

PROTEIN FOLDING

Global analysis of protein folding using massively parallel design, synthesis, and testing

Gabriel J. Rocklin,¹ Tamuka M. Chidyausiku,^{1,2} Inna Goreshnik,¹ Alex Ford,^{1,2} Scott Houliston,^{3,4} Alexander Lemak,³ Lauren Carter,¹ Rashmi Ravichandran,¹ Vikram K. Mulligan,¹ Aaron Chevalier,¹ Cheryl H. Arrowsmith,^{3,4,5} David Baker^{1,6*}

Proteins fold into unique native structures stabilized by thousands of weak interactions that collectively overcome the entropic cost of folding. Although these forces are “encoded” in the thousands of known protein structures, “decoding” them is challenging because of the complexity of natural proteins that have evolved for function, not stability. We combined computational protein design, next-generation gene synthesis, and a high-throughput protease susceptibility assay to measure folding and stability for more than 15,000 de novo designed miniproteins, 1000 natural proteins, 10,000 point mutants, and 30,000 negative control sequences. This analysis identified more than 2500 stable designed proteins in four basic folds—a number sufficient to enable us to systematically examine how sequence determines folding and stability in uncharted protein space. Iteration between design and experiment increased the design success rate from 6% to 47%, produced stable proteins unlike those found in nature for topologies where design was initially unsuccessful, and revealed subtle contributions to stability as designs became increasingly optimized. Our approach achieves the long-standing goal of a tight feedback cycle between computation and experiment and has the potential to transform computational protein design into a data-driven science.

The key challenge to achieving a quantitative understanding of the sequence determinants of protein folding is to accurately and efficiently model the balance among the many energy terms that contribute to the free energy of folding (1–3). Minimal protein domains (30 to 50 amino acids in length), such as the villin headpiece and WW domain, are commonly used to investigate this balance because they are the simplest protein folds found in nature (4). The primary experimental approach used to investigate this balance has been mutagenesis (5–12), but the results are context-dependent and do not provide a global view of the contributions to stability. Molecular dynamics simulations on minimal proteins have also been used to study folding (13–15), but these do not reveal which interactions specify and stabilize the native structure, and in general they cannot determine whether a given sequence will fold into a stable structure.

De novo protein design has the potential to reveal the sequence determinants of folding for minimal proteins by charting the space of non-natural sequences and structures to define what can and cannot fold. Protein sequence space (16) is vastly larger than the set of natural proteins that currently form the basis for nearly all models of protein stability (9, 12, 17–19) and is unbiased by selection for biological function. However, only two minimal proteins (<50 amino acids, stabilized exclusively by noncovalent interactions) have been computationally designed to date: FSD-1 (20) and DSI19 (21). In part, this is attributable to the cost of gene synthesis, which has limited such studies to testing tens of designs at most—a minuscule fraction of design space. Because of the small sample sizes, design experiments are typically unable to determine why some designs are stable and others are unstructured, resemble molten globules, or form aggregates (22).

Here, we present a new synthetic approach to examining the determinants of protein folding by exploring the space of potential minimal proteins using de novo computational protein design, with data generated by parallel DNA synthesis and protein stability measurements. To encode our designs, we use oligo library synthesis technology (23, 24), which was originally developed for transcriptional profiling and large gene assembly applications and is now capable of parallel synthesis of 10^4 to 10^5 arbitrarily specified DNA sequences long enough to encode short proteins (fig. S1). To assay designs for stabil-

ity, we express these libraries in yeast so that every cell displays many copies of one protein sequence on its surface, genetically fused to an expression tag that can be fluorescently labeled (25) (Fig. 1A). Cells are then incubated with varying concentrations of protease, those displaying resistant proteins are isolated by fluorescence-activated cell sorting (Fig. 1B), and the frequencies of each protein at each protease concentration are determined by deep sequencing (Fig. 1C; for reproducibility of the assay, see fig. S2). We then infer protease EC_{50} values (the protease concentration at which one-half of the cells pass the collection threshold) for each sequence from these data by modeling the complete selection procedure (Fig. 1D) (26). Finally, each design is assigned a “stability score” (Fig. 1E): the difference between the measured EC_{50} (on a \log_{10} scale) and the predicted EC_{50} in the unfolded state, according to a sequence-based model parameterized using EC_{50} measurements of scrambled sequences (figs. S3 and S4). A stability score of 1 corresponds to an EC_{50} value that exceeds the predicted EC_{50} in the unfolded state by a factor of 10. The complete experimental procedure costs less than \$7000 in reagents (mainly from DNA synthesis and sequencing) and requires ~10 hours of sorting per protease for each library.

Massively parallel measurement of folding stability

Proteolysis assays have been used to select for stable sequences (27–29) and to quantify stability for individual proteins (30) and proteins from cellular proteomes (31), but to date they have not been used to quantify stability for all sequences in a constructed library. To evaluate the ability of the assay to measure stability on a large scale, we obtained a synthetic DNA library encoding four small proteins [Pin1 WW domain (32), hYAP65 WW domain (5, 10), villin HP35 (7, 11), and BBL (8)] and 116 mutants of these proteins whose stability has been characterized in experiments on purified material. The library also contained 19,610 unrelated sequences (a fourth-generation designed protein library; see below), and all sequences were assayed for stability simultaneously. Although the stability score is not a direct analog of a thermodynamic parameter, stability scores measured with trypsin and separately measured with chymotrypsin were each well correlated with folding free energies (or melting temperatures) for all four sets of mutants, with r^2 values ranging from 0.63 to 0.85 (Fig. 1, F to I). Most mutants in this data set were predicted to have unfolded-state EC_{50} values similar to those of their parent sequences, so the relative stability scores of the mutants are very similar to their relative EC_{50} values. However, in the case of villin assayed with chymotrypsin, the unfolded-state model improved the correlation between protease resistance and folding free energy from $r^2 = 0.46$ (using raw EC_{50} values) to the reported $r^2 = 0.77$ by correcting for the effect of mutations such as Lys⁷⁰ → Met and Phe⁵¹ → Leu on intrinsic chymotrypsin cleavage rates. The mutual agreement among trypsin results, chymotrypsin results,

¹Department of Biochemistry and Institute for Protein Design, University of Washington, Seattle, WA 98195, USA.

²Graduate Program in Biological Physics, Structure, and Design, University of Washington, Seattle, WA 98195, USA.

³Princess Margaret Cancer Centre, Toronto, Ontario M5G 1L7, Canada. ⁴Structural Genomics Consortium, University of Toronto, Toronto, Ontario M5G 1L7, Canada.

⁵Department of Medical Biophysics, University of Toronto, Toronto, Ontario M5G 1L7, Canada. ⁶Howard Hughes Medical Institute, University of Washington, Seattle, WA 98195, USA.

*Corresponding author. Email: dabaker@u.washington.edu

and experiments on purified protein indicates that the assay provides a robust measure of folding stability for small proteins.

Massively parallel testing of designed mini-proteins

We selected four protein topologies ($\alpha\alpha\alpha$, $\beta\alpha\beta\beta$, $\alpha\beta\beta\alpha$, and $\beta\beta\alpha\beta$) as design targets. These topologies have increasing complexity: The $\alpha\alpha\alpha$ topology features only two loops and exclusively local secondary structure (helices); the $\beta\beta\alpha\beta$ fold requires four loops and features a mixed parallel/antiparallel β sheet bridging the N and C termini. Of these topologies, only $\alpha\alpha\alpha$ proteins have been found in nature within the target size range of 40 to 43 residues; no proteins have been previously designed in any of the four topologies at this size [excluding designed $\alpha\alpha\alpha$ and $\beta\alpha\beta\beta$ proteins stabilized by multiple disulfide linkages

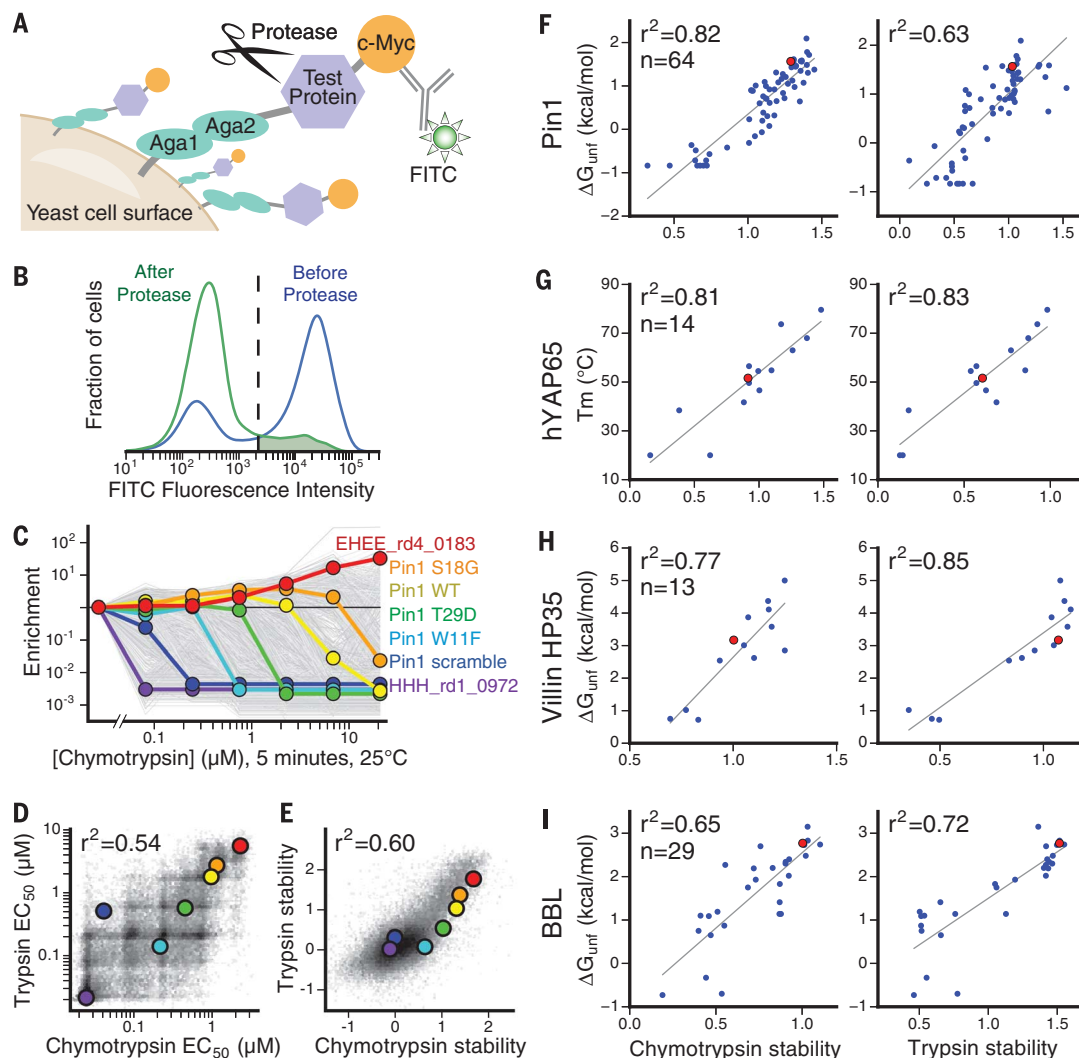
(33)]. For each topology, we first designed between 5000 and 40,000 de novo proteins using a blueprint-based approach described in (34). Each design has a unique three-dimensional main-chain conformation and a unique sequence predicted to be near-optimal for that conformation. We then selected 1000 designs per topology for experimental testing by ranking the designs by a weighted sum of their computed energies and additional filtering terms (26). The median sequence identity between any pair of tested designs of the same topology ranged from 15% to 35%, and designs were typically no more than 40 to 65% identical to any other design. This diversity is due to the different backbone conformations possible within a topology, along with the vast sequence space available even for small proteins (fig. S5). For each design, we also included two control sequences in our library: one made

by scrambling the order of amino acids in that design (preserving the overall amino acid composition), and a second made by scrambling the order while preserving both the composition and the hydrophobic or polar character at each position (35–37). The library comprised 12,459 different sequences in total: 4153 designed proteins and 8,306 control sequences. The designed proteins are named according to their secondary structure topology (using H for α helix and E for β strand), their design round, and a design number.

We assayed the sequence library for stability using both chymotrypsin and trypsin. To stringently identify stable designs, we ranked sequences by the lower of their trypsin or chymotrypsin stability scores, referred to simply as their (overall) stability score from here on. The fully scrambled sequences and patterned scrambled sequences had similar stability score distributions; most

Fig. 1. Yeast display enables massively parallel measurement of protein stability.

(A) Each yeast cell displays many copies of one test protein fused to Aga2. The C-terminal c-Myc tag is labeled with a fluorescent antibody. Protease cleavage of the test protein (or other cleavage) leads to loss of the tag and loss of fluorescence. FITC, fluorescein isothiocyanate. (B) Libraries of 10^4 unique sequences are sorted by flow cytometry. Most cells show high protein expression (measured by fluorescence) before proteolysis (blue). Only some cells retain fluorescence after proteolysis; those above a threshold (shaded green region) are collected for deep sequencing analysis. (C) Sequential sorting at increasing protease concentrations separates proteins by stability. Each sequence in a library of 19,726 proteins is shown as a gray line tracking its change in population fraction relative to that in the preselection library (enrichment). Enrichment traces for seven proteins at different stability levels are highlighted in color. (D) EC_{50} values for the seven highlighted proteins in (C) are plotted on top of the overall density of the 46,187 highest-confidence EC_{50} measurements from design rounds 1 to 4. (E) Same data as at left, showing that stability scores (EC_{50} values corrected for intrinsic proteolysis rates) correlate better than raw EC_{50} values between the proteases. (F to I) Stability scores measured in high-throughput correlate with individual folding stability measurements for mutants of four small proteins. The wild-type sequence in each set is highlighted as a red circle. Credible intervals for all EC_{50} measurements are provided in (26). (F) Pin1 ΔG_{unf} data at 40°C from (32) by thermal denaturation. (G) hYAP65 melting temperature (T_m) data from (5, 10). (H) Villin HP35 ΔG_{unf} data at 25°C from (7, 11) by urea denaturation. (I) BBL ΔG_{unf} data at 10°C from (8) by thermal denaturation.



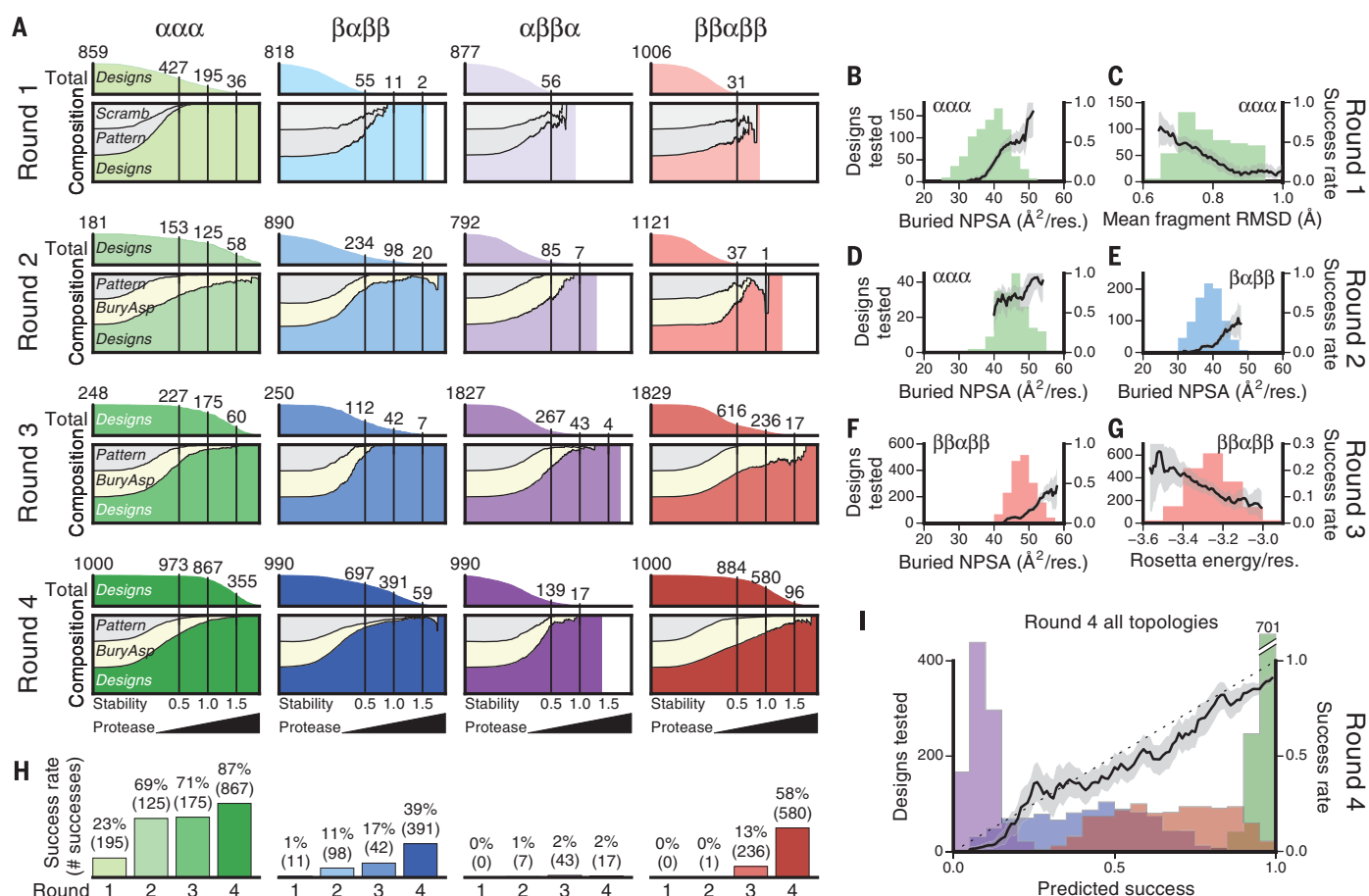


Fig. 2. Iterative, high-throughput computational design generates thousands of stable proteins and reveals stability determinants.

(A) Stability data for designs and control sequences, separated by topology ($\alpha\alpha\alpha$, $\beta\alpha\beta\beta$, $\alpha\beta\alpha\alpha$, and $\beta\beta\alpha\beta$) and by design round (1 to 4). For each round and topology, the upper plot shows the total number of designed proteins (y axis) exceeding a given stability score threshold (x axis; stability increases from left to right). The number of designs tested (top left) may be lower than the number originally ordered (described in the text) because low-confidence data were removed (26). Lower plots show the relative amounts of the three categories of sequences (y axis) exceeding a given stability score threshold (x axis), as above. Round 1 categories were designed sequences (colors), fully scrambled sequences ("Scramb," light gray), and hydrophobic-polar pattern-preserving scrambled sequences ("Pattern," dark gray). Categories in rounds 2 to 4 were designs, patterned scrambles, and point mutants of designs, with single Asp mutations expected to be destabilizing ("BuryAsp," yellow). (B to G) Determinants of stability from rounds 1 to 3 [as labeled in (A)]. Colored histograms show the number

of tested designs (left y axis) in each bin for the structural metric on the x axis. Black lines show the success rate (fraction of designs tested with stability score > 1.0, right y axis) within a moving window the size of the histogram bin width, with a shaded 95% confidence interval from bootstrapping. Design success is shown as a function of NPSA from hydrophobic residues [(B), (D), (E), and (F)]; as a function of geometric agreement between 9-residue fragments of similar sequences in the design models and natural proteins [see text and (26)], measured in average RMSD (C); and as a function of Rosetta total energy (G). (H) Overall success rate and number of successful designs per round (stability score > 1.0 with both proteases) for all topologies across all rounds. (I) Design success as a function of predicted success according to the topology-specific logistic regression models used to select round-4 designs for testing (trained on data from rounds 1 to 3). As in (B) to (G), colored histograms indicate the number of tested designs at each level of predicted success (left y axis), and the black line indicates the success rate (right y axis). See fig. S8 for individual success rates for each topology.

of these controls had stability scores below 0.5, and only one had a score greater than 1.0 (Fig. 2A, round 1). In contrast, 206 designed sequences had stability scores above 1.0 (Fig. 2A, round 1). Most of these (195 of 206) were $\alpha\alpha\alpha$ designs (both left-handed and right-handed bundles); the remaining 11 were $\beta\alpha\beta\beta$. The clustering of the 206 most stable designs around the $\alpha\alpha\alpha$ topology, and the high stability of designed sequences relative to control sequences with chemically identical compositions, strongly suggest that these stable designs fold into their designed structures.

To examine this further, we selected six stable designs (four $\alpha\alpha\alpha$ and two $\beta\alpha\beta\beta$) for *Escherichia coli* expression, purification, and further characterization by size exclusion chromatography (SEC) and circular dichroism (CD) spectroscopy. All six designs eluted from SEC as expected for a 5- to 7-kDa monomer, and the CD spectra were consistent with the designed secondary structure (fig. S6A and table S1). Five of the six designs had clear, cooperative melting transitions, refolded reversibly, and were highly stable for minimal proteins: All had melting temperatures above 70°C, and the $\beta\alpha\beta\beta$ design EHEE_rd1_0284 had

only partially melted at 95°C (free energy of unfolding $\Delta G_{\text{unf}} = 4.7$ kcal/mol; Fig. 3D). The sixth design, HHH_rd1_0005, did not refold and showed signs of aggregation (fig. S6A). We determined solution structures for EHEE_rd1_0284 and the left-handed $\alpha\alpha\alpha$ design HHH_rd1_0142 by nuclear magnetic resonance (NMR); each structure closely matched the design model [average backbone root mean square deviation (RMSD) = 2.2 Å for each NMR ensemble member against the design model] (Fig. 3A; see table S2 for NMR data summary). In sum, both high-throughput control experiments and

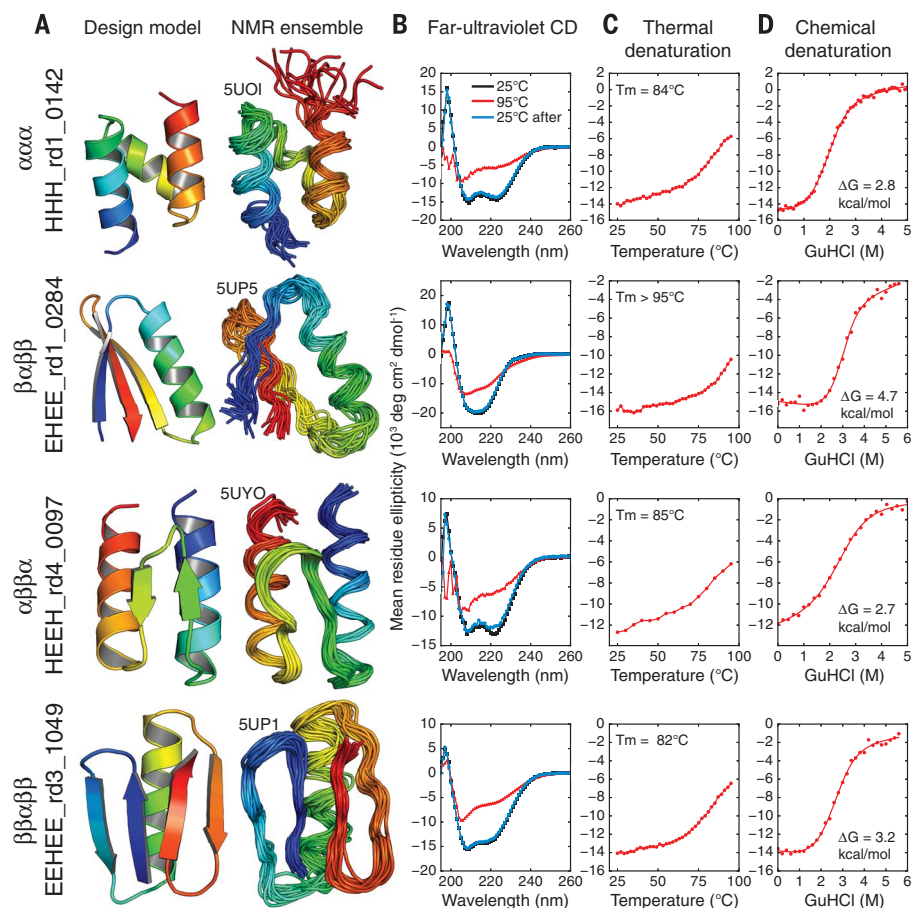


Fig. 3. Biophysical characterization of designed minimal proteins. (A) Design models and NMR solution ensembles for designed minimal proteins. PDB codes are given above each NMR ensemble. (B) Far-ultraviolet CD spectra at 25°C (black), 95°C (red), and 25°C after melting (blue). (C) Thermal melting curves measured by CD at 220 nm. Melting temperatures were determined using the derivative of the curve. (D) Chemical denaturation in GuHCl measured by CD at 220 nm and 25°C. Unfolding free energies were determined by fitting to a two-state model (red solid line). CD data for all 22 purified proteins are given in table S1 and fig. S6.

low-throughput characterization of individual proteins indicate that the protease-resistant designs folded as designed.

Global determinants of stability

This large set of stable and unstable minimal proteins with varying physical properties enabled us to quantitatively examine which protein features correlated with folding. We computed more than 60 structural and sequence-based metrics and examined which metrics differed between the 195 most stable $\alpha\alpha\alpha$ designs (stability score > 1.0 , considered to be design successes) and the 664 remaining $\alpha\alpha\alpha$ designs (considered to be failures) using the Kolmogorov-Smirnov two-sample test. Significant differences indicate that a particular metric captures an important contribution to protein stability and that this contribution was poorly optimized among the tested designs.

The dominant difference between stable and unstable $\alpha\alpha\alpha$ designs was the total amount of buried nonpolar surface area (NPSA) from hydrophobic amino acids (Fig. 2B). Stable designs buried more NPSA than did unstable designs ($P < 5 \times$

10^{-38} ; fig. S7A), and none of the 95 designs below 32 \AA^2 per residue were stable. Above this threshold, the success rate (ratio of successful designs to tested designs) steadily increased as buried NPSA increased (Fig. 2B). Stable designs also had better agreement between their sequences and their local structures, as assessed by quantifying the geometric similarity (in \AA of RMSD) between 9-residue fragments of the designs and 9-residue fragments of natural proteins similar in local sequence to the designed fragment (Fig. 2C) (26). Fragments of stable designs were more geometrically similar to fragments of natural proteins of similar local sequence, whereas fragments of unstable designs were more geometrically distant from the fragments of natural proteins matching their local sequence ($P < 2 \times 10^{-26}$; fig. S7B). Other metrics were only weakly correlated with success despite substantial variability among designs, including different measures of amino acid packing density and the total Rosetta energy itself. Although local sequence structure agreement and especially buried NPSA are well known to be important for protein

stability (1, 9), it is very challenging to determine the precise strength of these contributions at a global level in the complex balance of all the energetic contributions influencing protein structure. Our results directly demonstrate how specific imbalances (underweighting buried NPSA and local sequence structure agreement in the Rosetta energy model and the design procedure) led to hundreds of design failures, and our data and approach provide a new route to refining this balance in biophysical modeling.

Iterative, data-driven protein design

We sought to use these findings to increase the success rate of protein design by (i) changing the design procedure to increase buried NPSA, and (ii) reweighting the metrics used to select designs for testing (26). Using the improved design and ranking procedure, we built a second generation of 4150 designs, along with two control sequences per design: a pattern-preserving scrambled sequence as before (now also preserving Gly and Pro positions), and a second control identical to the designed sequence, but with the most buried side chain (according to the design model) replaced with Asp. As in round 1, almost no scrambled sequences had stability scores greater than 1 (our cutoff defining success) despite the increased hydrophobicity of the scrambled sequences (Fig. 2A, round 2). However, a much larger proportion of second-generation designs proved stable: Success for $\alpha\alpha\alpha$ designs improved from 23% to 69%, $\beta\alpha\beta\beta$ designs improved from 1% to 11% successful, and we also obtained seven stable $\alpha\beta\beta\alpha$ designs and one stable $\beta\beta\alpha\beta$ design (Fig. 2H). These increases demonstrate how iterative, high-throughput protein design can make concrete improvements in design and modeling. Nearly all stable designs were destabilized via the single buried Asp substitution: The median drop in stability score for these designs was 1.1, and only 33 buried Asp controls had stability scores greater than 1.0, compared with 271 designs (Fig. 2A, round 2). This substantial destabilization from a single designed substitution provides further large-scale evidence that the stable designs fold into their designed structures. We purified and characterized seven second-generation proteins by SEC and CD, all of which (including three $\alpha\beta\beta\alpha$ designs and one $\beta\beta\alpha\beta$ design) were monomeric, displayed their designed secondary structure in CD, and folded cooperatively and reversibly after thermal denaturation (fig. S6B and table S1). Although the $\alpha\beta\beta\alpha$ and $\beta\beta\alpha\beta$ designs were only marginally stable, the second-generation $\beta\alpha\beta\beta$ design EHEE_rd2_0005 is, to our knowledge, the most thermostable minimal protein ever found (lacking disulfides or metal coordination): Its CD spectrum is essentially unchanged at 95°C, and its denaturation midpoint concentration (C_m) is above 5 M guanidine hydrochloride (GuHCl) (fig. S6B).

The amount of buried NPSA was the strongest observed determinant of folding stability for second-generation $\beta\alpha\beta\beta$ designs (Fig. 2E) and continued to show correlation with stability for second-generation $\alpha\alpha\alpha$ designs (Fig. 2D). The success rate for $\alpha\alpha\alpha$ designs improved in round

2 at all levels of buried NPSA (compare Fig. 2D with Fig. 2B), indicating that improvement of design properties unrelated to buried NPSA (mainly local sequence structure compatibility) contributed to the increase in success rate along with the increase in NPSA. This also illustrates the coupling between different contributions to stability. Although analyzing single terms makes it possible to identify key problems with the design procedure and imbalances in the energy model, the specific success rates shown in Fig. 2 depend on the overall protein context and are not, on their own, fully general.

To improve the stability of the other two topologies, we built a third generation of designs with even greater buried NPSA, at the cost of increased exposure of hydrophobic surface. This might decrease the solubility of the designs, highlighting one of the limits of our approach aimed at optimizing stability. To increase buried NPSA in the $\beta\alpha\beta\beta$ topology, we expanded the architecture from 41 to 43 residues. This led to a large increase in the $\beta\alpha\beta\beta$ success rate (~0% to 13%; Fig. 2H) and 236 newly discovered stable $\beta\alpha\beta\beta$ designs (Fig. 2A, round 3). We purified four third-generation designs (fig. S6C and table S1) and found the $\beta\alpha\beta\beta$ design EEHEE_rd3_1049 to be very stable (Fig. 3). We determined the solution structure of this design by NMR, revealing that it folds into its designed structure, which is not found in nature at this size range (average backbone RMSD = 1.5 Å; Fig. 3). Buried NPSA remained the dominant determinant of stability within the tested $\beta\alpha\beta\beta$ designs (Fig. 2F). We also observed that a newly improved Rosetta energy function [optimized independently from this work (19)] provided significant discrimination between stable and unstable designs, both for the $\beta\alpha\beta\beta$ topology (Fig. 2G) and for other topologies.

Having accumulated nearly 1000 examples of stable designs from rounds 1 to 3, we asked whether more systematic use of these data could result in the selection of better designs. We designed 2000 to 6000 new proteins per topology (using the improved energy function) and then selected 1000 designs each for experimental testing by ranking the designs using topology-specific linear regression, logistic regression, and gradient-boosting regression models trained on the structural features and experimental stabilities of the 10,000 designs from rounds 1 to 3. Many designs selected for testing were predicted to have a low likelihood of folding but were included to increase sequence diversity and because better designs could not be found (26). Despite this, an even larger fraction of designs proved stable than before; notably, the success rate for $\beta\alpha\beta\beta$ designs increased from 17% to 39%, and the success rate for $\beta\alpha\beta\beta$ designs increased from 13% to 58% (Fig. 2H). Although the success rate for designing the $\alpha\beta\alpha$ topology remained low (as predicted by the models), five purified fourth-generation designs in this topology possessed the highest stability yet observed for the fold by CD (fig. S6D and table S1). We solved the structure of one of these (HEEH_rd4_0097) by NMR and found that it adopts the designed

structure in solution (average backbone RMSD = 1.5 Å; Fig. 3). The overall increase in success across the four rounds (Fig. 2H)—from 200 stable designs in round 1 (nearly all in a single topology) to more than 1800 stable designs in round 4 spread across all four topologies—demonstrates the power of our massively parallel approach to drive systematic improvement in protein design.

Of the models used to rank designs, logistic regression was the most successful and was quite accurate: When designs were binned according to their predicted success probability, the number of successes in each bin was close to that predicted beforehand by the logistic regressions (Fig. 2I and fig. S8A). The accuracy of the regression models demonstrates that large-scale analysis of stable and unstable designed proteins can be used to build predictive models of protein stability. Although the models we built are limited by their training data and not fully general, the inputs to the models were global features of all proteins, such as buried NPSA and total hydrogen bonding energy. This gives these models greater potential for generality than other models used in iterative protein engineering that are typically specific to particular protein families (38, 39), although those approaches have their own advantages. Retrospectively, we found that a single logistic regression trained on data from all topologies from rounds 1 to 3 performed comparably to the topology-specific regressions at ranking round-4 designs within each topology (fig. S8B). Ultimately, continued application of our approach should greatly expand and broaden the available training data, which can be integrated with other sources of physical, chemical, and biological information (19, 40) to build a new generation of general-purpose protein energy functions (22).

Sequence determinants of stability

We next examined determinants of stability at the individual-residue level by constructing a library containing every possible point mutant of 14 designs, as well as every point mutant in three paradigm proteins from decades of folding research: villin HP35, Pin1 WW domain, and hYAP65 WW domain L30K (Leu³⁰ → Lys) mutant. This library of 12,834 point mutants is comparable in size to the 12,561 single mutants found in the entire ProTherm database (41) and is unbiased toward specific mutations. We assayed this library for stability using trypsin and chymotrypsin, and determined an overall stability effect for each mutation by using the independent results from each protease to account for dynamic range of the assay (fig. S9) (26). The mutational effects were qualitatively consistent with the designed structures for 13 of 14 designs (fig. S10, A to N). As expected, the positions on the designs that were most sensitive to mutation were the core hydrophobic residues, including many Ala residues, which indicates that the designed cores are tightly packed (Fig. 4A and fig. S10, A to N). Mutations to surface residues had much smaller effects, highlighting the potential of these proteins

as stable scaffolds whose surfaces can be engineered for diverse applications.

To examine the mutability of protein surfaces in greater detail and to probe more subtle contributions to stability, we divided the 260 surface positions in 12 of the designs into categories based on secondary structure and calculated the average stability effect of each amino acid for each category using the ~5000 stability measurements at these positions (Fig. 4, E to L) (26). We observed specific, although weak, amino acid preferences within helices (Fig. 4E), helix N-caps (Fig. 4F), the first and last turns of helices (Fig. 4, G and H), middle strands and edge strands (Fig. 4, I and J), and loop residues (Fig. 4, K and L). Asp, Ser, Thr, and Asn were favorable for capping helices, but were, except for Asn, as unfavorable as Gly when inside helices (Fig. 4, E and F). Hydrophobic side chains were stabilizing even when located on the solvent-facing side of a β sheet, and this effect was stronger at middle strand positions than at edge strand positions (Fig. 4, I and J). Most notably, we observed stabilization from charged amino acids on the first and last turns of α helices when these charges counteracted the C-to-N negative-to-positive helical dipole; charges that enhanced the dipole were destabilizing (42). We isolated this effect by comparing the average stability of each amino acid on the first and last helical turns with the average stability of each amino acid at all helical sites (polar sites only in both cases; Fig. 4, G and H). The effect remained significant even when we restricted the analysis to positions that were Arg or Lys in the original designs to control for any bias in the designed structures favoring original, designed residues over mutant residues, although no significant effect was seen at Glu positions (fig. S11). We had not examined agreement with this dipolar preference during the four rounds of design, and after this observation, we found that the net favorable charge on the first and last helical turns (stabilizing charges minus destabilizing charges summed over all helices) discriminated between stable and unstable fourth-generation $\alpha\alpha\alpha$ designs better than any other metric we examined, explaining in part why the success rate had not reached 100%.

In the three naturally occurring proteins, mutations at conserved positions were generally destabilizing, although each natural protein possessed several highly conserved positions that we experimentally determined to be unimportant or deleterious to stability. In villin HP35, these were Trp⁶⁴, Lys⁷⁰, Lys⁷⁵, and Phe⁷⁶ (villin HP35 consists of residues 42 to 76), which are required for villin to bind F-actin (Fig. 4B and fig. S12) (43, 44). In Pin1, the highly conserved Ser¹⁶ is deleterious for stability but directly contacts the phosphate on phosphopeptide ligands of Pin1 (45), highlighting a stability/function trade-off in Pin1 (6, 46) discoverable without directly assaying function (Fig. 4C and fig. S12) (46). In hYAP65, the conserved residues His³², Thr³⁷, and Trp³⁹ are relatively unimportant for stability, but these residues form the peptide recognition pocket in

YAP-family WW domains (Fig. 4D and fig. S12) (47, 48). These examples illustrate how our approach enables high-throughput identification of functional residues, even without a functional assay or a protein structure [as in computational approaches (49)], via comparison between stability data and residue conservation.

Stability measurement of all known small protein domains

How stable are these designed proteins relative to naturally occurring proteins? To examine this, we synthesized DNA encoding (i) all 472 sequences in the Protein Data Bank (PDB) between 20 and 50 residues in length and containing only the 19 non-Cys amino acids, and (ii) one repre-

sentative for all 706 domains meeting these criteria in the Pfam protein family database. These DNA sequences were prepared by reverse translation in an identical manner to the designs (26). We included this DNA (and DNA for all stable designs from rounds 1 to 3) in the library containing our fourth-generation designs to facilitate a head-to-head comparison. The large majority of these natural proteins successfully displayed on yeast (92% each for PDB and Pfam sequences), which was comparable to the fourth-generation buried Asp mutants (also 92%) but lower than fourth-generation scrambled sequences (96%) and fourth-generation designs (99%). The most resistant overall sequence (measured by stability score) was a C-terminal coiled-coil domain from a TRP

channel (3HRO, stability score 1.93). This protein is likely stabilized by intersubunit interactions made possible by assembly on the yeast surface (50). Of the 100 unique, monomeric sequences with PDB structures, the most protease-resistant was a peripheral subunit binding domain (*aaa* topology) from the thermophile *Bacillus stearothermophilus* (2PDD, stability score 1.48), which has been studied as an ultrafast-folding protein (4, 8). A total of 774 designed proteins had higher stability scores than this most protease-resistant natural monomeric protein. As illustrated in Fig. 5, the number of stable proteins we discovered exceeds the number of natural proteins in the PDB (monomeric or not) in this size range by a factor of 50.

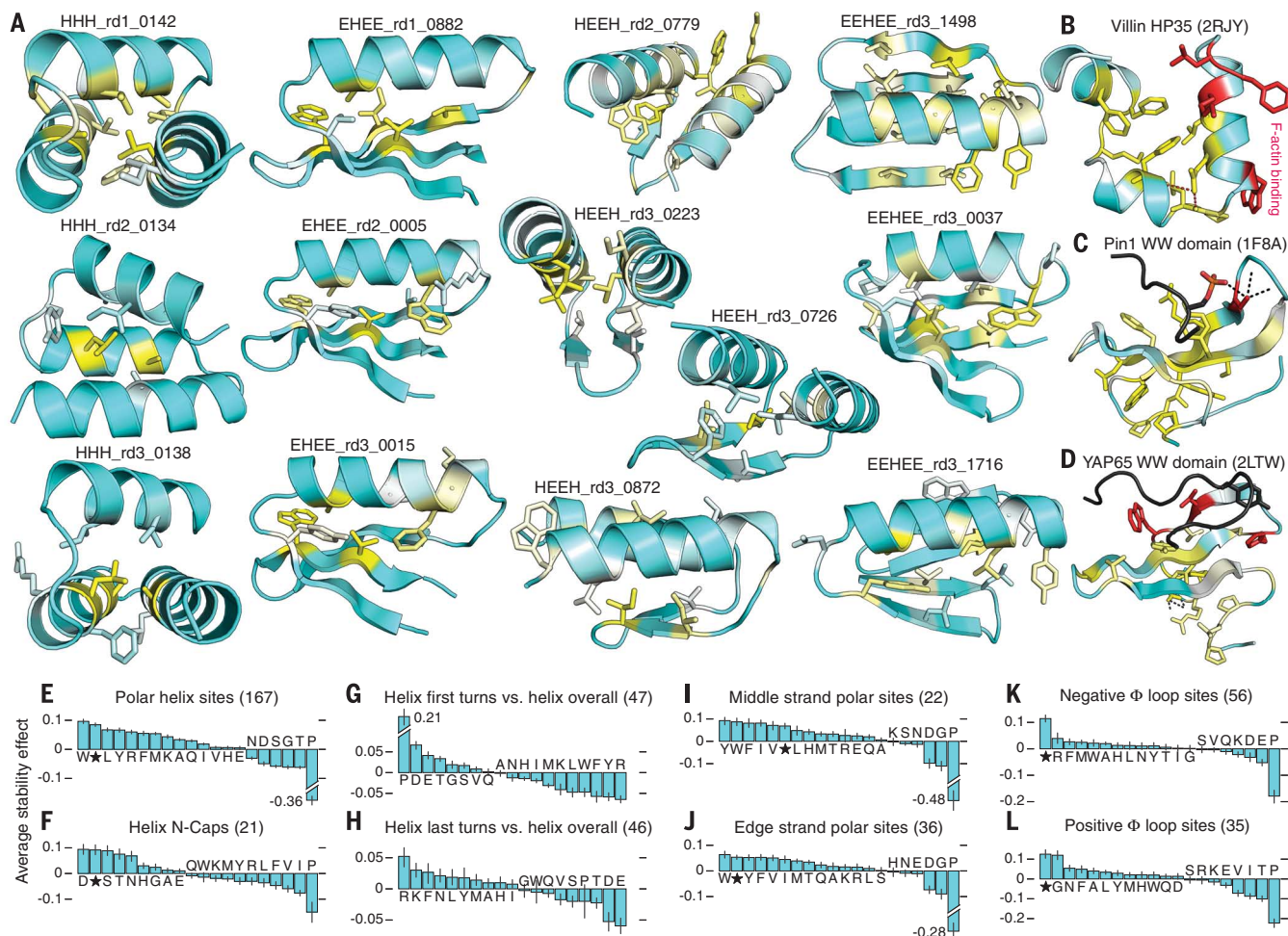


Fig. 4. Comprehensive mutational analysis of stability in designed and natural proteins. (A) Average change in stability due to mutating each position in 13 designed proteins, depicted on the design model structures.

Positions where mutations are most destabilizing are colored yellow and shown in stick representation; positions where mutations have little effect are colored blue. Each protein's color scale is different to emphasize the relative importance of positions; see fig. S10 for full data for all proteins. (B to D) As in (A) for native proteins, with conserved residues not contributing to stability colored red. (B) Villin HP35. In red, Trp⁶⁴, Lys⁷⁰, Lys⁷⁵, and Phe⁷⁶ (HP35 consists of residues 42 to 76) have little effect on stability but are conserved for function (F-actin binding). (C) Pin1 WW domain, shown bound to a doubly phosphorylated peptide. In red, Ser¹⁶ is conserved

and critical for function but is destabilizing relative to mutations at that position. (D) hYAP65 L30K, shown bound to a Smad7-derived peptide. In red, His³², Thr³⁷, and Trp³⁹ form the peptide recognition motif and are conserved but unimportant for stability. (E to L) Average stability effect of each amino acid at different categories of surface positions, in units of stability score (positive = stabilizing, negative = destabilizing). The average stability of all amino acids in each panel was set to zero. The number of individual positions examined in each category is listed in parentheses with the category name. The average stability effect of the original "wild-type" designed residue (unique to each particular site within a category) is shown by a black star. Error bars indicate the 50% confidence interval for the average stability effect, calculated using bootstrapping. See (26) for a full description of the analysis.

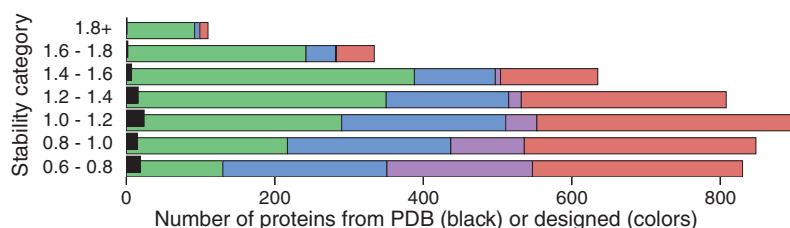


Fig. 5. Comparison of naturally occurring and designed protein stability. Designed and naturally occurring proteins are separated into bins by stability score (y axis). The total number of designed proteins in each bin is shown by the colored bar, subdivided by topology from left to right as follows: $\alpha\alpha\alpha$ (green), $\beta\alpha\beta\beta$ (blue), $\alpha\beta\beta\alpha$ (violet), $\beta\beta\alpha\beta\beta$ (red). The total number of naturally occurring proteins with PDB structures (lacking disulfides) in each bin is shown by a black bar.

Conclusion

We have shown that proteins can be computationally designed and assayed for folding thousands at a time, and that high-throughput design experiments can provide quantitative insights into the determinants of protein stability. Large libraries can be designed in a relatively unbiased manner (as in our first generation) to maximize the protein property space examined, or properties can be tuned to increase the design success rate at the cost of diversity. The power of our iterative learning approach to progressively home in on more subtle contributions to stability is highlighted by the progression of our $\alpha\alpha\alpha$ design sets from early rounds, in which design failures were caused by insufficient buried nonpolar surface area, to the last round, where helix-side chain electrostatics had the greater effect. The large numbers of folded and not-folded designs will also provide stringent tests of molecular dynamics simulation approaches that have successfully reproduced structures (13, 15) and some thermodynamic measurements (14, 51) of natural proteins, but have not yet been challenged with plausible but unstable protein structures like our design failures.

The four solution structures, saturation mutagenesis data on 13 of 14 designs, and more than 30,000 negative control experiments indicate that the large majority of our stable sequences are structured as designed. These 2788 designed proteins, stable without disulfides or metal coordination, should have numerous applications in bioengineering and synthetic biology. Many are more stable than any comparably sized monomeric proteins found in the PDB, making them ideal scaffolds for engineering inhibitors of intracellular protein-protein interactions. Their small size may also help to promote membrane translocation and endosomal escape (52, 53). As DNA synthesis technology continues to improve, high-throughput protein design will become possible for larger proteins as well, revealing determinants of protein stability in more complex structures. We have entered a new era of iterative, data-driven de novo protein design and modeling.

REFERENCES AND NOTES

- K. A. Dill, *Biochemistry* **29**, 7133–7155 (1990).
- A. D. Robertson, K. P. Murphy, *Chem. Rev.* **97**, 1251–1268 (1997).
- C. N. Pace, J. M. Scholtz, G. R. Grimsley, *FEBS Lett.* **588**, 2177–2184 (2014).
- H. Gelman, M. Gruebele, *Q. Rev. Biophys.* **47**, 95–142 (2014).
- X. Jiang, J. Kowalski, J. W. Kelly, *Protein Sci.* **10**, 1454–1465 (2001).
- M. Jäger *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **103**, 10648–10653 (2006).
- S. Xiao, Y. Bi, B. Shan, D. P. Raleigh, *Biochemistry* **48**, 4607–4616 (2009).
- H. Neuweiler *et al.*, *J. Mol. Biol.* **390**, 1060–1073 (2009).
- C. N. Pace *et al.*, *J. Mol. Biol.* **408**, 514–528 (2011).
- C. L. Araya *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **109**, 16858–16863 (2012).
- S. Xiao *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **110**, 11337–11342 (2013).
- C. N. Pace *et al.*, *Protein Sci.* **23**, 652–661 (2014).
- K. Lindorff-Larsen, S. Piana, R. O. Dror, D. E. Shaw, *Science* **334**, 517–520 (2011).
- S. Piana, K. Lindorff-Larsen, D. E. Shaw, *Proc. Natl. Acad. Sci. U.S.A.* **109**, 17845–17850 (2012).
- H. Nguyen, J. Maier, H. Huang, V. Perrone, C. Simmerling, *J. Am. Chem. Soc.* **136**, 13959–13962 (2014).
- P.-S. Huang, S. E. Boyken, D. Baker, *Nature* **537**, 320–327 (2016).
- C. A. Rohl, C. E. M. Strauss, K. M. S. Misura, D. Baker, *Methods Enzymol.* **383**, 66–93 (2004).
- T. J. Magliery, *Curr. Opin. Struct. Biol.* **33**, 161–168 (2015).
- H. Park *et al.*, *J. Chem. Theory Comput.* **12**, 6201–6212 (2016).
- B. I. Dahiyat, S. L. Mayo, *Science* **278**, 82–87 (1997).
- H. Liang *et al.*, *Angew. Chem. Int. Ed.* **48**, 3301–3303 (2009).
- Z. Li, Y. Yang, J. Zhan, L. Dai, Y. Zhou, *Annu. Rev. Biophys.* **42**, 315–335 (2013).
- S. Kosuri, G. M. Church, *Nat. Methods* **11**, 499–507 (2014).
- M. G. F. Sun, M.-H. Seo, S. Nim, C. Corbi-Verge, P. M. Kim, *Sci. Adv.* **2**, e1600692 (2016).
- E. T. Boder, K. D. Wittrup, *Nat. Biotechnol.* **15**, 553–557 (1997).
- See supplementary materials.
- V. Sieber, A. Plückthun, F. X. Schmid, *Nat. Biotechnol.* **16**, 955–960 (1998).
- M. D. Finucane, M. Tuna, J. H. Lees, D. N. Woolfson, *Biochemistry* **38**, 11604–11612 (1999).
- C. Park, S. Zhou, J. Gilmore, S. Marqusee, *J. Mol. Biol.* **368**, 1426–1437 (2007).
- C. Park, S. Marqusee, *Nat. Methods* **2**, 207–212 (2005).
- P. Leuenberger *et al.*, *Science* **355**, eaai7825 (2017).
- M. Jäger, M. Dendle, J. W. Kelly, *Protein Sci.* **18**, 1806–1813 (2009).
- G. Bhardwaj *et al.*, *Nature* **538**, 329–335 (2016).
- N. Koga *et al.*, *Nature* **491**, 222–227 (2012).
- S. Kamtekar, J. M. Schiffer, H. Xiong, J. M. Babik, M. H. Hecht, *Science* **262**, 1680–1685 (1993).
- A. R. Davidson, R. T. Sauer, *Proc. Natl. Acad. Sci. U.S.A.* **91**, 2146–2150 (1994).
- M. H. Hecht, A. Das, A. Go, L. H. Bradley, Y. Wei, *Protein Sci.* **13**, 1711–1723 (2004).

- R. J. Fox *et al.*, *Nat. Biotechnol.* **25**, 338–344 (2007).
- P. A. Romero, A. Krause, F. H. Arnold, *Proc. Natl. Acad. Sci. U.S.A.* **110**, E193–E201 (2013).
- A. Leaver-Fay *et al.*, *Methods Enzymol.* **523**, 109–143 (2013).
- M. D. S. Kumar *et al.*, *Nucleic Acids Res.* **34**, D204–D206 (2006).
- E. G. Baker *et al.*, *Nat. Chem. Biol.* **11**, 221–228 (2015).
- D. S. Doering, P. Matsudaira, *Biochemistry* **35**, 12677–12685 (1996).
- J. Meng *et al.*, *Biochemistry* **44**, 11963–11973 (2005).
- M. A. Verdecia, M. E. Bowman, K. P. Lu, T. Hunter, J. P. Noel, *Nat. Struct. Biol.* **7**, 639–643 (2000).
- B. K. Shoichet, W. A. Baase, R. Kuroki, B. W. Matthews, *Proc. Natl. Acad. Sci. U.S.A.* **92**, 452–456 (1995).
- P. A. Chong, H. Lin, J. L. Wrana, J. D. Forman-Kay, *J. Biol. Chem.* **281**, 17069–17075 (2006).
- E. Aragón *et al.*, *Structure* **20**, 1726–1736 (2012).
- A. H. Elcock, *J. Mol. Biol.* **312**, 885–896 (2001).
- E. T. Boder, J. R. Bill, A. W. Nields, P. C. Marrack, J. W. Kappler, *Biotechnol. Bioeng.* **92**, 485–491 (2005).
- S. Piana, J. L. Klepeis, D. E. Shaw, *Curr. Opin. Struct. Biol.* **24**, 98–105 (2014).
- J. S. Appelbaum *et al.*, *Chem. Biol.* **19**, 819–830 (2012).
- J. R. LaRochelle, G. B. Cobb, A. Steinauer, E. Rhoades, A. Schepartz, *J. Am. Chem. Soc.* **137**, 2536–2541 (2015).

ACKNOWLEDGMENTS

Supported by the Howard Hughes Medical Institute (D.B.) and the Natural Sciences and Engineering Research Council of Canada (C.H.A.). G.J.R. is a Merck Fellow of the Life Sciences Research Foundation. C.H.A. holds a Canada Research Chair in Structural Genomics. We thank S. Rettie for mass spectrometry support; C. Lee for deep sequencing support; S. Ovchinnikov for assistance quantifying sequence conservation; V. Nguyen, A. Yehdego, T. Howard, and K. Lau for assistance with protein purification; and H. Gelman and many other members of the Baker lab for helpful discussions. This work was facilitated by the Hyak supercomputer at the University of Washington and by donations of computing time from Rosetta@Home participants. The Structural Genomics Consortium is a registered charity (number 1097737) that receives funds from AbbVie; Bayer Pharma AG; Boehringer Ingelheim; Canada Foundation for Innovation; Eshelman Institute for Innovation; Genome Canada through Ontario Genomics Institute grant OGI-055; Innovative Medicines Initiative (EU/EFPIA) through ULTRA-DD grant 115766; Janssen Pharmaceuticals; Merck & Co.; Novartis Pharma AG; Ontario Ministry of Research, Innovation and Science (MRIS); Pfizer; São Paulo Research Foundation-FAPESP; Takeda; and the Wellcome Trust. The RosettaScripts code and blueprint files used for protein design are provided in the supplementary materials. The data for this work (designed sequences and structures, deep sequencing counts, EC₅₀ values, stability scores, and structural analysis of the designed models) are also provided in supplementary materials. The python code for inferring EC₅₀ values and for fitting the unfolded state model is provided at https://github.com/asford/protease_experimental_analysis. G.J.R. and D.B. are inventors on provisional patent application no. 62/491,518 filed 28 April 2017 by the University of Washington that covers (i) the method described in this work for computationally designing and experimentally verifying stable miniproteins, and (ii) the 4000 most stable protein sequences designed in the work. Author contributions: G.J.R. designed the research, the experimental approach, and the proteins; G.J.R., T.M.C., I.G., S.H., L.C., R.R., and A.C. performed experiments; all authors analyzed data; G.J.R., A.F., and V.K.M. contributed new computational tools; C.H.A. and D.B. supervised research; and G.J.R. and D.B. wrote the manuscript.

SUPPLEMENTARY MATERIALS

www.sciencemag.org/content/357/6347/168/suppl/DC1
Materials and Methods
Supplementary Text
Figs. S1 to S12
Tables S1 to S3
References (54–88)

28 February 2017; accepted 9 June 2017
10.1126/science.aan0693

Global analysis of protein folding using massively parallel design, synthesis, and testing

Gabriel J. Rocklin, Tamuka M. Chidyausiku, Inna Goreshnik, Alex Ford, Scott Houliston, Alexander Lemak, Lauren Carter, Rashmi Ravichandran, Vikram K. Mulligan, Aaron Chevalier, Cheryl H. Arrowsmith and David Baker

Science **357** (6347), 168-175.
DOI: 10.1126/science.aan0693

Exploring structure space to understand stability

Understanding the determinants of protein stability is challenging because native proteins have conformations that are optimized for function. Proteins designed without functional bias could give insight into how structure determines stability, but this requires a large sample size. Rocklin *et al.* report a high-throughput protein design and characterization method that allows them to measure thousands of miniproteins (see the Perspective by Woolfson *et al.*). Iterative rounds of design and characterization increased the design success rate from 6 to 47%, which provides insight into the balance of forces that determine protein stability.

Science, this issue p. 168; see also p. 133

ARTICLE TOOLS

<http://science.sciencemag.org/content/357/6347/168>

SUPPLEMENTARY MATERIALS

<http://science.sciencemag.org/content/suppl/2017/07/12/357.6347.168.DC1>

RELATED CONTENT

<http://science.sciencemag.org/content/sci/357/6347/133.full>

REFERENCES

This article cites 86 articles, 15 of which you can access for free
<http://science.sciencemag.org/content/357/6347/168#BIBL>

PERMISSIONS

<http://www.sciencemag.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of Service](#)