

1 **Full title: *De novo* protein fold families expand the designable ligand binding site space**

2 **Short title: *De novo* protein families expand the ligand binding site space**

3 Xingjie Pan^{*1,2}, Tanja Kortemme^{*1,2,3,4}

4

5 ¹Department of Bioengineering and Therapeutic Sciences, University of California San
6 Francisco, San Francisco, CA, USA.

7 ²UC Berkeley – UCSF Graduate Program in Bioengineering, University of California San
8 Francisco, San Francisco, CA, USA.

9 ³Quantitative Biosciences Institute (QBI), University of California San Francisco, San Francisco,
10 CA, USA

11 ⁴Chan Zuckerberg Biohub, San Francisco, CA, USA.

12 *Correspondence to: xingjiepan@gmail.com; tanjakortemme@gmail.com.

13

14 **Abstract**

15 A major challenge in designing proteins *de novo* to bind user-defined ligands with high
16 specificity and affinity is finding backbone structures that can accommodate a desired binding
17 site geometry with high precision. Recent advances in methods to generate protein fold families
18 *de novo* have expanded the space of accessible protein structures, but it is not clear to what
19 extend *de novo* proteins with diverse geometries also expand the space of designable ligand
20 binding functions. We constructed a library of 25,806 high-quality ligand binding sites and
21 developed a fast protocol to place (“match”) these binding sites into both naturally occurring and
22 *de novo* protein families with two fold topologies: Rossmann and NTF2. 5,896 and 7,475 binding
23 sites could be matched to the Rossmann and NTF2 fold families, respectively. *De novo*
24 designed Rossmann and NTF2 protein families can support 1,791 and 678 binding sites that
25 cannot be matched to naturally existing structures with the same topologies, respectively. While
26 the number of protein residues in ligand binding sites is the major determinant of matching
27 success, ligand size and primary sequence separation of binding site residues also play
28 important roles. The number of matched binding sites are power law functions of the number of
29 members in a fold family. Our results suggest that *de novo* sampling of geometric variations on
30 diverse fold topologies can significantly expand the space of designable ligand binding sites for
31 a wealth of possible new protein functions.

32

33 **Author summary**

34 *De novo* design of proteins that can bind to novel and highly diverse user-defined small
35 molecule ligands could have broad biomedical and synthetic biology applications. Because
36 ligand binding site geometries need to be accommodated by protein backbone scaffolds at high
37 accuracy, the diversity of scaffolds is a major limitation for designing new ligand binding
38 functions. Advances in computational protein structure design methods have significantly

39 increased the number of accessible stable scaffold structures. Understanding how many new
40 ligand binding sites can be accommodated by the *de novo* scaffolds is important for designing
41 novel ligand binding proteins. To answer this question, we constructed a large library of ligand
42 binding sites from the Protein Data Bank (PDB). We tested the number of ligand binding sites
43 that can be accommodated by *de novo* scaffolds and naturally existing scaffolds with same fold
44 topologies. The results showed that *de novo* scaffolds significantly expanded the ligand binding
45 space of their respective fold topologies. We also identified factors that affect difficulties of
46 binding site accommodation, as well as the relationship between the number of scaffolds and
47 the accessible ligand binding site space. We believe our findings will benefit future method
48 development and applications of ligand binding protein design.
49
50

51 **Introduction**

52 Ligand binding is a major class of protein functions, and the ability to design ligand binding *de*
53 *novo* has many important applications(1) such as engineering of biosensors and ligand-
54 controlled protein functions(2, 3). Naturally occurring proteins recognize their cognate ligands
55 with high affinity and specificity using defined three-dimensional geometries of binding sites with
56 high shape complementarity between ligands and proteins. For the formation of favorable
57 hydrophobic and polar interactions, the chemical groups on the protein must be placed at
58 specific spatial positions relative to the ligand(4, 5). Designing new ligand binding proteins
59 therefore requires the ability to build binding sites with defined geometries into stable protein
60 scaffolds that can accommodate the desired interaction geometry with high accuracy. While this
61 approach has led to the successful design of enzymatic activity(6, 7), ligand binding proteins(8,
62 9), and biosensors(2, 3, 10), it has been limited by both the availability of defined binding site
63 geometries and stable protein scaffolds into which these binding sites can be placed(3).

64
65 Several methods have recently been developed to address the first problem, increasing the
66 number of potential ligand binding sites one could generate. The RIF docking method(11)
67 generates ensembles of billions of side chain placements that make defined hydrogen-bonding
68 and non-polar interactions with a target ligand. Other methods(12, 13) use statistics from the
69 protein data bank (PDB) to find three-dimensional placements of amino acid residues that form
70 favorable interactions with fragments of a ligand, which can then be assembled into complete
71 binding site geometries. Protein-ligand interactions defined by these methods have been built
72 successfully into a *de novo* designed beta barrel(11), and a parametrically designed helical
73 bundle(13).

74

75 Naturally occurring proteins solve the second problem, finding a suitable protein backbone to
76 accommodate a specific binding site geometry, not by using a different fold for each function but
77 instead by evolving structural variation in existing protein fold families. This variation allows
78 proteins with the same fold topology (identity and connectivity of secondary structure elements)
79 to tune the precise geometry of binding sites to recognize diverse ligands(14). This strategy has
80 recently been mimicked by advances in computational protein design methods. These methods
81 have generated *de novo* designed protein fold families with large numbers of diverse
82 geometries(15, 16), which have significantly expanded the accessible designable protein
83 structure space. The resulting *de novo* proteins might be able to support binding sites that
84 cannot be built onto naturally occurring proteins in the PDB, but the extent to which *de novo* fold
85 families could improve binding site design has not been explored. Understanding the
86 relationship between the space occupied by protein structures, and the space available to
87 support different functions, is important for developing methods to design proteins *de novo* that
88 can bind to novel and highly diverse user-defined ligands.

89
90 Here, we studied the ability of native and *de novo* fold families to support a large number of
91 different ligand binding sites. We built a high-quality ligand binding site library from high
92 resolution protein crystal structures. We then matched the binding site library to members of
93 protein folding families using two protocols: a newly developed “fast matching” protocol and the
94 standard method for matching in the Rosetta program for structure modeling and design(5). We
95 calculated the number of matched binding sites for four fold families with two different
96 topologies. We studied the effects of binding site sizes, ligand sizes and primary sequence
97 separation of binding site residues on the matching success rates and determined the increase
98 of numbers of matches with increasing the sizes of fold families. Together, we show that *de*
99 *nov*o fold family design is a promising approach to broaden the scope of designable ligand
100 binding sites.

101 Results

102 We first constructed a library of ligand binding sites from native proteins in the PDB. We
103 extracted 25,806 ligands that have at most 100 heavy atoms as well as the ligand binding site
104 residues from 23,238 cluster representative structures from the PDB95 database(17) where
105 chains from the protein data bank are clustered at 95% identity (**Methods**). The extracted
106 ligands have between 1 and 93 heavy atoms (**Fig 1A,B**). 80.6% percent of the ligands have 13
107 or fewer heavy atoms, and 7,335 (28.4%) of the ligands have only 1 heavy atom. There are
108 2,461 unique ligand types in the 25,806 binding sites. The distribution of ligand type frequencies
109 has a long tail (**S1 Table**). There are 33 frequent ligand types that appear in over a hundred
110 binding sites, while 1,817 ligand types only appear in single binding sites. The frequent ligand
111 types include common crystallographic additives such as glycerol; 1,2-ethanediol; ions such as
112 SO_4^{2-} and Mg^{2+} ; and cofactors such as heme and flavin adenine dinucleotide (FAD). Ligands
113 that appear in multiple binding sites are seen as vertical stripes in **Fig 1B**. Binding sites have
114 between 2 and 41 residues, with 81.2% of the binding sites having 7 or fewer binding site
115 residues. The number of protein residues in binding sites scales linearly with the number of
116 ligand heavy atoms, with a slope 0.35 (**Fig 1B**). The frequencies of amino acid types in binding
117 sites are different from those for whole proteins reported by UniProtKB/Swiss-Prot (**Fig 1C**). We
118 defined the enrichment ratios of amino acids as their frequencies in ligand binding sites divided
119 by their frequencies in whole proteins. The large aromatic side chains Trp, Tyr and Phe are the
120 top 1, top 3 and top 6 enriched amino acid residues, respectively. His, characterized by its
121 ability to coordinate metal ions and to catalyze chemical reactions, is the second most enriched
122 amino acid residue. Asp and Arg are the 4th and 5th enriched amino acid residues, which may
123 play important roles in interacting with charged ligands. Binding sites with single heavy atom
124 ligands have different amino acid preference than those binding to ligands with at least two
125 heavy atoms (**Fig 1D**). For the binding sites with single heavy atom ligands, the negatively

126 charged residues Asp and Glu, which can form favorable electrostatic interactions with
127 positively charged metal ions, are highly enriched. The enrichment ratios of Asp and Glu are 4.6
128 and 2.2, respectively. The top 5 enriched residues that bind to ligands with at least 2 heavy
129 atoms are Trp, His, Tyr, Phe and Arg.

130
131 The binding site library is useful for testing the ability of protein fold families to support ligand
132 binding sites. A protein scaffold can in principle support a ligand binding site if the binding site
133 residues can be built onto the scaffold such that the key interactions between the ligand and
134 binding site protein residues are preserved. The Rosetta matcher protocol(5) has been shown to
135 be successful in matching ligand binding sites to protein scaffolds(8). However, the Rosetta
136 matcher is too slow to match tens of thousands of binding sites to thousands of scaffolds
137 because it samples all possible side chain rotamers of binding site residues. To perform all-
138 against-all matching between the library of ligand binding sites and the sets of scaffolds, we
139 developed a new fast match protocol (**Fig 2A**). In the fast match protocol, the binding site is
140 anchored and matched as a rigid body (**Methods**). This rigid body approximation drastically
141 improved the matching speed. We tested the run time by matching the binding site library to the
142 native NTF2 fold family (CATH superfamily 3.10.450.50)(18). The mean time to find a
143 successful standard Rosetta match is 706 s while the mean time of a successful fast match is
144 3.1s. As a trade-off, the rigid body approximation of the fast match method may discard binding
145 sites that can be matched by the Rosetta matcher using alternative side chain rotamers.
146 Therefore, in this study we focused on matching ligand binding sites using the side chain
147 rotamers present in the original ligand binding site in the PDB. Using these original rotamers
148 also let us directly compare the backbone geometries in the native binding sites and the
149 backbone geometries in our scaffold libraries.

150

151 We matched the binding site library to backbone scaffolds of *de novo* designed Rossmann and
152 NTF2 protein fold families generated by the loop-helix-loop unit combinatorial sampling (LUCS)
153 method(15), as well as the two native fold families with the same topology from the CATH
154 database(18) (**Fig 2B, Methods**). To determine if a fold family can support a given ligand
155 binding site, we first used fast match to match the ligand binding site to all protein scaffolds in
156 the family. Then we used the Rosetta matcher to match the binding site to the scaffolds that
157 passed the fast match (**Methods**). To limit computational time, once the Rosetta matcher found
158 a match for a given binding site, we skipped matching to further scaffolds in the same family.
159 Since we used stringent matching criteria (**Methods**), the matched binding sites in the scaffold
160 closely recapitulated the interactions between the ligands and binding site residues in the
161 original protein structures from which the binding sites were derived (**Fig 2C**).

162
163 Between 5896 and 7548 binding sites could be successfully matched by the Rosetta matcher to
164 each fold family when considering all binding sites (**Table 1**). The number of binding site
165 residues was the major determinant of the matching success rate (**Fig 3**). For the *de novo*
166 Rossmann fold family, the success rates for 2, 3 and 4 protein residue binding sites were
167 93.8%, 33.4% and 6.5%, respectively. Only 13 binding sites with 5 or 6 residues could be
168 matched. There was no match for binding sites with more than 6 protein residues. Similar
169 dependencies on binding site sizes were observed across the 4 different protein fold families
170 (**Table 2**). Because almost all 2-residue binding sites could be matched and the matching
171 success rates were low for binding sites with more than 3 residues, we used 3-residue binding
172 sites to further study properties of successful matches. We constructed a new library of binding
173 sites that all have 3 protein residues (**Methods**) and matched the binding sites to the scaffold
174 libraries using the same protocol. The number of successfully matched 3-residue binding sites
175 ranged from 2,142 to 3,715 (**Table 1**).

176

177 For the successfully matched 3-residue binding sites, we first analyzed the positions of matches
178 relative to the surface of the scaffolds. For each scaffold, we used the Rosetta Layer residue
179 selector(19) to assign layers to all of its residues in a side chain independent manner
180 (**Methods**). Residues on the surfaces of scaffolds were assigned to the surface layer; deeply
181 buried residues were assigned to the core layer; and the rest of residues were assigned to the
182 boundary layer (**Fig 4A**). In all of the fold families, surface layer residues were most abundant,
183 which accounted for 47%-63% of all residues. 29%-39% residues were in the boundary layer
184 and 6%-20% residues were in the core layer (**Fig 4B**). NTF2 fold proteins had more surface
185 layer residues which was likely due to the pocket of this fold. We defined the layer of each
186 residue in a matched binding site as the layer of its matched scaffold residue position. The
187 frequencies of matched residue layers are similar to the frequencies of scaffold layers (**Fig 4C**).
188 To evaluate the positions of matches at the binding site level, we defined a depth score for each
189 matched binding site. The depth score of a matched binding site is the number of boundary
190 residues plus two times the number of core residues. The depth scores for binding sites
191 matched to different fold families had similar distributions (**Fig 4D**). 20%-27% binding sites were
192 entirely matched to protein surfaces and had depth scores of 0. The remainder of matched
193 binding sites were buried to some extent. The majority of binding sites were in shallow pockets
194 with depth scores ranging from 2 to 4. Only 8%-12% binding sites were matched to deeply
195 buried positions with depth scores of 5 or 6.

196
197 We then tested factors that affect the “matchability” of 3-residue binding sites. We compared the
198 number of overlapping binding sites that were matched to both of two fold families to the
199 expected number of overlapping binding sites (**Fig 5A**). If matching to one fold family is
200 independent from matching to another fold family, the probability of overlapping binding sites
201 should be the product of the probabilities of matching to each fold family. We compared 4 pairs
202 of scaffold libraries (**Fig 5A**): *de novo* designed Rossmann folds versus *de novo* designed NTF2

203 folds (top left) or versus native Rossmann folds(top right), and native NTF2 folds versus *de novo*
204 designed NTF2 folds (bottom left) or versus native Rossmann folds (bottom right). For all 4 pairs
205 of scaffold libraries, the observed number of overlapping binding sites was significantly higher
206 than the number of expected overlapping binding sites (chi-squared test p-value < 10^{-300}). This
207 result indicates that some binding sites had higher matchabilities (probabilities to be matched to
208 multiple scaffold libraries). We investigated the contribution of ligand sizes to binding site
209 matchabilities. As expected, the matching success rates for 3-residue binding sites decreased
210 with an increase of the number of ligand heavy atoms (**Fig 5B**), likely because larger ligands are
211 more likely to clash with the scaffold backbones. We also hypothesized that binding sites whose
212 residues have larger separations in primary sequences are more difficult to match. To confirm
213 that non-local binding sites are harder to match, we calculated the mean inter-residue primary
214 sequence distances for each 3-residue binding site and plotted the mean distances against the
215 matching success rates (**Fig 5C**). When the 3-residues in a binding site were consecutive in
216 primary sequence, the mean primary sequence distance was 1.33, and the matching success
217 rates were higher than 80%. The success rates dropped rapidly with the increase of mean
218 distance and reached a plateau at low match success rates when the mean distance reached
219 70.

220
221 Next, we studied how the number of matched 3-residue binding sites grew with an increase of
222 the number of scaffolds in fold families (**Methods**). The log of the number of matched binding
223 sites scaled linearly with the log of the number of scaffolds (**Fig 6A-D**). This power law
224 relationship was valid for both the number of fast matches and Rosetta matches across the 4
225 different fold families. The powers of the power law functions (slopes of the log-log plots) ranged
226 from 0.184 to 0.298. Since the powers were small, the increase of matches progressively
227 diminished as the number of scaffolds got large. Because there is a limited number of
228 designable structures for each fold family, the power law relationship cannot continue

229 indefinitely, but it can still provide a reasonable estimation of the upper bound of the number of
230 matches. Extrapolating the *de novo* fold family power law relationships to the number of
231 representative structures from the PDB95 database, i.e., 23,238 structures, the numbers of
232 expected Rosetta matches for the Rossmann fold family and the NTF2 fold family would be
233 7,346 and 6,640. Based on this analysis, the extrapolated numbers of matchable binding sites
234 are still much smaller than the number of total binding sites, highlighting the importance of
235 having diverse fold topologies for different functions.

236

237 Finally, to understand how *de novo* scaffolds expand protein function space, we studied the
238 binding sites that can be matched to *de novo* scaffolds but not to the native scaffolds of the
239 same topology. We plotted the number of binding sites that were matched to only *de novo* fold
240 families versus the number of *de novo* scaffolds (**Fig 6E,F**). For each topology, there are more
241 than 1,000 binding sites that are exclusively matched to *de novo* scaffolds. These relationships
242 also follow power law functions. The slopes are larger than the slopes of the total matches (**Fig**
243 **6A-D**), indicating that binding sites that can match to both native and *de novo* fold families
244 saturate quickly.

245

246 **Discussion**

247 Advances in computational protein structure sampling methods(20) have expanded the
248 accessible structure space of *de novo* designed proteins. In particular, two recently developed
249 computational methods(15, 16) are capable of engineering *de novo* protein families that contain
250 defined variations in geometry of proteins that share the same overall fold topology. We probed
251 the functional implications of *de novo* protein fold families generated by the LUCS method(15)
252 by matching known ligand binding sites to both native and *de novo* fold families. We found that
253 thousands of ligand binding sites that cannot be matched to native fold families can be matched
254 to LUCS-generated members of *de novo* fold families of the same topology, showing that LUCS
255 generated structures expand both the accessible protein structure and the accessible protein
256 function space. The number of matched binding sites increased as a power law function of the
257 number of scaffolds. This relationship allowed us to estimate the upper bound of matches as the
258 number of scaffolds grew and showed that, in addition to geometric variation, different fold
259 topologies are necessary to support diverse functions.

260
261 Previous studies have shown that computationally generated artificial (ART) compact
262 homopolypeptide structures can match virtually every native ligand binding pocket(21, 22). In
263 contrast, the native and *de novo* fold families we studied here can only be matched to a limited
264 fraction of native binding sites. A likely reason is that we only used structures with two
265 topologies while the ART structures are generated using secondary structure preferences from
266 thousands of random PDB structures with many different fold topologies. These two behaviors
267 together support that the diversity of topologies is important for the repertoire of native ligand
268 binding functions. Additionally, the *de novo* designed structures we used were subjected to
269 filters for a set of physical properties such as core packing, hydrogen bonding and surface
270 exposed hydrophobic patches(15). These filters are designed to eliminate structures that are not

271 likely to fold, whereas the ART structures are model polyleucine homopolypeptides.
272 Requirements for folding places diverse additional constraints on the accessible conformational
273 space of protein structures.
274
275 Using the new fast match protocol introduced here as well as the Rosetta matcher, we were
276 able to match a library of high-quality binding sites to *de novo* protein fold families. To engineer
277 new ligand binding proteins, the matching step is typically followed by sequence design(3, 8) to
278 optimize the binding site protein environment. Ligand binding site design is a challenging
279 problem because the designed sequence must simultaneously be compatible with the protein
280 fold and precisely place binding site residues in their desired geometries for favorable
281 interactions with the ligand. Given the typically high stability of *de novo* designed protein
282 families(15, 23), matches generated by the protocol described here could be good model
283 systems for testing binding site design algorithms.
284
285 Another advantage of using *de novo* fold families for ligand binding site design is that the
286 systematic sampling of diverse geometries could provide an ensemble of negative states. Using
287 negative states in design has been shown to improve accuracy in protein stability prediction(24).
288 Thus, a *de novo* ensemble of negative states may increase success rates of ligand design
289 where high accuracy in both sampling and scoring designs is required. Ensembles of different
290 conformational states in *de novo* fold families also pave the way to engineer ligand binding-
291 induced conformational changes. Small molecule-induced switches could be designed by
292 building a ligand binding site in one of the structures in the *de novo* fold family and tuning the
293 free energy gaps between the ligand binding state and the other states. We envision that *de*
294 *nov*o designed protein fold families will play an important role in designing functions such as
295 ligand binding and protein switches.
296

297 **Methods**

298 **Binding site library construction**

299 Ligand binding sites were extracted from the PDB95(17) database. The representative pdb
300 structures for each cluster, which were listed in the pdb_95.cod file, were used for binding site
301 extraction. The representative structures were filtered by resolution. Only crystal structures
302 whose resolutions were better than 2 Å were kept. Ligand residues were identified by built-in
303 functions in PyRosetta(25). In this study, we focused on ligands that had at most 100 heavy
304 atoms. Ligands that had average heavy atom B-factors greater than 60 Å² were filtered out.
305 Ligands that did not have protein residues within 5 Å were also excluded from subsequent
306 processing. We calculated the Rosetta 2-body energy scores(26, 27) between ligands and
307 protein residues that have at least one heavy atom within 5 Å from any ligand heavy atom.
308 Ligand binding site residues were defined as protein residues that had favorable van der Waals,
309 electrostatic or hydrogen bond interactions with the ligand. A residue was included in a binding
310 site if the sum of its Rosetta energy(27) terms fa_atr, fa_elec, hbond_bb_sc and hbond_sc was
311 less than -1 Rosetta energy units (REU). We excluded protein residues from consideration that
312 had total Rosetta scores greater than 50 to avoid poorly modeled residues, such as those who
313 have severe clashes with the protein environment. We also excluded all residues with missing
314 heavy atoms in the PDB file. We only kept ligand binding sites that have at least two protein
315 residues. To prevent overcounting ligands in structures which had multiple chains of the same
316 protein in their asymmetric units, only one binding site was extracted for the same ligand in a
317 given structure.

318

319 **Fast match protocol**

320 We developed a new fast match protocol to rapidly match the library of binding sites to the sets
321 of protein scaffolds. During the fast match, a ligand binding site is treated as a rigid body. When

322 the fast matcher matches a ligand binding site to a scaffold, it first iterates through all pairs of
323 binding site protein residues and scaffold residues. For each pair of residues, the protocol
324 superimposes the N, Ca and C atoms of the binding site residue to the corresponding atoms in
325 the scaffold residue. The remainder of the binding site is transformed as a rigid body. Then the
326 matcher finds the closest scaffold residues to each binding site protein residue. The distances
327 between residues are defined as the Ca-Ca distances. If all distances between binding site
328 protein residues and their closest scaffold residues are within 2 Å, the backbone N, Ca and C
329 atoms of the binding site protein residues are superimposed to the N, Ca and C atoms of their
330 closest scaffold residues. The superimposition minimizes the root mean squared deviation
331 (RMSD) between the corresponding atoms. If the RMSD is within 1 Å, the cosine of angles
332 between the vectors pointing from Ca to Cb of corresponding residues are calculated. If all the
333 cosine values are greater than 0.7, clashes between the matched binding site and the scaffold
334 backbone are checked. Two atoms are defined to clash when the distance between them is less
335 than the sum of their Lennard-Jones radii times a scale factor of 0.6. The match is accepted if
336 the ligand and protein side chains from the binding site do not clash with the scaffold backbone
337 atoms that are not matched to binding site residues.

338

339 **Standard Rosetta matcher**

340 For each binding site successfully matched to a scaffold using fast match, we ran the standard
341 Rosetta matcher(5). We made mol2 files for ligands using Open Babel(28) and generated ligand
342 parameter files with the molfile_to_params.py script distributed with Rosetta. The relative
343 positions of a ligand and a binding site protein residue are defined by 6 heavy atoms. On the
344 ligand side, the heavy atom closest to the protein residue and the two ligand heavy atoms
345 closest to the first ligand heavy atom are defined as the anchor atoms. On the protein residue
346 side, the heavy atom closest to the ligand and two protein atoms closest to the first protein
347 heavy atom are defined as the anchor atoms. For each binding site, we generated a constraint

348 file where the relative positions between the ligand anchor atoms and protein residue anchor
349 atoms were constrained. We used stringent matching criteria similar to those used in previous
350 work(8, 12). The relative distances between ligands and binding site residues are sampled at
351 ideal values; the relative angles and torsions are sampled at the ideal values and $\pm 10^\circ$ from the
352 ideal values. The binding sites were matched using the standard Rosetta matcher with the
353 following command:

```
354  
355 match.linuxgccrelease -match:output_format PDB -match:match_grouper  
356 SameSequenceGrouper -match:consolidate_matches -match:output_matches_per_group 1 -  
357 use_input_sc -in:ignore_unrecognized_res -ex1 -ex2 -enumerate_ligand_rotamers false -  
358 match::lig_name LIG_NAME -match:geometric_constraint_file CST_FILE -s SCAFFOLD_PDB -  
359 match::scaffold_active_site_residues POS_FILE
```

360
361 where LIG_NAME is the 3-letter name of the ligand, CST_FILE is the constraint file,
362 SCAFFOLD_PDB is the pdb file of the scaffold structure and POS_FILE is the file that stores
363 the matchable residues. In this study, all residues on a scaffold are matchable.

364

365 **Construction of scaffold libraries**

366 The *de novo* Rossmann and NTF2 fold families were reported in ref.(15). The scaffolds in these
367 fold families were generated by the LUCS method and filtered by a set of designability
368 filters(15). We randomly selected 1,000 scaffolds from each *de novo* fold family as the scaffold
369 set for ligand binding site matching. The native fold families of Rossmann and NTF2 folds were
370 obtained from the CATH database(18). The native Rossmann fold scaffolds were extracted from
371 the CATH 3.40.50.1980 superfamily and the native NTF2 family structures were from the CATH
372 3.10.450.50 superfamily. Because the automatic classification algorithm of the CATH database
373 did not force all structures in a CATH superfamily to have a same topology, we manually

374 excluded the CATH structures that have different topologies from the *de novo* designed
375 scaffolds. As a result, the native Rossmann fold scaffold set had 20 structures and the native
376 NTF2 fold scaffold set had 103 structures. The C-terminal helices in *de novo* NTF2 scaffolds
377 occluded the ligand binding pocket. In contrast, only 35 out of 103 native NTF2 scaffolds had C-
378 terminal helices. Among these native C-terminal helices, 31 helices pointed away from pocket
379 entrances, and thus, did not affect the accessibility of ligand binding sites, leaving only 4
380 scaffolds with pocket occluding C-terminal helices. We therefore trimmed the C-terminal helices
381 in *de novo* NTF2 proteins to expose the ligand binding pocket.

382

383 **Construction of a library of 3-residue binding sites**

384 The 3-residue binding site library was constructed from the library of all binding sites. We
385 eliminated binding sites with fewer than 3 residues. The binding sites with 3 protein residues
386 were kept unchanged. For binding sites with more than 3 protein residues, we scored the total
387 Rosetta two-body energy(26) between the ligand and each protein residue. We kept the 3
388 protein residues with lowest total two-body energies and removed the remainder of the binding
389 site residues.

390

391 **Assignment of layers to scaffold residues**

392 The Rosetta Layer selector(19) with the default settings was applied to assign layers to each
393 scaffold residue. The layer of a residue was determined by a weighted count of the number of
394 neighbor amino acid residues in a cone extending along its Ca-Cb vector. A residue is assigned
395 to the surface layer if the weighted count is less than 2; a residue is assigned to the core layer if
396 the weighted count is greater than 5.2; all other residues are assigned to the boundary layer.

397

398 **Calculation of the numbers of matches for subsets of fold families**

399 During the process of matching a binding site to a fold family, we recorded the number of
400 scaffolds in the fold family that we tested to find the first successful fast match and called this
401 number the first-fast-match-encounter-number. The number of fast matches for a subset of a
402 fold family with N scaffolds was defined as the number of binding sites with first-fast-match-
403 encounter-numbers smaller than or equal to N. The number of Rosetta matches for subsets of
404 fold families were calculated in the same way.

405

406 **Data availability**

407 All relevant data are available in the manuscript and supporting information data files. Rosetta
408 source code is available from rosettacommons.org. Scripts, the binding site library and the
409 scaffold sets are available at [https://github.com/Kortemme-](https://github.com/Kortemme-Lab/match_ligand_binding_sites/releases/tag/v1)
410 [Lab/match_ligand_binding_sites/releases/tag/v1](https://github.com/Kortemme-Lab/match_ligand_binding_sites/releases/tag/v1) .

411

412 **Acknowledgments**

413 This work was supported by grants from the National Institutes of Health (NIH) (R01-
414 GM110089) and the National Science Foundation (NSF) (DBI-1564692) to TK. We additionally
415 acknowledge a UCSF Discovery Fellowship to XP. TK is a Chan Zuckerberg Biohub
416 Investigator.

417

418 **Author contributions**

419 XP conceived the idea for the project and developed the approach, with contributions from TK.
420 XP developed the computational methods and performed the simulations. TK provided
421 guidance, mentorship and resources. XP and TK wrote the manuscript.

422

423 **Competing interests:** The authors declare no competing interests.

424 **References**

- 425 1. Feldmeier K, Hocker B. Computational protein design of ligand binding and catalysis. *Curr*
426 *Opin Chem Biol.* 2013;17(6):929-33.
- 427 2. Feng J, Jester BW, Tinberg CE, Mandell DJ, Antunes MS, Chari R, et al. A general strategy
428 to construct small molecule biosensors in eukaryotes. *eLife.* 2015;4.
- 429 3. Glasgow AA, Huang YM, Mandell DJ, Thompson M, Ritterson R, Loshbaugh AL, et al.
430 Computational design of a modular protein sense-response system. *Science (New York,*
431 *NY.* 2019;366(6468):1024-8.
- 432 4. Yang W, Lai L. Computational design of ligand-binding proteins. *Current opinion in*
433 *structural biology.* 2017;45:67-73.
- 434 5. Zanghellini A, Jiang L, Wollacott AM, Cheng G, Meiler J, Althoff EA, et al. New algorithms
435 and an in silico benchmark for computational enzyme design. *Protein Sci.*
436 2006;15(12):2785-94.
- 437 6. Jiang L, Althoff EA, Clemente FR, Doyle L, Rothlisberger D, Zanghellini A, et al. De novo
438 computational design of retro-aldol enzymes. *Science (New York, NY.*
439 2008;319(5868):1387-91.
- 440 7. Rothlisberger D, Khersonsky O, Wollacott AM, Jiang L, DeChancie J, Betker J, et al. Kemp
441 elimination catalysts by computational enzyme design. *Nature.* 2008;453(7192):190-5.
- 442 8. Tinberg CE, Khare SD, Dou J, Doyle L, Nelson JW, Schena A, et al. Computational design
443 of ligand-binding proteins with high affinity and selectivity. *Nature.* 2013;501(7466):212-6.
- 444 9. Polizzi NF, Wu Y, Lemmin T, Maxwell AM, Zhang SQ, Rawson J, et al. De novo design of a
445 hyperstable non-natural protein-ligand complex with sub-Å accuracy. *Nat Chem.*
446 2017;9(12):1157-64.
- 447 10. Bick MJ, Greisen PJ, Morey KJ, Antunes MS, La D, Sankaran B, et al. Computational
448 design of environmental sensors for the potent opioid fentanyl. *eLife.* 2017;6.

- 449 11. Dou J, Vorobieva AA, Sheffler W, Doyle LA, Park H, Bick MJ, et al. De novo design of a
450 fluorescence-activating beta-barrel. *Nature*. 2018;561(7724):485-91.
- 451 12. Lucas JE, Kortemme T. New computational protein design methods for de novo small
452 molecule binding sites. *PLoS computational biology*. 2020;16(10):e1008178.
- 453 13. Polizzi NF, DeGrado WF. A defined structural unit enables de novo design of small-
454 molecule-binding proteins. *Science (New York, NY)*. 2020;369(6508):1227-33.
- 455 14. Orengo CA, Pearl FM, Bray JE, Todd AE, Martin AC, Lo Conte L, et al. The CATH
456 Database provides insights into protein structure/function relationships. *Nucleic acids
457 research*. 1999;27(1):275-9.
- 458 15. Pan X, Thompson MC, Zhang Y, Liu L, Fraser JS, Kelly MJS, et al. Expanding the space of
459 protein geometries by computational design of de novo fold families. *Science (New York,
460 NY)*. 2020;369(6507):1132-6.
- 461 16. Basanta B, Bick MJ, Bera AK, Norn C, Chow CM, Carter LP, et al. An enumerative
462 algorithm for de novo design of proteins with diverse pocket structures. *Proceedings of the
463 National Academy of Sciences of the United States of America*. 2020;117(36):22135-45.
- 464 17. Eswar N, John B, Mirkovic N, Fiser A, Ilyin VA, Pieper U, et al. Tools for comparative
465 protein structure modeling and analysis. *Nucleic acids research*. 2003;31(13):3375-80.
- 466 18. Sillitoe I, Dawson N, Lewis TE, Das S, Lees JG, Ashford P, et al. CATH: expanding the
467 horizons of structure-based functional annotations for genome sequences. *Nucleic acids
468 research*. 2019;47(D1):D280-D4.
- 469 19. Lemán JK, Weitzner BD, Lewis SM, Adolf-Bryfogle J, Alam N, Alford RF, et al.
470 Macromolecular modeling and design in Rosetta: recent methods and frameworks. *Nat
471 Methods*. 2020;17(7):665-80.
- 472 20. Huang PS, Boyken SE, Baker D. The coming of age of de novo protein design. *Nature*.
473 2016;537(7620):320-7.

- 474 21. Skolnick J, Gao M. Interplay of physics and evolution in the likely origin of protein
475 biochemical function. *Proceedings of the National Academy of Sciences of the United*
476 *States of America*. 2013;110(23):9344-9.
- 477 22. Skolnick J, Gao M, Zhou H. How special is the biochemical function of native proteins?
478 *F1000Res*. 2016;5.
- 479 23. Baker D. What has de novo protein design taught us about protein folding and biophysics?
480 *Protein Sci*. 2019;28(4):678-83.
- 481 24. Davey JA, Damry AM, Euler CK, Goto NK, Chica RA. Prediction of Stable Globular
482 Proteins Using Negative Design with Non-native Backbone Ensembles. *Structure*.
483 2015;23(11):2011-21.
- 484 25. Chaudhury S, Lyskov S, Gray JJ. PyRosetta: a script-based interface for implementing
485 molecular modeling algorithms using Rosetta. *Bioinformatics*. 2010;26(5):689-91.
- 486 26. Park H, Bradley P, Greisen P, Jr., Liu Y, Mulligan VK, Kim DE, et al. Simultaneous
487 Optimization of Biomolecular Energy Functions on Features from Small Molecules and
488 Macromolecules. *J Chem Theory Comput*. 2016;12(12):6201-12.
- 489 27. Alford RF, Leaver-Fay A, Jeliazkov JR, O'Meara MJ, DiMaio FP, Park H, et al. The Rosetta
490 all-atom energy function for macromolecular modeling and design. *J Chem Theory Comput*.
491 2017.
- 492 28. O'Boyle NM, Banck M, James CA, Morley C, Vandermeersch T, Hutchison GR. Open
493 Babel: An open chemical toolbox. *J Cheminform*. 2011;3:33.
- 494
- 495

496 **TABLES**

497

498 **Table 1. Number of matched binding sites**

Binding site library	Match type	Native Rossmann	Native NTF2	De novo Rossmann	De novo NTF2
All binding sites	fast	6860 (248)*	8761 (795)	9034 [2442]**	8909 [943]
	Rosetta	5896 (212)	7450 (580)	7475 [1791]	7548 [678]
3 protein residue binding sites	fast	3556 (324)	5714 (1306)	6537 [3305]	6128 [1720]
	Rosetta	2142 (199)	3541 (807)	3715 [1772]	3686 [952]

499

500 * Numbers in parentheses are binding sites that cannot be matched to *de novo* scaffolds with
501 the same topology.

502 ** Numbers in square brackets are binding sites that cannot be matched to native scaffolds with
503 the same topology.

504

505 **Table 2. Dependency of matching success on binding site size (number of protein**
506 **residues)**

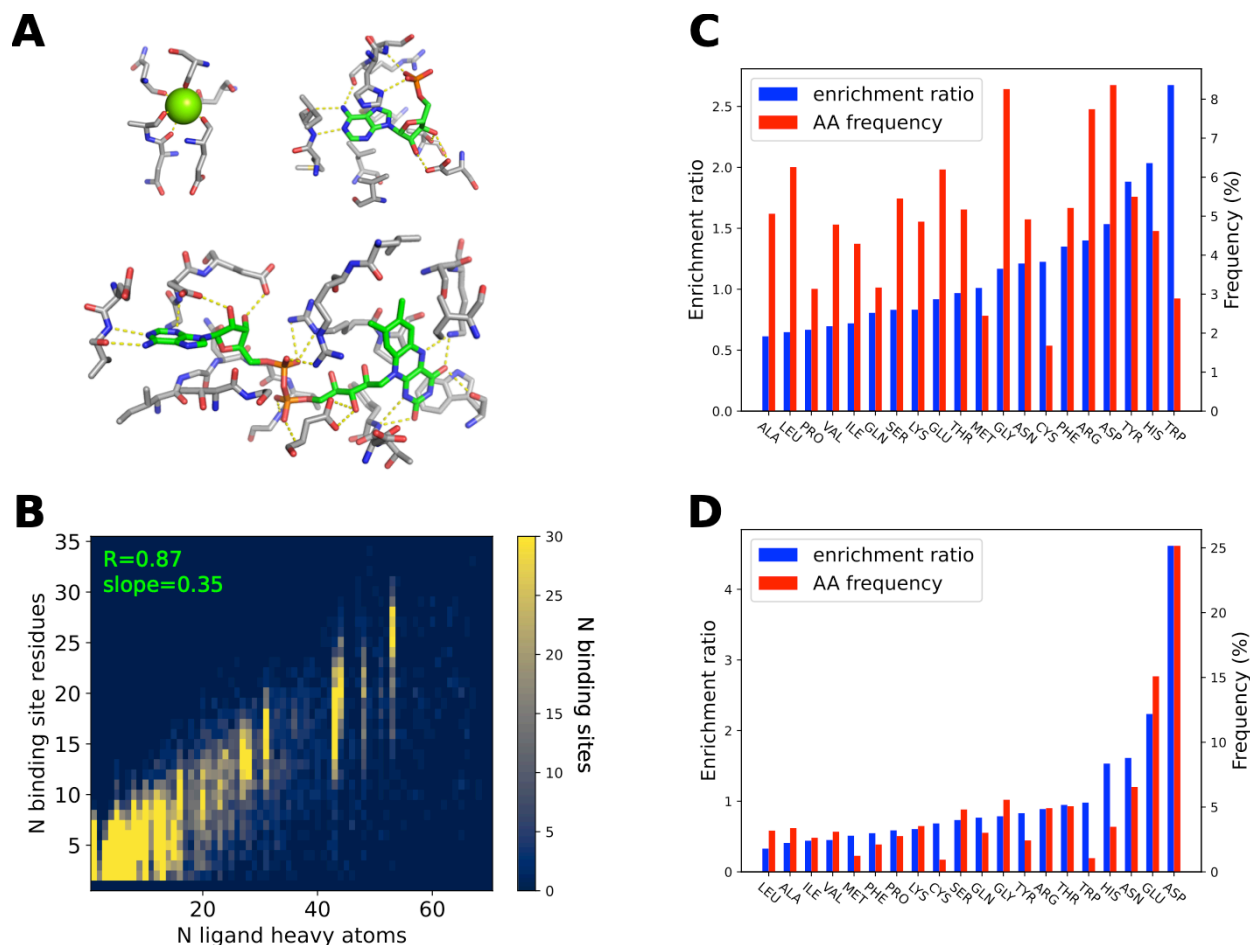
Binding site size	Native Rossmann		Native NTF2		De novo Rossmann		De novo NTF2	
	success count	success rate	success count	success rate	success count	success rate	success count	success rate
2	4590	80.9%	5340	94.2%	5328	93.8%	5359	94.4%
3	1182	21.4%	1792	32.5%	1853	33.4%	1882	33.9%
4	118	2.7%	272	6.3%	281	6.5%	276	6.4%
5	6	0.2%	38	1.4%	12	0.4%	27	1.0%
6	0	0	6	0.4%	1	0.06%	3	0.2%
7	0	0	2	0.2%	0	0	1	0.1%

507

508

509 **FIGURES**

510 **Figure 1**



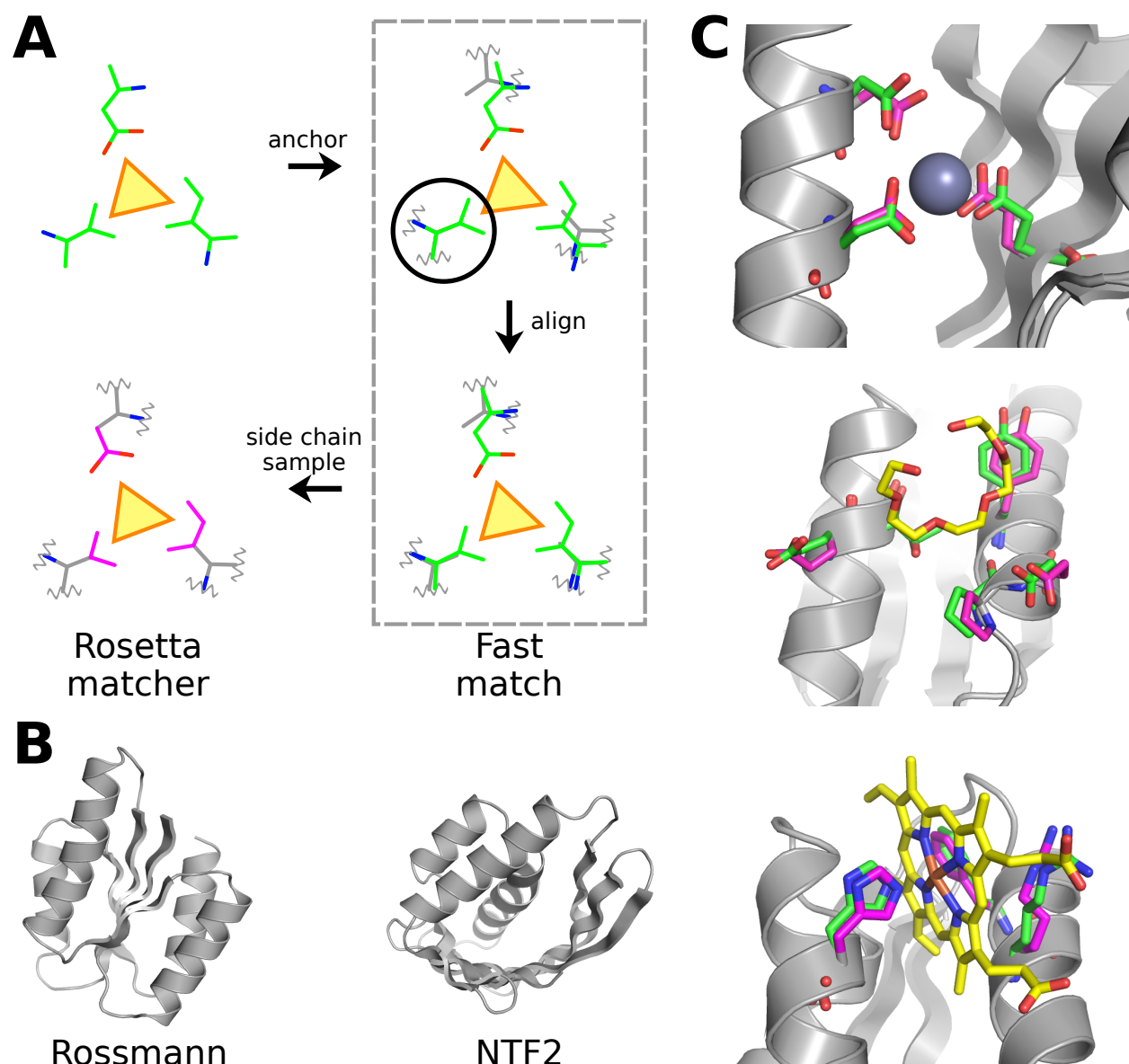
511

512 **Fig 1. The ligand-binding site library.**

513 **A.** Binding site examples. The Mg^{2+} ion is shown as a sphere; small molecules and protein
 514 residues are shown as sticks; carbon atoms are colored in green (small molecule) or grey
 515 (protein residues); oxygen atoms are colored in red; nitrogen atoms are colored in blue; polar
 516 interactions are shown as yellow dashed lines. **B.** Joint distribution of binding site sizes
 517 (numbers of binding site protein residues) and numbers of ligand heavy atoms. Binding site
 518 sizes are linearly correlated with the numbers of ligand heavy atoms. **C, D.** Amino acid (AA)
 519 frequencies (red, right y-axis) in ligand-binding sites and enrichment ratios (blue, left Y-axis) in
 520 ligand-binding sites compared to all residues in a protein. **C.** Distributions of all ligand binding
 521 sites. **D.** Distributions of single heavy atom ligand binding sites.

522

523 **Figure 2**



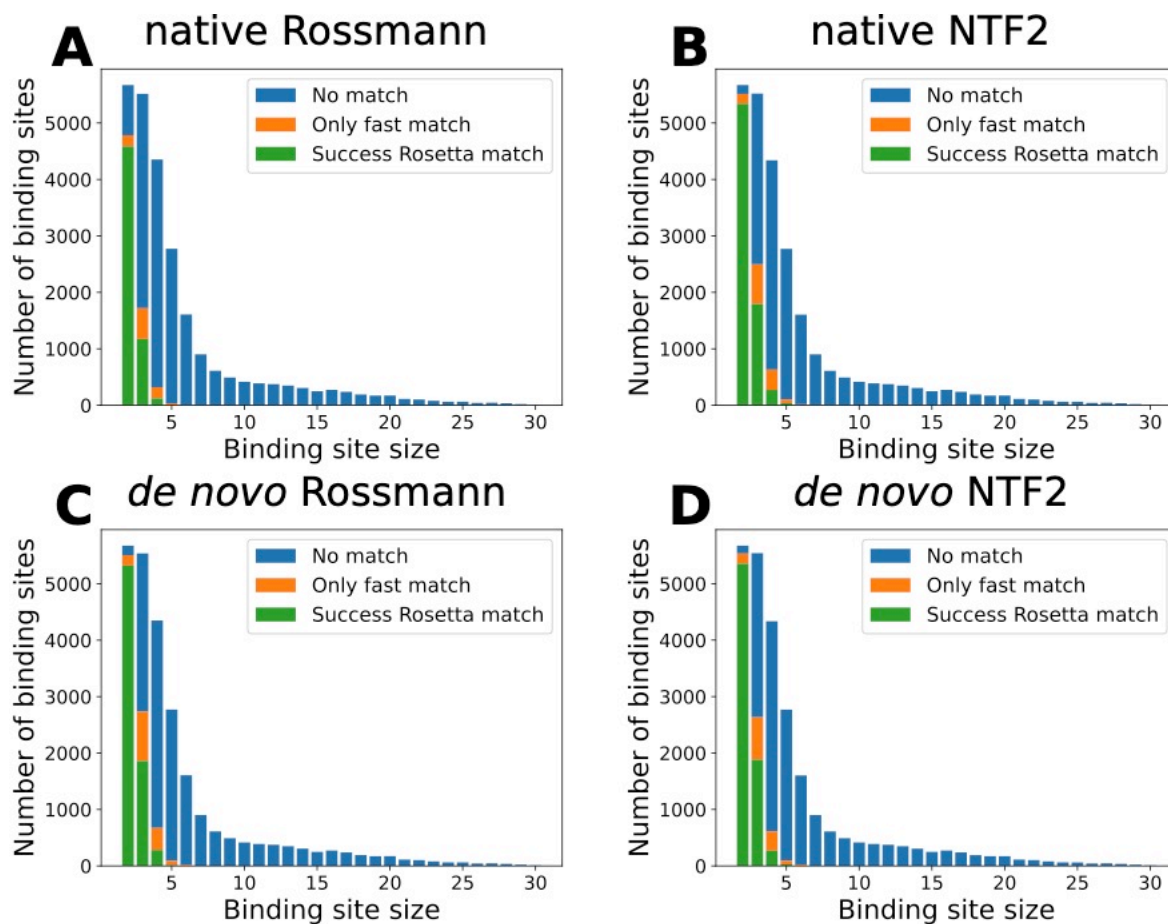
524

525 **Fig 2. Matching ligand binding sites to scaffold libraries.**

526 **A.** Schematic of the matching protocol. The ligand is represented as a yellow triangle. The
527 ligand-binding site as a rigid body (green) is first matched to the scaffold (grey) by anchoring to
528 a scaffold residue shown in the black circle. Then the binding site residues are aligned to the
529 corresponding scaffold residues. Finally, the standard Rosetta matcher is applied to build the
530 binding site side chains (magenta) onto the scaffold. **B.** The binding sites are matched to native
531 and *de novo* designed scaffold families with Rossmann or NTF2 fold topologies. **C.** Examples of
532 matches. The coloring scheme is the same as **A.**

533

534 **Figure 3**



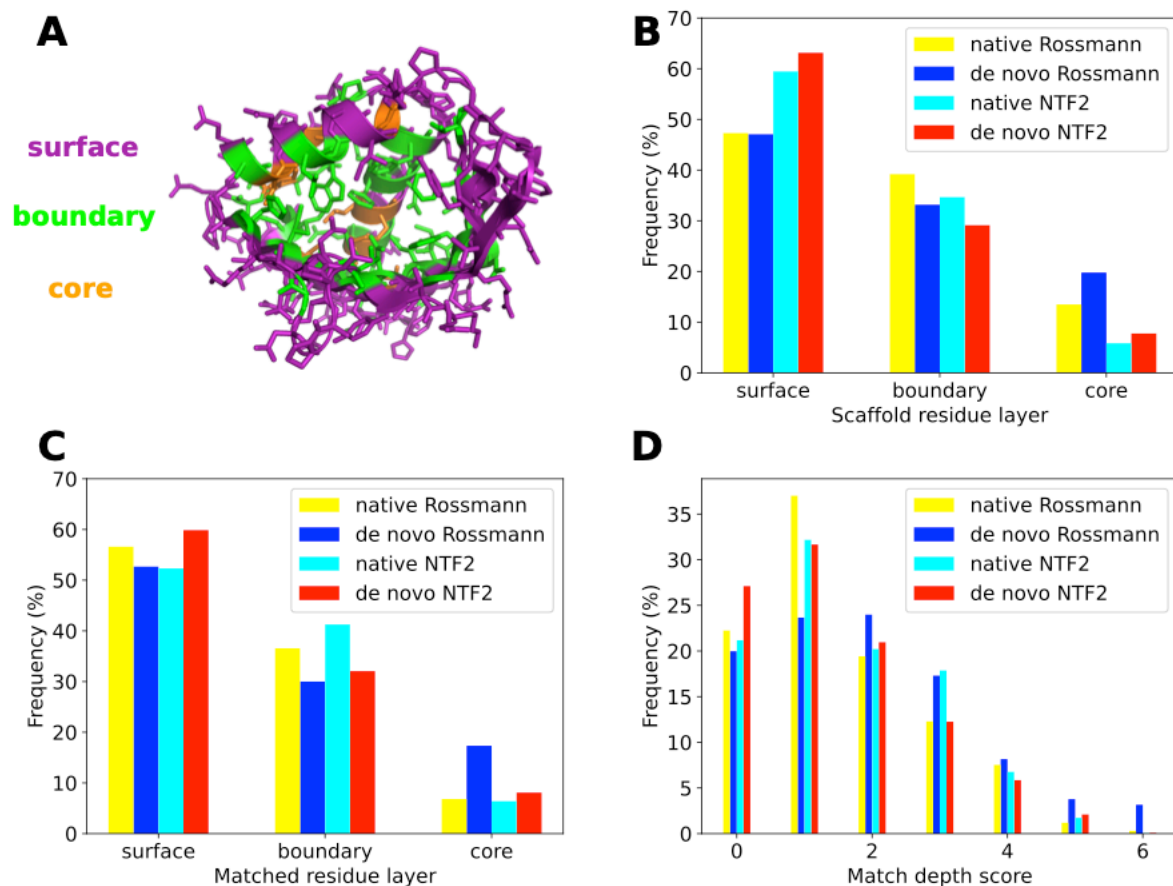
535

536 **Fig 3. Matchability of ligand binding sites depends on the binding site size.**

537 Histograms of numbers of matches vs binding site sizes (number of protein residues in the
538 binding site). Bindings sites that cannot be matched to any scaffold are shown in blue. Bindings
539 sites that can be matched to at least one scaffold by the fast match method but cannot be
540 matched by the standard Rosetta matcher are shown in orange. Binding sites that can be
541 matched to at least one scaffold by the standard Rosetta matcher are in green. **A-D.** Results for
542 4 scaffold libraries; scaffold sets are indicated in each panel title.

543

544 **Figure 4**



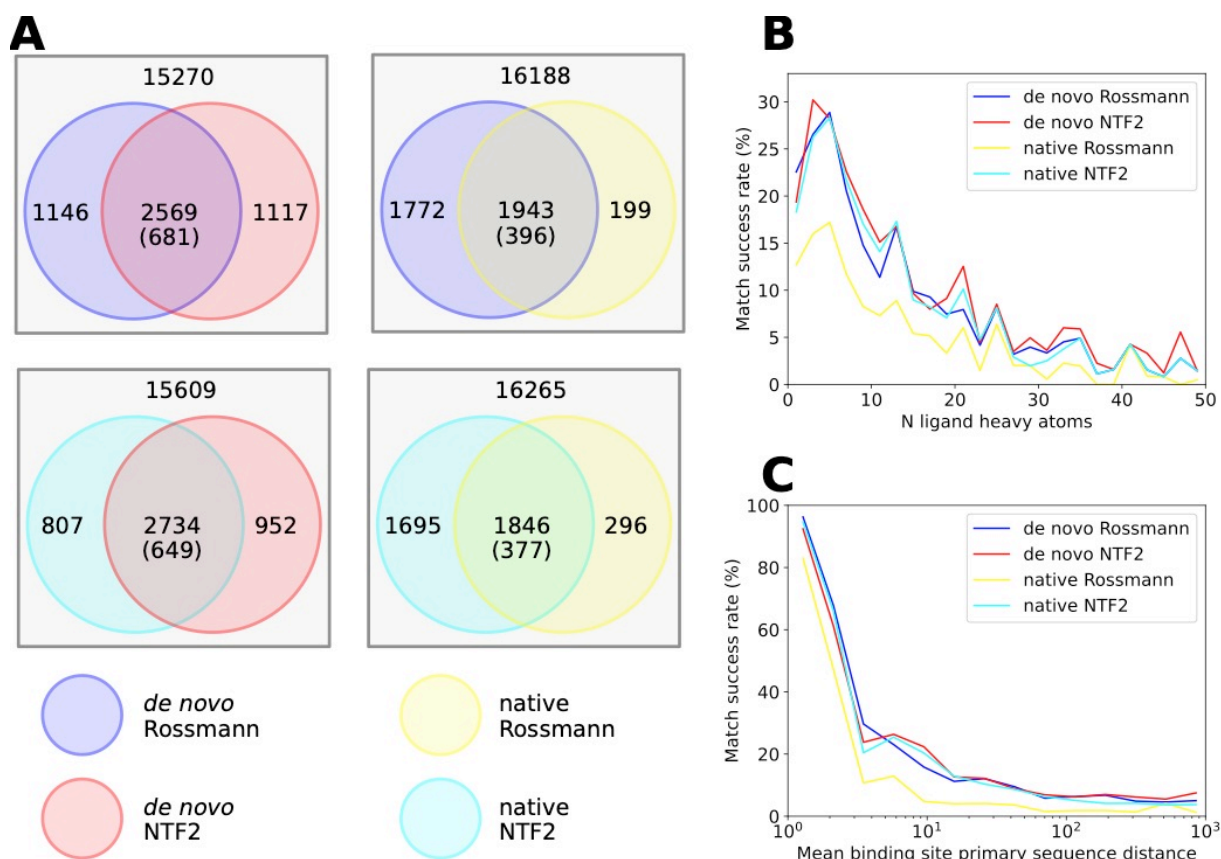
545

546 **Fig 4. Ligand binding sites are matched to all layers of scaffolds.**

547 **A.** An example of scaffold residue layers assigned to a scaffold (PDB:3FH1) from the native
 548 NTF2 fold family by the Rosetta Layer residue selector. The surface, boundary and core layers
 549 are colored in purple, green and orange, respectively. **B.** Distributions of residue layers in
 550 different scaffold libraries. **C.** Distributions of residue layers of binding sites matched to different
 551 scaffold libraries. **D.** Distributions of binding site depth scores matched to different scaffold
 552 libraries.

553

554 **Figure 5**



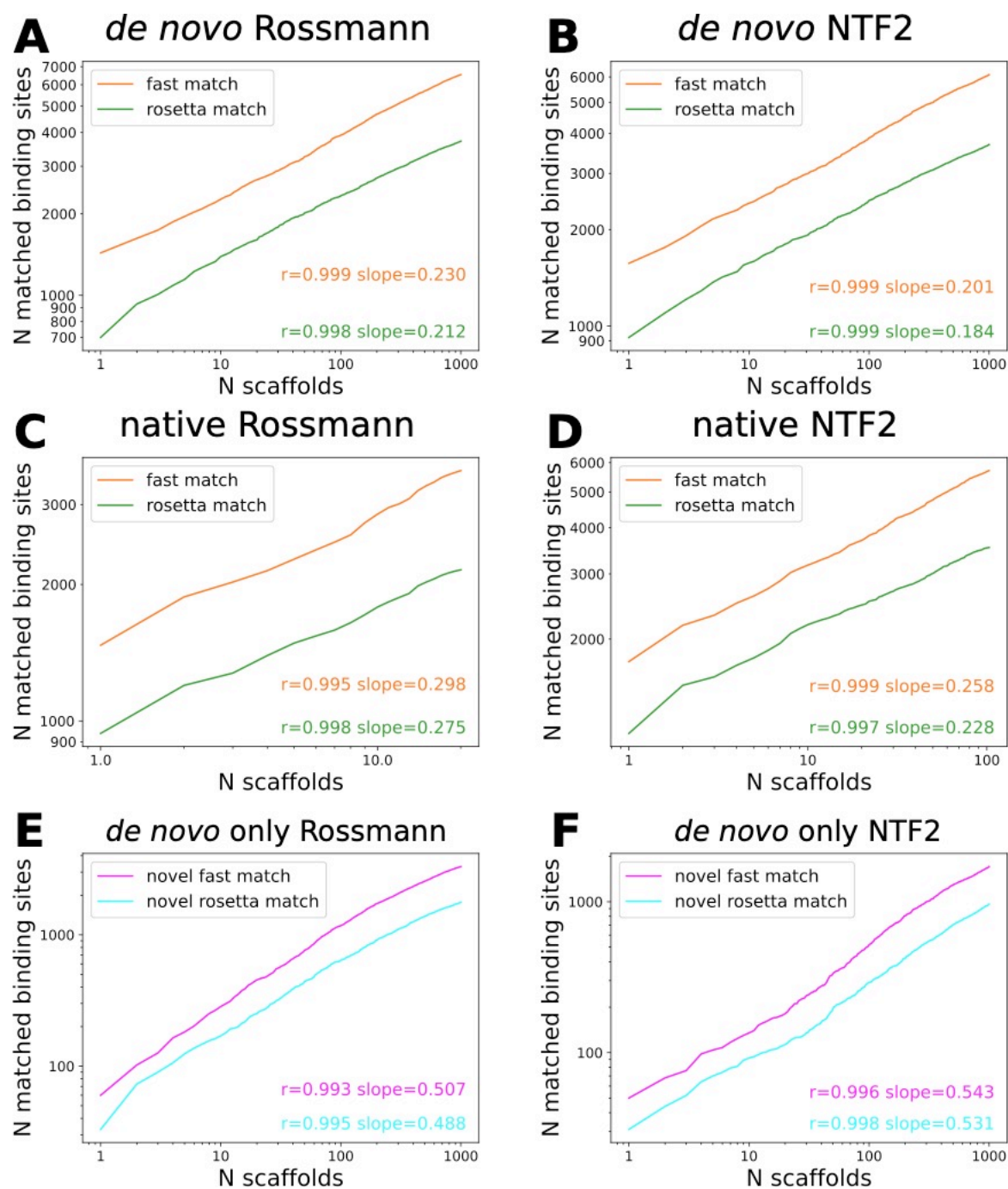
555

556 **Fig 5. Features affecting matching success rates of 3-residue ligand binding sites.**

557 **A.** Venn diagrams of the number of Rosetta-matched 3-residue binding sites between pairs of
558 scaffold sets. The number in the overlapping region is the observed number of binding sites that
559 can be matched to both scaffold sets, with the expected number in parentheses. The number in
560 the non-overlapping region within a circle denotes the binding sites that can only be matched to
561 this scaffold set. The number outside the circles denotes the binding sites that cannot be
562 matched to either of the two scaffold sets. **B.** The numbers of ligand heavy atoms are negatively
563 correlated with the match success rates. **C.** The mean primary sequence distances between
564 binding site residues are negatively correlated with match success rates.

565

566 **Figure 6**



567

568 **Fig 6. Numbers of matches scale as power-law functions of numbers of scaffolds in fold**
 569 **families.**

570 **A-D.** Log-log plots of the number of 3-residue matches vs the number of scaffolds. **E-F.** Log-log
 571 plots of the number of 3-residue binding sites that can only be matched to *de novo* scaffolds of
 572 specific topologies vs the number of scaffolds.

573 **Supporting information**

574 **S1 Table. Ligand type frequencies in the binding site library.**

575 **S1 File. Summary tables of matching results to all fold families.**