

Dissecting the stability determinants of a challenging de novo protein fold using massively parallel design and experimentation

Tae-Eun Kim^{1,2,3,9}, Kotaro Tsuboyama^{2,3,9}, Scott Houliston^{4,5}, Cydney M. Martell^{1,2,3,6}, Claire M. Phoumyvong^{1,2,3}, Hugh K. Haddox⁷, Cheryl H. Arrowsmith^{4,5}, Gabriel J. Rocklin^{2,3,6,8}

¹ Driskill Graduate Program in Life Sciences, Northwestern University Feinberg School of Medicine, Chicago, IL 60611

² Department of Pharmacology, Northwestern University Feinberg School of Medicine, Chicago, IL 60611

³ Center for Synthetic Biology, Northwestern University, Evanston, IL 60208

⁴ Structural Genomics Consortium, University of Toronto, Toronto, ON, M5G 1L7, Canada.

⁵ Princess Margaret Cancer Centre and Department of Medical Biophysics, University of Toronto, Toronto, ON, M5G 2M9, Canada.

⁶ Center for Life Processes Institute, Northwestern University, Evanston, IL 60208

⁷ Department of Biochemistry & Institute for Protein Design, University of Washington, Seattle, WA 98195

⁸ Robert H. Lurie Comprehensive Cancer Center, Northwestern University, Chicago, IL 60611

⁹ These authors contributed equally to this work

*To whom correspondence should be addressed:

Gabriel J. Rocklin

303 E Superior St

Simpson Querrey Building 11-517

Chicago, IL 60611

grocklin@gmail.com

Author Contributions: TEK, KT, and GJR designed the research, TEK, KT, SH, CM, CP performed research, TEK, KT, SH, and GJR analyzed the data, HKH developed computational tools, GJR and CHA supervised the research, and TEK and GJR wrote the paper. TEK and KT contributed equally to this work.

Competing interest: The authors declare no competing interests.

Abstract

Designing entirely new protein structures remains challenging because we do not fully understand the biophysical determinants of folding stability. Yet some protein folds are easier to design than others. Previous work identified the 43-residue $\alpha\beta\beta\alpha$ fold as especially challenging: the best designs had only a 2% success rate, compared to 39-87% success for other simple folds (1). This suggested the $\alpha\beta\beta\alpha$ fold would be a useful model system for gaining a deeper understanding of folding stability determinants and for testing new protein design methods. Here, we designed over ten thousand new $\alpha\beta\beta\alpha$ proteins and found over three thousand of them to fold into stable structures using a high-throughput protease-based assay. Nuclear magnetic resonance, hydrogen-deuterium exchange, circular dichroism, deep mutational scanning, and scrambled sequence control experiments indicated that our stable designs fold into their designed $\alpha\beta\beta\alpha$ structures with exceptional stability for their small size. Our large dataset enabled us to quantify the influence of universal stability determinants including nonpolar burial, helix capping, and buried unsatisfied polar atoms, as well as stability determinants unique to the $\alpha\beta\beta\alpha$ topology. Our work demonstrates how large-scale design and test cycles can solve challenging design problems while illuminating the biophysical determinants of folding.

Significance

Most computationally designed proteins fail to fold into their designed structures. This low success rate is a major obstacle to expanding the applications of protein design. In previous work, we discovered a small protein fold that was paradoxically challenging to design (only a 2% success rate) even though the fold itself is very simple. Here, we used a recently developed high-throughput approach to comprehensively examine the design rules for this simple fold. By designing over ten thousand proteins and experimentally measuring their folding stability, we discovered the key biophysical properties that determine the stability of these designs. Our results illustrate general lessons for protein design and also demonstrate how high-throughput stability studies can quantify the importance of different biophysical forces.

Introduction

Improving our understanding of the determinants of protein stability (2–4) would accelerate biological, biomedical, and biotechnology research. In particular, computational models of protein stability are commonly used for a range of applications, including protein design (5–7), stabilizing naturally occurring proteins (8, 9), and predicting the effects of point mutants (10–12). However, all of these models have important limitations. For example, most computationally designed proteins made by experts fail to fold and function (13–15). Non-experts avoid computational design techniques because they are not reliable. These challenges stem from our incomplete understanding of the biophysical determinants of folding stability, and therefore, improving our quantitative understanding of stability should improve the success and applicability of design.

Recently, we introduced a high-throughput approach to study protein folding stability that is particularly helpful for improving computational modeling and design. In our approach, we designed thousands of *de novo* proteins and measured their folding stabilities using a yeast display-based proteolysis assay coupled to next-generation sequencing (1). Several new studies have applied our methodology (16–19) as it has several advantages. First, measuring folding stability for thousands of proteins makes it possible to statistically quantify biophysical features that contribute to stability. Second, examining diverse sequences makes it easier to derive principles that are not specific to a particular protein context. Finally, assaying computationally designed proteins focuses the experimentation on the regions of sequence and structural space that are predicted to be low energy according to a particular computational model, which is especially useful for improving that model.

We previously used this approach to increase the success rate (i.e. fraction of designs that form stable, folded structures) of *de novo* miniprotein designs from 6% to 47% (1). Three different protein topologies could be designed very robustly (39-87% success), but a fourth topology ($\alpha\beta\beta\alpha$, 43 residues) proved very challenging. Only 2% of $\alpha\beta\beta\alpha$ designs folded into stable structures despite the simplicity of the structure and four repeated efforts to improve the design procedure (Fig. 1A). This suggested that our design procedure and stability model were missing something fundamental about the $\alpha\beta\beta\alpha$ topology, and that this particular fold could be a useful model system for building a deeper understanding of folding stability. Here, we investigated this by asking two main questions. First, how can we improve our design procedure to obtain a large number of stable $\alpha\beta\beta\alpha$ proteins for further analysis? Notably, there are no naturally occurring examples of the 43-residue $\alpha\beta\beta\alpha$ fold for us to learn from, although this architecture is similar to the unusual 55-residue $\alpha\beta\beta\alpha$ fold of the gpW protein from bacteriophage lambda (20). Second, how do the biophysical and topological features of different $\alpha\beta\beta\alpha$ designs combine to determine each protein's folding stability? We investigated these questions by designing and experimentally testing over ten thousand new $\alpha\beta\beta\alpha$ miniproteins using our high-throughput approach. We also examined whether the structure prediction model AlphaFold 2 (21) could be applied to differentiate stable and unstable designs.

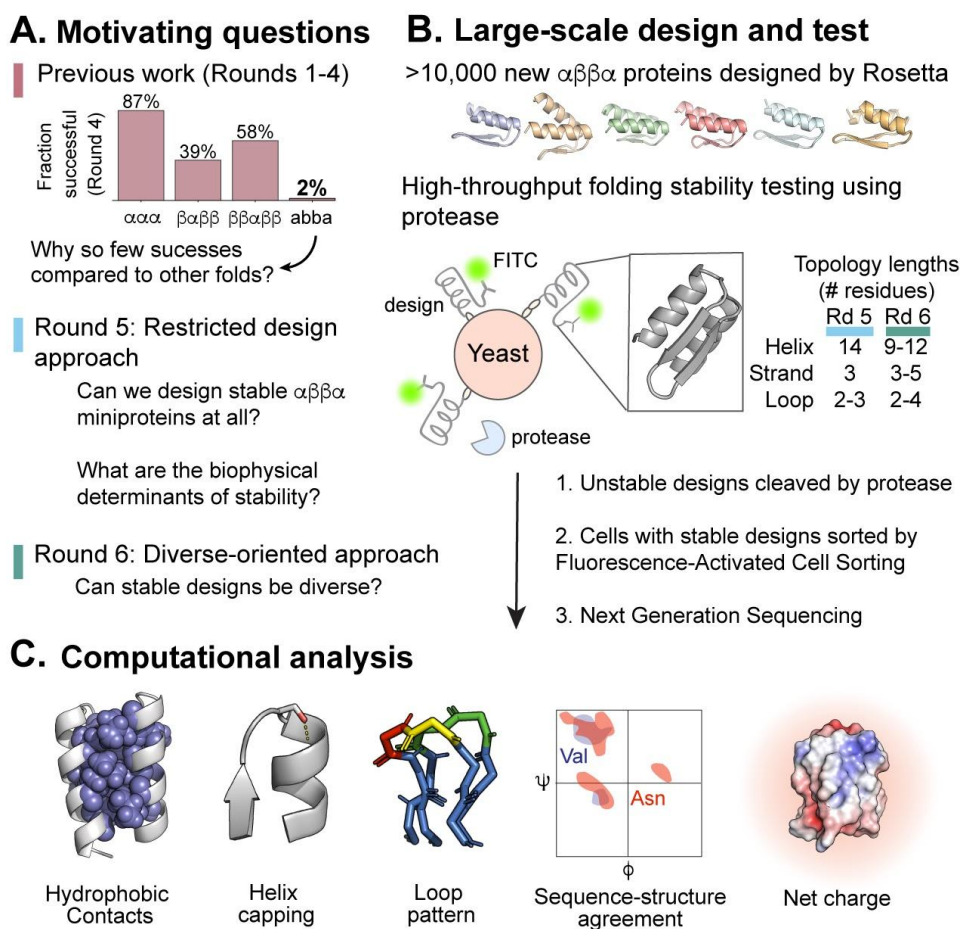


Fig. 1. Design strategy for generating and testing $\alpha\beta\beta\alpha$ miniproteins. (A) Previously, we performed four iterative design-test-analysis cycles to generate stable $\alpha\beta\beta\alpha$ miniprotein designs, but only achieved a 2% success rate (1). (B) Here, we designed thousands of new $\alpha\beta\beta\alpha$ miniproteins using Rosetta and experimentally tested them for their folding stability using a combined yeast display and protease sensitivity assay. (C) We then performed computational analysis to identify and understand the relative importance of key stability determinants (e.g. hydrophobic contacts, helix capping, loop patterning, local sequence-structure agreement, and net charge).

Results

Designing $\alpha\beta\beta\alpha$ miniproteins using a restricted design strategy

We first computationally designed thousands of new $\alpha\beta\beta\alpha$ miniproteins (“Round 5”) based on lessons learned from our previous four rounds of design (1). All designs were based on a single protein architecture (22) that previously led to the greatest number of stable designs (Fig. S1A). This architecture restricted our new $\alpha\beta\beta\alpha$ miniproteins to 14-residue α -helices, 3-residue β -strands, and a specific loop structure (Fig. 1B). In addition, we ensured our designs met strict criteria for buried nonpolar surface area, Rosetta energy, and predicted secondary structure (Fig. S1B). Finally, we required the middle loop to have a hydrophobic residue, required solvent-facing residues on the β -strands to be polar or charged, set a minimum threshold for the total number of hydrophobic residues, and eliminated Gly, Thr, and Val in helices (Fig. S1C-D) (see Methods). We expected these restrictions to increase the success of our new designs, although they would reduce the potential sequence and structural diversity.

Based on this “restricted” design strategy, we generated 28,000 $\alpha\beta\beta\alpha$ miniproteins using an improved version of the Rosetta score function. This score function was previously parameterized to correlate with our earlier high-throughput data on miniprotein folding stability (23). In addition, we used an improved sequence sampling procedure that minimizes over-compaction and produces more native-like protein cores containing bulky residues (24). Our final set of 6,000 $\alpha\beta\beta\alpha$ designs were chosen by ranking all 28,000 $\alpha\beta\beta\alpha$ designs using a linear regression model trained on previous large-scale $\alpha\beta\beta\alpha$ stability data (1). After we ranked our designs, we eliminated designs that were more than 31/43 residues identical to a higher-ranking design. Each design based on this restricted strategy is named HEEH_TK_rd5_#####, where HEEH indicates the pattern of α -helices (H) and β -strands (E), TK indicates the designer (author TEK), rd5 indicates these designs follow our four previous efforts (1), and ##### is the design number.

Biophysical characterization of $\alpha\beta\beta\alpha$ miniproteins using a restricted design strategy

We measured the folding stabilities of our newly designed $\alpha\beta\beta\alpha$ miniproteins using the high-throughput protease sensitivity assay introduced previously (Fig. 1B) (1). Briefly, all sequences were synthesized as DNA oligonucleotides in a pooled library. We then used *S. cerevisiae* to express and display our sequences on their cell surface, along with a c-terminal myc tag. Next, we subjected the yeast cells to varying concentrations of trypsin and chymotrypsin (tested separately) (Fig. S2A-B) and fluorescently labeled the cells displaying protease-resistant sequences. Finally, we sorted the fluorescently labeled cells by flow cytometry and identified the protease-resistant sequences by deep sequencing (Fig. 1B). As previously, we assigned each design a “stability score”, defined as the difference between that sequence’s observed protease sensitivity and the predicted sensitivity of that sequence in its unfolded state. Each one-unit increase in stability score indicates a 10-fold higher amount of protease required to cleave that sequence under assay conditions, compared with the predicted protease concentration required to cleave that sequence in its unfolded state (1). To conservatively identify stable designs, each design’s overall stability score is the minimum of the stability scores observed separately with trypsin and chymotrypsin. We previously observed that sequences of scrambled amino acids (not designed sequences) rarely have stability scores above 1, and so we classify designs as stable when their stability score exceeds 1.

We found 38% of our new $\alpha\beta\beta\alpha$ designs to be stable with a mean stability score of 0.81 (Fig. 2A). This greatly exceeded our previous success rate of 2% (Fig. 1A) (1). We also included control sequences in our library whose residue compositions matched our $\alpha\beta\beta\alpha$ designs, but with the ordering of the residues scrambled in a specific manner: polar residues remained polar, nonpolar residues remained nonpolar, and proline and glycine residues remained in their identical positions. In contrast to our designs, almost all scrambled sequences had stability scores < 1 with a mean stability score of -0.86 (Fig. 2A). This suggests that the protease resistance observed for a subset of designs can be attributed to the folding stability of their designed structures, rather than generic properties of their sequences such as residue composition or patterning. In addition, stability scores measured using trypsin and chymotrypsin were correlated with each other despite the differing specificities of the proteases (Fig. S2A-B). This further indicates that our measured stability scores reflect folding stability rather than protease-specific factors.

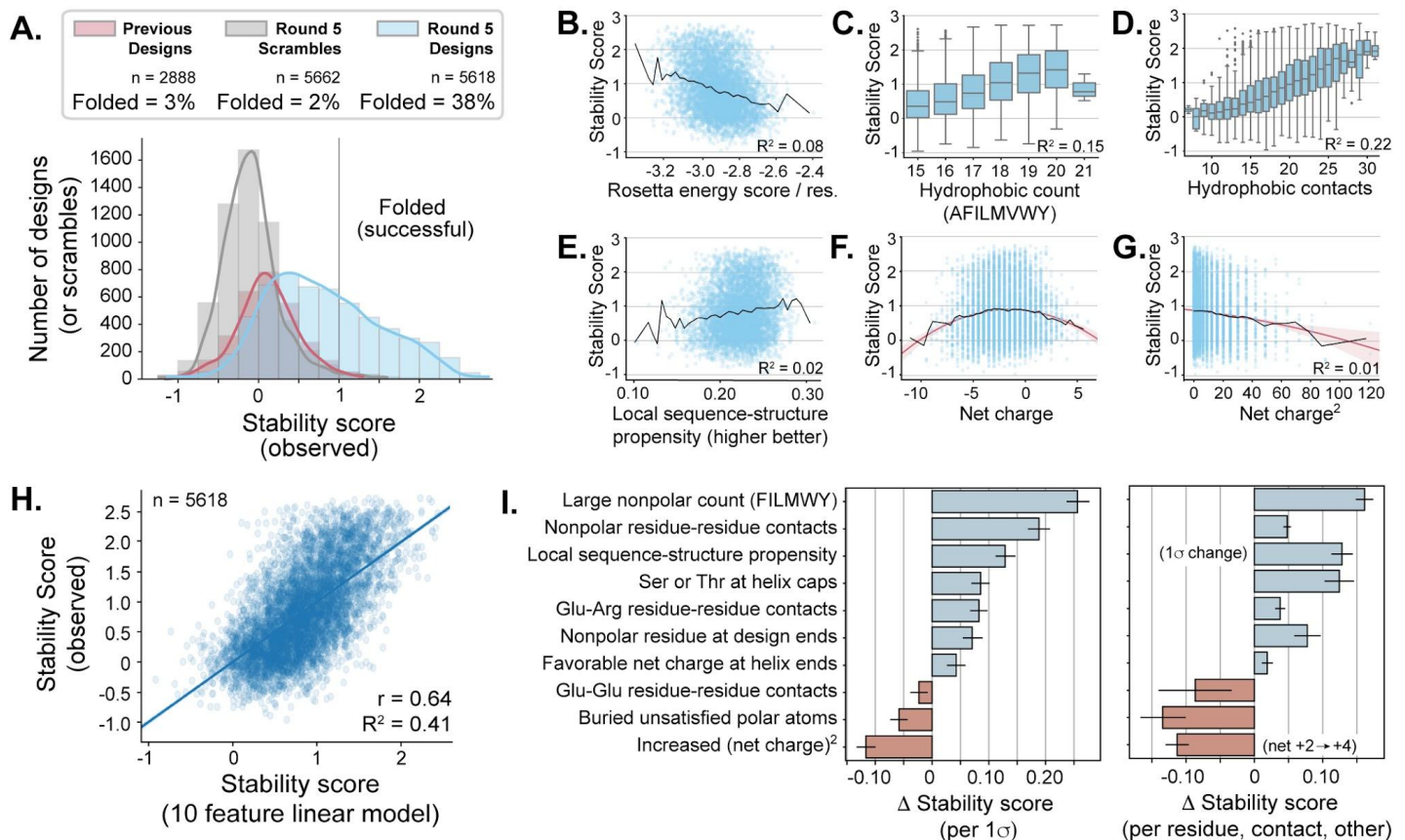


Fig. 2. Experimental testing and analysis of $\alpha\beta\alpha$ stability determinants from a restricted design strategy. (A) The stability score distributions of designed $\alpha\beta\alpha$ miniproteins (blue), scrambled sequences (gray) and previously published $\alpha\beta\alpha$ miniproteins (red) (1); the vertical line at stability score = 1 denotes the threshold above which we consider a design to be stable. (B-G) The relation between Individual protein features and stability score. For Rosetta energy, lower values indicate favorable energies, and for local sequence-structure propensity, higher values indicate favorable propensity. Black lines show moving averages; red lines show fits to quadratic (F) and linear (G) models. (H) A ten-feature linear regression model was built using normalized data, and the experimental stability scores are compared to the model's predicted stability scores. (I) The magnitudes of the coefficients from the model based on their importance in the dataset (left) and their biophysical strength (right). Error bars indicate 95% confidence intervals from bootstrapping.

We next sought to verify that stable $\alpha\beta\alpha$ miniproteins folded as designed using several orthogonal approaches. We first selected six stable $\alpha\beta\alpha$ designs with varying hydrophobicity values (25) and individually purified them from *E. coli* (Fig. 3A, Table S1). Size-exclusion chromatography indicated a mixture of monomeric and oligomeric species for each design, and the monomeric fraction was selected for analysis by circular dichroism (CD) spectroscopy. All designs exhibited helical secondary structure and reversible folding after heating to 95°C (Fig. 3B). None of the designs showed a clear melting transition, although designs HEEH_TK_rd5_0958 and HEEH_TK_rd5_3711 lost much of their helical character at 95°C. In contrast, design HEEH_TK_rd5_0420 was minimally perturbed during melting (Fig. 3C), indicating extreme thermostability.

Next, to spot-check the accuracy of our designed structures, we solved the structure of design HEEH_TK_rd5_0341 by nuclear magnetic resonance (NMR) (Table S2). The NMR ensemble revealed that the design folds into the expected $\alpha\beta\alpha$ structure with slightly more space between the α -helices and the β -hairpin

compared to the design model (average backbone root mean squared deviation (RMSD) = 2.25 Å) (Fig. 3E). As a comparison, the AlphaFold 2 (21) predicted structure has an RMSD of 1.16 Å to the Rosetta model, and an average RMSD of 1.58 Å to the NMR ensemble (Fig. 3E).

We also examined the local stability of design HEEH_TK_rd5_0341 by hydrogen deuterium exchange (HDX) NMR. The HDX opening free energies revealed differences in local stability in different regions of the topology. The most stable secondary structure was Helix 2, with opening energies around 4 kcal/mol at 15°C. The highest opening energy was 4.5 kcal/mol, observed at I36 (Fig. 3F). This highest opening energy typically indicates the global stability of the protein (26), making HEEH_TK_rd5_0341 almost 2 kcal/mol more stable than the previous highest stability observed for a designed $\alpha\beta\beta\alpha$ structure (1). Helix 1 was less stable ($\Delta G_{\text{open}} \sim 3$ kcal/mol), and the central β -hairpin was the least stable structure. Four residues in this hairpin (I21, G23, I24, and V26) form intramolecular hydrogen bonds that should protect those amides from exchange, but three of these residues exchanged too quickly to be measured by NMR. The fourth residue (V26) had an opening energy of 2.1 kcal/mol (Fig. 3F). This hierarchy of stabilities across the different secondary structures suggests the folding energy landscape is not fully cooperative.

Stability determinants of $\alpha\beta\beta\alpha$ designs from a restricted design strategy

We next investigated which design features correlated with folding stability. To this end, we computed over a thousand structural and sequence-based metrics for each design and analyzed whether particular metrics correlated with stability. Several of the strongest individual correlations are shown in Fig. 2. Designs were generally more stable if their Rosetta energy scores were lower (Fig. 2B) and had more hydrophobic residues and hydrophobic sidechain contacts (Fig. 2C-D). Hydrophobic residue count correlated more strongly with stability than Rosetta energy. Stability also increased if a design's sequence was highly compatible with its local backbone structure (see Methods and Fig. 2E). Finally, increased net charge destabilized our designs, although the optimal net charge was slightly negative (Fig. 2F-G). This stability change was approximately linear with the square of the net charge, as expected (27).

We also explored whether specific residues could individually have large influences on the stabilities of the designs. Because all designs are based on an identical architecture, each position in the sequence shares an identical structural role in all designs. Using the binomial test, we identified positions where specific amino acid identities had large and significant changes on the success rates of the designs (Fig. S3). Two positions near the N- and C-termini stood out as particularly important. Positions 2 and 39 are near the tips of each helix and contact each other in space. Across the design set, leucine residues at these positions increased the success rate of the designs by 25-39%, whereas other residues such as glutamate and tryptophan decreased the success rate by similar amounts. These differences in success rates were highly significant (adjusted p-value $< 10^{-18}$) (Fig. S3). The importance of these residues suggests that termini of the helices play an especially important role in the overall stability of designed $\alpha\beta\beta\alpha$ miniproteins.

To further examine individual residue contributions to stability, we performed deep mutational scanning analyses on the six $\alpha\beta\beta\alpha$ designs whose structures we verified by CD (Fig. 3A). Using our protease sensitivity assay, we measured the folding stability changes for all single mutants of each design (Fig. S2C-D, Fig. S4). The residues most sensitive to mutations were typically nonpolar residues buried in the protein core, although

critical polar and charged residues were observed as well. For example, designs HEEH_TK_rd5_0341 and HEEH_TK_rd5_3711 were destabilized if charged/polar residues lost hydrogen bonding interactions with the backbone (Fig. 3D, Fig. S4). The least hydrophobic design, HEEH_TK_rd5_0958, had a salt bridge that contributed to its stability (Glu6 and Arg10) (Fig. 3D).

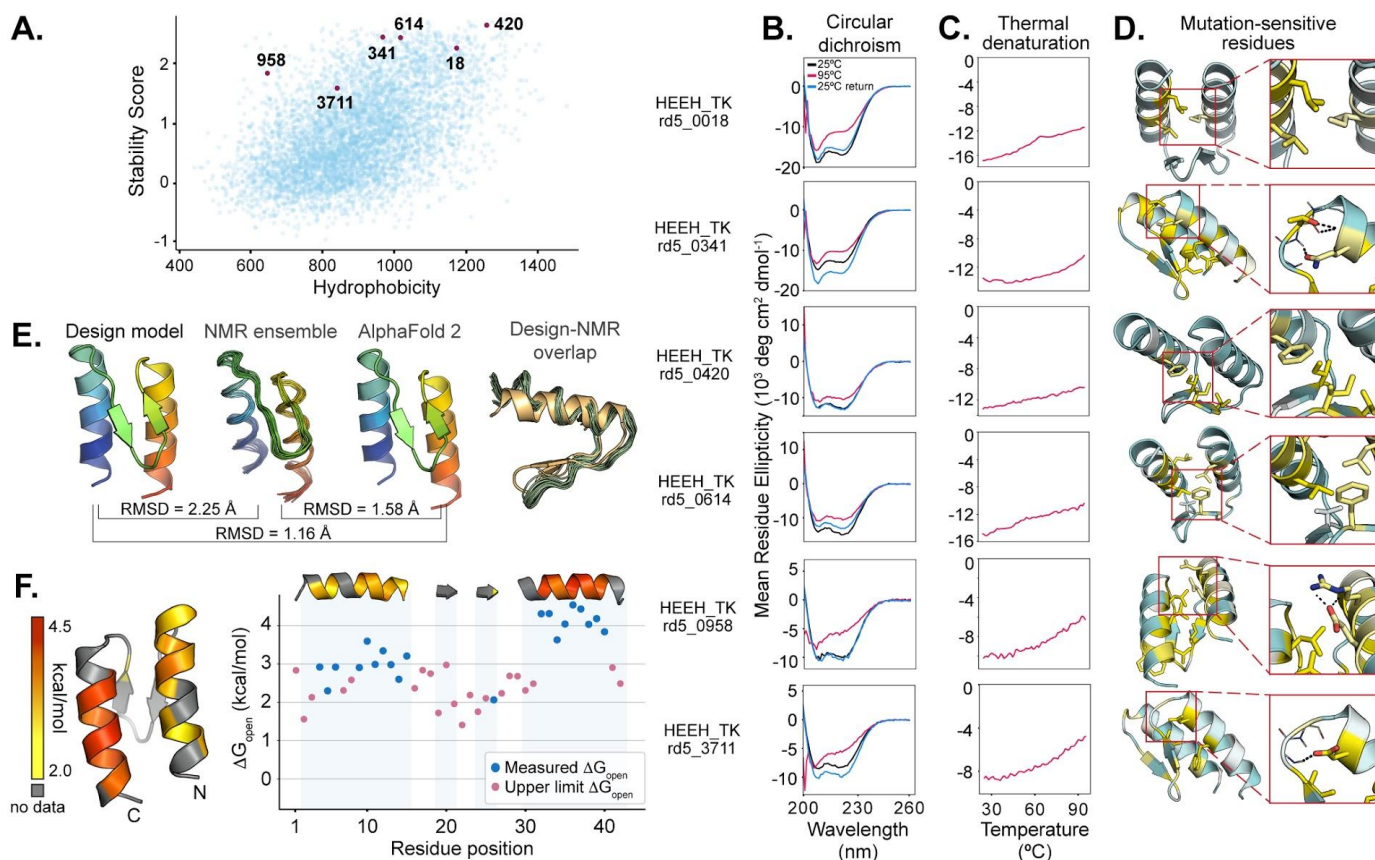


Fig. 3. Biophysical characterization of $\alpha\beta\alpha$ mini-proteins made using a restricted design strategy. (A) The stability scores of all $\alpha\beta\alpha$ mini-proteins made using a restricted design strategy are plotted by their hydrophobicity values (27). We selected six mini-proteins (red dots) with varying hydrophobicity for (B) far-ultraviolet circular dichroism measured at 25 $^{\circ}\text{C}$ (black), 95 $^{\circ}\text{C}$ (red), and 25 $^{\circ}\text{C}$ after melting (blue), and (C) thermal denaturation analysis measured at 220 nm at every 1 $^{\circ}\text{C}$ from 25 $^{\circ}\text{C}$ to 95 $^{\circ}\text{C}$. (D) Design models highlighting positions that are most tolerant (teal) or least tolerant (yellow) to mutations. Key designed residues that stabilize the protein are shown in stick representation. Each mini-protein's color scale is different to highlight the relative stabilizing or destabilizing effects within each protein; see Fig. S4 for complete data. (E) Comparison of HEEH_TK_rd5_341 Rosetta design model, NMR ensemble, and AlphaFold 2-predicted structures; overlay of the Rosetta design model (yellow) and NMR ensemble (green). (F) Opening energies determined by hydrogen-deuterium exchange for HEEH_TK_rd5_341. Observed measurements are colored red-yellow on a cartoon model and plotted in blue. For residues that exchanged too quickly to measure, the upper limit of ΔG_{open} is plotted in red.

The overall patterns of mutational sensitivity were consistent with the designed $\alpha\beta\alpha$ structures and provided further evidence that the stable sequences folded as designed. However, our mutational data also revealed some unexpectedly stable mutants (Fig. S4). For example, we expected that mutants to G23 would be highly destabilizing because G23 should be critical for forming the central β -hairpin. However, in four of the six designs, mutants to G23 could actually increase folding stability (Fig. S4). To investigate this, we predicted the structures of all mutant sequences using AlphaFold 2 (21). Although most mutants were predicted to have

similar structures to the original design, some predictions (including mutants of G23) suggest the possibility of alternative, compact structures (Fig. S5).

Modeling relative contributions of biophysical determinants on folding stability

Our previous analysis identified individual determinants of stability without considering how various features relate to each other. Hence, we next analyzed which protein features were the most important contributors to stability and how they compared to each other. Instead of prioritizing predictive accuracy, we used linear regression to build a parsimonious, interpretable, low-resolution model. Our moderately accurate model ($r = 0.64$, $r^2 = 0.41$; Fig. 2H, Table S3) included ten features chosen for either their large individual contributions to stability or their biophysical interest. Adding all 25 additional Rosetta energy terms provides only a minimal improvement to this low-resolution model (Table S4).

To analyze the strengths of the different features, we compared the different coefficients both in terms of their importance within our dataset (e.g. the impact of a one standard deviation change in each term, Fig. 2I left) and in terms of their biophysical strength (e.g. the impact of one additional residue, contact, charge, etc., Fig. 2I right). By representing the features in these two ways, we were able to observe how each feature contributes to a design's stability while holding all other features constant. Relative to the variance in the features, the count of large nonpolar residues is the largest contributor to folding stability (Fig. 2I). Additional biophysical determinants known to stabilize globular proteins (2, 28–30), such as contacts between adjacent nonpolar residues and Ser/Thr helix capping, contribute to folding stability as well (Fig. 2I). However, our model also points to the stabilizing role of nonpolar residues at the design ends, which is a feature specific to the $\alpha\beta\beta\alpha$ topology (Fig. 2I). Whereas previous studies on the relative importance of stability determinants were based on assays that changed one feature on individual proteins (31, 32), our large-scale testing enabled us to analyze over a thousand protein features on several thousand proteins in parallel. This, in turn, allowed us to develop a model that offers criteria for designing even more stable $\alpha\beta\beta\alpha$ miniproteins.

Designing $\alpha\beta\beta\alpha$ miniproteins using a diversity-oriented design strategy

Our restricted-design strategy (Round 5) focused on improving the success rate of designing stable $\alpha\beta\beta\alpha$ miniproteins but at the cost of reducing their structural diversity. Because we were now able to successfully generate stable $\alpha\beta\beta\alpha$ designs, we next investigated whether we could loosen the design restrictions that we had imposed, increase the diversity of our $\alpha\beta\beta\alpha$ miniproteins, and identify additional determinants of stability. Hence, we designed a new round of “diversity-oriented” (Round 6) $\alpha\beta\beta\alpha$ miniproteins based on fourteen different protein architectures instead of one. This allowed designs to have a greater variety of helix, β -strand, and loop lengths, while keeping the overall size of the protein to 43 residues (Fig. 1B). In addition, we did not impose residue restrictions on β -strands or in the middle loop and permitted a greater number of hydrophobic residues.

Importantly, we also used our Round 5 stability data to directly re-weight the Rosetta energy function. Using ridge regression, we adjusted the weights on the Rosetta energy terms to create the best correlation with our measured Round 5 $\alpha\beta\beta\alpha$ stabilities, while regularizing the regression to penalize large deviations from the original weights. With this approach, we created three new energy functions labeled “Minor,” “Medium,” and

“Heavy” based on how much the weights deviated from the original weights. We used these three energy functions (and the original weights) to design our Round 6 designs (Fig. S6).

We generated ~ 20,000 designs and chose our final set of over five thousand $\alpha\beta\beta\alpha$ designs for experimental testing by identifying designs that had the greatest structural diversity, varied sequence identity (no closer than 28/43 residues), and an $\alpha\beta\beta\alpha$ topology as determined by the computer program “Define Secondary Structure of Proteins” (DSSP) (33, 34) (See Methods). Notably, we prioritized structural diversity in our final selection instead of prioritizing the expected success rate. Each design is named HEEH_KT_rd6_####, in which KT indicates the designer (author K.T.), rd6 indicates these designs constitute a new “Round 6” following the previous rounds of $\alpha\beta\beta\alpha$ design, and #### is a design number.

Stability determinants of $\alpha\beta\beta\alpha$ designs based on diversity-oriented design strategy

We tested the stabilities of our “diversity-oriented” $\alpha\beta\beta\alpha$ miniproteins (and matching scrambled sequences) using the high-throughput protease sensitivity assay (1). Surprisingly, 12% of our scrambled sequences had stability scores above 1, compared to 2% or fewer in previous rounds (Fig. 4A). We further found that scrambled sequences were most likely to be stable when they were very hydrophobic and when their sequences had high helical propensity as determined by DSSP (33, 34) (Fig. S7). This suggested that designed sequences might also be stabilized by these properties alone, even if they did not fold into their designed structures. To remove these potential “false positive” designs from our analysis, we restricted our analysis to designs with a lower nonpolar residue count and lower helical propensity (Fig. S7). Restricting our analysis in this way removed 25% of our total designs, while lowering the fraction of stable scrambles from 12% to 6% (Fig. 4B). The overall fraction of stable designs was 26% - still substantially above the “success rate” of the scrambled sequences (Fig. 4B).

We first analyzed the impact of differently weighted Rosetta energy functions on folding stability. On average, designs made using the reweighted energy functions had higher stability than designs made with the default energy function (Fig. 4C-D). However, some regularization (restraining the weights near their original values) was critical to successful re-weighting: the “Heavy” energy function, where the changes to the weights were the largest, performed much more poorly than the energy functions with “Minor” and “Moderate” changes to the weights (Fig. 4C-D). The success of the re-weighted energy functions suggests that empirical re-weighting could be an efficient practical tool for protein design in situations where large-scale data is available for a specific task. The designs created by the re-weighted energy functions would not have been favored under our previous design procedure, with larger changes to the weights leading to designs that appear less and less favorable according to the default energy function (Fig. 4E).

Next, we investigated how topological features (loop, β -strand, helix) of the designs affect folding stability. We selected the seven most common loop structures found in our designs (represented using ABEGO notation) (35) and the three most common β -strand lengths as inputs to another linear regression model (Fig. S8, Fig. 4F). The explanatory strength of this model is weak (95% conf. int. from bootstrapping, mean $r = 0.167$, mean $R^2 = 0.028$). This is due to the simplicity of the model and because the topology-only model excludes critical stability determinants such as hydrophobic residue count. Despite these shortcomings, this model still enables

us to examine the relative importance of different topological components. The largest structural contributors to stability are the lengths of β -strands and helices, with shorter β -strands (and corresponding longer helices) as the most favorable topological parameter (β -strand and helix lengths are inversely related because all designs have a fixed length of 43 residues) (Fig. 4F). Secondly, particular structures in loops 2 and 3 influenced folding stability as well. A loop structure of GBB in the first loop, GG in the second loop, and AB in the third loop increases the stability of a design more than other loop structures (Fig. 4F).

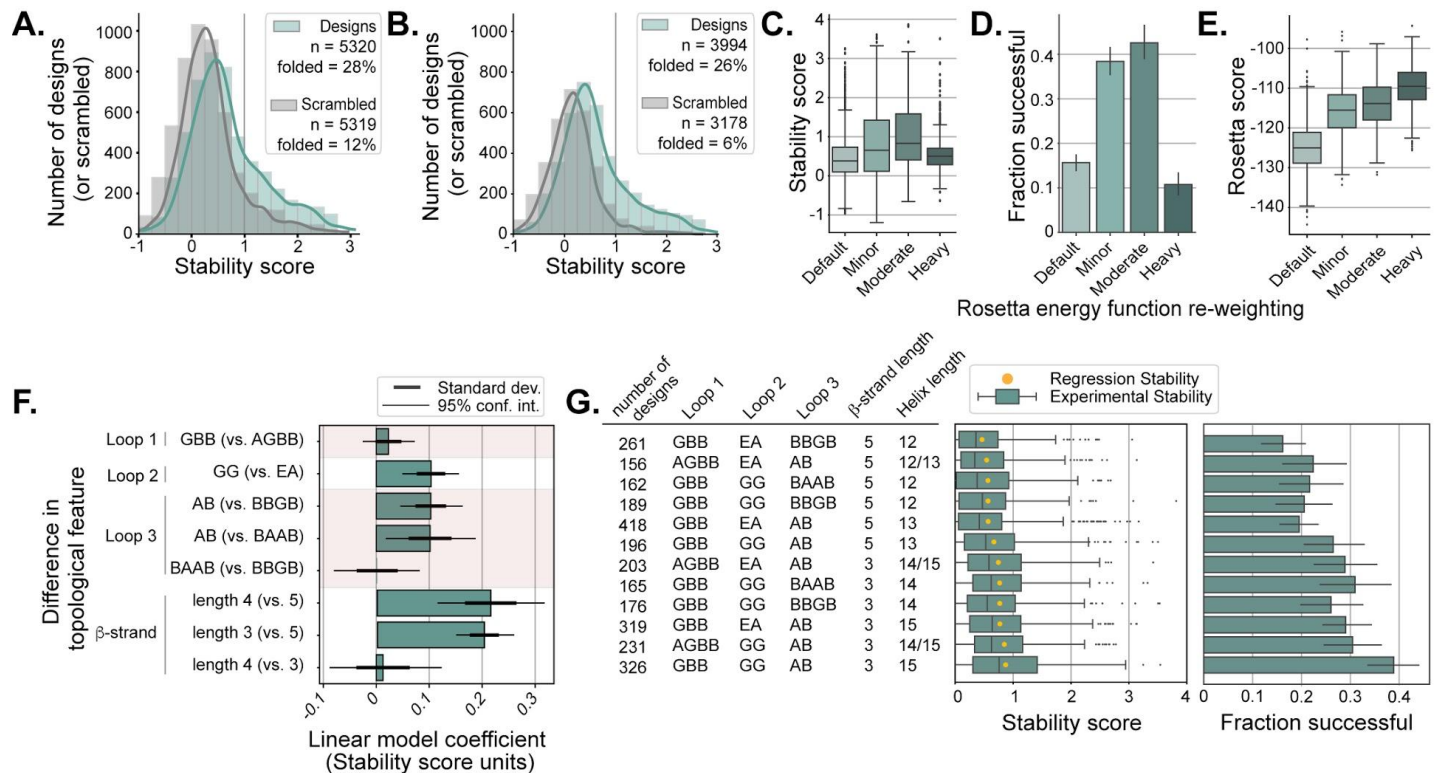


Fig. 4. Experimental testing and analysis of $\alpha\beta\alpha$ stability determinants from a diversity-oriented design strategy. (A) Stability score distribution of $\alpha\beta\alpha$ mini-proteins (green) and scrambled sequences (gray). (B) As in A, filtered to eliminate designed and scrambled sequences that may fold into non-designed structures; see text and Fig. S7. The vertical line at stability score = 1 denotes the threshold above which we consider a design to be stable. (C-D) Stability scores and success frequencies of designs made with differently-weighted Rosetta energy functions; “Heavy” indicates the largest amount of reweighting. (E) Rosetta scores (using the unmodified score function) of designs made using different weighting; the more positive scores of the designs from the re-weighted energy functions indicate these designs are less favorable according to the default energy function. (F) Stability contribution of the most common loop patterns (using ABEGO notation) and β -strand lengths based on a linear regression model. (G) The most common unique structure combinations (loop pattern, β -strand and helix lengths) are listed (left) followed by the distribution of observed stability scores (middle, with the expected stability from the linear regression model as a yellow dot). At right, the fraction of stable designs for each unique structure. All error bars indicate 95% confidence intervals from bootstrapping.

Based on this topology-focused model, we would expect $\alpha\beta\beta\alpha$ miniproteins with a GBB-GG-AB loop patterning, β -strands that are 4 residues long, and helices that are 14-residues long to be more stable on average than $\alpha\beta\beta\alpha$ miniproteins with any other loop, strand, and helix combination (Fig. 4F). Although designs with a β -strand length of 4 residues were not common in our dataset, a very similar design structure (GBB-GG-AB with a β -strand length of 3 residues) had the highest average stability score and the highest success rate in our dataset (Fig. 4G), which is in agreement with a previous study on loop patterning and stability (36). In fact, this design pattern is the protein architecture that we used to generate all the Round 5 $\alpha\beta\beta\alpha$ miniproteins (see Fig. S3A). However, the high success of this architecture in Round 6 may be due to using re-weighted energy functions that were optimized based on Round 5 designs with this specific architecture. Nonetheless, when we subset our Round 6 designs to identify $\alpha\beta\beta\alpha$ miniproteins with a GBB-GG-AB loop pattern and features that we previously determined to promote stability, these designs are diverse in their sequence identity and highly stable (81% successful) (Fig. 5). This provides a “recipe” for designing new stable $\alpha\beta\beta\alpha$ miniproteins in the future.

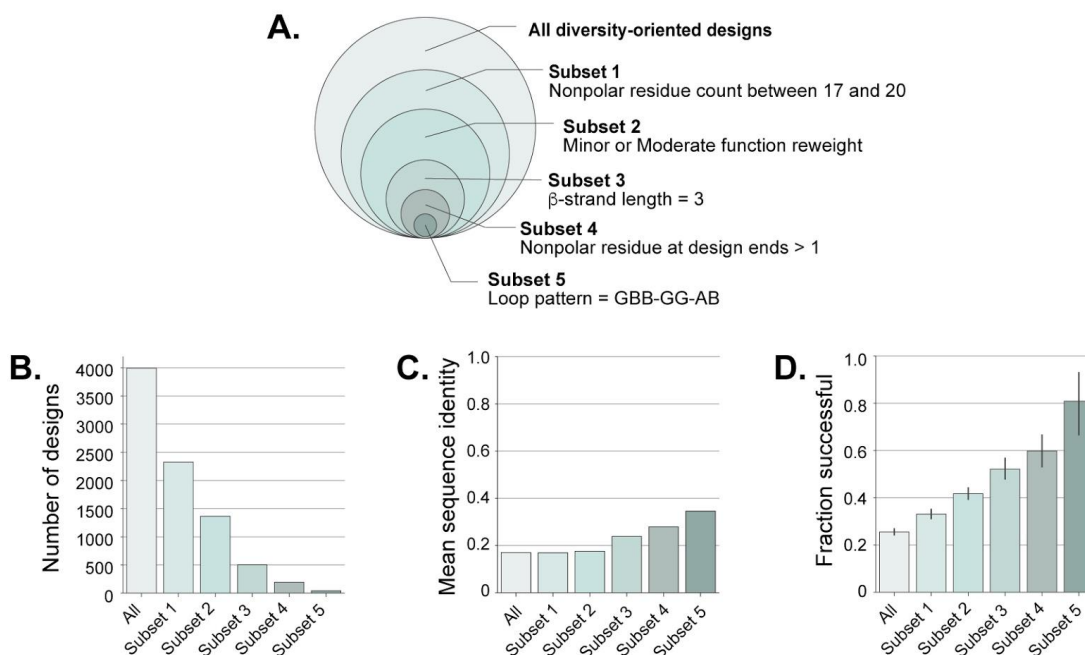


Fig 5. A recipe for building diverse high-stability $\alpha\beta\beta\alpha$ designs. (A) Designs made from a diversity-oriented strategy are grouped into subsets based on five features that we identified to be important for stability (Fig. 2I, Fig. 4F). (B) The number of designs that comprise each subset; (C) the mean sequence identity between any two designs in each subset; (D) the fraction of successful designs in each subset, with error bars indicating 95% confidence intervals from bootstrapping. Ideal designs (those with the parameters of Subset 5) are 80% successful with under 40% sequence identity between pairs of designs.

Predicting stable de novo $\alpha\beta\beta\alpha$ miniproteins by AlphaFold 2

When we designed and tested $\alpha\beta\beta\alpha$ miniproteins for their folding stability, AlphaFold 2 was not yet available. With its recent release (21), we wondered whether AlphaFold 2 could discriminate between stable and unstable miniproteins. We explored this possibility even though AlphaFold 2 is intended for structure prediction and not

stability prediction. Out of the ~5,600 and ~4,000 restricted and diversity-oriented designs, respectively, we found that 78% of the former and 20% of the latter had at least one predicted structure within 2 Å RMSD to the designed model. These predictions were equally in agreement with design models regardless of whether a design was experimentally unstable, moderately stable, or stable, indicating that AlphaFold 2 did not discriminate stable from unstable designs (Fig. S9 and S10A). We also examined whether the Rosetta energy scores of the AlphaFold 2-predicted models were better correlated with experimental stability scores than the scores of the original design models. The AlphaFold 2-predicted models did not improve the correlation with experiment for the Round 5 design set, but provided a small improvement for Round 6 (Fig. S10B-C). Neither RMSD nor AlphaFold 2's average confidence measure (pLDDT) showed much ability to enrich for stable designs (Fig. S10E), indicating that AlphaFold 2 is currently unable to determine the folding stability of these designed miniproteins.

Discussion

Understanding the biophysical determinants that enable proteins to fold and remain stable is important in protein design, drug development, and other areas. Here, we examined the stability determinants of the $\alpha\beta\beta\alpha$ miniprotein fold, which was previously identified as unusually challenging to design (1). We took advantage of an improved Rosetta design protocol (23, 24) to design over ten thousand $\alpha\beta\beta\alpha$ miniproteins using both restrictive and diversity-oriented design strategies. Our two design strategies led to over three thousand new stable designs (~2,100 restricted and ~1,000 diversity-oriented designs) and a much higher success rate (38%, Fig. 2A) than the 2% success previously reported (1). Selected $\alpha\beta\beta\alpha$ miniproteins folded as designed according to CD, NMR, and deep mutational scanning analysis (Fig. 3).

Our large dataset of stable designs enabled us to quantify determinants of stability for the previously-challenging $\alpha\beta\beta\alpha$ fold (Figs. 2I, 4F-G). Most of the stability determinants were common across globular proteins (2, 28, 29, 36, 32, 37–39), and similar to those previously observed in large-scale de novo design experiments (1). We also identified that designing hydrophobic residues near the termini was especially important for the $\alpha\beta\beta\alpha$ miniprotein fold (Fig. 2I). Our design success rate improved substantially when we used our large dataset to re-weight the Rosetta energy function specifically for $\alpha\beta\beta\alpha$ design (Fig. 4C-E). These observations largely explain the low success of previously designed $\alpha\beta\beta\alpha$ proteins: previous designs frequently employed non-optimal loop patterns, helix capping residues, and residues near the design termini, and typically had 13-16 nonpolar residues rather than the 17-20 used here (Fig. S1). Notably, the total number of nonpolar residues in each design is influenced by the design energy function and by parameters that restrict the amino acids that are sampled at each position according to the solvent accessibility of that position (40, 1). These restrictions are manually tuned to balance stability and solubility, as well as to reduce the search space of sequences. Designing proteins with too few nonpolar residues can thus be considered a failure of manual tuning as well as a failure of the design energy function.

Finally, our modeling enabled us to examine determinants of stability as independent and modular components, each of which independently contributes to the overall stability of a design. With this perspective, we suggest a “recipe” for designing diverse stable $\alpha\beta\beta\alpha$ miniproteins (Fig. 5). This viewpoint, however, assumes that determinants are additive, which is an approximation of the biophysics of protein folding (41). Despite this limitation, our $\alpha\beta\beta\alpha$ designs may be especially well-suited for examining folding stability using an additive

model due to their small size and simple structure. Simplified low-resolution models like our linear regression are valuable for building biophysical intuition about the strengths of different interactions (42, 43) as well as for guiding the construction of high-resolution models like the Rosetta energy function, which is also an additive model (44).

Overall, this study demonstrates the value of combining computational *de novo* protein design with large-scale experimental testing to solve a previously challenging design problem. As a result, the $\alpha\beta\beta\alpha$ topology expands the collection of known stable topologies (1, 36, 40, 45) and may also be valuable as binders for therapeutic, diagnostic, and synthetic biology applications (13, 46).

Materials and Methods

Computational protein design

We designed $\alpha\beta\beta\alpha$ miniproteins using Rosetta based on our previous work (1). Briefly, we used fragment assembly to build backbones according to protein architectures specified in a blueprint file (22), as in (40, 1). For the restricted design strategy (Round 5), we chose the protein architecture that previously led to generating the greatest number of stable (defined only here as stability score ≥ 0.8) $\alpha\beta\beta\alpha$ miniproteins (Fig. S1A). This architecture restricted the $\alpha\beta\beta\alpha$ miniprotein structure to have two helices that are 14-residues long, two β -strands that are 3-residues long, and three loops with an ABEGO pattern of GBB, GG, and AB, respectively. We also applied several design constraints. We forced the first residue in the middle loop (position 22) to be nonpolar (AFILMVWY), the solvent-facing residues in the β -strands (positions 20 and 25) to be polar or charged (QNSTDEHKR), and any helical positions were prevented from being designed as Gly, Thr, or Val as they are known to destabilize helix formation (47–49). We also required all designs to possess at least 15 hydrophobic residues (AFILMVWY) and no more than 21 hydrophobic residues. Finally, we filtered out designs with low Rosetta total energy scores or low buried nonpolar surface area (Fig. S1B-D).

Sequence design was performed using the Rosetta protocol FastDesign (24), the beta_nov16_protease version of the full-atom energy function, and a recently-improved sampling method designed to prevent over-compactness (50). In order to select $\alpha\beta\beta\alpha$ miniproteins for experimental testing, we ranked each $\alpha\beta\beta\alpha$ design by their predicted stability scores, which was determined by a Lasso regression model that we built using previous $\alpha\beta\beta\alpha$ miniprotein structural metrics and experimental stability scores (1). Based on this ranking, we selected the top $\sim 5,600$ designs with a threshold of 67% sequence identity for experimental testing.

Round 6 designs were designed as above with several changes. First, we utilized 14 different protein architectures. Moreover, we removed the hydrophobic restriction in the middle loop, were more permissive on non-helical residues (GDNST) inside the helices, and allowed hydrophobic residues to appear on the protein surface. We further specified a penalty for a protein's net charge outside the range of -5 or 3. Upon generating $\sim 20,000$ $\alpha\beta\beta\alpha$ designs, we took several steps to select over 7,000 designs for experimental testing. We first built Lasso and XGBoost regression models (51) using experimental data from Round 5 to identify $\sim 3,000$ designs with significantly different predictions between the two models (predicted stability scores were at least 0.25 scores away from the best-fit line between the models). We next independently performed principal component analysis to identify $\sim 9,000$ designs that were most distant from each other. From the combined $\sim 12,000$ $\alpha\beta\beta\alpha$ designs, we selected $\sim 7,400$ designs for experimental testing whose sequence identity was no closer than 66% to any other design.

Although all $\sim 7,400$ designs were experimentally tested, we determined afterwards that many of these structures either diverged away from the $\alpha\beta\beta\alpha$ topology during design or were not predicted to fold into an $\alpha\beta\beta\alpha$ structure by Rosetta's ab initio algorithm. We further found that scrambled sequences could form secondary structure according to psipred (52). To focus our analysis on designed $\alpha\beta\beta\alpha$ structures, we restricted our analysis to $\sim 5,300$ designs meeting these criteria: distance between the C-terminus to the middle loop < 22 Å, distance between the N- and C-termini < 20 Å, β -strand lengths according to

DSSP ≤ 5 residues, loop lengths ≤ 5 residues, and unbroken $\alpha\beta\beta\alpha$ secondary structural elements according to DSSP (33, 34).

All blueprint files, Rosetta input commands, design files are provided in the supporting information.

Energy function re-weighting

The Rosetta energy function is a weighted sum of individual, independent score terms (44). To test whether our experimental data could directly optimize the energy function for $\alpha\beta\beta\alpha$ miniprotein design, we sought to re-weight these terms in Round 6 to produce the best correlation with our experimentally measured stability scores from Round 5. In re-weighting, we also sought to bias our new weights to be as close to the original weights as possible by using ridge regression (51). However, because the L2 regularization in ridge regression biases coefficients to be near zero, we used ridge regression to identify optimal *perturbations* to our original weights, rather than directly optimizing the weights themselves. To determine the appropriate perturbations, we first regressed our set of experimentally measured stability scores against the original Rosetta (computational) total scores of the designs. We then used the residuals from this regression (i.e. the error in the prediction of experimental stability score) as the target values for our ridge regression. We used scikit-learn's implementation of Ridge regression (51) to determine new weights on the 25 unweighted Rosetta score terms that best fit the residuals of the first regression (Fig. S7). The coefficients in this second regression are effectively perturbations to the original Rosetta weights that minimize the error in predicting experimental stability scores (subject to the regularization constraint). After performing ridge regression, the new score function weights were determined based on the formula:

$$\text{NewWeight}_i = \text{OriginalWeight}_i * (1 + \text{Coefficient}_i)$$

where NewWeight_i is the new weight on score term i , OriginalWeight_i is the original weight for term i in the beta_nov16_protease energy function, and Coefficient_i is the coefficient on score term i in the ridge regression.

We tested three new weight sets (“Minor,” “Medium”, and “Heavy”) in addition to the default weights. These new weight sets were determined using three different regularization strengths in the ridge regression and are named based on the magnitude of the change. The “Minor” set used regularization strength $\alpha=200,000$; “Moderate” used $\alpha=20,000$, and “Heavy” used $\alpha=0.1$. “Heavy” corresponds to the value of α from a cross-validated ridge regression using scikit-learn's RidgeCV method (51). In all weight sets, the score term `fa_intra_rep_xover4` was maintained at its default value to avoid favoring extended structures. These weight sets are all provided in the supporting information.

Miniprotein library generation

We reverse translated the residue sequences and optimized the codons (based on *E. coli* codon frequencies) of all $\alpha\beta\beta\alpha$ miniproteins that we selected for experimental testing using DNAsworks 2.0 (53). We also included scrambled sequences (while preserving locations of P and G residues as well as nonpolar/polar patterning) for

each corresponding $\alpha\beta\alpha$ sequence (following (54)). Both oligo libraries (Round 5, and Round 6 + mutational scanning) were purchased from Agilent.

Yeast display and protease stability assay

DNA amplification, yeast display proteolysis, sorting, and next-generation sequencing were all performed by research contract to the University of Washington BioFab (55) according to the protocol of (1). Yeast display was performed using a display vector with improved protease resistance (24).

Computing stability scores

We calculated a “stability score” for each design based on a probabilistic model described previously (1). The model determines the EC_{50} (the protease concentration at which half of the yeast cells pass selection during flow cytometry) for each design. The difference between the experimental EC_{50} in the folded state and predicted EC_{50} in the unfolded state (based on the identical model from (1)) is what we call a “stability score.” The overall stability score for each sequence is the minimum of the independent stability scores measured by trypsin and chymotrypsin. As previously (1), Round 5 data were filtered based on the confidence interval of the EC_{50} estimate: only sequences where the 95% confidence interval was smaller than 2.0 (meaning the equivalent of two selection rounds, or 9x protease concentration) were retained for analysis. However, Round 6 data was not filtered based on the EC_{50} confidence intervals. As a result, 38/3178 scrambles (1%) but no designed sequences (Fig. 4B) with low-confidence stability estimates were included in the analysis. The mutational scanning data were also not filtered based on the EC_{50} confidence interval; however, only 6/4650 (0.1%) sequences had low confidence stability estimates.

Computing metrics and Regression modeling

The Rosetta models used for computing structural features were the lowest energy structures from at least 1,000 ab initio trajectories and 200 relax trajectories starting from the design models. Rosetta design models were scored using the Rosetta score function, and we computed structural and biophysical features pertaining to secondary structure, dipeptides, hydrophobicity, hydrogen-bonding, and fragment quality using the `score_monomeric_designs` package (https://github.com/Haddock/score_monomeric_designs).

For regression modeling, we performed linear regression by bootstrapping (sampling with replacement 1000 times) using Python scikit learn (51) and selected the 95% confidence interval for each variable’s coefficient for analysis. For the restricted design strategy, we first used stepwise linear regression to identify eight features (large nonpolar count, nonpolar residue-residue contacts, local sequence-structure propensity, Ser/Thr at helix caps, Glu-Arg residue-residue contacts, nonpolar residue at design ends (which we define as positions 1, 2, 42, and 43 of the 43-residue-long protein structure), Glu-Glu residue-residue contacts, and increased net charge) that increased the correlation coefficient between predicted and experienced stability scores. We also selected two features (favorable net charge at helix ends and buried unsatisfied polar atoms) to determine their relative contributions to stability. For the diversity-oriented topology-focused linear regression model, we selected the seven most common loop

patterns and the three most common β -strand lengths found in our dataset as inputs to a linear regression model.

Calculation of local sequence-structure agreement

The compatibility of each protein sequence with its local backbone structure (Fig. 2E) was computed using the `abego_res_profile` method from (1).

Protein expression and purification

We purchased six $\alpha\beta\beta\alpha$ designs whose nucleotide sequences were optimized for *E. coli* expression and encoded in the pET-28a(+) vector (that contains an N-terminal His-tag and thrombin cleavage) from Twist Bioscience. The plasmid vectors were transformed in BL21(DE3) competent cells (Invitrogen or Sigma Aldrich) and grown overnight in a starter culture of 50 mL LB media (Fisher Bioreagents) and 50 $\mu\text{g}/\text{mL}$ kanamycin at 37°C while shaking at 225 rpm. 16-18 hrs later, we inoculated 500 mL of LB media and 50 $\mu\text{g}/\text{mL}$ kanamycin with 10 mL of the starter culture and allowed the competent cells to grow until $\text{OD}_{600} \sim 0.6-0.8$.

In preparation for NMR analysis, we transformed one $\alpha\beta\beta\alpha$ design encoded in pET-28a(+) into BL21(DE3) competent cells (Sigma Aldrich) and grown in an LB media starter culture (as stated above). After 16-18 hrs, we pelleted the cells by centrifugation, replaced the LB media with M9 media (40 mM Na_2HPO_4 , 8.5 mM NaCl, 20 mM KH_2PO_4 , 60 mM d-Biotin, 55 mM Thiamine, 0.1 mM CaCl_2 , 0.01 mM ZnSO_4 , 2 mM MgSO_4 , 50 $\mu\text{g}/\text{mL}$ kanamycin) that included 15 mM $^{15}\text{NH}_4\text{Cl}$ and 10 mM ^{13}C glucose (Cambridge Isotopes) and resuspended the pellet. We then inoculated 500 mL of LB media with M9 media (including 15 mM $^{15}\text{NH}_4\text{Cl}$, 10 mM ^{13}C glucose and 50 $\mu\text{g}/\text{mL}$ kanamycin) with 10 mL of the resuspended starter culture and allowed the competent cells to grow until $\text{OD}_{600} \sim 0.6-0.8$.

Afterwards, for both labeled and unlabeled competent cells, we induced protein expression by adding a final concentration of 500 mM Isopropyl β -D-1-thiogalactopyranoside (IPTG) (Fisher Bioreagents) to the LB media and allowing the cells to grow overnight at 15°C while shaking at 225 rpm. We then harvested the cells by centrifugation at 4°C and lysed the cells in 30 mL of lysis buffer (20 mM Tris, 500 mM NaCl, 30mM imidazole, 0.25% CHAPS, 1mM PMSF, pH 8.0), which included 60 mg of chicken lysozyme (Sigma), 1.5 μL of benzonase nuclease (Sigma Millipore), and 1 tablet of Pierce protease Inhibitor EDTA-free (ThermoFisher) followed by sonication (QSonica SL-18).

Next, we separated insoluble bacterial material by centrifugation (10,000 x g for 30 min) and purified the $\alpha\beta\beta\alpha$ miniproteins by immobilized metal-affinity chromatography (IMAC), which involved transferring the supernatant onto Econo-pac columns (Bio-Rad) that were previously prepared with Ni-NTA (Qiagen), washing the column with 15 mL Wash Buffer (20 mM Tris, 500 mM NaCl, 30 mM imidazole, 0.25% CHAPS, 5% glycerol, pH 8.0), and eluting the samples in 10 mL of Elution Buffer (20 mM Tris, 300 mM NaCl, 500 mM imidazole, 5% glycerol, pH 8.0). We initially verified the size and purification of the miniproteins by SDS-PAGE electrophoresis and Coomassie stain gel analysis. Then, we concentrated them by a centrifugal filtration system (Amicon Ultracel-15) (we did not concentrate the protein for NMR analysis).

We further purified both labeled and unlabeled miniproteins by size-exclusion chromatography (Bio-Rad NGC Chromatography System) using a Superdex 75 10/300 GL column (GE Healthcare) and eluted in PBS buffer. Miniprotein size and purification were verified first by Coomassie stain gel analysis and then by mass spectrometry (Synapt G2 Si, Waters).

Circular dichroism

Far-ultraviolet circular dichroism measurements were performed on six $\alpha\beta\beta\alpha$ designs (HEEH_rd5_0018, HEEH_rd5_0341, HEEH_rd5_0420, HEEH_rd5_0614, HEEH_rd5_0958, and HEEH_rd5_3711) using a Jasco J-815 spectrophotometer. All analysis was performed on the unmodified expression constructs including a 21-residue N-terminal linker (MGSSHHHHHSSGLVPRGSHM). We measured the concentration samples with a Qubit 4 Fluorometer (Invitrogen) and diluted them to a final concentration of ~0.1-0.4 mg/mL in PBS buffer. Wavelength scan measurements were made using a 1 mm path-length cuvette from 195 to 260 nm at 25°C and 95°C. We also measured temperature melts at 220 nm for every 1°C from 25°C to 95°C. For temperature melt analysis, we smoothed the data with a Savitsky-Golay filter of polyorder = 3.

Nuclear Magnetic Resonance

NMR spectra for HEEH_TK_rd5_0341 structure determination were acquired at 288 K, on Bruker spectrometers operating at 600 or 800 MHz, equipped with TCI cryoprobes with the protein buffered in 20 mM sodium phosphate (pH 7.5, 150 mM NaCl) at concentrations of ~ 1.0 mM. Resonance assignments were determined for $^{15}\text{N}/^{13}\text{C}$ -labeled protein using FMCGUI (56) based on a standard suite of 3D triple and double-resonance NMR experiments collected as described previously (57). All 3D spectra were acquired with non-uniform sampling in the indirect dimensions and were reconstructed by the multi-dimensional decomposition software qMDD (58), interfaced with NMRPipe (59). Peak picking was performed manually using NMRFAM-Sparky (60). Torsion angle restraints were derived from TALOS+ (61). Automated NOE assignments and structure calculations were performed using CYANA 2.1 (62). The best 20 of 100 CYANA-calculated structures were refined with CNSSOLVE (63) by performing a short restrained molecular dynamics simulation in explicit solvent (64). The final 20 refined structures comprise the NMR ensemble. Structure quality scores were performed using Procheck analysis (65) and the PSVS server (66).

Hydrogen-deuterium exchange and analysis

NH/D amide exchange rates were determined for HEEH_TK_rd5_0341 (including the 21aa N-terminal linker MGSSHHHHHSSGLVPRGSHM) by performing an exchange series at 288 K (at 600 MHz), by monitoring the decay rate of amide peak intensities in ^1H - ^{15}N HSQC spectra collected over the course of 24 hrs. Exchange was initiated by mixing phosphate buffered HEEH_TK_rd5_0341 (at ~ 1.0 mM) with D_2O at a ratio of 1:19. Each HSQC time point was acquired in ~ 16 minutes; the first time point was started ~ 5 minutes following the initiation of the reaction. Peak intensities were fitted to a single exponential decay, with an offset due to the presence of residual 5% H_2O . Opening free energies were calculated from these rates as previously described (67–69). For residues where exchange was too fast to

quantify, we calculated an “upper limit” ΔG_{open} based on an exchange rate of 0.1 min^{-1} (the fastest quantifiable rate was 0.066 min^{-1}).

AlphaFold analysis

All $\alpha\beta\beta\alpha$ miniprotein structures were predicted from their primary sequences using AlphaFold 2 (21) without using multiple sequence alignment (MSA) information because $\alpha\beta\beta\alpha$ miniproteins have low similarity to natural proteins. Five models were generated for each sequence and the lowest RMSD model to the designed structure was used for the analysis in Figs. S5, S9, and S10.

Acknowledgements

This work was supported, in part, by the National Institute of General Medical Sciences through award number 1DP2GM140927 and award number 5T32GM105538. K.T. was supported by JSPS KAKENHI grant number (19J30003) and is currently supported by a Human Frontier Science Program Long-Term Fellowship. Computational work was supported in part through the resources and staff contributions provided for the Quest high performance computing facility at Northwestern University (which is jointly supported by the Office of the Provost, the Office for Research, and Northwestern University Information Technology). NMR work was performed by the Structural Genomics Consortium, a registered charity (no: 1097737) that receives funds from Bayer AG, Boehringer Ingelheim, Bristol Myers Squibb, Genentech, Genome Canada through Ontario Genomics Institute [OGI-196], EU/EFPIA/OICR/McGill/KTH/Diamond Innovative Medicines Initiative 2 Joint Undertaking [EUBOPEN grant 875510], Janssen, Merck KGaA (aka EMD in Canada and US), Pfizer and Takeda. Yeast display selections and next generation sequencing were performed by the University of Washington BioFab. CD spectroscopy was performed using Northwestern University’s Keck Biophysics Facility. We further thank the members of the Rocklin lab for discussions and comments on this manuscript.

References

1. G. J. Rocklin, *et al.*, Global analysis of protein folding using massively parallel design, synthesis, and testing. *Science* **357**, 168–175 (2017).
2. K. A. Dill, Dominant forces in protein folding. *Biochemistry* **29**, 7133–7155 (1990).
3. A. Goldenzweig, S. J. Fleishman, Principles of Protein Stability and Their Application in Computational Design. *Annu. Rev. Biochem.* **87**, 105–129 (2018).
4. M. Arai, Unified understanding of folding and binding mechanisms of globular and intrinsically disordered proteins. *Biophys. Rev.* **10**, 163–181 (2018).
5. P.-S. Huang, S. E. Boyken, D. Baker, The coming of age of de novo protein design. *Nature* **537**, 320–327 (2016).
6. S. E. Boyken, *et al.*, De novo design of protein homo-oligomers with modular hydrogen-bond network-mediated specificity. *Science* **352**, 680–687 (2016).
7. X. Pan, *et al.*, Expanding the space of protein geometries by computational design of de novo fold families. *Science* **369**, 1132–1136 (2020).
8. A. Goldenzweig, *et al.*, Automated Structure- and Sequence-Based Design of Proteins for High Bacterial Expression and Stability. *Mol. Cell* **63**, 337–346 (2016).
9. J. G. Wiese, S. Shanmugaratnam, B. Höcker, Extension of a de novo TIM barrel with a rationally designed

- secondary structure element. *Protein Sci. Publ. Protein Soc.* (2021) <https://doi.org/10.1002/pro.4064>.
10. C. J. Lalaurie, *et al.*, The de novo design of a biocompatible and functional integral membrane protein using minimal sequence complexity. *Sci. Rep.* **8**, 14564 (2018).
 11. A. Nisthal, C. Y. Wang, M. L. Ary, S. L. Mayo, Protein stability engineering insights revealed by domain-wide comprehensive mutagenesis. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 16367–16377 (2019).
 12. A. Broom, K. Trainor, Z. Jacobi, E. M. Meiering, Computational Modeling of Protein Stability: Quantitative Analysis Reveals Solutions to Pervasive Problems. *Structure* **28**, 717–726.e3 (2020).
 13. A. Chevalier, *et al.*, Massively parallel de novo protein design for targeted therapeutics. *Nature* **550**, 74–79 (2017).
 14. L. Cao, *et al.*, De novo design of picomolar SARS-CoV-2 miniprotein inhibitors. *Science* **370**, 426–431 (2020).
 15. C. M. Bryan, *et al.*, Computational design of a synthetic PD-1 agonist. *Proc. Natl. Acad. Sci.* **118** (2021).
 16. J. Dou, *et al.*, De novo design of a fluorescence-activating β -barrel. *Nature* **561**, 485–491 (2018).
 17. B. Basanta, *et al.*, An enumerative algorithm for de novo design of proteins with diverse pocket structures. *Proc. Natl. Acad. Sci.* **117**, 22135–22145 (2020).
 18. T. Linsky, *et al.*, Sampling of Structure and Sequence Space of Small Protein Folds. *bioRxiv*, 2021.03.10.434454 (2021).
 19. J. M. Singer, *et al.*, Large-scale design and refinement of stable proteins using sequence-only models. *bioRxiv*, 2021.03.12.435185 (2021).
 20. K. L. Maxwell, *et al.*, The solution structure of bacteriophage lambda protein W, a small morphogenetic protein possessing a novel fold. *J. Mol. Biol.* **308**, 9–14 (2001).
 21. J. Jumper, *et al.*, Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
 22. P.-S. Huang, *et al.*, RosettaRemodel: A Generalized Framework for Flexible Backbone Protein Design. *PLoS ONE* **6**, e24109 (2011).
 23. H. Park, *et al.*, Simultaneous Optimization of Biomolecular Energy Functions on Features from Small Molecules and Macromolecules. *J Chem Theory Comput*, **12** (2016).
 24. J. B. Maguire, *et al.*, Perturbing the energy landscape for improved packing during computational protein design. *Proteins Struct. Funct. Bioinforma.* **89**, 436–449 (2020).
 25. O. D. Monera, T. J. Sereda, N. E. Zhou, C. M. Kay, R. S. Hodges, Relationship of sidechain hydrophobicity and α -helical propensity on the stability of the single-stranded amphipathic α -helix. *J. Pept. Sci.* **1**, 319–329 (1995).
 26. B. M. P. Huyghues-Despointes, J. M. Scholtz, C. N. Pace, Protein conformational stabilities can be determined from hydrogen exchange rates. *Nat. Struct. Biol.* **6**, 910–912 (1999).
 27. R. S. Negin, J. D. Carbeck, Measurement of Electrostatic Interactions in Protein Folding with the Use of Protein Charge Ladders. *J. Am. Chem. Soc.* **124**, 2911–2916 (2002).
 28. L. Serrano, Fersht, Alan R, Capping and α -helix stability. *Nature* **342**, 296–299 (1989).
 29. W.-Y. Wan, E. J. Milner-White, A recurring two-hydrogen-bond motif incorporating A serine or threonine residue is found both at α -helical N termini and in other situations. *J. Mol. Biol.* **286**, 1651–1662 (1999).
 30. C. Nick Pace, J. M. Scholtz, G. R. Grimsley, Forces stabilizing proteins. *FEBS Lett.* **588**, 2177–2184 (2014).
 31. C. N. Pace, *et al.*, Contribution of hydrogen bonds to protein stability. *Protein Sci.* **23**, 652–661 (2014).
 32. C. N. Pace, *et al.*, Contribution of Hydrophobic Interactions to Protein Stability. *J. Mol. Biol.* **408**, 514–528 (2011).
 33. W. Kabsch, C. Sander, Dictionary of protein secondary structure: Pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* **22**, 2577–2637 (1983).
 34. W. G. Touw, *et al.*, A series of PDB-related databanks for everyday needs. *Nucleic Acids Res.* **43**, D364–D368 (2015).
 35. R. T. Wintjens, M. J. Rooman, S. J. Wodak, Automatic Classification and Analysis of $\alpha\alpha$ -Turn Motifs in Proteins. *J. Mol. Biol.* **255**, 235–253 (1996).
 36. Y.-R. Lin, *et al.*, Control over overall shape and size in de novo designed proteins. *Proc. Natl. Acad. Sci.* **112**, E5478–E5485 (2015).
 37. A. Chakrabarty, A. J. Doig, R. L. Baldwin, Helix capping propensities in peptides parallel those in proteins.

- Proc. Natl. Acad. Sci.* **90**, 11332–11336 (1993).
38. M. Kurnik, L. Hedberg, J. Danielsson, M. Oliveberg, Folding without charges. *Proc. Natl. Acad. Sci.* **109**, 5705–5710 (2012).
 39. Y. Gavrillov, S. Dagan, Y. Levy, Shortening a loop can increase protein native state entropy. *Proteins* **83**, 2137–2146 (2015).
 40. N. Koga, *et al.*, Principles for designing ideal protein structures. *Nature* **491**, 222–227 (2012).
 41. K. A. Dill, Additivity Principles in Biochemistry. *J. Biol. Chem.* **272**, 701–704 (1997).
 42. G. Bellesia, A. I. Jewett, J. Shea, Relative stability of de novo four-helix bundle proteins: Insights from coarse grained molecular simulations. *Protein Sci.* **20**, 818–826 (2011).
 43. T. Ha-Duong, “Coarse-Grained Models of the Proteins Backbone Conformational Dynamics” in *Protein Conformational Dynamics*, Advances in Experimental Medicine and Biology., K. Han, X. Zhang, M. Yang, Eds. (Springer International Publishing, 2014), pp. 157–169.
 44. R. F. Alford, *et al.*, The Rosetta All-Atom Energy Function for Macromolecular Modeling and Design. *J. Chem. Theory Comput.* **13**, 3031–3048 (2017).
 45. E. Marcos, *et al.*, De novo design of a non-local β -sheet protein with high stability and accuracy. *Nat. Struct. Mol. Biol.* **25**, 1028–1034 (2018).
 46. L. Cao, *et al.*, Robust de novo design of protein binding proteins from target structural information alone. *bioRxiv*, 2021.09.04.459002 (2021).
 47. K. T. O’Neil, W. F. DeGrado, A Thermodynamic Scale for the Helix-Forming Tendencies of the Commonly Occurring Amino Acids. *Science* **250**, 646–651 (1990).
 48. P. C. Lyu, M. I. Liff, L. A. Marky, N. R. Kallenbach, Side Chain Contributions to the Stability of Alpha-Helical Structure in Peptides. *Science* **250**, 669–673 (1990).
 49. S. Padmanabhan, S. Marqusee, T. Ridgeway, T. M. Laue, R. L. Baldwin, Relative helix-forming tendencies of nonpolar amino acids. *Nature* **344**, 268–270 (1990).
 50. R. E. Pavlovicz, H. Park, F. DiMaio, Efficient consideration of coordinated water molecules improves computational protein-protein and protein-ligand docking discrimination. *PLOS Comput. Biol.* **16**, e1008103 (2020).
 51. F. Pedregosa, *et al.*, Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
 52. D. T. Jones, Protein secondary structure prediction based on position-specific scoring matrices. *J. Mol. Biol.* **292**, 195–202 (1999).
 53. D. M. Hoover, J. Lubkowski, DNAWorks: an automated method for designing oligonucleotides for PCR-based gene synthesis. *Nucleic Acids Res.* **30**, e43 (2002).
 54. S. Kamtekar, J. M. Schiffer, H. Xiong, J. M. Babik, M. H. Hecht, Protein design by binary patterning of polar and nonpolar amino acids. *Science* **262**, 1680–1685 (1993).
 55. J. Vrana, *et al.*, Aquarium: open-source laboratory software for design, execution and data management. *Synth. Biol.* **6**, ysab006 (2021).
 56. A. Lemak, C. A. Steren, C. H. Arrowsmith, M. Llinás, Sequence specific resonance assignment via Multicanonical Monte Carlo search using an ABACUS approach. *J. Biomol. NMR* **41**, 29 (2008).
 57. A. Lemak, *et al.*, A novel strategy for NMR resonance assignment and protein structure determination. *J. Biomol. NMR* **49**, 27–38 (2011).
 58. K. Kazimierczuk, V. Yu. Orekhov, Accelerated NMR Spectroscopy by Using Compressed Sensing. *Angew. Chem. Int. Ed.* **50**, 5556–5559 (2011).
 59. F. Delaglio, *et al.*, NMRPipe: A multidimensional spectral processing system based on UNIX pipes. *J. Biomol. NMR* **6**, 277–293 (1995).
 60. W. Lee, M. Tonelli, J. L. Markley, NMRFAM-SPARKY: enhanced software for biomolecular NMR spectroscopy. *Bioinformatics* **31**, 1325–1327 (2015).
 61. Y. Shen, F. Delaglio, G. Cornilescu, A. Bax, TALOS+: a hybrid method for predicting protein backbone torsion angles from NMR chemical shifts. *J. Biomol. NMR* **44**, 213–223 (2009).
 62. P. Güntert, “Automated NMR Structure Calculation With CYANA” in *Protein NMR Techniques*, Methods in Molecular Biology™., A. K. Downing, Ed. (Humana Press, 2004), pp. 353–378.
 63. A. T. Brünger, *et al.*, Crystallography & NMR System: A New Software Suite for Macromolecular Structure Determination. *Acta Crystallogr. D Biol. Crystallogr.* **54**, 905–921 (1998).

64. J. P. Linge, M. A. Williams, C. A. E. M. Spronk, A. M. J. J. Bonvin, M. Nilges, Refinement of protein structures in explicit solvent. *Proteins Struct. Funct. Bioinforma.* **50**, 496–506 (2003).
65. R. A. Laskowski, M. W. MacArthur, D. S. Moss, J. M. Thornton, PROCHECK: a program to check the stereochemical quality of protein structures. *J. Appl. Crystallogr.* **26**, 283–291 (1993).
66. A. Bhattacharya, R. Tejero, G. T. Montelione, Evaluating protein structures determined by structural genomics consortia. *Proteins Struct. Funct. Bioinforma.* **66**, 778–795 (2007).
67. Y. Bai, J. S. Milne, L. Mayne, S. W. Englander, Primary structure effects on peptide group hydrogen exchange. *Proteins Struct. Funct. Bioinforma.* **17**, 75–86 (1993).
68. G. P. Connelly, Y. Bai, M.-F. Jeng, S. W. Englander, Isotope effects in peptide group hydrogen exchange. *Proteins Struct. Funct. Bioinforma.* **17**, 87–92 (1993).
69. D. Nguyen, L. Mayne, M. C. Phillips, S. Walter Englander, Reference Parameters for Protein Hydrogen Exchange Rates. *J. Am. Soc. Mass Spectrom.* **29**, 1936–1939 (2018).