

1   **De novo protein design by citizen scientists**

2  
3   Brian Koepnick, Jeff Flatten, Tamir Husain, Alex Ford, Daniel-Adriano Silva, Matthew J. Bick, Aaron Bauer,  
4   Gaohua Liu, Yojiro Ishida, Alexander Boykov, Roger D. Estep, Susan Kleinfelter, Linda Wei, Foldit Players,  
5   Gaetano T. Montelione, Zoran Popovic, Firas Khatib, Seth Cooper, David Baker\*

6   **Abstract**

7  
8   Online citizen science projects such as GalaxyZoo<sup>1</sup>, Eyewire<sup>2</sup> and Phylo<sup>3</sup> have been very  
9   successful for data collection, annotation, and processing, but for the most part have harnessed  
10   human pattern recognition skills rather than human creativity. An exception is the game  
11   EteRNA<sup>4</sup>, in which game players learn to build new RNA structures by exploring the discrete  
12   two-dimensional space of Watson-Crick base pairing possibilities. Building new proteins,  
13   however, is a more challenging task to present in a game, as both the representation and  
14   evaluation of a protein structure are intrinsically three-dimensional. We posed the challenge of  
15   *de novo* protein design in the online protein folding game Foldit<sup>5</sup>. Players were presented with a  
16   fully extended peptide chain and challenged to craft a folded protein structure with an amino  
17   acid sequence encoding that structure. After many iterations of player design, analysis of the  
18   top scoring solutions, and subsequent game improvement, Foldit players can now, starting from  
19   an extended polypeptide chain, generate a diversity of protein structures and sequences which  
20   encode them *in silico*. 146 Foldit player designs with sequences unrelated to naturally occurring  
21   proteins were encoded in synthetic genes; 58 were found to be expressed in *E. coli* with good  
22   solubility and to adopt stable monomeric folded structures in solution. The diversity of these  
23   structures is unprecedented in *de novo* protein design, representing 21 different folds—including  
24   a new fold that is unobserved in natural proteins. High resolution structures were determined for  
25   three of the designs, and are nearly identical to the player models. This work makes explicit the  
26   considerable implicit knowledge contributing to success in *de novo* protein design, and shows  
27   that citizen scientists can discover creative new solutions to outstanding scientific challenges,  
28   such as the protein design problem.

29  
30   **Main Text**

31  
32   The principle underlying *de novo* protein design is that proteins fold to their lowest free energy  
33   state<sup>6</sup>; hence, designing a new protein structure requires finding an amino acid sequence whose  
34   lowest energy state is the prescribed structure. In practice, this challenge can be divided into  
35   two subproblems: first, crafting a protein backbone that is designable (i.e. that could be the  
36   lowest energy state of some sequence); and second, finding a sequence whose lowest energy  
37   state is the crafted structure. One of the challenges of protein design is the exponentially  
38   increasing number of conformations available to a polypeptide chain, which is astronomical  
39   even for a modestly-sized protein of 60-100 residues. Thus, the first subproblem of crafting a  
40   plausible backbone is extremely open-ended, and the second subproblem is difficult because it  
41   is not tractable to explicitly check that a designed sequence has lower energy in the crafted  
42   structure than in any other structure. There has been considerable progress in *de novo* protein  
43   design in recent years<sup>7-10</sup>, but it is unclear whether all of the contributions to this success have  
44   been made explicit in the protocols used to design proteins, and how much implicit knowledge  
45   have been used.

46 resides in the expertise of the designers. Disentangling the role of expert knowledge is  
47 particularly difficult for the extremely open-ended challenge posed by the first subproblem (i.e.  
48 crafting a plausible backbone), for which there are a practically unlimited number of solutions.  
49 Because full computer enumeration of backbones is not possible, there is considerable room for  
50 human creativity and intuition in generating and designing new protein structures.

51  
52 To investigate how crowd-based creativity could contribute to solving the long-standing protein  
53 design problem, we incorporated *de novo* protein design tools into the protein folding game  
54 Foldit. Foldit is a free online computer game developed to crowdsource problems in protein  
55 modeling, and offers full control over the three-dimensional structure of a protein model<sup>5</sup> (Figure  
56 1). Players compete to build a model with the lowest free energy, as calculated by the Rosetta  
57 energy function<sup>11</sup>. In the past, Foldit has been primarily applied to protein structure prediction  
58 problems, in which Foldit players were presented with an unstructured amino acid sequence  
59 and challenged to determine its native conformation<sup>5,12</sup>. Foldit players in one case redesigned a  
60 loop region of an already folded structure<sup>13</sup>, but the *de novo* design of an entire protein is a far  
61 more expansive challenge.

62 We repeatedly challenged Foldit players to design stably folded proteins from scratch, and  
63 iteratively improved the game based on their results. In each challenge, players were provided  
64 with a poly-isoleucine backbone in a fully extended conformation (60-100 residues in length),  
65 and were given seven days to fold the backbone into a compact structure and identify a  
66 sequence specifying this backbone. Initially, most top-scoring (low energy) Foldit player designs  
67 were highly extended, lacked a solvent-inaccessible core, and were composed entirely of polar  
68 residues (Figure 1F). Such extended, fully  $\alpha$ -helical structures have more favorable hydrogen  
69 bonding, electrostatic, and local torsional energies than collapsed structures, which must contort  
70 to create a buried core. While poly-lysine and other extended polar sequences resembling these  
71 initial Foldit solutions are often  $\alpha$ -helical in solution<sup>15,16</sup>, the lack of long-range interactions  
72 precludes specific folding into a single stable structure<sup>17</sup>. This highlights a limitation of using  
73 absolute energy as an optimization criterion for protein design: a low energy design does not  
74 guarantee structural specificity, which arises only if all other alternative conformations have  
75 higher energy. To favor the design of globular solvent-excluding protein folds, with sequences  
76 that uniquely encode them, we introduced three supplementary design rules into Foldit: a “Core  
77 Exists” rule that requires a minimum proportion of residues (e.g. 30%) to be solvent-inaccessible  
78 in the designed structure; a “Secondary Structure Design” rule that prohibits glycine and alanine  
79 in all secondary structure elements; and a “Residue Interaction Energy” rule to penalize large  
80 residues that make insufficient intramolecular interactions in the designed structure. With the  
81 addition of these rules to Foldit, subsequent top-scoring designs from Foldit players were  
82 compact globular proteins.

83  
84 We obtained custom synthetic genes encoding 12 player designs for which structure prediction  
85 calculations converged on the player designed conformation<sup>14</sup>. The sequences of these proteins  
86 have no homology to any known protein (Sup. Table 1). The *de novo* designs were expressed in  
87 *E. coli* and purified by metal affinity and size exclusion chromatography. Chromatography and  
88 circular dichroism (CD) spectroscopy indicated that 7 of the 12 designs were monomeric and

89 folded in solution, with helical secondary structure consistent with the players' models (Sup. Fig.  
90 1). All of the experimentally tested proteins described in this paper, are entirely the work of  
91 Foldit players.

92  
93 During gameplay, the Foldit application uploads the player's latest model to the Foldit server  
94 every 2-5 minutes; from these snapshots we can reconstruct the process by which a Foldit  
95 player develops a protein design (Figure 2). Foldit players employ more varied and complex  
96 exploration strategies than standard Rosetta automated design protocols, and frequently revert  
97 to a previous iteration of their model to explore an alternative path, resulting in a highly-  
98 branched search tree. A typical automated design protocol, by contrast, includes only two  
99 branch points<sup>18</sup>. In addition, Foldit players regularly sample much higher energy states than the  
100 automated protocol (Figure 2), which has only a limited ability to escape local energy minima.

101  
102 Encouraged by the success of Foldit players in designing stable proteins from scratch, we made  
103 additions to the game encouraging players to explore more diverse protein structures. Up until  
104 this point, all top-scoring Foldit designs had consisted of either three or four  $\alpha$ -helices connected  
105 by minimal loops. Indeed, Foldit players had determined that designs with  $\beta$ -sheets did not  
106 score on par with  $\alpha$ -helical bundles, and competitive players had abandoned any attempt to  
107 design more varied folds. (This has an interesting parallel to protein design by practicing  
108 scientists, which has also focused much more on helical bundles than other classes of protein  
109 folds<sup>19-22</sup>.) To encourage the design of a wider variety of folds, we introduced a "Secondary  
110 Structure" rule stipulating that no more than 50% of residues may form  $\alpha$ -helices. Foldit players  
111 responded by designing a multitude of mixed  $\alpha/\beta$  proteins, which were indistinguishable from  
112 expert designs upon visual inspection. However, structure prediction calculations for these  $\alpha/\beta$   
113 design sequences showed poor sampling close to the target design structure, suggesting that  
114 the designed sequences did not strongly encode their local structure<sup>14</sup>. Further analysis showed  
115 that these player designs contained many residues with locally strained backbone  
116 conformations (backbone phi and psi torsions in unfavored regions of the Ramachandran  
117 plot<sup>23,24</sup>). That such designs had very low energies revealed a problem in the Rosetta energy  
118 function at the time: since Rosetta users typically sampled backbones starting from fragments of  
119 native proteins, unfavorable local conformations were rarely encountered—hence it had not  
120 been discovered that the energies associated with local backbone strain were being  
121 underestimated. We addressed this flaw in the Rosetta model by increasing the steepness of  
122 the energetic penalties associated with strained local backbone geometry; this is now standard  
123 in the latest Rosetta energy function<sup>11</sup>. We also added to Foldit an "Ideal Loops" rule restricting  
124 players to a set of 19 unstrained reverse-turn conformations<sup>7</sup>, and incorporated new tools to aid  
125 generation of unstrained backbones: a fragment lookup-based loop-closure tool, an interactive  
126 Ramachandran map, and a protein Blueprint scheme for drag-and-drop assembly of secondary  
127 structure elements and common loop conformations. Together, these upgrades brought about a  
128 marked improvement in the local backbone quality of Foldit player-designed proteins (Sup. Fig.  
129 3).

130  
131 The importance of reducing local backbone strain was borne out in experimental  
132 characterization. Prior to the backbone modeling improvements described in the previous

133 paragraph, only 5 of 37 Foldit  $\alpha/\beta$  designs tested (14%) were monomeric and structured in  
134 solution. Following the backbone modeling additions, 46 of 97 (47%) were monomeric and  
135 exhibited the expected secondary structure in solution. Most showed exceptional stability in  
136 thermal and chemical denaturation experiments, with free energies of unfolding ( $\Delta G_{unf}$ ) up to  
137 >15 kcal/mol; indeed, 32 designed proteins remained completely folded at 95°C (Figure 3;  
138 Extended Data Fig. 1). This success rate surpasses previous reports of designed  $\alpha/\beta$   
139 proteins<sup>7,12</sup>.

140  
141 Overall, the 58 successful Foldit designs are diverse in structure, representing 21 different  
142 protein folds (Figure 3; Extended Data Fig. 2)—one of which is a new fold previously  
143 unobserved in natural proteins. The success of Foldit designs is not attributed to just one or two  
144 exceptional Foldit players, but is shared broadly by the Foldit community. The 58 successful  
145 designs were created by 37 different Foldit players (the most prolific player authored 11  
146 successful designs); 19 designs were created collaboratively by at least two cooperating  
147 players, and 5 designs were not top-scoring, but regardless were flagged by players as personal  
148 favorites.

149  
150 We succeeded in solving structures of three Foldit player-designed proteins. X-ray crystal  
151 structures of two designed proteins, named by their designers as foldit1 and Peak6, closely  
152 match the designed conformations, with Ca-RMSD of 1.1 and 0.9 Å, respectively (Figure 4).  
153 Well-resolved electron density in the protein core shows that most sidechains adopt the  
154 intended rotamers and preserve the designed packing interactions. The solution NMR structure  
155 of a third design, foldit3, also closely matches the design conformation, with a with a Ca-RMSD  
156 of 1.1 Å between the design model and a representative structure (i.e., the medoid conformer<sup>25</sup>)  
157 of the ensemble.

158  
159 We can draw several general conclusions about scientific models, citizen science, and the  
160 interplay between the two. First, a scientific model which holds within the domain space  
161 considered by practicing scientists may not hold outside of this domain. This is most vividly  
162 illustrated by the highly extended structures generated by Foldit players in their first *de novo*  
163 design efforts, and later by the structures with strained local geometry not previously sampled  
164 by Rosetta users. Second, for citizen scientists to make essential and creative scientific  
165 contributions through online gaming, the scoring function of the game must be an accurate  
166 representation of the science. In our initial iterations, Foldit did not present to players a  
167 sufficiently accurate and general model to allow them to robustly design new proteins, even  
168 though the underlying Rosetta software had been used for protein design by practicing  
169 scientists. Third and most important, citizen science offers a powerful way to systematically  
170 improve a scientific model, through iterations of model trial and model improvement. Human  
171 game players are exceptionally capable at finding and exploiting unanticipated solutions that are  
172 otherwise unexplored by experienced scientists, whose focus is not on getting a high score, but  
173 rather on solving their specific scientific problem.

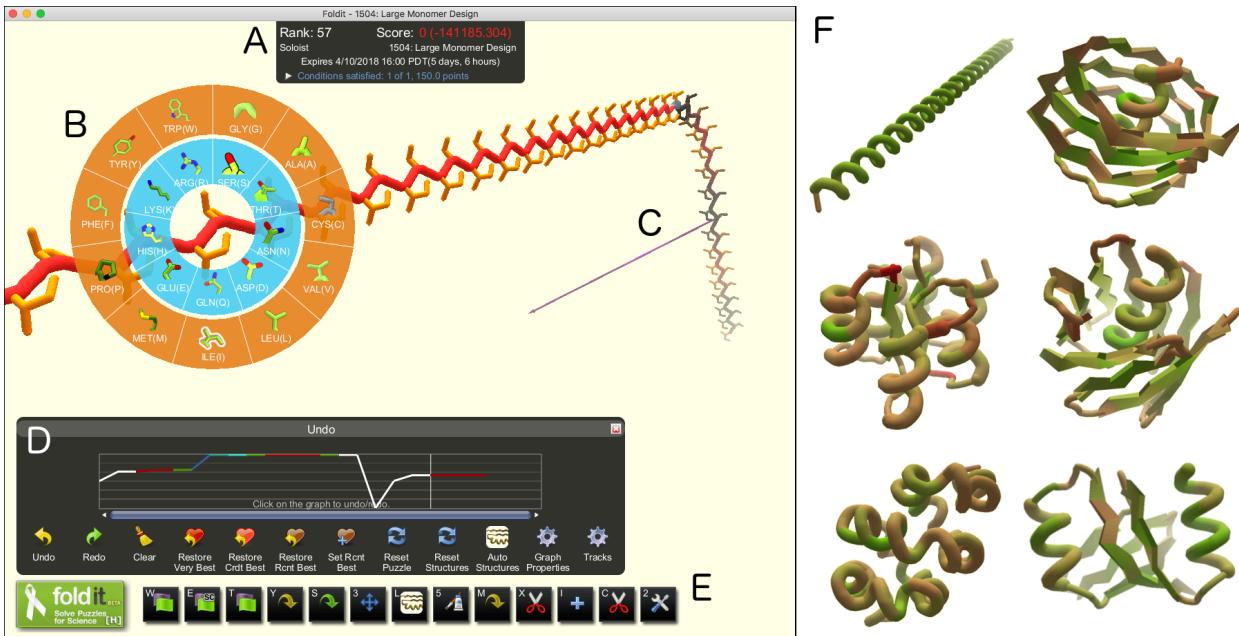
174  
175 We have demonstrated that non-expert citizen scientists, playing the online computer game  
176 Foldit, can accurately design completely new protein structures from scratch. Locally, players'

177 solutions are physically plausible and resemble natural proteins, but globally, they are creative  
178 and diverse. Proteins designed by citizen-scientist Foldit players are by no measure inferior to  
179 those of expert protein designers: they fold accurately to the intended conformation, show  
180 exceptional folding stability, and span a wide diversity of structures. This result is all the more  
181 impressive given that *de novo* protein design was an almost completely unsolved problem just a  
182 few years ago, and the diversity in protein folds spanned by the successful Foldit players'  
183 models considerably exceeds that in any previous protein design report. The sustained success  
184 of Foldit players over a wide diversity of protein folds highlights the power of human creativity  
185 when guided by scientific understanding presented in a readily comprehensible form.

186

187

188



189  
190 Figure 1. The Foldit user interface. (A) The Foldit score is the Rosetta energy with a negative  
191 multiplier, so that better models yield higher scores. (B) The design palette allows players to  
192 change the amino-acid residue identity at any position of the model. (C) The Pull tool allows  
193 players to manipulate the three-dimensional structure of the model. (D) The Undo graph tracks  
194 the score as a model is developed, and allows players to backtrack and load previous versions  
195 of a model. (E) Additional Foldit tools (from left to right): full structure minimization, sidechain  
196 minimization, backbone minimization, auto-design sidechains, repack sidechains,  
197 translate/rotate model, secondary structure assignment, idealize secondary structure, manually  
198 design sidechains, delete residues, insert residues, insert cutpoint, idealize peptide bond  
199 geometry. (F) Foldit players explore diverse structures that have no sequence or structural  
200 homology to natural proteins.  
201  
202  
203

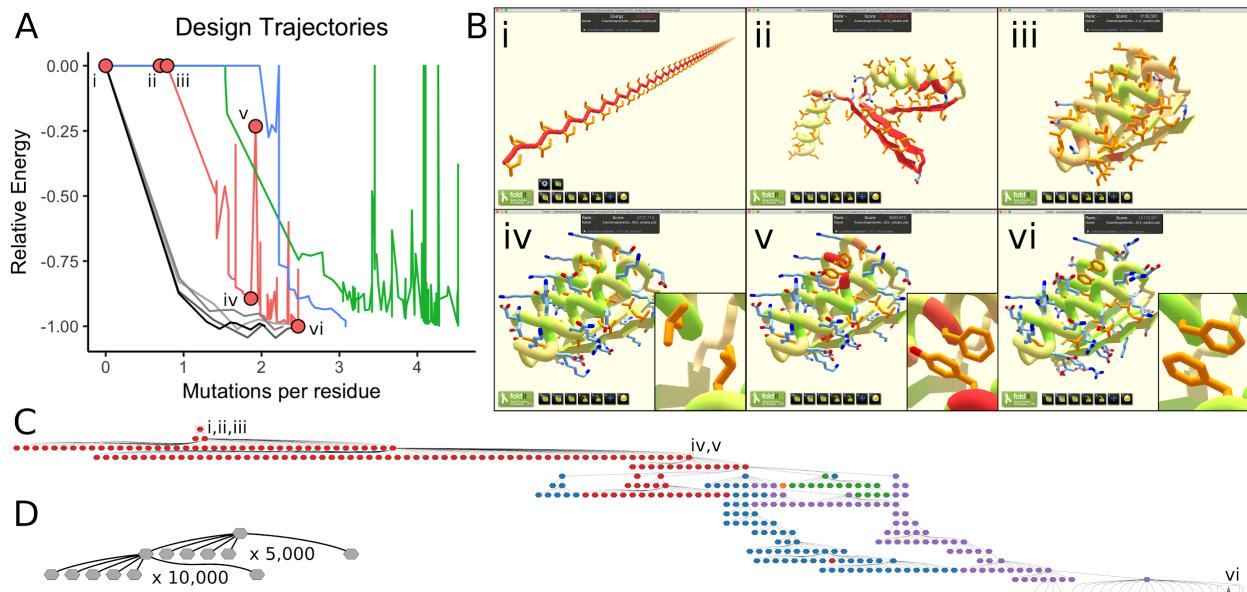
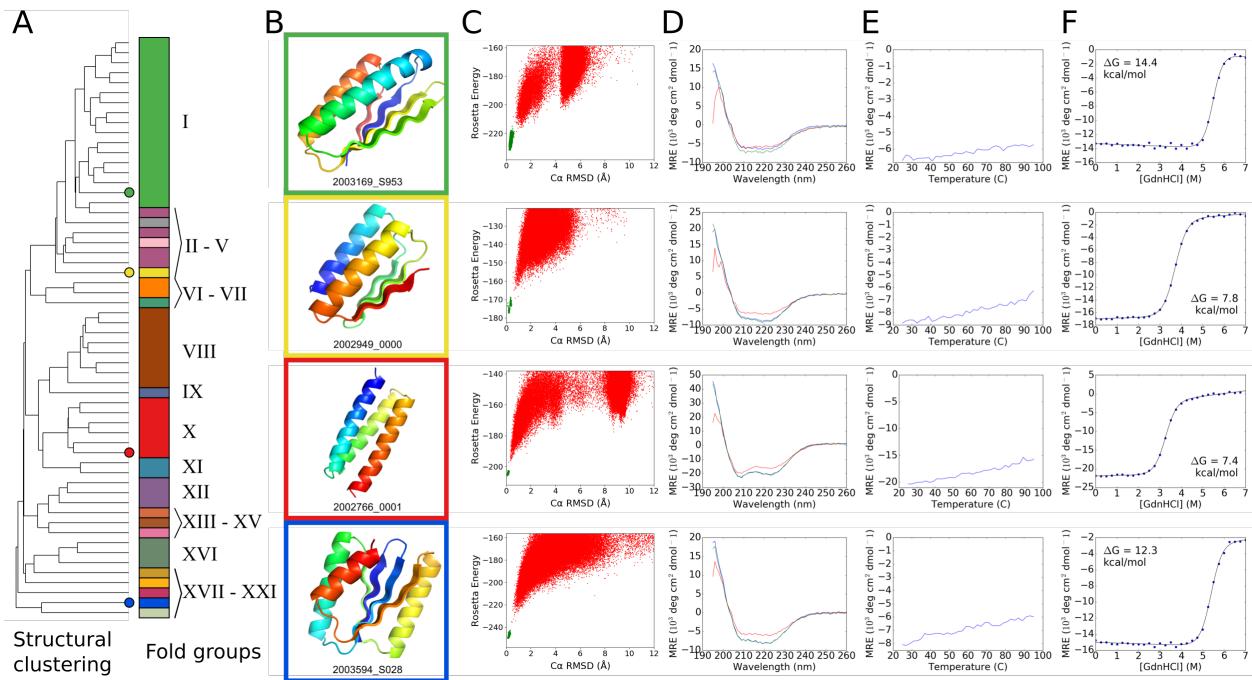
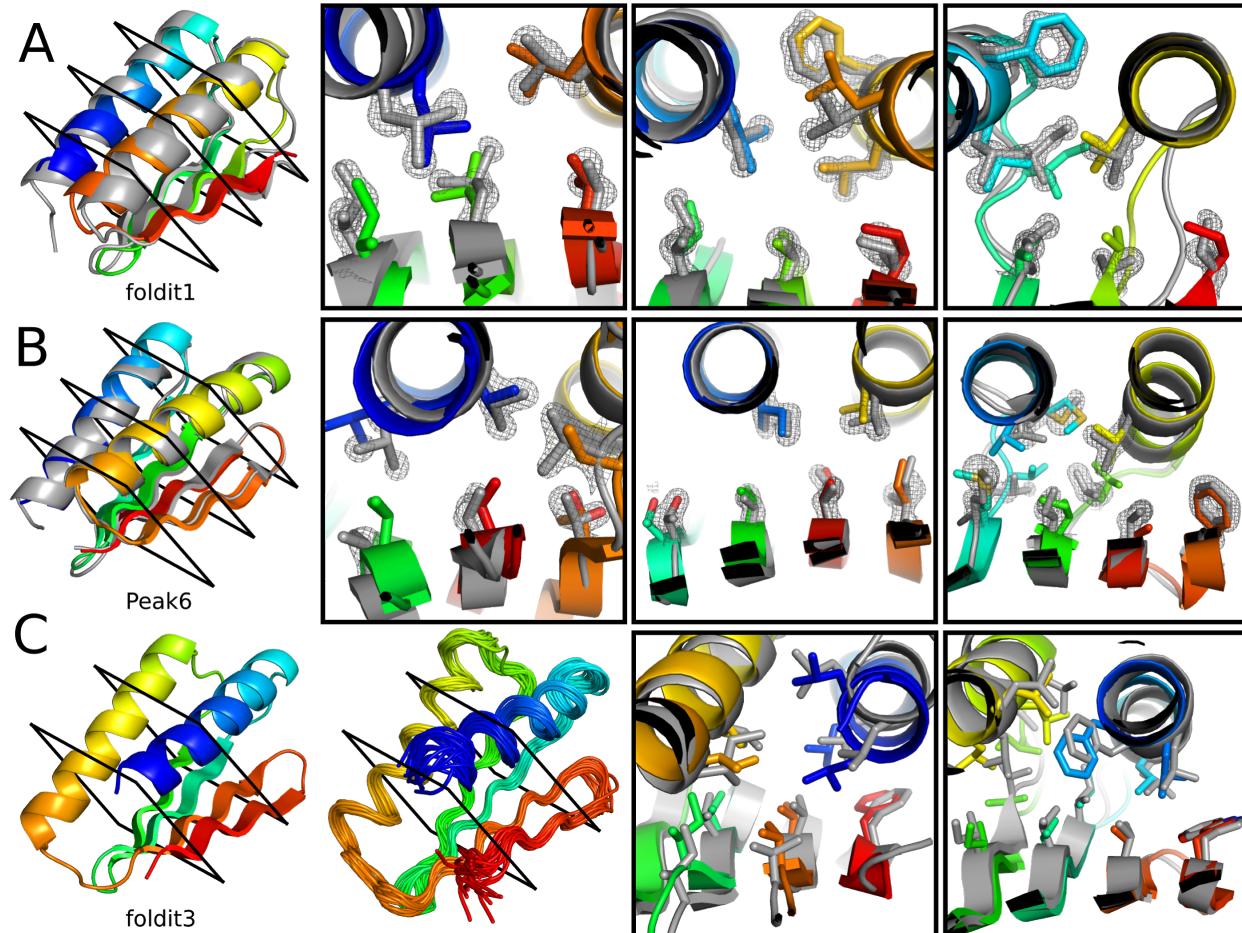


Figure 2. Comparison of Foldit player and automated design sampling strategies. (A) Single trajectories (ignoring abandoned branches) for three Foldit player-designed proteins in red (foldit1), blue (Peak6), and green (2003169\_S953); and design trajectories for four Rosetta-designed proteins in gray. The y-axis is the Rosetta energy rescaled so that the final design has a value of -1.00, and positive energies are shown as zero. Foldit players are willing to undergo large increases in energy to explore new regions; the Rosetta protocol in contrast has a limited ability to escape local energy minima. Red circles correspond to structures shown in (B). (B) Snapshots from the design trajectory of foldit1: (i) the initial extended chain of poly-isoleucine; (ii) development of secondary structure; (iii) development of folded tertiary structure; (iv) sequence design of folded structure, with inset showing favorable packing between two Leu sidechains at positions 13 and 45; (v) high-energy intermediate design, with inset showing redesign at positions 13 and 45, which results in steric clashes with the protein backbone; (vi) the final refined design, with inset showing favorable interactions between two Phe sidechains at positions 13 and 45. (C) The design strategy for foldit1 represented as a graph, showing all branch points where multiple design trajectories were spawned from a single intermediate. The final design was reached only after 17 branch points. Node colors correspond to five different cooperating Foldit players, and the final design is marked as a star. (D) Similar representation of a Rosetta design trajectory; there are only two branch points.



225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244

Figure 3. Structural characterization of Foldit player designed proteins. (A) Dendrogram showing all 58 folded Foldit player designs clustered by structural similarity ( $\text{TM-align}^{26}$ ), with colored circles highlighting the four designs characterized in (B-F). The stacked bars show the 21 different folds among the clustered designs (Extended Data Fig. 2). Fold XX (see design 2003594\_S028) is a new fold, previously unobserved in natural proteins. (B) Cartoon depiction of four select Foldit designs. (C) Rosetta@home *ab initio* calculations show that the sequence for each design has an energy landscape that is strongly funneled toward the design structure. Rosetta energy is on the y-axis and  $\text{Ca-RMSD}$  to the designed structure on the x-axis; points represent lowest energy structures sampled starting from an extended chain (red points), and starting from the Foldit design model (green points). (D) Circular dichroism (CD) spectra indicate that the designs adopt the expected secondary structure content in solution at 25°C (blue trace), when heated to 95°C (red trace), and when cooled again to 25°C (green trace). (E) CD mean residue ellipticity at 220 nm as temperature is increased from 25°C to 95°C; the designs do not temperature denature. (F) Cooperative unfolding during titration with guanidinium hydrochloride. Blue circles show CD mean residue ellipticity at 220 nm with increasing concentration of denaturant, and the black curve shows a two-state unfolding model fit to the data.  $\Delta G_{\text{unf}}$  values were determined by linear extrapolation using the fit model parameters<sup>27</sup>.



245  
 246 Figure 4. High-resolution structures of Foldit player designed proteins. (A) The foldit1 design  
 247 (fold V in Fig 3: 3  $\beta$ -strands with sheet order 1-2-3) model backbone (rainbow) aligns to the  
 248 crystal structure (gray) with Ca-RMSD of 1.1 Å. (B) The Peak6 design (fold III: 4 strands, sheet  
 249 order 1-2-4-3) model backbone (rainbow) aligns to the crystal structure (gray) with Ca-RMSD of  
 250 0.9 Å. Cross-sections show core residue sidechains, with the composite omit 2mFo-DFc map  
 251 contoured at 2.0  $\sigma$ . (C) The foldit3 design model (fold XVIII: 4 strands, sheet order 2-1-3-4) and  
 252 NMR ensemble. The design model aligns to the representative (medoid) NMR model with a Ca-  
 253 RMSD of 1.1 Å. Cross sections compare core side chains in the design model (rainbow) and a  
 254 representative NMR model (gray).  
 255  
 256  
 257

258     **Methods**  
259  
260     **Foldit protein design puzzles**  
261     Foldit puzzles were set up with a model poly-isoleucine in fully extended conformation, with  
262     fixed length ranging from 60 to 100 residues. Each puzzle was posted online for seven days,  
263     during which Foldit players competed to develop a protein model with the lowest energy, as  
264     calculated by the Rosetta energy function. Foldit puzzles used the talaris2013\_cart  
265     scorefunction with the following modifications: (1) the cart\_bonded scoreterm was upweighted  
266     (increased from 0.5 to 2.0) to ensure realistic bond lengths and angles as players cut and splice  
267     the backbone chain; (2) a penalty-only envsmooth scoreterm (weighted at 2.0) was added to  
268     supplement the Rosetta solvation treatment, and to discourage the design of buried polar and  
269     exposed nonpolar residues; (3) the reference energy of alanine was modified (increased to 3.0)  
270     to discourage the excessive design of alanine. See Supplementary Information for configuration  
271     files for all Foldit puzzles. Each Foldit puzzle was accompanied by a brief description, along with  
272     an explanation of any supplementary rules enforced in the puzzle. Design puzzles were  
273     accessible to all Foldit users; Foldit user registration is free and open to the public, at  
274     <http://fold.it>. Models were collected continuously as Foldit players worked on the puzzles, since  
275     the Foldit application automatically uploads the user's latest model to a server every 2-5  
276     minutes. This study was approved by the University of Washington Institutional Review Board,  
277     and informed consent for research participation was obtained from all Foldit users at the time of  
278     user registration.  
279  
280     **Protein design selection**  
281     After the end of each puzzle, we selected player models for further analysis as follows: First we  
282     selected the lowest-energy model from each of the 10 top-ranked groups, where independent  
283     players were treated as individual groups (designs named with suffix "0000-9"). Second, we  
284     selected the lowest-energy model from the 10 top-ranked solo players, which includes  
285     independent players as well as group members that developed a model without assistance from  
286     their group (suffix "s000-9"). Third, we visually inspected models that were flagged by Foldit  
287     players for special consideration, and selected any models that appeared plausible (suffix  
288     "S\*\*\*"). Last, we ranked and pruned the set of remaining models, by removing any models that  
289     align to a better-scoring model with Ca-RMSD less than 2.5 Å. We visually inspected the 50 top-  
290     ranked models in the pruned set and selected any models that appeared plausible (suffix "1001-  
291     50"). Models deemed "implausible" typically lacked secondary structure, contained buried polar  
292     residues, or included long stretches of completely polar residues. The sequences of selected  
293     models were subjected to Rosetta *ab initio* structure prediction<sup>14</sup>, using the distributed  
294     computing platform Rosetta@home. If *ab initio* predictions identified any decoy structures with  
295     energy comparable to (or lower than) the designed structure, or if *ab initio* predictions were  
296     unable to sample the designed structure, the design was rejected. All other designs were  
297     selected for experimental characterization. The majority of experimentally tested designs (96 of  
298     146) were top-ranked group or solo designs, which were selected "blindly" (without visual  
299     inspection). See supplementary data for structures and FASTA sequences of all tested designs.  
300  
301     **Protein expression and purification**

302 A 6x-His tag with TEV-cleavable linker (sequence 'MGHHHHHHGWSENLYFQGS') was  
303 prepended to the N-terminus of each design selected for experimental characterization.  
304 Plasmids containing the encoded genes were ordered from Genscript in pET15 (designs with  
305 prefix between 997258 and 1998925), or in pET21 (1998555-2002990), or from Twist in pET28  
306 (2003048-2003594) vectors. Plasmids were transformed into *E. coli* BL21 Star (DE3) cells  
307 (Invitrogen), and grown overnight in 4 mL Luria-Bertani medium (LB) with 50 µg/mL carbenicillin  
308 (for pET15, pET21 vectors) or 30 µg/mL kanamycin (for pET29). Overnight cultures were used  
309 to inoculate 0.5 L auto-induction media, and grown at 37 °C for 18 hours. Cultures were pelleted  
310 and resuspended in 25 mL lysis buffer (20 mM Tris pH 8.0, 300 mM NaCl, 1 mg/mL lysozyme,  
311 0.1 mg/mL DNase, 1 mM PMSF), and lysed by microfluidization. The cell lysate was pelleted  
312 and supernatant was filtered with a 0.22 µm filter before loading onto a 2 mL nickel affinity  
313 gravity column. Protein bound to the column was washed with 20 mL wash buffer (20 mM Tris  
314 pH 8.0, 500 mM NaCl, 30 mM imidazole) and eluted in 10 mL elution buffer (20 mM Tris pH 8.0,  
315 500 mM NaCl, 250 mM imidazole). Purified protein was dialyzed into TBS (20 mM Tris pH 8.0,  
316 300 mM NaCl) at 4°C overnight to remove imidazole and further purified by gel filtration on an  
317 AKTAexpress (GE Healthcare) with a Superdex S75 10/300 GL column (GE Healthcare). For  
318 proteins containing cysteine, dialysis and gel filtration were carried out in TBS with 1 mM TCEP.  
319 Protein expression and solubility was confirmed by SDS-PAGE of samples from lysate pellet  
320 and supernatant. All purified proteins were verified by mass spectrometry.

321

### 322 **Circular dichroism**

323 Purified protein was dialyzed into 50 mM sodium phosphate pH 7.4 at 4°C overnight (plus 500  
324 µM TCEP for proteins containing cysteine). All circular dichroism data were collected on an  
325 AVIV Model 420 spectrometer. Far UV spectra and temperature melts were measured with 11-  
326 62 µM protein in a quartz cuvette with path length of 1 mm. Protein concentration was  
327 determined by absorbance at 280 nm using a NanoDrop spectrophotometer (Thermo Scientific),  
328 using predicted extinction coefficients. Wavelength spectra were measured between 195 and  
329 260 nm at 25°C, 95°C, and again after cooling to 25°C. For temperature melts, ellipticity at 220  
330 nm was monitored as temperature increased from 25°C to 95°C, in increments of 2°C. Chemical  
331 titrations were carried out with 1.0-21 µM protein in a quartz cuvette with path length of 10 mm.  
332 Ellipticity at 220 nm was monitored at concentrations of guanidinium chloride increasing from 0  
333 to 7 M, in increments of 0.25 M. Denaturation curves were fitted with non-linear regression to  
334 two-state unfolding model with six parameters: the folding free energy, m-value, and slope and  
335 y-intercept for baseline curves<sup>28</sup>.

336

### 337 **X-ray crystallography**

338 Prior to x-ray crystallography, the N-terminal 6x-His tag was cleaved from protein samples by  
339 incubation with 250 µg TEV protease at 25°C for four hours in 20 mM Tris pH 8.0, 300 mM  
340 NaCl, 1 mM DTT. The reaction product was dialyzed into TBS overnight at 4°C to remove DTT  
341 and flowed over a 2 mL metal affinity gravity column to remove TEV protease and residual  
342 histidine tag. The cleaved protein was further purified by gel filtration as described above.  
343 Purified protein was concentrated to 20-100 mg/mL in 20 mM Tris pH 8.0, 300 mM NaCl.  
344 Crystallization screening was carried out with a variety of 96-condition spare matrix suites  
345 available from Qiagen or Hampton Research. A Mosquito Crystal nanoliter robot (TTP Labtech)

346 was used to prepare screens in 3-well sitting drop plates, with 200 nL drops and  
347 protein:precipitant ratios of 1:1, 1:2, and 2:1.

348

349 foldit1 was crystallized at 20 mg/mL in 50 mM HEPES pH 7.5, 0.2 M potassium chloride, 35%  
350 v/v pentaerythritol propoxylate. Crystals were flash-frozen in liquid nitrogen without further cryo-  
351 protection.

352

353 Peak6 was crystallized at 40 mg/mL in 0.1 M sodium acetate pH 4.5, 0.2 M lithium sulfate, 50%  
354 w/v PEG 400. Crystals were briefly soaked in mother liquor plus 20% PEG 200, then flash  
355 frozen in liquid nitrogen.

356

357 X-ray diffraction datasets were collected at the Advanced Light Source (Berkeley, CA). Data  
358 was processed with HKL2000<sup>31</sup>. Crystal structures were solved by molecular replacement with  
359 Phaser<sup>29</sup>, using the backbone of the original designed model with sidechains truncated to the  
360 beta carbon. Models were built and refined in iterative cycles using Coot and PHENIX<sup>30,31</sup>.  
361 Diffraction data and refinement statistics are listed in Extended Data Table 2.

362

363 **NMR spectroscopy**

364 NMR studies were performed using uniformly <sup>15</sup>N, <sup>13</sup>C-enriched protein samples. Synthetic  
365 genes were obtained from Genscript already incorporated into plasmid pET15TEV\_NESG,  
366 which includes a N-terminal 6xHis purification tag, followed by a TEV protease cleavage site  
367 (sequence 'MGHHHHHGWSSENLYFQGS'). *E. coli* BL21(DE3) cells harboring plasmid  
368 pET15TEV\_NESG-foldit3 were grown in 1L MJ9 minimal media<sup>32</sup>, supplemented with 100 µg/ml  
369 ampicillin at 37 °C. In order to produce uniformly <sup>15</sup>N and <sup>13</sup>C enriched protein samples, 1g / L  
370 <sup>15</sup>NH<sub>4</sub>-salts and 2g / L U-<sup>13</sup>C glucose were added as sole a nitrogen and a carbon sources,  
371 respectively. When O.D.<sub>600</sub> reached around 0.5 units, the culture was transferred to 18 °C, and  
372 the protein production was induced by addition of 1 mM IPTG. After overnight incubation, the  
373 cells were collected and resuspended in 20 ml binding buffer (20 mM Tris-HCl pH 8.0, 500 mM  
374 NaCl and 20 mM imidazole). After passing the cells through 900-1000 psi French press twice,  
375 cell debris were removed by 10,000 rpm for 30 min. The supernatant was further spun down at  
376 40,000 rpm for 1hr. The obtained supernatant (soluble fraction) was mixed with 1 ml of Ni-resin  
377 and incubated at 4 °C for 1 hr. The non-specific binding proteins were removed by 20 mL  
378 binding buffer and washing buffer (20 mM Tris-HCl pH 8.0, 500 mM NaCl and 50 mM imidazole)  
379 and the target protein was eluted by 5 mL elution buffer (20 mM Tris-HCl pH 8.0, 500 mM NaCl  
380 and 300 mM imidazole). The protein was dialyzed against GF buffer (20 mM Tris-HCl pH 8.0,  
381 100 mM NaCl) for overnight and gel filtration was carried out using AKTA express with high-load  
382 26/600 Superdex 200 pg column. Homogeneity (> 97%) was validated by SDS polyacrylamide  
383 gel electrophoresis. The purified protein was dialyzed against 20 mM potassium phosphate (pH  
384 6.5), and the protein concentration was adjusted to between 0.3-0.4 mM for NMR studies.

385

386 All NMR spectra were recorded at 25 °C using cryogenic NMR probes. All NMR data were  
387 collected on the Bruker AVANCE III 600 MHz spectrometers and processed using the program  
388 NMRPipe<sup>33</sup>, and analyzed using the programs SPARKY and XEASY<sup>34</sup>. Spectra were referenced  
389 to external DSS. Sequence-specific resonance assignments were determined using AutoAssign

software together with interactive manual analysis, as described previously<sup>35</sup>. Backbone dihedral angle constraints were derived from the chemical shifts using the program TALOS\_N<sup>36</sup> for residues located in well-defined secondary structure elements. The programs ASDP<sup>37</sup> and CYANA<sup>38,39</sup> were used to automatically assign NOEs and to calculate structures. RPF analysis<sup>37,40</sup> was used in parallel to guide iterative cycles of noise/artifact peak removal, peak picking, and NOESY peak assignments. The 20 conformers with the lowest target CYANA function value were then refined in explicit water<sup>41</sup> using the program CNS<sup>42</sup>. The structural statistics and global structure quality factors (Extended Data Table 3) including Verify3D<sup>43</sup>, ProsaII<sup>44</sup>, PROCHECK<sup>45</sup>, and MolProbity<sup>46</sup> raw and statistical Z-scores were computed using the PSVS<sup>47</sup> 1.5 and PDBStat<sup>48</sup> software packages. The global goodness-of-fit of the final structure ensembles with the NOESY peak list data, the NMR DP score, was determined using the RPF analysis program<sup>40</sup>.

402

### 403 **Code Availability**

Because Foldit crowdsourcing relies on regulated, fair competition between participants, the source code of the Foldit user interface is not open. The underlying Rosetta macromolecular modeling suite (<https://www.rosettacommons.org>) is freely available to academic and non-commercial users, and commercial licenses are available via the University of Washington CoMotion Express License Program. Analysis scripts used in this paper are available in the Supplementary Information.

410

### 411 **Data Availability**

The atomic coordinates of foldit1 and Peak6 crystal structures, and the foldit3 NMR structure, have been deposited in the RCSB Protein Database with accession numbers 6MRR, 6MRS, and 6MSP, respectively. Chemical shift and NOESY peak list data for foldit3 were deposited in the Biological Magnetic Resonance Bank (BMRB ID 30527).

416

### 417 **References**

1. Lintott, C. J. *et al.* Galaxy Zoo: morphologies derived from visual inspection of galaxies from the Sloan Digital Sky Survey. *Monthly Notices of the Royal Astronomical Society* **389**, 1179–1189 (2008).
2. Kim, J. S. *et al.* Space-time wiring specificity supports direction selectivity in the retina. *Nature* **509**, 331–336 (2014).
3. Kawrykow, A. *et al.* Phylo: A Citizen Science Approach for Improving Multiple Sequence Alignment. *PLoS ONE* **7**, (2012).
4. Lee, J. *et al.* RNA design rules from a massive open laboratory. *Proceedings of the National Academy of Sciences* **111**, 2122–2127 (2014).
5. Cooper, S. *et al.* Predicting protein structures with a multiplayer online game. *Nature* **466**, 756–760 (2010).
6. Epstein, C. J., Goldberger, R. F. & Anfinsen, C. B. The Genetic Control of Tertiary Protein Structure: Studies With Model Systems. *Cold Spring Harbor Symposia on Quantitative Biology* **28**, 439–449 (1963).
7. Lin, Y.-R. *et al.* Control over overall shape and size in de novo designed proteins. *Proceedings of the National Academy of Sciences* (2015).

- 434 8. Huang, P.-S., Boyken, S. E. & Baker, D. The coming of age of de novo protein design.  
435 *Nature* **537**, 320–327 (2016).
- 436 9. Marcos, E. *et al.* Principles for designing proteins with cavities formed by curved  $\beta$   
437 sheets. *Science* **355**, 201–206 (2017).
- 438 10. Dou, J. *et al.* De novo design of a fluorescence-activating  $\beta$ -barrel. *Nature* (2018).
- 439 11. Alford, R. F. *et al.* The Rosetta All-Atom Energy Function for Macromolecular Modeling  
440 and Design. *J Chem Theory Comput* **13**, 3031–3048 (2017).
- 441 12. Khatib, F. *et al.* Crystal structure of a monomeric retroviral protease solved by protein  
442 folding game players. *Nat Struct Mol Biol* **18**, 1175–1177 (2011).
- 443 13. Eiben, C. B. *et al.* Increased Diels-Alderase activity through backbone remodeling  
444 guided by Foldit players. *Nature Biotechnology* **30**, 190–192 (2012).
- 445 14. Rohl, C. A., Strauss, C. E. M., Misura, K. M. S. & Baker, D. Protein structure prediction  
446 using Rosetta. *Meth. Enzymol.* **383**, 66–93 (2004).
- 447 15. Blout, E. R. & Idelson, M. Compositional Effects on the Configuration of Water-soluble  
448 Polypeptide Copolymers of L-Glutamic Acid and L-Lysine. *Journal of the American  
449 Chemical Society* **80**, 4909–4913 (1958).
- 450 16. Doty, P., Imahori, K. & Klemperer, E. The solution properties and configurations of a  
451 polyampholytic polypeptide: copoly-L-lysine-L-glutamic acid. *Proceedings of the National  
452 Academy of Sciences* **44**, 424–431 (1958).
- 453 17. Ghosh, K. & Dill, K. A. Theory for Protein Folding Cooperativity: Helix Bundles. *J. Am.  
454 Chem. Soc.* **131**, 2306–2312 (2009).
- 455 18. Koga, N. *et al.* Principles for designing ideal protein structures. *Nature* **491**, 222–227  
456 (2013).
- 457 19. Regan, L. & DeGrado, W. Characterization of a helical protein designed from first  
458 principles. *Science* **241**, 976–978 (1988).
- 459 20. Harbury, P. B., Plecs, J. J., Tidor, B., Alber, T. & Kim, P. S. High-resolution protein  
460 design with backbone freedom. *Science* **282**, 1462–1467 (1998).
- 461 21. Thomson, A. R. *et al.* Computational design of water-soluble alpha-helical barrels.  
462 *Science* **346**, 485–488 (2014).
- 463 22. Jacobs, T. M. *et al.* Design of structurally distinct proteins using strategies inspired by  
464 evolution. *Science* **352**, 687–690 (2016).
- 465 23. Ramachandran, G. N. & Sasisekharan, V. Conformation of Polypeptides and Proteins.  
466 *Advances in Protein Chemistry* **23**, 283–437 (1968).
- 467 24. Chen, V. B. *et al.* MolProbity: all-atom structure validation for macromolecular  
468 crystallography. *Acta Crystallogr. D Biol. Crystallogr.* **66**, 12–21 (2010).
- 469 25. Montelione, G. T. *et al.* Recommendations of the wwPDB NMR Validation Task Force.  
470 *Structure* **21**, 1563–1570 (2013).
- 471 26. Zhang, Y. & Skolnick, J. TM-align: a protein structure alignment algorithm based on the  
472 TM-score. *Nucleic Acids Res.* **33**, 2302–2309 (2005).
- 473 27. Santoro, M. M. & Bolen, D. W. Unfolding free energy changes determined by the linear  
474 extrapolation method. 1. Unfolding of phenylmethanesulfonyl  $\alpha$ -chymotrypsin using  
475 different denaturants. *Biochemistry* **27**, 8063–8068 (1988).
- 476 28. Otwinowski, Z. & Minor, W. Processing of X-ray diffraction data collected in oscillation  
477 mode. *Methods Enzymol.* **276**, 307–326 (1997).

- 478 29. McCoy, A. J. *et al.* Phaser crystallographic software. *J Appl Crystallogr* **40**, 658–674  
479 (2007).
- 480 30. Emsley, P., Lohkamp, B., Scott, W. G. & Cowtan, K. Features and development of Coot.  
481 *Acta Crystallogr. D Biol. Crystallogr.* **66**, 486–501 (2010).
- 482 31. Afonine, P. V. *et al.* Towards automated crystallographic structure refinement with  
483 phenix.refine. *Acta Crystallogr. D Biol. Crystallogr.* **68**, 352–367 (2012).
- 484 32. Jansson M, Li YC, Jendeberg L, Anderson S, Montelione GT, Nilsson B. High-level  
485 production of uniformly  $^{15}\text{N}$ - and  $^{13}\text{C}$ -enriched fusion proteins in Escherichia coli. *J  
486 Biomol NMR* **7**, 131–141 (1996).
- 487 33. Delaglio, F. *et al.* Nmrpipe - a Multidimensional Spectral Processing System Based on  
488 Unix Pipes. *J Biomol NMR* **6**, 277–293 (1995).
- 489 34. Bartels, C., Xia, T. H., Billeter, M., Guntert, P. & Wuthrich, K. The Program Xeasy for  
490 Computer-Supported Nmr Spectral-Analysis of Biological Macromolecules. *J Biomol  
491 NMR* **6**, 1–10 (1995).
- 492 35. Liu, G. H. *et al.* NMR data collection and analysis protocol for high-throughput protein  
493 structure determination. *Proceedings of the National Academy of Sciences of the United  
494 States of America* **102**, 10487–10492 (2005).
- 495 36. Shen, Y., Delaglio, F., Cornilescu, G. & Bax, A. TALOS+: a hybrid method for predicting  
496 protein backbone torsion angles from NMR chemical shifts. *J Biomol NMR* **44**, 213–223,  
497 doi:10.1007/s10858-009-9333-z (2009).
- 498 37. Huang, Y. J., Tejero, R., Powers, R. & Montelione, G. T. A topology-constrained  
499 distance network algorithm for protein structure determination from NOESY data.  
500 *Proteins* **62**, 587–603, doi:10.1002/prot.20820 (2006).
- 501 38. Guntert, P., Mumenthaler, C. & Wuthrich, K. Torsion angle dynamics for NMR structure  
502 calculation with the new program DYANA. *Journal of Molecular Biology* **273**, 283–298  
503 (1997).
- 504 39. Herrmann, T., Guntert, P. & Wuthrich, K. Protein NMR structure determination with  
505 automated NOE assignment using the new software CANDID and the torsion angle  
506 dynamics algorithm DYANA. *Journal of Molecular Biology* **319**, 209–227 (2002).
- 507 40. Huang, Y. J., Powers, R. & Montelione, G. T. Protein NMR recall, precision, and F-  
508 measure scores (RPF scores): Structure quality assessment measures based on  
509 information retrieval statistics. *Journal of the American Chemical Society* **127**, 1665–  
510 1674 (2005).
- 511 41. Linge, J. P., Williams, M. A., Spronk, C. A., Bonvin, A. M. & Nilges, M. Refinement of  
512 protein structures in explicit solvent. *Proteins* **50**, 496–506, doi:10.1002/prot.10299  
513 (2003).
- 514 42. Brunger, A. T. *et al.* Crystallography & NMR system: A new software suite for  
515 macromolecular structure determination. *Acta Crystallographica Section D-Biological  
516 Crystallography* **54**, 905–921 (1998).
- 517 43. Luthy, R., Bowie, J. U. & Eisenberg, D. Assessment of protein models with three-  
518 dimensional profiles. *Nature* **356**, 83–85, doi:10.1038/356083a0 (1992).
- 519 44. Sippl, M. J. Recognition of errors in three-dimensional structures of proteins. *Proteins*  
520 **17**, 355–362, doi:10.1002/prot.340170404 (1993).

- 521        45. Laskowski, R. A., Macarthur, M. W., Moss, D. S. & Thornton, J. M. Procheck - a  
522        Program to Check the Stereochemical Quality of Protein Structures. *Journal of Applied*  
523        *Crystallography* **26**, 283-291 (1993).
- 524        46. Word, J. M., Bateman, R. C., Presley, B. K., Lovell, S. C. & Richardson, D. C. Exploring  
525        steric constraints on protein mutations using MAGE/PROBE. *Protein Science* **9**, 2251-  
526        2259 (2000).
- 527        47. Bhattacharya, A., Tejero, R. & Montelione, G. T. Evaluating protein structures  
528        determined by structural genomics consortia. *Proteins* **66**, 778-795,  
529        doi:10.1002/prot.21165 (2007).
- 530        48. Tejero, R., Snyder, D., Mao, B., Aramini, J.M., Montelione, G.T. PDBStat: A universal  
531        restraint converter and restraint analysis software package for protein NMR. *J. Biomol.*  
532        *NMR* **56**, 337-351 (2013).
- 533

534 **Acknowledgements**

535 We thank all Foldit players for their gameplay contributions, and for feedback offered on the  
536 https://fold.it website. We thank A. Kang, S.A. Rettie, C. Chow, and L. Carter for help with  
537 experiments; D. Alonso, L. Goldschmidt, P. Vecchiato, D. Kim for computer support; and  
538 Rosetta@home (<https://boinc.bakerlab.org>) volunteers for computing resources. We thank G.  
539 Rocklin, V. Mulligan, and other members of the Baker lab for discussions. This material is based  
540 upon work supported by the National Science Foundation Graduate Research Fellowship under  
541 Grant No. DGE-1256082, and National Institutes of Health Grant 1R01 GM120574 (to G.T.M.).  
542 The ALS-ENABLE beamlines are supported in part by the National Institutes of Health, National  
543 Institute of General Medical Sciences, grant P30 GM124169-01. The Advanced Light Source is  
544 a DOE User Facility under Contract No. DE-AC02-05CH11231. This material is based upon  
545 work supported by the National Science Foundation under grant no. 1629879. This work was  
546 supported by National Institutes of Health grant 1UH2CA203780.

547

548 **Author contributions**

549 B.K., Z.P., F.K., S.C., and D.B. designed the study.  
550 B.K., J.F., T.H., A.F., D.A.S., and S.C. developed Foldit software tools.  
551 A. Boykov, R.D.E., S.K., L.W., and Foldit Players designed all proteins.  
552 B.K., F.K., A.F., and A. Bauer analyzed Foldit player designs.  
553 B.K. performed biophysical characterization.  
554 B.K. and M.J.B. determined crystal structures.  
555 G.L., Y.I., and G.T.M. determined the NMR structure.  
556 B.K. and D.B. wrote the manuscript with input from all authors.

557

558 **Competing interests**

559 The authors declare no competing interests. G.T.M. is a co-founder of Nexomics Biosciences,  
560 Inc.

561

Extended Data Table 1. Foldit player-designed proteins selected for experimental testing

Design ID	Designers	Nearest Sequence Homolog			Nearest Structural Homolog		Experimental Characterization			
		BLAST Score	BLAST E-value	BLAST Hit	TM-align Score	TM-align Hit	Expressed	Soluble	Monomeric	Structured
997258_0001	PLAYER_2,MaartenDesnouck,MurloW	-	-	-	0.806	5cwoA	Yes	Yes	Yes	No
997258_0004	Timo van der Laan	-	-	-	<b>0.814</b>	<b>5cwpA</b>	Yes	Yes	Yes	Yes
997383_S346	frood66,PLAYER_16,PLAYER_10	35.4	5.5	KFV13184.1	0.685	3vf0A	Yes	Yes	No	-
997523_0000	PLAYER_13,MurloW	-	-	-	0.877	4tqlA	Yes	Yes	No	-
997523_0003	caglar	<b>48.9</b>	<b>2.00E-04</b>	<b>XP_015806396.1</b>	<b>0.844</b>	<b>4e40A</b>	Yes	Yes	Yes	Yes
997523_0005	vakobo,Grom,PLAYER_14	38.9	0.67	XP_015806396.1	0.811	5j0lE	Yes	Yes	No	-
997523_0006	Timo van der Laan	40	0.25	GBE60965.1	0.819	4tqlA	Yes	Yes	No	-
997523_0008	eikem	<b>35.4</b>	<b>8.7</b>	<b>WP_044163677.1</b>	<b>0.782</b>	<b>5xqjA</b>	Yes	Yes	Yes	Yes
997523_1003	nemo7731	<b>43.5</b>	<b>0.014</b>	<b>WP_055520104.1</b>	<b>0.740</b>	<b>5j0lA</b>	Yes	Yes	Yes	Yes
997523_1040	PLAYER_6	-	-	-	<b>0.811</b>	<b>5k7vA</b>	Yes	Yes	Yes	Yes
997791_1027	PLAYER_6	-	-	-	<b>0.820</b>	<b>2ojqA</b>	Yes	Yes	Yes	Yes
997915_0000	Galaxie,PLAYER_2,karstenw	<b>36.2</b>	<b>2.4</b>	<b>YP_009282838.1</b>	<b>0.813</b>	<b>5k7vA</b>	Yes	Yes	Yes	Yes
1998469_0000	PLAYER_15,retiredmichael	-	-	-	0.686	4p2fA	No	-	-	-
1998520_0002	Galaxie,PLAYER_23	35.8	2.8	XP_011658777.1	0.703	3w40A	Yes	Yes	No	-
1998555_1041	frood66	35	7.4	WP_102588688.1	0.631	2kptA	Yes	No	-	-
1998925_s005	PLAYER_9	-	-	-	0.642	3dyjB	No	-	-	-
1998925_s008	PLAYER_18	<b>36.6</b>	<b>1.4</b>	<b>WP_013842421.1</b>	<b>0.789</b>	<b>4fsxA</b>	Yes	Yes	Yes	Yes
2000240_0002	retiredmichael,PLAYER_15,LociOilng,smilingone,PLAYER_21	37.4	0.78	WP_047150276.1	0.719	1vdwA	No	-	-	-
2000240_s003	spvincnt	33.9	5.8	XP_003761380.1	0.600	4x00A	Yes	Yes	No	-
2000485_1070	spvincnt	-	-	-	0.595	4dlqA	No	-	-	-
2000518_0000	MurloW,PLAYER_9	-	-	-	0.634	4xevD	No	-	-	-
2000518_0003	PLAYER_22	-	-	-	0.610	m572A	Yes	No	-	-
2000518_S468	eusair	-	-	-	0.593	4g0hA	No	-	-	-
2000665_1003	mbinfield,Bruno Kestemont	-	-	-	0.569	4acjA	Yes	Yes	No	-
2001044_0000	MurloW	<b>35.8</b>	<b>2.2</b>	<b>XP_018018081.1</b>	<b>0.746</b>	<b>3d6kA</b>	Yes	Yes	Yes	Yes
2002089_0001	Galaxie,Susume	35	5.7	WP_117139598.1	0.705	4p1xF	No	-	-	-
2002089_1029	Galaxie,Susume	-	-	-	0.686	4ok4KA	Yes	Yes	No	-
2002243_1016	Susume	<b>36.2</b>	<b>2.1</b>	<b>PKR98267.1</b>	<b>0.656</b>	<b>2qzgA</b>	Yes	Yes	Yes	Yes
2002290_S122	fiendish_ghoul	79.7†	1.00E-17†	XP_001631727.1†	0.914†	2pig8†	No	-	-	-
2002308_0000	Mark-,PLAYER_3,Bletchley Park,PLAYER_1	<b>34.3</b>	<b>4.2</b>	<b>KKU76692.1</b>	<b>0.654</b>	<b>3rk0C</b>	Yes	Yes	Yes	Yes
2002308_0005	frood66,actiasluna,Mike Lewis	36.2	2.6	XP_022288518.1	0.622	2kt9A	Yes	No	-	-
2002308_S695	Susume	38.9	0.29	XP_012556696.1	0.689	3pg5A	Yes	No	-	-
2002334_0005	fiendish_ghoul	37.4	1	XP_016366754.1	0.637	5f1cB	Yes	Yes	No	-
2002376_0000	Mark-,Bletchley Park	36.2	4.3	CX16261.1	0.683	1fsaA	No	-	-	-
2002376_0003	PLAYER_21,LociOilng	45.4	0.002	XP_013073713.1	0.650	4he8G	Yes	Yes	Yes	No
2002469_0001	Bruno Kestemont,gloverd,Scopper	37	1.8	PIW76571.1	0.657	1xioA	No	-	-	-
2002469_S848	fiendish_ghoul	45.1	0.003	WP_075688920.1	0.695	5kilA	No	-	-	-
2002486_0006	fiendish_ghoul	-	-	-	0.685	2q1fA	Yes	No	-	-
2002486_1012	Mark-	34.7	7.9	WP_068417976.1	0.646	4zg4E	No	-	-	-
2002486_1048	fiendish_ghoul	35	4.8	WP_096386863.1	0.617	2bkaA	Yes	Yes	No	-
2002544_0000	Mark-,PLAYER_3	-	-	-	0.645	5ms2A	Yes	Yes	Yes	No
2002553_0000	Mark-,Bletchley Park	-	-	-	0.628	3dcpA	Yes	Yes	No	-
2002553_s003	Susume	<b>36.2</b>	<b>2</b>	<b>PYT68698.1</b>	<b>0.744</b>	<b>2g0lA</b>	Yes	Yes	Yes	Yes
2002565_0002	Galaxie,tokens	35.4	8.9	XP_012770361.1	0.666	3m1cA	Yes	Yes	-	-
2002590_0001	Galaxie,tokens	34.3	9	WP_077438279.1	0.726	3dd0A	Yes	Yes	-	-
2002590_S567	tokens	36.2	1.8	WP_074948256.1	0.738	4wyab	No	-	-	-
2002613_0004	PLAYER_25	36.2	2.7	XP_021195829.1	0.752	5iz3A	No	-	-	-
2002713_0000	retiredmichael,smilingone	-	-	-	<b>0.862</b>	<b>3rh3A</b>	Yes	Yes	Yes	Yes
2002713_0004	Mark-,gitwut,Bletchley Park	-	-	-	<b>0.742</b>	<b>3fajA</b>	Yes	Yes	Yes	Yes
2002713_1006	Mark-,PLAYER_3	<b>38.9</b>	<b>0.19</b>	<b>WP_044473091.1</b>	<b>0.799</b>	<b>5nxqA</b>	Yes	Yes	Yes	Yes
2002745_0000	Mark-,Bletchley Park,georg137	<b>34.7</b>	<b>5.8</b>	<b>GBB96165.1</b>	<b>0.819</b>	<b>3ripA</b>	Yes	Yes	Yes	Yes
2002745_0001	PLAYER_4	-	-	-	<b>0.873</b>	<b>5aqtB</b>	Yes	Yes	Yes	Yes
2002745_0003	Galaxie,PLAYER_2	-	-	-	<b>0.831</b>	<b>5cwiA</b>	Yes	Yes	Yes	Yes
2002745_0004	mirp,Bruno Kestemont,Paulo Roque	-	-	-	<b>0.807</b>	<b>4uy3A</b>	Yes	Yes	Yes	Yes
2002745_0008	Madde,kabubi	-	-	-	<b>0.849</b>	<b>2okuA</b>	Yes	Yes	Yes	Yes
2002766_0000	Mark-,Bletchley Park	<b>37.7</b>	<b>0.2</b>	<b>WP_116244417.1</b>	<b>0.806</b>	<b>1y4cA</b>	Yes	Yes	Yes	Yes
2002766_0001	LociOilng	<b>36.2</b>	<b>2</b>	<b>XP_013096440.1</b>	<b>0.839</b>	<b>1s94A</b>	Yes	Yes	Yes	Yes
2002766_0002	PLAYER_25	-	-	-	<b>0.813</b>	<b>4tqlA</b>	Yes	Yes	Yes	Yes
2002766_0003	Galaxie,tokens	-	-	-	0.846	4iv6A	Yes	Yes	Yes	No
2002766_0004	actiasluna,PLAYER_16,Blipperman	-	-	-	0.903	5cwmA	Yes	Yes	No	-
2002766_0006	PLAYER_17	-	-	-	0.853	4hwha	Yes	Yes	No	-
2002787_0005	fiendish_ghoul	35	6.2	WP_083480128.1	0.813	4kyzA	Yes	Yes	Yes	No
2002877_S005	gitwut	35.8	5.2	KKS40574.1	0.693	3tv9A	Yes	Yes	No	-
2002922_1013	Mark-,PLAYER_11,Bletchley Park	-	-	-	0.690	5ms2A	Yes	No	-	-
2002922_1018	Hollinas,Bruno Kestemont,Scopper	-	-	-	0.701	6eqtD	Yes	No	-	-
2002922_s004	tokens	<b>36.2</b>	<b>1.3</b>	<b>EXM13361.1</b>	<b>0.710</b>	<b>2ejxA</b>	Yes	Yes	Yes	Yes
2002949_0000	Galaxie,Susume,PLAYER_12	<b>35.4</b>	<b>1.6</b>	<b>WP_019366325.1</b>	<b>0.712</b>	<b>3rf0A</b>	Yes	Yes	Yes	Yes
a.k.a. foldit1										
2002949_0007	fiendish_ghoul	34.3	8.4	PJE68342.1	0.859	2ebbA	Yes	No	-	-
2002990_0006	fiendish_ghoul	<b>37.4</b>	<b>1.7</b>	<b>RDI84383.1</b>	<b>0.803</b>	<b>4hhua</b>	Yes	Yes	Yes	Yes
2002990_1031	fiendish_ghoul	-	-	-	0.734	4clfA	Yes	Yes	No	-
2002990_1039	fiendish_ghoul	<b>37</b>	<b>2</b>	<b>WP_113960216.1</b>	<b>0.713</b>	<b>4zhqD</b>	Yes	Yes	Yes	Yes
2002990_s006	retiredmichael	-	-	-	0.717	3ej0A	Yes	No	-	-
2003048_0003	PLAYER_11,gitwut	-	-	-	<b>0.738</b>	<b>5adkX</b>	Yes	Yes	Yes	Yes
2003048_0005	Bruno Kestemont,Scopper	35.4	3.8	WP_054386055.1	0.749	2bjIA	No	-	-	-
2003048_1024	actiasluna,PLAYER_5	-	-	-	0.689	3lg0A	Yes	Yes	No	-
2003048_1050	fiendish_ghoul	-	-	-	0.715	1k8kf	Yes	Yes	No	-
2003048_S697	tokens	37.4	0.62	WP_071355392.1	0.683	5izvA	No	-	-	-
2003048_s009	ZeroLeak7	-	-	-	0.744	513sD	No	-	-	-
2003169_S953	tokens	-	-	-	<b>0.767</b>	<b>4wjba</b>	Yes	Yes	Yes	Yes
2003169_s001	tokens	-	-	-	<b>0.721</b>	<b>1e2sP</b>	Yes	Yes	Yes	Yes
2003169_s008	Susume	-	-	-	0.694	2125A	Yes	No	-	-
2003205_0000	PLAYER_12,tokens	<b>36.6</b>	<b>2.3</b>	<b>WP_107152233.1</b>	<b>0.735</b>	<b>5cw9A</b>	Yes	Yes	Yes	Yes
2003205_0002	fiendish_ghoul	-	-	-	<b>0.888</b>	<b>4kyzA</b>	Yes	Yes	Yes	Yes
2003205_0003	actiasluna,PLAYER_26,PLAYER_5,Blipperman	37	1.4	XP_020607139.1	0.660	2b8wA	No	-	-	-
2003205_0006	Vinara	-	-	-	0.701	4p0ea	No	-	-	-

2003205_0008	kabubi	-	-		0.674	3cniA	Yes	No	-	-
2003205_1035	actiasluna	34.7	9.2	RHR75271.1	0.715	3c6kD	Yes	No	-	-
2003205_S506	fiendish_ghoul	45.8	2.00E-04	2MQ8_A	0.766	2ddzA	Yes	Yes	No	-
2003205_S722	Susume	36.6	2.4	OVF09666.1	0.722	4ky3A	Yes	No	-	-
2003205_s002	markm457	37	1.7	WP_095044620.1	0.820	2kl8A	Yes	Yes	Yes	Yes
2003245_0004	fiendish_ghoul	36.6	3.3	OQR79684.1	0.827	4wjB	Yes	Yes	Yes	Yes
2003245_S383	fiendish_ghoul	39.3	0.36	WP_113146696.1	0.860	4wjB	Yes	Yes	Yes	Yes
2003245_S385	fiendish_ghoul	35	9.9	WP_008166005.1	0.844	4pxbA	Yes	Yes	Yes	Yes
2003265_0007	fiendish_ghoul	-	-		0.916	4neyB	Yes	Yes	Yes	Yes
2003265_1034	fiendish_ghoul	35.8	4.8	OBT67522.1	0.818	4pxdA	Yes	Yes	No	-
2003265_S115	fiendish_ghoul	-	-		0.808	6cdA	Yes	Yes	Yes	Yes
2003265_S714	Susume	-	-		0.695	5zrbB	Yes	Yes	No	-
2003265_s003	MurloW	-	-		0.635	5i4mA	No	-	-	-
2003265_s005	markm457	37.4	1.2	XP_018401807.1	0.802	4pxdA	Yes	Yes	Yes	Yes
2003265_s008	Susume	41.2	0.057	WP_005505876.1	0.748	6gyA	Yes	Yes	Yes	Yes
a.k.a. foldit3										
2003285_0000	Galaxie,tokens	35.8	2.3	XP_013764557.1	0.794	4clIA	Yes	Yes	Yes	Yes
2003285_s000	tokens	43.9	0.003	ATA67140.1	0.791	4clIA	‡	‡	‡	‡
2003285_s005	Galaxie	35.4	3.1	WP_105677077.1	0.668	2pozA	Yes	Yes	Yes	Yes
2003285_s008	PLAYER_12	-	-		0.718	3jyyA	Yes	Yes	No	-
2003308_0004	spvincent,gitwut,Bletchley Park	36.6	2.9	OAA53803.1	0.678	5n9jW	Yes	Yes	Yes	Yes
2003308_0005	PLAYER_17	38.9	0.46	OJT21624.1	0.712	3ejoA	Yes	No	-	-
2003308_0009	kabubi	36.2	3.6	XP_012773543.1	0.674	6eqtA	Yes	No	-	-
2003308_1010	actiasluna,PLAYER_8	-	-		0.704	3ejoA	Yes	Yes	Yes	Yes
2003333_0000	markm457	37.7	0.71	XP_018574486.1	0.758	4clIA	Yes	No	-	-
2003333_0005	PLAYER_19	35	7.4	WP_006978974.1	0.813	4nezA	Yes	Yes	Yes	Yes
2003333_0006	fiendish_ghoul	-	-		0.819	4neyB	Yes	Yes	Yes	Yes
a.k.a. Peak6										
2003333_1006	MurloW	34.7	2.3	WP_075637384.1	0.651	5hzyA	Yes	Yes	Yes	Yes
2003333_1013	retiredmichael	35	6.4	WP_011364167.1	0.663	1ukxA	Yes	Yes	No	-
2003333_s001	PLAYER_15	-	-		0.739	5lywA	Yes	No	-	-
2003333_s003	Susume	35.8	3.9	WP_027956627.1	0.718	4znIB	Yes	Yes	No	-
2003360_0005	fiendish_ghoul	35	4.2	WP_044185747.1	0.885	5tp4B	Yes	Yes	Yes	Yes
2003360_1013	Bruno Kestemont,PLAYER_20	37	2.3	XP_003196708.1	0.644	2fhY	Yes	Yes	Yes	Yes
2003360_s000	markm457	36.2	5	ATY58566.1	0.725	3n5fA	Yes	Yes	Yes	Yes
2003360_s002	PLAYER_7	36.2	3.1	WP_057511752.1	0.744	4wjB	Yes	Yes	Yes	Yes
2003360_s003	MurloW	-	-		0.635	2ixnA	Yes	Yes	Yes	Yes
2003360_s004	LociOilng	38.9	0.58	XP_019623596.1	0.762	5eq7A	Yes	Yes	Yes	Yes
2003382_0002	Bruno Kestemont,ZeroLeak7	34.7	8.9	WP_051171461.1	0.708	3lwtx	Yes	Yes	No	-
2003382_0004	Vinara	-	-		0.672	5cqCA	Yes	Yes	No	-
2003382_0009	fiendish_ghoul	35.8	3.1	ABW09484.1	0.884	2n3zA	Yes	Yes	No	-
2003382_1021	gitwut	33.1	9.7	WP_037858059.1	0.719	4nogB	Yes	No	-	-
2003382_s005	markm457	35.8	3.1	XP_018821081.1	0.766	2ln3A	Yes	Yes	No	-
2003382_s008	tokens	39.3	0.16	WP_086637674.1	0.720	2nzcA	Yes	Yes	Yes	Yes
2003414_1018	PLAYER_7	37.7	1.1	WP_056892421.1	0.671	1k8fK	Yes	Yes	No	-
2003455_0001	Galaxie,tokens	40	0.085	XP_010698604.1	0.752	4hhuA	Yes	Yes	Yes	Yes
2003455_0002	ZeroLeak7	-	-		0.625	1pp0A	No	-	-	-
2003455_0009	Vinara	-	-		0.698	5ov5A	Yes	Yes	Yes	No
2003455_1023	PLAYER_2	-	-		0.795	1yz7A	No	-	-	-
2003455_S886	Bruno Kestemont	34.3	9.1	WP_039196110.1	0.667	3c66B	No	-	-	-
2003455_S943	Susume	35	5.4	WP_067214428.1	0.667	4nezA	Yes	Yes	Yes	No
2003455_s008	Susume	36.2	1.9	XP_020230923.1	0.709	5o85C	‡	‡	‡	‡
2003485_0000	Galaxie,markm457	-	-		0.864	4pxdA	Yes	Yes	Yes	Yes
2003485_0002	Bruno Kestemont	38.9	0.35	WP_006459835.1	0.667	4zivB	Yes	Yes	Yes	Yes
2003485_1017	ZeroLeak7	-	-		0.681	1pp0A	Yes	Yes	No	-
2003485_1029	fiendish_ghoul	35.4	7	WP_091182063.1	0.832	6cd0A	Yes	Yes	No	-
2003485_1036	Vinara	-	-		0.704	4akrA	No	-	-	-
2003485_S412	Susume	36.6	2.5	WP_067248652.1	0.712	2vcgA	Yes	No	-	-
2003532_0000	actiasluna	35	1.1	YP_007675131.1	0.651	2kt9A	No	-	-	-
2003532_1020	actiasluna	35.8	4.4	WP_012982902.1	0.785	5ae2D	Yes	Yes	No	-
2003532_1022	kabubi	35	8.3	CCA74651.1	0.680	5nj5A	Yes	No	-	-
2003594_0000	Galaxie,tokens	36.6	4.4	SDB06526.1	0.682	4wjB	Yes	Yes	Yes	Yes
2003594_S028	tokens	37.4	2.1	CCX30340.1	0.672	4irxA	Yes	Yes	Yes	Yes
2003594_S603	tokens	38.1	1.4	PYN54536.1	0.638	4zhqD	Yes	Yes	Yes	No
2003594_s008	Susume	-	-		0.635	6cfwk	Yes	Yes	No	-

Successful designs are shown in bold. Foldit player usernames are shown only for players who consented to be named in print; non-consenting players are listed anonymously as PLAYER\_1, PLAYER\_2, etc.

BLAST<sup>49</sup> search was conducted against the nr database of non-redundant protein sequences; scores are omitted where BLAST was unable to find a significant sequence alignment.

TM-align<sup>26</sup> search was conducted against all non-redundant protein chains in the PDB.

† Design 2002290\_S122 has high sequence homology with, and is structurally similar to, a family of bacterial transferases with an unusual β-solenoid fold. Unfortunately, the design failed to express.

‡ Because the DNA was never delivered, designs 2003285\_s000 and 2003455\_s008 were not tested experimentally.

562 **Extended Data**

563

564 Extended Data Table 1. Foldit player-designed proteins selected for experimental testing  
565 (See Extended Data Table 1 PDF)

566 Extended Data Table 2. Crystallographic data collection and refinement statistics

	<b>foldit1</b>	<b>Peak6</b>
<b>Wavelength</b>	1	1
<b>Resolution range</b>	28.92 - 1.18 (1.222 - 1.18)	26.21 - 1.541 (1.596 - 1.541)
<b>Space group</b>	P 1 21 1	P 31 2 1
<b>Unit cell</b>	24.045 43.584 29.276 90 98.998 90	52.414 52.414 56.086 90 90 120
<b>Total reflections</b>	60389 (6169)	129411 (4118)
<b>Unique reflections</b>	18574 (1830)	12866 (860)
<b>Multiplicity</b>	3.3 (3.4)	10.1 (4.8)
<b>Completeness (%)</b>	92.67 (88.38)	94.86 (65.00)
<b>Mean I/sigma(I)</b>	25.65 (9.97)	18.52 (1.34)
<b>Wilson B-factor</b>	10.36	17.88
<b>R-merge</b>	0.02508 (0.1209)	0.0872 (0.7896)
<b>R-meas</b>	0.03015 (0.1439)	0.09186 (0.878)
<b>R-pim</b>	0.01654 (0.07738)	0.02847 (0.3694)

<b>CC1/2</b>	0.999 (0.991)	0.999 (0.714)
<b>CC*</b>	1 (0.998)	1 (0.913)
<b>Reflections used in refinement</b>	18574 (1749)	12861 (860)
<b>Reflections used for R-free</b>	1829 (174)	1282 (85)
<b>R-work</b>	0.1464 (0.1278)	0.1682 (0.2761)
<b>R-free</b>	0.1819 (0.1755)	0.1975 (0.3091)
<b>CC(work)</b>	0.963 (0.982)	0.967 (0.830)
<b>CC(free)</b>	0.956 (0.956)	0.953 (0.806)
<b>Number of non-hydrogen atoms</b>	690	755
<b>macromolecules</b>	574	646
<b>ligands</b>		20
<b>solvent</b>	116	89
<b>Protein residues</b>	68	77
<b>RMS(bonds)</b>	0.008	0.007
<b>RMS(angles)</b>	0.83	1.03
<b>Ramachandran favored (%)</b>	100.00	100.00
<b>Ramachandran allowed (%)</b>	0.00	0.00
<b>Ramachandran outliers (%)</b>	0.00	0.00

<b>Rotamer outliers (%)</b>	0.00	0.00
<b>Clashscore</b>	2.60	3.75
<b>Average B-factor</b>	16.37	24.96
<b>macromolecules</b>	14.54	22.82
<b>ligands</b>		47.36
<b>solvent</b>	25.39	35.49
<b>Number of TLS groups</b>		3

567 Statistics for the highest-resolution shell are shown in parentheses.

568

569

570 Extended Data Table 3: NMR data and refinement statistics for foldit3<sup>a</sup>

<b>Distance restraints</b>	
Total NOE-based restraints	2012
Intra-residue	553
Inter-residue	
Sequential ( $ i-j  = 1$ )	505
Medium-range ( $ i-j  \leq 4$ )	301
Long-range ( $ i-j  > 5$ )	653
Hydrogen bonds restraints	66
<b>Dihedral angle restraints</b>	118
phi	59
psi	59
<b>Restricting restraints / restrained residue</b>	23.0
<b>Restricting long range restraints / restrained residue</b>	6.2
<b>Structure quality statistics</b>	
<b>Restraint Violations</b>	
RMS of distance violation / restraint <sup>b</sup> (Å)	0.01
RMS of dihedral angle violation / restraint (°)	0.88

Max distance restraint violation (Å)	0.66
Max dihedral angle violation (degrees)	7.80
<b>Average r.m.s.d. to medoid conformer<sup>c</sup> (Å)</b>	
Backbone (N, Ca, C')	0.71 ± 0.11
Heavy atoms (all N, C, S, and O)	1.52 ± 0.11
<b>RPF Scores</b>	
Recall	0.912
Precision	0.936
F-measure	0.924
NMR DP-score	0.786
<b>Structure quality factors (raw score / Z-scores<sup>d</sup>)</b>	
Procheck G-factor (phi / psi only)	-0.09 / -0.04
Procheck G-factor (all dihedral angles)	-0.14 / -0.83
Verify3D	0.45 / -0.16
ProsaII (-ve)	0.91 / 1.08
MolProbity clashscore	17.51 / -1.48
<b>Ramachandran plot summary (Richardson statistics)</b>	
Most favored regions (%)	97.3
Allowed regions (%)	2.5

Disallowied regions (%)	0.1
-------------------------	-----

571  
572   <sup>a</sup>Analyzed for the ensemble of 20 lowest-energy structures, residues 1-97, using PDBStat<sup>50</sup> and  
573 PSVS<sup>51</sup> ver 1.5 software.

574  
575   <sup>b</sup>Calculated by using sum over  $r^{-6}$  averaging method.

576  
577   <sup>c</sup>Calculated among 20 structures for "well defined" residues, defined as those that have sum of  
578 phi and psi order parameters  $S(\phi)+S(\psi) > 1.8$ . The "well defined" residues are: 21-45, 48-54,  
579 58-76, 81-87, and 90-96.

580  
581   With respect to mean and standard deviation for a set of 252 X-ray structures with sequence  
582 lengths < 500, resolution  $\leq 1.80 \text{ \AA}$ , R-factor  $\leq 0.25$ , and R-free  $\leq 0.28$ ; a positive value indicates  
583 a 'better' score.

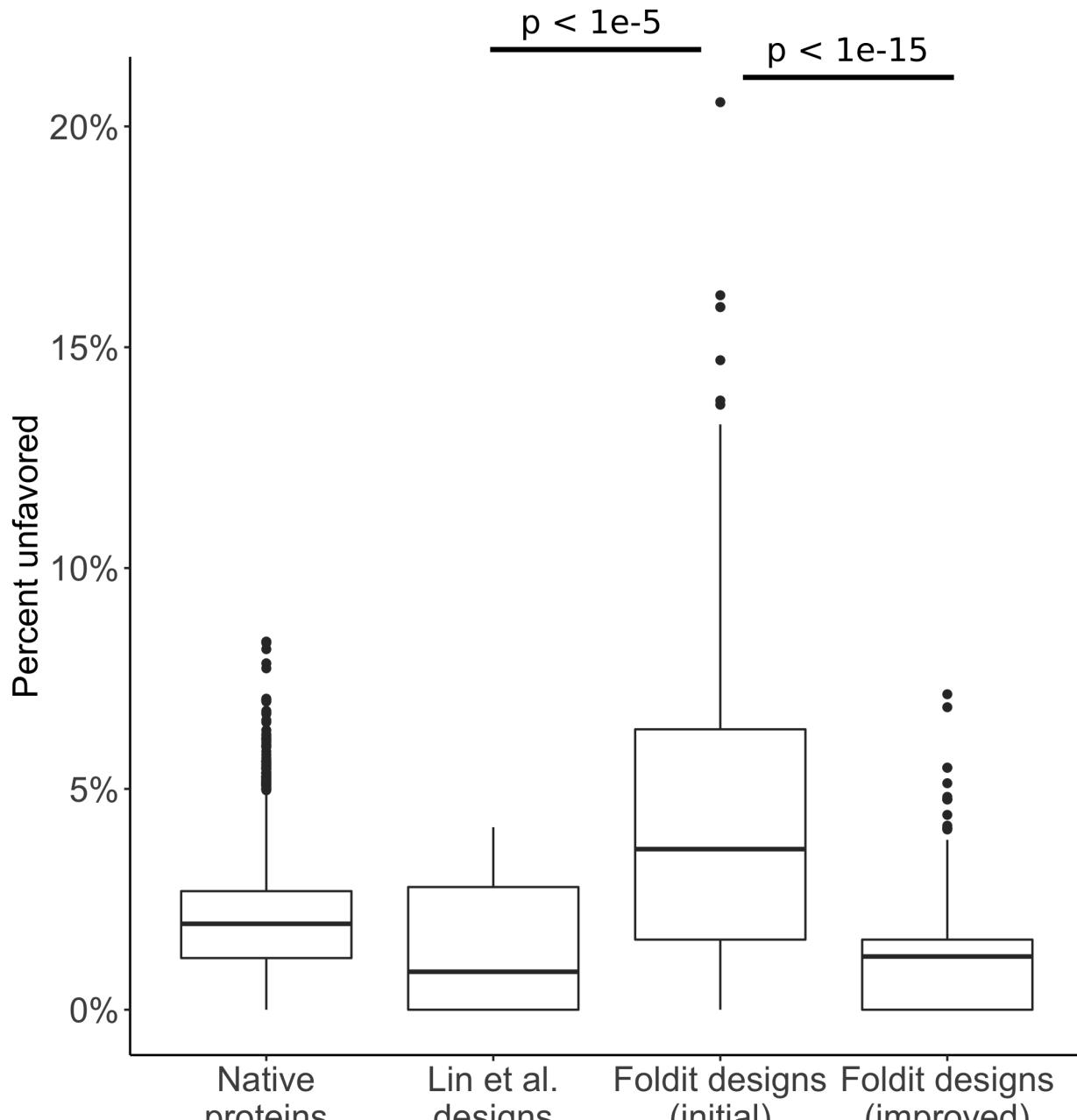
584

585 Extended Data Figure 1. Biophysical characterization of successful designs  
586 (See Supplementary Figure 1 PNG)  
587



588  
589  
590  
591  
592  
593  
594

Extended Data Figure 2. Protein folds represented by successful Foldit player designs. Each fold has a unique arrangement and connectivity of secondary structure elements, depicted in cartoon diagrams. Diagrams are labeled with Roman numerals as in Figure 3. Fold XX is a new fold, previously unobserved in natural proteins; TM-align<sup>26</sup> and DALI<sup>52</sup> alignments against the entire PDB found no structural homologs for design 2003594\_S028 with this fold.



595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605

Extended Data Figure 3. Improvement of backbone quality in Foldit designs. MolProbity<sup>24</sup> was used to calculate the proportion of residues with “unfavored” or “outlier” backbone torsions in: high-resolution crystal structures of native proteins ( $n = 6342$ ), successful *de novo* design models by Lin et al.<sup>7</sup> ( $n = 18$ ), and top-ranking Foldit player-designs from before ( $n = 717$ ) and after ( $n = 250$ ) improvements to Foldit backbone modeling tools. Initial Foldit player designs contained significantly more unfavored torsions than successful designs by Lin et al. ( $p < 1e-5$ , two-tailed t-test). Improvements to Foldit’s backbone modeling tools led Foldit players to produce designs with fewer unfavored torsions ( $p < 1e-15$ , two-tailed t-test). Boxplots show: center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, outliers.

606 Extended Data Foldit player testimonials:  
607  
608 Susan Kleinfelter, a.k.a. Susume (designer of foldit3, and collaborating designer of foldit1):  
609 *I almost always start with a pencil sketch of the protein I am going to build; if I don't make a*  
610 *sketch I look at either an old design of mine that I am modifying, or a protein in the pdb or in*  
611 *some other foldit puzzle that I am modifying. I come up with the sketches either by generating*  
612 *shapes using the rules from the Koga & Koga 2012 paper<sup>[18]</sup>, or by modifying proteins I find in*  
613 *the pdb. The Koga rules all rely on very short loops between the sheets and helices, so*  
614 *adapting a natural fold includes shortening the loops to the few that are built in to foldit, and*  
615 *lengthening or shortening strands or helices to make those loops possible. The foldit1*  
616 *backbone came from a series of sketches I did of modifications to a ferredoxin fold - drop a*  
617 *strand (foldit1), add a strand, add a helix, etc., and/or reverse the termini of any of those*  
618 *designs, all following the Koga rules. The foldit3 backbone was from an exercise I did to find all*  
619 *the 4-strand 2-helix shapes that follow the Koga rules other than the two featured in the Koga*  
620 *paper (ferredoxin-like and IF3-like).*

621  
622 *I hand fold first (1-4 hours), then run scripts (2-4 days). The tools I use most to fold up the*  
623 *backbone are blueprint and dragging dots in the rama map (before blueprint came out I just*  
624 *used rama map). At the start I make all the strands ILE, the helices LEU, and the loop AAs*  
625 *according to Fig. S3 of the Lin 2015 paper<sup>[7]</sup>. I usually get the whole backbone pretty close to*  
626 *the planned shape without wiggling, then band the sheets and shake and wiggle a couple of*  
627 *times. Then I do a few rounds of mutate tool, manually fix AAs, shake and wiggle at various CI.*  
628 *By "fix" I mean manually fix any AAs that mutate has chosen that I don't like, such as buried*  
629 *serine or other hydrophilics, not enough hydrophobics on edge strands or helices, or non-*  
630 *hydrophobics in ideal loop positions where Lin 2015 says they should be hydrophobic.*

631  
632 *At this point I switch to running scripts: a remix and mutate script on all the loops plus a few*  
633 *extra residues on each side of them; one or two banders that alternate random collections of*  
634 *bands with mutate; idealize the backbone and switch to medium wiggle; another bander; and a*  
635 *bands plus local wiggle script at the end. I use only a few scripts and run each one for a long*  
636 *time (several hours or overnight) to let lots of nearby positions and mutations get tried. In*  
637 *between scripts I look for AAs I don't like and manually change them back to ones I like, which*  
638 *always lowers my score but I think makes successful folding more likely. Sometimes I will run*  
639 *all the scripts on one track letting foldit choose all sidechains, in order to get higher on the*  
640 *leaderboard, and run one or more other tracks where I choose some sidechains to share with*  
641 *scientists. When puzzle 1297 (the one in which foldit1 was designed) came out I was still using*  
642 *public versions of scripts, but later I modified some of my favorites to let me mark certain*  
643 *sidechains that I want to stay hydrophobic, since that is the most common manual "fix" I was*  
644 *making.*

645  
646 *Once I start running scripts I also take the primary sequence after each script and run it through*  
647 *either jpred<sup>[53]</sup> or psipred<sup>[54]</sup> to see if the secondary structure I want is considered likely to form. I*  
648 *often manually change AAs just with the purpose of improving the match between the psipred or*  
649 *jpred prediction and my design. I consider this to be both negative design (making it less likely*

650 *the protein will fold up some other way than designed) and gaming the system (because*  
651 *Rosetta uses psipred to inform its prediction, and a design only goes to the wet lab if Rosetta*  
652 *successfully matches the player's design). I try to come up with a compromise between foldit's*  
653 *score function and a good psipred/jpred prediction, but I will sometimes submit a design to*  
654 *scientists that has a good secondary structure prediction even if its foldit score is much lower*  
655 *than what I had before.*

656

657 *High score is not my biggest motivator, although in the long script runs the score function is*  
658 *choosing which modifications to keep. I place higher priority on following the rules from the*  
659 *Koga and Lin papers<sup>[7,18]</sup> (including the ones built into foldit in the form of the ideal loop filter),*  
660 *getting a good psipred prediction, and doing folds that I think are original and cool. I like the*  
661 *challenge of being creative within the somewhat constrained solution space defined by the ideal*  
662 *design rules and foldit's filters, and I think following those rules gives me a leg up in the contest I*  
663 *really care about - getting proteins to fold successfully in the wet lab and ultimately to be*  
664 *published.*

665

666 (See also Susume's instructional screencast at <https://youtu.be/-nizMbICCM0>)

667

668

669 Linda Wei, a.k.a. Galaxie (collaborating designer of foldit1):  
670 *To begin a design puzzle I look at the puzzle page, noting the filter requirements and puzzle*  
671 *comments. Usually I try to maximize the number of helix segments allowed and use the*  
672 *remaining segments for sheets and loops. Depending on time and computer usage required by*  
673 *other puzzles, several designs are produced, one with the minimum number of segments and*  
674 *one with the maximum amount. Tools used most often in this type of puzzle include blueprint,*  
675 *rama map, idealize ss, freeze, band and the move tool. Once the basic design is set, group and*  
676 *public shared recipes are used to refine the shape of the protein in low to auto wiggle,*  
677 *progressing from low to full ci. Mid to end game involves scripts that idealize and fine tune the*  
678 *protein in medium wiggle. "hand mutating" is used to correct low scoring residues. Generally I*  
679 *rely on the scoring function as an indication of strategy effectiveness. Often I will go back to an*  
680 *earlier save to try a different strategy in an attempt to achieve a higher score.*

681 Player testimonials (cont.)  
682  
683 fiendish\_ghoul (designer of Peak6):  
684 *Hand-folding of one protein may take 2-4 hours, and the most of this time is occupied by*  
685 *running mutate tool. Before i start folding i usually draw a layout to help determine SS lengths*  
686 *and types of loops. Hand-folding starts with designation of secondary structures, formation of*  
687 *loops with blueprint tool where it is possible, and pulling the resulting structure close to its*  
688 *intended shape. Most amino acids are initially set to valine, amino acids in loops are assigned*  
689 *according to their surroundings, loop shape and amino acid preferences shown in*  
690 *supplementary material from "Control over overall shape and size in de novo designed proteins"*  
691 *paper<sup>[7]</sup>. Then the protein is wiggled, shaken, mutated and shaken again, and this procedure is*  
692 *repeated with manual interventions between rounds to fix inappropriate mutation choices and*  
693 *other issues. Clashing importance is gradually increased between rounds from 0.01-0.1 to 1.*  
694 *After hand-folding stage, running recipes takes around 10 hours. TvdL DRemixW is run the*  
695 *majority of time and some idealizing recipes like Microidealize and Random Idealize are used at*  
696 *the end. I use only publicly available recipes since i don't make my own.*

697

698 One approach to designing proteins i use is to combine structures (mostly  $\beta\alpha\beta\beta$  and  $\beta\alpha\beta$  units)  
699 described in "Control over overall shape and size in de novo designed proteins" paper<sup>[7]</sup> in  
700 various ways. Some modifications can be applied to these standard structures, like replacing  
701 sheet-sheet loop with loop-helix-loop-helix-loop sequence.

702

703 Another approach is to modify existing folds. For example, proteins with N and C termini  
704 positioned near each other (like in ferredoxin-like fold from previously mentioned paper) can be  
705 imagined as continuous structure, and then a "cut" can be made in place of one of the loops of  
706 the protein to get a protein with similar shape but different SS sequence. Recently validated  
707 protein was designed this way.

708

709 These approaches have been yielding my best scoring solutions so far and are relatively easy  
710 to implement. More experimental folds tend to have much lower score. In general, simpler folds  
711 with 2 helices and 3-4 sheets have the highest score, and the more complex the protein gets,  
712 the more difficult it is to make it score well.

713

714 Regarding foldit tools, i find blueprint tool useful because it speeds up folding process and  
715 makes it possible to quickly sample different structures. Rama Map is useful for creating  
716 structures like beta bulges and loops that are not in blueprint library.

717

- 718      **Extended References**
- 719      49. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment  
720      search tool. *Journal of Molecular Biology* **215**, 403–410 (1990).
- 721      50. Tejero, R., Snyder, D., Mao, B., Aramini, J.M., Montelione, G.T. PDBStat: A universal  
722      restraint converter and restraint analysis software package for protein NMR. *J. Biomol.*  
723      *NMR* **56**, 337-351 (2013).
- 724      51. Bhattacharya, A., Tejero, R., Montelione, G.T. Evaluating protein structures determined  
725      by structural genomics consortia. *PROTEINS: Struct. Funct. Bioinformatics* **66**, 778 - 795  
726      (2007).
- 727      52. Holm, L. & Laakso, L. M. Dali server update. *Nucleic Acids Res.* **44**, W351–5 (2016).
- 728      53. Drozdetskiy, A., Cole, C., Procter, J. & Barton, G. J. JPred4: a protein secondary  
729      structure prediction server. *Nucleic Acids Res.* **43**, W389–94 (2015).
- 730      54. McGuffin, L. J., Bryson, K. & Jones, D. T. The PSIPRED protein structure prediction  
731      server. *Bioinformatics* **16**, 404–405 (2000).
- 732