

July 30, 2020

A Bayesian mutation-selection framework for detecting site-specific adaptive evolution in protein-coding genes

Nicolas Rodrigue¹, Thibault Latrille² and Nicolas Lartillot²

¹Department of Biology, Institute of Biochemistry, and School of Mathematics and Statistics,
Carleton University, Ottawa, Canada

²Université de Lyon, Université Lyon 1, CNRS; UMR 5558, Laboratoire de Biométrie et Biologie Évolutive, F-69622 Villeurbanne, France

Running head: Detecting site-specific adaptation with mutation-selection models

Keywords: Nearly-neutral evolution; fitness landscape; Dirichlet process; Markov chain Monte Carlo.

Correspondence: Nicolas Rodrigue

209 Nesbitt Biology Building,

1125 Colonel By Drive Ottawa, Ontario, CANADA

K1A 0C6

`nicolas.rodrigue@carleton.ca`

tel: +1 613 520 2600 x 4194

Abstract

In recent years, codon substitution models based on the mutation-selection principle have been extended for the purpose of detecting signatures of adaptive evolution in protein-coding genes. However, the approaches used to date have either focused on detecting global signals of adaptive regimes—across the entire gene—or on contexts where experimentally derived site-specific amino acid fitness profiles are available. Here, we present a Bayesian site-heterogeneous mutation-selection framework for site-specific detection of adaptive substitution regimes given a protein-coding DNA alignment. We offer implementations, briefly present simulation results, and apply the approach on a few real data sets. Our findings suggest that the new approach shows greater sensitivity than traditional methods, with a reasonably low false-positive rate. Finally, we and outline a potential research program with the framework.

Introduction

Codon substitution models (Goldman and Yang, 1994; Muse and Gaut, 1994) are among the important modern tools used for uncovering potential signals of molecular adaptation from protein-coding gene alignments. One set of broadly used models focuses on estimating the ratio of rates for non-synonymous (dN) and synonymous (dS) substitutions. These models introduce a multiplicative parameter, denoted ω , to entries in a codon substitution matrix corresponding to nonsynonymous events. Because ω is the only distinction between the rate specification of nonsynonymous and synonymous events, it directly corresponds to $\omega = dN/dS$.

Fitting a model with a single (global) nonsynonymous multiplicative parameter almost always leads to $\omega < 1$ (Yang, 2006), given the pervasive purifying selection that operates at most codon sites over most of evolutionary history. Many efforts were thus made to develop codon substitution models with distributions of ω values across sites and/or across the branches of a phylogeny (reviewed in Yang, 2019). A common objective of such developments is to uncover specific sites having evolved under an adaptive regime (i.e., with $\omega > 1$), perhaps along a particular branch of the phylogeny.

Meanwhile, another set of codon substitution models was proposed, with a focus on accounting for purifying selection at the amino acid level in a site-heterogeneous manner (Halpern and Bruno, 1998). Having nucleotide-level parameters controlling a mutational process, and amino acid fitness profiles controlling selection, they have come to be known as *mutation-selection* models (e.g., Yang and Nielsen, 2008; Rodrigue et al., 2010). In these models, the dN/dS ratio is not explicitly parameterized. Instead, it is an emerging quantity, induced by the interplay between mutation, selection, and drift. Spielman and Wilke (2015) have shown how to calculate the dN/dS induced by the mutation-selection framework—which we denote ω_0 (Rodrigue and Lartillot, 2017)—and found that, under specific conditions (i.e., a substitution process at equilibrium, without selection

on synonymous variants), it is always true that $\omega_0 \leq 1$, as expected from a model focused on purifying selection.

In the last few years, the mutation-selection framework has been extended for the purpose of detecting genes having evolved under an adaptive regime, in either a global (Rodrigue and Lartillot, 2017) or site-specific (Bloom, 2017) manner. Like their traditional predecessors, these recent mutation-selection models introduce a multiplicative parameter on nonsynonymous rates. However, because amino acid profiles are also involved in modulating nonsynonymous rates, such a multiplicative parameter—which we denote as ω_* (Rodrigue and Lartillot, 2017)—cannot be interpreted as the dN/dS ratio; we chose to emphasize this distinction with an asterisk in the notation. Given that the mutation-selection formulation itself induces a certain dN/dS ratio, ω_0 , the net overall dN/dS ratio, ω , can be thought of as $\omega = \omega_0 \times \omega_*$, which can be rearranged to $\omega_* = \omega/\omega_0$. The latter equation helps clarify the interpretation of ω_* as a measure of the deviation in nonsynonymous rates from the expectation under the pure mutation-selection equilibrium; in particular, $\omega_* > 1$ indicates that nonsynonymous rates are higher than expected, even though they might not be so high as to lead to $\omega > 1$.

New Approaches

Here, we conduct a first exploration of a Bayesian mutation-selection model with site-heterogeneous amino acid fitness profiles and site-heterogeneous ω_* values. The Bayesian nature of the model qualifies it as a *random-effects* approach, in contrast to the *fixed-effects* approach sometimes utilized in maximum-likelihood versions of mutation-selection models (Halpern and Bruno, 1998; Holder et al., 2008; Tamuri et al., 2014; Bloom, 2017).

Results and Discussion

Models With Global ω or ω_*

We first contrasted the difference in behaviour between a traditional codon substitution model inspired from Muse and Gaut (1994), with a global ω parameter (a traditional model we denote MG-M0, described in detail in the Materials and Methods section), and a mutation-selection model with a Dirichlet process prior on amino acid profiles across sites, and a global ω_* parameter (a model presented in Rodrigue and Lartillot, 2017, which we denote here as MutSel-M0 $_*$, and also described in the Materials and Methods section). Figure 1 shows results of the two models on data generated through a simulation approach explicitly allowing for fluctuating selection at some sites; for these sites, amino acid fitness profiles change along the branches of the phylogeny, as described in Rodrigue and Lartillot (2017). Such a system mimics an adaptive substitution process, where the simulated substitution history tracks a changing amino acid fitness optimum along the branches of the tree, and thus accrues more nonsynonymous substitutions than expected under a pure nearly-neutral regime (i.e., mutation-selection balance). An important distinction with Rodrigue and Lartillot (2017) is that here the simulated data set contains 90% of codon sites generated under a pure nearly-neutral mutation-selection formulation (Rodrigue et al., 2010) and 10% generated under the adaptive regime. We produced alignments of 300 codons in length, repeating the simulation thrice with different sets of empirically inferred amino acid profiles (see Lowe and Rodrigue, 2020, and the Materials and Methods section).

Results under the traditional MG-M0 model (red) reflect the overall purifying selection governing most of the data-generating processes, with a posterior mean ω values at 0.14, 0.15, 0.13 in replicates in panels 1A, 1B, and 1C. The fact that 10% of sites were produced under an adaptive regime is underwhelming to the MG-M0 model, and indeed little is generally expected of it in

practice. Results under the MutSel-M0_{*} model (blue) show a posterior distribution for ω_* situated above 1, with $p(\omega_* > 1 \mid D) \geq 0.99$ (where D refers to the data set) for the first two replicates (fig. 1A and 1B); indeed, the second replicate has a posterior mean that surpasses 2. For the third replicate, we find a slightly lower probability, at $p(\omega_* > 1 \mid D) \sim 0.93$.

Previous studies (Rodrigue and Lartillot, 2017; Lowe and Rodrigue, 2020) have shown that a simulation conducted with 100% of sites under the pure nearly-neutral mutation-selection formulation leads to a posterior distribution of ω_* situated around 1; here, however, 10% of sites have evolved with higher than expected nonsynonymous rates, which pulls the distribution to the right. Already with the use of single additional parameter ω_* , the mutation-selection framework allows us to detect adaptation where the traditional framework with a single ω parameter would not.

Figure 2 shows results of these models on a hand-full of real alignments. Panel 2A displays the results on the well-known β -GLOBIN alignment sampled across 17 vertebrates (Yang et al., 2000). As in the simulation, the MG-M0 model indicates that $\omega < 1$. In contrast, under MutSel-M0_{*}, the posterior mean of ω_* is around 1.8, with a high posterior probability in favor of a value greater than 1, $p(\omega_* > 1 \mid D) > 0.99$, indicating the presence of adaptive evolution in this gene. As suggested by the simulation experiment presented above, adaptive evolution on even a relatively small fraction of the sites of the gene could be sufficient to induce such a rightward shift in the posterior distribution of ω_* .

Panel 2B displays results on an alignment of the alcohol dehydrogenase (ADH) gene sampled across 23 species of *Drosophila*. Here again, the MG-M0 model indicates that $\omega < 1$, with a posterior mean ~ 0.13 . In contrast, with the MutSel-M0_{*} model we find a posterior mean $\omega_* \sim 1.2$, and $p(\omega_* > 1 \mid D) > 0.99$. No previous phylogenetic approach has ever found signals of adaptive evolution in this gene, in spite of the fact that population-genetic approaches have long suggested adaptation in many instances (e.g., McDonald and Kreitman, 1991; Matzkin and Eanes, 2003;

Matzkin, 2004). While a specific scenario of ADH adaptation in specific species has been refuted by Siddiq and Thornton (2019), their study nonetheless provides strong experimental evidence of major fitness effects of some mutations suggesting adaptive opportunities across *Drosophila*.

The four lower panels of figure 2 (2C, 2D, 2E, and 2F) show results on four genes sampled across placental mammals. Again, $\omega < 1$ in all four genes, whereas ω_* is either around 1, or slightly below, which does not suggest adaptive evolution in these genes across placental mammals. Previous simulations studies have pointed to epistasis (Rodrigue and Lartillot, 2017) or weak evolutionary signal (Lowe and Rodrigue, 2020) as potential reasons for $\omega_* < 1$. The conditions thus tend to make the model conservative in the detection of adaptive regimes.

Models with heterogeneous ω or ω_*

In spite of the potential of the MutSel-M0_{*} model—able to capture relatively subtle signals of adaptive evolution—it still does not directly allow us to pinpoint which sites are most responsible for such signals. This is one of the motivations of *site-models*. Classical site-models (Nielsen and Yang, 1998; Yang et al., 2000; Yang and Swanson, 2002) consider alignment sites as having been produced from a distribution of possible ω values. They are typically used in the context of an empirical Bayes approach for identifying sites with a strong statistical support for a $\omega > 1$; and they are more efficient at detecting positive selection than the simple MG-M0 model with a single ω for all sites. For instance, they do find sites under positive selection in the case of the β -GLOBIN gene (detailed below, but also see Yang et al., 2000). On the other hand, site-models might still miss those sites under weaker positive selection. In particular, an adaptive regime at a site could be sufficiently strong to increase the dN/dS ratio, but not to the point of driving it well above 1. In other words, at least in their current version, these models might present the same limitation as the classical MG-M0 model, as compared with MutSel-M0_{*} model, although now at the level

of the single site. This in turn suggests that the rationale of estimating ω_* in the context of a mutation-selection model should be explored not just globally over the whole gene (Rodrigue and Lartillot, 2017), but as a distribution across sites of the gene (Bloom, 2017).

To illustrate this point, and for simplicity here, we work with the classical MG-M3 model, inspired from Muse and Gaut (1994) and Yang et al. (2000), which invokes a finite mixture of three ω values—with their respective weights—jointly estimated with all other parameters from the data. We also study a new model referred to as MutSel-M3_{*}, which is built from a finite mixture of three ω_* values across sites, and respective weights, combined with the distribution of amino acid profiles across sites (the infinite mixture Dirichlet process), and global mutational parameters. The two forms of across-site heterogeneity are independent in the model construction, in that each site draws its amino acid profile and its ω_* independently from the two corresponding mixtures.

As a verification, figure 4 shows the results under the MG-M3 (red) and MutSel-M3_{*} (blue) models on three simulated data sets, this time generated entirely under the pure mutation-selection framework (i.e., no adaptive regimes were included). In accordance with the simulation, no sites have high probabilities of having $\omega_* > 1$ (or $\omega > 1$). One of the simulations (fig. 4B) suggests that the mixture of ω_* values may tend to be under-estimated, which would again tend to make the model conservative vis-à-vis inferences of adaptive evolution.

Figure 5 shows the results on the three simulated data sets studied in figure 1 (i.e., with 10% of sites simulated with an adaptive regime). The panels include vertical marks at the top, showing the 30 codon sites simulated under adaptive regimes. Sites evolving under an adaptive regime tend to accrue more nonsynonymous substitutions than under a nearly-neutral regime, which would shift ω_* to the right of the unit. With a threshold posterior probability of 0.95 for $p(\omega_* > 1 \mid D)$, the MutSel-M3_{*} model correctly identifies 23/30 sites (76%), calls 1 false positive, and misses 7 sites for the first and second replicates, whereas for the third replicate it correctly identifies 20/30,

with no false positives. Of note, a single false positive out of 24 discoveries, using a threshold of 0.95, corresponds to an accuracy of $\sim 96\%$, thus suggesting that the posterior probabilities are well-calibrated, correctly reflecting our actual rate of true discovery. The MG-M3 models identifies no sites at this threshold, although the plot suggests that it nonetheless faintly detects some adaptive signal. Interestingly, the sites leading to false positives under the MutSel-M3_{*} model also tempt the MG-M3 model; the simulations are stochastic processes, and can, from time to time, accumulate a disproportionately high number of nonsynonymous substitutions, even when the configuration of the simulating model is one of pure mutation-selection balance. In other words, this false positive may not come about solely as a result of a problem with MutSel-M3_{*} model itself, but rather, at least partly, from a chance occurrence in the simulation. It is particularly noteworthy that some of the sites correctly identified by MutSel-M3_{*} show virtually no signal under MG-M3 (e.g., sites 52, 103, 285 in the first replicate, panel 5A). All of the sites simulated with an adaptive regime, but missing the 0.95 threshold under MutSel-M3_{*}, nonetheless have relatively high probabilities of having $\omega_* > 1$. Overall, the MutSel-M3_{*} model seems to have considerably greater sensitivity than the traditional-style MG-M3, at the cost of a mildly increased risk of false positives.

Figure 6 displays the results obtained from analyzing the six real data sets mentioned above with the MG-M3 and MutSel-M3_{*} models. For the β -GLOBIN alignment (fig. 6A), our Bayesian version of the classic MG-M3 model leads to the same set of sites identified with these traditional models in the maximum likelihood contexts (Yang et al., 2000): at the 95% threshold, the sites are 7, 11, 42, 48, 50, 54, 67, 85, and 123. Under the MutSel-M3_{*} model, these same sites are also found, and the following three are added: 10, 74, and 84. The complete lists of sites identified at different thresholds are given in table 1. It is interesting to note that the MG-M3 model found $p(\omega > 1 \mid D) = 0.381$ for site 10, $p(\omega > 1 \mid D) = 0.244$ for site 74, and $p(\omega > 1 \mid D) = 0.074$ for site 84. These last three sites, and site 84 in particular, yield results compatible with the interpretation

of having evolved under a mild adaptive regime, of changing amino acid fitness profiles over time, leading to an increase in nonsynonymous rate; the increase is not to the point where $\omega > 1$ at a site in question, although it is enough for $\omega_* > 1$. Sites 10 and 74 are known to be involved in oxygen affinity, which could indeed make them a target of adaptive evolution.

Results of the analysis of ADH (fig. 6B) suggest several sites under adaptive evolution under the MutSel-M3_{*} model, whereas the MG-M3 yields posterior probabilities of $\omega > 1$ at all sites that are numerically indistinguishable from 0. Given that most studies suggesting adaptation in this gene have relied on population-genetic methodologies, which pool the statistics across all sites, a comparison of sites uncovered by the MutSel-M3_{*} model with previous results is not possible. Much more work will be required to determine the plausibility of results on this gene.

Our analysis of the mammalian gene VWF also suggests several sites with adaptive signatures under the MutSel-M3_{*} model, and none under the MG-M3 model. A previous study, utilizing branch heterogeneous models, has suggested adaptive evolution conferring venom resistance to oposoms that prey on pitvipers (Jansa and Voss, 2011). Moreover, variants of this gene have been found to have dramatic effects on its own expression levels in mice (Lemmerhirt et al., 2006). Again, however, more work is required to assess these results.

Of the remaining mammalian gene alignments studied with the MutSel-M3_{*} model, two suggest very few sites having evolved under adaptive regimes (ADORA3 and S1PR1, in fig. 6A and 6C respectively), and one (RBP3, fig. 6B) with none. The traditional MG-M3 model suggests no sites under adaptive evolution.

Future directions

The traditional codon models based on ω have become increasingly well understood thanks to decades of empirical applications and simulation studies. We believe that a similar project should

be considered within the mutation-selection framework.

In particular, more sophisticated simulations, at a larger scale and with increasingly realistic conditions, will be needed to better characterize the behavior of the proposed mutation-selection-based approaches to potential model violations, including context-dependent mutational features (e.g., Laurin-Lemay et al., 2018), uneven codon usage (Yang and Nielsen, 2008), variable effective population sizes over the phylogeny (Platt et al., 2018), and epistatic effects (Pollock et al., 2012; Shah et al., 2015). Some of these model violations have been shown to decrease estimates of ω_* (e.g., epistasis, in Rodrigue and Lartillot, 2017), while others may very well lead to increased estimates. Suppose, for instance, a context with uneven codon usage, known to be pronounced and be highly variable across *Drosophila* (Powell and Moriyama, 1997), such that the codon TTG is virtually the only one used to encode leucine, and GTG is strongly favored to encode valine. Also suppose that leucine and valine are of equivalent fitness at a given site. In such a context, nonsynonymous substitutions between TTG and GTG accumulate more readily than synonymous substitutions. If such features were to be present to a high extent could mislead the MutSel-M3* model to infer $\omega_* > 1$, in this case suggesting adaptive evolution where the regime is in fact one of strict purifying selection on codon usage. Other violations on the assumption of synonymous neutrality and homogeneity, known to be false in several contexts (Wisotsky et al., in press), could have similar effects.

Also, a more detailed examination, ideally combined with experimental corroborations, of the sites uncovered by the model is pressing, and would ideally be based on far more than the handful of data sets of the present study. This would help build our empirical understanding how the model behaves in different contexts. We hope to apply the model on a few thousand genes from the OrthoMamm database (Scornavacca et al., 2019) in a first step, before broader applications across varied taxonomic contexts. It will be important to contrast results with those based on other

methods (Moutinho et al., 2019).

While we have outlined the modeling strategy with a three-component finite mixture of ω_* values, in combination with a Dirichlet process prior on amino acid profiles, many other possibilities could be considered: various parametric families on ω_* (as did Yang et al., 2000, with ω), non-parametric approaches on ω_* (as proposed for ω by Huelsenbeck et al., 2006), grids of predetermined ω_* values (in the spirit of Murrell et al., 2013), along with similar choices on modeling amino acid fitness heterogeneity (e.g., Rodrigue et al., 2010; Rodrigue, 2013; Rodrigue and Lartillot, 2014), and the potentially complex interactions between the numerous combinations. We propose these modeling ideas in two independent software packages (see below). One of our Markov chain Monte Carlo implementations can run under fixed topology as well as sample over trees, and thus enable studies of the impact of phylogenetic uncertainty in inferences of adaptive evolution, utilizing both traditional and mutation-selection codon substitution models; this also suggests more extensive studies on the potential of such models for phylogenetic inference *per se*. Another implementation we offer lends itself to integrative modeling objectives, with a wide suite of potential research avenues utilizing the mutation-selection-based approaches. Foreseeable directions with the latter implementation include modeling the evolution of effective population size over the phylogeny, along with joint inferences of continuous-trait evolution, as formalized by Lartillot and Poujol (2011).

Materials and methods

Data

Add data list and brief descriptions. For convenience, all data sets (simulated and real) are included in a supplementary information file.

Substitution models

Write out MG-M0 and MutSel-M0_{*} substitution matrix entries.

Priors

List all priors, mentioning the difference between the chronograms in BayesCode versus exponential prior and branch lengths in PB-MPI-2.

Implementations

The models presented have been implemented in an experimental version (2) of PhyloBayes-MPI (<https://github.com/bayesiancook/pbmpi2>), allowing for a joint sampling across parameter space, auxiliary variables, and tree topology space. We have also implemented the models in a new software called BayesCode, which is focused on integrative comparative methods under fixed topology (<https://github.com/bayesiancook/bayescode>). With the PhyloBayes-MPI implementation, running an MCMC sampling from the posterior distribution under MutSel-M3_{*} model, on the β -globin gene (originally obtained from Z. Yang's website at <http://abacus.gene.ucl.ac.uk/ziheng/data/YNGP2000data.tgz>), can be called with the following command: `mpirun -np 4 ./pb_mpi -mutsel -freeomega -omegafinite -nomega 3 -d bglobin.phy -T bglobin.tre -s -sbdp -x 1 1100 mutselm3`, which will iterate over 1100 cycles, with each cycle having a complex set of many update mechanisms, and with computations parallelized over 4 cores. Note that the option `-T bglobin.tre` specifies a sampling under a fixed tree topology (as specified in the `bglobin.tre` file), but using a lowercase `-t` will only use the provided tree as a starting point, and removing the option altogether will start from a random topology, with the MCMC run performing updates on the tree, along with other parameters. Once the run complete, the command `./readpb_mpi -siteomegagto -x 100 mutselm3` will produce a file of site-specific posterior prob-

abilities of $\omega_* > 1$, skipping over the first 100 draws from the sample. The MG-M3 model can be obtained by invoking a finite-mixture on amino acid profiles (rather than the Dirichlet process), and forcing a single, flat profile of amino acids (which thus cancel out of the mutation-selection formulation); this is done by replacing the `-sbdp` option with `-catfix uniform`. The same follow-up command as before can be used to produce the site-specific posterior probabilities of $\omega > 1$. MG-M0 and MutSel-M0_{*} can be obtained by setting `-nomega 1`. With the BayesCode implementation, the MutSel-M3 model can be called with the following command: `./mutselomega -a bglobin.phy -t bglobin.tre --freeomega --omegancat 3 mutselm3`. Adding the option `--flatfitness` will yield the MG-M3 model, and setting `--omegancat 1` can be used to obtain the MutSel-M0_{*} and MG-M0 models.

Acknowledgements

The work was funded by the Natural Sciences and Engineering Research Council of Canada (NR), Carleton University (NR), and ... ask Nico and Thibault if they'd like to mention anything here.

References

- Bloom, J. D. 2017. Identification of positive selection in genes is greatly improved by using experimentally informed site-specific models. *Biol. Direct* 12:1.
- Goldman, N., and Z. Yang. 1994. A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Mol. Biol. Evol.* 11:725–736.
- Halpern, A. L., and W. J. Bruno. 1998. Evolutionary distances for protein-coding sequences: modeling site-specific residue frequencies. *Mol. Biol. Evol.* 15:910–917.
- Holder, M. T., D. J. Zwickl, and C. Dessimoz. 2008. Evaluating the robustness of phylogenetic

- methods to among-site variability in substitution processes. *Phil. Tran. R. Soc. B* 363:4013–4021.
- Huelsenbeck, J. P., S. Jain, S. W. D. Frost, and S. L. K. Pond. 2006. A Dirichlet process model for detecting positive selection in protein-coding DNA sequences. *Proc. Natl. Acad. Sci. USA* 103:6263–6268.
- Jansa, S. A., and R. S. Voss. 2011. Adaptive evolution of the venom-targeted vwf protein in opossums that eat pitvipers. *PLoS One* 6:e20997.
- Lartillot, N., and R. Poujol. 2011. A phylogenetic model for investigating correlated evolution of substitution rates and continuous phenotypic characters. *Mol. Biol. Evol.* 28:729–744.
- Laurin-Lemay, S., H. Philippe, and N. Rodrigue. 2018. Multiple factors confounding phylogenetic detection of selection on codon usage. *Mol. Biol. Evol.* 35:1463–1472.
- Lemmerhirt, H. L., J. A. Shavit, G. G. Levy, S. M. Cole, J. C. Long, and D. Ginsburg. 2006. Enhanced VWF biosynthesis and elevated plasma VWF due to a natural variant in the murine Vwf gene. *Blood* 108:3061–3067.
- Lowe, C., and N. Rodrigue. 2020. Detecting adaptation from multi-species protein-coding dna sequence alignments alignments. *Phylogenetics in the Genomic Era* 4–5.
- Matzkin, Luciano M. 2004. Population Genetics and Geographic Variation of Alcohol Dehydrogenase (Adh) Paralogs and Glucose-6-Phosphate Dehydrogenase (G6pd) in *Drosophila mojavensis*. *Mol. Biol. Evol.* 21:276–285.
- Matzkin, Luciano M., and Walter F. Eanes. 2003. Sequence variation of alcohol dehydrogenase (adh) paralogs in cactophilic *drosophila*. *Genetics* 163:181–194.

- McDonald, J. H., and M. Kreitman. 1991. Adaptive protein evolution at the *adh* locus in *drosophila*. *Nature* 351:652–654.
- Moutinho, A. F., F. F. Trancoso, and J. Y. Dutheil. 2019. The impact of protein architecture on adaptive evolution. *Mol. Biol. Evol.* 36:2013–2028.
- Murrell, Ben, Sasha Moola, Amandla Mabona, Thomas Weighill, Daniel Sheward, Sergei L Kosakovsky Pond, and Konrad Scheffler. 2013. Fubar: a fast, unconstrained Bayesian approximation for inferring selection. *Mol. Biol. Evol.* 30:1196–1205.
- Muse, S. V., and B. S. Gaut. 1994. A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with application to the chloroplast genome. *Mol. Biol. Evol.* 11:715–724.
- Nielsen, R., and Z. Yang. 1998. Likelihood models for detecting positively selected amino acid sites and applications to the HIV-1 envelope gene. *Genetics* 148:929–936.
- Platt, A., C. C. Weber, and D. A. Liberles. 2018. Protein evolution depends on multiple distinct population size parameters. *BMC Evol. Biol.* 18:17.
- Pollock, D. D., G. Thiltgen, and R. A. Goldstein. 2012. Amino acid coevolution induces an evolutionary Stokes shift. *Proc. Natl. Acad. Sci. USA* 109:E1352–9.
- Powell, J. R., and E. N. Moriyama. 1997. Evolution of codon usage bias in *drosophila*. *Proc. Natl. Acad. Sci. USA* 94:7784–7790.
- Rodrigue, N. 2013. On the statistical interpretation of site-specific variables in phylogeny-based substitution models. *Genetics* 193:557–564.
- Rodrigue, N., and N. Lartillot. 2014. Site-heterogeneous mutation-selection models within the phylobayes-mpi package. *Bioinformatics* 30:1020–1021.

- Rodrigue, N., and N. Lartillot. 2017. Detecting adaptation in protein-coding genes using a Bayesian site-heterogeneous mutation-selection codon substitution model. *Mol. Biol. Evol.* 34:204–214.
- Rodrigue, N., H. Philippe, and N. Lartillot. 2010. Mutation-selection models of coding sequence evolution with site-heterogeneous amino acid fitness profiles. *Proc. Natl. Acad. Sci. USA* 107:4629–4634.
- Scornavacca, C., K. Belkhir, J. Lopez, R. Dernat, F. Delsuc, E. J. P. Douzery, and V. Ranwez. 2019. OrthoMaM v10: Scaling-Up Orthologous Coding Sequence and Exon Alignments with More than One Hundred Mammalian Genomes. *Mol. Biol. Evol.* 36:861–862.
- Shah, P., D. M. McCandlish, and J. B. Plotkin. 2015. Contingency and entrenchment in protein evolution under purifying selection. *Proc. Natl. Acad. Sci. USA* 112:E3226–35.
- Siddiq, M. A., and J. W Thornton. 2019. Fitness effects but no temperature-mediated balancing selection at the polymorphic *adh* gene of *drosophila melanogaster*. *Proc. Natl. Acad. Sci. USA* 116:21634–21640.
- Spielman, S. J., and C. O. Wilke. 2015. The relationship between dN/dS and scaled selection coefficients. *Mol. Biol. Evol.* 32:1097–1108.
- Tamuri, A. U., N. Goldman, and M. dos Reis. 2014. A Penalized Likelihood Method for Estimating the Distribution of Selection Coefficients from Phylogenetic Data. *Genetics* 197:257–271.
- Wisotsky, S. R., S. L. Kosakovsky Pond, S. D. Shank, and S. V. Muse. in press. Synonymous site-to-site substitution rate variation dramatically inflates false positive rates of selection analyses: ignore at your own peril. *Mol. Biol. Evol.* .
- Yang, Z. 2006. *Computational Molecular Evolution*. Oxford Series in Ecology and Evolution.

- Yang, Z. 2019. Adaptive molecular evolution. In *Handbook of statistical genomics, vol. i*, ed. D. J. Balding, I. Moltke, and J. Marioni. John Wiley & Sons, Ltd.
- Yang, Z., and R. Nielsen. 2008. Mutation-selection models of codon substitution and their use to estimate selective strengths on codon usage. *Mol. Biol. Evol.* 25:568–579.
- Yang, Z., R. Nielsen, N. Goldman, and A. M. Pedersen. 2000. Codon-substitution models for heterogeneous selection pressure at amino acid sites. *Genetics* 155:431–449.
- Yang, Z., and W. J. Swanson. 2002. Codon-substitution models to detect adaptive evolution that account for heterogeneous selective pressures among site classes. *Mol. Biol. Evol.* 19:49–57.

Table 1. Amino acid sites under positive selection.

Data	Model	Sites
β -GLOBIN	MG-M3	7 , 11, 42 , 48 , 50 54 , 67 , 85 , 123
	MutSel-M3 _*	7 , 10, 11, <i>14</i> , 42 , 48 , 50 , 54 , 67 , 74 , 84, 85 , <i>110</i> , <i>113</i> , 123
ADH	MG-M3	-
	MutSel-M3 _*	9, 39, 49, 57 , 68 , 69, <i>72</i> , <i>81</i> , 85, 98, 133, 163, 165, <i>170</i> , 185 , <i>187</i> , <i>197</i> , 201, <i>205</i> , 208, 216 , 229, 253
VWF	MG-M3	-
	MutSel-M3 _*	5, 9, 26, 41 , <i>82</i> , 85, <i>103</i> , 108 , 125, <i>147</i> , 148, <i>158</i> , <i>177</i> , 182, <i>197</i> , 226, 227, 235 , 239 , 241, 242 , 247, 288, <i>291</i> , 307, 313 , <i>318</i> , <i>324</i> , <i>339</i> , 371, <i>379</i> , 390
ADORA3	MG-M3	-
	MutSel-M3 _*	2, <i>4</i> , <i>93</i> , 96
RBP3	MG-M3	-
	MutSel-M3 _*	-
S1PR1	MG-M3	-
	MutSel-M3 _*	<i>1</i> , <i>58</i> , <i>144</i> , <i>145</i> , <i>146</i> , <i>148</i>

Note.—Numbers in *italic* font are at the 0.9 level, in plain font at the 0.95 level, and in **bold** font at 0.99 level.

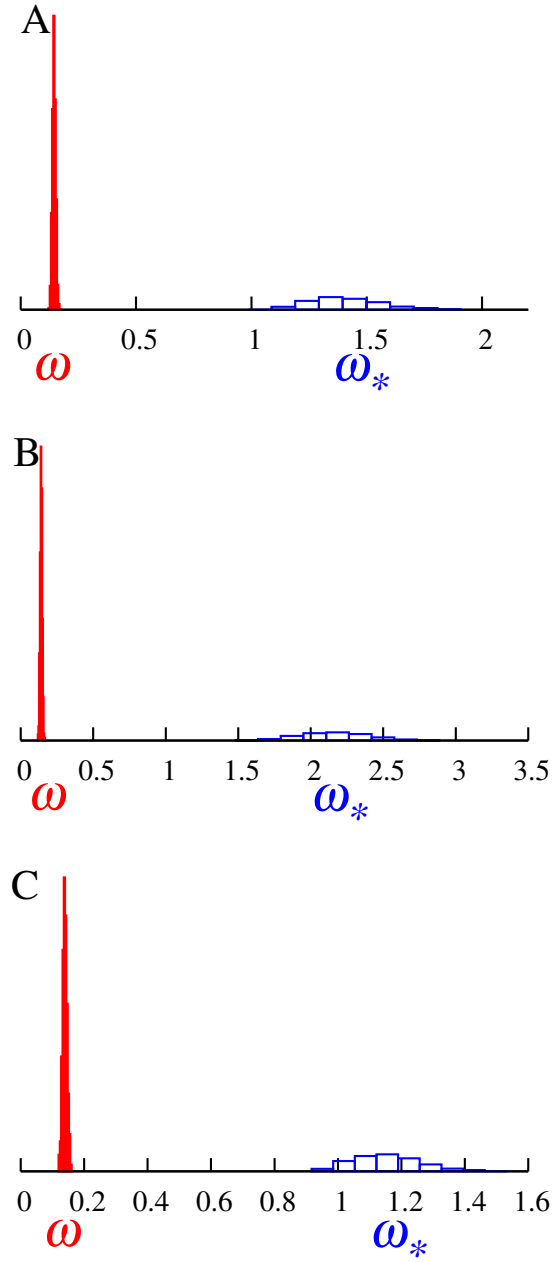


Figure 1. Posterior distributions of ω (red, under MG-M0) and ω_* (blue, under MutSel-M0_{*}) on simulated data sets with 10% of sites evolved under adaptive evolution (see methods).

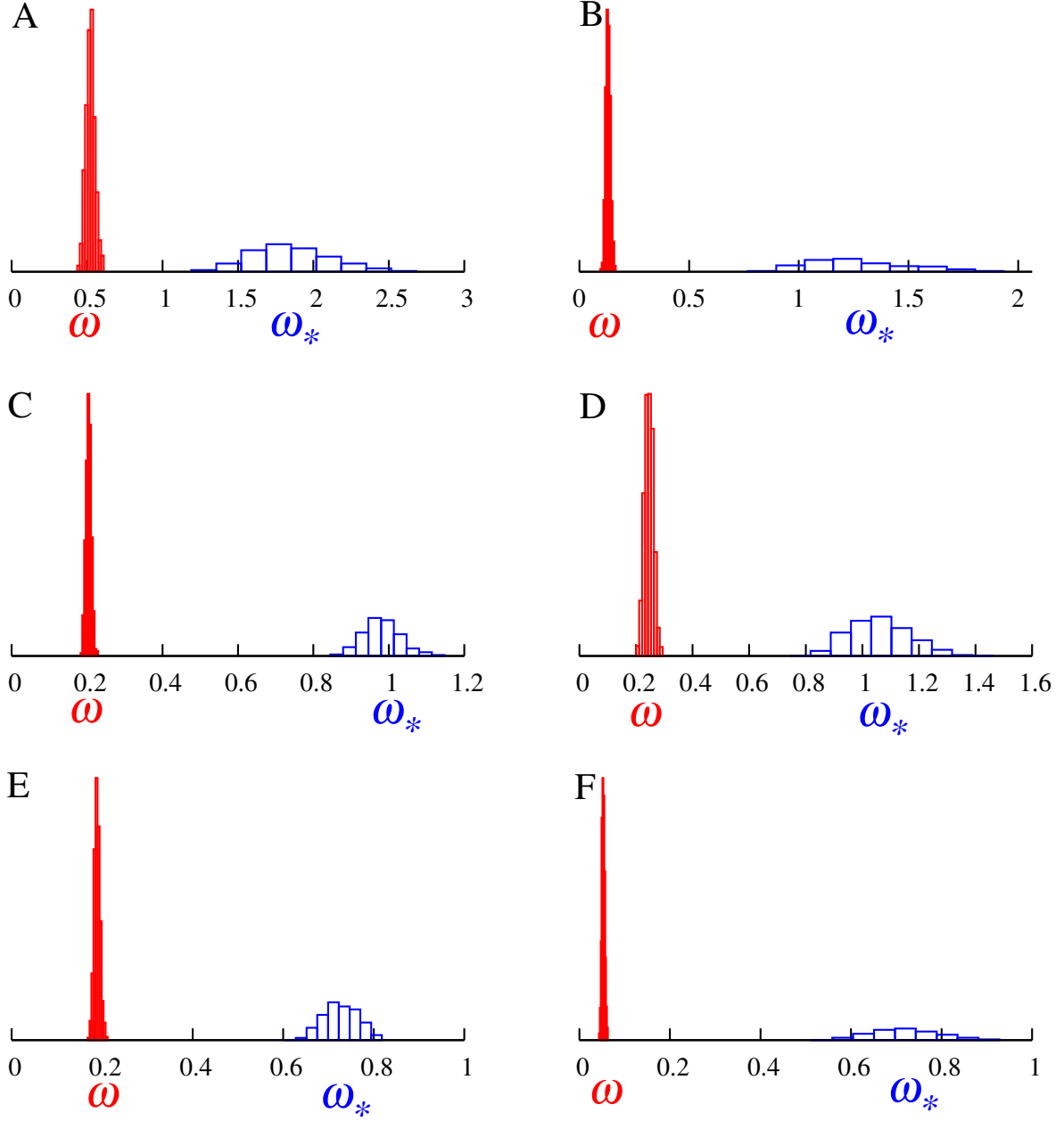


Figure 2. Posterior distributions of ω (red, under MG-M0) and ω_* (blue, under MutSel-M0_{*}) on β -globin17-144, adh, vwf, adora3, rbp3, slpr1 data sets (see methods).

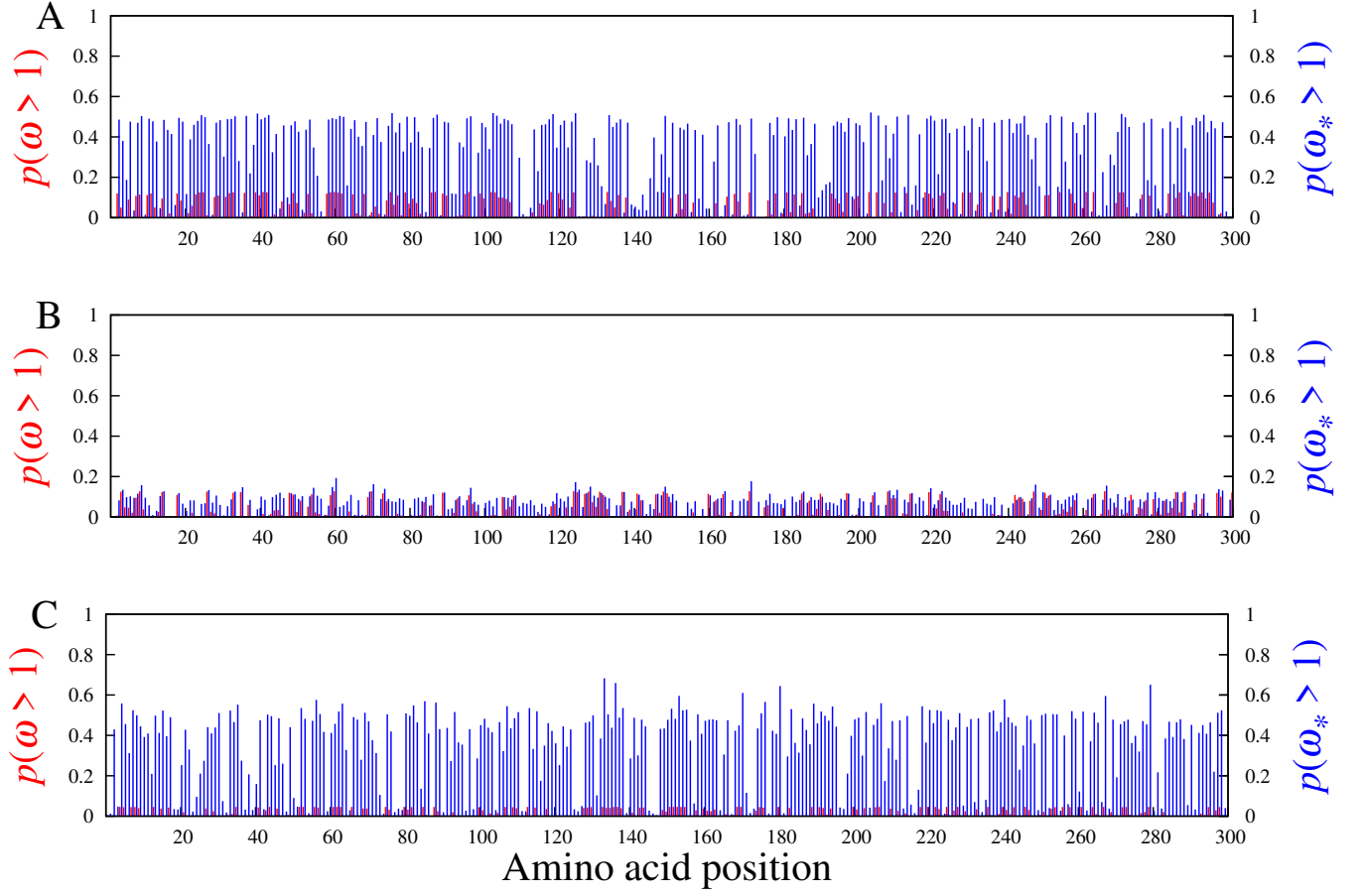


Figure 3. Site-specific posterior probabilities of ω (red, under MG-M3) and ω_* (blue, under MutSel-M3_{*}) being greater than 1 on data sets simulated under the pure mutation-selection framework.

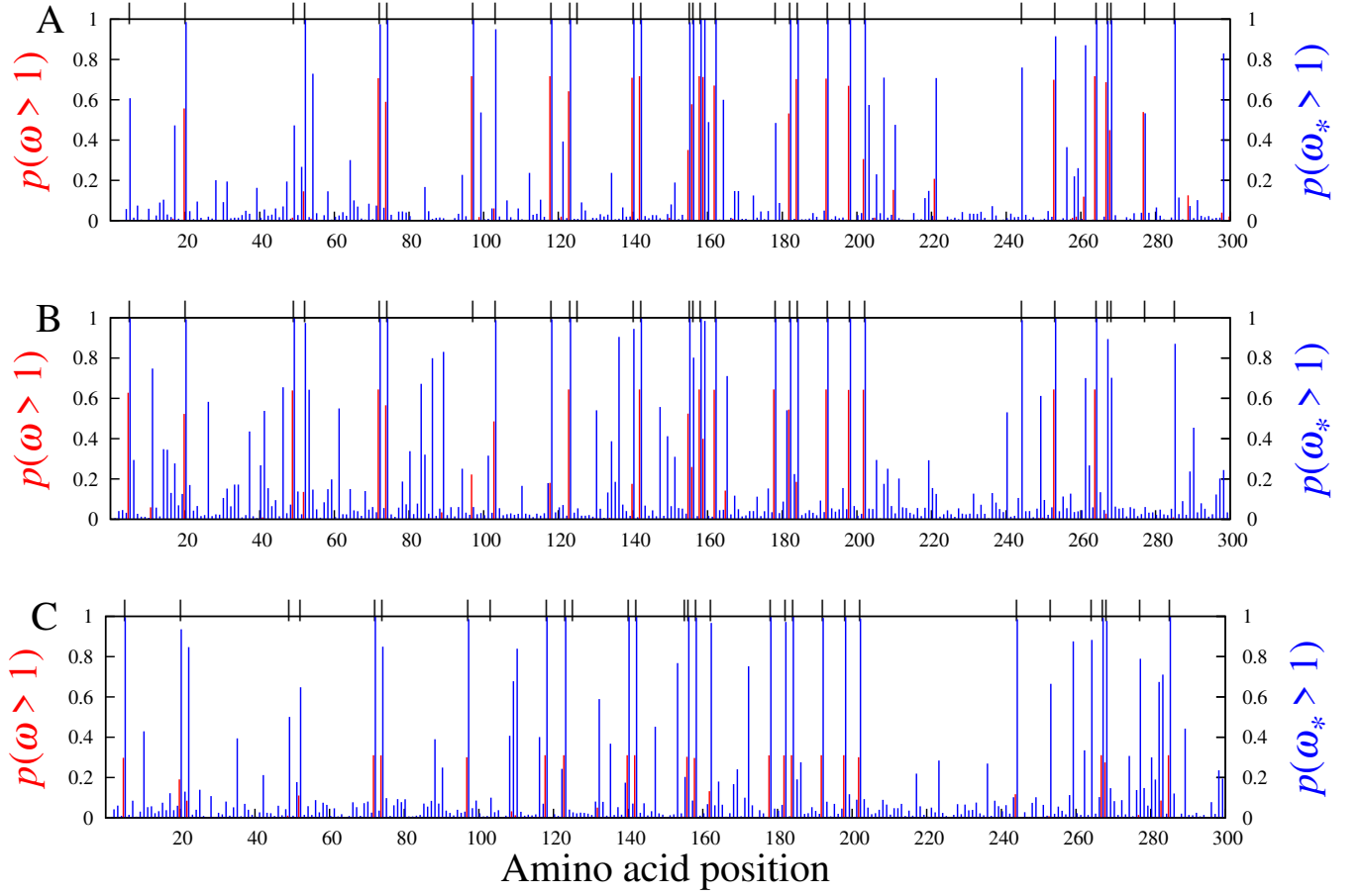


Figure 4. Site-specific posterior probabilities of ω (red, under MG-M3) and ω_* (blue, under MutSel-M3_{*}) being greater than 1 on data sets simulated with 30 sites (marked with at top of panels) under an adaptive regime, and the remaining 270 site under the pure mutation-selection framework.

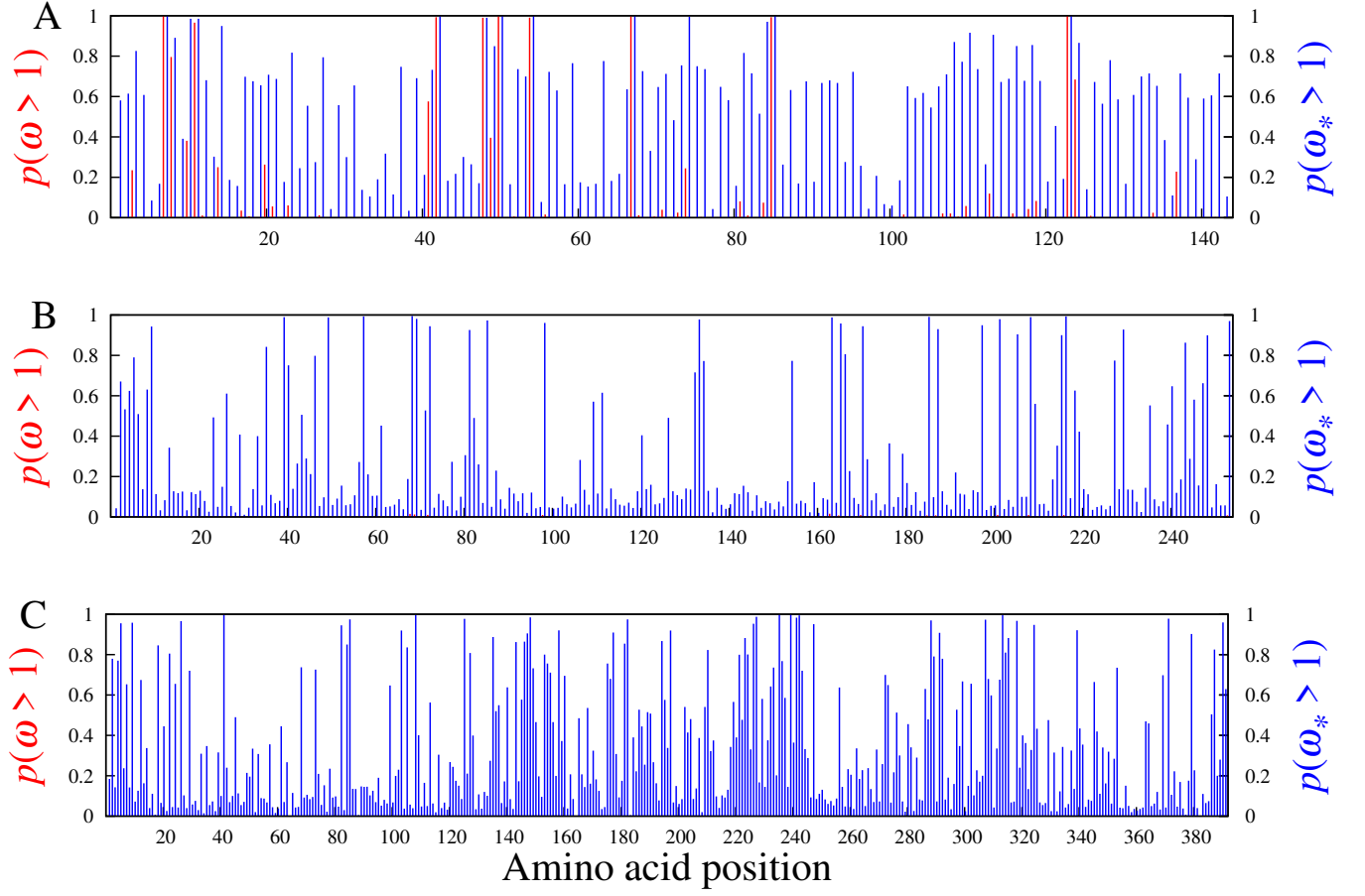


Figure 5. Site-specific posterior probabilities of ω (red, under MG-M3) and ω_* (blue, under MutSel-M3_{*}) being greater than 1 on β -globin, adh, and vwf.

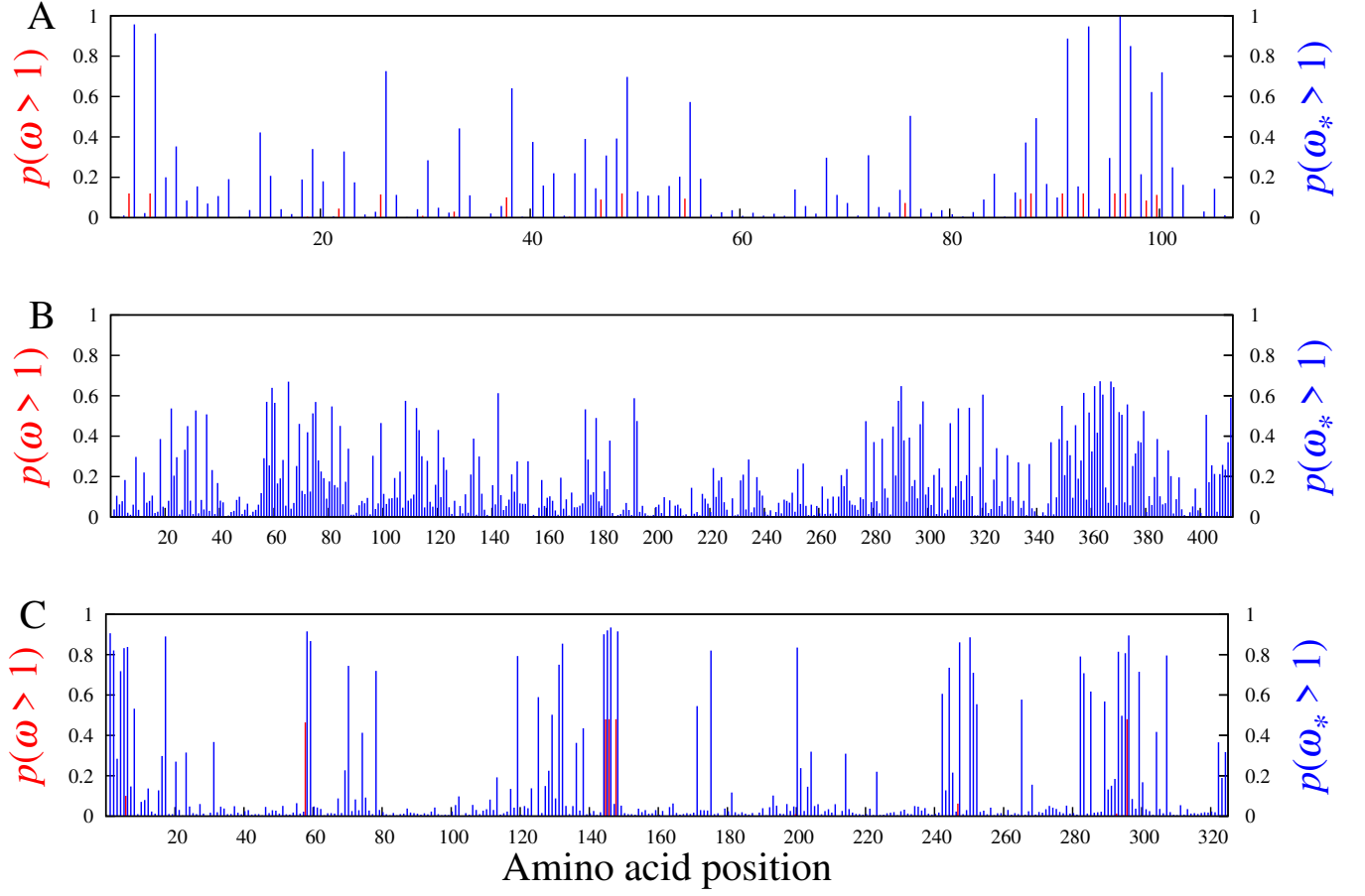


Figure 6. Site-specific posterior probabilities of ω (red, under MG-M3) and ω_* (blue, under MutSel-M3*) being greater than 1 on *adora3*, *rbp3*, *slpr1*.