*, **42**, D917–D921.
- Murrell,B. *et al.* (2012) Detecting individual sites subject to episodic diversifying selection. *PLoS Genet.*, **8**, e1002764–e1002764.
- Muse,S.V. and Gaut,B.S. (1994) A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with application to the chloroplast genome. *Mol. Biol. Evol.*, **11**, 715–724.
- Nguyen,L.T. *et al.* (2014) IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.*, msu300.
- Petrov,T. *et al.* (2012) Model decomposition and stochastic fragments. *Electr. Notes Theor. Comput. Sci.*, **284**, 105–124.
- Phillips,M.J. *et al.* (2004) Genome-scale phylogeny and the detection of systematic biases. *Mol. Biol. Evol.*, **21**, 1455–1458.
- Proux,E. *et al.* (2009) Selectome: a database of positive selection. *Nucleic Acids Res.*, **37**, D404–D407.
- Regnier,E.D. and Shechter,S.M. (2013) State-space size considerations for disease-progression models. *Stat. Med.*, **32**, 3862–3880.
- Ren,F. *et al.* (2005) An empirical examination of the utility of codon-substitution models in phylogeny reconstruction. *Syst. Biol.*, **54**, 808–818.
- Rodrigue,N. *et al.* (2008) Uniformization for sampling realizations of Markov processes: applications to Bayesian implementations of codon substitution models. *Bioinformatics*, **24**, 56–62.
- Ronquist,F. (1998) Fast Fitch-parsimony algorithms for large data sets. *Cladistics*, **14**, 387–400.
- Rubinstein,N.D. *et al.* (2011) Evolutionary models accounting for layers of selection in protein-coding genes and their impact on the inference of positive selection. *Mol. Biol. Evol.*, **28**, 3297–3308.
- Schabauer,H. *et al.* (2012) SlimCodeML: an optimized version of CodeML for the branch-site model. In: *2012 IEEE 26th International Parallel and Distributed Processing Symposium Workshops & PbD Forum*, pp. 706–714. IEEE.
- Smith,M.D. *et al.* (2015) Less is more: an adaptive branch-site random effects model for efficient detection of episodic diversifying selection. *Mol. Biol. Evol.*, **32**, 1342–1353.
- Stamatakis,A. (2014) RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics*, **30**, 1312–1313.
- Stamatakis,A. *et al.* (2005) RAxML-III: a fast program for maximum likelihood-based inference of large phylogenetic trees. *Bioinformatics*, **21**, 456–463.
- Stamatakis,A.P. *et al.* (2004). New fast and accurate heuristics for inference of large phylogenetic trees. In: *Parallel and Distributed Processing Symposium, 2004. Proceedings. 18th International*, p. 193. IEEE.
- Storey,J.D. *et al.* (2004) Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)*, **66**, 187–205.
- Susko,E. and Roger,A.J. (2007) On reduced amino acid alphabets for phylogenetic inference. *Mol. Biol. Evol.*, **24**, 2139–2150.
- Swofford,D.L. and Olsen,G.J. (1990). Phylogeny reconstruction. In: Hillis,D. and Moritz,C. (eds.) *Molecular Systematics*, pp. 411–501. Sinauer Associates, Sunderland, MA.
- Ullah,G. *et al.* (2012) Simplification of reversible Markov chains by removal of states with low equilibrium occupancy. *J. Theor. Biol.*, **311**, 117–129.
- Valle,M. *et al.* (2014) Optimization strategies for fast detection of positive selection on phylogenetic trees. *Bioinformatics*, btt760.
- Vera-Ruiz,V.A. *et al.* (2014) Statistical tests to identify appropriate types of nucleotide sequence recoding in molecular phylogenetics. *BMC Bioinformatics*, **15**, S8.
- Vilella,A.J. *et al.* (2009) EnsemblCompara GeneTrees: complete, duplication-aware phylogenetic trees in vertebrates. *Genome Res.*, **19**, 327–335.
- Wertheim,J.O. *et al.* (2014) Relax: detecting relaxed selection in a phylogenetic framework. *Mol. Biol. Evol.*, msu400.
- Wong,W.S. *et al.* (2004) Accuracy and power of statistical methods for detecting adaptive evolution in protein coding sequences and for identifying positively selected sites. *Genetics*, **168**, 1041–1051.
- Yang,Z. (2007) PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.*, **24**, 1586–1591.
- Yang,Z. *et al.* (1998) Models of amino acid substitution and applications to mitochondrial protein evolution. *Mol. Biol. Evol.*, **15**, 1600–1611.
- Yeang,C.H. and Haussler,D. (2007) Detecting coevolution in and among protein domains. *PLoS Comput. Biol.*, **3**, e211.
- Zhang,J. *et al.* (2005) Evaluation of an improved branch-site likelihood method for detecting positive selection at the molecular level. *Mol. Biol. Evol.*, **22**, 2472–2479.