

# mGPFusion: Predicting protein stability changes with Gaussian process kernel learning and data fusion

Emmi Jokinen<sup>1</sup>, Markus Heinonen<sup>1,2</sup> and Harri Lähdesmäki<sup>1</sup>

<sup>1</sup>Department of Computer Science, Aalto University, 02150 Espoo, Finland

<sup>2</sup>Helsinki Institute for Information Technology, Finland

## Abstract

**Motivation:** Proteins are commonly used by biochemical industry for numerous processes. Refining these proteins' properties via mutations causes stability effects as well. Accurate computational method to predict how mutations affect protein stability are necessary to facilitate efficient protein design. However, accuracy of predictive models is ultimately constrained by the limited availability of experimental data.

**Results:** We have developed mGPFusion, a novel Gaussian process (GP) method for predicting protein's stability changes upon single and multiple mutations. This method complements the limited experimental data with large amounts of molecular simulation data. We introduce a Bayesian data fusion model that re-calibrates the experimental and *in silico* data sources and then learns a predictive GP model from the combined data. Our protein-specific model requires experimental data only regarding the protein of interest and performs well even with few experimental measurements. The mGPFusion models proteins by contact maps and infers the stability effects caused by mutations with a mixture of graph kernels. Our results show that mGPFusion outperforms state-of-the-art methods in predicting protein stability on a dataset of 15 different proteins and that incorporating molecular simulation data improves the model learning and prediction accuracy.

**Availability:** Software implementation and datasets are available at [github.com/emmijokinen/mgpfusion](https://github.com/emmijokinen/mgpfusion)

**Contact:** [emmi.jokinen@aalto.fi](mailto:emmi.jokinen@aalto.fi)

## 1 Introduction

Proteins are used in various applications by pharmaceutical, food, fuel, and many other industries and their usage is growing steadily (Kirk *et al.*, 2002; Sanchez and Demain, 2010). Proteins have important advantages over chemical catalysts, as they are derived from renewable resources, are biodegradable and are often highly selective (Cherry and Fidantsef, 2003). Protein engineering is used to further improve the properties of proteins, for example to enhance their catalytic activity, modify their substrate specificity or to improve their thermostability (Rapley and Walker, 2000). Increasing the stability is an important aspect of protein engineering, as the proteins used in industry should be stable in the industrial process conditions, which often involve higher than ambient temperature and non-aqueous solvents (Bommarius *et al.*, 2011). The properties of a protein are modified by introducing alterations to its amino acid sequence. Mutations in general tend to be destabilising, and if too many destabilising mutations are implemented, the protein may not remain functional without compensatory stabilising mutations (Tokuriki and Tawfik, 2009).

The stability of a protein can be defined as the difference in Gibbs energy  $\Delta G$  between the folded and unfolded (or native and denatured) state of the protein. More precisely, the Gibbs energy difference determines the thermodynamic stability  $\Delta G_t$  of the protein, as it does not take into account the kinetic stability  $\Delta G_k$  which determines the energy needed for the transition between the folded and unfolded states (Anslyn and Dougherty, 2006) (see Supplementary Figure S1). Here we will consider only the thermodynamic stability and from now on it will be referred to merely as stability  $\Delta G$ .

The effect of mutations can be defined by the change they cause to the Gibbs energy  $\Delta G$ , denoted as  $\Delta\Delta G$  (Pace and Scholtz, 1997). To comprehend the significance of stability changes upon mutations, we can consider globular proteins, the most common type of enzymes, whose polypeptide chain is folded up in a compact ball-like shape with an irregular surface (Alberts *et al.*, 2007). These proteins are only marginally stable and the difference in Gibbs energy between the folded and unfolded state is only about 5–15 kcal/mol, which is not much more than the energy of a single hydrogen bond that is about 2–5 kcal/mol (Branden and Tooze, 1999). Therefore, even one mutation that breaks a hydrogen bond can prevent a protein from folding properly.

The protein stability can be measured with many techniques, including thermal, urea and guanidinium chloride (GdmCl) denaturation

curves that are determined as the fraction of unfolded proteins at different temperatures or at different concentrations of urea or GdmCl (Pace and Shaw, 2000). Some of the experimentally measured stability changes upon mutations have been gathered in thermodynamic databases such as Protherm (Kumar *et al.*, 2006).

A variety of computational methods have been introduced to predict the stability changes upon mutations more effortlessly than through experimental measurements. These methods utilise physics or knowledge-based potentials (Leaver-Fay *et al.*, 2011), their combinations, or different machine learning methods. The machine learning methods utilise support vector machines (SVM) (Capriotti *et al.*, 2005b, 2008; Chen *et al.*, 2013; Cheng *et al.*, 2006; Folkman *et al.*, 2014; Liu and Kang, 2012; Pires *et al.*, 2014a), random forests (Tian *et al.*, 2010; Wainreb *et al.*, 2011), neural networks (Dehouck *et al.*, 2009; Giollo *et al.*, 2014), and Gaussian processes (Pires *et al.*, 2014b). However, it has been assessed that although on average many of these methods provide good results, they tend to fail on details (Potapov *et al.*, 2009). In addition, many of these methods are able to predict the stability effects only for single-point mutations.

We introduce mGPFusion (mutation Gaussian Processes with data fusion), a method for predicting stability effects of both point and multiple mutations. mGPFusion is a protein-specific model – in contrast to earlier stability predictors that aim to estimate arbitrary protein structure or sequence stabilities – and achieves markedly higher accuracy while utilising data only from a single protein at a time. In contrast to earlier works that only use experimental data to train their models, we also combine exhaustive Rosetta (Leaver-Fay *et al.*, 2011) simulated point mutation *in silico* stabilities to our training data.

A key part of mGPFusion is the automatic scaling of simulated data to better match the experimental data distribution based on those variants that have both experimental and simulated stability values. Furthermore, we estimate a variance resulting from the scaling, which places a higher uncertainty on very destabilising simulations. Our Gaussian process model then utilises the joint dataset with their estimated heteroscedastic variances and uses a mixture of graph kernels to assess the stability effects caused by changes in amino acid sequence according to 21 substitution models. Our experiments on a novel 15 protein dataset show a state-of-the-art stability prediction performance, which is also sustained when there is access only to a very few experimental stability measurements.

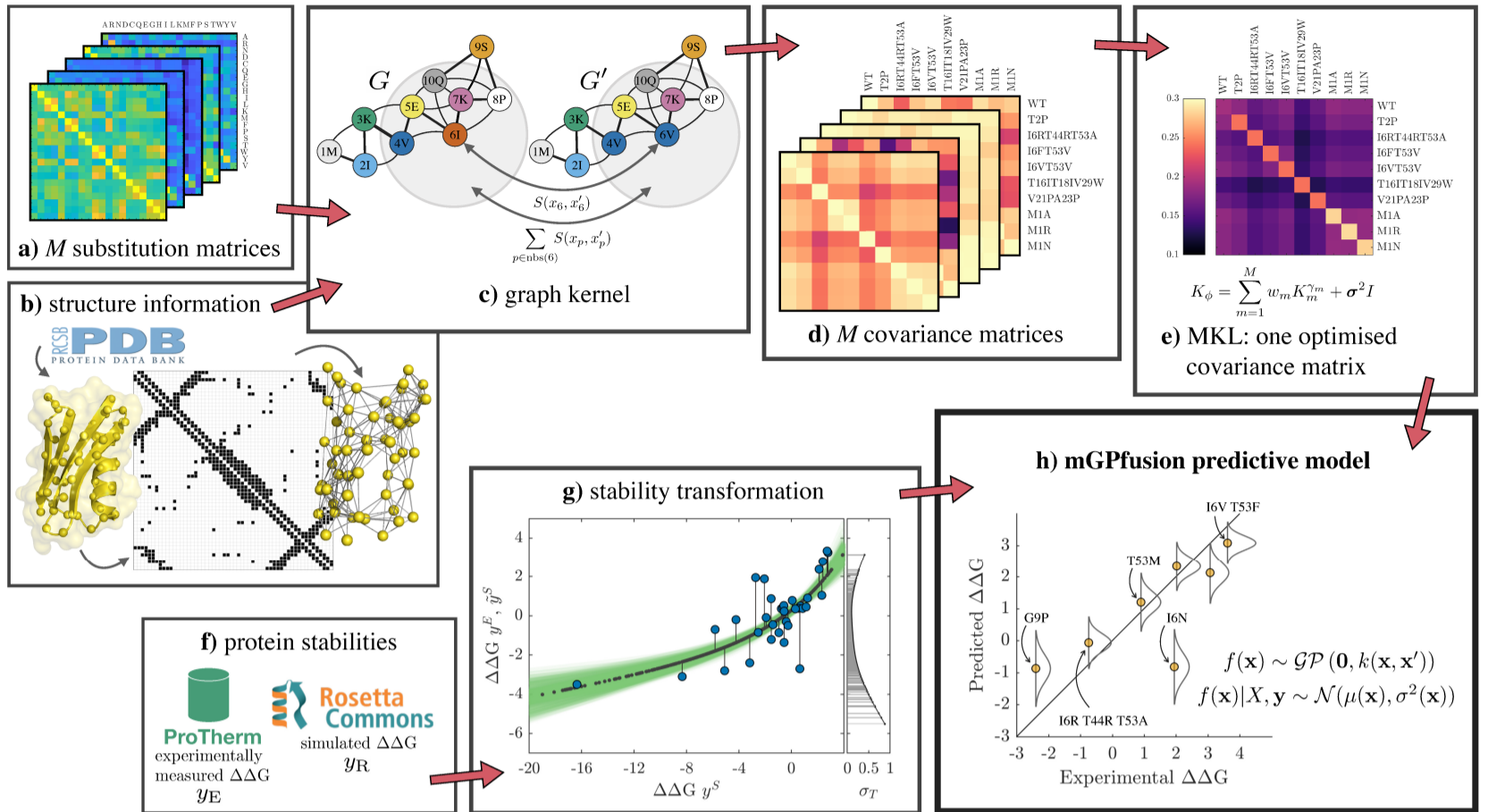


Figure 1: Pipeline illustration for mGPFusion. **a)**  $M = 21$  substitution matrices utilise different information sources and give scores to pairwise amino acid substitutions. **b)** The wild-type structures from Protein Data Bank are modelled as contact graphs. **c)** The graph kernel measures similarity of two sequences by a substitution model  $S$  over all positions  $p$  and their neighbourhoods  $\text{nbs}(p)$  in the contact graph. **d)** Each substitution matrix is used to create a separate covariance matrix. **e)** Multiple kernel learning (MKL) is used for finding the optimal combination of the base kernels. The kernel matrix measures variant similarities. **f)** Experimentally measured  $\Delta\Delta G$  values  $y^E$  are gathered from ProTherm and Rosetta’s ddg monomer application is used to simulate the stability effects  $y^S$  for all single point mutations. **g)** Bayesian scaling for the simulated values  $y^S$  at the x-axis. Possible scalings are coloured with green and the chosen scaling from  $y^S$  into scaled values  $\tilde{y}^S$  is marked by black dots. The scaling is fitted to a subset of experimentally measured stabilities  $y^E$  (circles). **h)** The stability predictive GP model is trained using experimental and simulated data through the kernel matrix.

## 2 Methods

Following Pires *et al.* (2014b) we choose a Bayesian model family of Gaussian processes for prediction of mutation effects on protein stability due to its inherent ability to handle uncertainty in a principled way. Bayesian modelling is a natural approach for combining the experimental and simulated data distribution, while it is also suitable for learning the underlying mixture of substitution models that governs the mutational process.

The pipeline for mGPFusion is presented in Figure 1. The first part of mGPFusion consists of collection of *in silico* and experimental datasets discussed in Section 2.1, the scaling of the *in silico* dataset in Section 2.2 and the fusion of these two datasets in Section 2.3. The second part consists of the Gaussian process model described in Section 2.4 with detailed description of the graph kernels in Sections 2.5-2.6 and model inference in Section 2.7. Finally, the evaluation criteria used are described in Section 2.8.

### 2.1 Experimental and *in silico* data

ProTherm is a database of numerical thermodynamic parameters for proteins and their mutants (Kumar *et al.*, 2006). From ProTherm we gathered all proteins with at least 50 unique mutations whose  $\Delta\Delta G$  has been measured by thermal denaturation, and where a PDB code for a 3D structure of the protein was reported. We required the proteins to have at least 50 unique mutations, so that we would have a representative test set and get sufficiently reliable estimates of prediction accuracy on individual proteins and examine how the amount of experimental training data affects the accuracy of the model. The 3D structures are necessary for obtaining the connections between residues. We collected the 3D structures with the reported PDB codes from the Protein Databank, [www.rcsb.org](http://www.rcsb.org) (Berman *et al.*, 2000). The 15 pro-

teins that fulfilled these requirements are listed in Table 1. We averaged the stability values for proteins with multiple measurements and ignored mutations to residues not present in their 3D structures. These data sets are available at [github.com/emmijokinen/mgpfusion](https://github.com/emmijokinen/mgpfusion).

| Protein (organism)  | PDB  | mutations |       |             |
|---|------|-----------|-------|-------------|
|   |      | all       | point | point (sim) |
| T4 Lysozyme ( <i>Enterobacteria phage T4</i> )              | 2LZM | 349       | 264   | 3116        |
| Barnase ( <i>Bacillus amyloliquefaciens</i> )               | 1BNI | 182       | 163   | 2052        |
| Gene V protein ( <i>Escherichia virus m13</i> )             | 1VQB | 124       | 92    | 1634        |
| Glycosyltransferase A ( <i>Homo sapiens</i> )               | 1LZI | 116       | 114   | 2470        |
| Chymotrypsin inhibitor 2 ( <i>Hordeum vulgare</i> )         | 2C12 | 98        | 77    | 1235        |
| Protein G ( <i>Streptococcus sp. gx7805</i> )               | 1PGA | 89        | 34    | 1064        |
| Ribonuclease H ( <i>Escheria coli</i> )                     | 2RN2 | 83        | 65    | 2945        |
| Cold shock protein B ( <i>Bacillus subtilis</i> )           | 1CSP | 80        | 50    | 1273        |
| Apomyoglobin ( <i>Physeter catodon</i> )                    | 1BVC | 80        | 56    | 2907        |
| Hen egg white lysozyme ( <i>Gallus gallus</i> )             | 4LYZ | 63        | 50    | 2451        |
| Ribonuclease A ( <i>Bos taurus</i> )                        | 1RTB | 57        | 50    | 2356        |
| Peptidyl-prolyl cis-trans isomerase ( <i>Homo sapiens</i> ) | 1PIN | 56        | 56    | 2907        |
| Ribonuclease T1 isozyme ( <i>Aspergillus oryzae</i> )       | 1RN1 | 53        | 48    | 1957        |
| Ribonuclease ( <i>Streptomyces auerofaciens</i> )           | 1RCG | 54        | 45    | 1824        |
| Bovine pancreatic trypsin inhibitor ( <i>Bos taurus</i> )   | 1BPI | 53        | 47    | 1102        |
| total   |      | 1537      | 1211  | 31293       |

Table 1: The 15 protein data from ProTherm database with counts of point mutations, all mutations, and of simulated point mutation stability changes.

We also generate simulated data of the stability effects of all possible single mutations of the proteins. Our method can utilise any simulated stability values. We used the “ddG monomer” application of Rosetta 3.6 (Leaver-Fay *et al.*, 2011) using the high-resolution backrub-based protocol 16 recommended in Kellogg *et al.* (2011). The predictions  $y^S$  made with Rosetta are given in Rosetta Energy Units (REU). Kellogg *et al.* (2011) suggest transformation  $0.57y^S$  for converting the predictions into physical units. The simulated data scaled this way is not as accurate as the experimental data, the correlation and root mean

square error (rmse) with respect to the experimental data are shown for all proteins in Table 2 and for individual proteins in Supplementary Table S2, on rows labelled Rosetta. For this reason, we use instead a Bayesian scaling described in the next section and different noise models for the experimental and simulated data, described in Section 2.3.

For each of the 15 proteins, let  $\mathbf{x}_i = (x_{i1}, \dots, x_{iM})$  denote its  $M$ -length variant  $i$  with positions  $p$  labelled with residues  $x_{ip} \in \{A, R, N, \dots, V\}$ . We denote the wild-type protein as  $\mathbf{x}_0$ . We collect 15 separate sets of simulated and experimental data. We denote the  $N_E$  experimental variants of each protein as  $X_E = (\mathbf{x}_1^E, \dots, \mathbf{x}_{N_E}^E)^T$  with the corresponding experimental stability values  $\mathbf{y}_E = (y_1^E, \dots, y_{N_E}^E)^T \in \mathbb{R}^{N_E}$ . Similarly, we denote the  $N_S$  simulated observations as  $X_S = (\mathbf{x}_1^S, \dots, \mathbf{x}_{N_S}^S)^T$  and  $\mathbf{y}_S = (y_1^S, \dots, y_{N_S}^S)^T \in \mathbb{R}^{N_S}$ .

## 2.2 Bayesian scaling of *in silico* data

The described transformation from REU to physical units may not be optimal for all proteins. We therefore applied instead a linear-exponential scaling function to obtain *scaled* Rosetta simulated stabilities  $\tilde{y}^S$ ,

$$\tilde{y}^S = g(y^S | \theta_j) = a_j e^{c_j y^S} + b_j y^S + d_j. \quad (1)$$

This scaling transforms the Rosetta simulations  $y^S$  for each protein  $j = 1, \dots, 15$  to correspond better to the experimental data. The parameters  $\theta_j = (a_j, b_j, c_j, d_j)$  define the weight  $a_j$  and steepness  $c_j$  of the exponential term, while the linear term has slope  $b_j$  and intercept  $d_j$ . To avoid overfitting, we perform Bayesian linear regression and start by defining parameter prior  $p(\theta_j) = p(a_j)p(b_j)p(c_j)p(d_j)$  that reflects our beliefs about realistic scalings having only moderate steepness:

$$\begin{aligned} p(a_j) &= \text{Gamma}(a_j | \alpha_a, \beta_a) \\ p(b_j) &= \text{Beta}(1/2 \cdot b_j | \alpha_b, \beta_b) \\ p(c_j) &= \text{Beta}(10/3 \cdot c_j | \alpha_c, \beta_c) \\ p(d_j) &= \mathcal{N}(d_j | \mu_d, \sigma_d^2). \end{aligned} \quad (2)$$

The empirically selected hyperparameter values are listed in Supplementary Table S1 and the priors are illustrated in Figure S2.

We compute the posterior for  $\theta_j$  using the subset of simulated data that have corresponding experimentally measured data:

$$p(\theta_j | \mathbf{y}_E, \mathbf{y}_S) \propto \prod_{i: \mathbf{x}_i \in X_E \cap X_S} \mathcal{N}(y_i^E | g(y_i^S | \theta_j), \sigma_n^2) p(\theta_j).$$

The product iterates over all  $N_{E \cap S}$  simulated  $\Delta\Delta G$ 's that have a matching experimentally observed value. The  $\sigma_n^2$  is the scaling error variance, which was set to  $\sigma_n^2 = 0.5$ . The parameters  $\theta$  for each protein were sampled using a random walk Metropolis-Hastings MCMC algorithm (the `mhsample` function in Matlab) for  $N_{MC} = 10000$  samples with a burn-in set to 500. The proposal distribution was selected to be a symmetric uniform distribution such that  $[a_{s+1}, b_{s+1}, c_{s+1}, d_{s+1}] \sim U(a_s \pm 0.4, b_s \pm 0.04, c_s \pm 0.04, d_s \pm 0.4)$ . Given the sample of scaling parameters  $(\theta_j^{(s)})_{s=1}^{N_{MC}}$ , we define the scaled simulated data as the average scaling over the MCMC sample, and record also the sample scaling variance

$$\tilde{y}_i^S = \frac{1}{N_{MC}} \sum_{s=1}^{N_{MC}} g(y_i^S | \theta_j^{(s)}) \quad (3)$$

$$\sigma_T^2(i) = \frac{1}{N_{MC}} \sum_{s=1}^{N_{MC}} (g(y_i^S | \theta_j^{(s)}) - \tilde{y}_i^S)^2. \quad (4)$$

See Figure 1 g) for an illustration of the scaling. We collect the scaled simulated value and its variance from each simulated point into vectors  $\tilde{\mathbf{y}}_S = (\tilde{y}_1^S, \dots, \tilde{y}_{N_S}^S)$  and  $\boldsymbol{\sigma}_T^2 = (\sigma_T^2(1), \dots, \sigma_T^2(N_S)) \in \mathbb{R}^{N_S}$ .

## 2.3 Data fusion and noise models

For each protein  $j$ , we organise its experimental data  $(X_E, \mathbf{y}_E)$  and transformed simulated data  $(X_S, \tilde{\mathbf{y}}_S)$  along with the wild-type information  $(\mathbf{x}_0, y_0)$  into a single joint dataset of variants  $X = (\mathbf{x}_0, X_E, X_S)$

and stabilities  $\mathbf{y} = (y_0, \mathbf{y}_E, \tilde{\mathbf{y}}_S)$  of size  $\mathbb{R}^N$  where  $N = 1 + N_E + N_S$  is the total number of simulated and experimental data points, including the wild-type. We assume heteroscedastic additive noise models for the three information sources

$$\begin{aligned} y_0 &= f(\mathbf{x}_0) + \varepsilon_0, & \varepsilon_0 &\sim \mathcal{N}(0, \sigma_0^2) \\ y_i^E &= f(\mathbf{x}_i^E) + \varepsilon_i^E, & \varepsilon_i^E &\sim \mathcal{N}(0, \sigma_E^2) \\ \tilde{y}_i^S &= f(\mathbf{x}_i^S) + \varepsilon_i^S, & \varepsilon_i^S &\sim \mathcal{N}(0, (\sigma_E + \sigma_S + t\sigma_T(i))^2), \end{aligned} \quad (5)$$

where the observed values are noisy versions of the underlying ‘true’ stability function  $f(\mathbf{x})$  corrupted by zero-mean noise with data source specific variances. We learn a Gaussian process based stability function  $f(\mathbf{x})$  in the next Section.

The Equations (5) encode that the experimental data are corrupted by a global experimental noise variance  $\sigma_E^2$ . The simulated stabilities are additionally corrupted by a global Rosetta simulator error variance  $\sigma_S^2$ , and by the value-dependent transformation variance  $t\sigma_T^2(i)$  scaled by parameter  $t$ . The model then encapsulates that we trust the Rosetta data less than the experimental data. By definition, the  $\Delta\Delta G$  of the wild-type is zero ( $y_0 = 0$ ) with very small assumed error,  $\sigma_0 = 10^{-6}$ . Note that  $\sigma_T^2$  are fixed by equation (4), while we infer the optimal values for the remaining three free parameters  $(\sigma_E, \sigma_S, t)$  (See Section 2.4). The parameters  $\sigma_E^2$  and  $\sigma_S^2$  are assigned priors

$$\begin{aligned} \sigma_E &\sim \text{Gamma}(\sigma_E | \alpha_E, \beta_E) \\ \sigma_S &\sim \text{Gamma}(\sigma_S | \alpha_S, \beta_S). \end{aligned} \quad (6)$$

The values of these hyperparameters are shown in Supplementary Table S1.

## 2.4 Gaussian processes

We use a Gaussian process (GP) function  $f$  to predict the stability  $f(\mathbf{x}) \in \mathbb{R}$  of a protein variant  $\mathbf{x}$ . Gaussian processes are a family of non-parametric, non-linear Bayesian models (Rasmussen and Williams, 2006). A zero-mean GP prior

$$f(\mathbf{x}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}')),$$

defines a distribution over functions  $f(\mathbf{x})$  whose mean and covariance are

$$\begin{aligned} \mathbb{E}[f(\mathbf{x})] &= 0 \\ \text{cov}[f(\mathbf{x}), f(\mathbf{x}')] &= k(\mathbf{x}, \mathbf{x}'). \end{aligned}$$

For any collection of protein variants  $X = \mathbf{x}_1, \dots, \mathbf{x}_N$ , the function values follow a multivariate normal distribution  $\mathbf{f} \sim \mathcal{N}(\mathbf{0}, K_{XX})$ , where  $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))^T \in \mathbb{R}^N$ , and where  $K_{XX} \in \mathbb{R}^{N \times N}$  with  $[K_{XX}]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ . The key property of Gaussian processes is that they encode functions that predict similar stability values  $f(\mathbf{x}), f(\mathbf{x}')$  for protein variants  $\mathbf{x}, \mathbf{x}'$  that are similar, as encoded by the kernel  $k(\mathbf{x}, \mathbf{x}')$ . The key part of GP modelling is then to infer a kernel that measures the mutation’s effects to the stability.

Let a dataset of noisy stability values from two sources be  $\mathbf{y} \in \mathbb{R}^N$ , the corresponding protein structures  $X = (\mathbf{x}_i)_{i=1}^N$ , and a new protein variant  $\mathbf{x}_*$  whose stability we wish to predict. A Gaussian process defines a joint distribution over the observed values  $\mathbf{y}$  of variants  $X$ , and the unknown function value  $f(\mathbf{x}_*)$  of the unseen variant  $\mathbf{x}_*$ ,

$$\begin{bmatrix} \mathbf{y} \\ f(\mathbf{x}_*) \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} K_{XX} + \text{diag}(\boldsymbol{\sigma}^2) & \mathbf{k}_{X*} \\ \mathbf{k}_{*X} & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right),$$

where  $\mathbf{k}_{X*} = \mathbf{k}_{*X}^T \in \mathbb{R}^N$  is a kernel vector with elements  $k(\mathbf{x}_i, \mathbf{x}_*)$  for all  $i = 1, \dots, N$ , and where  $\boldsymbol{\sigma}^2 = (\sigma_0^2, \sigma_E^2 \mathbf{1}^T, (\sigma_E \mathbf{1}^T + \sigma_S \mathbf{1}^T + t\boldsymbol{\sigma}_T^T)^2)^T$  collects final variances of the data points from equations (5). Here the exponents are elementwise. The conditional distribution gives the posterior distribution of the stability prediction as

$$f(\mathbf{x}_*) | (X, \mathbf{y}) \sim \mathcal{N}(\mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*)),$$

where the prediction mean and variance are

$$\begin{aligned} \mu(\mathbf{x}_*) &= \mathbf{k}_{*X} (K_{XX} + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \mathbf{y}, \\ \sigma^2(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_{*X} (K_{XX} + \text{diag}(\boldsymbol{\sigma}^2))^{-1} \mathbf{k}_{X*}. \end{aligned}$$

Hence, in GP regression the stability predictions  $\mu(\mathbf{x}_*) \pm \sigma(\mathbf{x}_*)$  will come with uncertainty estimates.

## 2.5 Graph kernel

Next, we consider how to compute the similarity function  $k(\mathbf{x}, \mathbf{x}')$  between two variants of the same protein structure. We will encode the 3D structural information of the two protein variants as a contact map and measure their similarity by the formalism of graph kernels (Vishwanathan *et al.*, 2010).

We consider two residues to be in contact if their closest atoms are within 5 Å of each other in the PDB structure, which is illustrated in Figure 1 b). All variants of the same protein have the same length, with only different residues at mutating positions. Furthermore, we assume that all variants share the wild-type protein contact map.

To compare protein variants, we construct a weighted decomposition kernel (WDK) (Menchetti *et al.*, 2005) between two protein variants  $\mathbf{x} = (x_1, \dots, x_M)$  and  $\mathbf{x}' = (x'_1, \dots, x'_M)$  of length  $M$ ,

$$k(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^M \left( S(x_p, x'_p) \sum_{l \in \text{nbs}(p)} S(x_l, x'_l) \right), \quad (7)$$

where  $\text{nbs}(p)$  defines the set of neighbouring positions to position  $p$ , and  $S$  is a substitution matrix. The kernel iterates over all positions  $p$  and compares for each of them their residues through a substitution matrix  $S(x_p, x'_p)$ . Furthermore, the similarity of the residues at each position is multiplied by the average similarity of the residues at its neighbouring positions  $S(x_l, x'_l)$ . Hence, the kernel defines the similarity of two protein variants as the average position and neighbourhood similarity over all positions. The kernel matrix is normalised so that for two data points, the normalised kernel is  $\hat{k}(x_p, x'_p) = k(x_p, x'_p) / \sqrt{k(x_p, x_p)k(x'_p, x'_p)}$ , as defined by Shawe-Taylor and Cristianini (2004). The kernel is illustrated in Figure 1 c).

The above WDK kernel allows us to compare the effects of multiple simultaneous mutations. However, as the wild type protein structure is used for all of the protein variants, changes that the mutations may cause to the protein structure are not taken into consideration. This may cause problems if mutations that alter the protein structure significantly are introduced – especially if many of them are introduced simultaneously. On the other hand, substitution matrices that have their basis in sequence comparisons, should take these effects into account to some extent as these kinds of mutations are usually highly destabilising and do not occur often in nature. In the next section, we will discuss how we utilise different substitution matrices with multiple kernel learning.

## 2.6 Substitution matrices and multiple kernel learning

The BLOSUM substitution models have been a common choice for protein models (Giguere *et al.*, 2013), while mixtures of substitution models were proposed by Cichonska *et al.* (2017). BLOSUM matrices score amino acid substitutions by their appearances throughout evolution, as they compare the frequencies of different mutations in similar blocks of sequences (Henikoff and Henikoff, 1992). However, there are also different ways to score amino acids substitutions, such as chemical similarity and neighbourhood selectivity (Tomii and Kanehisa, 1996). When the stability effects of mutations are evaluated, the frequency of an amino acid substitution in nature may not be the most important factor.

To take into account different measures of similarity between amino acids, we employed a set of 21 amino acid substitution matrices gathered from AAindex2<sup>1</sup> (Kawashima *et al.*, 2008). AAindex2 currently contains 94 substitution matrices. From these we selected those that had no gaps concerning substitutions between the 20 naturally occur-

ring amino acids and scaled them between zero and one as

$$S = \frac{S_0 - \min(S_0) + 1}{\max(S_0) - \min(S_0) + 1}. \quad (8)$$

Out of these matrices, we only chose those 23 matrices that were positive semidefinite. Furthermore, there were two pairs of matrices that were extremely similar, and we only selected one matrix from each pair, ending up with 21 substitution matrices. These substitution matrices are used together with Equation 7 for computing 21 base kernel matrices. Finally, MKL is used to find an optimal combination of the base kernels of form

$$K_\phi = \sum_{m=1}^{21} w_m K_m^{(\gamma_m)}, \quad (9)$$

where  $w_m$  is a kernel specific weight,  $\gamma_m$  is an (elementwise) exponent. The elementwise exponent retains the SDP property of  $K_\phi$  (Shawe-Taylor and Cristianini, 2004). We observe empirically that the optimal kernel weights  $w_m$  tend to be sparse (See Figure 2).

The selected substitution matrices are listed in Figure 2. These matrices have different basis and through multiple kernel learning (MKL) our model learns which of these are important for inferring the stability effects that mutations cause on different proteins. The figure illustrates this by showing the average weights of the base kernel matrices obtained via the multiple kernel learning.

## 2.7 Parameter inference

The complete model has five parameters  $\phi = (\sigma_E, \sigma_S, t, \mathbf{w}, \boldsymbol{\gamma})$  to infer, of which the variance parameters  $(\sigma_E, \sigma_S, t)$  parameterise the joint data variance  $\sigma_\phi^2$ , while the MKL parameters  $\mathbf{w} = (w_1, \dots, w_{21})$  and  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_{21})$  parameterise the kernel matrix  $K_\phi$ . In a Gaussian process model these can be jointly optimised by the marginal (log) likelihood with priors

$$\begin{aligned} \log p(\mathbf{y}|\phi)p(\sigma_E)p(\sigma_S) &= \log \int p(\mathbf{y}|\mathbf{f}, \phi)p(\mathbf{f}|\phi)p(\sigma_E)p(\sigma_S)d\mathbf{f} \\ &\propto -\frac{1}{2}\mathbf{y}^T(K_\phi + \text{diag}(\sigma_\phi^2))^{-1}\mathbf{y} - \frac{1}{2}\log |K_\phi + \text{diag}(\sigma_\phi^2)| \\ &\quad + \log \text{Gamma}(\sigma_E|\alpha_E, \beta_E) + \log \text{Gamma}(\sigma_S|\alpha_S, \beta_S), \end{aligned} \quad (10)$$

which automatically balances model fit (the square term) and the model complexity (the determinant) to avoid overfitting (Rasmussen and Williams, 2006). The parameters can be optimised by maximising the marginal log likelihood (10) using gradient ascent, since the marginal likelihood can be differentiated analytically (see Supplementary Equations S1 and S2). We utilised a limited-memory projected quasi-Newton algorithm (minConf.TMP<sup>2</sup>), described by (Schmidt *et al.*, 2009).

## 2.8 Evaluation criteria

We chose to evaluate the accuracy of our predictions using the same metrics that have been used by many others – correlation  $\rho$  between the predicted and experimentally measured  $\Delta\Delta G$  values (Capriotti *et al.*, 2005a; Dehouck *et al.*, 2009; Kellogg *et al.*, 2011; Pires *et al.*, 2014b; Potapov *et al.*, 2009) and the root mean square error (rmse) (Dehouck *et al.*, 2009; Pires *et al.*, 2014a,b), which are determined in the Supplementary Equations S3 and S4. We use marginal likelihood maximisation to infer model parameters and perform cross-validation to evaluate the model performance on test data. Below we only report evaluation metrics obtained from the test sets not used at any stage of the model learning or data transformation sampling.

## 3 Results

In this section we evaluate the performance of mGPFusion on predicting stability effects of mutations, and compare it to the state-of-the-art prediction methods mCSM, PoPMuSiC and Rosetta. Rosetta is a

<sup>1</sup><http://www.genome.jp/aaindex/>

<sup>2</sup><http://www.cs.ubc.ca/~schmidtm/Software/minConf.html>

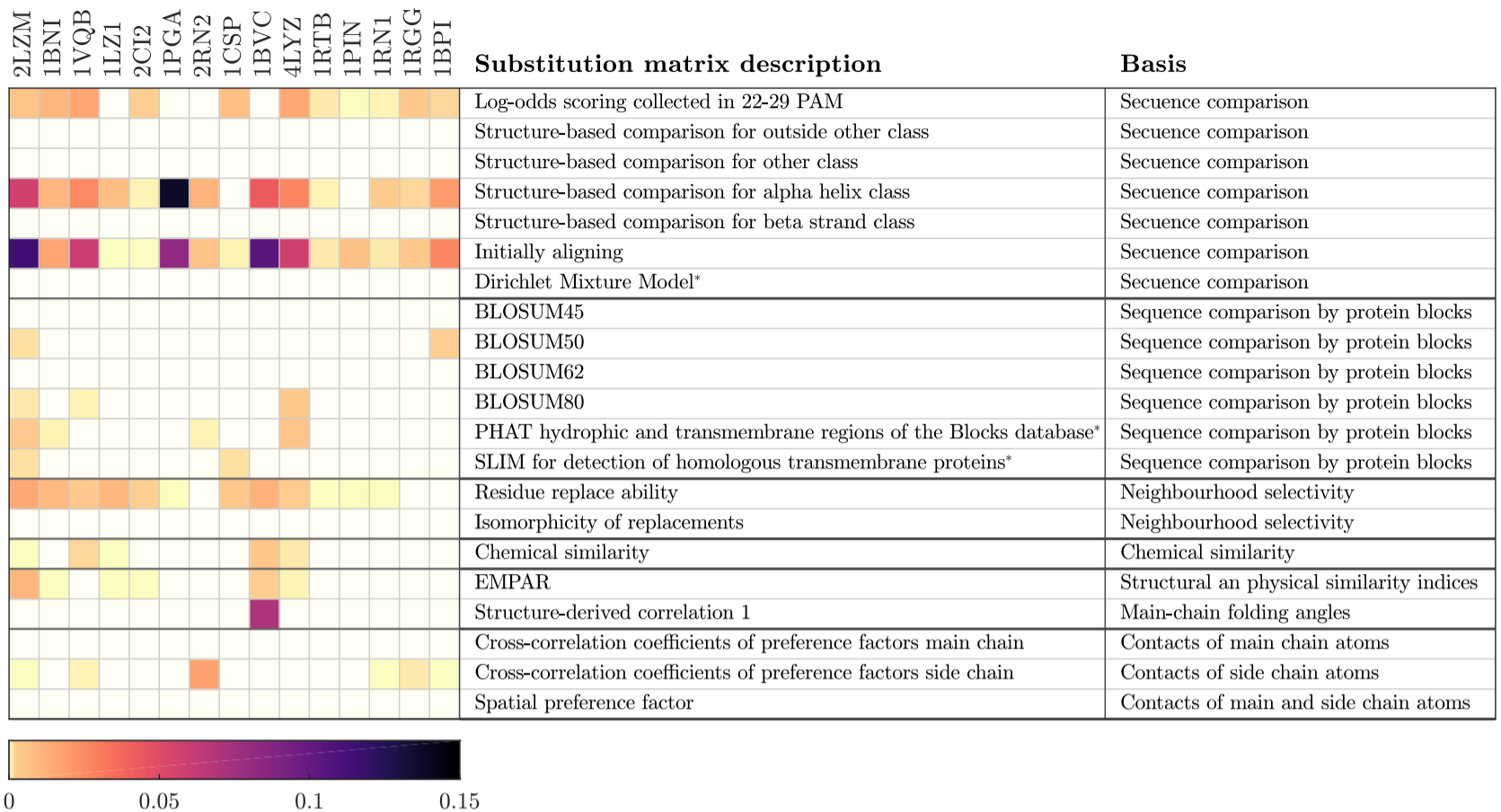


Figure 2: Average weights for kernels utilising the described substitution matrices from AAindex2, when GP models were trained with mutation level cross-validation. Basis for the substitution matrices are obtained from (Tomii and Kanehisa, 1996). \* were added to AAindex2 in a later release, and their basis were not determined by Tomii and Kanehisa (1996).

molecular modelling software whose `ddg_monomer` module can directly simulate the stability changes  $\Delta\Delta G$  of a protein upon mutations. PoP-MuSiC and mCSM are machine learning models that predict stability based on protein variant features. We run Rosetta locally, and use mCSM and PoPMuSiC models through their web servers<sup>3</sup>. This may give these methods an advantage over mGPFusion since parts of our testing data were likely used within their training data.

We compare four different variants of our method: mGPFusion that uses both simulated data and MKL, “mGPFusion, only B62” that uses simulated data but incorporates only one kernel matrix (BLOSUM62 substitution matrix), mGP model that uses MKL but does not use simulated data, and “mGP, only B62” that uses only the base GP model but does not incorporate simulated data and uses only the BLOSUM62 substitution matrix. In addition, we experiment on transforming Rosetta predictions with the Bayesian scaling. We perform the experiments for the 15 proteins separately using either position or mutation level (leave-one-out) cross-validation regarding the methods mGP, mGPFusion and the Bayesian scaling of Rosetta. Pires *et al.* (2014b) used protein and position level cross-validation to evaluate their model. In protein level cross-validation all mutations in a protein are either in the test or training set exclusively. When we train our model using protein level cross-validation, we use no experimental data and rely only on the simulated data. Position level cross-validation is defined so that all mutations in a position are either in the test or training set exclusively. However, datasets in Pires *et al.* (2014b) contained only point mutations and therefore we had to extend the definition to also include multiple mutations. In position level cross-validation we train one model for each position using only the part of data that has a wild-type residue in that position. Therefore, in position level cross-validation we construct a test set that contains all protein variants that have a mutation at position  $p$  and use as training set all the protein variants that have a wild-type residue at that position. Dehouck *et al.* (2009) evaluated their models by randomly selecting training and test sets so that each mutation was exclusively in one of the sets, but both sets could contain mutations from the same position of the same pro-

tein. We call this mutation level cross-validation. When we use all available experimental data with mutation level cross-validation, this corresponds to leave-one-out cross-validation.

### 3.1 Predicting point mutations

Table 2 summarises the average prediction performance over all 15 proteins for all compared methods, types of mutations and cross-validation types. We first compare the performances on single point mutations, where mGPFusion and mGP achieve the highest performance with  $\rho = 0.81$  and  $\text{rmse} = 1.07$  kcal/mol, and  $\rho = 0.81$  and  $\text{rmse} = 1.04$  kcal/mol, respectively with mutation level cross-validation. With only one kernel utilising the BLOSUM62 matrix instead of MKL, the performance decreases slightly, but the competing methods are still outperformed, as mCSM achieves  $\rho = 0.64$  and  $\text{rmse} = 1.37$  kcal/mol, PoPMuSiC  $\rho = 0.61$  and Rosetta  $\rho = 0.55$ . Applying Bayesian scaling on Rosetta simulator improves the performance of standard Rosetta from  $\rho = 0.55$  to  $\rho = 0.65$  and decreases the  $\text{rmse}$  from 1.63 kcal/mol to 1.35 kcal/mol, which is interestingly even slightly better than the performances of mCSM and PoPMuSiC.

With position level cross-validation mGPFusion achieves the highest performance of  $\rho = 0.70$  and  $\text{rmse} = 1.26$  kcal/mol, likely due to having still access to simulated variants from that position, since they are always available to the learner. Without simulation data, the baseline machine learning model mGP performance decreases to  $\rho = 0.51$  and  $\text{rmse} = 1.54$  kcal/mol, thus demonstrating the importance of the data fusion. Cross-validation could not be performed for the off-the-shelf methods mCSM and PoPMuSiC. Even still, mGPFusion (trained with one or multiple kernels) outperforms competing state-of-the-art methods and achieves markedly higher prediction performance as quantified by both mutation and position level cross-validations. Also mGP outperforms these methods when quantified by mutation level cross-validation. With protein level cross-validation mGPFusion achieves slightly better results than Rosetta.

<sup>3</sup>[biosig.unimelb.edu.au/mcsm](http://biosig.unimelb.edu.au/mcsm) and [omictools.com/popmusic-tool](http://omictools.com/popmusic-tool)



| Method  | Correlation $\rho$     |             |       |                        |             |       |                        |             |       | rmse                   |             |       |                        |             |       |                        |             |       |
|---|------------------------|-------------|-------|------------------------|-------------|-------|------------------------|-------------|-------|------------------------|-------------|-------|------------------------|-------------|-------|------------------------|-------------|-------|
|   | Point mutations        |             |       | Multiple mutations     |             |       | All mutations          |             |       | Point mutations        |             |       | Multiple mutations     |             |       | All mutations          |             |       |
|   | cross-validation level |             |       | cross-validation level |             |       | cross-validation level |             |       | cross-validation level |             |       | cross-validation level |             |       | cross-validation level |             |       |
|   | mut.                   | pos.        | prot. | mut.                   | pos.        | prot. | mut.                   | pos.        | prot. | mut.                   | pos.        | prot. | mut.                   | pos.        | prot. | mut.                   | pos.        | prot. |
| mGPFusion   | <b>0.81</b>            | <b>0.70</b> | 0.56  | <b>0.88</b>            | 0.61        | 0.49  | <b>0.83</b>            | 0.64        | 0.52  | 1.07                   | <b>1.26</b> | 1.61  | <b>1.33</b>            | 2.45        | 2.53  | <b>1.13</b>            | <b>1.87</b> | 1.84  |
| mGPFusion, only B62   | 0.79                   | 0.69        | 0.56  | 0.86                   | <b>0.64</b> | 0.50  | 0.82                   | <b>0.66</b> | 0.52  | 1.11                   | 1.30        | 1.62  | 1.43                   | <b>2.40</b> | 2.50  | 1.18                   | <b>1.85</b> | 1.84  |
| mGP   | 0.81                   | 0.51        | -     | 0.86                   | 0.52        | -     | 0.83                   | 0.50        | -     | <b>1.04</b>            | 1.54        | -     | 1.44                   | 2.65        | -     | 1.14                   | 2.09        | -     |
| mGP, only B62   | 0.76                   | 0.34        | -     | 0.86                   | 0.55        | -     | 0.80                   | 0.49        | -     | 1.26                   | 1.95        | -     | 1.45                   | 2.56        | -     | 1.30                   | 2.23        | -     |
| Rosetta scaled  | 0.65                   | 0.63        | -     | 0.51                   | 0.39        | -     | 0.60                   | 0.48        | -     | 1.35                   | 1.38        | -     | 2.49                   | 2.99        | -     | 1.66                   | 2.22        | -     |
| Predictions from off-the-shelf implementations with no cross-validation |                        |             |       |                        |             |       |                        |             |       |                        |             |       |                        |             |       |                        |             |       |
| Rosetta   |                        | 0.55        |       |                        | 0.40        |       |                        | 0.49        |       |                        | 1.63        |       |                        | 2.74        |       |                        | 1.92        |       |
| mCSM  |                        | 0.61        |       |                        | -           |       |                        | -           |       |                        | 1.40        |       |                        | -           |       |                        | -           |       |
| PoPMuSiC  |                        | 0.64        |       |                        | -           |       |                        | -           |       |                        | 1.37        |       |                        | -           |       |                        | -           |       |

Table 2: Comparison of different methods on the 15 protein dataset with respect to  $\rho$  and rmse. Mutation, position, and protein are referred to as mut., pos., and prot., respectively. Predictions from off-the-shelf implementations of Rosetta, mCSM and PoPMuSiC are used directly without cross-validation.

### 3.2 Predicting multiple mutations

Next, we tested stability prediction accuracies for variants containing either single or multiple mutations. Figure 3 shows a scatter plot of mGPFusion predictions for all 1537 single and multiple mutation variants (covering all 15 proteins) against the experimental  $\Delta\Delta G$  values using the mutation level (leave-one-out) cross-validation. The points are coloured by the number of simultaneous mutations in the variants, with 326 variants having at least 2 mutations (See Table 1). Figure 3 illustrates the mGPFusion’s overall high accuracy of  $\rho = 0.83$  and  $\text{rmse} = 1.13$  kcal/mol on both single and multiple mutations (See Table 2). Scatter plots for the individual proteins can be found in Supplementary Figure S3. Dehouck *et al.* (2009) suggested that considering the predictive power after removal of most badly predicted stability effects of mutations may give more relevant evaluation, as some of the experimental measurements may have been made in non-physiological conditions or affected by significant error, associated with a poorly resolved structure, or indexed incorrectly in the database. They thus reported correlation and rmse of the predictions after excluding 10 % of the predictions with most negative impacts on the correlation coefficient. Pires *et al.* (2014b) also reported their accuracy after 10 % outlier removal. If we remove the 10% worst predicted stability effects from the combined predictions, we achieve correlation  $\rho$  of 0.92 and rmse of 0.67 kcal/mol. We report these results for all the methods in Supplementary Table S3 and also present the error distribution in Supplementary Figure S5.

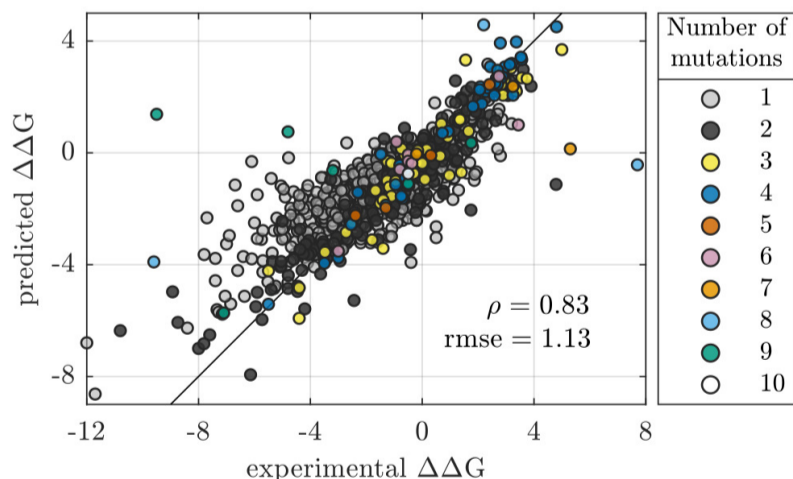


Figure 3: Scatter plot for the mutation level (leave-one-out) predictions made for all 15 proteins (See Table 1). The colour indicates the number of simultaneous mutations.

The high accuracy is retained for variants with multiple mutations as well ( $\rho = 0.88$  and  $\text{rmse} = 1.33$  kcal/mol, see Tables 2 and S2). Table 3 lists mGPFusion’s rmse for different number of simultaneous mutations. The model accuracy in fact improves up to 6 mutations. This is explained by the training set often containing the same single point mutations that appear in variants with multiple mutations. The model can then infer the combined effect of pointwise mutations. The model seems to fail when predicting the effects of 7-9 simultaneous mu-

tations. Most of these mutations (8/12) are for Ribonuclease (1RGG) and their effects seem to be exceptionally difficult to predict. This may be because only few of the point mutations that are part of the multiple mutations are present in the training data. However, these mutations seem to be exceptionally difficult to predict for Rosetta as well, which could indicate that the experimental measurements concerning these mutations are not quite accurate. PoPMuSiC and mCSM are unable to predict multiple mutations, while Rosetta supports them, but its rmse accuracy decreases already with two mutations.

| mutations occurrences | 1    | 2    | 3    | 4    | 5    | 6    | 7    | 8    | 9    | 10   |
|-----------------------|------|------|------|------|------|------|------|------|------|------|
| mGPFusion             | 1.07 | 1.06 | 0.80 | 0.51 | 0.40 | 1.01 | 3.02 | 5.89 | 5.16 | 0.25 |
| mGPFusion, only B62   | 1.11 | 1.12 | 0.77 | 0.59 | 0.29 | 1.14 | 3.00 | 6.78 | 5.56 | 0.11 |
| mGP                   | 1.04 | 1.03 | 0.61 | 0.50 | 0.18 | 0.92 | 3.23 | 6.18 | 6.75 | 0.08 |
| mGP, only B62         | 1.26 | 0.96 | 0.65 | 0.83 | 0.26 | 1.14 | 2.95 | 6.90 | 6.57 | 0.05 |
| Rosetta scaled        | 1.35 | 2.10 | 1.92 | 2.94 | 2.29 | 2.32 | 2.93 | 6.75 | 7.28 | 2.69 |
| Rosetta               | 1.63 | 2.27 | 2.11 | 3.78 | 2.93 | 2.21 | 2.92 | 5.80 | 7.45 | 3.42 |

Table 3: Root-mean-square errors for different number of simultaneous mutations for all 15 proteins, with models trained by leave-one-out cross-validation. Rosetta is added for comparison.

With multiple mutations, the decrease in performance between the position and mutation level cross-validations becomes clearer than with single mutations. With the position level cross-validation the stability effects of multiple mutations are predicted multiple times, which partly explains this loss of accuracy. For example, the effects of mutants with nine different simultaneous mutations, which were the most difficult cases in the mutation level cross-validation, are predicted nine times. Surprisingly, mGPFusion trained with protein level cross-validation achieves higher correlation and smaller errors than Rosetta; mGPFusion utilising simulated  $\Delta\Delta G$  values for only single mutations, can predict the stability effects of multiple mutations better than Rosetta.

### 3.3 Uncertainty of the predictions

Gaussian processes provide a mean  $\mu(\mathbf{x})$  and a standard deviation  $\sigma(\mathbf{x})$  for the stability prediction of a protein variant  $\mathbf{x}$ . The standard deviation allows estimation of the prediction accuracy even without test data. Figure 1 h) visualises the uncertainty of a few predictions made for the protein G (1PGA) when mutation level cross-validation is used. The estimated standard deviation allows a user to automatically identify low quality predictions that can appear e.g. in parts of the input protein space from which less data is included in model training. Conversely, in order to minimise the amount of uncertainty in the mGPFusion predictions, estimated standard deviation can be used to guide next experiments. The probabilistic nature of the predictions also admits an alternative error measure of negative log probability density (NLPD)  $\text{nlpd} = -\sum_{i=1}^N \log p(y_i|\mu(\mathbf{x}_i), \sigma^2(\mathbf{x}_i))$ , which can naturally take into account the prediction variance.

### 3.4 Effect of training set size

The results presented in Sections 3.1–3.3 used all available data for training with cross-validation to obtain unbiased performance measures. The inclusion of thousands of simulated variants allows the model

to learn accurate models with less experimentally measured variants. Hence, we study how the mGPFusion model with or without simulated data performs with reduced number of experimental observations. To facilitate this, we randomly selected subsets of experimental data of size 0, 10, 20, and so on. We learned the mGP and mGPFusion models with these reduced experimental data sets while always using the full simulated data sets. This also allows us to estimate how the models work with different number of cross-validation folds. For example, the point of a learning curve which utilises 2/3 or 4/5 of the training data correspond to an average of multiple 3-fold or 5-fold cross-validations, respectively.

The learning curve in Figure 4a) shows how the averaged correlation for protein 2LZM improves when the size of the experimental data set increases. The right-most values at  $N = 348$  are obtained with leave-one-out cross-validation. The inclusion of simulated data in mGPFusion (dark blue line) consistently improves the performance of mGP, which is trained without simulated data. Figure 4b) illustrate the difference in root mean square error. Learning curves for all proteins listed in Table 1 can be found from the Supplementary Figures S6-S8. When the number of experimental samples is zero, the mGPFusion model is trained solely using the simulated data with scaling  $0.57y^S$ , and the mGP model predicts the stability effect of every mutation as zero. The last point on the learning curves is obtained with mutation level cross-validation (see Tables 2 and S2).

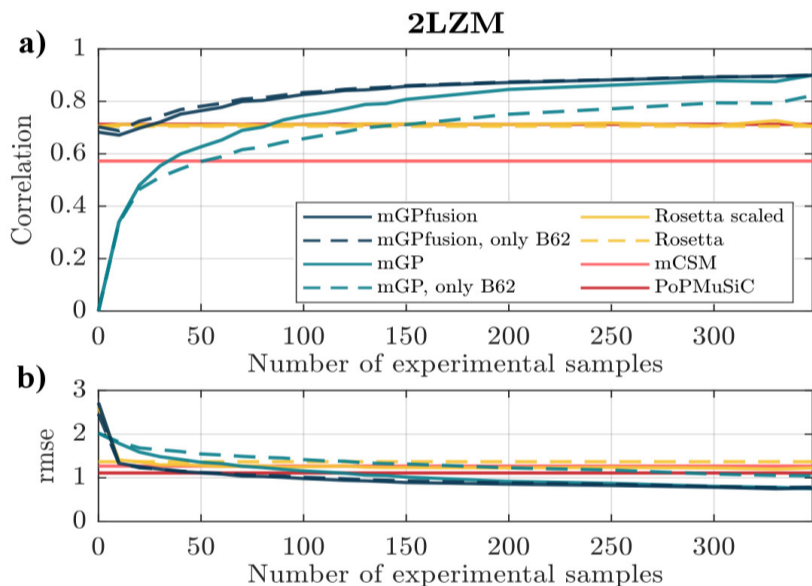


Figure 4: a) Correlation and b) root mean square error of predictions made by models with different number of experimental training samples for T4 Lysozyme (2LZM). The results of Rosetta, mCSM and PoPMuSiC are invariant to training data (because mCSM and PoPMuSiC are pre-trained), and are thus constant lines. For both figures, an average of 100 randomly selected training sets is taken at each point.

### 3.5 Effect of data fusion and multiple substitution matrices

In the beginning of the learning curves, when only little training data is available, mGPFusion quite consistently outperforms the mGP model, demonstrating that the additional simulated data improves the prediction accuracy. However, when more training data becomes available, the performance of mGP model is almost as good or sometimes even better than the performance of the mGPFusion model. This shows that if enough training data is available, it is not necessary to simulate additional data in order to obtain accurate predictions. Table 2 also shows, that the data fusion can compensate the lack of relevant training data – with the mGPFusion models that utilise the additional data, the decrease in accuracy is smaller when position level cross-validation is used instead of mutation level cross-validation, than with the mGP models.

The varying weights for the base kernels between different proteins (shown in Figure 2) already illustrated that different proteins benefit from different similarity measures for amino acid substitutions. The learning curves also support this observation – with some of the

proteins mGPFusion model trained with only one kernel that utilises BLOSUM62, provides approximately as good results as the mGPFusion model trained with multiple kernels. However, with many of the proteins, utilising just BLOSUM62 does not seem to be sufficient and the accuracy of the model can be improved by using different substitution matrices. Prior knowledge of appropriate substitution models for each protein could enable creation of accurate prediction models with just one substitution model, but the MKL seems to be a good tool for selecting suitable substitution models when such knowledge is not available. It seems that the data fusion and number or relevance of used substitution matrices can compensate each other – the learning curves show, that the difference between mGPFusion models trained with one or multiple kernels is smaller than the difference between the mGP models utilising one or multiple kernels. This indicates that if additional simulated data is exploited, the use of multiple or appropriate substitution models is not as important than without the data fusion. On the other hand, if data fusion is not applied, the use of MKL can more significantly improve the accuracy of the mGP model.

### 3.6 Effect of the Bayesian transformation on Rosetta

The Bayesian scaling of simulated Rosetta values, proposed in Section 2.2, improves the match of Rosetta simulated values to empirical  $\Delta\Delta G$  values even without using the Gaussian process framework. The Bayesian scaling improves the performance of standard Rosetta simulations from  $\rho = 0.55$  and  $\text{rmse} = 1.63$  kcal/mol to  $\rho = 0.65$  and  $\text{rmse} = 1.35$  kcal/mol (see Table 2 and Supplementary Table S2). This shows that the scaling proposed by Kellogg *et al.* (2011) indeed is not always the optimal scaling and significant improvements can be gained by optimising the scaling using a set of training data.

Figure 1 g) visualises the Bayesian scaling for protein 1PGA, where the very destabilising  $\Delta\Delta G$  values are dampened by the scaling (black dots) to less extreme values by matching the scaled simulated values to the experimental points (blue circles). The black dots along the scaling curve indicate the grid of point mutations after transformation. The scaling variance  $\sigma_T^2$  is indicated by the green region’s vertical width, and on the right panel. The scaling tends to dampen very small values into less extreme stabilities, while it also estimates higher uncertainties for stability values further away from  $\Delta\Delta G = 0$ . However, the scalings vary between different proteins, as can be seen from the transformations for each of the 15 proteins presented in Supplementary Figure S9.

## 4 Conclusions

We present a novel method mGPFusion for predicting stability effects of both single and multiple simultaneous mutations. mGPFusion utilises structural information in form of contact maps and integrates that with amino acid residues and combines both experimental and comprehensive simulated measurements of mutations’ stability effects. In contrast to earlier general-purpose stability models, mGPFusion model is protein-specific by design, which improves the accuracy but necessitates having a set of experimental measurements from the protein. In practise small datasets of 10–20 experimental observations were found to provide state-of-the-art accuracy models when supplemented by large simulation datasets.

An important advantage over most state-of-the-art machine learning methods is that mGPFusion is able to predict the effects of multiple simultaneous mutations in addition to single point mutations. Our experiments show that mGPFusion is reliable in predicting up to six simultaneous mutations in our dataset. Furthermore, the Gaussian process framework provide a way to estimate the (un)certainty of the predictions even without a separate test set. We additionally proposed a novel Bayesian scaling method to re-calibrate simulated protein stability values against experimental observations. This is a crucial part of the mGPFusion model, and also alone improved protein-specific Rosetta stability predictions by calibrating them using experimental data.

mGPFusion is best suited for a situation, where a protein is thoroughly experimented on and accurate predictions for stability effects

upon mutations are needed. It takes some time to set up the framework and train the model, but after that new predictions can be made in fractions of a second. The most time-consuming part is running the simulations with Rosetta, at least when the most accurate protocol 16 is used. Simulating all 19 possible point mutations for one position took about 12 hours, but simulations for different positions can be run on parallel. The time needed for training the prediction model depends on the amount of experimental and simulated training data. With no simulated data, the training time ranged from few seconds to few minutes. With data fusion and a single kernel, the training time was under an hour. With data fusion and MKL with 21 kernels, the training time was from a few minutes to a day.

## Acknowledgements

We acknowledge the computational resources provided by the Aalto Science-IT.

## Funding

This work has been supported by the Academy of Finland Center of Excellence in Systems Immunology and Physiology, the Academy of Finland grants no. 260403 and 299915, and the Finnish Funding Agency for Innovation Tekes (grant no 40128/14, Living Factories).

## References

- Alberts, B., Johnson, A., Lewis, J., Raff, M., Roberts, K., and Walter, P. (2007). *Molecular biology of the cell*. Garland Science, 5 edition.
- Anslyn, E. V. and Dougherty, D. A. (2006). *Modern physical organic chemistry*. University Science Books.
- Berman, H. M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T., Weissig, H., Shindyalov, I. N., and Bourne, P. E. (2000). The protein data bank. *Nucleic acids research*, **28**(1), 235–242.
- Bommarius, A. S., Blum, J. K., and Abrahamson, M. J. (2011). Status of protein engineering for biocatalysts: how to design an industrially useful biocatalyst. *Current opinion in chemical biology*, **15**(2), 194–200.
- Branden, C. and Tooze, J. (1999). *Introduction to protein structure*. Garland, 2 edition.
- Capriotti, E., Fariselli, P., and Casadio, R. (2005a). I-Mutant2.0: predicting stability changes upon mutation from the protein sequence or structure. *Nucleic acids research*, **33**(suppl 2), W306–W310.
- Capriotti, E., Fariselli, P., Calabrese, R., and Casadio, R. (2005b). Predicting protein stability changes from sequences using support vector machines. *Bioinformatics*, **21**(suppl 2), ii54–ii58.
- Capriotti, E., Fariselli, P., Rossi, I., and Casadio, R. (2008). A three-state prediction of single point mutations on protein stability changes. *BMC bioinformatics*, **9**(2).
- Chen, C.-W., Lin, J., and Chu, Y.-W. (2013). iStable: off-the-shelf predictor integration for predicting protein stability changes. *BMC bioinformatics*, **14**(2).
- Cheng, J., Randall, A., and Baldi, P. (2006). Prediction of protein stability changes for single-site mutations using support vector machines. *Proteins: Structure, Function, and Bioinformatics*, **62**(4), 1125–1132.
- Cherry, J. R. and Fidantsef, A. L. (2003). Directed evolution of industrial enzymes: an update. *Current opinion in biotechnology*, **14**(4), 438–443.
- Cichonska, A., Ravikumar, B., Parri, E., Timonen, S., Pahikkala, T., Airola, A., Wennerberg, K., Rousu, J., and Aittokallio, T. (2017). Computational-experimental approach to drug-target interaction mapping: A case study on kinase inhibitors. *PLoS computational biology*, **13**(8), e1005678.
- Dehouck, Y., Grosfils, A., Folch, B., Gilis, D., Bogaerts, P., and Rooman, M. (2009). Fast and accurate predictions of protein stability changes upon mutations using statistical potentials and neural networks: PoPMuSiC-2.0. *Bioinformatics*, **25**(19), 2537–2543.
- Folkman, L., Stantic, B., and Sattar, A. (2014). Feature-based multiple models improve classification of mutation-induced stability changes. *BMC genomics*, **15**(Suppl 4).
- Giguere, S., Marchand, M., Laviolette, F., Drouin, A., and Corbeil, J. (2013). Learning a peptide-protein binding affinity predictor with kernel ridge regression. *BMC bioinformatics*, **14**(1), 82.
- Giollo, M., Martin, A. J., Walsh, I., Ferrari, C., and Tosatto, S. C. (2014). NeEMO: a method using residue interaction networks to improve prediction of protein stability upon mutation. *BMC genomics*, **15**(4), 1.
- Henikoff, S. and Henikoff, J. G. (1992). Amino acid substitution matrices from protein blocks. *Proceedings of the National Academy of Sciences*, **89**(22), 10915–10919.
- Kawashima, S., Pokarowski, P., Pokarowska, M., Kolinski, A., Katayama, T., and Kanehisa, M. (2008). AAindex: amino acid index database, progress report 2008. *Nucleic acids research*, **36**(suppl 1), D202–D205.
- Kellogg, E. H., Leaver-Fay, A., and Baker, D. (2011). Role of conformational sampling in computing mutation-induced changes in protein structure and stability. *Proteins: Structure, Function, and Bioinformatics*, **79**(3), 830–838.
- Kirk, O., Borchert, T. V., and Fuglsang, C. C. (2002). Industrial enzyme applications. *Current opinion in biotechnology*, **13**(4), 345–351.
- Kumar, M. S., Bava, K. A., Gromiha, M. M., Prabakaran, P., Kitajima, K., Uedaira, H., and Sarai, A. (2006). ProTherm and ProNIT: thermodynamic databases for proteins and protein–nucleic acid interactions. *Nucleic Acids Research*, **34**(suppl 1), D204–D206.
- Leaver-Fay, A., Tyka, M., Lewis, S. M., Lange, O. F., Thompson, J., Jacak, R., Kaufman, K., Renfrew, P. D., Smith, C. A., Sheffler, W., et al. (2011). ROSETTA3: an object-oriented software suite for the simulation and design of macromolecules. *Methods in enzymology*, **487**, 545.
- Liu, J. and Kang, X. (2012). Grading amino acid properties increased accuracies of single point mutation on protein stability prediction. *BMC bioinformatics*, **13**(1), 1.
- Menchetti, S., Costa, F., and Frasconi, P. (2005). Weighted decomposition kernels. In *Proceedings of the 22nd international conference on Machine learning*, pages 585–592. ACM.
- Pace, C. N. and Scholtz, J. M. (1997). Measuring the conformational stability of a protein. *Protein structure: A practical approach*, **2**, 299–321.
- Pace, C. N. and Shaw, K. L. (2000). Linear extrapolation method of analyzing solvent denaturation curves. *Proteins: Structure, Function, and Bioinformatics*, **41**(S4), 1–7.
- Pires, D. E., Ascher, D. B., and Blundell, T. L. (2014a). DUET: a server for predicting effects of mutations on protein stability using an integrated computational approach. *Nucleic acids research*, page gku411.
- Pires, D. E., Ascher, D. B., and Blundell, T. L. (2014b). mCSM: predicting the effects of mutations in proteins using graph-based signatures. *Bioinformatics*, **30**(3), 335–342.
- Potapov, V., Cohen, M., and Schreiber, G. (2009). Assessing computational methods for predicting protein stability upon mutation: good on average but not in the details. *Protein Engineering Design and Selection*, **22**(9), 553–560.
- Rapley, R. and Walker, J. M. (2000). *Molecular Biology and Biotechnology*. Royal Society of Chemistry, 4 edition.
- Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian processes for machine learning*. The MIT Press.
- Sanchez, S. and Demain, A. L. (2010). Enzymes and bioconversions of industrial, pharmaceutical, and biotechnological significance. *Organic Process Research & Development*, **15**(1), 224–230.
- Schmidt, M. W., Berg, E., Friedlander, M. P., and Murphy, K. P. (2009). Optimizing costly functions with simple constraints: A limited-memory projected quasi-newton algorithm. In *International Conference on Artificial Intelligence and Statistics*, page None.
- Shawe-Taylor, J. and Cristianini, N. (2004). *Kernel methods for pattern analysis*. Cambridge university press.
- Tian, J., Wu, N., Chu, X., and Fan, Y. (2010). Predicting changes in protein thermostability brought about by single- or multi-site mutations. *BMC bioinformatics*, **11**(1), 1.
- Tokuriki, N. and Tawfik, D. S. (2009). Stability effects of mutations and protein evolvability. *Current opinion in structural biology*, **19**(5), 596–604.
- Tomii, K. and Kanehisa, M. (1996). Analysis of amino acid indices and mutation matrices for sequence comparison and structure prediction of proteins. *Protein Engineering, Design and Selection*, **9**(1), 27–36.
- Vishwanathan, S. V. N., Schraudolph, N. N., Kondor, R., and Borgwardt, K. M. (2010). Graph kernels. *The Journal of Machine Learning Research*, **11**, 1201–1242.
- Wainreb, G., Wolf, L., Ashkenazy, H., Dehouck, Y., and Ben-Tal, N. (2011). Protein stability: a single recorded mutation aids in predicting the effects of other mutations in the same amino acid site. *Bioinformatics*, **27**(23), 3286–3292.



## Supplementary material

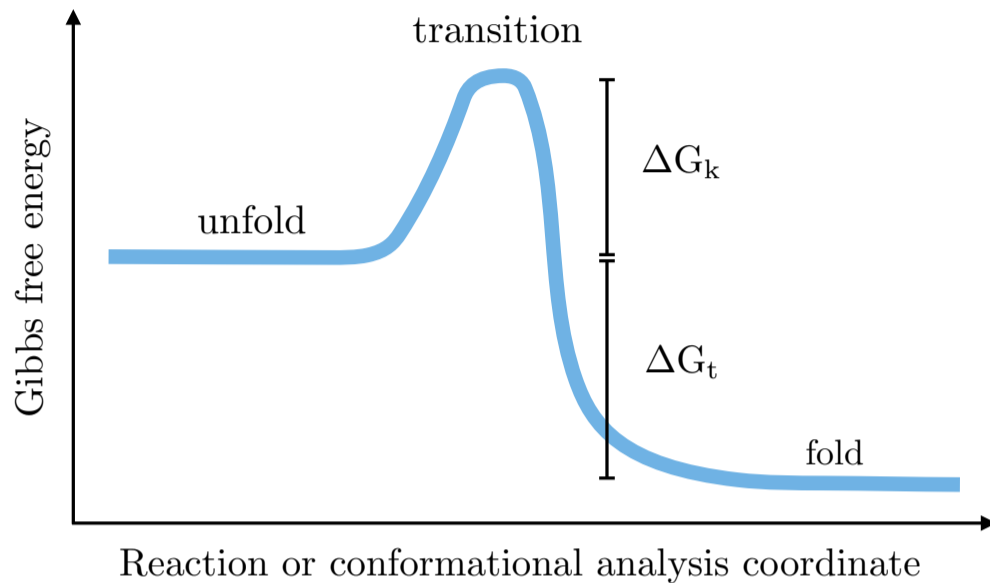


Figure S5: The stability of a protein is determined by the thermodynamic and kinetic stabilities,  $\Delta G_t$  and  $\Delta G_k$ , respectively. We only consider the thermodynamic stability.

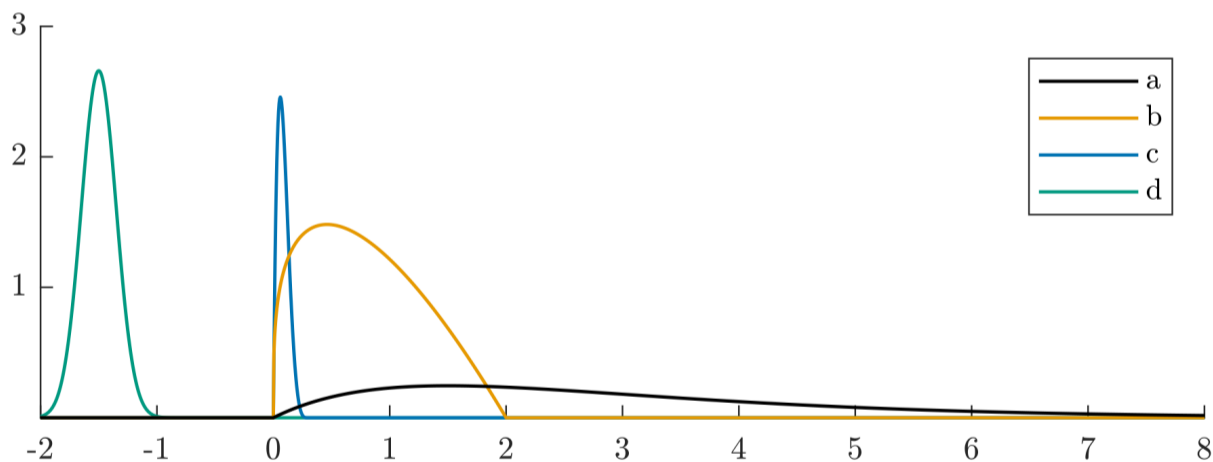


Figure S6: Priors presented by Equation 2. Here  $\mu_d = -1.5$ , the most likely value for  $-a$ . Other hyperparameter values are presented in Table S4.

Table S4: Values for the hyperparameters used in the priors of  $a, b, c, d, \sigma_E$  and  $\sigma_S$  presented in Equations 2 and 6, respectively.

| $a$             | $b$              | $c$            | $d$               | $\sigma_E$       | $\sigma_S$        |
|-----------------|------------------|----------------|-------------------|------------------|-------------------|
| $\alpha_a = 2$  | $\alpha_b = 1.3$ | $\alpha_c = 2$ | $\mu_d = -a$      | $\alpha_E = 2.5$ | $\alpha_S = 50$   |
| $\beta_a = 1.5$ | $\beta_b = 2$    | $\beta_c = 5$  | $\sigma_d = 0.15$ | $\beta_E = 0.02$ | $\beta_S = 0.007$ |

The partial derivatives of the marginal likelihood with respect to the parameters  $\phi$  are obtained from Equation (9) as follows:

$$\begin{aligned}\frac{\partial}{\partial\phi_j}\log p(\mathbf{y}|X,\phi) &= \frac{1}{2}\mathbf{y}^T K_\phi^{-1} \frac{\partial K_\phi}{\partial\phi_j} K_\phi^{-1} \mathbf{y} - \frac{1}{2}\text{tr}\left(K_\phi^{-1} \frac{\partial K_\phi}{\partial\phi_j}\right) \\ &= \frac{1}{2}\text{tr}\left(\left(\boldsymbol{\alpha}\boldsymbol{\alpha}^T - K_\phi^{-1}\right) \frac{\partial K_\phi}{\partial\theta_j}\right),\end{aligned}\tag{S11}$$

where  $\boldsymbol{\alpha} = K_\phi^{-1}\mathbf{y}$ ,  $K_\phi$  is determined as

$$K_\phi = \sum_{m=1}^M w_m K_m^{\gamma_m} + \text{diag}\left(\begin{array}{c} \sigma_0 \\ \sigma_E \mathbf{1}_{N_E} \\ \sigma_E \mathbf{1}_{N_E} + \sigma_S \mathbf{1}_{N_S} + t\boldsymbol{\sigma}_T \end{array}\right)^2$$

and the partial derivatives of  $K_\phi$  with respect to the optimised parameters are

$$\frac{\partial K_\phi}{\partial\sigma_E} = \text{diag}\left(\begin{array}{c} 0 \\ 2\sigma_E \mathbf{1}_{N_E} \\ 2(\sigma_E \mathbf{1}_{N_E} + \sigma_S \mathbf{1}_{N_S} + t\boldsymbol{\sigma}_T) \end{array}\right)\tag{S12}$$

$$\frac{\partial K_\phi}{\partial\sigma_R} = \text{diag}\left(\begin{array}{c} 0 \\ \mathbf{0}_{N_E} \\ 2(\sigma_E \mathbf{1}_{N_E} + \sigma_S \mathbf{1}_{N_S} + t\boldsymbol{\sigma}_T) \end{array}\right)\tag{S13}$$

$$\frac{\partial K_\phi}{\partial t} = \text{diag}\left(\begin{array}{c} 0 \\ \mathbf{0}_{N_E} \\ 2(\sigma_E \mathbf{1}_{N_E} + \sigma_S \mathbf{1}_{N_S} + t)\boldsymbol{\sigma}_T \end{array}\right)\tag{S14}$$

$$\frac{\partial K_\phi}{\partial w_m} = K_m^{\gamma_m}\tag{S15}$$

$$\frac{\partial K_\phi}{\partial\gamma_m} = w_m K_m^{\gamma_m} \log K_m\tag{S16}$$

Correlation  $\rho$  and root-mean-square error rmse for the predictions are determined as

$$\rho = \frac{\sum_{i=1}^{N_*} (y_i - \bar{y})(\mu(\mathbf{x}_i) - \bar{\mu})}{\sqrt{\sum_{i=1}^{N_*} (y_i - \bar{y})^2 \sum_{i=1}^{N_*} (\mu(\mathbf{x}_i) - \bar{\mu})^2}}\tag{S17}$$

$$\text{rmse} = \sqrt{\frac{1}{N_*} \sum_{i=1}^{N_*} (y_i - \mu(\mathbf{x}_i))^2},\tag{S18}$$

where  $\bar{y}$  is the mean of the experimentally measured values,  $\mu(\mathbf{x}_i)$  is prediction mean,  $\bar{\mu}$  is the average of all prediction means, and  $N_*$  is the number of predictions.

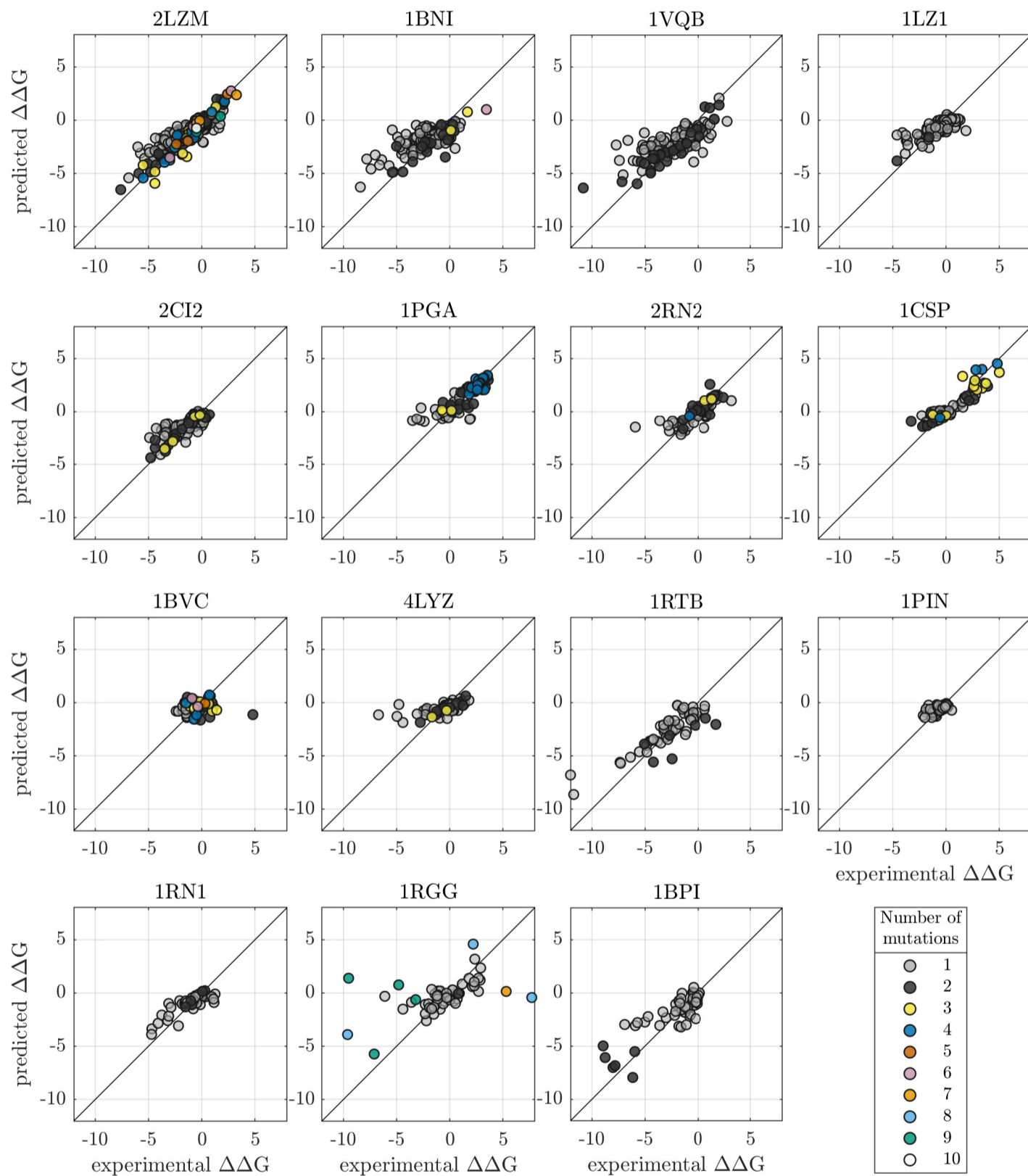


Figure S7: Mutation-level predictions for all 15 proteins presented in Table 1. The predictions are coloured by the number of simultaneous mutations.

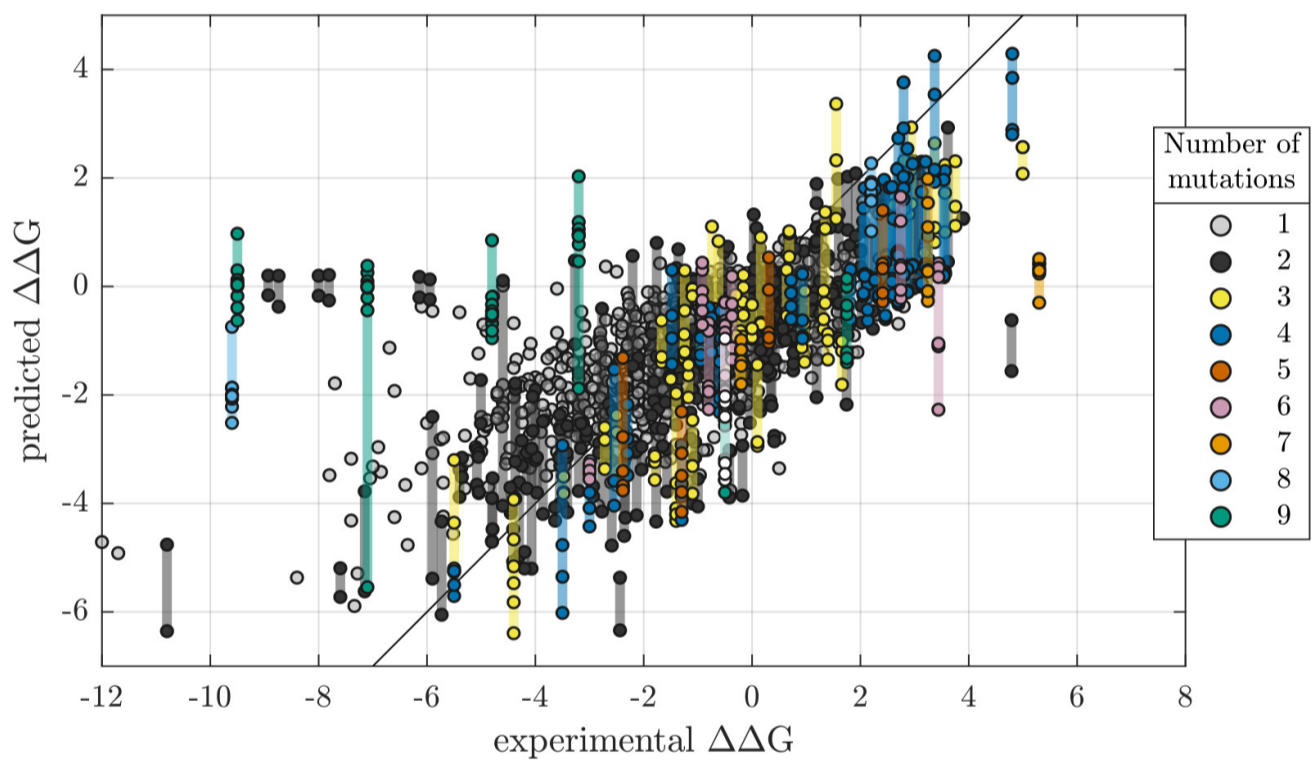


Figure S8: Position-level predictions for all 15 proteins. When the effects of a mutant are predicted multiple times, they are connected by a line.



Table S5: (Continues on the next page) Comparison of different methods on the 15 protein dataset with respect to  $\rho$  and rmse. Off-the-shelf implementations of Rosetta, mCSM and PoPMuSiC are used directly without cross-validation.

| Protein  | Method              | Correlation $\rho$ |       |       |                    |       |       |               |       |       | rmse            |      |       |                    |      |       |               |      |       |
|----------|---------------------|--------------------|-------|-------|--------------------|-------|-------|---------------|-------|-------|-----------------|------|-------|--------------------|------|-------|---------------|------|-------|
|          |                     | Point mutations    |       |       | Multiple mutations |       |       | All mutations |       |       | Point mutations |      |       | Multiple mutations |      |       | All mutations |      |       |
|          |                     | mut.               | pos.  | prot. | mut.               | pos.  | prot. | mut.          | pos.  | prot. | mut.            | pos. | prot. | mut.               | pos. | prot. | mut.          | pos. | prot. |
| 2LZM     | mGPFusion           | 0.87               | 0.80  | 0.75  | 0.96               | 0.83  | 0.64  | 0.90          | 0.76  | 0.68  | 0.82            | 1.02 | 1.12  | 0.57               | 1.53 | 2.14  | 0.76          | 1.30 | 1.43  |
|          | mGPFusion, only B62 | 0.86               | 0.77  | 0.75  | 0.96               | 0.87  | 0.64  | 0.90          | 0.82  | 0.69  | 0.84            | 1.08 | 1.12  | 0.59               | 1.16 | 2.11  | 0.79          | 1.13 | 1.42  |
|          | mGP                 | 0.86               | 0.59  | -     | 0.97               | 0.85  | -     | 0.90          | 0.72  | -     | 0.82            | 1.34 | -     | 0.48               | 1.24 | -     | 0.75          | 1.29 | -     |
|          | mGP, only B62       | 0.75               | 0.37  | -     | 0.94               | 0.77  | -     | 0.82          | 0.61  | -     | 1.12            | 1.93 | -     | 0.78               | 1.44 | -     | 1.05          | 1.70 | -     |
|          | Rosetta scaled      | 0.74               | 0.73  | -     | 0.68               | 0.66  | -     | 0.70          | 0.65  | -     | 1.05            | 1.06 | -     | 1.65               | 1.84 | -     | 1.23          | 1.51 | -     |
| 1BNI     | Rosetta             |                    | 0.75  |       |                    | 0.68  |       |               | 0.71  |       |                 | 1.13 |       |                    | 1.93 |       |               | 1.37 |       |
|          | mCSM                |                    | 0.57  |       |                    | -     |       |               | -     |       |                 | 1.27 |       |                    | -    |       |               | -    |       |
|          | PoPMuSiC            |                    | 0.71  |       |                    | -     |       |               | -     |       |                 | 1.11 |       |                    | -    |       |               | -    |       |
|          | mGPFusion           | 0.77               | 0.64  | 0.62  | 0.86               | 0.70  | 0.39  | 0.77          | 0.55  | 0.57  | 1.21            | 1.37 | 1.69  | 1.28               | 2.49 | 2.17  | 1.22          | 1.67 | 1.75  |
|          | mGPFusion, only B62 | 0.74               | 0.61  | 0.62  | 0.84               | 0.79  | 0.40  | 0.74          | 0.60  | 0.57  | 1.27            | 1.45 | 1.69  | 1.24               | 2.03 | 2.14  | 1.27          | 1.60 | 1.75  |
| 1VQB     | mGP                 | 0.81               | 0.65  | -     | 0.86               | 0.82  | -     | 0.80          | 0.63  | -     | 1.08            | 1.40 | -     | 1.32               | 2.00 | -     | 1.11          | 1.55 | -     |
|          | mGP, only B62       | 0.61               | 0.48  | -     | 0.87               | 0.85  | -     | 0.63          | 0.41  | -     | 1.61            | 2.32 | -     | 1.15               | 1.67 | -     | 1.57          | 2.20 | -     |
|          | Rosetta scaled      | 0.59               | 0.58  | -     | 0.17               | 0.29  | -     | 0.53          | 0.45  | -     | 1.58            | 1.58 | -     | 2.51               | 2.90 | -     | 1.70          | 1.94 | -     |
|          | Rosetta             |                    | 0.62  |       |                    | 0.18  |       |               | 0.56  |       |                 | 1.70 |       |                    | 2.33 |       |               | 1.77 |       |
|          | mCSM                |                    | 0.60  |       |                    | -     |       |               | -     |       |                 | 1.62 |       |                    | -    |       |               | -    |       |
| 1LZ1     | PoPMuSiC            |                    | 0.66  |       |                    | -     |       |               | -     |       |                 | 1.53 |       |                    | -    |       |               | -    |       |
|          | mGPFusion           | 0.67               | 0.50  | 0.49  | 0.93               | 0.83  | 0.75  | 0.76          | 0.69  | 0.60  | 1.71            | 1.94 | 2.25  | 1.15               | 1.62 | 2.06  | 1.59          | 1.82 | 2.20  |
|          | mGPFusion, only B62 | 0.65               | 0.53  | 0.50  | 0.91               | 0.82  | 0.73  | 0.75          | 0.69  | 0.58  | 1.75            | 1.94 | 2.25  | 1.35               | 1.82 | 2.43  | 1.66          | 1.89 | 2.30  |
|          | mGP                 | 0.79               | 0.12  | -     | 0.96               | 0.70  | -     | 0.85          | 0.50  | -     | 1.41            | 2.41 | -     | 0.70               | 1.97 | -     | 1.27          | 2.24 | -     |
|          | mGP, only B62       | 0.79               | 0.29  | -     | 0.97               | 0.75  | -     | 0.85          | 0.55  | -     | 1.50            | 2.89 | -     | 0.63               | 2.30 | -     | 1.33          | 2.66 | -     |
| 2CI2     | Rosetta scaled      | 0.47               | 0.46  | -     | 0.71               | 0.68  | -     | 0.57          | 0.59  | -     | 1.99            | 2.00 | -     | 1.90               | 1.96 | -     | 1.97          | 1.99 | -     |
|          | Rosetta             |                    | 0.49  |       |                    | 0.73  |       |               | 0.59  |       |                 | 2.26 |       |                    | 2.06 |       |               | 2.21 |       |
|          | mCSM                |                    | 0.53  |       |                    | -     |       |               | -     |       |                 | 2.24 |       |                    | -    |       |               | -    |       |
|          | PoPMuSiC            |                    | 0.51  |       |                    | -     |       |               | -     |       |                 | 2.29 |       |                    | -    |       |               | -    |       |
|          | mGPFusion           | 0.75               | 0.59  | 0.58  | 1.00               | 0.11  | -1.00 | 0.77          | 0.57  | 0.52  | 0.83            | 0.99 | 1.06  | 0.56               | 2.40 | 3.75  | 0.83          | 1.07 | 1.16  |
| 1PGA     | mGPFusion, only B62 | 0.73               | 0.56  | 0.59  | 1.00               | 0.05  | -1.00 | 0.76          | 0.56  | 0.53  | 0.86            | 1.03 | 1.06  | 1.10               | 2.41 | 3.51  | 0.87          | 1.10 | 1.14  |
|          | mGP                 | 0.75               | 0.39  | -     | 1.00               | 0.56  | -     | 0.78          | 0.47  | -     | 0.81            | 1.15 | -     | 0.13               | 1.65 | -     | 0.80          | 1.17 | -     |
|          | mGP, only B62       | 0.71               | -0.31 | -     | 1.00               | 0.42  | -     | 0.74          | 0.21  | -     | 0.91            | 1.43 | -     | 0.27               | 2.36 | -     | 0.90          | 1.47 | -     |
|          | Rosetta scaled      | 0.57               | 0.55  | -     | -1.00              | -1.00 | -     | 0.53          | 0.46  | -     | 0.99            | 1.01 | -     | 3.23               | 3.25 | -     | 1.07          | 1.15 | -     |
|          | Rosetta             |                    | 0.59  |       |                    | -1.00 |       |               | 0.55  |       |                 | 1.04 |       |                    | 3.41 |       |               | 1.12 |       |
| 2RN2     | mCSM                |                    | 0.67  |       |                    | -     |       |               | -     |       |                 | 0.97 |       |                    | -    |       |               | -    |       |
|          | PoPMuSiC            |                    | 0.64  |       |                    | -     |       |               | -     |       |                 | 0.95 |       |                    | -    |       |               | -    |       |
|          | mGPFusion           | 0.73               | 0.72  | 0.64  | 0.95               | 0.87  | 0.85  | 0.82          | 0.81  | 0.71  | 0.85            | 0.90 | 1.07  | 0.55               | 0.80 | 1.21  | 0.80          | 0.86 | 1.10  |
|          | mGPFusion, only B62 | 0.69               | 0.67  | 0.63  | 0.92               | 0.86  | 0.86  | 0.79          | 0.79  | 0.72  | 0.91            | 0.97 | 1.07  | 0.71               | 1.01 | 1.12  | 0.87          | 0.99 | 1.08  |
|          | mGP                 | 0.65               | 0.61  | -     | 0.92               | 0.79  | -     | 0.76          | 0.72  | -     | 0.95            | 1.02 | -     | 0.66               | 1.01 | -     | 0.90          | 1.02 | -     |
| 1CSP     | mGP, only B62       | 0.51               | 0.74  | -     | 0.92               | 0.71  | -     | 0.68          | 0.63  | -     | 1.16            | 1.39 | -     | 0.71               | 1.40 | -     | 1.08          | 1.39 | -     |
|          | Rosetta scaled      | 0.60               | 0.60  | -     | 0.61               | 0.61  | -     | 0.63          | 0.63  | -     | 1.00            | 1.00 | -     | 1.27               | 1.27 | -     | 1.06          | 1.11 | -     |
|          | Rosetta             |                    | 0.63  |       |                    | 0.62  |       |               | 0.65  |       |                 | 1.09 |       |                    | 1.30 |       |               | 1.13 |       |
|          | mCSM                |                    | 0.74  |       |                    | -     |       |               | -     |       |                 | 0.86 |       |                    | -    |       |               | -    |       |
|          | PoPMuSiC            |                    | 0.75  |       |                    | -     |       |               | -     |       |                 | 0.85 |       |                    | -    |       |               | -    |       |
| 1CSP     | mGPFusion           | 0.68               | 0.47  | 0.69  | 0.90               | 0.35  | 0.32  | 0.85          | 0.43  | 0.50  | 1.26            | 1.54 | 1.64  | 0.53               | 2.09 | 2.74  | 0.88          | 2.00 | 2.38  |
|          | mGPFusion, only B62 | 0.82               | 0.59  | 0.71  | 0.76               | 0.60  | 0.62  | 0.81          | 0.53  | 0.70  | 0.87            | 1.22 | 1.06  | 0.69               | 0.95 | 0.88  | 0.83          | 1.13 | 1.02  |
|          | mGP                 | 0.62               | 0.61  | -     | 0.93               | -0.24 | -     | 0.84          | -0.14 | -     | 1.40            | 1.58 | -     | 0.45               | 3.01 | -     | 0.94          | 2.81 | -     |
|          | mGP, only B62       | 0.57               | -0.46 | -     | 0.92               | -0.08 | -     | 0.81          | 0.06  | -     | 1.53            | 1.73 | -     | 0.47               | 2.07 | -     | 1.02          | 2.01 | -     |
|          | Rosetta scaled      | 0.69               | 0.59  | -     | 0.09               | 0.07  | -     | 0.24          | 0.11  | -     | 1.21            | 1.42 | -     | 2.81               | 3.09 | -     | 2.33          | 2.87 | -     |
| 1CSP     | Rosetta             |                    | 0.69  |       |                    | 0.03  |       |               | 0.28  |       |                 | 1.70 |       |                    | 3.51 |       |               | 2.95 |       |
|          | mCSM                |                    | -0.10 |       |                    | -     |       |               | -     |       |                 | 1.94 |       |                    | -    |       |               | -    |       |
|          | PoPMuSiC            |                    | 0.28  |       |                    | -     |       |               | -     |       |                 | 1.89 |       |                    | -    |       |               | -    |       |
|          | mGPFusion           | 0.79               | 0.58  | 0.71  | 0.78               | 0.60  | 0.61  | 0.79          | 0.53  | 0.70  | 0.91            | 1.21 | 1.05  | 0.67               | 1.01 | 0.91  | 0.86          | 1.14 | 1.02  |
|          | mGPFusion, only B62 | 0.82               | 0.59  | 0.71  | 0.76               | 0.60  | 0.62  | 0.81          | 0.53  | 0.70  | 0.87            | 1.22 | 1.59  | 0.69               | 0.95 | 1.18  | 0.83          | 1.13 | 1.51  |
| 1CSP     | mGP                 | 0.77               | 0.12  | -     | 0.75               | 0.42  | -     | 0.77          | 0.22  | -     | 0.93            | 1.45 | -     | 0.74               | 1.21 | -     | 0.89          | 1.36 | -     |
|          | mGP, only B62       | 0.83               | 0.09  | -     | 0.77               | 0.42  | -     | 0.82          | 0.23  | -     | 0.82            | 1.45 | -     | 0.68               | 1.20 | -     | 0.80          | 1.36 | -     |
|          | Rosetta scaled      | 0.66               | 0.64  | -     | 0.48               | 0.50  | -     | 0.62          | 0.57  | -     | 1.09            | 1.13 | -     | 1.20               | 1.08 | -     | 1.12          | 1.11 | -     |
|          | Rosetta             |                    | 0.70  |       |                    | 0.47  |       |               | 0.65  |       |                 | 1.07 |       |                    | 1.25 |       |               | 1.11 |       |
|          | mCSM                |                    | 0.71  |       |                    | -     |       |               | -     |       |                 | 1.04 |       |                    | -    |       |               | -    |       |
| 1CSP     | PoPMuSiC            |                    | 0.71  |       |                    | -     |       |               | -     |       |                 | 1.16 |       |                    | -    |       |               | -    |       |
|          | mGPFusion           | 0.85               | 0.23  | 0.33  | 0.92               | 0.73  | 0.48  | 0.92          | 0.75  | 0.38  | 0.64            | 1.04 | 1.10  | 0.91               | 1.66 | 2.58  | 0.75          | 1.45 | 1.80  |
|          | mGPFusion, only B62 | 0.86               | 0.22  | 0.34  | 0.91               | 0.69  | 0.59  | 0.91          | 0.72  | 0.54  | 0.65            | 1.04 | 1.10  | 0.96               | 1.87 | 2.13  | 0.78          | 1.60 | 1.57  |
|          | mGP                 | 0.88               | -0.06 | -     | 0.94               | 0.75  | -     | 0.94          | 0.77  | -     | 0.54            | 1.07 | -     | 0.76               | 1.51 | -     | 0.63          | 1.36 | -     |
|          | mGP, only B62       | 0.87               | -0.37 | -     | 0.92               | 0.71  | -     | 0.92          | 0.72  | -     | 0.60            | 1.12 | -     | 0.86               | 1.59 | -     | 0.71          | 1.43 | -     |
| 1CSP     | Rosetta scaled      | 0.23               | 0.20  | -     | 0.68               | 0.69  | -     | 0.59          | 0.64  | -     | 1.04            | 1.06 | -     | 2.19               | 2.29 | -     | 1.58          | 1.92 | -     |
|          | Rosetta             |                    | 0.33  |       |                    | 0.68  |       |               | 0.60  |       |                 | 1.11 |       |                    | 1.92 |       |               | 1.47 |       |
|          | mCSM                |                    | 0.42  |       |                    | -     |       |               | -     |       |                 | 1.02 |       |                    | -    |       |               | -    |       |
| PoPMuSiC |                     | 0.48               |       |       | -                  |       |       | -             |       |       | 0.99            |      |       | -                  |      |       | -             |      |       |

Table S5: (Continued) Comparison of different methods on the 15 protein dataset with respect to  $\rho$  and rmse. Off-the-shelf implementations of Rosetta, mCSM and PoPMuSiC are used directly without cross-validation.

| Protein | Method              | Correlation $\rho$ |       |       |                    |       |       |               |       |       | rmse            |      |       |                    |      |       |               |      |       |
|---------|---------------------|--------------------|-------|-------|--------------------|-------|-------|---------------|-------|-------|-----------------|------|-------|--------------------|------|-------|---------------|------|-------|
|         |                     | Point mutations    |       |       | Multiple mutations |       |       | All mutations |       |       | Point mutations |      |       | Multiple mutations |      |       | All mutations |      |       |
|         |                     | mut.               | pos.  | level | mut.               | pos.  | level | mut.          | pos.  | level | mut.            | pos. | level | mut.               | pos. | level | mut.          | pos. | level |
| 1BVC    | mGPFusion           | 0.41               | 0.43  | 0.48  | -0.28              | -0.25 | -0.63 | 0.08          | -0.09 | 0.09  | 0.74            | 0.72 | 1.65  | 1.64               | 1.48 | 2.70  | 1.09          | 1.21 | 2.02  |
|         | mGPFusion, only B62 | 0.50               | 0.48  | 0.48  | -0.25              | -0.25 | -0.66 | 0.14          | -0.07 | 0.14  | 0.70            | 0.71 | 1.65  | 1.62               | 1.47 | 2.43  | 1.06          | 1.20 | 1.92  |
|         | mGP                 | -0.05              | -0.12 | -     | -0.10              | -0.21 | -     | -0.13         | -0.23 | -     | 0.99            | 1.00 | -     | 1.30               | 1.35 | -     | 1.09          | 1.21 | -     |
|         | mGP, only B62       | -0.05              | 0.06  | -     | -0.11              | -0.21 | -     | -0.14         | -0.20 | -     | 0.99            | 0.99 | -     | 1.29               | 1.33 | -     | 1.09          | 1.20 | -     |
|         | Rosetta scaled      | 0.42               | 0.40  | -     | -0.63              | -0.57 | -     | 0.09          | -0.09 | -     | 0.75            | 0.76 | -     | 1.65               | 1.45 | -     | 1.10          | 1.20 | -     |
|         | Rosetta             |                    | 0.47  |       |                    | -0.65 |       |               | 0.14  |       |                 | 1.67 |       |                    | 2.45 |       |               | 1.94 |       |
| 4LYZ    | mGPFusion           | 0.61               | 0.27  | 0.34  | 0.95               | 0.57  | 0.66  | 0.65          | 0.34  | 0.35  | 1.47            | 1.65 | 3.46  | 0.69               | 1.02 | 2.35  | 1.35          | 1.45 | 3.26  |
|         | mGPFusion, only B62 | 0.59               | 0.30  | 0.34  | 0.96               | 0.64  | 0.64  | 0.64          | 0.40  | 0.35  | 1.48            | 1.64 | 3.50  | 0.57               | 0.93 | 2.17  | 1.35          | 1.43 | 3.27  |
|         | mGP                 | 0.65               | -0.05 | -     | 0.96               | 0.46  | -     | 0.68          | 0.18  | -     | 1.38            | 1.78 | -     | 0.36               | 1.15 | -     | 1.24          | 1.59 | -     |
|         | mGP, only B62       | 0.62               | 0.27  | -     | 0.97               | 0.47  | -     | 0.64          | 0.21  | -     | 1.48            | 1.84 | -     | 0.31               | 1.16 | -     | 1.33          | 1.63 | -     |
|         | Rosetta scaled      | 0.29               | 0.28  | -     | 0.70               | 0.68  | -     | 0.33          | 0.34  | -     | 2.26            | 2.22 | -     | 1.22               | 1.29 | -     | 2.09          | 1.94 | -     |
|         | Rosetta             |                    | 0.33  |       |                    | 0.71  |       |               | 0.35  |       |                 | 3.61 |       |                    | 2.05 |       |               | 3.35 |       |
| 1RTB    | mGPFusion           | 0.92               | 0.81  | 0.70  | 0.78               | 0.74  | 0.75  | 0.86          | 0.71  | 0.68  | 1.25            | 1.73 | 2.44  | 2.22               | 2.24 | 1.71  | 1.40          | 1.85 | 2.36  |
|         | mGPFusion, only B62 | 0.89               | 0.79  | 0.70  | 0.73               | 0.73  | 0.67  | 0.86          | 0.75  | 0.69  | 1.48            | 1.85 | 2.44  | 1.82               | 1.81 | 1.90  | 1.53          | 1.84 | 2.38  |
|         | mGP                 | 0.92               | 0.69  | -     | 0.78               | 0.30  | -     | 0.76          | 0.18  | -     | 1.26            | 2.62 | -     | 3.44               | 3.39 | -     | 1.69          | 2.81 | -     |
|         | mGP, only B62       | 0.91               | 0.58  | -     | 0.48               | 0.58  | -     | 0.85          | 0.16  | -     | 1.67            | 3.37 | -     | 2.12               | 2.11 | -     | 1.73          | 3.14 | -     |
|         | Rosetta scaled      | 0.65               | 0.61  | -     | 0.73               | 0.73  | -     | 0.67          | 0.65  | -     | 1.99            | 2.08 | -     | 1.63               | 1.62 | -     | 1.95          | 1.99 | -     |
|         | Rosetta             |                    | 0.69  |       |                    | 0.73  |       |               | 0.70  |       |                 | 2.45 |       |                    | 2.01 |       |               | 2.40 |       |
| 1PIN*   | mGPFusion           | 0.64               | 0.49  | 0.53  | -                  | -     | -     | 0.64          | 0.49  | 0.53  | 0.47            | 0.53 | 1.00  | -                  | -    | -     | 0.47          | 0.53 | 1.00  |
|         | mGPFusion, only B62 | 0.59               | 0.51  | 0.53  | -                  | -     | -     | 0.59          | 0.51  | 0.53  | 0.50            | 0.54 | 1.00  | -                  | -    | -     | 0.50          | 0.54 | 1.00  |
|         | mGP                 | 0.61               | 0.40  | -     | -                  | -     | -     | 0.61          | 0.40  | -     | 0.49            | 0.58 | -     | -                  | -    | -     | 0.49          | 0.58 | -     |
|         | mGP, only B62       | 0.23               | 0.25  | -     | -                  | -     | -     | 0.23          | 0.25  | -     | 0.70            | 0.76 | -     | -                  | -    | -     | 0.70          | 0.76 | -     |
|         | Rosetta scaled      | 0.50               | 0.49  | -     | -                  | -     | -     | 0.50          | 0.49  | -     | 0.54            | 0.54 | -     | -                  | -    | -     | 0.54          | 0.54 | -     |
|         | Rosetta             |                    | 0.53  |       |                    | -     |       |               | 0.53  |       |                 | 0.99 |       |                    | -    |       |               | 0.99 |       |
| 1RN1    | mGPFusion           | 0.83               | 0.62  | 0.67  | 0.97               | 0.38  | -0.19 | 0.84          | 0.55  | 0.65  | 0.85            | 1.17 | 1.18  | 0.19               | 0.91 | 1.28  | 0.81          | 1.13 | 1.19  |
|         | mGPFusion, only B62 | 0.81               | 0.58  | 0.67  | 0.96               | 0.42  | -0.01 | 0.81          | 0.51  | 0.66  | 0.92            | 1.24 | 1.18  | 0.35               | 0.77 | 0.99  | 0.88          | 1.17 | 1.17  |
|         | mGP                 | 0.74               | 0.04  | -     | 0.99               | 0.62  | -     | 0.75          | 0.04  | -     | 1.06            | 1.75 | -     | 0.12               | 0.58 | -     | 1.01          | 1.61 | -     |
|         | mGP, only B62       | 0.77               | 0.27  | -     | 0.99               | 0.50  | -     | 0.78          | 0.10  | -     | 1.00            | 1.78 | -     | 0.13               | 0.82 | -     | 0.95          | 1.66 | -     |
|         | Rosetta scaled      | 0.63               | 0.57  | -     | 0.21               | 0.18  | -     | 0.61          | 0.54  | -     | 1.12            | 1.18 | -     | 0.98               | 1.01 | -     | 1.10          | 1.15 | -     |
|         | Rosetta             |                    | 0.67  |       |                    | 0.20  |       |               | 0.65  |       |                 | 1.20 |       |                    | 1.24 |       |               | 1.20 |       |
| 1RGG    | mGPFusion           | 0.69               | 0.51  | 0.61  | 0.42               | 0.66  | 0.96  | 0.52          | 0.52  | 0.73  | 1.41            | 1.70 | 1.56  | 5.67               | 5.95 | 4.20  | 2.65          | 4.75 | 2.23  |
|         | mGPFusion, only B62 | 0.69               | 0.52  | 0.61  | 0.22               | 0.30  | 0.93  | 0.40          | 0.33  | 0.55  | 1.42            | 1.68 | 1.56  | 6.21               | 6.14 | 5.43  | 2.85          | 4.90 | 2.64  |
|         | mGP                 | 0.71               | -0.61 | -     | 0.19               | 0.08  | -     | 0.35          | -0.05 | -     | 1.38            | 2.24 | -     | 6.81               | 6.88 | -     | 3.05          | 5.53 | -     |
|         | mGP, only B62       | 0.72               | -0.53 | -     | 0.08               | 0.07  | -     | 0.30          | -0.00 | -     | 1.53            | 2.01 | -     | 6.89               | 6.77 | -     | 3.14          | 5.41 | -     |
|         | Rosetta scaled      | 0.58               | 0.43  | -     | 0.78               | 0.86  | -     | 0.42          | 0.51  | -     | 1.57            | 1.76 | -     | 7.08               | 7.24 | -     | 3.23          | 5.74 | -     |
|         | Rosetta             |                    | 0.60  |       |                    | 0.77  |       |               | 0.39  |       |                 | 1.58 |       |                    | 6.90 |       |               | 3.16 |       |
| 1BPI    | mGPFusion           | 0.68               | 0.67  | 0.51  | -0.39              | -0.02 | -0.06 | 0.85          | -0.13 | 0.12  | 1.27            | 1.28 | 1.83  | 2.17               | 7.66 | 7.69  | 1.40          | 3.64 | 3.11  |
|         | mGPFusion, only B62 | 0.64               | 0.65  | 0.51  | -0.53              | 0.33  | 0.21  | 0.81          | 0.05  | 0.12  | 1.32            | 1.32 | 1.82  | 2.86               | 7.07 | 7.70  | 1.57          | 3.40 | 3.11  |
|         | mGP                 | 0.56               | 0.57  | -     | -0.29              | 0.52  | -     | 0.83          | 0.84  | -     | 1.43            | 1.43 | -     | 1.72               | 4.19 | -     | 1.47          | 2.28 | -     |
|         | mGP, only B62       | 0.43               | 0.58  | -     | -0.55              | 0.46  | -     | 0.76          | 0.70  | -     | 1.68            | 1.73 | -     | 2.52               | 6.58 | -     | 1.80          | 3.34 | -     |
|         | Rosetta scaled      | 0.52               | 0.52  | -     | 0.23               | -0.00 | -     | 0.08          | -0.08 | -     | 1.52            | 1.52 | -     | 7.43               | 7.43 | -     | 2.88          | 3.62 | -     |
|         | Rosetta             |                    | 0.51  |       |                    | -0.00 |       |               | 0.13  |       |                 | 1.80 |       |                    | 7.68 |       |               | 3.09 |       |
| total   | mGPFusion           | 0.81               | 0.70  | 0.56  | 0.88               | 0.61  | 0.49  | 0.83          | 0.64  | 0.52  | 1.07            | 1.26 | 1.61  | 1.33               | 2.45 | 2.53  | 1.13          | 1.87 | 1.84  |
|         | mGPFusion, only B62 | 0.79               | 0.69  | 0.56  | 0.86               | 0.64  | 0.50  | 0.82          | 0.66  | 0.52  | 1.11            | 1.30 | 1.62  | 1.43               | 2.40 | 2.50  | 1.18          | 1.85 | 1.84  |
|         | mGP                 | 0.81               | 0.51  | -     | 0.86               | 0.52  | -     | 0.83          | 0.50  | -     | 1.04            | 1.54 | -     | 1.44               | 2.65 | -     | 1.14          | 2.09 | -     |
|         | mGP, only B62       | 0.76               | 0.34  | -     | 0.86               | 0.55  | -     | 0.80          | 0.49  | -     | 1.26            | 1.95 | -     | 1.45               | 2.56 | -     | 1.30          | 2.23 | -     |
|         | Rosetta scaled      | 0.65               | 0.63  | -     | 0.51               | 0.39  | -     | 0.60          | 0.48  | -     | 1.35            | 1.38 | -     | 2.49               | 2.99 | -     | 1.66          | 2.22 | -     |
|         | Rosetta             |                    | 0.55  |       |                    | 0.40  |       |               | 0.49  |       |                 | 1.63 |       |                    | 2.74 |       |               | 1.92 |       |
|         | mCSM                |                    | 0.61  |       |                    | -     |       | -             |       |       | 1.40            |      |       | -                  |      |       | -             | -    |       |
|         | PoPMuSiC            |                    | 0.64  |       |                    | -     |       | -             |       |       | 1.37            |      |       | -                  |      |       | -             | -    |       |

\* dataset for 1PIN contained no multiple mutations.

Table S6: Comparison of different methods on the 15 protein dataset with respect to  $\rho$  and rmse after removing 10% of predictions with largest errors. Mutation, position, and protein are referred to as mut., pos., and prot., respectively. Off-the-shelf implementations of Rosetta, mCSM and PoPMuSiC are used directly without cross-validation.

| Results after 10 %<br>outlier removal                  | Correlation $\rho$     |      |       |                        |      |       |                        |      |       | rmse                   |      |       |                        |      |       |                        |      |       |
|--|------------------------|------|-------|------------------------|------|-------|------------------------|------|-------|------------------------|------|-------|------------------------|------|-------|------------------------|------|-------|
|  | Point mutations        |      |       | Multiple mutations     |      |       | All mutations          |      |       | Point mutations        |      |       | Multiple mutations     |      |       | All mutations          |      |       |
|  | cross-validation level |      |       | cross-validation level |      |       | cross-validation level |      |       | cross-validation level |      |       | cross-validation level |      |       | cross-validation level |      |       |
| Method   | mut.                   | pos. | prot. | mut.                   | pos. | prot. | mut.                   | pos. | prot. | mut.                   | pos. | prot. | mut.                   | pos. | prot. | mut.                   | pos. | prot. |
| mGPFusion  | 0.87                   | 0.77 | 0.75  | 0.97                   | 0.84 | 0.63  | 0.92                   | 0.81 | 0.71  | 0.69                   | 0.85 | 0.93  | 0.62                   | 1.55 | 1.90  | 0.67                   | 1.15 | 1.13  |
| mGPFusion, only B62                                    | 0.86                   | 0.73 | 0.75  | 0.96                   | 0.84 | 0.65  | 0.91                   | 0.80 | 0.72  | 0.69                   | 0.85 | 0.93  | 0.64                   | 1.43 | 1.85  | 0.67                   | 1.10 | 1.12  |
| mGP  | 0.89                   | 0.57 | -     | 0.98                   | 0.76 | -     | 0.93                   | 0.74 | -     | 0.65                   | 1.01 | -     | 0.49                   | 1.66 | -     | 0.62                   | 1.23 | -     |
| mGP, only B62  | 0.86                   | 0.27 | -     | 0.97                   | 0.76 | -     | 0.91                   | 0.66 | -     | 0.73                   | 1.30 | -     | 0.54                   | 1.48 | -     | 0.69                   | 1.37 | -     |
| Rosetta scaled   | 0.80                   | 0.78 | -     | 0.73                   | 0.71 | -     | 0.78                   | 0.75 | -     | 0.84                   | 0.87 | -     | 1.67                   | 1.99 | -     | 1.00                   | 1.31 | -     |
| Off-the-shelf implementations with no cross-validation |                        |      |       |                        |      |       |                        |      |       |                        |      |       |                        |      |       |                        |      |       |
| Rosetta  | 0.75                   |      |       | 0.67                   |      |       | 0.73                   |      |       | 0.94                   |      |       | 1.85                   |      |       | 1.11                   |      |       |
| mCSM   | 0.71                   |      |       | -                      |      |       | -                      |      |       | 0.89                   |      |       | -                      |      |       | -                      |      |       |
| PoPMuSiC   | 0.73                   |      |       | -                      |      |       | -                      |      |       | 0.86                   |      |       | -                      |      |       | -                      |      |       |

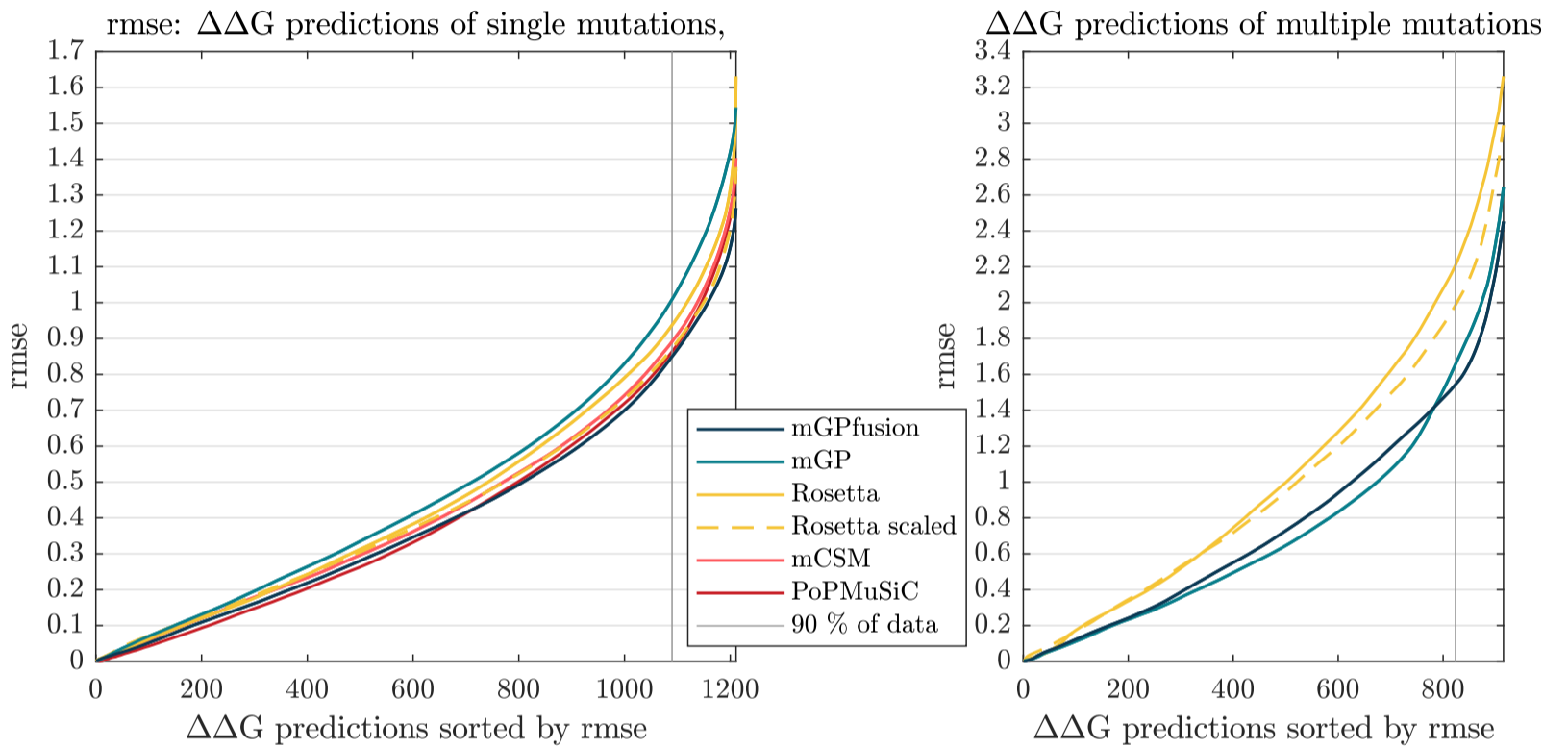


Figure S9: rmse with different amount of predictions, when predictions are sorted by the error. Position level cross-validation was used for mGPFusion, mGP and Rosetta scaled.

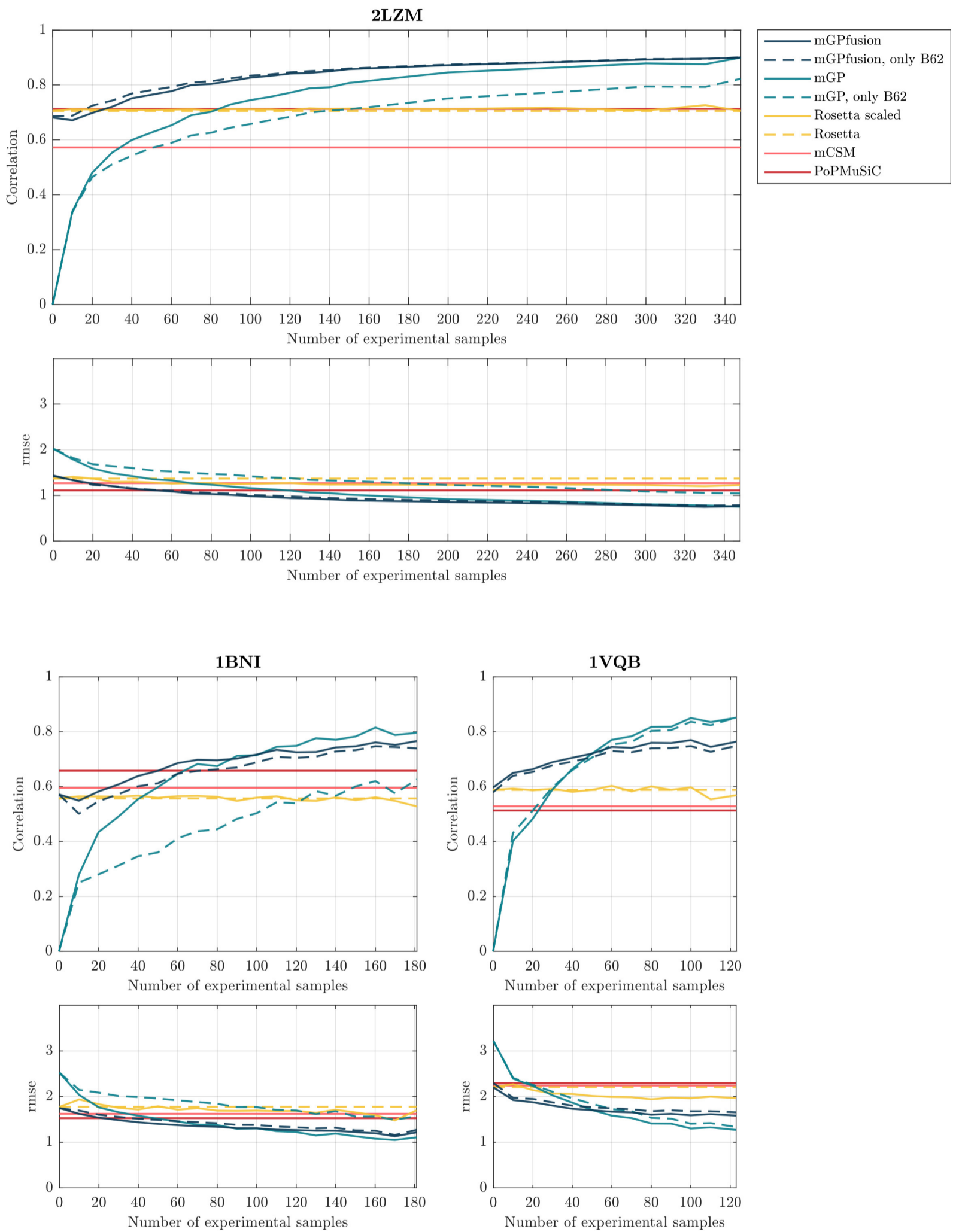


Figure S10: Learning curves.



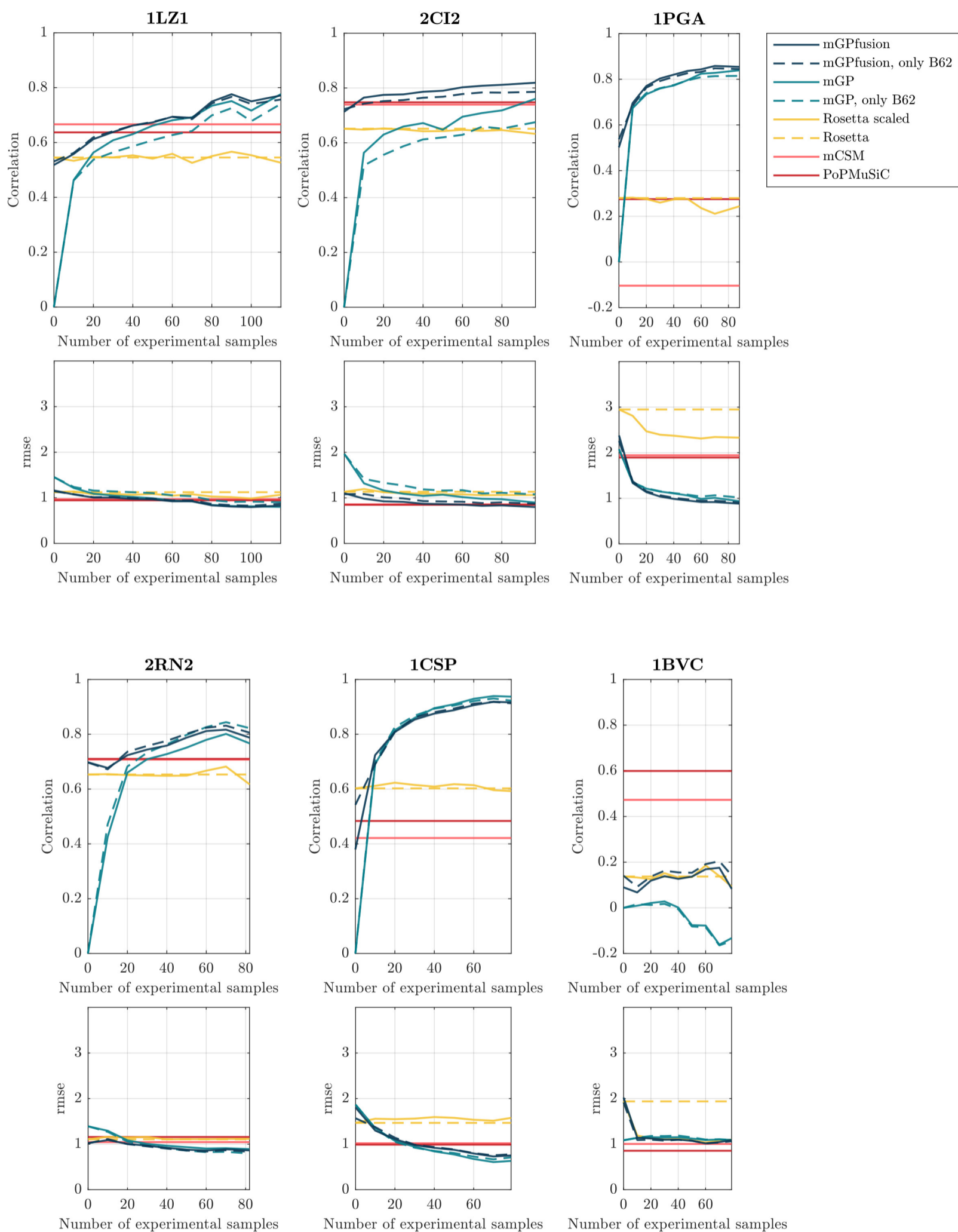


Figure S11: Learning curves.

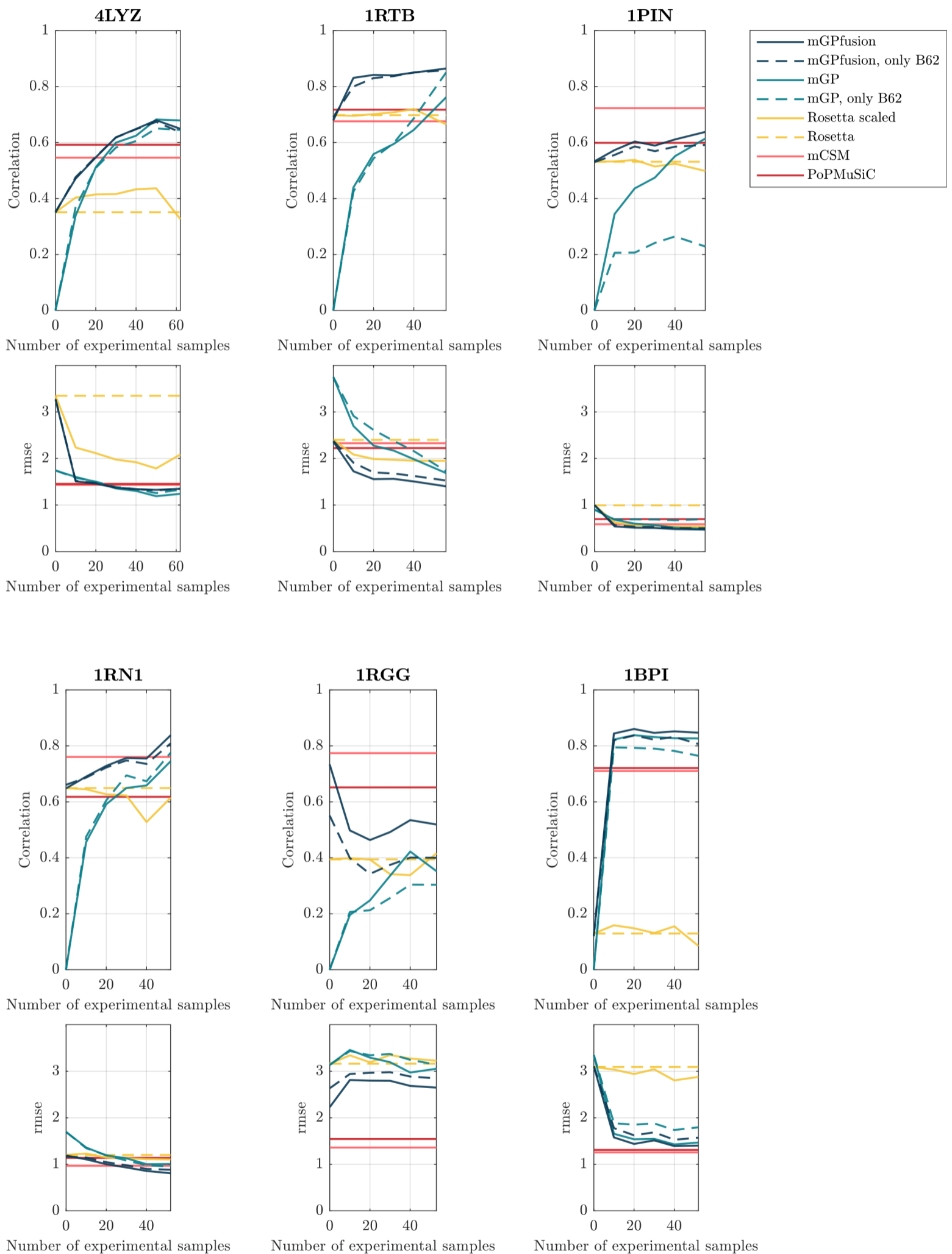


Figure S12: Learning curves.

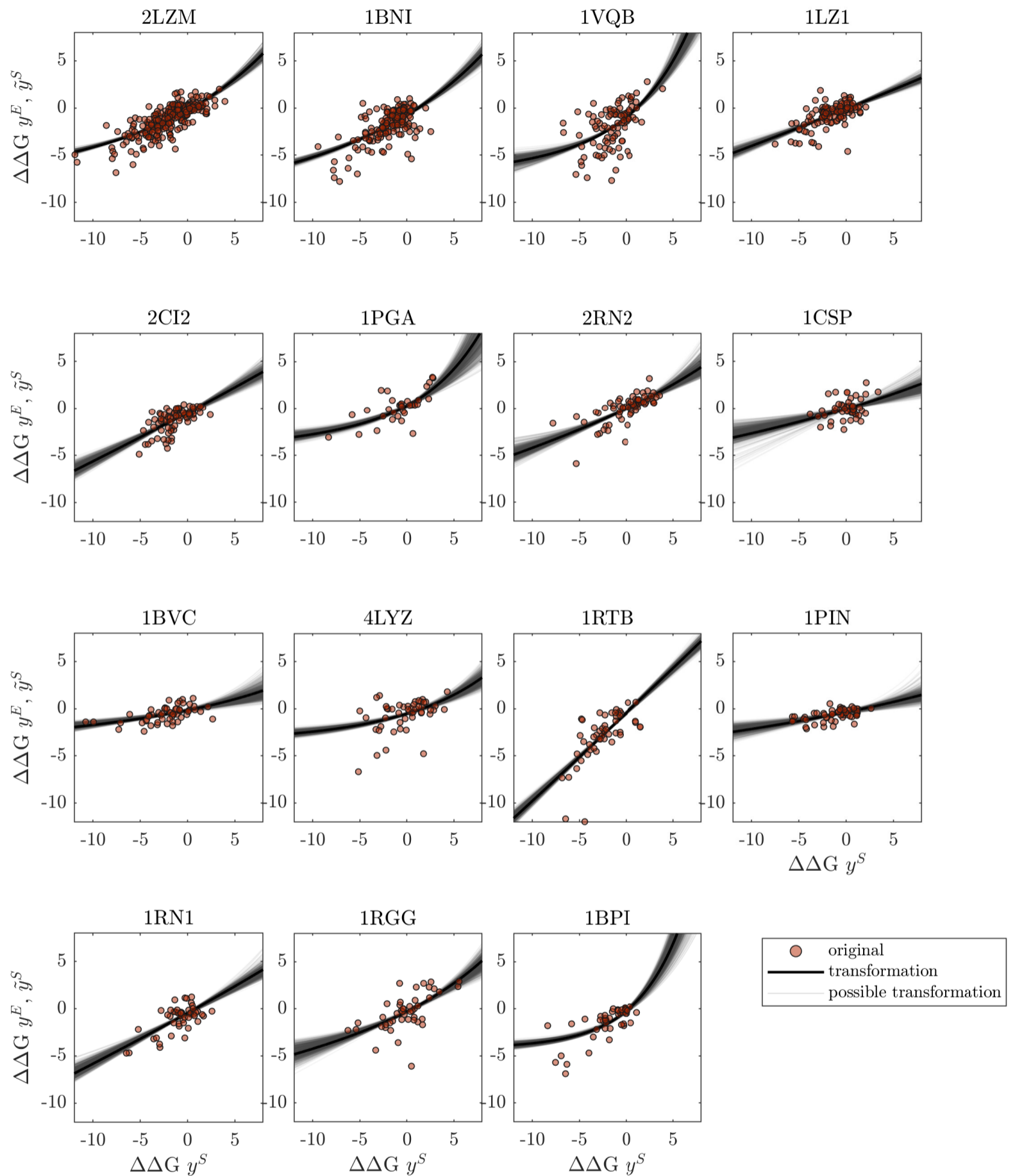


Figure S13: Transformations for all 15 proteins presented in Table 1. The red circles mark the simulated  $\Delta\Delta G$ -values  $y^S$  with respect to the experimental measured  $\Delta\Delta G$ -values  $y^E$ . Thin black lines show possible transformations for  $y^S$ , whereas the thick black line shows the selected transformation from  $y^S$  to  $\tilde{y}^S$ .