



Published in final edited form as:

Phys Biol. ; 12(2): 025002. doi:10.1088/1478-3975/12/2/025002.

## Relationship between protein thermodynamic constraints and variation of evolutionary rates among sites

Julian Echave<sup>1</sup>, Eleisha L. Jackson<sup>2</sup>, and Claus O. Wilke<sup>2</sup>

Julian Echave: jechave@unsam.edu.ar; Claus O. Wilke: wilke@austin.utexas.edu

1

2

### Abstract

Evolutionary-rate variation among sites within proteins depends on functional and biophysical properties that constrain protein evolution. It is generally accepted that proteins must be able to fold stably in order to function. However, the relationship between stability constraints and among-sites rate variation is not well understood. Here, we present a biophysical model that links the thermodynamic stability changes due to mutations at sites in proteins ( $\Delta G$ ) to the rate at which mutations accumulate at those sites over evolutionary time. We find that such a “stability model” generally performs well, displaying correlations between predicted and empirically observed rates of up to 0.75 for some proteins. We further find that our model has comparable predictive power as does an alternative, recently proposed “stress model” that explains evolutionary-rate variation among sites in terms of the excess energy needed for mutants to adopt the correct active structure ( $\Delta G^*$ ). The two models make distinct predictions, though, and for some proteins the stability model outperforms the stress model and vice versa. We conclude that both stability and stress constrain site-specific sequence evolution in proteins.

### Keywords

protein evolution; rate variation among sites; biophysical model; thermodynamics; stability; stress

## 1. Introduction

The evolution of protein-coding genes is shaped by functional and biophysical constraints on the expressed proteins (Pal et al. 2006, Thorne 2007, Worth et al. 2009, Wilke & Drummond 2010, Grahn et al. 2011, Liberles et al. 2012). These constraints create patterns of rate variation *among* and *within* proteins. Among proteins, the primary determinant of rate variation is gene expression level (Drummond & Wilke 2008), though many other factors have been identified that also contribute to rate variation (Lemos et al. 2005, Xia et al. 2009, Liao et al. 2010, Pang et al. 2010). Within proteins, the primary determinants of rate variation seem to be linked to geometrical properties of the folded protein, in particular the Relative Solvent Accessibility (RSA) (Bustamante et al. 2000, Dean et al. 2002, Franzosa & Xia 2009, Ramsey et al. 2011, Shahmoradi et al. 2014) and the Local Packing Density (LPD) (Liao et al. 2005, Franzosa & Xia 2009, Yeh et al. 2014a, Yeh et al. 2014b) of sites in the three-dimensional protein structure.

To develop a mechanistic understanding of the causes that link geometrical properties, such as RSA and LPD, with site-specific rates of evolution, we need to develop explicit models of protein evolution. For example, recently a mechanistic “stress model” was proposed to explain the LPD–rate relationship (Huang et al. 2014). According to this stress model, LPD is a proxy of the stress energy  $G^*$ , a thermodynamic quantity that is a measure of the excess free energy needed for a folded mutant protein to adopt the correct active conformation. The stress model considers the effect of the stress free energy difference  $G^*$  but not that of possible mutational changes on global stability  $G$ . However, most proteins will function properly if they have folded *stably* into the correct conformation. To what extent stability constraints shape site-specific sequence evolution is not known.

Recent work has shown that describing protein evolution from the perspective of thermodynamic stability provides a wealth of insight into important aspects of protein evolution, such as the evolution of mutational robustness (Bloom et al. 2007), the origin of epistatic interactions (Bershtein et al. 2006, Gong et al. 2013), lethal mutagenesis (Chen & Shakhnovich 2009), determinants of evolutionary rate at protein level (Drummond & Wilke 2008, Serohijos et al. 2012), the evolution of novel function (Bloom et al. 2006, Tokuriki et al. 2008), and the expected equilibrium distributions of stability and the explanation of marginal stability (Taverna & Goldstein 2002, Goldstein 2011, Wylie & Shakhnovich 2011). Moreover, some studies suggest that  $G$ -based models are useful to study site-specific constraints. For example, Bloom & Glassman (2009) have shown that changes in stability upon mutation ( $G$  values) are intimately linked to the patterns of amino-acid substitutions observed over evolutionary divergence, to the extent that  $G$  values can actually be inferred with accuracy comparable to state-of-the art structure-based methods solely from an alignment of diverged protein sequences. More recently, Arenas et al. (2013) have used stability-based models to predict site-specific amino acid distributions. Despite the recognized importance of folding stability, stability-based models have not been used to predict the variation of evolutionary rates among sites.

Here, we investigate the relationship between mutational changes of stability and the site-dependency of rates of substitution. Following Bloom & Glassman (2009), we derived a neutral “stability model” of evolution which relates the  $G$ s due to mutations at a site with the site’s rate of substitution. For a diverse set of more than 200 enzymes, we compare the predicted rates with empirical rates (inferred from multiple sequence alignments) and with predictions of the stress model. The  $G$ -based and  $G^*$ -based predictions have on average similar correlations with empirical rates. However, the two models make significant independent contributions, which suggests that both stability and stress mould sequence divergence.

## 2. Stability Model: $G$ -based rates

Our stability model is based on earlier work by Bloom and coworkers (Bloom & Glassman 2009, Bloom et al. 2005). The core idea of Bloom’s model is that there is a stability threshold  $G^{\text{threshold}}$  such that all proteins more stable than the threshold are neutral (i.e. have all the same fitness) whereas all proteins less stable than the threshold are inviable

(have fitness = 0). Thus, if  $G$  is the stability of a protein, then its fitness  $f(G)$  is assumed to be:

$$f(\Delta G) = \begin{cases} 1 & \text{if } \Delta G \leq \Delta G^{\text{threshold}}, \\ 0 & \text{if } \Delta G > \Delta G^{\text{threshold}}. \end{cases} \quad (1)$$

It is convenient to define

$$\Delta G^{\text{extra}} = \Delta G - \Delta G^{\text{threshold}}, \quad (2)$$

so that

$$f(\Delta G^{\text{extra}}) = \begin{cases} 1 & \text{if } \Delta G^{\text{extra}} \leq 0, \\ 0 & \text{if } \Delta G^{\text{extra}} > 0. \end{cases} \quad (3)$$

We further assume that the mutational effect on stability of a mutation  $i \rightarrow j$  at site  $k$  is independent of the sequence background. We refer to this stability change as  $\Delta\Delta G_{ij}^k$ . Because of the assumption of sequence independence, the stability difference between two sequences can be written as

$$\Delta G(j_1, j_2, \dots) - \Delta G(i_1, i_2, \dots) = \sum_k \Delta\Delta G_{i_k j_k}^k \quad (4)$$

where  $i_1, i_2, \dots$  and  $j_1, j_2, \dots$  represent the amino acids of the two sequences, respectively. While this assumption cannot strictly be true, in practice it has worked well in several applications (e.g. Bloom et al. 2005, Bloom & Glassman 2009). The assumption is further supported by the observation that mutational effects on stability are frequently additive (Wells 1990, Serrano et al. 1993, Zhang et al. 1995) and tend to be conserved during evolution (Ashenberg et al. 2013).

Next we describe the evolutionary process. Throughout this work, we assume that the product of the protein-wide mutation rate  $\mu$  and the effective population size  $N_e$  is small,  $\mu N_e \ll 1$ . As a consequence, our populations are monomorphic, and we only have to track the evolution of a single representative sequence over time. We further assume that at most a single mutation arises at each time step.

The probability that a substitution  $i \rightarrow j$  occurs at site  $k$  in a single time step,  $Q_{ij}^k$ , can be written as the product of the probability that the mutation  $i \rightarrow j$  occurs,  $M_{ij}$ , and the probability it goes to fixation

$$Q_{ij}^k = M_{ij} \times p_{\text{fix}}. \quad (5)$$

Here, we have assumed that all sites experience the same mutational process, so that  $M_{ij}$  does not depend on  $k$ . Note that  $M_{ij}$  scales with the effective population size  $N_e$ , since all sequences in the population may mutate in one time step, and  $p_{\text{fix}}$  scales with  $1/N_e$ , because

we are modeling the case of neutral evolution (Eq. 3). Thus  $N_e$  cancels, and we can set it equal to 1 without loss of generality.

Under the assumption of neutral evolution, the fixation probability is either one or zero, depending on whether the mutation keeps the extra stability in the negative or not. Because we have previously assumed that stability effects are independent of the sequence background (Eq. 4), they are fully specified by  $i, j$ , and  $k$ . (In other words, a mutation from  $i$  to  $j$  at site  $k$  always has the same stability effect  $\Delta\Delta G_{ij}^k$ .) However, the extra stability after the mutation,  $\Delta G^{\text{extra}} + \Delta\Delta G_{ij}^k$ , depends on the sequence background through the value of  $G^{\text{extra}}$  before the mutation. From Eq. 3 we find the conditional fixation probability

$$p_{\text{fix}}(\Delta\Delta G_{ij}^k | \Delta G^{\text{extra}}) = \begin{cases} 1 & \text{if } \Delta G^{\text{extra}} + \Delta\Delta G_{ij}^k \leq 0, \\ 0 & \text{if } \Delta G^{\text{extra}} + \Delta\Delta G_{ij}^k > 0. \end{cases} \quad (6)$$

If  $\Delta G^{\text{extra}} + \Delta\Delta G_{ij}^k \leq 0$ , then the mutated protein is viable, and hence it fixes with probability 1. (Recall that we set  $N_e = 1$ .) By contrast, if  $\Delta G^{\text{extra}} + \Delta\Delta G_{ij}^k > 0$ , then the mutated protein is not viable and will not fix.

To proceed, we could write down a Markov process that keeps track of the extra stability at all time points (Bloom et al. 2007, Raval 2007). Instead, here we employ the “mean field” approximation of Bloom & Glassman (2009), in which we assume that  $G^{\text{extra}}$  before mutation is drawn randomly from the steady-state distribution of  $G^{\text{extra}}$  values,  $p_0(G^{\text{extra}})$ , so that we can write the unconditional fixation probability as

$$p_{\text{fix}}(\Delta\Delta G_{ij}^k) = \int p_{\text{fix}}(\Delta\Delta G_{ij}^k | \Delta G^{\text{extra}}) p_0(\Delta G^{\text{extra}}) d\Delta G^{\text{extra}}. \quad (7)$$

For  $p_0(G^{\text{extra}})$ , Bloom & Glassman (2009) make the ansatz that it has an exponential probability-density function  $p_0(G^{\text{extra}})$ :

$$p_0(\Delta G^{\text{extra}}) = \begin{cases} \alpha e^{\alpha \Delta G^{\text{extra}}} & \text{if } \Delta G^{\text{extra}} \leq 0, \\ 0 & \text{if } \Delta G^{\text{extra}} > 0. \end{cases} \quad (8)$$

where  $\alpha > 0$  is a free parameter. This form cannot be derived from first principles, but it is justified by visual inspection of the probability density functions obtained under simulations (Bloom et al. 2007) (but see Wylie & Shakhnovich 2011).

Inserting Eq. 6 and Eq. 8 into Eq. 7, we obtain

$$p_{\text{fix}}(\Delta\Delta G_{ij}^k) = \begin{cases} 1 & \text{if } \Delta\Delta G_{ij}^k \leq 0, \\ \int_{-\infty}^{-\Delta\Delta G_{ij}^k} \alpha e^{\alpha \Delta G^{\text{extra}}} & \text{if } \Delta\Delta G_{ij}^k > 0. \end{cases} \quad (9)$$

After taking the integral, we find

$$p_{\text{fix}}(\Delta\Delta G_{ij}^k) = \begin{cases} 1 & \text{if } \Delta\Delta G_{ij}^k \leq 0, \\ e^{-\alpha\Delta\Delta G_{ij}^k} & \text{if } \Delta\Delta G_{ij}^k > 0. \end{cases} \quad (10)$$

The stability model is completely specified by Eq. 5 and Eq. 10.

Next we consider the calculation of site-specific substitution rates. The substitution process at site  $k$  is described by a rate matrix  $\mathbf{Q}^k$  with elements

$$Q_{ij}^k = \begin{cases} M_{ij} \times p_{\text{fix}}(\Delta\Delta G_{ij}^k) & \text{if } i \neq j, \\ -\sum_{j \neq i} Q_{ij}^k & \text{if } i = j. \end{cases} \quad (11)$$

The stationary distribution  $\pi_i^k$  of the substitution process is given by the left null eigenvector of  $\mathbf{Q}^k$ , normalized such that  $\sum_i \pi_i^k = 1$ . The rate of substitution at site  $k$ ,  $K_{\text{stability}}^k$ , follows as

$$K_{\text{stability}}^k = \sum_i \sum_{j \neq i} \pi_i^k Q_{ij}^k = -\sum_i \pi_i^k Q_{ii}^k. \quad (12)$$

The subscript “stability” emphasizes that this rate estimate is calculated using the stability model.

In the case of symmetric mutations,  $M_{ji} = M_{ij}$ , the equilibrium frequencies can be expressed as

$$\pi_i^k = \frac{e^{-\alpha\Delta\Delta G_{0i}^k}}{\sum_j e^{-\alpha\Delta\Delta G_{0j}^k}}, \quad (13)$$

where  $\Delta\Delta G_{0i}^k$  is the stability change relative to an arbitrarily chosen reference amino acid at site  $k$ . In the limit of unbiased mutations,  $M_{ij} = \text{const}$  for  $i \neq j$ , the rate can be simplified to

$$K_{\text{stability}}^k = 2 \sum_i \pi_i^k [\text{rank}(-\pi_i^k) - 1]. \quad (14)$$

Here,  $\text{rank}(-\pi_i^k)$  represents the rank order of  $\pi_i^k$ , from largest to smallest. (The advantage of using Eq. 14 instead of Eq. 12 is that the latter contains a double-sum and hence is slower to evaluate.)

### 3. The Stress Model: $G^*$ -based rates

The stability model is based on the assumption that fitness depends on whether the protein is stable enough to fold, so that the probability of fixation of a mutation will depend on the difference of folding free energy between the mutant and the wild-type, each in their respective equilibrium conformations. A different mechanistic model, the “stress model,” was recently derived based on the idea that, to be viable, a mutant must not only be stable, it must also be able to adopt a correct active conformation (Huang et al. 2014). Following this idea, the fixation probability of a mutant was modeled as the mutant’s probability of

adopting the active conformation. According to this model, the rate of substitution for site  $k$  is

$$K_{\text{stress}}^k = a + b \langle \Delta \Delta G^* \rangle^k, \quad (15)$$

where  $G^* = G_{\text{mutant}}(\mathbf{r}_{\text{active}}) - G_{\text{wt}}(\mathbf{r}_{\text{active}})$  is the free energy difference between mutant and wild-type when both adopt the active conformation and  $\langle G^* \rangle^k$  is its average over random mutations at site  $k$ . Since in general the active conformation will not necessarily be the relaxed equilibrium conformation,  $G^*$  represents the energy needed to stress the protein into adopting the right active conformation.

Further assuming that the active conformation is the wild type's equilibrium conformation and approximating the free energy landscape using the parameter-free Anisotropic Network Model of Yang et al. (2009), it can be shown that  $\langle G^* \rangle^k \propto \text{WCN}^k$ , where

$\text{WCN}^k = \sum_{j \neq k} d_{kj}^{-2}$  is the Weighted Contact Number introduced by Lin et al. (2008) and found to be among the best structural predictors of site-dependent evolutionary rates (Yeh et al. 2014a, Yeh et al. 2014b). Because of the proportionality between  $\text{WCN}^k$  and  $\langle G^* \rangle^k$ , we can also write Eq. 15 as

$$K_{\text{stress}}^k = a + \tilde{b} \text{WCN}^k. \quad (16)$$

In practice, we obtain rates  $K_{\text{stress}}^k$  by calculating the  $\text{WCN}^k$  for each site  $k$  in a protein structure, fitting the linear expression  $a + \tilde{b} \text{WCN}^k$  to a set of empirically estimated rates, and then using Eq. 16 to calculate a predicted rate at each site.

It is worthwhile to keep in mind that while the stability model takes into account whether the mutant is able to fold, the stress model takes into account the probability that the mutant adopts the right conformation. In principle both factors can affect fitness independently and therefore may both have an influence on substitution rates. If this is the case, both models are incomplete: the stability model does not consider the effect of possible conformational changes as long as the mutant is stable and the stress model takes stability for granted and considers only the destabilization of the active structure.

## 4. Comparing the theoretical models with empirical data

### 4.1. Data set and calculation of empirical and predicted evolutionary rates

We tested our theory on the data set of Huang et al. (2014), which consists of 213 monomeric enzymes of known structure covering diverse structural and functional classes. Each structure is accompanied by up to 300 homologous sequences. In our analysis, we omitted four structures (1bbs, 1bs0, 1din, 1hpl) that had missing data at insertion sites. We aligned the homologous sequences for each structure with MAFFT (Multiple Alignment using Fast Fourier Transform) (Katoh et al. 2005, Katoh & Standley 2013). Using the resulting alignments as input, we inferred Maximum Likelihood phylogenetic trees with RAxML (Randomized Axelerated Maximum Likelihood), using the LG substitution matrix (named after Le and Gacuel) and the CAT model of rate heterogeneity (Stamatakis 2014).

For each structure, we then used the respective sequence alignment and phylogenetic tree to infer site-specific substitution rates with Rate4Site, using the empirical Bayesian method and the amino-acid Jukes-Cantor mutational model (aaJC) (Mayrose et al. 2004). The aaJC model poses equal probabilities for all amino-acid mutations, so that it is consistent with the theory presented in Section 2 and with the assumption of modeling amino-acid mutations as completely random perturbations made in the derivation of the stress model (Huang et al. 2014). Site-specific *relative* rates were obtained by dividing site-specific rates by their average over all sites of the protein, so that the mean relative rate of all sites in a protein was 1. In the following, we will refer to the rates inferred by Rate4Site as *empirical* rates, and will denote them by  $K_{R4S}$ . We will refer to the rates calculated according to the stability model ( $K_{\text{stability}}$ ) or the stress model ( $K_{\text{stress}}$ ) as *predicted* rates. If necessary, we will distinguish between the predictions of the stability and stress model using the terms  $G$ -predicted rates and  $G^*$ -predicted rates, respectively.

We calculated  $G$  values with the program FoldX, following the default protocol (Guerois et al. 2002, Schymkowitz et al. 2005). Specifically, we first optimized the energy for each structure, using the RepairPDB method. We then calculated a  $\Delta\Delta G_{0j}^k$  value for all possible 19 amino-acid substitutions at all sites in all proteins, using the PositionScan method, and considering the amino acid present in the PDB structure at each site as the reference amino acid at that site.

Rates predicted by the stability model were obtained using Eq. 13 and Eq. 14 either with  $\alpha = 1$  or with  $\alpha$  chosen specifically for each protein. To determine the appropriate scale factor  $\alpha$  for each protein, we maximized the correlation coefficient between the predicted site-specific rates as given by Eqs. 13 and 14 and the empirical site-specific rates as calculated by Rate4Site. To calculate the rates predicted by the stress model, we performed a linear fit between the site-dependent  $K_{R4S}$  and WCN for each protein, and then used Eq. 16 to calculate  $K_{\text{stress}}$  at each site.

All statistical analysis was carried out with R (R Core Team 2014). To fit the stability model to the data, we used the built-in function `optimize()` with default parameter settings. To fit the stress model to the data, we used the built-in function `lm()`. Correlation coefficients between predicted and empirical rates were calculated using `cor()` and partial correlations were obtained using the function `pcor.test()` of package `ppcor`.

All data and analysis scripts necessary to reproduce this work are available at: [https://github.com/wilkelab/therm\\_constraints\\_rate\\_variation/](https://github.com/wilkelab/therm_constraints_rate_variation/).

## 4.2. Relationship between empirical and predicted evolutionary rates

We found that rates predicted by the stability model correlate significantly with the empirical rates. Correlation coefficients ranged between 0.25 and 0.75, with a median of 0.57 (Figure 1A). Scale values  $\alpha$  fell between 0.52 and 2.63, with a median of 1.19. We further found that correlation coefficients and scale values were not correlated ( $r = 0.05$ ,  $P = 0.47$ ). To determine to what extent optimizing  $\alpha$  for each protein affected the resulting correlation coefficients, we also calculated correlation coefficients with  $\alpha = 1$  for all

proteins. We found that adjusting  $\alpha$  made only a small difference, resulting on average in an increase in correlation coefficient of 0.007 (Figure 1A).

We next investigated the functional relationship between empirical rates and rates predicted by the stability model. We pooled the data from all sites in all 209 proteins and calculated the joint distribution of the two rates. We also grouped sites into 20 bins of similar number of points using quantile breaks along the predicted rates axis. Figure 2 shows the joint distribution as well as the mean empirical rates and the 25% and 75% quantiles for each bin. The mean empirical rates fall nearly on top of the  $x = y$  line (which represents a perfect fit), with only a small amount of curvature around the mean predicted rate. The correlation between average empirical and predicted rates is  $r = 0.995$ , consistent with a very good linear fit. Despite the good fit of average rates, there is significant variation around  $x = y$ , as can be seen from the dispersion of the joint distribution around the  $x = y$  line and the error bars in Figure 2. The overall square correlation between  $G$ -predicted rates and empirical rates is  $r^2 = 0.31$ , so that 69% of the variance of empirical rates is not explained by the stability model.

Next, we compared the predictions of the stability model with those from the stress model (Huang et al. 2014), which describes site-specific evolutionary rates in terms of the increased stress that results in the protein's active conformation due to mutation ( $G^*$ ). In a protein-by-protein comparison, the stability model is somewhat better (dots above the  $x = y$  line in Figure 3) for 127 of the 209 proteins, a proportion significantly larger than 50% (binomial test: 61%,  $P = 0.002$ ). When considering all sites together, the two models perform comparably. The correlations between empirical and predicted rates for all sites are 0.56 with  $G$ -based predictions and 0.55 with  $G^*$ -based predictions. However, even though the two models perform comparably on average, there is substantial variation around the mean trend (Figure 3). For some proteins, the  $G$  model clearly outperforms the stress model and vice versa. Also, considering all sites, the partial correlations between empirical rates and predicted rates for one model controlling for the predictions of the other are 0.33 and 0.31 for the  $G$  model and the stress model, respectively. These values are large and highly significant ( $P \ll 10^{-3}$ ), showing that the predictions of the two models are quite independent and may be accounting for different constraints.

The relative independence of stress and stability as determinants of site-specific evolutionary rates suggests that considering both factors should improve predictions. To verify this hypothesis, we fit empirical rates to a linear combination of rates predicted from  $G^*$  and  $G$ . Considering all sites of all proteins, the two-variable model results in a square correlation  $R^2 = 0.38$ , approximately a 23% improvement over  $R^2 = 0.31$  of the stability model and  $R^2 = 0.30$  of the stress model. Both predictors in the two-variable model are highly significant ( $P < 10^{-15}$ ). These results further support the idea that stability and stress provide significant independent constraints to evolutionary divergence at site level.

All  $G$ -based predictions presented above used  $G$  values calculated by FoldX. It is possible that a different  $G$  predictor would yield substantially different results. In particular, even though FoldX is a state-of-the-art  $G$  predictor its predictions explain only 25% of the variance in measured  $G$  values (Potapov et al. 2009, Thiltgen & Goldstein



2012), indicating a substantial need for improved  $G$  prediction methods with higher accuracy. Therefore, we also asked to what extent our results depended on the method by which we calculated  $G$  values. We calculated a second set of  $G$  values, using the ddg monomer application in Rosetta (Kellogg et al. 2011). Because this application runs approximately 500 times slower than FoldX, we could not run it on all proteins in our data set. Instead, we arbitrarily selected five proteins (PDB IDs 1bp2, 1lba, 1ljl, 1pyl, and 2acy) as a test case. We found that FoldX performs similarly or better than ddg monomer (Figure 4). Thus, in our application here, we could not identify any major differences between predictions obtained from FoldX and those obtained from Rosetta ddg monomer.

## 5. Conclusion

We have developed a biophysical model linking stability changes  $G$  due to mutations at individual sites in proteins to site-specific evolutionary rates. This stability model predicts site-specific rates in very good agreement with empirical rates. Indeed, the overall correlation between empirical rates and  $G$ -based predictions is similar to the correlation with the best structural determinant, the packing density measure WCN, which, according to a recent mechanistic stress model, is a measure of the local stress introduced by mutations into the active protein structure  $G^*$  (Yeh et al. 2014b, Huang et al. 2014). However, despite the similar performance, large partial correlations show that the two factors  $G$  and  $G^*$  result in largely independent predictions. Moreover, there are proteins for which the stability model performs significantly better than the stress model, while for other proteins the reverse is true. Consistently, a two-variable model that combines stability and stress significantly improves predictions. Therefore, both the overall stability  $G$  and the stress  $G^*$  seem to capture distinct thermodynamic constraints on protein evolution.

The stability model presented here is a neutral model in which mutations are either neutral or completely deleterious according to whether the mutant's stability is above a certain threshold (Taverna & Goldstein 2002, Bloom et al. 2005, Bloom & Glassman 2009). A presumably more sophisticated model is based on posing a continuous dependence between fitness and  $G$  (Tokuriki & Tawfik 2009, Chen & Shakhnovich 2009, Goldstein 2011, Wylie & Shakhnovich 2011). However, even though the continuous fitness models appear to be more realistic than the neutral stability-threshold models, in a recent study Arenas et al. (2013) found that the neutral model leads to better predictions of site-specific amino-acid distributions. This finding provides additional support for our choice of using a neutral  $G$ -based model. In future work, it will be worthwhile to explore the site-dependency of substitution rates using continuous fitness-stability models.

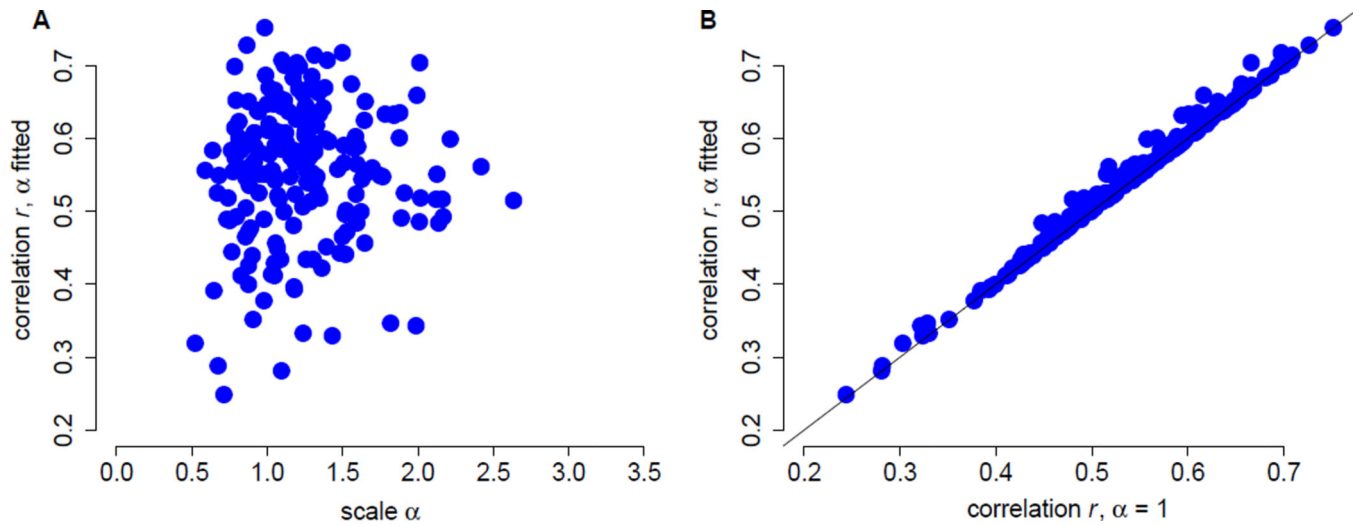
## Acknowledgements

We would like to thank Stephanie Spielman for help with setting up the evolutionary-rate calculations. J.E. is a researcher of CONICET. E.L.J. is funded by an NSF Graduate Research Fellowship, grant number DGE-1110007. C.O.W. is supported by NIH grant R01 GM088344, DTRA grant HDTRA1-12-C-0007, ARO grant W911NF-12-1-0390, and the BEACON Center for the Study of Evolution in Action (NSF Cooperative Agreement DBI-0939454). The Texas Advanced Computing Center (TACC) provided high-performance computing resources.

## References

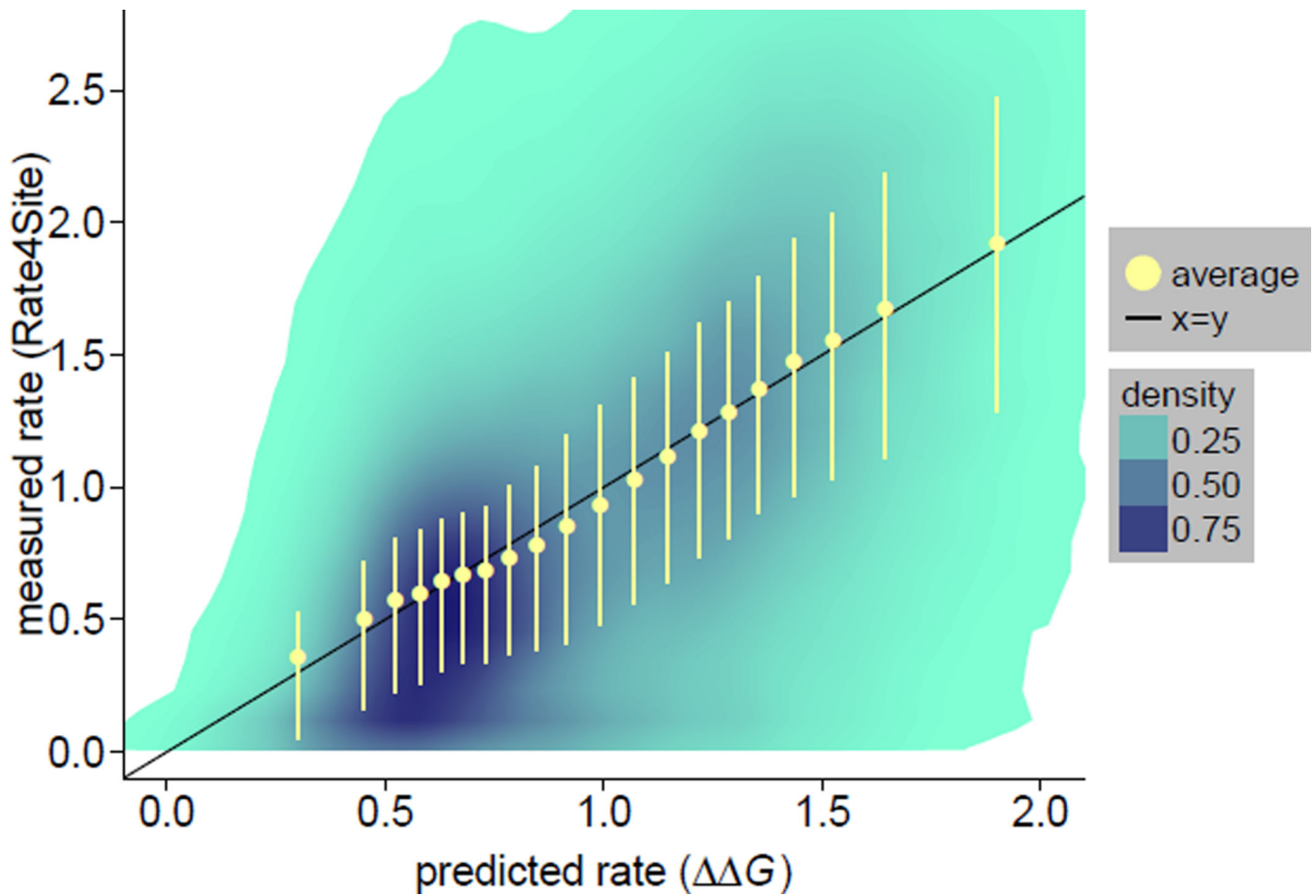
- Arenas M, Dos Santos HG, Posada D, Bastolla U. *Bioinformatics*. 2013; 29:3020–3028. [PubMed: 24037213]
- Ashenberg O, Gong LI, Bloom JD. *Proc. Natl. Acad. Sci. USA*. 2013; 110:21071–21076. [PubMed: 24324165]
- Bershtein S, Segal M, Bekerman R, Tokuriki N, Tawfik DS. *Nature*. 2006; 444:929–932. [PubMed: 17122770]
- Bloom JD, Glassman MJ. *PLOS Comp. Biol.* 2009; 5:e1000349.
- Bloom JD, Labthavikul ST, Otey CR, Arnold FH. *Proc. Natl. Acad. Sci. USA*. 2006; 103:5869–5874. [PubMed: 16581913]
- Bloom JD, Raval A, Wilke CO. *Genetics*. 2007; 175:255–266. [PubMed: 17110496]
- Bloom JD, Silberg JJ, Wilke CO, Drummond DA, Adami C, Arnold FH. *Proc. Natl. Acad. Sci. USA*. 2005; 102:606–611. [PubMed: 15644440]
- Bustamante CD, Townsend JP, Hartl DL. *Mol. Biol. Evol.* 2000; 17:301–308. [PubMed: 10677853]
- Chen P, Shakhnovich EI. *Genetics*. 2009; 183:639–650. [PubMed: 19620390]
- Dean AM, Neuhauser C, Grenier E, Golding GB. *Mol. Biol. Evol.* 2002; 19:1846–1864. [PubMed: 12411594]
- Drummond DA, Wilke CO. *Cell*. 2008; 134:341–335. [PubMed: 18662548]
- Franzosa EA, Xia Y. *Mol. Biol. Evol.* 2009; 26:2387–2395. [PubMed: 19597162]
- Goldstein RA. *Proteins*. 2011; 79:1396–1407. [PubMed: 21337623]
- Gong LI, Suchard MA, Bloom JD. *eLife*. 2013; 2:e00631. [PubMed: 23682315]
- Grahnen JA, Nandakumar P, Kubelka J, Liberles DA. *BMC Evol Biol*. 2011; 11:361. [PubMed: 22171550]
- Guerois R, Nielsen JE, Serrano L. *J. Mol. Biol.* 2002; 320:369–387. [PubMed: 12079393]
- Huang TT, Marcos ML, Hwang JK, Echave J. *BMC Evol. Biol.* 2014; 14:78. [PubMed: 24716445]
- Katoh K, Kuma KI, Toh H, Miyata T. *Nucl. Acids Res.* 2005; 33:511–518. [PubMed: 15661851]
- Katoh K, Standley DM. *Mol. Biol. Evol.* 2013; 30:772–780. [PubMed: 23329690]
- Kellogg EH, Leaver-Fay A, Baker D. *Proteins: Structure, Function, Bioinformatics*. 2011; 79:830–838.
- Lemos B, Bettencourt BR, Meiklejohn CD, Hartl DL. *Mol. Biol. Evol.* 2005; 22:1345–1354. [PubMed: 15746013]
- Liao BY, Weng MP, Zhang J. *Genome Biol. Evol.* 2010; 2:39–43. [PubMed: 20333223]
- Liao H, Yeh W, Chiang D, Jernigan RL, Lustig B. *Protein Eng. Des. Sel.* 2005; 59–64. [PubMed: 15788422]
- Liberles DA, Teichmann SA, Bahar I, Bastolla U, Bloom J, BornbergBauer E, Colwell LJ, de Koning APJ, Dokholyan NV, Echave J, Elofsson A, Gerloff DL, Goldstein RA, Grahnen JA, Holder MT, Lakner C, Lartillot N, Lovell SC, Naylor G, Perica T, Pollock DD, Pupko T, Regan L, Roger A, Rubinstein N, Shakhnovich E, Sjölander K, Sunyaev S, Teufel AI, Thorne JL, Thornton JW, Weinreich DM, Whelan S. *Protein Sci.* 2012; 21:769–785. [PubMed: 22528593]
- Lin CP, Huang SW, Lai YL, Yen SC, Shih CH, Lu CH, Huang CC, Hwang JK. *Proteins*. 2008; 72:929–935. [PubMed: 18300253]
- Mayrose I, Graur D, Ben-Tal N, Pupko T. *Mol. Biol. Evol.* 2004; 21:1781–1791. [PubMed: 15201400]
- Pal C, Papp B, Lercher MJ. *Nat Rev Genet.* 2006; 7:337–348. [PubMed: 16619049]
- Pang K, Cheng C, Xuan Z, Sheng H, Ma X. *BMC Systems Biol.* 2010; 4:179.
- Potapov V, Cohen M, Schreiber G. *Protein Eng. Des. Sel.* 2009; 22:553–560. [PubMed: 19561092]
- R Core Team. *R: A Language and Environment for Statistical Computing* R Foundation for Statistical Computing. Vienna, Austria: 2014. URL: <http://www.R-project.org>
- Ramsey DC, Scherrer MP, Zhou T, Wilke CO. *Genetics*. 2011; 188:479–488. [PubMed: 21467571]
- Raval A. *Phys. Rev. Lett.* 2007; 99:138104. [PubMed: 17930643]
- Schymkowitz JW, Rousseau F, Martins IC, Ferkinghoff-Borg J, Stricher F, Serrano L. *Proc. Natl. Acad. Sci. USA*. 2005; 102:10147–10152. [PubMed: 16006526]

- Serohijos A WR, Rimas Z, Shakhnovich EI. *Cell Reports*. 2012; 2:1–8. [PubMed: 22840390]
- Serrano L, Day AG, Fersht AR. *J. Mol. Biol.* 1993; 233:305–312. [PubMed: 8377205]
- Shahmoradi A, Sydykova DK, Spielman SJ, Jackson EL, Dawson ET, Meyer AG, Wilke CO. *J. Mol. Evol.* 2014 in press.
- Stamatakis A. *Bioinformatics*. 2014; 30:1312–1313. [PubMed: 24451623]
- Taverna DM, Goldstein RA. *Proteins*. 2002; 46:105–109. [PubMed: 11746707]
- Thiltgen G, Goldstein RA. *PLoS ONE*. 2012; 7:e46084. [PubMed: 23144695]
- Thorne JL. *Curr. Opin. Struct. Biol.* 2007; 17:337–341. [PubMed: 17572082]
- Tokuriki N, Stricher F, Serrano L, Tawfik DS. *PLoS Comp. Biol.* 2008; 4:e1000002.
- Tokuriki N, Tawfik DS. *Curr. Opin. Struct. Biol.* 2009; 19:596–604. [PubMed: 19765975]
- Wells JA. *Biochemistry*. 1990; 29:8509–8517. [PubMed: 2271534]
- Wilke CO, Drummond DA. *Current Opinion in Structural Biology*. 2010; 20:385–389. [PubMed: 20395125]
- Worth CL, Gong S, Blundell TL. *Nat. Rev. Mol. Cell. Biol.* 2009; 10:709–720. [PubMed: 19756040]
- Wylie CS, Shakhnovich EI. *Proc. Natl. Acad. Sci. USA*. 2011; 108:9916–9921. [PubMed: 21610162]
- Xia Y, Franzosa EA, Gerstein MB. *PLoS Comput. Biol.* 2009; 5
- Yang L, Song G, Jernigan RL. *Proc Natl Acad Sci USA*. 2009; 106:12347–12352. [PubMed: 19617554]
- Yeh SW, Huang TT, Liu JW, Yu SH, Shih CH, Hwang JK, Echave J. *Biomed Res. Int.* 2014b; 2014:572409. [PubMed: 25121105]
- Yeh SW, Liu JW, Yu SH, Shih CH, Hwang JK, Echave J, et al. *Mol. Biol. Evol.* 2014a; 31:135–139. [PubMed: 24109601]
- Zhang XJ, Baase WA, Shoichet BK, Wilson KP, Matthews BW. *Protein Eng.* 1995; 8:1017–1022. [PubMed: 8771182]



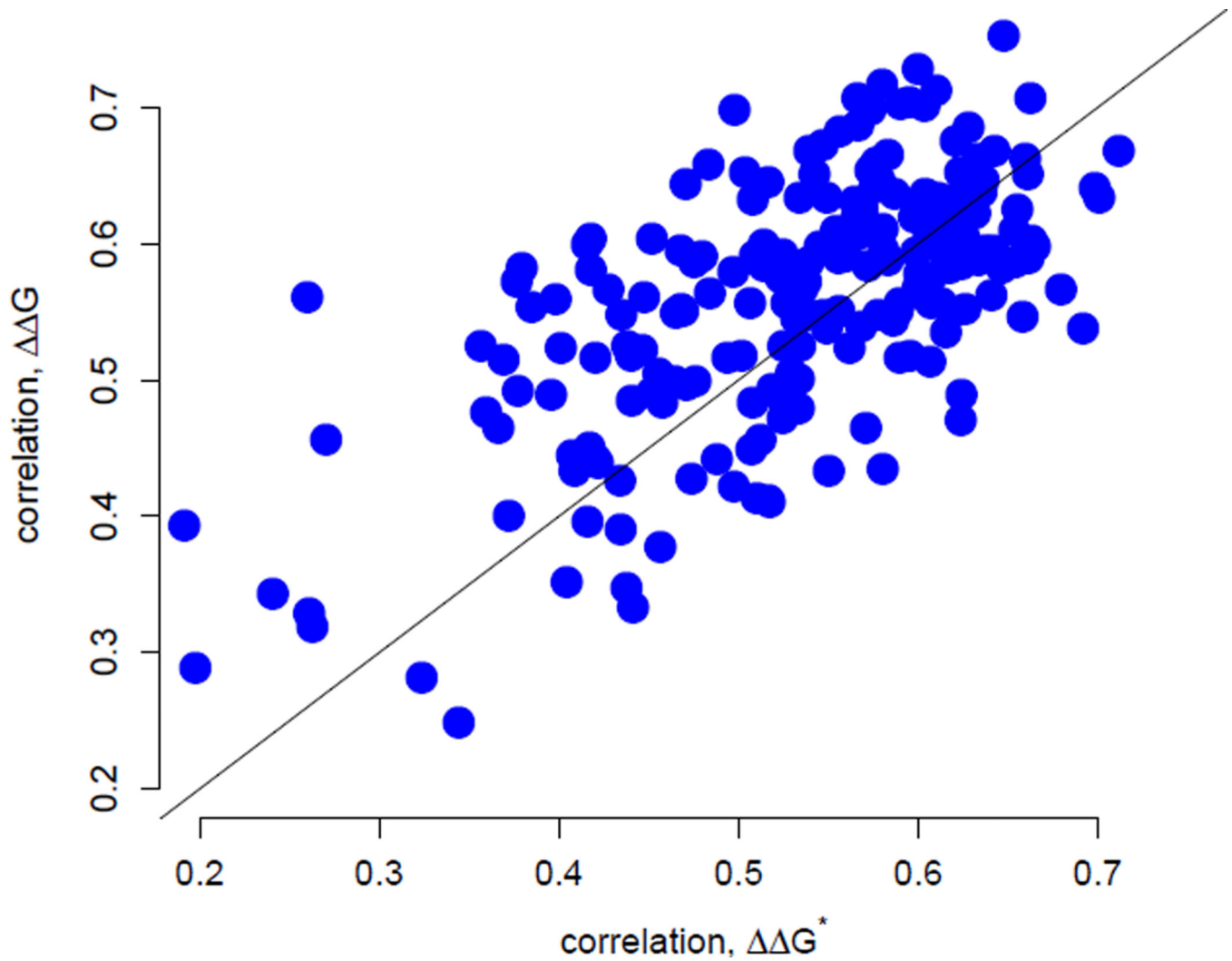
**Figure 1.**

Correlations between rates predicted from  $G$  and rates inferred by Rate4Site. (A) Correlation coefficients vs. the fitted, protein-specific scales  $\alpha$ . Each dot represents data for one protein. There is no relationship between the correlation coefficients and  $\alpha$  ( $r = 0.10$ ,  $P = 0.16$ ). (B) Fitted  $\alpha$  values provide only a small benefit over  $\alpha = 1$ . Fitting  $\alpha$  to each protein increases correlation coefficients, on average, by 0.007 (paired  $t$ -test, mean difference  $\bar{d} = 0.007$ ,  $df = 208$ ,  $P < 10^{-10}$ ).



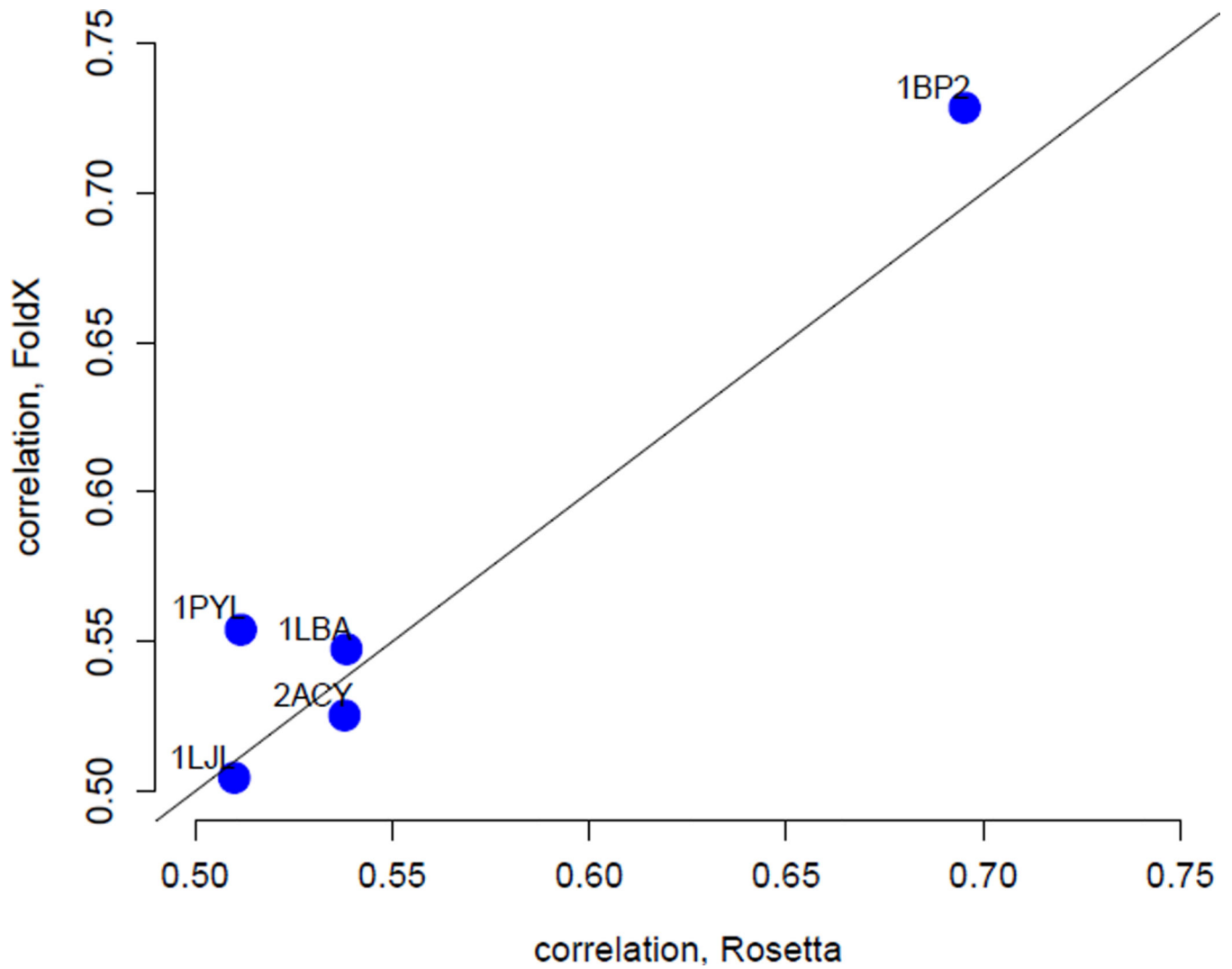
**Figure 2.**

The relationship between rates predicted from  $\Delta\Delta G$  and rates inferred by Rate4Site is nearly linear. The joint distribution of empirical vs. predicted rates is shown using shaded areas. All sites were grouped into 20 bins of approximately equal number of sites using quantile breaks on the predicted rate axis. Yellow dots are the mean rates obtained by averaging over sites within a bin. Yellow error bars correspond to the 25% and 75% quantiles for each bin. Average empirical rates (yellow circles) are very close to the  $x = y$  line that corresponds to a perfect empirical-predicted fit (the correlation coefficient between mean empirical and predicted rates is 0.995). However, there is substantial variation around the mean trend, as can be seen from shaded areas and yellow error bars (correlation between non-averaged empirical and predicted rates is 0.558).



**Figure 3.**

Correlations between rates inferred by Rate4Site and rates predicted by either the stress  $G^*$ -based model (shown along the  $x$  axis) or the stability  $G$ -based model (shown along the  $y$  axis). The correlation coefficients from the two models are significantly correlated ( $r = 0.64$ ,  $P < 10^{-10}$ ). Correlations have similar magnitudes, with the  $G$ -based model giving slightly better results on average (paired t-test, mean difference  $d = 0.026$ ,  $df = 208$ ,  $P < 0.001$ ). For 127 of the 209 proteins the stability model gives better correlations while for 82 of the 207 proteins the stress model gives better results.



**Figure 4.**

Rates predicted using  $G$  values obtained from FoldX perform as well as or better than the ones obtained from the ddg monomer protocol in Rosetta. Shown are the correlation coefficients of measured rates with rates predicted using the stability model with FoldX

$G$  values (y axis) vs. Rosetta  $G$  values (x axis) for five proteins.