

Protein–Protein Interfaces: Properties, Preferences, and Projections

Jeffrey J. Headd, Y. E. Andrew Ban, Paul Brown, Herbert Edelsbrunner,
Madhuwanti Vaidya, and Johannes Rudolph*

Department of Biochemistry, Duke University Medical Center, Durham, North Carolina 27710

Received January 12, 2007

Herein, we study the interfaces of a set of 146 transient protein–protein interfaces in order to better understand the principles of their interactions. We define and generate the protein interface using tools from computational geometry and topology and then apply statistical analysis to its residue composition. In addition to counting individual occurrences, we evaluate pairing preferences, both across and as neighbors on one side of an interface. Likelihood correction emphasizes novel and unexpected pairs, such as the His–Cys pair found in most complexes of serine proteases with their diverse inhibitors and the Met–Met neighbor pair found in unrelated protein interfaces. We also present a visualization of the protein interface that allows for facile identification of residue–residue contacts and other biochemical properties.

Keywords: protein–protein interactions • protein interfaces • geometric topology • visualization

Introduction

Protein–protein interactions play a significant role in the majority of intracellular processes, and understanding how proteins transiently form complexes is essential for grasping the nuances of biological systems. While we have yet to derive universal rules that allow us to identify interaction sites a priori or to reliably predict protein docking, formation of protein complexes and their subsequent stability have been linked to a few specific interfacial residues. These residues, commonly called hotspots, contribute the bulk of the binding energy between proteins, whereas a majority of other residues apparently serve as tolerant bystanders. Thus, the composition of surface residues involved in transient interactions is important to their function. With the ever increasing number of available high-resolution structures of protein complexes, studying the residue composition and pairing preferences of known protein–protein interfaces allows a better understanding of the fundamentals of protein–protein association.

Many prior studies have examined the composition of protein–protein recognition sites from a number of different perspectives.^{1–21} The work of Chakrabarti and Janin¹ has shown these sites have an amino acid composition similar to the overall protein surface. They break the sites down to core and rim regions, with the core having a composition different from the rest of the protein surface, suggesting that protein interfaces possess some unique characteristics. Glaser et al.² studied the pairing preferences for a much larger number of protein–protein complexes by including homo-dimeric complexes in addition to transient hetero-complexes. Though statistically more rigorous, their work does not take into account the

inherent interfacial differences between transient hetero-complexes and obligate homo-dimers, with the latter more resembling protein interiors. Ofra and Burkhard³ have accounted for these differences by classifying protein–protein interfaces into six categories based on interaction type and studying them individually. Mintz et al.²⁰ have developed a method for clustering similar interfaces based on shape and location of chemical functional groups from the entire PDB wherein the overall data set is again biased toward homo-dimers. Ma et al.¹⁵ have focused on elucidating the structurally conserved residues at protein–protein interfaces. Although these prior statistical studies have given us a foundation for thinking about protein–protein interactions, with a few exceptions,^{19,22} they have not found significant applications to the prediction of protein interaction sites, docking of proteins, and the identification of hotspot residues.

This lack of progress in employing the results of previous statistical studies can be attributed to two key drawbacks. First, the limited number of structurally characterized transient protein–protein complexes introduces significant bias, and the resulting statistics have potentially limited applications.^{1,8} As noted above, larger data sets contain disproportionately more homo-dimeric complexes^{2,3,20} or are assembled through automated processes that may include nonphysiological protein–protein pairs such as crystal contacts. Second, the lack of a rigorous definition of what constitutes a protein interface makes reliable automation for broader scale analysis difficult. Frequently, studies have relied on either a distance cutoff (e.g., at most 6 Å between atoms across the interface)^{2,3} and/or an area cutoff (e.g., at least 0.1 Å² increase in solvent accessible area upon separation).^{1,2,5–14} The choice of a cutoff, which does not arise from an intrinsic property of protein–protein complexes, can greatly influence the size and composition of the protein interface and frequently gives rise to artificial holes, 80

* To whom correspondence should be addressed: Johannes Rudolph, Department of Biochemistry, Mailstop 3813 LSRC, C125 Duke University Medical Center, Durham, NC 27710. Tel, (919) 668-6188; fax, (919) 613-8642; e-mail, johannes@alum.mit.edu.

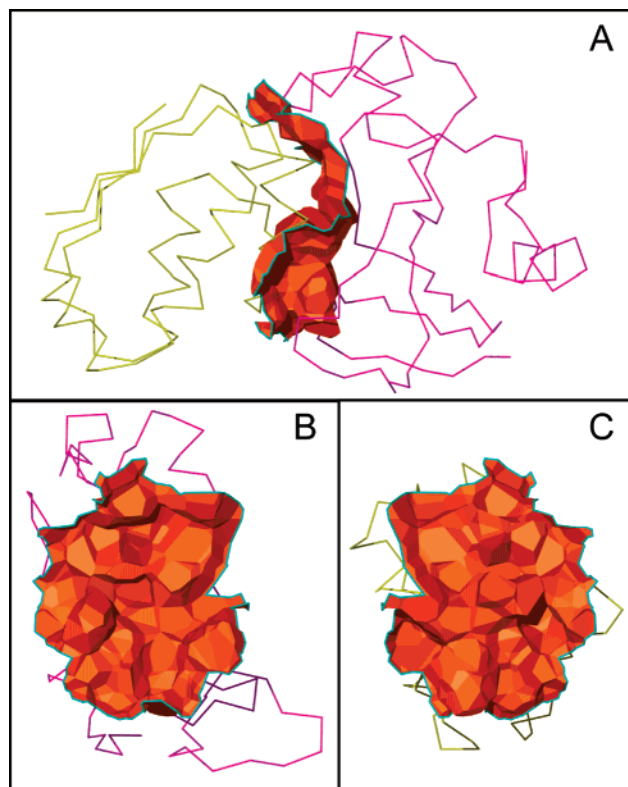


Figure 1. Three-dimensional interface. (A) Interface surface for barnase/barstar (1BRS) complex, with barnase colored purple and barstar colored yellow; (B) image from panel A rotated to show the barstar side of the interface against the barnase protein; (C) image from panel A rotated to show the barnase side of the interface against the barstar protein.

81 overlaps, or extensions. Thus, we and others have developed
 82 alternative definitions of the interface between two proteins
 83 based on geometric topology,²¹ “halfway points”,¹⁶ or molecular
 84 contacts.^{17,18}

85 To address the limitations of previous methods and handle
 86 a wider range of protein–protein complex types, we base our
 87 work on a concrete and unambiguous definition of the
 88 interface. Our interface surface is a subset of the Voronoi
 89 diagram of the entire complex. As seen in Figure 1, the surface
 90 separates two chains in a protein complex and resembles a
 91 wrinkled sheet of paper. Its polygonal assembly captures
 92 features of the contacts between two chains in more detail than
 93 previously described methods, better representing the complex-
 94 ity and complementarity of protein–protein interactions.
 95 Its hierarchical construction allows for a retraction of the
 96 interface surface to a core region that is highly enriched in
 97 hotspot residues.²¹

98 In this study, we analyze the amino acid composition of the
 99 interface of a set of 135 manually selected high-resolution
 100 protein–protein complexes yielding 146 protein interfaces. Our
 101 database of transient protein interfaces is more than twice as
 102 large as previously assembled, hand-culled data sets, and is
 103 significantly more diverse. Using our definition of protein
 104 interfaces, we not only derive a more rigorous statistical
 105 analysis, but are also able to directly measure pairing prefer-
 106 ences across and on one side of the surface. Our analysis reveals
 107 a number of intriguing results, including a higher than expected
 108 contribution of backbone atoms. Finally, we provide a novel
 109 flattened view of the interface surface, allowing for facile

110 visualization of interfacial residue composition, residue con-
 111 tacts, and comparison of homologous structures.

Experimental Section

Interface Construction. We define a protein–protein inter-
 113 face surface as described previously.²¹ Briefly, the surface
 114 formed by two or more proteins is a subset of the Voronoi
 115 diagram of the set of spheres making up the space-filling
 116 diagrams of the involved proteins. We center a space-filling ball
 117 at each atom, and grow the balls simultaneously such that the
 118 Voronoi cells do not change. As a ball grows, we clip it to within
 119 its Voronoi cell, and each time two or more such clipped balls
 120 intersect, we add the convex hull of their centers (generically
 121 a simplex) to K , which is the dual complex of the Voronoi
 122 diagram. This process creates a filtration of dual complexes,
 123 adding simplices to K until K equals the Delaunay triangulation
 124 (the dual of the Voronoi diagram). Simplices that are formed
 125 at the same time enter the complex at the same time t , which
 126 defines the ordering of the filtration. This ordering gives each
 127 simplex a rank value, which is assigned chronologically.
 128 Participating atoms at the interface are those that share Voronoi
 129 polygons in K with an atom on a complementary protein chain.
 130 The collection of Voronoi polygons that belong to a single residue
 131 are denoted a *tile*. Tiles are defined on both sides of the
 132 surface, defining two tilings, one for each of the two proteins
 133 separated by the surface. Overlapping tiles correspond to inter-
 134 chain pairs and adjacent tiles correspond to intrachain pairs.
 135

136 In our analysis, we use the ordered filtration to unambigu-
 137 ously delineate between boundary and core regions of the
 138 interface surface through a retraction process that isolates those
 139 polygons whose rank value is less than or equal to the median
 140 rank value for a given surface. We define this subset as the *core*
 141 of an interface.

142 Since our surface defines atom pairs across the interface,
 143 residue pairs easily follow. In addition, the total number of
 144 atom pairs per residue pair can be determined, clarifying the
 145 extent of the interaction. We can also identify residue neighbor
 146 pairs on the same side of the interface, defined by atoms whose
 147 Voronoi polygons share an adjacent edge in the interface
 148 surface. Finally, the concrete definition allows for normalization
 149 by area and perimeter contributions of residues.

Interface Flattening. For the purpose of visualization, we
 150 fully triangulate the interface surface and then map the
 151 triangulation into the plane, effectively flattening the surface.
 152 Almost all surfaces defined by only two proteins are simply
 153 connected and can therefore be flattened to a round disk in
 154 the plane. While the disk does not necessarily represent the
 155 general shape of the surface, it is easy to view, and it lends
 156 itself to comparing different interfaces. Importantly, the flat-
 157 tening process preserves all connectivity information, both
 158 across the surface and between neighboring tiles. If a complex
 159 consists of three or more proteins, we flatten the sheets defined
 160 by the pairs separately. Flattening is only performed once per
 161 interface, and retracted regions are removed from display,
 162 leaving the remnants of the original disk to represent the
 163 remaining interface surface. Finally, in the uncommon case in
 164 which an interface surface has nonzero genus, we have to cut
 165 the surface to remove the genus before flattening. For an
 166 example of an interface with nonzero genus, see the neurotoxic
 167 vipoxin complex from Western Sand Viper, PDB code 1JLT, as
 168 depicted in Figure 8 in Ban et al.²¹

169 The algorithm we use for flattening is based on a theorem
 170 by Tutte²³ which states that a convex mapping of a simple,
 171

172 3-connected, planar graph is a valid straight-line embedding
 173 of the graph. Specifically, if we draw the boundary of the surface
 174 as a convex curve in the plane and we express the image of
 175 each interior vertex as a convex combination of the images of
 176 its neighbor vertices, then Tutte's theorem guarantees that we
 177 indeed have a flattening of the triangulation. Let N be the total
 178 number of vertices and $n < N$ the number of interior vertices.
 179 The most straightforward implementation of Tutte's theorem
 180 maps the $N - n$ boundary vertices to equally spaced points in
 181 sequence along a unit circle, and it solves a system of linear
 182 equations to compute the image of the interior vertices. Letting
 183 v_1 to v_N be the vertices of the surface triangulation and μ the
 184 map to the plane, the equation for the i th vertex is

$$\mu(v_i) = \sum_{j=1}^N \lambda_{i,j} \mu(v_j)$$

185 where $\lambda_{i,j} = 0$ if (v_i, v_j) is not an edge in the triangulation, $\lambda_{i,j} =$
 186 $1/d_i$ if (v_i, v_j) is an edge, d_i is the number of neighbor vertices
 187 of v_i , and v_i is assumed to be an interior vertex. We have n
 188 linear equations in n unknowns and therefore a unique
 189 solution. Although the system of equations can be big, it is
 190 necessarily sparse and therefore permits efficient computation.

191 We refer to the above implementation of Tutte's theorem
 192 as the *uniform method* because it treats the neighbors of each
 193 vertex the same way. While it produces topologically accurate
 194 flattened images, it introduces a significant amount of distortion.
 195 To reduce the distortion, we use a more sophisticated
 196 implementation of Tutte's theorem referred to as the *mean*
 197 *value coordinates method* as described by Floater.²⁴ The
 198 boundary vertices are again mapped in sequence to points on
 199 a unit circle, but the spacing along the circle is chosen
 200 according to the lengths of the boundary edges. The second
 201 difference to the uniform method is in the choice of the $\lambda_{i,j}$. To
 202 describe the new weights, we assume v_j is a neighbor of v_i and
 203 let $\alpha_{i,j}$ and $\beta_{i,j}$ be the angles at v_i of the two triangles that share
 204 the edge (v_i, v_j) . Setting

$$w_{i,j} = \frac{\tan(\alpha_{i,j}/2) + \tan(\beta_{i,j}/2)}{\|v_i - v_j\|}$$

205 the corresponding weight used in the linear systems is

$$\lambda_{i,j} = \frac{w_{i,j}}{\sum_{l=1}^N w_{i,l}}$$

206 assuming $w_{i,j} = 0$ if v_i and v_l are not neighbors. For a further
 207 description of how the weights $w_{i,j}$ are calculated and affect
 208 the flattening, see Floater.²⁴ It is easy to see that the $\lambda_{i,j}$ are
 209 non-negative and add up to 1, if summed over all j . In words,
 210 the weights express each interior vertex as a convex combina-
 211 tion of its neighbors; hence, Tutte's theorem applies.

212 To assess the difference between the two implementations,
 213 we measure area distortion as a weighted sum of square ratios
 214 of normalized areas of triangles. Specifically, we normalize such
 215 that the (one-sided) area of the interface surface is 1 and the
 216 area of its image (the disk in the plane) is 1. Denoting a triangle
 217 in the surface by t and its image by $\mu(t)$, we define

$$D = \sum_{\text{triangles } t} \left[\frac{\text{area}(\mu(t))}{\text{area}(t)} \right]^2 \text{area}(t)$$

Table 1. The Data Set Used for All Statistical Analysis, Consisting of 135 Complexes Hand-Culled from the PDB

A00	1BKD	1DN2	1FSS	1JHL	1MLC	1SG1	1WEJ	2KAI
1A0R	1BQL	1DQJ	1FVC	1JQJ	1NCA	1SPB	1Y8R	2MTA
1A22	1BRC	1DVF	1G3N	1JRH	1NFD	1STF	1YCQ	2PCC
1A2K	1BRS	1DX5	1GC1	1JTG	1NMB	1TAB	1YCS	2PTC
1A4Y	1BTH	1EER	1GG2	1KB5	1NSN	1TBQ	1YDR	2SIC
1ACB	1BUH	1EFN	1GLA	1KF6	1NVU	1TCO	1YYM	2SNI
1AGR	1BVK	1EFU	1GOT	1KKL	1OMW	1TGS	1Z92	2TEC
1AHW	1BXI	1F47	1GUA	1KXV	1OSP	1TOC	1ZJD	2TRC
1AIP	1C4Z	1FBI	1HIA	1LOY	1P22	1TT5	2ASS	3HFL
1AK4	1CBW	1FC2	1HWG	1LDK	1PPE	1UDI	2B4J	3HFM
1A07	1CHO	1FDL	1IIR	1LM8	1PPF	1UGH	2B4S	3HHR
1ATN	1CSE	1FIN	1IAI	1MCT	1PVH	1US7	2B5I	3SGB
1AVW	1DAN	1FLE	1IGC	1MDA	1QFU	1USU	2BTF	3TPI
1AVZ	1DFJ	1FQ1	1JDH	1MEL	1RZK	1UUG	2C2V	4CPA
1BI8	1DHK	1FS1	1JEL	1MKW	1SEB	1VFB	2JEL	4HTC

218 It is not difficult to see that $D \geq 1.0$ and $D = 1.0$ if, and only
 219 if, the area of each triangle is the same as the area of its image.
 220 More generally, the smaller D is, the closer μ is to an equiareal
 221 map. We have measured the distortion of flattenings as
 222 computed by the uniform and the mean value coordinates
 223 methods. As shown in Supplemental Table 1 on Supporting
 224 Information, the latter method has $D \leq 2.0$ for almost all cases
 225 and thus introduces significantly less area distortion than the
 226 uniform method with values of D between 5.0 and 64.0.
 227 Because of this difference, the mean value coordinates method
 228 is used throughout this paper. A visual comparison between
 229 the 1BRS interface surface flattened using both methods is
 230 shown in Supplemental Figure 1 in Supporting Information.

231 **Data Set.** Our data set (Table 1) is a significant augmentation
 232 of the set of PDB complexes used by Chakrabarti and Janin.¹
 233 It consists of 135 high-resolution hetero-complexes (ranging
 234 from 1.2 to 3.0 Å resolution) containing 146 protein-protein
 235 interfaces and 10 296 interfacial amino acids. The distribution
 236 of number of interface residues and side-chains versus area of
 237 each complex is depicted in Figure 2.

238 **Single Residue Statistics.** Letting $\#(a)$ be the number of type
 239 a residues that contribute tiles to interface surfaces and $\#_{\text{total}}$
 240 $= \sum_a \#(a)$ the total number of contributing residues (of any
 241 type), we define the *relative frequency of residue type a* equal
 242 to

$$\text{Prob}[a] = \frac{\#(a)}{\#_{\text{total}}}$$

243 which is the probability of a randomly chosen contributing
 244 residue to be of type a . Similarly, letting $\text{area}(a)$ be the total
 245 area contributed by type a residues and $\text{area}_{\text{total}} = \sum_a \text{area}(a)$
 246 the total (two-sided) area of the interface surfaces, we define
 247 the *area-weighted relative frequency of residue type a* equal to

$$\text{Prob}_{\text{area}}[a] = \frac{\text{area}(a)}{\text{area}_{\text{total}}}$$

248 which is the probability that a randomly chosen point and side
 249 of an interface surface belongs to a tile of a type a residue.
 250 Both relative and area-weighted relative frequencies will be
 251 used to calculate the statistics for single residue occurrences
 252 (see Results).

253 **Residue Pair Statistics.** We distinguish between *interchain*
 254 pairs of residues that are separated by at least one shared
 255 Voronoi polygon in the interface surface, and *intrachain* or
 256 *neighbor* pairs of residues whose tiles belong to the same side
 257 and share at least one Voronoi edge in their boundaries. We
 258 note that residue pairs from the same chain are only counted

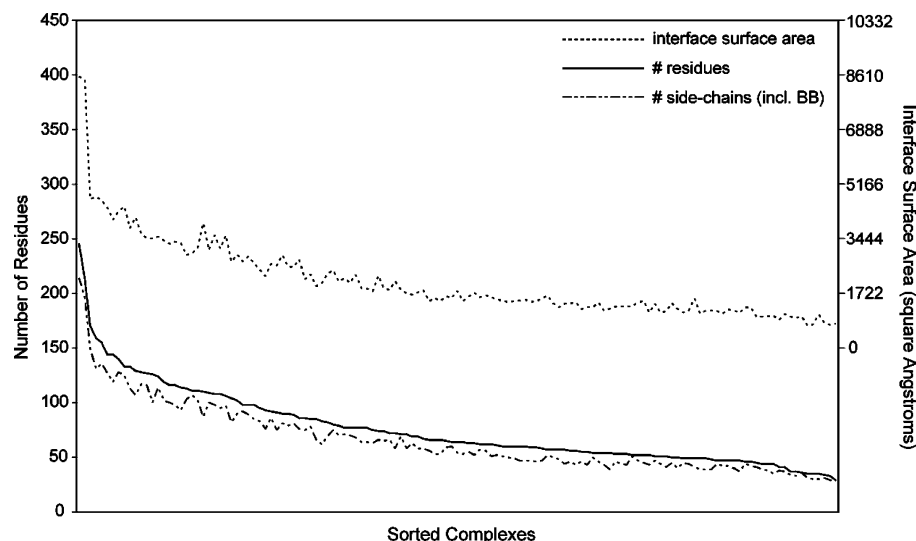


Figure 2. Graphical summary of data sets. From top to bottom: the surface area, the number of residues, and the number of side chains (including backbone as one category) of the complexes sorted by number of residues from left to right. Numbers of residues and side chains are marked on the left, and areas (in Å²) are marked on the right.

259 if they are adjacent to each other at the interface surface, not
 260 if they are adjacent elsewhere within the protein or they are
 261 contiguous along the backbone. It will be convenient to assume
 262 an arbitrary but fixed ordering among the residue types so we
 263 can write $a \leq b$ by which we mean that a is equal to b or a
 264 precedes b in this ordering.

265 Letting $\#C(a,b)$ be the number of interchain pairs in which
 266 one residue is of type a and the other of type b and $C_{\text{total}} =$
 267 $\sum_{a \leq b} \#C(a,b)$ the total number of interchain pairs, we define
 268 the *relative frequency of the pair a, b* to be equal to

$$\text{Prob}[a,b] = \frac{\#C(a,b)}{\#C_{\text{total}}}$$

269 which is the probability that a randomly chosen interchain pair
 270 consists of a type a and a type b residue. We note that the pairs
 271 are unordered, so $\text{Prob}[a,b] = \text{Prob}[b,a]$. To get area weighted
 272 formulas, we observe that $\text{area}_{\text{total}} = \sum_{a \leq b} \text{area}(a,b)$, where
 273 $\text{area}(a,b)$ is the total (two-sided) area contributed by pairs in
 274 which one residue is of type a and the other of type b . We
 275 define the *area-weighted relative frequency of the pair a, b*
 276 to be equal to

$$\text{Prob}_{\text{area}}[a,b] = \frac{\text{area}(a,b)}{\text{area}_{\text{total}}}$$

277 which is the probability that a randomly chosen point on an
 278 interface surface belongs to the tile of a type a residue on one
 279 side and to the tile of a type b residue on the other side of the
 280 interface surface.

281 Letting $\#N(a,b)$ be the number of neighbor pairs in which
 282 one residue is of type a and other of type b and $\#N_{\text{total}} = \sum_{a \leq b}$
 283 $\#N(a,b)$ the total number of neighbor pairs, we define the
 284 *relative frequency of the neighbor pair a, b* to be equal to

$$\text{Prob}[a,b] = \frac{\#N(a,b)}{N_{\text{total}}}$$

285 which is the probability that a randomly chosen neighbor pair
 286 consists of a type a and a type b residue. We note again that
 287 the pairs are unordered, so $\text{Prob}[a,b] = \text{Prob}[b,a]$. If two tiles
 288 share a common portion of their boundary, we can measure

the length or perimeter of that portion and use it as a weight 289
 in our statistics. Letting $\text{perim}(a,b)$ be the total shared perimeter 290
 of tiles contributed by pairs of type a and b and $\text{perim}_{\text{total}} =$ 291
 $\sum_{a \leq b} \text{perim}(a,b)$, we define the *perimeter-weighted relative* 292
frequency of the neighbor pair a, b to be equal to 293

$$\text{Prob}_{\text{perim}}[a,b] = \frac{\text{perim}(a,b)}{\text{perim}_{\text{total}}}$$

294 which is the probability that a randomly chosen point in the 294
 interior of an interface surface that does not belong to the 295
 interior of a tile belongs to the shared boundary of tiles 296
 contributed by a type a and type b residue. 297

298 We define a triplet as an intrachain pair whose members 298
 both form interchain pairs with a third residue on the other 299
 chain. Propensities for triplets are identified by visual inspec- 300
 tion of a representative sample (20–50% of observed cases) for 301
 each type. 302

303 **Likelihood Correction.** The probabilities for inter- and 303
 intrachain pairs depend on the probabilities of single residues. 304
 For example, in the absence of any bias, the relative frequency 305
 of the pair $a \neq b$ is $\text{Prob}[a,b] = 2\text{Prob}[a]\text{Prob}[b]$, where the 306
 factor of 2 accounts for the fact that the pair is unordered. On 307
 the other hand, $\text{Prob}[a,a] = \text{Prob}[a]^2$, without the factor of 2. 308
 It thus makes sense to consider the ratio of the left over the 309
 right side, which in the absence of any bias is one. Taking the 310
 logarithm, we obtain positive numbers for pairs that are more 311
 likely than warranted by the probabilities of its constituents 312
 and negative numbers for pairs that are less likely. The formulas 313
 are 314

$$\text{LogOdds}(a,b) = \log_2 \frac{\text{Prob}[a,b]}{2\text{Prob}[a]\text{Prob}[b]}, \quad \text{where } a \neq b$$

$$\text{LogOdds}(a,a) = \log_2 \frac{\text{Prob}[a,a]}{\text{Prob}[a]\text{Prob}[a]}$$

We use similar area- and perimeter-weighted formulas to study 315
 the bias for or against forming inter- and intrachain pairs. 316
 Previous studies^{2,3} have used the $\text{LogOdds}(a,a)$ formula to 317
 calculate all log odds values, regardless of pairing partners, but 318

319 we treat each pair uniformly as unordered pairs, giving a
320 consistent statistical result.

321 Results and Discussion

322 **Larger Data Set of Transient Hetero-Complexes.** A major
323 drawback of using the characteristics of known complexes to
324 study protein-protein interactions is the limited available data
325 set. Although the PDB holds more than 34 000 structures, only
326 a small fraction of these represents complexes of proteins that
327 can exist independently in a folded native state. The data set
328 is sufficiently limited that statistical filtering for redundancy is
329 often not applied. For example, the largest previous analysis
330 of hand-selected protein-protein recognition sites contains 75
331 transient complexes,⁸ seven of which are complexes of a serine
332 protease bound to a protein inhibitor. Thus, we have attempted
333 to increase the number and diversity of readily available
334 protein-protein complexes for study. Our current list, built on
335 the foundation of Chakrabarti and Janin,¹ adds complexes from
336 the protein docking benchmark data set,⁴ and others culled
337 manually from the literature. New complexes were selected to
338 capture as much interface diversity as possible, taking advan-
339 tage of recent additions to the PDB consisting of multi-protein
340 systems (e.g., ubiquitination). Our data set includes 135 hetero-
341 complexes, some with multiple interfaces (Table 1).

342 Applying our interface definition to the 135 complexes in
343 our database (Table 1), we generate 146 protein-protein
344 interfaces. There are $\#_{\text{total}} = 10\,292$ residues in total, $P_{\text{total}} =$
345 $25\,875$ interchain pairs, and $N_{\text{total}} = 23\,545$ neighbor pairs in
346 our data set. This yields an average of $2P_{\text{total}}/\#_{\text{total}} = 5.03$ pairs
347 and $2N_{\text{total}}/\#_{\text{total}} = 4.58$ neighbors per residue. The most
348 common residue is present $\#(\text{Glu}) = 705$ times and the least
349 common residue is found $\#(\text{Met}) = 193$ times, demonstrating
350 the significant size of our interface database. The interfaces
351 range in size from 709.4 to 8,544.8 Å², with an average \pm
352 standard deviation of 2094.2 ± 1263.3 Å², counting both sides
353 so as to compare with accessible surface area methods that
354 have previously shown an average of 1906 ± 759 Å².¹ The
355 distribution of area for our data set is depicted in Figure 2. As
356 a point of reference, there are 64,640 residues in the 135
357 complexes, yielding 54,348 non-interface residues, such that
358 the interface regions account for almost 16% of the residues
359 in the data set.

360 **Flattening.** Visualization is another significant hurdle faced
361 when attempting to understand protein-protein interfaces. A
362 common representation shows the separate surfaces of the two
363 contributing proteins, often as a GRASP view.²⁵ From these
364 disjointed views, it is difficult to see the relative area contribu-
365 tions of residues or residue interactions across the interface.
366 Gabdoulline and Wade²⁶ have previously described a method
367 that attempts to remedy this dilemma by projecting their
368 analytically defined interface onto a flat surface with approx-
369 imate conservation of area. We introduce here a flattened
370 view and associated web tool that takes advantage of our
371 interface definition to yield a voidless map independent of
372 embedding and without overlapping points. As described in
373 Experimental Section, we flatten the potentially complicated
374 surface embedded in three-dimensional space to a disk.
375 Critically, the flattening procedure retains all neighbor con-
376 nectivity, both across as well as on each side of the interface.
377 The flattened interface can be colored by selectable attributes.
378 For example, the default view in our MAPS web tool divides
379 and colors the flattened interface by contributing residues
380 (Figure 3A,B). Thick black lines separate the tiles (residue

contributions) that are colored by type according to the 381
indicated palette. Thin black lines separate the atom contribu- 382
tions within each tile, and both atom name and residue 383
identification are obtained using a mouse-over tool. Other 384
selectable attributes for viewing include atom type, electrostat- 385
ics, backbone versus side-chain, and distance between atoms 386
across the interface. For example, the large contribution of 387
backbone atoms is immediately apparent in Figure 3F, whereas 388
hydrogen bonds across the interface are seen as the regions of 389
closest distance in Figure 3E. Of course, any of the selected 390
attributes can originate from either of the two contributing 391
proteins, as seen for residue type in Figure 3A,B. 392

393 To directly compare the attributes from two complexed 393
proteins, we have created a merged view (Figure 3C,D) in which 394
the bottom side remains unchanged while the top side is 395
reduced to frames that outline its tiles. These frames have a 396
black edge to aid in viewing and are colored to indicate the 397
selected attribute (e.g., residue type, atom type, distance). Thus, 398
it is easy to see which residue/attribute from one protein sits 399
across from which residue/attribute from the other protein. 400

401 Our MAPS web tool ([http://biogeometry.cs.duke.edu/research/](http://biogeometry.cs.duke.edu/research/docking/index.html)
402 [docking/index.html](http://biogeometry.cs.duke.edu/research/docking/index.html)) contains interfaces for all 135 complexes
403 from our data set (Table 1) and displays both the 3D interface
404 between the generating protein chains as well as the flattened
405 view with functionality as described above.

406 **Comparison of Related Protein Interfaces.** A particularly 406
powerful application of flattened interfaces is the direct 407
comparison of similar protein complexes to reveal key similar- 408
ities and differences. As an example, we consider the pig RNase 409
inhibitor (Figure 4A) bound to bovine RNase (1DFJ) (Figure 4C) 410
compared to the human RNase inhibitor (Figure 4B) bound to 411
angiogenin (1A4Y) (Figure 4D), two well-characterized protein- 412
protein complexes. The pig and human RNase inhibitors are 413
highly homologous (about 77% identity), whereas bovine RNase 414
and angiogenin, despite their similar protein folds, are signifi- 415
cantly different in sequence (only about 36% identity) and 416
function (ribonuclease activity vs inducer of angiogenesis). The 417
similarity of the binding interfaces is readily seen from the side 418
of the inhibitor (Figure 4A,B). Clearly, the RNase inhibitors use 419
the corresponding binding sites, showing the same interfacial 420
residues, including the hotspot residues found in both com- 421
plexes (Tyr434 and Asp435). Some differences include the 422
presence of Trp375 in 1A4Y but not 1DFJ and the more subtle 423
change of the backbone-mediated interaction of residue 436, 424
which is an Ile in 1A4Y and a Thr in 1DFJ. Despite a few 425
similarities in the region of the hotspot interaction (His119, 426
Lys41, Gln11 in 1DFJ vs His114, Lys40, Gln117 in 1A4Y), the 427
protein interfaces are dramatically different from the other side. 428
The direct visual comparison of two or more protein interfaces 429
is expected to facilitate experimental investigations aimed 430
toward understanding how protein-protein interaction can 431
attain both specificity and flexibility. 432

433 **Single Amino Acid Statistics.** Understanding the composi- 433
tion of binding sites is essential to understanding how transient 434
protein-protein complexes form. Despite differences in the 435
definition of the interface and the database of protein com- 436
plexes, our amino acid composition for protein interfaces is 437
comparable to previous studies.¹ For example, four of the five 438
most (Ser, Glu, Gly, and Asp) and the five least (Met, Trp, Cys, 439
His, and Phe) represented amino acids are identical in the two 440
data sets. Also, our retraction process toward the protected core 441
of about 50% area (see Experimental Section) shows a trend 442
in enrichment of amino acids similar to the selection of the 443

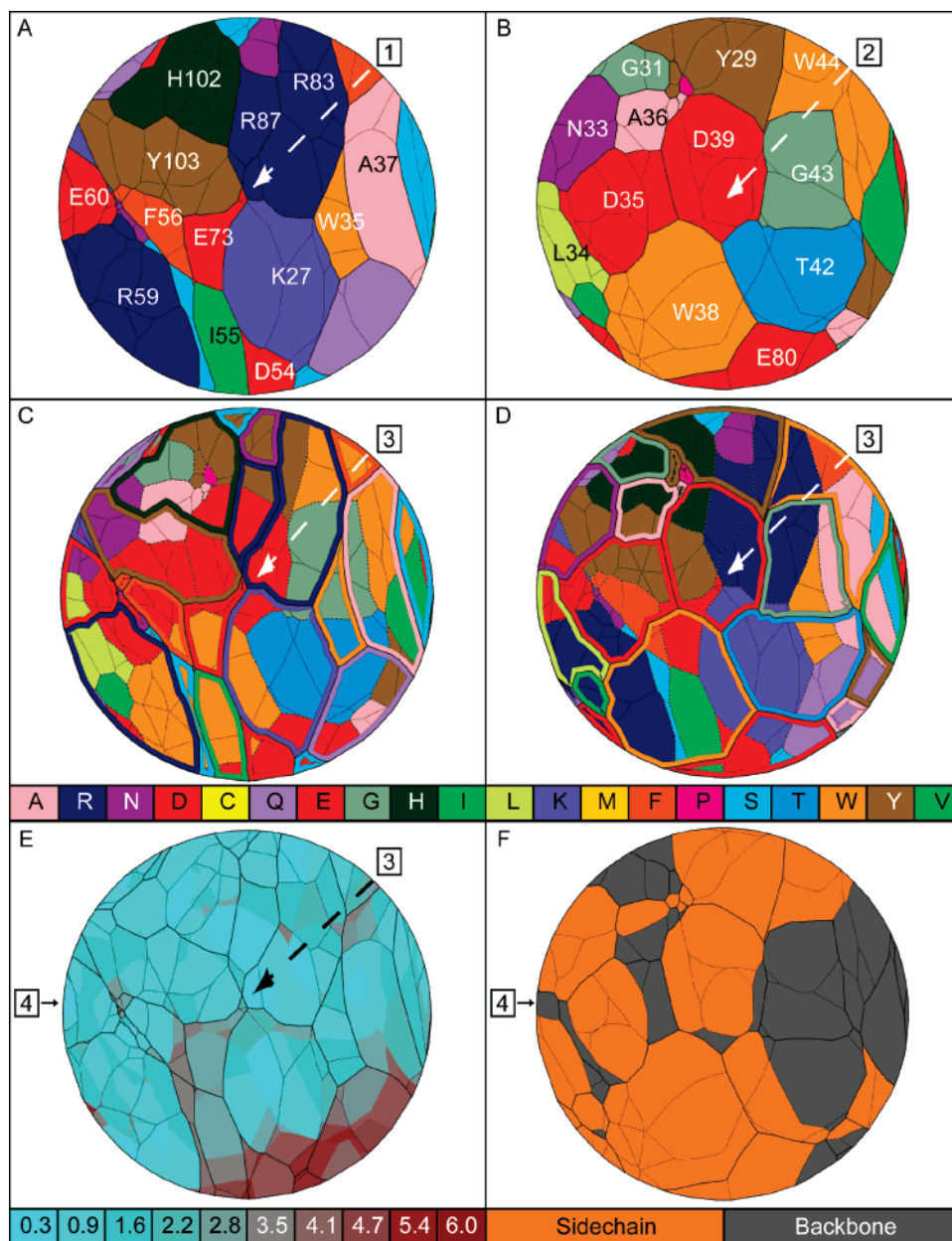


Figure 3. Flattened views of the interface of barnase/barstar (1BRS). (A) Interface surface of barnase (chain A), colored by residue type. The N_{γ2} atom of Arg 83 is indicated by 1. (B) Interface surface of barstar (chain D), colored by residue type. The O_{δ1} atom of Asp 39 is indicated by 2. (C and D) Merged view of the interface between barnase and barstar, colored by residue type, as viewed from barnase in C, and from barstar in D. The salt-bridge between N_{γ2} Arg 83 (barnase) and O_{δ1} Asp 39 (barstar) is indicated by 3. (E) Interface surface of barnase colored by distance gradient. The Arg 83–Asp 39 salt-bridge is indicated by 3, and a hydrogen bond between a backbone N from Leu 34 of barstar to O_{ε2} of Glu 60 of barnase is indicated by 4. (F) Interface surface of barstar colored by backbone/side chain. The Leu 34–Glu 60 hydrogen bond is indicated by 4.

444 53% core residues that contribute 72% of the interface area in
 445 Chakrabarti and Janin.¹ In particular, both methods and data
 446 sets reveal a substantial enrichment in aromatic and some
 447 hydrophobic residues, especially Leu and Ile, and a decrease
 448 in all charged residues, especially Asp and Glu (Table 2).

449 Because evolutionary selection of amino acids occurs via
 450 side-chain variation on an invariant backbone, we chose to
 451 analyze backbone contributions to protein interfaces separately.
 452 This was accomplished by designating backbone its own
 453 category, which consists of all backbone atoms of the residues
 454 making up a protein chain. The side chains are the residues
 455 without their backbone atoms. As a result, Gly is absorbed into

the backbone because it does not have a sidechain. In the case
 where both side chain and backbone atoms from one residue
 appear at the interface, both are counted in the appropriate
 categories. As previously noted by Lo Conte et al.,⁸ backbone
 atoms comprise a significant portion of the interface. By
 frequency, about 32% of all interfacial pairs are either backbone–
 backbone or backbone–side-chain contacts, with backbone
 atoms accounting for about 23% of total interface area (e.g.,
 Figure 3F). Though prior studies have shown that backbone
 carbonyl O atoms are commonly involved in hydrogen bonding
 at protein–protein interfaces,^{8,20} we find that all backbone
 atoms make a significant contribution. While not significantly

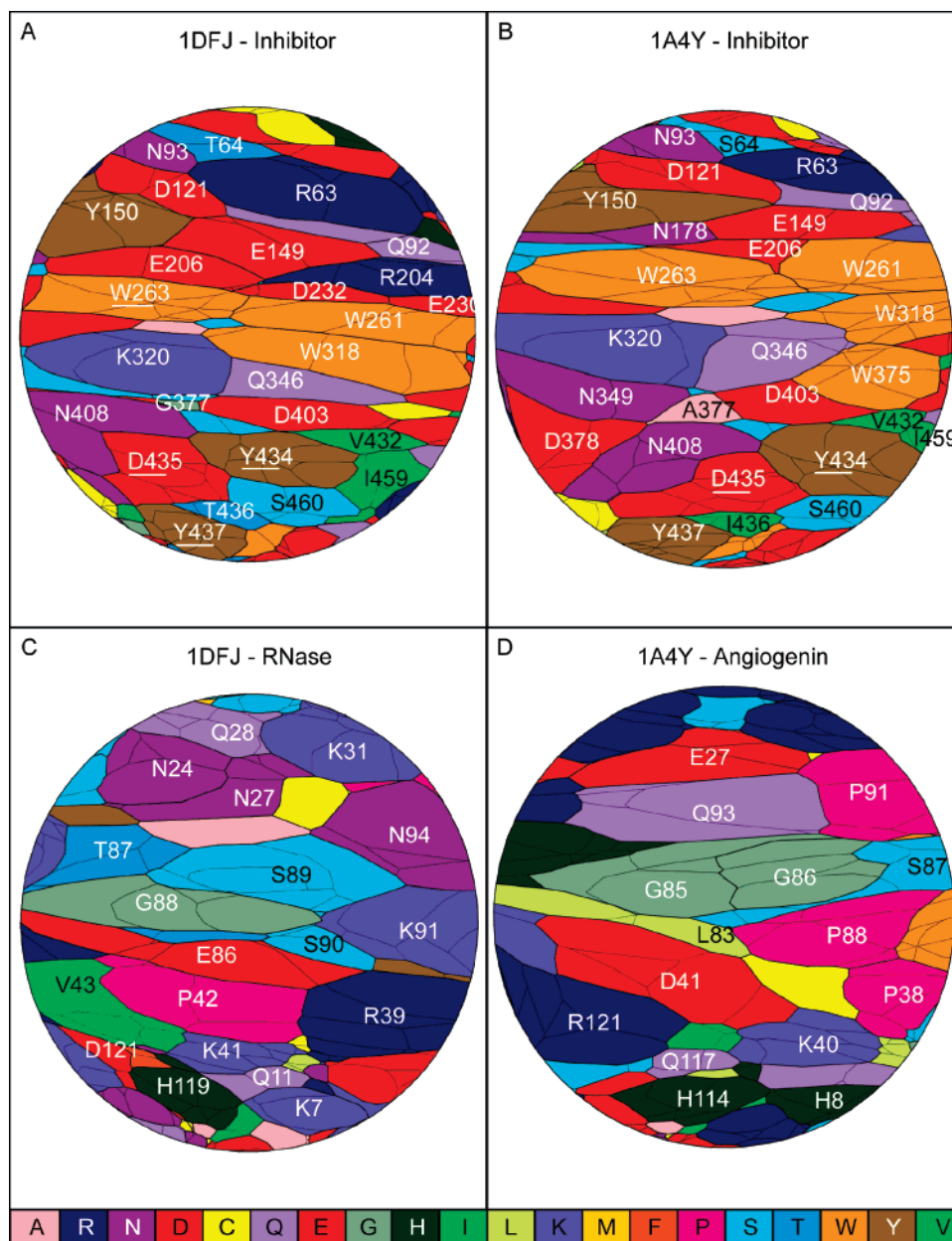


Figure 4. Comparison between the flattened views of 1DFJ and 1A4Y: (A) 1DFJ inhibitor chain, (B) 1A4Y inhibitor chain, (C) 1DFJ RNase chain, (D) 1A4Y angiogenin chain. Residues of interest are labeled in all four pictures, and the hotspot residues are underlined for the two inhibitors.

468 changing the rank order for amino acid prevalences at the
 469 interface, our classification serves as a noise filter, allowing the
 470 five most frequent side chains (Glu, Ser, Asp, Lys, and Arg) to
 471 each have more than 7% representation and the four least
 472 represented side chains (Met, Cys, Trp, and His) to each have
 473 less than 3.5% representation. In the halfway retracted inter-
 474 face, there are some changes in the ranked frequencies of the
 475 larger side chains. In particular, Tyr and Leu join Ser as the
 476 three most prevalent side chains. Additionally, our reclassifica-
 477 tion decreases the number from $\#_{total} = 10\ 292$ residues at the
 478 interface to $\#_{total} = 8934$ side chains (including the backbone
 479 category), with the number of pairs dropping from $P_{total} =$
 480 $25\ 875$ to $P_{total} = 16\ 603$ pairs. Accordingly, the average number
 481 of pairs drops from $2P_{total}/\#_{total} = 5.03$ per residue to $2P_{total}/$
 482 $\#_{total} = 3.72$ per side chain. This change is fairly uniform across
 483 the complexes (Figure 2). The large drop indicates that many

residues contribute only backbone atoms to the interface, often
 forming cross-interface interactions with other backbone-only
 contributors. The change in average pairs indicates that
 residues tend to interact with backbone atoms from multiple
 other residues, while not necessarily interacting with the
 corresponding side chains. Our observation suggests that some
 interfacial residues are selected for their contribution to folding
 and/or stability of the substituent proteins rather than for their
 contribution to transient complex formation. The high preva-
 lence of backbone at the interface also implies that a complete
 energetic characterization of interfacial contacts by experi-
 mentation remain elusive given the challenges in substitutions
 of backbone atoms.

To emphasize that most of the binding energy of protein-
 protein interactions is thought to be contributed by van der
 Waals interactions, we weigh amino acid occurrences by

Table 2. Amino Acid Occurrences by Interfacial Surface Area

	standard 20 AA ^a		side chain/backbone ^b			
	frequency		frequency		area	
	full	core	full	core	full	core
ALA	4.09	4.29	3.68	4.04	1.84	1.79
ASP	7.01	5.96	7.53	6.26	4.76	3.26
ARG	6.66	6.06	7.30	6.47	8.25	6.48
ASN	5.66	5.14	6.00	5.20	4.26	3.27
CYS	2.74	3.54	2.46	3.03	0.99	1.05
GLU	7.18	5.24	7.89	5.20	5.62	3.23
GLN	4.05	4.06	4.32	4.15	4.13	3.08
GLY	7.12	6.82	0.00	0.00	0.00	0.00
HIS	3.11	3.50	3.37	3.64	3.09	3.35
ILE	3.75	4.43	3.83	4.72	3.51	4.04
LYS	6.89	5.03	7.38	4.68	6.18	5.15
LEU	5.59	6.58	5.75	6.76	4.98	5.23
MET	2.09	2.63	2.16	2.63	2.34	2.51
PHE	3.50	4.47	3.68	4.68	3.87	5.00
PRO	4.20	3.84	4.05	3.51	2.86	2.49
SER	7.78	7.53	7.76	6.87	3.76	3.15
THR	6.08	5.16	6.12	4.78	3.59	2.96
TRP	2.60	3.58	2.75	3.88	3.64	4.97
TYR	5.75	7.34	6.09	7.70	6.36	7.80
VAL	4.16	4.81	4.16	4.74	2.98	3.21
BBA	na	na	3.70	7.05	22.98	27.97

^a Standard 20 amino acid definitions are used to calculate the frequency of each side chain type or backbone for the full and halfway retracted interface surface. ^b Frequency and area contribution of each side chain or backbone (see Results) for the full interface and the core.

Table 3. Most Common Interchain Pairs, Both in Area-Weighted Ranking and with Likelihood Correction for the Full as Well as the Halfway Retracted Interface Surface

	full		core	
	%	Log Odds	%	Log Odds
1 Glu-Arg	1.72	1.34	Glu-Arg	1.22
2 Asp-Arg	1.37	1.21	Asp-Lys	1.20
3 Asp-Lys	1.36	1.16	Arg-Trp	1.15
4 Glu-Lys	1.29	1.11	Arg-Tyr	1.06
5 Arg-Tyr	1.12	1.01	Leu-Phe	0.96
6 Asn-Tyr	0.72	0.89	Lys-Tyr	0.92
7 Lys-Tyr	0.71	0.89	Asn-Tyr	0.84
8 Glu-Tyr	0.70	0.89	Ile-Trp	0.81
9 Arg-Trp	0.70	0.89	Asp-Arg	0.75
10 Arg-Asn	0.70	0.88	Gln-Tyr	0.75

interfacial surface area throughout the rest of our analysis (Table 2). Not surprisingly, the biggest gainers from this procedure are the large amino acids Arg, Tyr, Trp, and Phe, whereas the biggest losers are Ser, Asp, Thr, and Ala. Despite some reordering, however, many of the same amino acids remain similarly ranked in this area-weighted analysis.

Interchain Pair Statistics. We next mine our data set for residue pairings across the interface surface. Pairwise statistics prove to be more informative than single residue or side-chain statistics because they more directly extract information about what interactions drive protein-protein association and/or prevent their disassociation. Interchain pairs are weighted by area to emphasize significant residue contacts across the interface. In addition, we use the log odds function to determine how different the probability of a pair is from uniformly random given the probabilities of its constituents (Table 3, Supplemental Table 2a-d in Supporting information). As for the single side chains, we are primarily interested in interactions that can be readily selected by evolution or tested by site-directed mutagenesis and thus ignore pairings that involve the

backbone. Area-weighted pairwise statistics yield a number of interesting results:

Not unexpectedly, the four highest scoring area-weighted pairs for the full interface are the four salt-bridge pairs (Glu-Arg, Asp-Arg, Asp-Lys, and Glu-Lys) (Table 3). Of these, the Glu-Arg pair is by far the most common, with 1.72% compared to 1.37%, 1.36%, and 1.29% for the Asp-Arg, Asp-Lys, and Glu-Lys pairs, respectively. Although some of these pairs arise from van der Waals contacts along the uncharged part of the side chain, many of them form salt-bridges. The prevalence of salt-bridges at the interface of transient protein-protein complexes has been previously noted by Ofra and Rost³ and emphasizes the importance of charge complementarity at protein interfaces, which tend to be protected from solvent.

The next most prevalent pairs are Tyr with the side chains of Arg, Asn, Lys, and Glu (Table 3). Of these, the Arg-Tyr pair shows the most interesting configurations, often as a hydrogen bond between the hydroxyl of Tyr and one of the three nitrogens (usually N_{η1} or N_{η2}) of Arg (about 40%). Almost as often, a classical cation- π interaction is observed (about 40%).²⁷ Of these, about two-thirds orient the amino group over the center of the ring, while about one-third orient the C_δ or C_γ atoms over the ring. Asn-Tyr pairs are seen most often as a hydrogen bond between the hydroxyl of Tyr and the O_{δ1} or N_{δ2} atom of the Asn residue (about 55%). Occasionally, Asn-Tyr pairs display an orientation similar to cation- π packing despite not being positively charged (about 15%). Lys-Tyr pairs are seen most often as hydrogen bonds between the hydroxyl of the Tyr and the N_ε atom of the Lys (about 55%). About one-third of these pack Lys carbon atoms against the Tyr ring, often with a hydrogen bond between the N_ε of Lys and the carbonyl of the Tyr backbone. Glu-Tyr pairs are seen most often as hydrogen bonds between the hydroxyl of the Tyr and either O_{ε1} or O_{ε2} of the Glu residue (about 65%).

Arg is part of the ninth and tenth ranked pairs. For the Arg-Trp pair, the cation- π interaction is most prevalent (about 75%). The core region is significantly enriched in these Arg-Trp pairs, which is consistent with their increased prevalence as hotspot residues. For the tenth ranked Arg-Asn pair, hydrogen bonds between the functional groups are by far the most common mode of interaction (about 90%).

Similar results are observed for the frequency statistics not weighted by area (data not shown).

Likelihood correction based on probabilities from individual occurrences serves to highlight two types of pairs. First, pairing preferences for residues that are rarely present at the interface (i.e., Cys, Met, His) are revealed and can be identified as recurring motifs. For example, the Cys-His pair arises from the proximity of a cysteine disulfide bridge that packs against the active site His of the catalytic triad in serine protease inhibitor complexes. This motif is particularly interesting because of the diversity observed in serine protease inhibitors. Although there is high homology among the serine-proteases, the inhibitors themselves are quite different with the exception of the cysteine disulfide positioned about 3.5 Å from the His. Although this disulfide is known to contribute significantly to the stability of these protease inhibitors,²⁸ our observations suggest that there may be other roles for these highly conserved Cys pairs, such as this interaction with the His of the serine protease. It is interesting to note that substitutions of this Cys-Cys pair have been performed that yield a protein of similar stability to wild-type yet retain a similar potency as trypsin

Protein-Protein Interfaces

research articles

583 inhibitors only if the substitutions are approximately size-
584 neutral (e.g., Gly-Leu but not Leu-Val).²⁹

585 The second class of pairs emphasized by likelihood correc-
586 tion includes those that remain high on the list despite their
587 intrinsic prevalence at the interface. These unusually prevalent
588 pairs fall into two categories. Leu-Val and Ile-Leu suggest a
589 hydrophobic component to transient protein complexes, as
590 noted previously,¹ whereas the charge pairs Asp-Lys, Glu-Arg,
591 Glu-Lys, and Asp-Arg again emphasize the importance of
592 charge complementarity at protein interfaces. Similar results
593 with only a slight reordering were observed for the core residues
594 following retraction (Table 3).

595 The 10 least common pairs (without likelihood correction)
596 almost all involve Cys (data not shown), in accord with its
597 infrequency at protein interfaces. Following likelihood correc-
598 tion, the self-pairs His-His, Lys-Lys, Arg-Arg, Tyr-Tyr, and
599 Asp-Asp are especially rare, with log odd ratios ranging from
600 -1.03 to -0.38, suggesting that the high steric cost of packing
601 like charges against each other is selected against by evolution.

602 These interchain pair statistics differ significantly from those
603 reported by Glaser et al.² In their study, Cys-Cys, Pro-Trp,
604 Asp-His and Arg-Trp are the most prevalent unweighted pairs,
605 and Arg-Trp, Pro-Trp, and Cys-Cys are the most prevalent
606 pairs as weighted by volume of the contributing amino acid.
607 The basis for the discrepancy between our results herein and
608 these prior studies is twofold. First, their data set of 621
609 interfaces is dominated by 404 homo-dimers, whose interfaces
610 more resemble protein interiors.³ Second, Glaser et al. use
611 overall residue volume, not area contribution to the interface,
612 to weight their amino acids, thus biasing their results toward
613 larger residues even if they contribute only a single atom to
614 the protein interface.

615 Although not our primary focus, we also examined the
616 extremely prevalent backbone-atom-backbone-atom (BBA-
617 BBA) pairs and BBA pairing to selected amino acids in more
618 detail. As about 70% of BBA-BBA contacts are at distance
619 between 3 and 6 Å, these appear to constitute van der Waals
620 packing or hydrogen bonding. In a number of complexes, we
621 observe series of hydrogen bonds between backbone O and N
622 atoms that mimic the hydrogen-bonding pattern of β -sheets.
623 These interactions, both parallel and antiparallel, occur in 2-3
624 residue stretches per side, such as the parallel β -sheet formation
625 between residues Gly 42, Val 43, and Met 44 of actin with
626 residues Tyr 65, Val 66, and Val 67 of DNase I in the complex
627 IATN as previously noted.³⁰ Additionally, these β -sheet-like
628 motifs are seen in a number of proteinase/inhibitor complexes
629 (1BTH, 1CBW, 1FLE, 1HIA, 2KAI, 2PTC, and 3TPI).^{31,32}

630 **Intrachain Pairs or Neighbor Statistics.** Unique to our
631 definition of the interface, neighbor information on each side
632 of a given protein complex is also captured. Neighbor pairing
633 preferences reveal the composition of common patches on a
634 protein that may be responsible for initial docking or subse-
635 quent stabilization of a transient interaction. Analogous to
636 weighting by area applied to the interchain pairs, we here
637 accumulate statistics in which each intrachain pair is weighted
638 by the length of the shared boundary between the contributed
639 regions (Supplemental Table 3a,b in Supporting Information).
640 We again exclude the prevalent BBA-pairings from our tabu-
641 lated analysis (Table 4).

642 As described by Jones and Thornton,⁷ surface patches that
643 correlate with protein docking sites in hetero-complexes show
644 a propensity for hydrophobic residues, particularly Ile, Leu,
645 Met, Phe, and Val, as well as Arg and the polar aromatic

Table 4. Most Common Intrachain or Neighbor Pairs, Both in Perimeter-Weighted Ranking and with Likelihood Correction for the Full as Well as the Halfway Retracted Interface Surface

	full			core				
	%		Log Odds	%		Log Odds		
1	Asp-Arg	1.00	Cys-Cys	3.85	Ile-Leu	0.78	Cys-Cys	3.68
2	Glu-Lys	0.92	Met-Met	1.80	Trp-Tyr	0.77	Met-Met	2.05
3	Glu-Arg	0.87	Trp-Trp	1.31	Asp-Tyr	0.70	Trp-Trp	1.37
4	Asp-Lys	0.60	Gln-His	1.09	Asp-Arg	0.70	Gln-His	1.33
5	Asp-Tyr	0.59	Ile-Leu	1.07	Leu-Tyr	0.70	Pro-Pro	1.33
6	Ile-Leu	0.57	Glu-Lys	1.06	Tyr-Tyr	0.69	Thr-Thr	1.22
7	Arg-Tyr	0.55	Asp-Arg	1.05	Ser-Tyr	0.68	His-Ser	1.21
8	Arg-Lys	0.54	Ser-Ser	0.93	Leu-Phe	0.66	Leu-Leu	1.20
9	Asn-Tyr	0.52	Met-Pro	0.92	Leu-Val	0.62	Met-Pro	1.19
10	Leu-Tyr	0.52	Ala-Ile	0.90	Ile-Tyr	0.58	Met-Val	1.17

646 residues Trp, Tyr, and His. Our observation of neighbor pair
647 preferences agree with these findings and complement them
648 by identifying specific neighbor contacts as well as their
649 interaction partners across the interface (selected triplets).

650 As with the interchain pairing preferences, the four op-
651 positively charged pairs (Asp-Arg, Glu-Lys, Glu-Arg, and Asp-
652 Lys) are the most prevalent neighbor pairs (Table 4). In contrast
653 to their interactions across the interface, however, these pairs
654 do not typically form salt-bridges. Instead, they appear to form
655 small dipoles, mostly on the periphery of interfaces. These
656 dipoles do not necessarily form salt-bridges across the interface
657 as they are paired with other charged residues (about 31%),
658 polar residues (about 30%), backbone atoms (about 24%), as
659 well as with hydrophobic residues (about 15%). How these
660 dipoles facilitate transient protein interactions remains to be
661 studied in more detail using electrostatic potentials. However,
662 these peripheral dipoles are reminiscent of the concept of
663 electrostatic steering pioneered by Fersht and Schreiber³³ with
664 an added element of local directionality. Asp-Arg dipoles are
665 enriched at protein interfaces versus noninterface in a ratio of
666 about 2:1, whereas Glu-Lys, Glu-Arg, and Asp-Arg dipoles
667 are about equally common at and outside interfaces.

668 The fifth most common neighbor pair is Asp-Tyr, which has
669 a notable preference for tripling with Arg and Lys (see for
670 example the Tyr 434-Asp 435 pair with Lys 40 in 1A4Y in Figure
671 4). This generates a salt-bridge across the interface flanked by
672 a Tyr residue. This configuration is consistent with the high
673 prevalence of Arg-Tyr and Lys-Tyr pairs across the interface
674 surface (Table 3).

675 The sixth most prevalent neighbor pair is the hydrophobic
676 Ile-Leu, which is enriched to be the most common neighbor
677 pair following halfway retraction. Most of the time, these Ile-
678 Leu pairs form triplets with other hydrophobic residues (about
679 50%), but we also see triplets with polar residues (about 26%),
680 backbone atoms (about 18%), and the occasional charged
681 amino acid (about 6%). The reason this hydrophobic pair is
682 more common than any of the others remains unclear, though
683 it is important to note that Ile and Leu are often found packed
684 near each other, primarily in hydrophobic protein interiors. In
685 fact, many hydrophobic pairs are seen in the core of the
686 interface (Table 4), suggesting that a pre-existing hydrophobic
687 patch can serve as a docking site for protein interactions. As
688 for Ile-Leu, these hydrophobic pairs are not necessarily across
689 from other hydrophobic residues. For example, Leu-Tyr forms
690 triplets with other hydrophobic residues (about 37%) and with
691 polar residues (about 30%), and Trp-Tyr pairs preferably form
692 triplets with polar (about 30%), backbone (about 30%), and

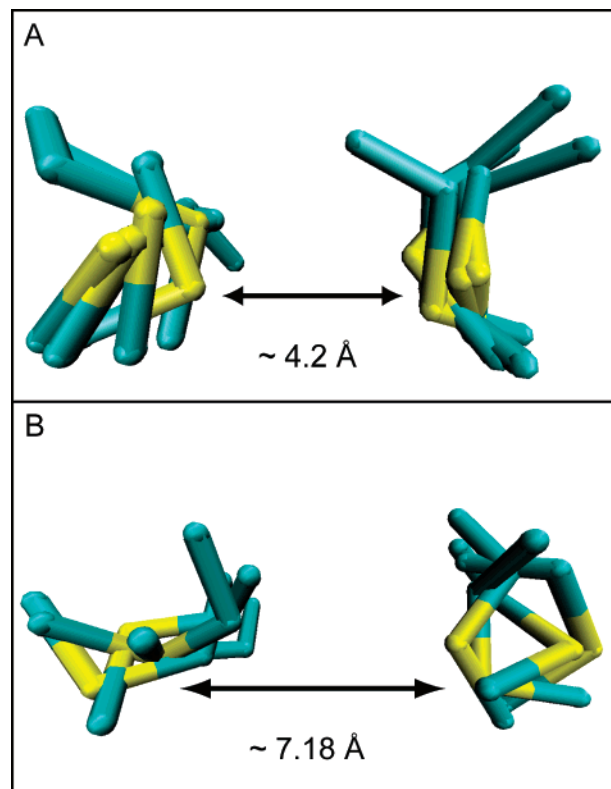


Figure 5. (A) Near Met–Met neighbor pairs from 1LDK, 1MDA, 1DN2, 1YCS, and 1AIP overlaid using rmsd alignment for side-chain atoms, excluding C_{β} . (B) Far Met–Met neighbor pairs from 1ATN, 1MDA, and 1AIP(2) overlaid using rmsd alignment for side-chain atoms, excluding C_{β} .

693 charged (about 24%) residues and relatively infrequently with
694 hydrophobic residues (about 16%).

695 Arg–Tyr and Lys–Tyr neighbor pairs are also often seen at
696 the interface (7th and 12th most frequently, respectively), which
697 is consistent with their propensity to form the cation– π motif.
698 Collectively, 7.6% of the observed Arg–Tyr and Lys–Tyr
699 neighbor pairs occur between consecutive residues in a protein
700 chain, which is consistent with the previously noted 7.3% for
701 all cation– π motifs, suggesting that many occur on α -helices.²⁷

702 Likelihood correction highlights a number of interesting
703 neighbor pairs (Table 4, Supplemental Table 3c,d in Supporting
704 Information). The Cys–Cys disulfide of trypsin protease inhibitors
705 noted above is detected again, which not surprisingly forms
706 triplets with His, as described above. Met–Met neighbor pairs
707 are the second most prevalent following likelihood correction.
708 We observe two similar yet different motifs, which we refer to
709 as *near* and *far* (Figure 5). The near Met–Met motif contains
710 S–S distances ranging from 3.5 to 4.7 Å, with an average \pm
711 standard deviation of 4.2 ± 0.4 Å, and contains only noncon-
712 secutive pairs (i.e., not consecutive along the protein chain).
713 The far Met–Met motif contains S–S distances ranging from
714 5.54 to 8.81 Å, with an average \pm standard deviation of $7.18 \pm$
715 1.2 Å, and contains both consecutive Met–Met residue pairs
716 and nonconsecutive pairs. Both Met–Met motifs pack against
717 hydrophobic residues or hydrophobic regions of charged
718 residues, moving close to a C_{β} and C_{γ} atom in all but one
719 observed case.

720 The Gln–His pair, which is the fourth most prevalent
721 neighbor pair in the perimeter normalized statistics, is seen
722 most frequently forming hydrogen bonds, either as Gln $O_{\epsilon 1}$ /

$N_{\epsilon 2}$ with His $N_{\delta 1}/N_{\epsilon 2}$ (about 65%) or Gln $O_{\epsilon 1}/N_{\epsilon 2}$ with His
backbone O/N (about 20%). Spurious contacts (about 15%)
include Gln–His neighbors that occur consecutively on the
protein chain, where steric requirements preclude H-bond
formation.

Conclusions

728 We have presented here a statistical analysis of protein–
729 protein interfaces from a large and diverse data set using a
730 reliable and consistent definition of the interface. Our analysis
731 serves as a foundation for prediction problems in protein–
732 protein docking. For example, our pairing and neighbor prefer-
733 ences can be used as weights in scoring functions to distinguish
734 between true and false predictions. Previously generated lists
735 of 2000–10 000 possible docking configurations containing one
736 or more correct answers³⁴ and data sets of native and decoy
737 docking configurations^{4,35} will serve as useful test sets for such
738 implementations. Additionally, residue frequencies and neigh-
739 bor preferences can be used to predict probable binding sites
740 for proteins whose 3D coordinates are available but whose
741 interaction sites remain unclassified. Identification of these
742 binding sites will allow the potential identification of novel
743 protein–protein pairs, leading to a greater understanding of
744 the networks of interactions in the proteome. We have also
745 provided a novel visualization that facilitates the analysis of
746 the intrinsic complexity of protein–protein interfaces. Our
747 simplified view allows easier recognition of interfacial residue
748 contacts and other biochemical characteristics. Insights derived
749 from such visual inspections will aid in the design of experi-
750 ments toward elucidating the specificity of protein–protein
751 association. Also, comparative studies of related interfaces are
752 made easier by having a single independent and simplified
753 entity. Combined, our statistical analysis and visualization serve
754 as a novel toolset for biochemists interested in the fundamen-
755 tals of protein–protein interactions. 756

Supporting Information Available: Tables listing the
area distortion measures for a sampling of complexes in the
data set; the complete area-weighted statistics and complete
likelihood corrected statistics for interchain pairs for the full
interface, complete area-weighted statistics and complete
likelihood corrected statistics for interchain pairs for the core
of the interface; complete perimeter-weighted statistics and
complete likelihood corrected perimeter-weighted statistics for
the intrachain neighbors for the full interface and complete
perimeter-weighted statistics and complete likelihood corrected
perimeter-weighted statistics for the intrachain neighbors for
the core of the interface; and figure of the comparison between
Uniform and Mean Value Coordinate flattening methods for
1BRS. This material is available free of charge via the Internet
at <http://pubs.acs.org>.

References

- 773 (1) Chakrabarti, P.; Janin, J. Dissecting protein-protein recognition
774 sites. *Proteins* **2002**, *47*, 334–343.
- 775 (2) Glaser, F.; Steinberg, D. M.; Vakser, I. A.; Ben-Tal, N. Residue
776 frequencies and pairing preferences at protein-protein interfaces.
777 *Proteins* **2001**, *43*, 89–102.
- 778 (3) Ofra, Y.; Rost, B. Analyzing six types of protein-protein interfaces.
779 *J. Mol. Biol.* **2003**, *325*, 377–387.
- 780 (4) Chen, R.; Mintseris, J.; Janin, J.; Weng, Z. A protein-protein
781 docking benchmark. *Proteins* **2003**, *52*, 88–91.
- 782 (5) Janin, J.; Chothia, C. The structure of protein-protein recognition
783 sites. *J. Biol. Chem.* **1990**, *265*, 16027–16030.

- 784 (6) Hubbard, S. J.; Argos, P. Cavities and packing at protein interfaces. *Protein Sci.* **1994**, *3*, 2194–2206. 830
- 785 (7) Jones, S.; Thornton, J. M. Analysis of protein-protein interaction 831
- 786 sites using surface patches. *J. Mol. Biol.* **1997**, *272*, 121–132. 832
- 787 (8) Lo Conte, L.; Chothia, C.; Janin, J. The atomic structure of protein- 833
- 788 protein recognition sites. *J. Mol. Biol.* **1999**, *285*, 2177–2198. 834
- 789 (9) Bahadur, R. P.; Chakrabarti, P.; Rodier, F.; Janin, J. A dissection 835
- 790 of specific and non-specific protein-protein interfaces. *J. Mol.* 836
- 791 *Biol.* **2004**, *336*, 943–955. 837
- 792 (10) Nooren, I. M. A.; Thornton, J. M. Structural characterization and 838
- 793 functional significance of transient protein-protein interactions. 839
- 794 *J. Mol. Biol.* **2003**, *325*, 981–1018. 840
- 795 (11) Jones, S.; Thornton, J. M. Principles of protein-protein interac- 841
- 796 tions. *Proc. Natl. Acad. Sci. U.S.A.* **1996**, *93*, 13–20. 842
- 797 (12) Tsai, C.-J.; Lin, S. L.; Wolfson, H. J.; Nussinov, R. Studies of 843
- 798 protein-protein interfaces: a statistical analysis of the hydro- 844
- 799 phobic effect. *Protein Sci.* **1997**, *6*, 53–64. 845
- 800 (13) Xu, D.; Tsai, C. J.; Nussinov, R. Hydrogen bonds and salt bridges 846
- 801 across protein-protein interfaces. *Protein Eng.* **1997**, *10*, 999– 847
- 802 1012. 848
- 803 (14) Rodier, F.; Bahadur, R. P.; Chakrabarti, P.; Janin, J. Hydration of 849
- 804 protein-protein interfaces. *Proteins* **2005**, *60*, 36–45. 850
- 805 (15) Ma, B.; Elkayam, T. W. H.; Nussinov, R. Protein-protein interac- 851
- 806 tions: structurally conserved residues distinguish between bind- 852
- 807 ing sites and exposed protein surfaces. *Proc. Natl. Acad. Sci. U.S.A.* 853
- 808 **2003**, *100*, 5772–5777. 854
- 809 (16) Duncan, B. S.; Olson, A. J. Texture mapping parametric molecular 855
- 810 surfaces. *Mol. Graphics* **1995**, *13*, 258–264. 856
- 811 (17) Varshney, A. F. P.; Brooks, J.; Richardson, D. C.; Wright, W. V.; 857
- 812 Manocha, D. Defining, computing, and visualizing molecular 858
- 813 interfaces. *Proc. IEEE Visualization* **1995**, 36–43. 859
- 814 (18) Word, J. M.; Lovell, S. C.; LaBean, T. H.; Taylor, H. C.; Zalias, M. 860
- 815 E.; Presley, B. K.; Richardson, J. S.; Richardson, D. C. Visualizing 861
- 816 and quantifying molecular goodness-of-fit: small-probe contact 862
- 817 dots with explicit hydrogen atoms. *J. Mol. Biol.* **1999**, *285*, 1711– 863
- 818 1733. 864
- 819 (19) Neuvirth, H.; Raz, R.; Schreiber, G. ProMate: A structure based 865
- 820 prediction program to identify the location of protein-protein 866
- 821 binding sites. *J. Mol. Biol.* **2004**, *338*, 181–199. 867
- 822 (20) Mintz, S.; Shulman-Peleg, A.; Wolfson, H. J.; Nussinov, R. Genera- 868
- 823 tion and analysis of a protein-protein Interface data set with 869
- 824 similar chemical and spatial patterns of interactions. *Proteins* 870
- 825 **2005**, *61*, 6–20. 871
- 826 (21) Ban, Y. E.; Edelsbrunner, H.; Rudolph, J. Interface surfaces for 872
- 827 protein-protein complexes. *J. Assoc. Comput. Mach.* **2006**, *53* (3), 873
- 828 361–378. 874
- 829 (22) Gottschalk, K.-E.; Neuvirth, H.; Schreiber, G. A novel method for 875
- scoring of docked protein complexes using predicted protein- 876
- protein binding sites. *Protein Eng. Des. Sel.* **2004**, *17* (2), 183– 877
189. 878
- (23) Tutte, W. How to draw a graph. *Proc. London Math. Soc.* **1963**, 879
- 743–768. 880
- (24) Floater, M. Mean value coordinates. *Comput. Aided Geometric* 881
- Des.* **2003**, *20*, 19–27. 882
- (25) Nicholls, A. J. GRASP: graphical representation and analysis of 883
- surface properties. *Biophys. J.* **1993**, *64*, A116. 884
- (26) Gabdoulline, R. R.; Wade, R. C. Analytically defined surfaces to 885
- analyze molecular interaction properties. *J. Mol. Graphics* **1996**, 886
- 14*, 341–353. 887
- (27) Gallivan, J. P.; Dougherty, D. A. Cation-pi interactions in structural 888
- biology. *Proc. Natl. Acad. Sci. U.S.A.* **1999**, *96*, 9459–9464. 889
- (28) Schwarz, H.; Hinz, H.; Mehlich, A.; Tschesche, H.; Wenzel, H. 890
- Stability studies of derivatives of the bovine pancreatic trypsin 891
- inhibitor. *Biochemistry* **1987**, *26* (12), 3544–3551. 892
- (29) Hagihara, Y.; Shiraki, K.; Nakamura, T.; Uegaki, K.; Takagi, M.; 893
- Imanaka, T.; Yumoto, N. Screening for stable mutants with amino 894
- acid pairs substituted for the disulfide between residues 14 and 895
- 38 of bovine pancreatic trypsin inhibitor (BPTI). *J. Biol. Chem.* 896
- 2002**, *277* (52), 51043–51048. 897
- (30) Kabsch, W.; Mannherz, H. G.; Suck, D.; Pai, E. F.; Holmes, K. C. 898
- Atomic structure of the actin:DNase I complex. *Nature* **1990**, *347*, 899
- 37–44. 900
- (31) Bode, W.; Huber, R. Natural protein proteinase inhibitors and 901
- their interaction with proteases. *Eur. J. Biochem.* **1992**, *204* (2), 902
- 433–451. 903
- (32) van de Locht, A.; Bode, W.; Huber, R.; Le Bonniec, B. F.; Stone, 904
- S. R.; Esmon, C. T.; Stubbs, M. T. The thrombin E192Q-BPTI 905
- complex reveals gross structural rearrangements: implications 906
- for the interaction with antithrombin and thrombomodulin. 907
- EMBO J.* **1997**, *16* (11), 2977–2984. 908
- (33) Schreiber, G.; Fersht, A. R. Rapid, electrostatically assisted as- 909
- sociation of proteins. *Nat. Struct. Biol.* **1996**, *3* (5), 427–431. 910
- (34) Wang, Y.; Agarwal, P. K.; Brown, P.; Edelsbrunner, H.; Rudolph, 911
- J. Coarse and reliable geometric alignment for protein docking. 912
- Proc. Pac. Symp. Biocomput.* **2005**, 64–75. 913
- (35) Gray, J. J.; Moughon, S.; Wang, C.; Schueler-Furman, O.; Kuhlman, 914
- B.; Rohl, C. A.; Baker, D. Protein-protein docking with simulta- 915
- neous optimization of Rigid-body Displacement and Side-chain 916
- conformations. *J. Mol. Biol.* **2003**, *331*, 281–299. 917
- PR070018+ 873