



Protein stability engineering insights revealed by domain-wide comprehensive mutagenesis

Alex Nisthal^{a,1,2}, Connie Y. Wang^b, Marie L. Ary^b, and Stephen L. Mayo^{a,c,1}

^aDivision of Biology and Biological Engineering, California Institute of Technology, Pasadena, CA 91125; ^bProtabit, LLC, Pasadena, CA 91106; and ^cDivision of Chemistry and Chemical Engineering, California Institute of Technology, Pasadena, CA 91125

Contributed by Stephen L. Mayo, July 3, 2019 (sent for review March 6, 2019; reviewed by Elizabeth M. Meiering and Timothy Whitehead)

The accurate prediction of protein stability upon sequence mutation is an important but unsolved challenge in protein engineering. Large mutational datasets are required to train computational predictors, but traditional methods for collecting stability data are either low-throughput or measure protein stability indirectly. Here, we develop an automated method to generate thermodynamic stability data for nearly every single mutant in a small 56-residue protein. Analysis reveals that most single mutants have a neutral effect on stability, mutational sensitivity is largely governed by residue burial, and unexpectedly, hydrophobics are the best tolerated amino acid type. Correlating the output of various stability-prediction algorithms against our data shows that nearly all perform better on boundary and surface positions than for those in the core and are better at predicting large-to-small mutations than small-to-large ones. We show that the most stable variants in the single-mutant landscape are better identified using combinations of 2 prediction algorithms and including more algorithms can provide diminishing returns. In most cases, poor *in silico* predictions were tied to compositional differences between the data being analyzed and the datasets used to train the algorithm. Finally, we find that strategies to extract stabilities from high-throughput fitness data such as deep mutational scanning are promising and that data produced by these methods may be applicable toward training future stability-prediction tools.

thermodynamic stability | mutagenesis | protein engineering | protein stability prediction | protein G

Thermodynamic stability is a fundamental property of proteins that significantly influences protein structure, function, expression, and solubility. Efforts to identify the molecular determinants of protein stability and to engineer improvements have thus been crucial in the development and optimization of a wide range of biotechnology products, including industrial-grade enzymes, antibodies, and other protein-based therapeutics and reagents (1–3). The ability to reliably predict the effect of mutations on protein stability would greatly facilitate engineering efforts, and much research has been devoted to developing computational tools for this purpose (4–10). Understanding how mutations affect stability can also shed light on various biological processes, including disease and drug resistance (11). Advances in genotyping and next-generation sequencing allow for the identification of significant numbers of missense mutations associated with human disease (12). Fast and accurate protein stability prediction could be used to quickly identify which of these mutations are likely to lead to disease phenotypes.

However, the accurate prediction of the impact of an amino acid substitution on protein stability remains an unsolved challenge in protein engineering. Correlation studies have shown that computational techniques can capture general trends, but fail to precisely predict the magnitude of mutational effects (13, 14). The success of these techniques is dependent on the quality of the input structure, conformational sampling, the free-energy function used to evaluate the mutant sequences, and importantly, the data used for training and testing (8, 13, 15). Traditionally, protein stability data are collected by generating and

purifying a small set of selected protein variants for characterization via calorimetry or spectroscopically measured chemical or thermal denaturation experiments. Values typically determined include the chemical or thermal denaturation midpoint (C_m or T_m , respectively), the free energy of unfolding (ΔG), and the change in ΔG relative to wild type (WT) ($\Delta\Delta G$). Although low-throughput, the widespread use of these methods has generated a wealth of protein stability data over time, which has shaped our current understanding of protein structure–function relationships (16–19). Much of this work has been aggregated in the ProTherm (20) database, commonly used as a training data resource. Until recently, ProTherm was the largest public source of thermodynamic protein stability data, containing over 25,000 entries from 1,902 scientific articles. The database has been critical to the development of a variety of computational tools, from knowledge-based potentials exclusively trained on experimental data (6) to physics-based potentials with atomic resolution (7) and everything in between. Unfortunately, the ProTherm website is no longer being supported. The ProTherm data are still available, however, in ProtBank (21), a recently developed online database for protein engineering data (<https://www.protabank.org/>).

Although training and validation datasets from ProTherm have been widely used, ProTherm data suffer from 3 flaws: 1) experimental conditions vary widely among entries, requiring manual filtering to obtain comparable data, which results in smaller

Significance

Using liquid-handling automation, we constructed and measured the thermodynamic stability of almost every single mutant of protein G (G β 1), a small domain. This self-consistent dataset is the largest of its kind and offers unique opportunities on 2 fronts: 1) insight into protein domain properties such as positional sensitivity and incorporated amino acid tolerance, and 2) service as a validation set for future efforts in protein stability prediction. As G β 1 is a model system for protein folding and design, and its single-mutant landscape has been measured by deep mutational scanning, we expect our dataset to serve as a reference for studies aimed at extracting stability information from fitness data or developing high-throughput stability assays.

Author contributions: A.N. and S.L.M. designed research; A.N. performed research; A.N. and C.Y.W. contributed new reagents/analytic tools; A.N., C.Y.W., and M.L.A. analyzed data; and A.N., C.Y.W., M.L.A., and S.L.M. wrote the paper.

Reviewers: E.M.M., University of Waterloo; and T.W., University of Colorado Boulder.

The authors declare no conflict of interest.

Published under the PNAS license.

Data deposition: The data reported in this paper have been deposited in ProtBank, <https://www.protabank.org/> (ID no. gwo52haU3).

¹To whom correspondence may be addressed. Email: nisthal@caltech.edu or steve@mayo.caltech.edu.

²Present address: Protein Engineering, Xencor, Inc., Monrovia, CA 91016.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1903888116/-DCSupplemental.

Published online August 1, 2019.

datasets; 2) little information is included on unfolded or alternatively folded sequences, precluding training on this type of mutational data; and 3) results from alanine (Ala) scanning mutagenesis are overrepresented, biasing the dataset toward large-to-small mutations. Thus, training or testing on ProTherm data may mask deficiencies in computational algorithms or result in predictions that are biased toward particular features of the dataset. As many of the stability-prediction tools available today rely on experimental data from ProTherm, it is perhaps not surprising that none is very accurate and all perform about the same (8, 13, 14, 22).

Comprehensive mutagenesis studies, with stabilities measured under fixed experimental conditions, could provide better training data. The low-throughput nature of traditional methods, however, makes the collection of stability data for large numbers of protein variants unfeasible. Several strategies have been devised to improve this process, including the use of genetic repressor systems (23), plate-based fluorescence assays (24, 25), differential scanning fluorimetry (26), and, more recently, yeast-displayed proteolysis (27). Unfortunately, these approaches generally make compromises by either 1) tying an easy-to-measure but indirect protein stability readout to large variant libraries, or 2) addressing the throughput of stability determination, but not the laborious nature of variant generation and purification.

Here, we develop an automated method that addresses both of these issues and apply it to obtain thermodynamic stability data from the comprehensive mutagenesis of an entire protein domain—the 56-residue β 1 domain of streptococcal protein G (G β 1). G β 1 was chosen for its small size, high amount of secondary structure, and well-behaved WT sequence. Drawing both inspiration and methodology from structural genomics, we couple automated molecular biology procedures with a high-throughput plate-based stability determination method, resulting in a 20-fold increase in throughput over traditional benchtop methods. We applied our experimental pipeline to G β 1 to produce a dataset that maintains constant experimental conditions, includes data on nonfolded sequences, and features an unbiased mutational distribution over 935 unique variants covering nearly every single mutant of G β 1. Data in hand, we examine positional

sensitivity and amino acid tolerance, and evaluate several protein stability-prediction algorithms and engineering strategies. Finally, we compare our dataset against one derived by deep mutational scanning (DMS), a technique that can generate large mutational datasets via functional selections and deep sequencing (28, 29), and explore whether stability data from DMS studies are applicable toward training future protein stability-prediction tools.

Results and Discussion

Automated Site-Directed Mutagenesis and Stability Determination Pipeline Increases Throughput 20-Fold. Using laboratory automation, we constructed, expressed, and purified nearly every single mutant of G β 1 fused to an N-terminal His6 purification tag. The automated pipeline is illustrated in Fig. 1A. Each variant was constructed explicitly instead of by saturation mutagenesis so that mutants not found in the first pass could be more easily recovered. Variants were constructed using a megaprimer method that requires only 1 mutagenic oligonucleotide, thereby halving oligonucleotide costs. After expression and purification, the thermodynamic stabilities of the generated variants were then determined using an improved version of our previously described plate-based chemical denaturation assay (24) (Fig. 1B). Enhancements include adaptation to automated liquid handling for increased speed and precision, and doubling the number of data points collected per curve to improve accuracy. Although the intent was to collect data on 19 amino acids at 56 positions for a total of 1,064 variants, a trade-off was made in which mutations at the buried tryptophan (Trp) at position 43 (W43) were excluded to preserve the integrity of the Trp-based fluorescence assay. Also, mutations incorporating cysteine (Cys) or Trp were omitted to avoid oligomerization by disulfide formation and potential interference with W43, respectively. Thus, mutations to 17 of 19 possible amino acids were made at 55 of 56 positions, for a total of 935 single mutants.

Each step of the workflow was developed as an independent module, allowing for optimization outside the full experimental pipeline. Modularization also permits flexible scheduling and parallelization, allowing modules to run multiple times per day.

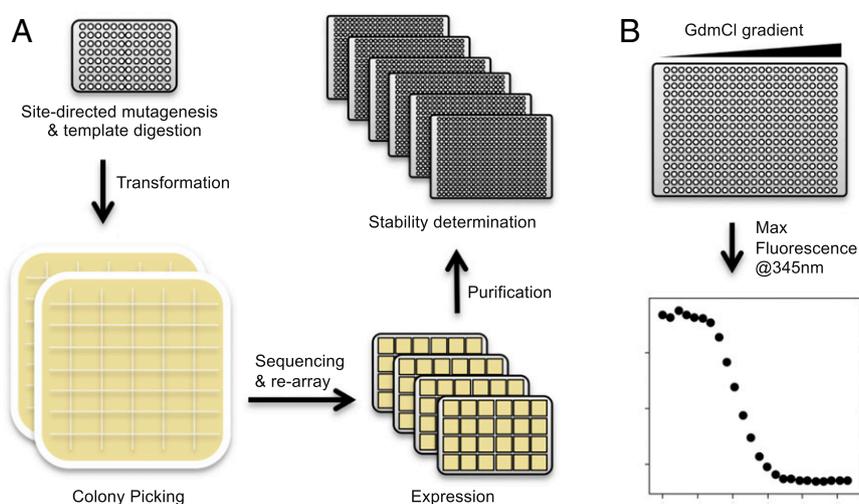


Fig. 1. Automated site-directed mutagenesis and stability determination pipeline. (A) Modular protocols enable the rapid construction, sequence verification, and stability determination of single-mutant variants. For illustrative purposes, each step shows the number of plates required for 96 individual reactions. Oligonucleotides direct mutagenesis reactions on 96-well PCR plates followed by bacterial transformation and plating onto 48-well agar trays. Individual colonies are picked, cultured, and rearrayed after successful sequence validation. Confirmed variants are expressed in 24-well culture blocks and NiNTA purified. Each reaction was tracked via a database throughout the pipeline, allowing for method optimization. (B) Thermodynamic stability was determined by measuring Trp fluorescence in response to a 24-point GdmCl gradient. Each row of a 384-well plate is one protein stability curve from which the concentration of denaturant at the midpoint of the unfolding transition (C_m) is directly measured. After estimating the slope of the curve (m value), the change in the free energy of unfolding ($\Delta\Delta G$) of each variant relative to WT is calculated by taking the difference between the WT and mutant C_m values and multiplying by their mean m value (*SI Appendix, Methods*).

For comparison, 8 d is a reasonable estimate for traditional procedures to construct, verify, express, purify, and measure the thermodynamic stability of 8 single mutants. Extrapolating to 935 variants (the number in this study), traditional procedures would take 935 d, or 2.5 y. In contrast, our platform can generate data on 935 variants in 5 to 6 wk, a speedup of at least 20-fold.

Stability Determination of G β 1 Single Mutants. We measured the Trp fluorescence of each variant in response to a 24-point guanidinium chloride (GdmCl) gradient 2 to 6 times as biological replicates. Each denaturation experiment generates an unfolding curve (Fig. 1*B*) from which we determined the concentration of denaturant at the midpoint of the unfolding transition (C_m), the $\Delta G(H_2O)$, and the slope (m value) using the linear extrapolation method (30, 31). While $\Delta\Delta G$ can be calculated in multiple ways (*SI Appendix, Methods*), we prefer the method that takes the difference between the mutant and WT C_m values and multiplies it by their mean m value (\bar{m}) (32) as shown in the following equation:

$$\Delta\Delta G = \bar{m} \times (C_m \text{ mutant} - C_m \text{ WT}),$$

where the individual mutant and WT m values are obtained by the linear extrapolation method. This method was selected because it is more precise yet in excellent agreement with the method that takes the difference between $\Delta G(H_2O)$ estimates for WT G β 1 and each mutant protein (*SI Appendix, Fig. S1*). Using this equation, stabilizing mutations have positive $\Delta\Delta G$ values, and destabilizing mutations have negative values. Of the 935 variants analyzed, 105 failed the assumptions of the linear extrapolation method (reversibility of folding/unfolding and 2-state behavior) due to poor stability, presence of a folding intermediate, or no expression (*SI Appendix, Fig. S2*). The 830 variants that passed these criteria are referred to as the quantitative dataset, and the remaining 105 are referred to as the qualitative dataset. We observed minimal effects of the N-terminal His6 purification tag on G β 1 stability, reflected by the strong correlation ($r = 0.95$; $P < 0.001$) between literature values of untagged G β 1 variants and corresponding values from our dataset (*SI Appendix, Fig. S3*). The single-mutant stabilities ($\Delta\Delta G$ s) for the entire dataset (33) are shown as a heat map in Fig. 2.

Stability Distribution of G β 1 Single Mutants Is Primarily Neutral. The $\Delta\Delta G$ distribution of G β 1 single mutants is primarily neutral ($\Delta\Delta G$ of 0 ± 1 kcal/mol) with a long tail of destabilizing variants (Fig. 3*A*). The median of the quantitative dataset is 0.05 kcal/mol with an interquartile range of 1.0 kcal/mol (Fig. 3*C*), and the fraction of positive, neutral, and negative mutations is 3%, 68%, and 29%, respectively. If we assume the qualitative data contains only negative mutations, then our complete dataset shifts the fractions to 3%, 60%, and 37%, respectively. Summing the positive and neutral mutations, almost two-thirds of the tested single mutants (63%) have at worst no effect on G β 1 stability. The fraction of destabilizing mutations (37%) is on the low end compared with an experimental dataset of 1,285 mutants from ProTherm, which shows that $\sim 50\%$ of single mutants are destabilizing ($\Delta\Delta G < 1$ kcal/mol) (34, 35). The destabilizing fraction we obtained for G β 1 would likely increase, however, upon making mutations to W43 and including Trp and Cys scanning variants as these residues are generally difficult to substitute in or out (36). Also, the G β 1 domain itself may skew mutational outcomes as its small size results in a large surface-to-buried area ratio. This ratio likely contributes to fewer destabilizing mutations than larger proteins with larger cores, assuming that most core mutations are destabilizing (17, 23, 37, 38).

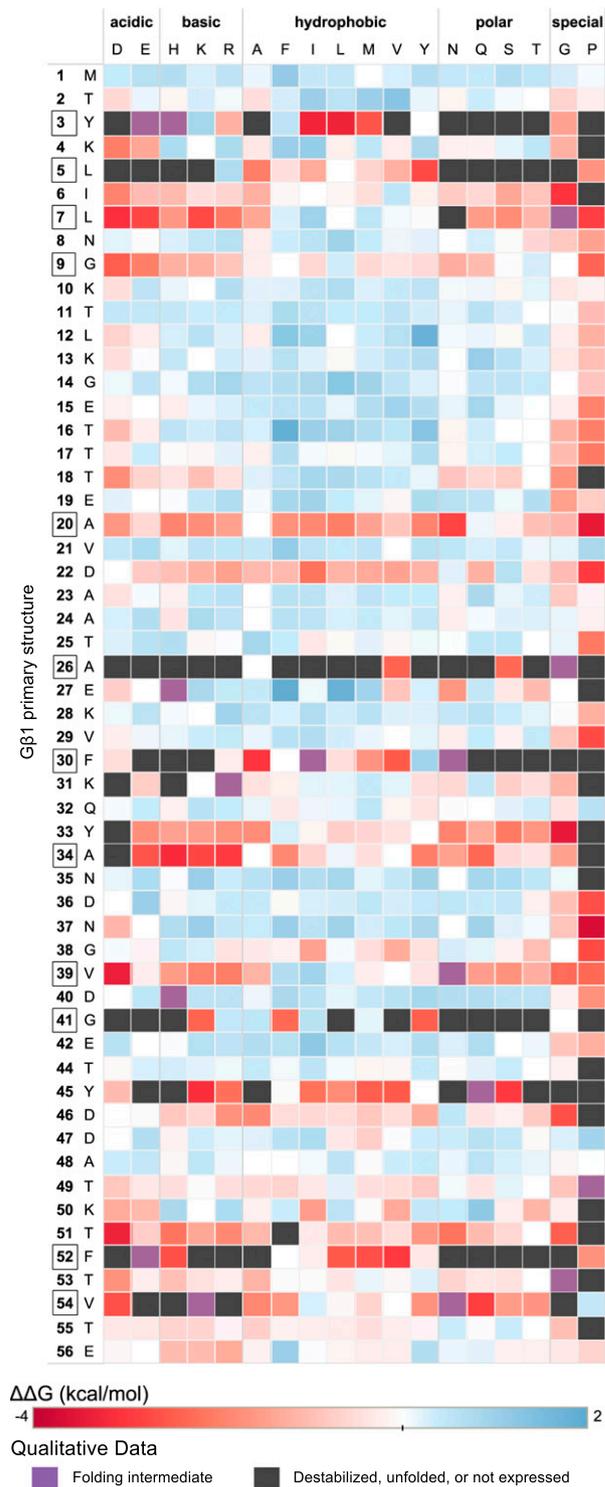


Fig. 2. Single-mutant thermodynamic stability landscape of G β 1. The vertical axis of the mutational matrix depicts the primary structure of G β 1, with the position and WT amino acid as indicated. Core positions as determined by RESCLASS (4) are boxed. The horizontal axis depicts mutant amino acids examined in the study, grouped by amino acid type. Variants from the quantitative dataset are colored by their determined $\Delta\Delta G$ value where red is destabilizing, blue is stabilizing, and white is neutral. Self-identity mutations such as M1M have $\Delta\Delta G = 0$ and thus are colored white. Variants from the qualitative dataset are colored based on whether a folding intermediate was detected (purple) or whether the mutant did not express, was unfolded, or was too destabilized to collect a quantitative measurement (black) (*SI Appendix, Fig. S2*).

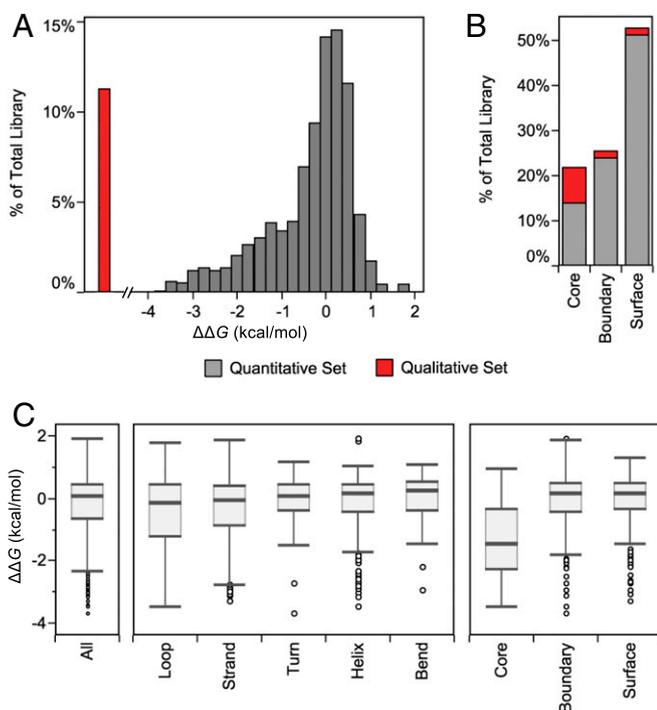


Fig. 3. Stability distribution of G β 1 single mutants. (A and B) The 935-member dataset is split between variants with quantitative data (gray) and those with only qualitative data (red) due to poor stability or misfolding. (A) The $\Delta\Delta G$ distribution is split into 0.25 kcal/mol bins. Variants belonging to the qualitative dataset are shown to the *Left* of the distribution, indicating values of $\Delta\Delta G$ less than -4 kcal/mol. (B) Variants are binned into core, boundary, or surface using RESCLASS (4). (C) Box-and-whisker plots of the quantitative dataset describe the median, the quartile cutoffs, and the outlier cutoffs of the $\Delta\Delta G$ distribution for all of the residues (All), binned into secondary structure classifications as defined by DSSP (63), or binned by RESCLASS. Outliers are shown as unfilled circles and are defined as points that are $1.5\times$ interquartile range above or below the third quartile or first quartile, respectively.

Positional Sensitivity Is Governed by Residue Burial. The heat map in Fig. 2, which is organized by primary structure, allows for a granular look at the distribution of mutational stability. We observe 2 clear trends: 1) the mutational sensitivity ($\Delta\Delta G$) of the domain is largely determined by the position of the mutation, not the amino acid identity, unless 2) the mutations are to glycine (Gly) or proline (Pro), for which most mutations are deleterious. Positions 3, 5, 26, 30, 41, 45, 52, and 54 are particularly sensitive to mutation. If we map the positional sensitivity (median $\Delta\Delta G$ at each position) onto the G β 1 structure (Fig. 4), we see that residues in the interior of the protein are more susceptible to destabilization. This is also observed when analyzing the distribution by tertiary structure, but not by secondary structure (Fig. 3C). That is, classifying residues into core, boundary, or surface with the RESCLASS algorithm (4) shows that the median $\Delta\Delta G$ for core residues is ~ 1.5 kcal/mol lower than that of the rest of the protein. In addition, the qualitative dataset, which contains mutants whose stabilities are difficult to measure or are fully unfolded, includes 5-fold more core variants than boundary or surface positions, adding further support to this observation (Fig. 3B). Although this relationship has been observed with other datasets using a variety of proxies for protein stability (17, 23, 37, 38), this study provides a comprehensive analysis at the whole domain level with direct thermodynamic stability measurements.

As seen in Fig. 2, however, not all core positions behave the same, as some are more sensitive to mutation than others. For engineering purposes, it would be useful to identify specific

protein attributes that could serve as quantitative predictors of positional sensitivity. We therefore performed linear regression with 10-fold cross-validation on a large number of attributes that might impact protein stability. Attributes tested included measures of residue burial, secondary structure type/propensity, structural flexibility, and the change upon mutation of residue descriptors such as hydrophobicity, volume, and charge. The best individual predictors were measures of residue burial: depth of the C β atom (39, 40) and occluded surface packing (OSP) (41, 42), with Pearson correlation coefficients (r) of 0.82 and 0.76, respectively (both P s < 0.001). This demonstrates that not all core positions are created equal, and that there is a direct relationship between how buried a position is and its sensitivity to mutation. Flexibility descriptors such as root-mean-squared fluctuations (RMSFs) (from molecular dynamics simulations) or secondary structure descriptors such as α -helix propensity performed less well ($r = 0.42$ and 0.06 , respectively). We repeated these analyses with sequence entropy (43) as an alternative metric of positional sensitivity, and the conclusions remain the same; C β depth and OSP were the 2 best predictors, with $r = 0.81$ and 0.78 , respectively (both P s < 0.001). Combinations of attributes were also tested, but these did not substantially improve predictability. Given the strong correlation between positional sensitivity and residue burial indicators like OSP and C β depth, calculation of these measures should be among the first tools employed when evaluating positions for substitution, provided structural information is available.

Hydrophobics Are the Best Tolerated Amino Acid Type. A common practice in protein redesign and optimization is to restrict core residues to nonpolar amino acids and only allow polar amino acids at the surface. We tested the validity of this strategy with our quantitative dataset by calculating median $\Delta\Delta G$ by incorporating amino acid and ranking the amino acids from worst tolerated to best tolerated across the entire domain (Fig. 5A). In general, the 2 worst amino acids for incorporation are Pro and Gly, which is not surprising given their vastly different Ramachandran preferences compared with all other amino acids. Beyond secondary structure-breaking amino acids, the third worst tolerated amino acid, interestingly, is aspartic acid (Asp), which may be rationalized by the fact that it is very hydrophilic (44) and has one of the highest charge densities among the amino acids (45). Unexpectedly, hydrophobic amino acids, particularly isoleucine (Ile) and phenylalanine (Phe), are among the best tolerated residues across all G β 1 positions. Even among surface positions, which make up over 50% of the dataset, Ile is the most favored individual residue, and hydrophobic amino acids as a whole are favored equally or better than the other amino acid types (Fig. 5B). The preference for hydrophobic

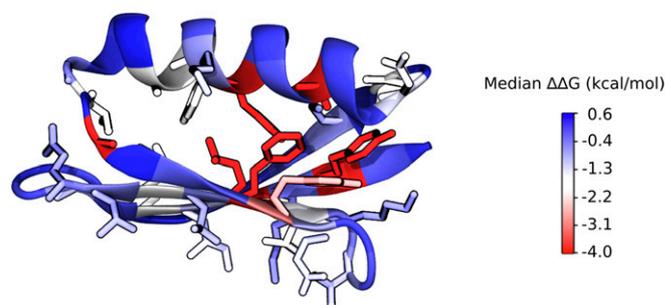


Fig. 4. Positional sensitivity (median $\Delta\Delta G$ at each position) of G β 1. G β 1 X-ray crystal structure (PDB ID: 1PGA) is colored by the positional sensitivity at each position. Side-chain atoms are shown for residues with a positional sensitivity score less than zero (destabilized).

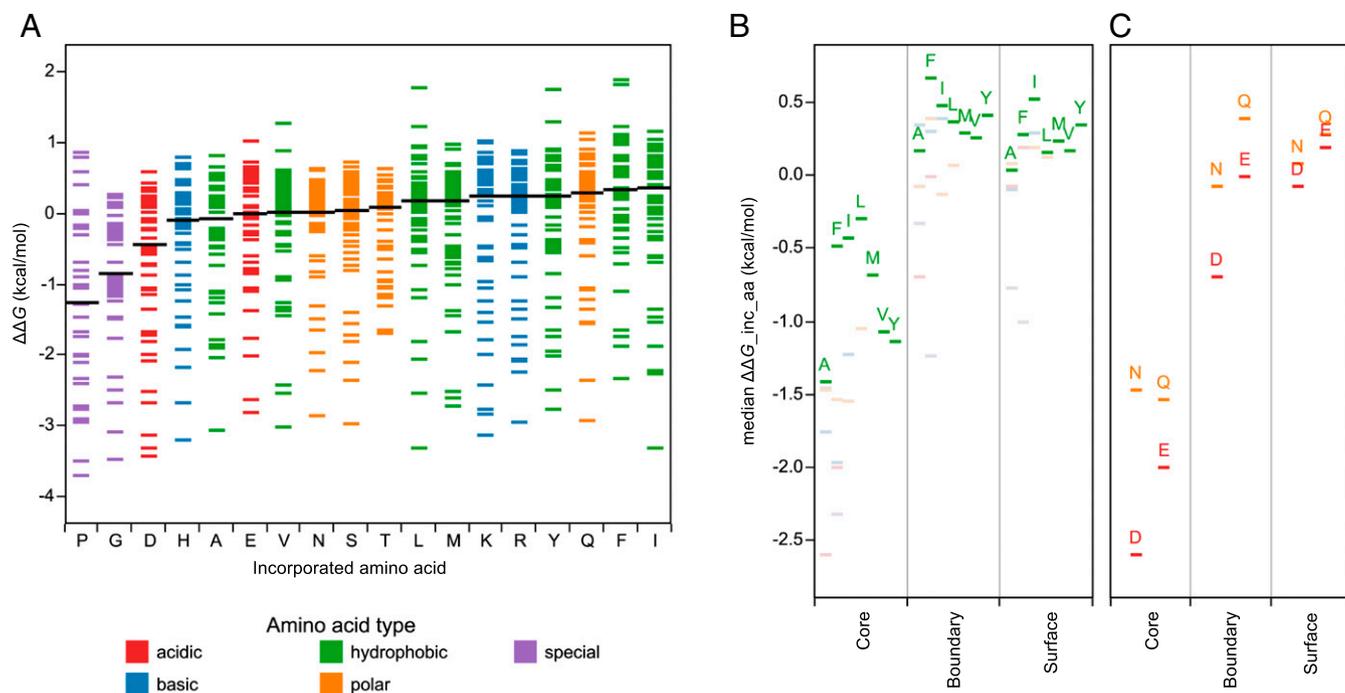


Fig. 5. $\Delta\Delta G$ distribution by incorporated amino acid. Amino acids are colored by physiochemical type. (A) Individual variants are shown as Gantt lines and distributed by the incorporated amino acid. The amino acid bins are ordered from *Left to Right* by the median $\Delta\Delta G$ of each distribution (black lines). (B) Median $\Delta\Delta G$ values of the incorporated amino acid distribution grouped by RESCLASS (4). For clarity, only hydrophobic amino acids are labeled. (C) Median $\Delta\Delta G$ values of chemically similar pairs (D/E and N/Q), grouped by RESCLASS.

amino acids extends to the chemically similar amino acid pairs, Asp/glutamic acid (Glu) and asparagine (Asn)/glutamine (Gln), where the pair member containing the extra methylene is better tolerated across the domain (Fig. 5A) and in almost every RESCLASS environment (Fig. 5C). To determine whether this observation is unique to G β 1, we performed domain-wide *in silico* stability predictions (6, 46) on 5 compositionally diverse proteins, including G β 1 as a control (*SI Appendix, Table S1*). Remarkably, the calculations recapitulated our observations for G β 1 and produced similar results for the other proteins, even across different RESCLASS types (*SI Appendix, Fig. S4*).

Several other experimental studies have also found that hydrophobic amino acids are well tolerated on the surface (47–51). The investigators attributed these findings to unique amino acid properties or structural contexts that enable these nonpolar mutations to stabilize the mutation site. However, our results suggest that non-position-specific increases in nonpolar surface area and volume are well tolerated, and the more the better. Larger hydrophobic amino acids like Ile, Phe, and tyrosine (Tyr) are consistently ranked as the best tolerated, and smaller

hydrophobics like Ala or valine (Val) do much worse across all 3 residue classes, including surface residues (Fig. 5B). Although multiple nonpolar mutations to the surface are still likely deleterious to protein stability and solubility (48), single mutations to hydrophobic amino acids should not be categorically excluded for stability optimization.

Benchmarking Protein Stability-Prediction Algorithms. We evaluated the ability of 3 stability-prediction algorithms, PoPMuSiC (46), FoldX (7), and Rosetta (5, 8), to recapitulate the 830 $\Delta\Delta G$ values in our quantitative dataset. While there are many protein stability predictors, we focus on these 3 algorithms because they have shown success in engineering point mutations with experimentally verified improved stability (48), and they represent a range of diverse energy function methodologies. To better understand the effect of training data on each algorithm's performance, we compare the mutational composition of $\Delta\Delta G$ datasets used in the development of each algorithm (Table 1). PoPMuSiC is a simplified-representation statistical energy function trained on a very large experimental dataset from ProTherm. FoldX is similarly

Table 1. Mutational composition of protein stability datasets

$\Delta\Delta G$ dataset	Total no.	Surface, %	Boundary, %	Core, %	+Vol Δ , %	-Vol Δ , %	Ala, %
FoldX training set (7)	339	32	27	36	3	97	61
FoldX test set (7)	625	32	30	35	5	95	54
PoPMuSiC training set (6)	2,644	26	32	40	33	67	28
PoPMuSiC test set (6)	350	21	27	48	40	60	26
Rosetta test set (8)	1,210	32	30	38	16	84	47
$\Delta\Delta G_{lit}$ (54)	82	71	21	8	52	48	20
This dataset	935	53	25	22	56	44	5
Top 175 variants of this dataset	175	63	32	5	78	22	3

Residues were classified as core, boundary, or surface with the RESCLASS (4) algorithm. Mutations with mislabeled or nonstandard PDB data (<5%) were omitted from residue classification. +Vol Δ , small-to-large mutations; -Vol Δ , large-to-small mutations.

trained, albeit with a smaller and more Ala-biased dataset, and mixes all-atom physical potentials with weighted statistical terms. Rosetta also mixes statistical and all-atom physical potentials, but is trained to recover native sequence compositions for protein design. A recent study systematically explored the effect of 19 different Rosetta parameter sets on single-mutant stability prediction (8), 4 of which are evaluated here. Three of the tested parameter sets use identical weights and terms but allow increasing amounts of backbone flexibility. That is, after side-chain repacking, the structure either undergoes no energy minimization, constrained backbone minimization, or unconstrained backbone minimization. Initially described as row 3, row 16, and row 19 (8), we refer to these parameter sets here as NoMin, SomeMin, and FullMin, respectively. The fourth Rosetta parameter set evaluated here (SomeMin_ddg) combines constrained minimization with optimized amino acid reference energies trained on single-mutant $\Delta\Delta G$ data from ProTherm, similar to FoldX and PoPMuSiC. Pearson correlation coefficients were used to evaluate algorithm performance overall (all mutations) as well as performance based on tertiary structure (RESCLASS) and volume change (Table 2). As energies from physical potentials can be dramatically skewed by atomic clashes, we excluded mutations with exceptionally high clash energies (clash outliers, see *SI Appendix*, Table S3). As might be expected, clash outliers typically occur in core positions where volume is increased upon substitution. All correlations here have values of $P < 0.001$ except for those in the core, as detailed in Table 2.

The Rosetta SomeMin method is the best-performing algorithm overall with a Pearson correlation coefficient of 0.64. The other methods perform less well at $r = 0.56$ (PoPMuSiC) and $r = 0.51$ (FoldX). All of the algorithms scored lower on our dataset than previously reported on independent test sets, where r values of 0.69 (Rosetta SomeMin) (8), 0.67 (PopMuSiC) (46), and 0.64 (FoldX) (7) were obtained. Comparing the different Rosetta methods, we observe that increasing backbone flexibility decreases the number of clash outliers, but does not necessarily improve overall performance. The constrained minimization in SomeMin considerably improves the correlation over NoMin, but unconstrained minimization in FullMin shows diminishing returns in allowing increased flexibility, as observed previously (8). Notably, the Rosetta SomeMin_ddg method performed worse than the SomeMin method ($r = 0.54$ and 0.64, respectively), demonstrating a limitation of training all-atom potentials with small, biased experimental datasets (Table 1).

If we look at the Pearson correlation coefficient by residue class, we find a general performance trend of boundary >

surface > core. Except for Rosetta NoMin, which performs poorly across all categories, the all-atom algorithms exhibit very strong correlations in the boundary ($r \sim 0.7$), with weaker correlations on the surface ($r \sim 0.5$). This performance bump could be explained by the presence of explicitly modeled electrostatic and hydrogen bonding interactions between protein atoms in the boundary, whereas at the surface, all interactions with solvent (water) become aggregated and averaged out due to the nature of the algorithms. In contrast, PoPMuSiC likely performs more similarly across these 2 residue classes ($r = 0.56$ and $r = 0.51$, respectively) due to several terms in the statistical function that are weighted by surface accessibility. All algorithms do a poor job at predicting core mutations (r values range from 0.13 to 0.37), possibly because these mutations are more likely to lead to structural rearrangements that are not well captured by the algorithms (6–8). The observed differences in correlation accuracy by residue class likely do not stem from differences in training data, as all 3 algorithms were trained on datasets that have very similar fractions of surface, boundary, and core residues (Table 1).

The data were also analyzed by mutations that either reduce side-chain volume (large to small, $-Vol\Delta$) or increase side-chain volume (small to large, $+Vol\Delta$). Overall, across all methods, large-to-small mutations are better predicted than the inverse, which correlates with the composition of the training sets used in algorithm development (Table 1).

All algorithms were also evaluated by the Spearman correlation coefficient to minimize penalties on skewed energies and instead reward correct rank ordering. The differences found with the Pearson method on the overall dataset are no longer observed (*SI Appendix*, Table S2). PoPMuSiC and all of the Rosetta methods perform about the same, with FoldX performing less well. However, the performance trend between residue classes is retained with boundary > surface > core, and the performance edge for large-to-small mutations is widened when evaluated by the Spearman coefficient. Because mutations that remove substantial volume often create a destabilizing cavity (52), the direction of the stability change of large-to-small mutations is more easily predicted and indeed captured by all of the algorithms equally well. The small-to-large mutation type can have very different outcomes (stabilized backbone accommodation or underpacked/overpacked destabilization) and thus is harder to rank, much less predict accurately, as observed here. The trend in the volume change data subset demonstrates why stability predictors often feature favorable correlation coefficients on their test sets, which nearly always contain a bias toward mutations to small amino acids like Ala, as observed in Table 1.

Table 2. Algorithm performance by Pearson correlation

Algorithm	Backbone minimization*	Clash outliers [†]	Pearson correlation coefficient (r)					
			Overall	Surface [‡]	Boundary [‡]	Core [‡]	+Vol Δ	-Vol Δ
PoPMuSiC		0	0.56 [#] (830)	0.51 [#] (477)	0.56 [#] (224)	0.33 [#] (129)	0.43 [#] (492)	0.64 [#] (338)
FoldX		17	0.51 [#] (813)	0.42 [#] (477)	0.68 [#] (224)	0.17 (112)	0.46 [#] (475)	0.56 [#] (338)
Rosetta [§]								
NoMin	None	22	0.33 [#] (801)	0.29 [#] (472)	0.26 [#] (222)	0.13 ^{**} (107)	0.38 [#] (468)	0.28 [#] (333)
SomeMin	Constrained	17	0.64 [#] (813)	0.53 [#] (476)	0.73 [#] (222)	0.37 [#] (115)	0.56 [#] (477)	0.66 [#] (336)
SomeMin_ddg [¶]	Constrained	6	0.54 [#] (824)	0.49 [#] (474)	0.68 [#] (224)	0.15 ^{††} (126)	0.46 [#] (487)	0.66 [#] (337)
FullMin	Unconstrained	3	0.60 [#] (827)	0.52 [#] (476)	0.69 [#] (223)	0.24 ^{††} (128)	0.48 [#] (491)	0.69 [#] (336)

Predicted $\Delta\Delta G$ s from stability algorithms were compared with experimental $\Delta\Delta G$ values for G β 1 single mutants in the quantitative dataset. Mutations with exceptionally high clash energies (clash outliers) were excluded when calculating correlation coefficients. Number of mutations is shown in parentheses. +Vol Δ , small-to-large mutations; -Vol Δ , large-to-small mutations. [#] $P < 0.001$; ^{||} $P = 0.07$; ^{**} $P = 0.193$; ^{††} $P = 0.099$; and ^{†††} $P = 0.005$.

*Level of backbone minimization after repacking for Rosetta methods.

[†]Number of mutations with a calculated repulsive energy >2 SDs above the mean.

[‡]Residues are classified as core, boundary, or surface using RESCLASS (4).

[§]Rosetta parameter sets NoMin, SomeMin, and FullMin were initially described as row 3, row 16, and row 19, respectively (8).

[¶]Combines constrained backbone minimization with optimized reference energies trained on ProTherm single-mutant $\Delta\Delta G$ data.

binding to IgG Fc, Olson et al. found that fitness values obtained using binding affinity enrichment ($\ln W$) had no correlation ($r = 0.013$) with $\Delta\Delta G$ values reported in the literature for 82 single mutants ($\Delta\Delta G_{lit}$). When we compared $\ln W$ with the $\Delta\Delta G$ values from our larger set of 830 single mutants, we found a better, but still poor correlation ($r = 0.19$) (SI Appendix, Fig. S7A).

To address this issue, Olson et al. devised a strategy to estimate single mutant stabilities from their DMS fitness data. This approach requires identifying destabilized mutational backgrounds using double-mutant fitness data so that the functional effect of a second mutation in these backgrounds could be used to compute single-mutant $\Delta\Delta G$ s. They identified 5 background mutations that produced a good correlation ($r = 0.91$) with $\Delta\Delta G_{lit}$ and later demonstrated an approach [see Wu et al. (55)] that avoids the need for preexisting stability data. In Fig. 7A, we plot our experimental $\Delta\Delta G$ s versus those predicted using the Wu et al. method ($\Delta\Delta G_{Wu}$) for 794 single mutants. The correlation ($r = 0.60$; $P < 0.001$) is dramatically lower than the value obtained using the smaller $\Delta\Delta G_{lit}$ dataset ($r = 0.91$). A closer look at the 82 mutants in $\Delta\Delta G_{lit}$ reveals a relatively small percentage of mutations in the core and a bias toward alanine substitutions, resulting in a dataset that does not reflect the breadth of possible mutations in the entire domain (Table 1). As seen in SI Appendix, Fig. S7B, the limited number of mutants in $\Delta\Delta G_{lit}$ masks the lower correlation between $\Delta\Delta G$ and $\Delta\Delta G_{Wu}$ by serendipitously avoiding off-diagonal single mutants.

A recent report by Otwinowski (56) reanalyzed the Olson et al. fitness data with a method based on a thermodynamic model describing 3 states (bound-folded, unbound-folded, and unfolded) that avoids the need for preexisting mutational or structural data. The method calculates distinct energies for folding ($E_{folding}$) and binding ($E_{binding}$). We compare the $E_{folding}$ energy ($\Delta\Delta G_{Otwinowski}$) with our experimental $\Delta\Delta G$ values in Fig. 7B, which shows an improved correlation ($r = 0.72$; $P < 0.001$) over the Wu et al. method ($r = 0.60$). SI Appendix, Table S4 analyzes the correlations for the 2 methods by residue class, volume change, and polarity change. The $\Delta\Delta G_{Otwinowski}$ energy yields better correlations across the board, with the core continuing to show a significantly lower correlation. Thus, although DMS fitness data are poorly correlated with thermodynamic stability, simple biophysical models can be constructed that lead to significantly improved correlations. We expect that large, comprehensive datasets containing thermodynamic measurements

such as those provided here will facilitate the development of improved methods to extract biophysical quantities (e.g., stability and binding) from fitness data, thus greatly expanding the utility of DMS and other deep-sequencing techniques.

Discussion

We described an automated chemical denaturation methodology that produces high-quality thermodynamic stability data at a throughput that enables the near-total site saturation mutagenesis of small protein domains. Although other low-cost methods such as thermal challenge assays or differential scanning fluorimetry can also provide useful data, and deep-sequencing approaches such as DMS can streamline the entire process, these methods do not directly report thermodynamic information. The automated pipeline described here makes gathering accurate thermodynamic stability data at a large scale feasible. The broad, unbiased nature of our near-complete G β 1 single-mutant study provides an important dataset for examining domain-wide trends, evaluating stability-prediction tools, and validating methods to extract stability values from DMS results. In addition, our analysis highlights the impact that training sets can have on computational predictors of stability.

We found that while the stability distribution of our G β 1 dataset features a long tail of destabilizing variants, most mutations (68%) are neutral. However, if variants without quantitative data and those omitted for technical reasons are assigned negative outcomes, destabilizing variants make up 45% of the 1,064 possible single mutants of G β 1, approaching predicted published values (34). Other trends followed conventions, with mutations to Gly, Pro, and core positions almost always being deleterious. However, not all core positions (as determined by RESCLASS) show the same degree of sensitivity, as measures of residue burial such as C β atom depth and OSP were found to best correlate with median $\Delta\Delta G$ at each position. Although the correlation of residue burial with individual $\Delta\Delta G$ measurements was previously reported for a collection of variants across many proteins (39), our domain-wide dataset allows the position-specific nature of the relationship to be fully observed. Similarly, using our unique dataset to calculate median $\Delta\Delta G$ by incorporated amino acid reveals an unexpected tolerance for large hydrophobic amino acids. This preference extended across tertiary structure, and stability predictions on 4 other proteins confirmed this trend.

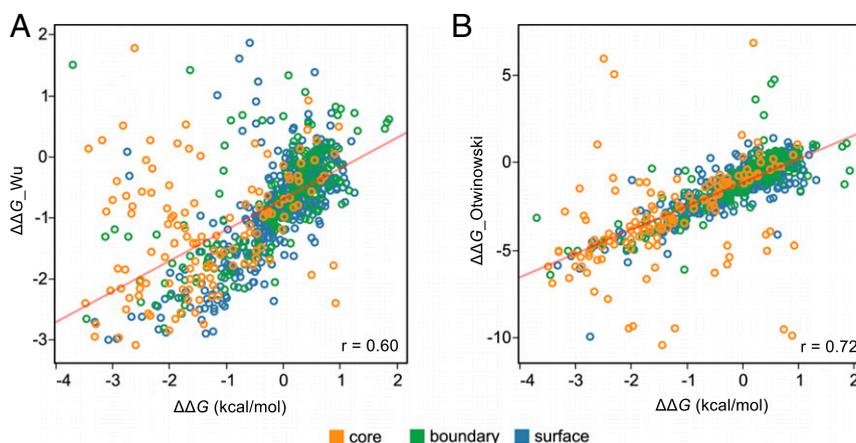


Fig. 7. Comparing experimental $\Delta\Delta G$ s with predictions obtained from DMS fitness data. G β 1 single-mutant stabilities from our experimental quantitative dataset ($\Delta\Delta G$) are plotted against (A) $\Delta\Delta G$ values predicted from the DMS data using the Wu et al. method (55) ($\Delta\Delta G_{Wu}$) ($n = 794$) or (B) $E_{folding}$ values predicted from the DMS data using Otwinowski's 3-state thermodynamic model (56) ($\Delta\Delta G_{Otwinowski}$) ($n = 812$). In both cases, DMS data are from Olson et al. (54). Points are colored by RESCLASS (4) values. A linear regression line is shown in red, and the Pearson correlation coefficient ($P < 0.001$) is shown in the Lower Right of each plot.

was calculated over a 20-ns molecular-dynamics trajectory in full solvent using NAMD (60). The depth of the C β atom was calculated by the RESCLASS algorithm (4) to decide core, boundary, and surface residues. Linear regression with 10-fold cross-validation was performed with scikit-learn (61) to identify attributes that correlate highly with positional sensitivity. Recursive feature elimination was also performed with scikit-learn using a ridge estimator, and 5-fold cross-validation was performed to evaluate combinations of attributes. Recursive feature elimination was also performed with scikit-learn to evaluate combinations of attributes.

Stability-Prediction Algorithms. The crystal structure of G β 1 (PDB ID: 1PGA) was used as the input structure for all algorithms. The webserver for PoPMuSIC, version 3, located at <http://www.dezyme.com>, was used to perform a "Systematic" command on the G β 1 crystal structure. A copy of FoldX (Version 3.0, beta 5) was retrieved from <http://foldxsuite.crg.eu/>. The crystal structure was prepared by using the "RepairPDB" command to perform Asn, Gln, and His flips, alleviate small van der Waals clashes, and optimize WT rotamer packing. Every mutant in the dataset was constructed through the "BuildModel" command, and the difference in energy between the WT reference and the corresponding mutant was averaged over 5 trials. A copy of Rosetta (version 3.3) was retrieved from <https://www.rosettacommons.org/>. The *ddg_monomer* application was used to generate single-mutant stability

data from the G β 1 crystal structure. We followed the available online documentation to prepare all necessary input files. Option sets described in the documentation pertain to the various Rosetta iterations tested in this paper (NoMin, low-resolution protocol; SomeMin, high-resolution protocol; FullMin, high-resolution protocol with an empty distance restraints file).

Statistical Visualization and Analysis. All plots were generated using Tableau. Custom Python scripts were developed to calculate the large number of thermodynamic stability curve fits. Correlation coefficients (Pearson's and Spearman's) were calculated either in Tableau or in the software package R (version 3.2.2). The ROC package for R was used for classification and ROC analysis (62).

Data Availability. The $\Delta\Delta G$ distribution of G β 1 single mutants generated during this work is publicly available in ProtaBank (<https://www.protabank.org/>), a protein engineering data repository, under the ID gwoS2haU3.

ACKNOWLEDGMENTS. A.N. thanks Jost Vielmetter for advice and feedback on the automated platform. S.L.M. acknowledges grants from the National Security Science and Engineering Faculty Fellowship program and the Defense Advanced Research Projects Agency Protein Design Processes program.

1. A. S. Bommaris, M. F. Paye, Stabilizing biocatalysts. *Chem. Soc. Rev.* **42**, 6534–6565 (2013).
2. A. Goldenzweig, S. J. Fleishman, Principles of protein stability and their application in computational design. *Annu. Rev. Biochem.* **87**, 105–129 (2018).
3. R. Rouet, D. Lowe, D. Christ, Stability engineering of the human antibody repertoire. *FEBS Lett.* **588**, 269–277 (2014).
4. B. I. Dahiya, S. L. Mayo, De novo protein design: Fully automated sequence selection. *Science* **278**, 82–87 (1997).
5. R. Das, D. Baker, Macromolecular modeling with Rosetta. *Annu. Rev. Biochem.* **77**, 363–382 (2008).
6. Y. Dehouck *et al.*, Fast and accurate predictions of protein stability changes upon mutations using statistical potentials and neural networks: PoPMuSIC-2.0. *Bioinformatics* **25**, 2537–2543 (2009).
7. R. Guerois, J. E. Nielsen, L. Serrano, Predicting changes in the stability of proteins and protein complexes: A study of more than 1000 mutations. *J. Mol. Biol.* **320**, 369–387 (2002).
8. E. H. Kellogg, A. Leaver-Fay, D. Baker, Role of conformational sampling in computing mutation-induced changes in protein structure and stability. *Proteins* **79**, 830–838 (2011).
9. S. M. Malakauskas, S. L. Mayo, Design, structure and stability of a hyperthermophilic protein variant. *Nat. Struct. Biol.* **5**, 470–475 (1998).
10. F. Zheng, G. Grigoryan, Sequence statistics of tertiary structural motifs reflect protein stability. *PLoS One* **12**, e0178272 (2017).
11. D. E. V. Pires, J. Chen, T. L. Blundell, D. B. Ascher, In silico functional dissection of saturation mutagenesis: Interpreting the relationship between phenotypes and changes in protein stability, interactions and activity. *Sci. Rep.* **6**, 19848 (2016).
12. D. J. Burgess, Disease genetics: Network effects of disease mutations. *Nat. Rev. Genet.* **16**, 317 (2015).
13. V. Potapov, M. Cohen, G. Schreiber, Assessing computational methods for predicting protein stability upon mutation: Good on average but not in the details. *Protein Eng. Des. Sel.* **22**, 553–560 (2009).
14. S. Khan, M. Vihinen, Performance of protein stability predictors. *Hum. Mutat.* **31**, 675–684 (2010).
15. J. A. Davey, R. A. Chica, Optimization of rotamers prior to template minimization improves stability predictions made by computational protein design. *Protein Sci.* **24**, 545–560 (2015).
16. T. Alber, Mutational effects on protein stability. *Annu. Rev. Biochem.* **58**, 765–798 (1989).
17. J. U. Bowie, J. F. Reidhaar-Olson, W. A. Lim, R. T. Sauer, Deciphering the message in protein sequences: Tolerance to amino acid substitutions. *Science* **247**, 1306–1310 (1990).
18. B. W. Matthews, Structural and genetic analysis of protein stability. *Annu. Rev. Biochem.* **62**, 139–160 (1993).
19. A. R. Fersht, L. Serrano, Principles of protein stability derived from protein engineering experiments. *Curr. Opin. Struct. Biol.* **3**, 75–83 (1993).
20. M. D. S. Kumar *et al.*, ProTherm and ProNIT: Thermodynamic databases for proteins and protein-nucleic acid interactions. *Nucleic Acids Res.* **34**, D204–D206 (2006).
21. C. Y. Wang *et al.*, ProtaBank: A repository for protein design and engineering data. *Protein Sci.* **27**, 1113–1124 (2018).
22. T. J. Magliery, Protein stability: Computation, sequence statistics, and new experimental methods. *Curr. Opin. Struct. Biol.* **33**, 161–168 (2015).
23. P. Markiewicz, L. G. Kleina, C. Cruz, S. Ehret, J. H. Miller, Genetic studies of the lac repressor. XIV. Analysis of 4000 altered *Escherichia coli lac* repressors reveals essential and non-essential residues, as well as "spacers" which do not require a specific sequence. *J. Mol. Biol.* **240**, 421–433 (1994).
24. B. D. Allen, A. Nisthal, S. L. Mayo, Experimental library screening demonstrates the successful application of computational protein design to large structural ensembles. *Proc. Natl. Acad. Sci. U.S.A.* **107**, 19838–19843 (2010).
25. J. P. Aucamp, A. M. Cosme, G. J. Lye, P. A. Dalby, High-throughput measurement of protein stability in microtiter plates. *Biotechnol. Bioeng.* **89**, 599–607 (2005).
26. J. J. Lavinder, S. B. Hari, B. J. Sullivan, T. J. Magliery, High-throughput thermal scanning: A general, rapid dye-binding thermal shift screen for protein engineering. *J. Am. Chem. Soc.* **131**, 3794–3795 (2009).
27. G. J. Rocklin *et al.*, Global analysis of protein folding using massively parallel design, synthesis, and testing. *Science* **357**, 168–175 (2017).
28. C. L. Araya, D. M. Fowler, Deep mutational scanning: Assessing protein function on a massive scale. *Trends Biotechnol.* **29**, 435–442 (2011).
29. D. M. Fowler, S. Fields, Deep mutational scanning: A new style of protein science. *Nat. Methods* **11**, 801–807 (2014).
30. C. N. Pace, Determination and analysis of urea and guanidine hydrochloride denaturation curves. *Methods Enzymol.* **131**, 266–280 (1986).
31. M. M. Santoro, D. W. Bolen, Unfolding free energy changes determined by the linear extrapolation method. 1. Unfolding of phenylmethanesulfonyl alpha-chymotrypsin using different denaturants. *Biochemistry* **27**, 8063–8068 (1988).
32. J. K. Myers, C. N. Pace, J. M. Scholtz, Denaturant m values and heat capacity changes: Relation to changes in accessible surface areas of protein unfolding. *Protein Sci.* **4**, 2138–2148 (1995).
33. A. Nisthal, Protein stability engineering insights revealed by domain-wide comprehensive mutagenesis. Protabank. https://www.protabank.org/study_analysis/gwoS2haU3/. Deposited 9 October, 2018.
34. N. Tokuriki, F. Stricher, J. Schymkowitz, L. Serrano, D. S. Tawfik, The stability effects of protein mutations appear to be universally distributed. *J. Mol. Biol.* **369**, 1318–1332 (2007).
35. N. Tokuriki, D. S. Tawfik, Stability effects of mutations and protein evolvability. *Curr. Opin. Struct. Biol.* **19**, 596–604 (2009).
36. S. Henikoff, J. G. Henikoff, Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. U.S.A.* **89**, 10915–10919 (1992).
37. W. Huang, J. Petrosino, M. Hirsch, P. S. Shenkin, T. Palzkill, Amino acid sequence determinants of beta-lactamase structure and activity. *J. Mol. Biol.* **258**, 688–703 (1996).
38. D. Rennell, S. E. Bouvier, L. W. Hardy, A. R. Poteete, Systematic mutation of bacteriophage T4 lysozyme. *J. Mol. Biol.* **222**, 67–88 (1991).
39. S. Chakravarty, R. Varadarajan, Residue depth: A novel parameter for the analysis of protein structure and stability. *Structure* **7**, 723–732 (1999).
40. K. P. Tan, T. B. Nguyen, S. Patel, R. Varadarajan, M. S. Madhusudhan, Depth: A web server to compute depth, cavity sizes, detect potential small-molecule ligand-binding cavities and predict the pK_a of ionizable residues in proteins. *Nucleic Acids Res.* **41**, W314–W321 (2013).
41. N. Pattabiraman, K. B. Ward, P. J. Fleming, Occluded molecular surface: Analysis of protein packing. *J. Mol. Recognit.* **8**, 334–344 (1995).
42. P. J. Fleming, F. M. Richards, Protein packing: Dependence on protein size, secondary structure and amino acid composition. *J. Mol. Biol.* **299**, 487–498 (2000).
43. W. S. J. Valdar, Scoring residue conservation. *Proteins* **48**, 227–241 (2002).
44. J. Kyte, R. F. Doolittle, A simple method for displaying the hydrophobic character of a protein. *J. Mol. Biol.* **157**, 105–132 (1982).
45. D. A. Dixon, W. N. Lipscomb, Electronic structure and bonding of the amino acids containing first row atoms. *J. Biol. Chem.* **251**, 5992–6000 (1976).
46. Y. Dehouck, J. M. Kwasiogoch, D. Gillis, M. Rومان, PoPMuSIC 2.1: A web server for the estimation of protein stability changes upon mutation and sequence optimality. *BMC Bioinformatics* **12**, 151 (2011).
47. S. Ayuso-Tejedor, O. Abián, J. Sancho, Underexposed polar residues and protein stabilization. *Protein Eng. Des. Sel.* **24**, 171–177 (2011).
48. A. Broom, Z. Jacobi, K. Trainor, E. M. Meiering, Computational tools help improve protein stability but with a solubility tradeoff. *J. Biol. Chem.* **292**, 14349–14361 (2017).
49. M. H. Cordes, R. T. Sauer, Tolerance of a protein to multiple polar-to-hydrophobic surface substitutions. *Protein Sci.* **8**, 318–325 (1999).

50. M. Machius, N. Declerck, R. Huber, G. Wiegand, Kinetic stabilization of *Bacillus licheniformis* alpha-amylase through introduction of hydrophobic residues at the surface. *J. Biol. Chem.* **278**, 11546–11553 (2003).
51. D. Poso, R. B. Sessions, M. Lorch, A. R. Clarke, Progressive stabilization of intermediate and transition states in protein folding reactions by introducing surface hydrophobic residues. *J. Biol. Chem.* **275**, 35723–35726 (2000).
52. W. A. Baase, L. Liu, D. E. Tronrud, B. W. Matthews, Lessons from the lysozyme of phage T4. *Protein Sci.* **19**, 631–641 (2010).
53. O. Buß, J. Rudat, K. Ochsenreither, FoldX as protein engineering tool: Better than random based approaches? *Comput. Struct. Biotechnol. J.* **16**, 25–33 (2018).
54. C. A. Olson, N. C. Wu, R. Sun, A comprehensive biophysical description of pairwise epistasis throughout an entire protein domain. *Curr. Biol.* **24**, 2643–2651 (2014).
55. N. C. Wu, C. A. Olson, R. Sun, High-throughput identification of protein mutant stability computed from a double mutant fitness landscape. *Protein Sci.* **25**, 530–539 (2016).
56. J. Otwinowski, Biophysical inference of epistasis and the effects of mutations on protein stability and function. arXiv:1802.08744v2 (30 March 2018).
57. A. Nisthal, "Accelerating the interplay between theory and experiment in protein design," PhD thesis, California Institute of Technology, Pasadena, CA (2012).
58. H. E. Klock, S. A. Lesley, The Polymerase Incomplete Primer Extension (PIPE) method applied to high-throughput cloning and site-directed mutagenesis. *Methods Mol. Biol.* **498**, 91–103 (2009).
59. W. Humphrey, A. Dalke, K. Schulten, VMD: Visual molecular dynamics. *J. Mol. Graph.* **14**, 33–38, 27–28 (1996).
60. J. C. Phillips et al., Scalable molecular dynamics with NAMD. *J. Comput. Chem.* **26**, 1781–1802 (2005).
61. F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
62. T. Sing, O. Sander, N. Beerenwinkel, T. Lengauer, ROCr: Visualizing classifier performance in R. *Bioinformatics* **21**, 3940–3941 (2005).
63. W. Kabsch, C. Sander, Dictionary of protein secondary structure: Pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* **22**, 2577–2637 (1983).