



Published in final edited form as:

*Nat Rev Genet.* 2016 February ; 17(2): 109–121. doi:10.1038/nrg.2015.18.

## Causes of evolutionary rate variation among protein sites

Julian Echave<sup>1,\*</sup>, Stephanie J. Spielman<sup>2</sup>, and Claus O. Wilke<sup>2,\*</sup>

<sup>1</sup>Escuela de Ciencia y Tecnología, Universidad Nacional de San Martín, San Martín, Buenos Aires, Argentina

<sup>2</sup>Department of Integrative Biology, Center for Computational Biology and Bioinformatics, and Institute for Cellular and Molecular Biology, The University of Texas at Austin, Austin, TX, USA

### Abstract

It has long been recognized that certain sites within a protein, such as sites in the protein core or catalytic residues in enzymes, are more conserved than are other sites. However, our understanding of rate variation among sites remains surprisingly limited. Recent progress to address this includes the development of a wide array of reliable methods to estimate site-specific substitution rates from sequence alignments. In addition, several molecular traits have been identified that correlate with site-specific rates, and novel mechanistic, biophysical models have been proposed to explain the observed correlations. Nonetheless, at best, current models explain approximately 60% of the observed variance, highlighting the limitations of current methods and models, and the need for new research directions.

### Introduction

Different protein-coding genes within the same species vary widely in their rates of evolution. For example, proteins that are highly expressed or that perform critical functions tend to evolve more slowly than will other proteins<sup>1</sup>. In addition to this gene-wide variation, and perhaps more interestingly, evolutionary rates vary among residues *within* a given protein. Although some of this variation is attributable to positive, diversifying selection, e.g. selection pressure triggering adaptation to environmental or other changes, there exists substantial rate heterogeneity even at sites not subject to such selection pressure. This heterogeneity likely emerges from the differing functional and/or biophysical constraints affecting different sites. Accurately modeling this among-site heterogeneity is critically important in evolutionary studies, particularly in phylogenetic inference<sup>2-8</sup>. Phylogenetic models which allow for among-site rate heterogeneity universally provide better fits to data than do models which assume constant rates across sites<sup>3;9-13</sup>. However, such models are largely phenomenological in nature and contain no information about the mechanistic source of among-site rate heterogeneity<sup>14</sup>. Although it is clear that substantial rate variation exists, the underlying mechanisms which generate the observed rate heterogeneity remain elusive.

Over the years, it has become apparent that site-specific evolutionary rates are influenced by a dynamic interplay between structural and functional constraints (Figure 1). In the 1960's,

\*Corresponding author: J. E. jechave@unsam.edu.ar, C. O. W. wilke@austin.utexas.edu.

Perutz et al.<sup>15</sup> investigated site-specific sequence variability in globin proteins and found that “internal sites” were generally more conserved than were “superficial sites.” They further reasoned that “special functions” had to be influencing sites which did not conform to this pattern<sup>15</sup>. Later, Kimura and Ohta built upon these observations by proposing the governing principle that “[f]unctionally less important molecules or parts of a molecule evolve (in terms of mutant substitutions) faster than more important ones”<sup>16</sup>. Kimura and Ohta additionally recognized that surface protein residues “are usually not very critical to maintaining the function or tertiary structure, and the evolutionary rates in these parts are expected to be much higher”<sup>16</sup>.

Following these early studies, most work on the sequence-structure-function relationship has been done from the perspective of structural biology. In general, such studies have not considered evolutionary *rates*, but have considered conservation only qualitatively or through conservation scores that do not take into account the nature of the evolutionary processes or the phylogenetic relationships. Therefore, our current understanding of how functional and structural constraints interact to shape evolutionary rate heterogeneity remains limited. To develop a complete picture of protein evolution, we need to identify the precise structural and functional properties which ultimately govern protein evolutionary rates, and we need to develop mechanistic explanatory models thereof.

In recent years, there has been significant progress on this front. Advances in computational evolutionary modeling have provided a variety of robust methods for estimating site-specific rates both from amino-acid sequences and from protein-coding DNA sequences. Further, numerous studies have discovered functional, structural, and dynamical molecular features that correlate with rates<sup>17-23</sup>, and biophysical models have been proposed that predict site-specific rates from protein thermodynamics<sup>24;25</sup>. Studies on the relationship between rate and structure have yielded consistent findings, suggesting that current data and methods provide a solid foundation upon which further knowledge can be built. Therefore, the time is ripe to review and synthesize our current understanding of site-specific evolutionary rates and the identified structural and functional aspects which influence them.

We focus here as much as possible on works that study site-specific rates estimated using state-of-the-art molecular evolution methods. To provide context or where we lack rate-based studies, we also discuss relevant work based on other measures of sequence conservation. In the following, we first describe current methods to estimate rates from sequence data. We then consider molecular traits related to site-specific conservation. Next, we discuss two recent mechanistic biophysical models that have been used to explain site-specific rates. We conclude with a discussion of challenges and future directions.

## Estimation of site-specific rates

Computing site-specific evolutionary rates requires two pieces of data: a multiple sequence alignment, with either codon or amino-acid data, and a corresponding phylogeny. Estimating the substitution rate at each individual site in a protein can be a computationally burdensome endeavor, much more so than estimating a mean rate for an entire protein sequence. Indeed, individual sites contain far less information than do entire protein alignments, and thus large

and diverse datasets are needed for reliable inference. In particular, most alignment sites must have experienced several substitution events for a reasonably accurate rate estimate to be made<sup>26</sup>.

Broadly speaking, site-specific inference methods follow one of two paradigms<sup>27-31</sup>: i) directly counting observed substitutions along a phylogeny, and ii) employing a Markov model of sequence evolution to infer evolutionary rate parameters, typically in a maximum-likelihood (ML) framework. Other methods, which compute sequence entropy or conservation scores<sup>32-34</sup>, are useful for assessing the tolerance of a given residue to mutation. However, they cannot be substituted for true measures of evolutionary rate, as they do not typically account for phylogeny, which represents the evolutionary relationships among sequences. Indeed, it is possible for a given site to have high entropy and a low evolutionary rate or vice versa.

### Inferring rates from codon data

In the context of protein-coding sequences, evolutionary rates are typically estimated with the ratio  $\omega=dN/dS$ , where  $dN$  is the rate of non-synonymous substitutions and  $dS$  is the rate of synonymous substitutions. To make  $dN$  and  $dS$  directly comparable, they are normalized to account for the approximately 3-fold higher likelihood that a random mutation is non-synonymous rather than synonymous<sup>35</sup>. The ratio  $\omega$  has been developed primarily to detect sites under adaptive evolution (for which  $\omega>1$ ), but it can also be used to estimate site-specific rates<sup>30;36</sup>.

Counting-based methods, the oldest class of  $dN/dS$  inference methods, calculate  $dN/dS$  simply by enumerating the observed changes either between pairs of sequences or along a phylogenetic tree<sup>5;29;37-39</sup>. While relatively fast, these methods do not adequately account for multiple substitutions, variation in branch lengths, and other biases, and therefore they tend to produce biased  $dN/dS$  estimates<sup>5;29;35</sup>.

Most modern-day inference approaches, on the other hand, estimate rates in a ML framework with an explicit Markov model of sequence evolution. By implicitly accounting for any hidden substitutions along branches, ML-based methods are more robust and less biased than are counting methods. Site-specific rates are obtained either by fitting a rate parameter individually to each site in the coding sequence (known as a “fixed-effects likelihood” or FEL approach)<sup>28;29;40</sup> or by considering the rate to be a random variable drawn from a distribution governing the entire protein (known as a “random-effects likelihood” or REL approach)<sup>9;28;29;41</sup>. In the REL approach, site-specific rates are calculated using a Bayes Empirical Bayes framework<sup>42</sup>. FEL lacks power for small datasets, and thus it is most appropriate for use on large datasets (at least 200 sequences in which most sites have experienced recurrent changes)<sup>29</sup>. On the other hand, REL is best suited for datasets of intermediate size (50–200 sequences), as its estimates on either very small or very large datasets are usually quite biased<sup>29</sup>. Importantly, while smaller datasets (e.g. 16–50 sequences) may suffice when detecting episodic and/or diversifying selection<sup>43</sup>, obtaining reliable site-specific  $dN/dS$  point estimates requires larger datasets. We note that FEL methods are primarily implemented in HyPhy<sup>44</sup> and its corresponding web-server DataMonkey<sup>45</sup>, whereas REL methods are implemented in both HyPhy<sup>44</sup> and PAML<sup>46</sup>.

There are two possible strategies for parameterizing  $dN/dS$  during rate inference. One can either fix  $dS$  across sites and simply estimate site-specific  $dN$  values, or one can estimate a separate  $dN$  and  $dS$  parameter for each site. Importantly, in the context of protein evolution,  $dN$  is the primary parameter of interest, and  $dS$  serves only as a normalization factor to determine into which selection regime (e.g. purifying, neutral, positive selection) a given residue falls. Given that site-specific rate estimates are inherently noisy, normalizing each site's inferred  $dN$  with a corresponding site-specific  $dS$  likely introduces substantial, and potentially confounding, error. Models which either fix  $dS$  to 1<sup>9;41;47</sup>, or similarly infer gene-wide  $dS$  estimates for normalization, therefore, may represent a more robust strategy for obtaining reliable rates of protein sequence evolution.

Although the inference approaches described above are generally implemented in a ML framework, several evolutionary rate inference approaches have recently emerged that use Bayesian, rather than frequentist, statistics. For example, a novel method known as renaissance counting<sup>30</sup>, implemented in the software package BEAST<sup>48</sup>, combines a counting-based approach with empirical Bayes regularization, thus leveraging the power of large datasets to produce site-specific estimates of accuracy comparable to FEL and REL. In addition, the inference method FUBAR (a Fast, Unconstrained Bayesian AppRoximation for inferring selection) adapts the REL framework to rapidly fit a large, pre-specified grid of evolutionary rates to the data in a hierarchical Bayesian framework<sup>49</sup>. This approach is exceptionally fast yet yields reliable rate estimates for data with sufficient divergence. Finally, an approach for estimating gene-wide  $dN/dS$  values using a Bayesian statistical framework has recently been described, and future development may see this method extended to site-specific estimation<sup>50</sup>.

### Inferring rates from amino-acid data

The primary approach for inferring rates from amino-acid data is implemented in the program Rate4Site<sup>51</sup>. Rate4Site estimates a per-site rate-scaling factor that indicates how rapidly each residue evolves relative to the mean protein rate. It is implemented in both ML-based and Bayesian frameworks, with the latter being the default<sup>52</sup>. Under the Bayesian framework, Rate4Site employs a random-effects approach, specifying either a single gamma distribution<sup>52</sup> or a mixture of gamma distributions<sup>11</sup> as the prior rate distribution. A Bayes Empirical Bayes approach is then used to calculate site-specific rates. Importantly, Rate4Site can only accommodate datasets with fewer than approximately 300 sequences<sup>51</sup>, and thus future research endeavors may seek to extend this method for use on larger datasets.

Alternatively, Fernandes and Atchley proposed a fixed-effects framework for estimating site-specific evolutionary rates from proteins<sup>53</sup>. Unlike the Bayes Empirical Bayes approach in Rate4Site, this method provides an independent rate estimate at each site, thus avoiding the confounding influences of mis-specified prior distributions on the rate. Finally, a relatively new method called GP4Rate uses a Gaussian process to infer site-specific evolutionary rates while taking protein tertiary structure into account<sup>54;55</sup>. This approach effectively accounts for non-independence among site-specific rates.

## Structural and environmental rate constraints

The among-sites rate variation observed in natural sequence alignments is, to a large extent, driven by the requirement that proteins fold properly and stably into their required, active conformation. In evolutionary terms, this requirement corresponds to purifying selection, such that sites at which mutations would disrupt folding or stability the most will be most conserved. In addition, proteins experience selection pressure to avoid disrupting their native environment, for example by forming non-specific protein-protein interactions<sup>56</sup>. In the following, we refer to these kinds of selective forces as structural and environmental constraints on sequence evolution.

Early studies of structural constraints established a basic paradigm dividing the protein into two general regions: the interior, which evolves slowly, and the surface, which evolves more rapidly<sup>15;16</sup> (Figure 1). This paradigm poses two questions: i) why is the protein interior more conserved than the surface? and ii) what are the salient structural differences between these two regions? Several biophysical measures have been proposed to explain the observed rate differences in structurally distinct regions. These measures include a residue's solvent accessibility, packing density, and flexibility. Although these measures are distinct, they all quantify the position of a given residue relative either to other nearby residues or to the protein as a whole.

### Solvent accessibility

The most obvious difference between surface and interior is that the surface is accessible to the external environment, e.g. water, while the interior is not. This observation defines solvent accessibility (ASA, accessible surface area, or SASA, solvent-accessible surface area), which indicates the surface area of a given residue that is accessible to water. ASA values are commonly normalized by the largest possible ASA for a given amino acid<sup>57</sup>, resulting in the relative measure RSA (Relative Solvent Accessibility). RSA ranges from 0 for completely buried residues to 1 for completely exposed residues.

Because early studies showed a relatively high conservation of protein cores, solvent accessibility formed the basis of research investigating the relationship between protein evolution and structure through the mid-1990s<sup>58-63</sup>. A broad consensus emerged that amino-acid substitution rates and properties differ between buried and exposed sites, with buried sites being more conserved and tending more towards hydrophobic residues, due to the local environment in a protein's core. Over time, it was generally assumed that solvent accessibility represented the dominant structural constraint on evolutionary rate. For example, one study partitioned residues according to RSA and secondary structure (e.g. helix, sheet, coil, turn, etc.) and found that exposed sites evolved more rapidly than buried ones, regardless of secondary structure<sup>64</sup>. A study of site-specific substitution rates for 25  $\alpha/\beta$  barrel enzymes found that a minimal model which considered RSA as the main factor could not be improved upon by adding other properties such as secondary structure or H-bonding information<sup>17</sup>. A third study showed that amino-acid properties, such as hydrophobicity or size, had little influence on site rates beyond the strong effect of solvent exposure<sup>65</sup>.

More recently, studies have leveraged the power of genomic datasets and sophisticated rate inference methods to perform more comprehensive analyses. For example, Franzosa and Xia<sup>18</sup> examined the correlation of site-specific  $dN/dS$  with several structural properties across nearly 1000 *Saccharomyces cerevisiae* proteins. They found that evolutionary rate increased linearly with RSA. Since then, several additional studies have reinforced the strong, positive relationship between RSA and site-specific rate<sup>66-68</sup>.

### Packing density

As we have seen in the previous subsection, solvent accessibility has become the *de facto* structural measurement to use in protein evolution studies. However, more recent work has called the central role of solvent accessibility into question, suggesting instead different structural measures which correlate more strongly with evolutionary rate.

In particular, instead of quantifying the extent to which a given residue comes into contact with solvent, as RSA and ASA do, we can also quantify the extent to which a residue comes into contact with other residues in the protein. This alternative concept, known as packing (or contact) density, indicates how densely packed a residue is within the protein tertiary structure. The two packing measures most commonly used in evolutionary studies are the contact number (CN) and the weighted contact number (WCN). For a given amino acid, CN simply counts the number of other residues within a local, structural neighborhood. By contrast, WCN considers all residues in the protein and weights them by the square of their inverse distance to the focal amino acid<sup>19;69</sup>.

Packing density was initially introduced into the protein evolutionary rate literature because theoretical calculations predicted that more densely packed proteins should be easier to design<sup>70</sup> and, as a consequence, evolve more rapidly<sup>71</sup>. The first studies relating contact density to evolutionary rate focused on whole-protein rather than site-specific rates, and broadly found that average protein rates are higher for proteins whose residues have, on average, higher packing densities<sup>71-73</sup>. Once packing density was established as a veritable predictor of evolutionary rate, several groups began investigating the packing–rate relationship on a site-specific basis. Franzosa and Xia found a modest but significant partial correlation of CN with site-specific rates while controlling for RSA, prompting them to conclude that CN influences rate independently from RSA<sup>18</sup>: More densely packed sites evolve more slowly. Importantly, Franzosa and Xia concluded that, while CN does predict evolutionary rate, RSA is a much stronger predictor<sup>18</sup>. However, subsequent works have challenged this finding, pointing out that using WCN instead of CN to estimate packing density results in stronger predictive power, and moreover that RSA's independent contribution becomes relatively small when WCN is controlled for<sup>19;23;74</sup>.

In all above-mentioned studies, WCN was calculated using  $C_{\alpha}$  carbons to represent residues. Recent work has proposed that using side chains rather than  $C_{\alpha}$  carbons may provide a more robust determinant<sup>75</sup>. Indeed, side-chain based WCN has consistently outperformed both  $C_{\alpha}$ -based WCN and RSA as a rate determinant (Figure 2). Therefore, it now appears that WCN is the main determinant of site-specific substitution rates, and RSA provides a comparatively minor independent contribution.

Even so, as can be seen in Figure 2, both the relative performance of WCN and RSA and the overall performance of either predictor varies widely among protein structures. While for some proteins we can explain over 60% of the observed rate variation with simple structural measures such as RSA and WCN, for other proteins we can only explain less than 10% with the same measures. Similarly, for some structures WCN outperforms RSA by over 15 percentage points, while for others the two measures perform comparably. The underlying cause of these discrepancies is not well understood, but there are several, not mutually exclusive possibilities: First, other predictors, e.g. related to protein function (see below), may be more important in some structures than in others. As an example, consider the case of structure 1AKO shown in Figure 2B. Second, in some cases alignments may be poor or contain insufficient or excessive divergence. In general, alignments need to be sufficiently diverged for accurate rate inference at individual sites<sup>26</sup> but not be saturated with mutations. Finally, standing polymorphisms, slightly deleterious mutations, or mutations hitchhiking on recent selective sweeps may cause biased rate estimates<sup>76</sup>, and as a consequence structural predictors may not work well on alignments in which any of these factors are highly prevalent.

### Flexibility

Proteins are not static structures; they are dynamic polymers that undergo constant conformational fluctuations. Such movements are frequently critical for protein function. For example, enzymes must shift their structural conformation to expose the active site before a substrate can be accommodated. Similarly, conformational changes could control the mutational tolerance of a site, such that a site in a highly flexible region of a protein structure would likely be more tolerant to mutations than a site in a less flexible region<sup>20-22;77</sup>.

At the site level, conformational dynamics can be quantified using measures of local flexibility, such as Mean Square Fluctuations (MSF) or B-factors. These quantities measure to what extent a given residue changes its position over time. Using these and similar quantities, several studies have found that site-specific sequence variation correlates with local flexibility<sup>20;22;24;77</sup>, such that flexible sites evolve more rapidly than rigid ones do.

That flexibility and evolutionary divergence correlate has been interpreted as evidence that protein dynamics imposes significant constraints on sequence evolution<sup>20;77</sup>. However, whether flexibility is the underlying causal factor in the observed flexibility–rate correlation is unclear. Local flexibility relates directly to packing density<sup>78</sup>, which, as discussed above, correlates strongly with rate. Therefore, it is possible that either flexibility or packing density represents the underlying causal factor affecting rate.

One of the first studies of the flexibility–evolution relationship used the inverse of the contact number as a proxy of flexibility, making the implicit assumption that flexibility was the causal factor<sup>79</sup>. Several later studies made similar arguments<sup>20;77</sup>. However, if flexibility were the actual determinant and packing just a proxy, site-specific rates would have higher correlations with flexibility than with packing, which is not the case. Instead, the reverse is true: Site-specific substitution rates correlate more strongly with measures of packing density (such as WCN) than with measures of flexibility (such as MSF)<sup>22;24</sup>. Moreover,

when packing density is controlled for, no residual correlation remains between rate and MSF<sup>24</sup>. Thus, it appears that flexibility correlates with rate simply because both quantities are determined by local packing density, and not because of a direct, causal relationship between flexibility and rate.

### **Other structural constraints: folding kinetics, protein expression, and cellular environment**

There are other quantities and constraints that broadly relate to the requirement of proper and stable protein expression. For example, structural factors that may constrain evolution at the site level include secondary structure, side-chain hydrogen bonds, unusual side-chain rotamers, nonplanar peptide bonds, strained main-chain conformations, or buried hydrophilic-charged residues<sup>80</sup>. However, the majority of these factors have little explanatory power for site-specific rates once solvent accessibility is controlled for<sup>17;64;81</sup>. One factor that does matter is structural disorder: sites in disordered regions tend to evolve more rapidly, and with less conservative amino-acid substitutions, than do ordered regions<sup>82-84</sup>.

In addition to being stable, proteins need to fold sufficiently rapidly. It is reasonable to expect this requirement for rapid folding to further constrain evolutionary divergence. Folding occurs via a transition state involving a small number of sites that assume their native conformation, a *folding nucleus*<sup>85</sup>. A classic study of the cytochrome c family speculated that those sites that were very conserved but not involved in activity could be the sites forming the folding nucleus<sup>86</sup>. In another study, the experimentally determined folding-nucleus sites of 9 proteins were found to be more conserved than average<sup>87</sup>. By contrast, a more exhaustive, systematic study found no significant evidence for extra conservation of folding-nucleus residues and, moreover, argued that the previously found special conservation was due to biases of the experimental data<sup>88</sup>. Note that the cited studies quantified conservation using entropy-based measures. A more recent study based on substitution rates found no significant differential conservation of folding-nucleus sites<sup>89</sup>.

The level at which a protein is expressed and the cellular location where it functions influence site variation as well. For example, it is well known that more highly expressed proteins evolve more slowly<sup>90</sup>. Analyses at the site level have shown that this evolutionary constraint is RSA dependent<sup>67;68</sup>. In yeast, the difference in mean evolutionary rates for lowly and highly expressed genes increases linearly with RSA. For the most buried residues (RSA=0), sites in highly expressed genes evolve approximately two times slower than sites in lowly expressed genes, while for the most solvent-exposed genes (RSA=1), the relative ratio in mean rates grows to above three<sup>67</sup>. An example of environment dependence can be seen in membrane proteins. In such proteins, the transmembrane regions are more evolutionarily conserved than the extramembrane regions are, and this effect seems to be separate from conservation due to solvent accessibility<sup>91;92</sup>.

### **Rate variation caused by protein function**

None of the structural and environmental constraints discussed in the preceding section are directly related to protein function. Indeed, while proteins need to fold stably into their active conformation, this requirement alone does not guarantee that they will function



properly, except perhaps for proteins whose sole purpose is structural, as building blocks of organs and tissues. All other proteins, including enzymes, transcription factors, molecular motors, and antibodies, have specific functional or active sites at which they experience additional selection pressures. These pressures may act in the form of purifying selection, causing increased evolutionary conservation, or in the form of positive, diversifying, or balancing selection, causing increased evolutionary variability.

### Purifying selection

In many cases, selection for function adds additional evolutionary constraints to the specific amino acids (such as catalytic sites) involved in a protein's function<sup>93;94</sup>. These sites—as well as their neighbors—are often particularly conserved<sup>93</sup> (as one example, consider the pattern of conservation near the active site in Figure 1). Curiously, selection for function seems to extend beyond just the active site and its immediate neighbors. Dean et al. found that distance to the active site correlated with site-specific variation in several enzymes<sup>17</sup>. Similarly, several experimental studies have observed that mutations far away from the active site can disable protein function by inducing protein-wide structural alterations (see also Box 1). These findings provide evidence for the presence of long-range, indirect interactions in protein structures, likely mediated by steric interactions among neighboring amino acids<sup>95</sup>.

Besides catalytically active sites, residues involved in protein–protein or protein–nucleic-acid interactions also experience added functional constraint and generally are more conserved than other surface sites<sup>18;96-98</sup>. The extent to which protein–protein interactions constrain site evolution seems to depend on the exact nature of the interaction. For example, obligate interactions, which often persist for the lifetime of the protein, are associated with lower rates compared to transient interactions that occur only occasionally<sup>96</sup>. Since residues involved in protein–protein interactions experience reduced solvent accessibility when the interacting protein partner is present, Franzosa and Xia asked whether this reduction could explain the added evolutionary constraint on interface residues<sup>18</sup>. Their answer was “not entirely.” They found that while the evolutionary constraint increases linearly with increasing amount of solvent-accessible surface area lost due to the interaction, there is also an additional, albeit low, fixed cost that can be attributed to the mere fact that a residue is participating in the interface interaction<sup>18</sup>. The fixed cost is independent of RSA. Its low magnitude is consistent with recent experiments and computer simulations showing that protein–protein interfaces can maintain function despite extensive divergence of one partner<sup>99</sup>.

Finally, ligand-binding sites tend to be more conserved than other sites, and this conservation is exploited in ligand-binding-site prediction methods<sup>100;101</sup>. However, the degree of conservation varies. For example, while catalytic sites are very conserved, allosteric sites, while more conserved than average, vary more than catalytic sites, because of either weaker constraints or positive selection on the regulatory mechanisms<sup>102</sup>.

## Positive selection

When organisms are faced with novel or changing environments, their protein-coding genes may experience positive selection, i.e., a selection pressure to change rather than remain the same. This selection pressure will primarily act on the sites directly involved in the specific function under selection. That individual sites can experience positive selection for specific function was first recognized over thirty years ago. Early examples include positive diversifying selection in the active sites of three related rodent protease inhibitors<sup>103</sup>, over dominance in the antigen recognition site of human and mouse class I MHC (major histocompatibility complex) genes<sup>104</sup>, and positive selection in the V3 region of the HIV-1 envelope gene<sup>41</sup>, likely reflecting immune escape or adaptation to cell tropism.

Positive selection is probably the most prevalent in viral surface proteins, which experience intense selection pressure to adapt to their host or to escape their host's immune response. For example, positive selection strongly shapes the evolution of the influenza hemagglutinin in protein<sup>105-108</sup>, which initiates fusion of the viral envelope with the cellular membrane. Indeed, in that protein, positive selection explains nearly as much rate variation as does RSA<sup>109</sup>.

Since positively selected sites can provide meaningful insight into the functioning of and selective constraints on a gene, the molecular evolution community is broadly interested in identifying such sites under many different scenarios. This interest has spurred the development of numerous tests for positive selection, based on the approaches and inference frameworks discussed in the previous section “Estimation of site-specific rates.” Importantly, these existing methods rarely consider the baseline structural constraints acting on most sites in a protein. Consequently, they tend to be overly conservative and likely miss many important sites<sup>36</sup>.

For example, the widely used test for  $dN/dS > 1$  makes the implicit assumption that sites evolve at  $dN/dS = 1$  in the absence of selection. However, for protein-coding sequences selection is virtually never absent; structural constraints will induce purifying selection that will push  $dN/dS$  to values much lower than 1 in nearly all cases. Thus, we can expect that positive selection at a site that otherwise would have been highly conserved may yield elevated  $dN/dS$  values that nevertheless remain below 1 (Figure 3). Purely statistical methods that consider only the value of  $dN/dS$  at individual sites will not be able to identify such site (Figure 3A). However, methods that incorporate appropriate baseline expectations derived from protein structure may be able to do so (Figure 3B). In general, functional constraints may be the reason why sites evolve faster or slower than expected from structural constraints (Figure 3B; see also Figure 2B). Finally, methods that incorporate both structural and functional information may accurately predict the rate of evolution at functionally selected sites, thus providing mechanistic insight into why a given site evolves at the rate it does (Figure 3C).

## Predicting rates from first principles

Correlations between rates and predictor variables allow us to identify factors that influence rate variation, but ultimately they do not provide explicit mechanistic insight into why a

given site is variable or conserved. To gain mechanistic insight, we need to develop biophysical models, grounded in first principles. Several biophysical models have been used to study issues such as marginal protein stability and site–site coevolution<sup>110;111</sup>. Here, we will focus on the models that have been used to study site-specific rates.

While phenomenological models depend directly on predictor variables such as RSA and WCN, biophysical models are based, essentially, on protein stability. This choice is reasonable because from a physicochemical perspective a protein's function is determined by its thermodynamics and kinetics, which are related to stability. In addition, stability is related to all the molecular features that correlate with evolutionary rates. For example, solvent accessibility (RSA) is related to stability via the energetic cost of burying a side chain into the core<sup>112;113</sup>. Similarly, packing density (CN and WCN) and flexibility (MSF) are related to the mean interaction energy of a site with the rest of the protein<sup>78</sup>. Therefore, molecular features such as RSA, WCN, and MSF could be mere proxies of stability, which would be the true determinant of protein fitness and, therefore, site-specific evolutionary rates.

Two distinct biophysical models have been proposed in the literature. The *stability-threshold* model, also referred to here as the *native-stability* model, links site-specific substitution rates to mutational changes of protein thermodynamic stability<sup>114</sup>. Specifically, it assumes that all proteins with sufficient stability in the native state function equally well and have identical fitness, and proteins that are not sufficiently stable have zero fitness. Thus, they impose a stability-threshold condition that the protein needs to meet at all times during its evolution. We note that variations of this model may employ a sigmoidal function instead of a hard threshold, but they show similar behavior and make the similar assumption that native stability is the critical factor in a protein's function<sup>115;116</sup>.

In the threshold model, the probability of fixation of a given mutation is equivalent to the probability that the mutation will push the stability below the threshold. This probability can be calculated under the assumption that the free-energy changes  $G_{ij,k}$  for mutations from amino acid  $i$  to amino acid  $j$  at site  $k$  are known<sup>25</sup>. ( $G$  measures the change in free energy between two protein variants, and  $G$  is a measure for the stability of the protein fold.) These free-energy changes can be estimated from atomic force fields such as FoldX<sup>117</sup> and subsequently converted into rate estimates<sup>25</sup>.

In contrast to the native-stability model, the *active-stability* or *stress* model assumes that a mutation will affect not just the native conformation of a protein but its whole energy landscape. For example, if a protein needs to adopt a certain active conformation to function, the stability of this active state likely affects fitness. The active-stability model postulates that the fixation probability is proportional to the probability of finding the mutant in the active conformation, which in turn is a function of the stability change of the active conformation,  $G^*$  (Figure 4). This stability change can be calculated analytically via perturbed elastic network models (ENMs)<sup>24;75</sup>. In particular, using the parameter-free Anisotropic Network Model (pfANM)<sup>118</sup>, one can show that the substitution rate should be proportional to the weighted contact number WCN. Thus, the active-stability model provides a mechanism for the observed rate–WCN correlation<sup>24;75</sup>.

The site-specific rates predicted by both the active-stability and the native-stability models are in good agreement with empirical rates<sup>24;25;75</sup>. However, the active-stability model tends to perform better, and there is little independent contribution from the native-stability model once the active-stability contribution is accounted for. Therefore, the empirical evidence so far favors the active-stability model. Mechanistically, the active-stability model can explain why native-stability predictions correlate with empirical rates: As shown in Figure 4,

$G^* = G + G^\ddagger$ , where the first term is the change in native-state stability and the second term is the change in activation free energy. From the perspective of the active-stability model,  $G$  affects evolutionary constraints via its effect on the stability of the active conformation. Importantly, a given mutation may destabilize the active state even if it increases native stability (Figure 4), and the native-state-model predictions for such mutations would be incorrect.

## Challenges and future directions

Our present-day understanding of site-specific rate variation broadly agrees with the initial picture developed over 40 years ago. However, whereas early work suggested definite structural regions (interior and surface), the emerging view is more nuanced (see also Table 1). Structural constraints decrease continuously from the solvent-inaccessible, tightly packed, and rigid protein interior towards the solvent-exposed, loosely packed, and flexible protein surface. Moreover, active sites and protein-protein interfaces exert additional evolutionary constraints, in the form of either positive or purifying selection, and these constraints seem to extend beyond the immediate residues involved in the protein's function.

Yet a complete and accurate predictive model of rate variation remains elusive. Our best current models can explain only ~60% of the observed rate variation, and only in some structures. Model performance varies widely among different proteins, for unknown reasons. Thus, while the field has made considerable progress, many important questions remain (Box 2). To make further progress, we will have to pursue three distinct research areas: First, we need to improve rate estimates by developing better inference methods and by quantifying the errors of these estimates for realistic data sets. Second, we must try to improve predictions by finding yet undiscovered relevant molecular traits. Third, to advance mechanistic understanding, we need further research on theoretical models. All of these efforts will likely benefit from stronger integration with experimental work, as discussed in Box 1. Some questions that could help orient future research in these three areas are listed in Box 2.

### The three main challenges

While the field of rate estimation is mature, we still see ample room for improvement. In particular, rate estimation is subject to both stochastic and systematic errors. Estimation methods are typically based on some model and assessed using data simulated using the same type of model. This practice serves to assess stochastic errors and their convergence with, for example, number of sequences and divergence<sup>29;51</sup>. It also serves to assess, for example, whether Bayesian or ML methods result in better estimates<sup>49;52</sup>. However, differences between the process that generated the actual data and the model used for analysis will lead to systematic errors<sup>119</sup>. Using the most rigorous statistical approaches and

increasing the amount of data cannot compensate for the use of incorrect models; on the contrary, it may lead to even more biased estimates<sup>120</sup>. Among the most important assumptions that may affect rate estimates are the codon or amino acid replacement model, prior rate distributions, the assumption that rates are constant over time and lineages, and independently evolving sites. Investigating the effects of violating such model assumptions is, we believe, the most important current challenge for improving rate inference methods.

Phenomenological models that combine predictors such as RSA and WCN are not in perfect agreement with observed rates, and the origin of the mismatch, especially of the wide variation of explanatory power among proteins, is unknown. RSA and WCN are the best currently-known predictors, but they are not the only ones. Even though other constraints, such as local flexibility, secondary structure, and side-chain hydrogen bonding, do not seem to have a large effect on determining the overall pattern of site-specific rates, these properties have been found to affect the evolutionary process of some sites, and further work in this area, in particular comparing and contrasting these quantities to RSA and WCN, may be worthwhile. More importantly, we will have to develop useful predictors that quantify functional constraints. Beyond the high conservation of a few sites directly involved in function, there is some evidence of functional constraints inducing longer range patterns. Measures such as the distance to the active site improve site-specific rate predictions in some cases, and these and similar functional predictors of site-specific rate variation are the obvious next direction for the field. Finally, phenomenological models can also be integrated directly into the rate-inference framework<sup>36;121</sup>, and such integrated models could provide both better rate estimates and novel insight into structural and functional evolutionary constraints.

Ultimately, we aim at a mechanistic understanding of protein evolution, derived as much as possible from first principles. We have described the two biophysical models that have been applied to the study of site-specific rates. These models are based on the idea that fitness depends on protein activity which, in turn, depends on stability changes. While one of the models depends on changes in the stability of the native conformation, the other depends on the (de)stabilization of an active conformation. Even if the active-stability model results in better predictions of site-specific rates, whether active-state stability or native-state stability or both are the primary drivers of site-specific evolution is not currently known. Moreover, there exists no framework currently to incorporate functional constraints into biophysical mechanistic models. For example, what is the biophysical origin of the increase of site-specific rates with distance to the active site? Including function explicitly in mechanistic biophysical models is one of the main challenges for further development of mechanistic models.

### **Other limitations of current work**

All the structural predictors of rate we have discussed here suffer from one important shortcoming: They ignore pairwise interactions between amino acids. While quantities such as solvent accessibility or packing density implicitly take into account the extent to which other amino acids are nearby, they cannot explicitly model the increased or decreased substitution rate at one site in response to a substitution at another site. Yet such co-

evolution among sites is well documented<sup>122-126</sup>, and it has been used to infer protein tertiary structure<sup>127-129</sup> and protein–protein interactions<sup>130;131</sup> from sequence alignments. How pairwise interactions among sites affect rate heterogeneity, however, is poorly understood. Indeed, studies have primarily focused on inferring structure from sequence alignments, and few attempts have been made to solve the inverse problem of predicting site co-variation from structure.

As a generic approach towards developing more accurate models of structural constraints, one could move away from simple summary statistics such as RSA or WCN and instead employ mechanistic, all-atom models of protein folding. In principle, we should be able to recover any evolutionary constraints attributable to protein folding stability from the detailed energetic models used in protein design algorithms, which naturally consider interactions among residues in the structure<sup>132</sup>. However, attempts to date have fallen short of expectations<sup>22;133;134</sup>, the simple quantities RSA and WCN perform much better in predicting site-specific rates than do sophisticated all-atom protein-design calculations. In particular, protein design underestimates the amount of co-variation among sites observed in natural protein sequences<sup>133</sup>. It also tends to overestimate the variability of buried sites and underestimate that of exposed sites<sup>134</sup>.

Importantly, we are limited in the extent to which we can compare findings from distinct existing studies, because they frequently employ widely diverging inference methods, data sets, or models. For example, while rates estimated with Rate4Site<sup>51</sup> should correlate with  $dN/dS$  estimates, the extent to which this relationship holds true has never been tested. Similarly, results found for highly diverged globular enzymes<sup>23;75</sup> may not be comparable to results found for viral proteins that have experienced little divergence<sup>22</sup>, and which may inherently possess distinct structural features relative to cellular proteins<sup>135</sup>. Also, while some studies are based on highly diverged data sets that include all major taxa<sup>23;75</sup>, others are based on a few closely related species<sup>18</sup>. Even though the mutational response of proteins of different major taxa seems to be universally distributed<sup>136</sup>, the extent to which patterns of rate variation among sites can be compared among different taxa has not been assessed. Thus, future work will have to test explicitly to what extent results obtained under one method, for one taxonomic group, or for one type of data carry over to very different methods, taxonomic groups, or data types.

Finally, we have focused here on the variation of evolutionary rates among sites while implicitly assuming that rates are constant. However, rates vary also over evolutionary time, a phenomenon called *heterotachy*<sup>137</sup>. If, as we saw, site-specific rates depend on structural and functional constraints, the rate of a site should change due to evolutionary divergence of protein structure and function. Indeed, heterotachy has been observed in relationship with functional<sup>26;138;139</sup> and nonfunctional divergence<sup>137</sup>. In addition, even under constant structural or functional constraints, changes in the amino acids of the local environment of a site also impact its substitution rate<sup>140</sup>. Addressing the challenges proposed here, especially the development of mechanistic biophysical models of evolution, coupled with studies of structural and functional divergence, should advance our understanding of the variation of evolutionary rates over time.

## Acknowledgments

J. E. is Principal Investigator of CONICET. This work was also supported in part by NIH grant F31 GM113622-01 to S. J. S and by NIH grant R01 GM088344, NSF Cooperative agreement DBI-0939454 (BEACON Center), and ARO grant W911NF-12-1-0390 to C. O. W.

## References

1. Zhang J, Yang JR. Determinants of the rate of protein sequence evolution. *Nat Rev Genet.* 2015; 16:409–420. [PubMed: 26055156]
2. Yang Z. Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: approximate methods. *J Mol Evol.* 1994; 39:306–314. [PubMed: 7932792]
3. Yang Z. Among-site rate variation and its impact on phylogenetic analyses. *Trends Ecol Evol.* 1996; 11:367–372. [PubMed: 21237881]
4. Lartillot N, Phillipe HA. Bayesian mixture model for across-site heterogeneities in the amino-acid replacement process. *Mol Biol Evol.* 2004; 21:1095–1109. [PubMed: 15014145]
5. Yang, Z. *Computational Molecular Evolution.* Oxford University Press; 2006.
6. Holder MT, Zwickl DJ, Dessimoz C. Evaluating the robustness of phylogenetic methods to among-site variability in substitution processes. *Phil Trans Royal Soc B.* 2008; 363:4013–4021.
7. Wang HC, Li K, Susko E, Roger AJ. A class frequency mixture model that adjusts for site-specific amino acid frequencies and improves inference of protein phylogeny. *BMC Evol Biol.* 2008; 8:331. [PubMed: 19087270]
8. Le SQ, Dang CC, Gascuel O. Modeling protein evolution with several amino acid replacement matrices depending on site rates. *Mol Biol Evol.* 2012; 29:2921–2936. [PubMed: 22491036]
9. Yang ZH, Nielsen R, Goldman N, Pedersen AMK. Codon-substitution models for heterogeneous selection pressure at amino acid sites. *Genetics.* 2000; 155:431–449. [PubMed: 10790415]
10. Buckley TR, Simon C, Chambers GK. Exploring among-site rate variation models in a maximum likelihood framework using empirical data: Effects of model assumptions on estimates of topology, branch lengths, and bootstrap support. *Syst Biol.* 2001; 50:67–86. [PubMed: 12116595]
11. Mayrose I, Friedman N, Pupko T. A gamma mixture model better accounts for among site rate heterogeneity. *Bioinformatics.* 2005; 21:ii151–ii158. [PubMed: 16204095]
12. Delpont W, Scheffler K, Gravenor MB, Muse SV, Kosakovsky Pond SL. Benchmarking multi-rate codon models. *PLOS ONE.* 2010; 5:e11587. [PubMed: 20657773]
13. Lartillot N. Probabilistic models of eukaryotic evolution: time for integration. *Phil Trans R Soc London B.* 2015; 370:20140338. [PubMed: 26323768]
14. Liberles DA, Teufel AI, Liu L, Stadler T. On the need for mechanistic models in computational genomics and metagenomics. *Genome Biol Evol.* 2013; 5:2008–2018. [PubMed: 24115604]
15. Perutz MF, Kendrew JC, Watson HC. Structure and function of haemoglobin: II. Some relations between polypeptide chain configuration and amino acid sequence. *J Mol Biol.* 1965; 13:669–678.
16. Kimura M, Ota T. On some principles governing molecular evolution. *Proc Natl Acad Sci USA.* 1974; 71:2848–52. [PubMed: 4527913]
17. Dean AM, Neuhauser C, Grenier E, Golding GB. The pattern of amino acid replacements in  $\alpha/\beta$ -barrels. *Mol Biol Evol.* 2002; 19:1846–1864. One of the first studies considering both structural and functional determinants of site-specific amino-acid substitution rates. [PubMed: 12411594]
18. Franzosa EA, Xia Y. Structural determinants of protein evolution are context-sensitive at the residue level. *Mol Biol Evol.* 2009; 26:2387–2395. This landmark study found that that site-specific rate (dN/dS) increases linearly with solvent accessibility in yeast. [PubMed: 19597162]
19. Shih CH, Hwang JK. Evolutionary information hidden in a single protein structure. *Proteins.* 2012; 80:1647–1657. [PubMed: 22454236]
20. Nevin Gerek Z, Kumar S, Banu Ozkan S. Structural dynamics flexibility informs function and evolution at a proteome scale. *Evol Appl.* 2013; 6:423–433. [PubMed: 23745135]

21. Marsh JA, Teichmann SA. Parallel dynamics and evolution: Protein conformational fluctuations and assembly reflect evolutionary changes in sequence and structure. *BioEssays*. 2014; 36:209–218. [PubMed: 24272815]
22. Shahmoradi A, et al. Predicting evolutionary site variability from structure in viral proteins: Buriedness, packing, flexibility, and design. *J Mol Evol*. 2014; 79:130–142. [PubMed: 25217382]
23. Yeh SW, et al. Site-specific structural constraints on protein sequence evolutionary divergence: local packing density versus solvent exposure. *Mol Biol Evol*. 2014; 31:135–139. First study showing that site-specific rates correlate more strongly with WCN than with RSA. [PubMed: 24109601]
24. Huang TT, Del Valle Marcos ML, Hwang JK, Echave J. A mechanistic stress model of protein evolution accounts for site-specific evolutionary rates and their relationship with packing density and flexibility. *BMC Evol Biol*. 2014; 14:78. This paper introduces the stress model of protein evolution, a biophysical model based on mutational changes of active-state stability. [PubMed: 24716445]
25. Echave J, Jackson EL, Wilke CO. Relationship between protein thermodynamic constraints and variation of evolutionary rates among sites. *Phys Biol*. 2015; 12:025002. Study of rate variation among sites using the native-stability threshold biophysical model. [PubMed: 25787027]
26. Meyer AG, Spielman SJ, Bedford T, Wilke CO. Time dependence of evolutionary metrics during the 2009 pandemic influenza virus outbreak. *Virus Evolution*. 2015; 1:vev006–10. [PubMed: 26770819]
27. Nielsen R. Mapping mutations on phylogenies. *Syst Biol*. 2002; 51:729–739. [PubMed: 12396587]
28. Kosakovsky Pond SL, Frost SDW. A simple hierarchical approach to modeling distributions of substitution rates. *Mol Biol Evol*. 2004; 22:223–234. [PubMed: 15483327]
29. Kosakovsky Pond SL, Frost SDW. Not so different after all: A comparison of methods for detecting amino acid sites under selection. *Mol Biol Evol*. 2005; 22:1208–1222. Landmark paper benchmarking different methods of site-specific rate inference. [PubMed: 15703242]
30. Lemey P, Minin VN, Bielejec F, Kosakovsky Pond SL, Suchard MA. A counting renaissance: combining stochastic mapping and empirical Bayes to quickly detect amino acid sites under positive selection. *Bioinformatics*. 2012; 28:3248–3256. [PubMed: 23064000]
31. Rodrigue N. On the statistical interpretation of site-specific variables in phylogeny-based substitution models. *Genetics*. 2013; 193:557–564. [PubMed: 23222651]
32. Valdar WS. Scoring residue conservation. *Proteins*. 2002; 48:227–241. [PubMed: 12112692]
33. Capra JA, Singh M. Predicting functionally important residues from sequence conservation. *Bioinformatics*. 2007; 23:1875–1882. [PubMed: 17519246]
34. Johansson F, Toh H. A comparative study of conservation and variation scores. *BMC Bioinformatics*. 2010; 11:311–388. [PubMed: 20534127]
35. Muse SV. Estimating synonymous and nonsynonymous substitution rates. *Mol Biol Evol*. 1996; 13:105–114. [PubMed: 8583885]
36. Meyer AG, Wilke CO. Integrating sequence variation and protein structure to identify sites under selection. *Mol Biol Evol*. 2013; 30:36–44. [PubMed: 22977116]
37. Li WH, Wu CI, Luo CC. A new method for estimating synonymous and nonsynonymous rates of nucleotide substitution consider the relative likelihood of nucleotide and codon changes. *Mol Biol Evol*. 1985; 2:150–174. [PubMed: 3916709]
38. Nei M, Gojobori T. Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions. *Mol Biol Evol*. 1986; 3:418–426. [PubMed: 3444411]
39. Yang Z, Nielsen R. Estimating synonymous and nonsynonymous substitution rates under realistic evolutionary models. *Mol Biol Evol*. 2000; 17:32–42. [PubMed: 10666704]
40. Meyer S, von Haeseler A. Identifying site-specific substitution rates. *Mol Biol Evol*. 2003; 20:182–189. [PubMed: 12598684]
41. Nielsen R, Yang Z. Likelihood models for detecting positively selected amino acid sites and applications to the HIV-1 envelope. *Genetics*. 148:929–936. [PubMed: 9539414]
42. Yang Z, Wong WSW, Nielsen R. Bayes Empirical Bayes inference of amino acid sites under positive selection. *Mol Biol Evol*. 2005; 22:1107–1118. [PubMed: 15689528]



43. Murrell B, et al. Detecting individual sites subject to episodic diversifying selection. *PLOS Genet.* 2012; 8:e1002764. [PubMed: 22807683]
44. Kosakovsky Pond SL, Frost SDW, Muse SV. HyPhy: hypothesis testing using phylogenies. *Bioinformatics.* 2005; 21:676–679. [PubMed: 15509596]
45. Delpont W, Poon AFY, Frost SDW, Kosakovsky Pond SL. Datamonkey 2010 a suite of phylogenetic analysis tools for evolutionary biology. *Bioinformatics.* 2010; 26:2455–2457. [PubMed: 20671151]
46. Yang Z. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol.* 2007; 24:1586–1591. [PubMed: 17483113]
47. Goldman N, Yang Z. A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Mol Biol Evol.* 1994; 11:725–736. [PubMed: 7968486]
48. Drummond AJ, Rambaut A. BEAST : Bayesian evolutionary analysis by sampling trees. *BMC Evol Biol.* 2007; 7:214. [PubMed: 17996036]
49. Murrell B, et al. FUBAR: A Fast, Unconstrained Bayesian AppRoximation for inferring selection. *Mol Biol Evol.* 2013; 30:1196–1205. This paper introduces an extremely rapid yet accurate method to infer dN/dS. [PubMed: 23420840]
50. Angelis K, dos Reis M, Yang Z. Bayesian estimation of nonsynonymous/synonymous rate ratios for pairwise sequence comparisons. *Mol Biol Evol.* 2014; 31:1902–1913. [PubMed: 24748652]
51. Pupko T, Bell RE, Mayrose I, Glaser F, Ben-Tal N. Rate4Site: an algorithmic tool for the identification of functional regions in proteins by surface mapping of evolutionary determinants within their homologues. *Bioinformatics.* 2002; 18 Suppl 1:S71–S77. This paper introduced the Rate4Site method that is now widely used to calculate site-specific rates from amino-acid sequence data. [PubMed: 12169533]
52. Mayrose I, Graur D, Ben-Tal N, Pupko T. Comparison of site-specific rate-inference methods for protein sequences: empirical Bayesian methods are superior. *Mol Biol Evol.* 2004; 21:1781–1791. [PubMed: 15201400]
53. Fernandes AD, Atchley WR. Site-specific evolutionary rates in proteins are better modeled as non-independent and strictly relative. *Bioinformatics.* 2008; 24:2177–2183. [PubMed: 18662926]
54. Huang YF, Golding GB. Phylogenetic Gaussian process model for the inference of functionally important regions in protein tertiary structures. *PLOS Comput Biol.* 2014; 10:e1003429–12. [PubMed: 24453956]
55. Huang YF, Golding GB. FuncPatch : A web server for the fast bayesian inference of conserved functional patches in protein 3D structures. *Bioinformatics.* 2015; 31:523–531. [PubMed: 25322839]
56. Yang JR, Liao BY, Zhuang SM, Zhang J. Protein misinteraction avoidance causes highly expressed proteins to evolve slowly. *Proc Natl Acad Sci USA.* 2012; 109:E831–E840. [PubMed: 22416125]
57. Tien MZ, Meyer AG, Sydykova DK, Spielman SJ, Wilke CO. Maximum allowed solvent accessibilities of residues in proteins. *PLOS ONE.* 2013; 8:e80635. This paper provides accurate normalization constants required for the calculation of relative solvent accessibility. [PubMed: 24278298]
58. Hubbard TJ, Blundell TL. Comparison of solvent-inaccessible cores of homologous proteins: definitions useful for protein modelling. *Protein Eng.* 1987; 1:159–171. [PubMed: 3507702]
59. Lim WA, Sauer RT. Alternative packing arrangements in the hydrophobic core of lambda repressor. *Nature.* 1989; 339:31–36. [PubMed: 2524006]
60. Overington J, Johnson MS, Sali A, Blundell TL. Tertiary structural constraints on protein evolutionary diversity: templates, key residues and structure prediction. *Proc Biol Sci.* 1990; 241:132–145. [PubMed: 1978340]
61. Topham CM, et al. Fragment ranking in modelling of protein structure. Conformationally constrained environmental amino acid substitution tables. *J Mol Biol.* 1993; 229:194–220. [PubMed: 8421300]
62. Wako H, Blundell TL. Use of amino acid environment-dependent substitution tables and conformational propensities in structure prediction from aligned sequences of homologous proteins. I. solvent accessibility classes. *J Mol Biol.* 1994; 238:682–692. [PubMed: 8182743]

63. Koshi JM, Goldstein RA. Context-dependent optimal substitution matrices. *Protein Eng.* 1995; 8:641–645. [PubMed: 8577693]
64. Goldman N, Thorne JL, Jones DT. Assessing the impact of secondary structure and solvent accessibility on protein evolution. *Genetics.* 1998; 149:445–458. [PubMed: 9584116]
65. Conant GC, Stadler PF. Solvent exposure imparts similar selective pressures across a range of yeast proteins. *Mol Biol Evol.* 2009; 26:1155–1161. [PubMed: 19233963]
66. Ramsey DC, Scherrer MP, Zhou T, Wilke CO. The relationship between relative solvent accessibility and evolutionary rate in protein evolution. *Genetics.* 2011; 188:479–488. [PubMed: 21467571]
67. Scherrer MP, Meyer AG, Wilke CO. Modeling coding-sequence evolution within the context of residue solvent accessibility. *BMC Evol Biol.* 2012; 12:179. [PubMed: 22967129]
68. Franzosa EA, Xia Y. Independent effects of protein core size and expression on residue-level structure–evolution relationships. *PLOS ONE.* 2012; 7:e46602. [PubMed: 23056364]
69. Lin CP, et al. Deriving protein dynamical properties from weighted protein contact number. *Proteins.* 2008; 72:929–935. [PubMed: 18300253]
70. England JL, Shakhnovich E. Structural determinant of protein designability. *Phys Rev Lett.* 2003; 90:218101. [PubMed: 12786593]
71. Bloom JD, Drummond DA, Arnold FH, Wilke CO. Structural determinants of the rate of protein evolution in yeast. *Mol Biol Evol.* 2006; 23:1751–1761. [PubMed: 16782762]
72. Shakhnovich B, Deeds E, Delisi C, Shakhnovich E. Protein structure and evolutionary history determine sequence space topology. *Genome Res.* 2005; 15:385–392. [PubMed: 15741509]
73. Zhou T, Drummond DA, Wilke CO. Contact density affects protein evolutionary rate from bacteria to animals. *J Mol Evol.* 2008; 66:395–404. [PubMed: 18379715]
74. Yeh SW, et al. Local packing density is the main structural determinant of the rate of protein sequence evolution at site level. *BioMed Res Int.* 2014:572409. [PubMed: 25121105]
75. Marcos ML, Echave J. Too packed to change: side-chain packing and site-specific substitution rates in protein evolution. *PeerJ.* 2015; 3:e911. [PubMed: 25922797]
76. Mugal CF, Wolf JBW, Kaj I. Why time matters: Codon evolution and the temporal dynamics of dN/dS. *Mol Biol Evol.* 2014; 31:212–231. [PubMed: 24129904]
77. Liu Y, Bahar I. Sequence evolution correlates with structural dynamics. *Mol Biol Evol.* 2012; 29:2253–2263. Study of the correlation between flexibility and site-specific sequence entropy. [PubMed: 22427707]
78. Halle B. Flexibility and packing in proteins. *Proc Natl Acad Sci U S A.* 2002; 99:1274–1279. [PubMed: 11818549]
79. Liao H, Yeh W, Chiang D, Jernigan RL, Lustig B. Protein sequence entropy is closely related to packing density and hydrophobicity. *Protein Eng Des Sel.* 2005; 18:59–64. [PubMed: 15788422]
80. Worth CL, Gong S, Blundell TL. Structural and functional constraints in the evolution of protein families. *Nat Rev Mol Cell Biol.* 2009; 10:709–720. [PubMed: 19756040]
81. Bustamante CD, Townsend JP, Hartl DL. Solvent accessibility and purifying selection within proteins of *Escherichia coli* and *Salmonella enterica*. *Mol Biol Evol.* 2000; 17:301–308. [PubMed: 10677853]
82. Brown CJ, et al. Evolutionary rate heterogeneity in proteins with long disordered regions. *J Mol Evol.* 2002; 55:104–110. [PubMed: 12165847]
83. Brown CJ, Johnson AK, Daughdrill GW. Comparing models of evolution for ordered and disordered proteins. *Mol Biol Evol.* 2010; 27:609–621. [PubMed: 19923193]
84. Tóth-Petróczy A, Tawfik DS. Slow protein evolutionary rates are dictated by surface-core association. *Proc Natl Acad Sci U S A.* 2011; 108:11151–11156. Systematic study of the distributions of site-specific rates for yeast proteins. [PubMed: 21690394]
85. Finkelstein AV, Ivankov DN, Garbuzynskiy SO, Galzitskaya OV. Understanding the folding rates and folding nuclei of globular proteins. *Curr Protein Pept Sci.* 2007; 8:521–536. [PubMed: 18220841]
86. Ptitsyn OB. Protein folding and protein evolution: Common folding nucleus in different subfamilies of c-type cytochromes? *J Mol Biol.* 1998; 278:655–666. [PubMed: 9600846]

87. Mirny L, Shakhnovich E. Evolutionary conservation of the folding nucleus. *J Mol Biol.* 2001; 308:123–129. [PubMed: 11327757]
88. Larson SM, Ruczinski I, Davidson AR, Baker D, Plaxco KW. Residues participating in the protein folding nucleus do not exhibit preferential evolutionary conservation. *J Mol Biol.* 2002; 316:225–233. Study that shows that sites involved in the folding nucleus are not particularly conserved. [PubMed: 11851333]
89. Tseng YY, Liang J. Are residues in a protein folding nucleus evolutionarily conserved? *J Mol Biol.* 2004; 335:869–880. [PubMed: 14698285]
90. Drummond DA, Bloom JD, Adami C, Wilke CO, Arnold FH. Why highly expressed proteins evolve slowly. *Proc Natl Acad Sci USA.* 2005; 102:14338–14343. [PubMed: 16176987]
91. Franzosa EA, Xue R, Xia Y. Quantitative residue-level structure–evolution relationships in the yeast membrane proteome. *Genome Biol Evol.* 2013; 5:734–744. [PubMed: 23512408]
92. Spielman SJ, Wilke CO. Membrane environment imposes unique selection pressures on transmembrane domains of G protein-coupled receptors. *J Mol Evol.* 2013; 76:172–182. [PubMed: 23355009]
93. Bartlett GJ, Porter CT, Borkakoti N, Thornton JM. Analysis of catalytic residues in enzyme active sites. *J Mol Biol.* 2002; 324:105–121. [PubMed: 12421562]
94. Chelliah V, Chen L, Blundell TL, Lovell SC. Distinguishing structural and functional restraints in evolution in order to identify interaction sites. *J Mol Biol.* 2004; 342:1487–1504. [PubMed: 15364576]
95. McLaughlin RN Jr, Poelwijk FJ, Raman A, Gosal WS, Ranganathan R. The spatial architecture of protein function and adaptation. *Nature.* 2012; 490:138–142. [PubMed: 23041932]
96. Mintseris J, Weng Z. Structure, function, and evolution of transient and obligate protein–protein interactions. *Proc Natl Acad Sci USA.* 2005; 102:10930–10935. This paper shows that sites that participate in obligate protein-protein interactions are more conserved than those involved in transient interactions. [PubMed: 16043700]
97. Kim PM, Lu LJ, Xia Y, Gerstein MB. Relating three-dimensional structures to protein networks provides evolutionary insights. *Science.* 2006; 314:1938–1941. [PubMed: 17185604]
98. Huang YW, Chang CM, Lee CW, Hwang JK. The conservation profile of a protein bears the imprint of the molecule that is evolutionarily coupled to the protein. *Proteins.* 2015; 83:1407–1413. [PubMed: 25846748]
99. Kachroo AH, et al. Systematic humanization of yeast genes reveals conserved functions and genetic modularity. *Science.* 2015; 348:921–925. [PubMed: 25999509]
100. Glaser F, Morris RJ, Najmanovich RJ, Laskowski RA, Thornton JM. A method for localizing ligand binding pockets in protein structures. *Proteins Struct Funct Genet.* 2006; 62:479–488. [PubMed: 16304646]
101. Capra JA, Laskowski RA, Thornton JM, Singh M, Funkhouser TA. Predicting protein ligand binding sites by combining evolutionary sequence conservation and 3D structure. *PLoS Comput Biol.* 2009; 5:e1000585. [PubMed: 19997483]
102. Yang JS, Seo SW, Jang S, Jung GY, Kim S. Rational engineering of enzyme allosteric regulation through sequence evolution analysis. *PLoS Comput Biol.* 2012; 8:e1002612–10. [PubMed: 22807670]
103. Hill RE, Hastie ND. Accelerated evolution in the reactive centre regions of serine protease inhibitors. *Nature.* 1987; 326:96–99. [PubMed: 3493437]
104. Hughes AL, Nei M. Pattern of nucleotide substitution at major histocompatibility complex class I loci reveals overdominant selection. *Nature.* 1988; 335:167–170. [PubMed: 3412472]
105. Bush RM, Fitch WM, Bender CA, Cox NJ. Positive selection on the H3 hemagglutinin gene of human influenza virus A. *Mol Biol Evol.* 1999; 16:1457–1465. [PubMed: 10555276]
106. Shih AC, Hsiao T, Ho M, Li W. Simultaneous amino acid substitutions at antigenic sites drive influenza a hemagglutinin evolution. *Proc Natl Acad Sci USA.* 2007; 104:6283–6288. [PubMed: 17395716]
107. Pan K, Deem MW. Quantifying selection and diversity in viruses by entropy methods, with application to the haemagglutinin of H3N2 influenza. *J Roy Soc Interface.* 2011; 8:1644–1653. [PubMed: 21543352]

108. Tusche C, Steinbrück L, McHardy AC. Detecting patches of protein sites of influenza A viruses under positive selection. *Mol Biol Evol.* 2012; 29:2063–2071. [PubMed: 22427709]
109. Meyer AG, Wilke CO. Geometric constraints dominate the antigenic evolution of influenza H3N2 hemagglutinin. *PLOS Pathog.* 2015; 11:e1004940. [PubMed: 26020774]
110. Liberles DA, et al. The interface of protein structure, protein biophysics, and molecular evolution. *Protein Sci.* 2012; 21:769–785. [PubMed: 22528593]
111. Harms MJ, Thornton JW. Evolutionary biochemistry: revealing the historical and physical causes of protein properties. *Nat Rev Genet.* 2013; 14:559–571. [PubMed: 23864121]
112. Zhou H, Zhou Y. Quantifying the effect of burial of amino acid residues on protein stability. *Proteins.* 2004; 322:315–322. [PubMed: 14696193]
113. Shaytan AK, Shaitan KV, Khokhlov AR. Solvent accessible surface area of amino acid residues in globular proteins: Correlation of apparent transfer free energies with experimental hydrophobicity scales. *Biomacromolecules.* 2009; 10:1224–1237. [PubMed: 19334678]
114. Bloom JD, Glassman MJ. Inferring stabilizing mutations from protein phylogenies: Application to influenza hemagglutinin. *PLOS Comput Biol.* 2009; 5:e1000349. [PubMed: 19381264]
115. Wylie SC, Shakhnovich EI. A biophysical protein folding model accounts for most mutational fitness effects in viruses. *Proc Natl Acad Sci USA.* 2011; 108:9916–9921. [PubMed: 21610162]
116. Wylie SC, Shakhnovich EI. Mutation induced extinction in finite populations: Lethal mutagenesis and lethal isolation. *PLOS Comput Biol.* 2012; 8:e1002609. [PubMed: 22876168]
117. Guerois R, Nielsen JE, Serrano L. Predicting changes in the stability of proteins and protein complexes: a study of more than 1000 mutations. *J Mol Biol.* 2002; 320:369–387. [PubMed: 12079393]
118. Yang L, Song G, Jernigan RL. Protein elastic network models and the ranges of cooperativity. *Proc Natl Acad Sci USA.* 2009; 106:12347–52. [PubMed: 19617554]
119. Spielman SJ, Wilke CO. The relationship between  $dN/dS$  and scaled selection coefficients. *Mol Biol Evol.* 2015:1097–1108. This paper establishes a mathematical relationship between mutation–selection models and  $dN/dS$  ratios. [PubMed: 25576365]
120. Kolaczowski B, Thornton JW. Performance of maximum parsimony and likelihood phylogenetics when evolution is heterogeneous. *Nature.* 2004; 431:980–984. [PubMed: 15496922]
121. Kleinman CL, Rodrigue N, Lartillot N, Philippe H. Statistical potentials for improved structurally constrained evolutionary models. *Mol Biol Evol.* 2010; 27:1546–1560. [PubMed: 20159780]
122. Pagel M. Detecting correlated evolution on phylogenies: A general method for the comparative analysis of discrete characters. *Proc R Soc B-Biological Sci.* 1994; 255:37–45.
123. Muse SV. Evolutionary analyses of DNA sequences subject to constraints on secondary structure. *Genetics.* 1995; 139:1429–1439. [PubMed: 7768450]
124. Poon AFY, Lewis FI, Kosakovsky Pond SL, Frost SDW. An evolutionary-network model reveals stratified interactions in the V3 loop of the HIV-1 envelope. *PLoS Comput Biol.* 2007; 3:2279–2290.
125. Carlson JM, et al. Phylogenetic dependency networks: inferring patterns of CTL escape and codon covariation in HIV-1 Gag. *PLoS Comput Biol.* 2008; 4:e1000225. [PubMed: 19023406]
126. Kryazhimskiy S, Dushoff J, Bazykin GA, Plotkin JB. Prevalence of Epistasis in the Evolution of Influenza A Surface Proteins. *PLoS Genet.* 2011; 7:e1001301. [PubMed: 21390205]
127. Burger L, van Nimwegen E. Disentangling direct from indirect co-evolution of residues in protein alignments. 2010
128. Morcos F, et al. Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proc Natl Acad Sci USA.* 2011; 108:E1293–E1301. [PubMed: 22106262]
129. Jones DT, Buchan DW, Cozzetto D, Pontil M. PSICOV: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments. *Bioinformatics.* 2012; 28:184–190. [PubMed: 22101153]
130. Skerker JM, et al. Rewiring the specificity of two-component signal transduction systems. *Cell.* 2008; 133:1043–1054. [PubMed: 18555780]

131. Cheng RR, Morcos F, Levine H, Onuchic JN. Toward rationally redesigning bacterial two-component signaling systems using coevolutionary information. *Proc Natl Acad Sci USA*. 2014; 111:E563–E571. [PubMed: 24449878]
132. Leaver-Fay A, et al. ROSETTA3: an object-oriented software suite for the simulation and design of macromolecules. *Methods in Enzymol*. 2011; 487:545–574. [PubMed: 21187238]
133. Ollikainen N, Kortemme T. Computational protein design quantifies structural constraints on amino acid covariation. *PLOS Comput Biol*. 2013; 9:e1003313. [PubMed: 24244128]
134. Jackson EL, Ollikainen N, Covert AW, Kortemme T, Wilke CO. Amino-acid site variability among natural and designed proteins. *PeerJ*. 2013; 1:e211. [PubMed: 24255821]
135. Tokuriki N, Oldfield CJ, Uversky VN, Berezovsky IN, Tawfik DS. Do viral proteins possess unique biophysical features? *Trends Biochem Sci*. 2009; 34:53–59. [PubMed: 19062293]
136. Faure G, Koonin EV. Universal distribution of mutational effects on protein stability, uncoupling of protein robustness from sequence evolution and distinct evolutionary modes of prokaryotic and eukaryotic proteins. *Phys Biol*. 2015; 12:035001. [PubMed: 25927823]
137. Lopez P, Casane D, Philippe H. Heterotachy, an important process of protein evolution. *Mol Biol Evol*. 2002; 19:1–7. [PubMed: 11752184]
138. Gu X, Gu X. Statistical methods for testing functional divergence after gene duplication. *Mol Biol Evol*. 1999; 16:1664–74. [PubMed: 10605109]
139. Gu X. A simple statistical method for estimating type-ii (cluster-specific) functional divergence of protein sequences. *Mol Biol Evol*. 2006; 23:1937–1945. [PubMed: 16864604]
140. Pollock DD, Thiltgen G, Goldstein RA. Amino acid coevolution induces an evolutionary stokes shift. *Proc Natl Acad Sci*. 2012; 109:E1352–E1359. This paper introduces the concept of evolutionary Stokes shift: when an amino-acid substitution occurs at a site, its neighbors evolve more rapidly to accommodate the substitution. [PubMed: 22547823]
141. Leferink NGH, et al. Impact of residues remote from the catalytic centre on enzyme catalysis of copper nitrite reductase. *Nature Commun*. 2014; 5:4395. [PubMed: 25022223]
142. Fowler DM, et al. High-resolution mapping of protein sequence-function relationships. *Nat Methods*. 2010; 7:741–746. [PubMed: 20711194]
143. Fowler DM, Fields S. Deep mutational scanning: A new style of protein science. *Nat Methods*. 2014; 11:801–807. [PubMed: 25075907]
144. Romero PA, Tran TM, Abate AR. Dissecting enzyme function with microfluidic-based deep mutational scanning. *Proc Natl Acad Sci USA*. 2015; 112:7159–7164. [PubMed: 26040002]
145. Bloom JD. An experimentally informed evolutionary model improves phylogenetic fit to divergent lactamase homolog. *Mol Biol Evol*. 2014; 31:2753–2769. [PubMed: 25063439]
146. Abriata LA, Palzkill T, Dal Peraro M. How structural and physicochemical determinants shape sequence constraints in a functional enzyme. *PLOS ONE*. 2015; 10:e0118684. This paper shows one example (TEM lactamase) for which functional constraints relax slowly with distance to the active site. [PubMed: 25706742]
147. Bloom JD. An experimentally determined evolutionary model dramatically improves phylogenetic fit. *Mol Biol Evol*. 2014; 31:1956–1978. One of the first studies to parameterize a phylogenetic model with experimentally measured, site-specific parameters. [PubMed: 24859245]
148. Doud MB, Ashenberg O, Bloom JD. Site-specific amino acid preferences are mostly conserved in two closely related protein homologs. *Mol Biol Evol*. 2015; 32:2944–2960. [PubMed: 26226986]

## Biographies

**Julian Echave.** Escuela de Ciencia y Tecnología, Universidad de San Martín, San Martín, Buenos Aires, Argentina.

Julian Echave is Professor of Physical Chemistry at the National University of San Martín (Argentina) and Principal Investigator of the National Scientific and Technical Research Council (Argentina). He received his PhD in Physical Chemistry from the National

University of La Plata (Argentina) in 1991 for work on algebraic methods to solve the time-dependent Schrödinger equation. From 1991 to 1993, he was a postdoc in the group of David C. Clary at Cambridge University (UK) where he developed the quantum theory of four-atom chemical reactions. From 1993, he has studied the dynamics of various physical, chemical, and biological processes. He is currently interested in the biophysical mechanisms that shape the evolutionary divergence of protein sequence, structure, and dynamics.

**Stephanie J. Spielman.** Department of Integrative Biology, The University of Texas at Austin.

Stephanie Spielman is a PhD Candidate in the Ecology, Evolution, and Behavior program at The University of Texas at Austin, working under the supervision of Claus Wilke. She received her Sc.B. in Biology, with Honors, from Brown University in 2010. Her PhD research focuses on the development of and relationships among statistical models of coding-sequence evolution, and she is broadly interested in molecular evolution and phylogenetics.

**Claus O. Wilke.** Department of Integrative Biology, The University of Texas at Austin. Claus Wilke is Professor of Integrative Biology at The University of Texas at Austin. He received his PhD in Theoretical Physics from the University of Bochum in Germany in 1999, for work on quantitative models of evolutionary processes. From 2000 to 2005, he was a postdoc in the Adami lab at Caltech, where he received training in biological physics, evolutionary biology, and artificial life. Claus Wilke's current research interests focus on molecular evolution, structural biology, systems biology, and viral evolution.

## Glossary

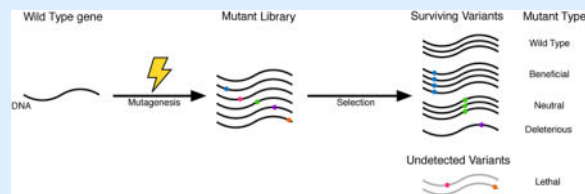
<b>Substitution</b>	Mutation that has spread to all members of the population (i.e., has fixed), substituting the ancestral variant.
<b>evolutionary rate</b>	Number of substitutions (fixed mutations) per unit of evolutionary time.
<b>non-synonymous mutation</b>	DNA mutation that changes from a codon that codes for one amino acid to a codon that codes for a different amino acid.
<b>synonymous mutation</b>	DNA mutation that changes from a codon that codes for one amino acid to a codon that codes for the same amino acid.
<b>non-synonymous evolutionary rate (<math>dN</math>)</b>	Rate at which non-synonymous substitutions (fixed mutations) occur per unit of evolutionary time.
<b>synonymous evolutionary rate (<math>dS</math>)</b>	Rate at which synonymous substitutions (fixed mutations) occur per unit of evolutionary time.
<b><math>dN/dS</math></b>	Ratio of non-synonymous to synonymous evolutionary rates.
<b>positive selection</b>	Fixation of mutations that increase fitness (adaptive).

<b>purifying selection</b>	Loss of mutations that decrease fitness (deleterious).
<b>structural constraint</b>	Structural feature that correlates with sequence conservation (e.g. solvent accessibility).
<b>functional constraint</b>	Functional features that correlates with sequence conservation (e.g. involvement in the active site).
<b>Rate4Site</b>	Popular software to estimate relative site-specific rates from amino-acid sequence data.
<b>solvent accessible surface area (SASA)</b>	surface area of a given residue that is accessible to water.
<b>accessible surface area (ASA)</b>	same as SASA.
<b>relative solvent accessibility (RSA)</b>	Measures the proportion of an amino acid's surface that is accessible to solvent (i.e., water) in the folded protein structure, from 0 (completely inaccessible) to 1 (completely accessible). Calculated as the ratio of the SASA of a given residue in the protein structure and the maximum SASA of that residue in a fully solvent-accessible conformation.
<b>contact number (CN)</b>	Number of neighboring residues present in a protein structure within a given distance (e.g., 10Å) from a focal residue.
<b>weighted contact number (WCN)</b>	Similar to contact number, but the neighboring residues are weighted by their inverse square distance to the focal residue, and all residues in a structure are considered to be neighboring residues.
<b>mean square fluctuation (MSF)</b>	Time-average of the square norm of the vector that connects the instantaneous coordinates of a site to its equilibrium coordinates; measures the amount of movement a residue undergoes over time.
<b>B factor</b>	Quantity that measures the amount of thermal motion of an atom in a protein crystal structure; also referred to as “temperature factor”,
<b><math>G</math></b>	Mutational change of stability; the folding free energy difference between mutant and wild type when each is in its own native conformation.
<b><math>G^*</math></b>	Mutational change of stability of the active conformation; free energy difference between the active conformation of the mutant and the active conformation of the wild type.
<b><math>G^\ddagger</math></b>	Mutational change of the activation free energy; difference between mutant and wild type of the free energy needed to deform the protein from the native into the active conformation.

**Box 1****Experimental approaches to measure site-specific variation**

Most of the work quantifying site-specific rate variation has been conducted computationally; however, a growing body of literature has emerged that takes an experimental approach to the question. Specifically, several experimental lines of research have sought to determine site-specific amino-acid preferences and/or tolerance to mutations. These quantities are intimately tied to evolutionary rate: sites which are more tolerant to mutation, or at which more amino-acids are selectively tolerated, will generally evolve more rapidly. Conversely, sites with low mutational tolerance will evolve more slowly<sup>66;119</sup>. Results from experimental work sampling all possible mutations across all residues in a given protein have supported these theoretical predictions. For example, McLaughlin, Jr., et al.<sup>95</sup> have shown that functionally important residues are generally less mutationally tolerant than are residues with less stringent functional constraint. Leferink et al.<sup>141</sup> demonstrated that mutations which increase solvent accessibility at an active site have strong influences on an enzyme's catalytic efficiency, demonstrating a tight relationship between evolutionary rate and function.

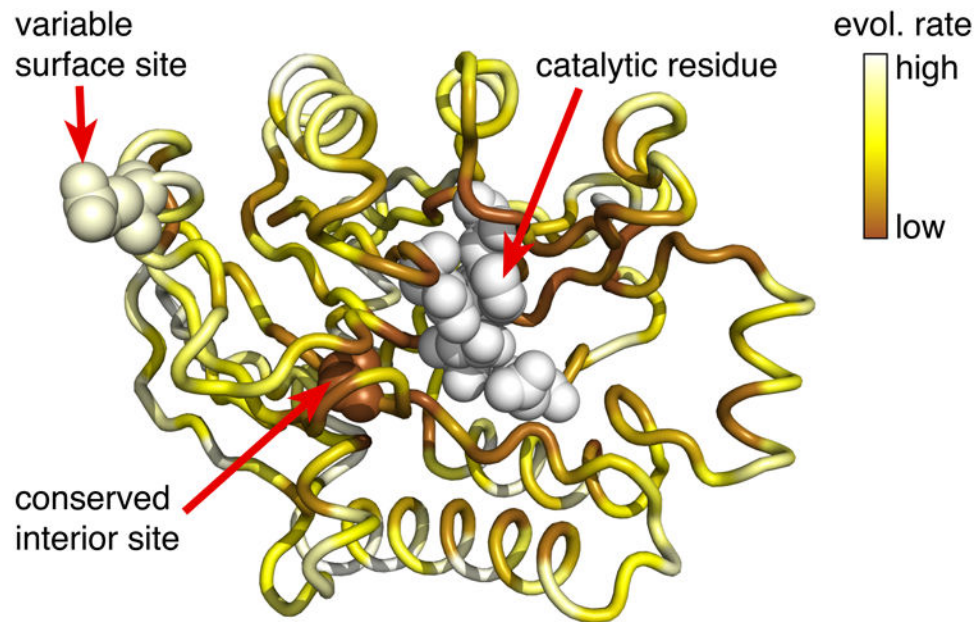
One new and powerful approach to addressing site-specific properties in proteins is deep mutational scanning, an experimental approach which samples as many as one million protein variants at a time<sup>142-144</sup>. Under deep mutational scanning, many different variants of a given gene are subjected to selection in a high-throughput procedure. By measuring the relative enrichment or depletion of variants after selection, this procedure allows for precise quantification of the gene variants' relative fitnesses (see figure). Deep mutational scanning studies conducted on proteins from bacteria<sup>145;146</sup> and viruses<sup>147;148</sup> have revealed extensive heterogeneity in mutational tolerance within a given protein, and that tolerance generally correlates with solvent accessibility. Moreover, deep mutational scanning on TEM lactamase has shown that residues near active sites can sustain very few substitutions<sup>146</sup>. Finally, Bloom linked the experimentally measured mutational tolerances to evolutionary models with site heterogeneity and showed that these experimentally informed models better account for observed variation in natural sequences than do standard phylogenetic models<sup>145;147</sup>.



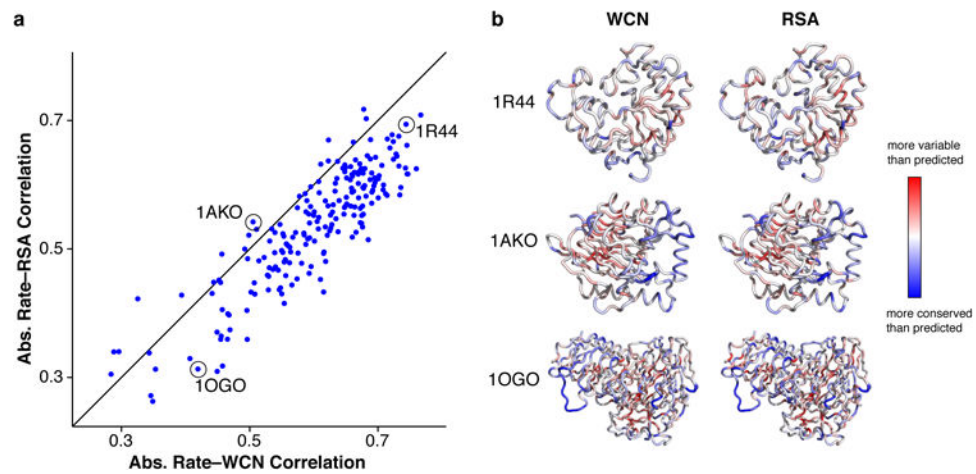


**Box 2****Open questions**

1. How accurate are current rate estimates, and how can we improve them?
  - Can we quantify the expected errors and biases of rate estimates?
  - How do rate estimates depend on the number of sequences and the degree of sequence divergence?
  - How robust are rate estimates with respect to violation of model assumptions, such as prior rate distributions and site independence?
2. Do we know all the molecular determinants of site rates?
  - How much of the variation of rates among sites do the factors we currently know actually explain?
  - Is the unexplained variation due to unidentified factors or errors of rate estimates or unexplainable variation (i.e., noise and biases of estimates)?
  - What other molecular features affect evolutionary rates? How can we incorporate the effect of specific functional features, such as an active site, into quantitative predictors of rates?
3. What are the mechanisms that produce site-specific evolutionary rate variation?
  - Does natural selection favor more stable proteins?
  - Is there a stability optimum due to stability–activity trade-offs?
  - Is there a stability threshold above which all mutants are neutral?
  - What is more important, stability or the correct active-site conformation?
  - How can we incorporate protein function explicitly into mechanistic models?

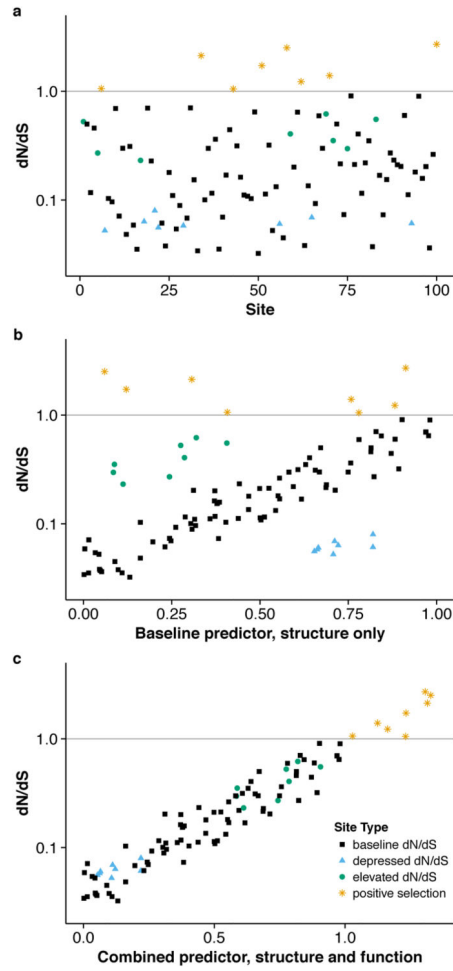


**Figure 1. Structural and functional constraints shape site-specific evolutionary divergence**  
Substitution rates for Exonuclease III of *Escherichia coli* are mapped onto its structure (pdb code 1AKO) using a divergence scale that goes from darker low rates to lighter high rates. Due to structural constraints, substitution rates are low in the protein interior and high on the surface. Residues close to the catalytic sites (visualized in gray) also evolve slowly, likely because of functional constraints. Evolutionary rate data are taken from Ref. 75.

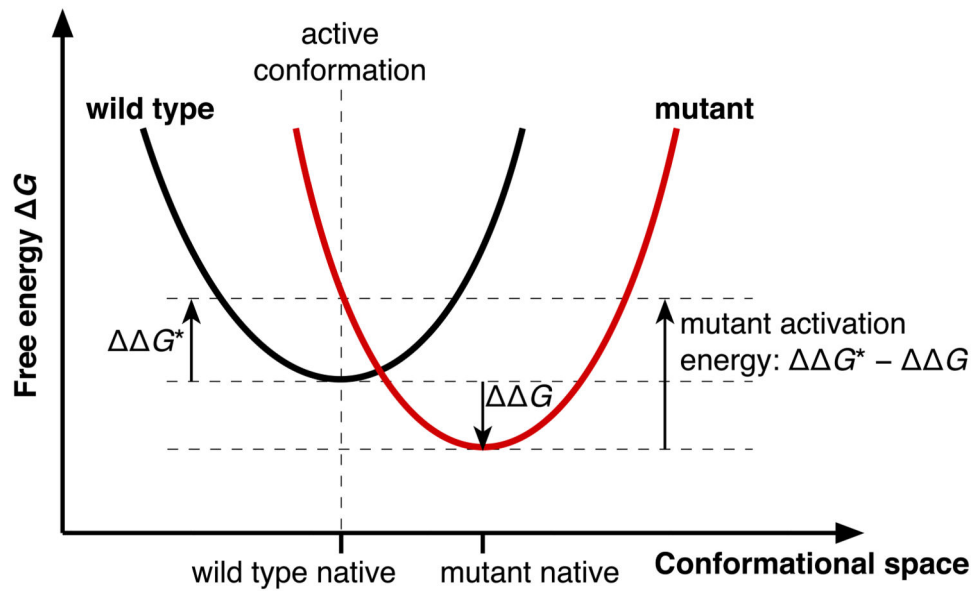


**Figure 2. Weighted contact number (WCN) correlates more strongly with site-specific rate than relative solvent accessibility (RSA) does**

(a) Absolute rate–RSA correlation vs. absolute rate–WCN correlation for 209 enzyme structures<sup>75</sup>. The solid line represents the  $x=y$  line. The rate–WCN correlations are systematically stronger than the rate–RSA correlations. However, for some proteins RSA performs better than WCN, as can be seen for example, for highlighted structure 1AKO. (b) Observed vs. predicted rate, mapped onto the backbone of three structures. Rate predictions were obtained from either WCN or RSA. The structures were chosen to represent low (1OGO), moderate (1AKO), and strong (1R44) structure–rate correlations, as highlighted in part (a). Colors represent the differences between observed and predicted rates at each site, with white representing a perfectly accurate prediction. As can be seen for structures 1AKO and 1OGO, poor predictions often coincide with surface loops that are more conserved than predicted from structure alone. These surface loops likely experience additional purifying selection due to function—compare e.g. the location of the conserved surface loops in 1AKO to the location of the protein's active site, as shown in Fig. 1.



**Figure 3. Predictors of evolutionary variation can help identify important sites in a protein**  
 (a) When plotted against the linear position (site) in the protein, site-specific evolutionary rates appear random. We can identify sites under positive selection ( $dN/dS > 1$ , indicated in yellow) but we cannot easily identify other important sites (here indicated in green and blue). (b) If we can identify a baseline predictor that captures the effect of protein structure on site-specific evolutionary rate, then sites that deviate from this baseline expectation clearly stand out (green: sites that evolve more rapidly than expected, blue: sites that are more conserved than expected). Such sites are likely functionally important. (c) If we can develop a predictor that can capture both the effects of structure and the effects of functional importance on evolutionary rate, then previously outlying sites appear to now follow the overall trend. This result indicates that we have identified the proper underlying reasons for why blue sites evolve slower and green or yellow sites faster than expected under the model of (b).



**Figure 4. Trade-off between native stability and active stability**

A mutation shifts the free energy landscape from that of the wild-type protein (black curve) to that of the mutant (red curve). The mutant has a different equilibrium conformation and its stability differs by an amount of  $\Delta\Delta G$  from that of the wild type (difference between red and black minima). The stability of the active conformation changes by an amount of  $\Delta\Delta G^*$  due to the mutation (difference between the intersections of the vertical “active conformation” line with the red and black curves). To function, the mutant protein must deform from its equilibrium conformation to the active conformation, which requires an activation energy  $G^\ddagger = G^* - G$ . We assume here that the wild type native structure is the active conformation. In this scenario, a mutation may stabilize the native state ( $\Delta\Delta G < 0$ ) yet destabilize the active state ( $\Delta\Delta G^* > 0$ ). Thus, even a stabilizing mutation can create an energy barrier that may reduce or eliminate proper protein function.

**Table 1**  
**Summary of quantities and/or biophysical effects observed to influence site-specific rate**

Rates in the cited studies have been estimated either from codon or from amino-acid data, as discussed in section “Estimation of site-specific rates”.

Quantity/physical effect	Effect on rate	References
<i>Structural constraints</i>		
Contact Number (CN, WCN)	Decreases with increasing CN/WCN	18;19;22-24;74;75;79
Relative Solvent Accessibility (RSA)	Increases with increasing RSA	17;18;23;24;64;66-68;74;75;81
Structural flexibility	Increases with increasing flexibility	20-22;24
Structural disorder	Increased in disordered regions	82-84
<i>Functional constraints</i>		
Protein–protein interfaces	Depressed in interface regions	18;96;97
Protein–nucleic acid interfaces	Depressed in interface regions	98
Catalytic sites	Depressed at and near catalytic sites	17
<i>Environmental constraints</i>		
Gene expression level	Decreases with increasing expression level, in particular at surface sites	67;68