

Shelling the Voronoi interface of protein-protein complexes predicts residue activity and conservation

Benjamin Bouvier^{a,c,*} Raik Grünberg^{b,*} Michael Nilges^c

Frederic Cazals^{a,**}

^a*INRIA Sophia-Antipolis, Project Geometrica, F-06902 Sophia-Antipolis, France*

^b*EMBL-CRG Systems Biology Unit, CRG-Centre de Regulacio Genomica, Dr.*

Aiguader 88, 08003 Barcelona, Spain

^c*Unit de Bioinformatique Structurale, Institute Pasteur, 75724 Paris Cedex 15,*

France

Abstract

The accurate description of protein-protein interfaces remains a challenging task. Traditional criteria, based on atomic contacts or changes in solvent accessibility, tend to over or underpredict the interface itself and cannot discriminate active from less relevant parts. A recent molecular dynamics simulation study by Mihalek and co-authors concluded that active residues tend to be ‘dry’, that is, insulated from water fluctuations. We show that patterns of ‘dry’ residues can, to a large extent, be predicted by a fast, parameter-free and purely geometric analysis of protein interfaces. We introduce the shelling order of Voronoi facets as a straightforward quantitative measure of an atom’s depth inside an interface. We analyze the correlation between Voronoi shelling order, dryness, and conservation on a set of 54 protein-protein complexes. Residues with high shelling order tend to be dry; evolutionary conservation also correlates with dryness and shelling order but, perhaps not surprisingly, is a much less accurate predictor of either property. Voronoi shelling order thus seems a meaningful and efficient descriptor of protein interfaces. Moreover, the strong correlation with dryness suggests that water dynamics within protein interfaces may, in first approximation, be described by simple diffusion models.

Key words:

Protein-protein complex, interface activity, hotspots, conservation, Voronoi models.

PACS:

* both authors contributed equally

**Corresponding author

Email address: `Frederic.Cazals@sophia.inria.fr` (Frederic Cazals).

1 INTRODUCTION

2 Specific recognition between proteins plays a crucial role in almost all cellular
3 processes and most proteins are embedded in highly connected (and dynami-
4 cally changing) networks of interaction partners [1]. Despite much progress [2],
5 identifying the exact interface between two proteins remains difficult. On the
6 one hand, exact predictions are hindered by the complex and dynamic nature
7 of proteins [3,4]; on the other hand, the descriptors we employ to study the
8 interface may be flawed or ill-chosen.

9 A protein-protein interface is traditionally defined by the ‘geometric footprint’,
10 which refers to all atoms within a given distance of the interaction partner.
11 Somewhat more precise definitions rely on the loss of solvent accessibility (SA)
12 upon binding [5]. Yet, as much as half of this footprint can seemingly be irrel-
13 evant to binding [6]. As contributions to specificity and affinity appeared very
14 unevenly distributed, substantial effort has been spent on the identification
15 of areas or residue patches that are actively involved in molecular recognition
16 [7–10]. This lead to the definition of ‘hotspot’ residues [11,12]. Hotspots refer
17 to the usually very small number [12] of ‘key’ residues in a protein-protein
18 interface, the mutation of which causes large changes in the binding free en-
19 ergy. Contrary to this focus on isolated residues, more recent studies have
20 revealed strong non-additive, collective effects [13] which point to a modular
21 organization of interfaces into interaction clusters [14].

22 Also the evolutionary record seems of limited use for distinguishing relevant
23 from irrelevant. The sequence conservation of protein-protein interfaces is
24 hardly statistically significant and depends heavily on surface-patch selection

25 techniques [15]. A commonly adopted view states that, unlike catalytic sites
26 that are highly unlikely to transform in a series of discrete steps without com-
27 plete loss of activity [16], the assembly of proteins involves a continuous scale
28 of binding modes, from transient to stable, leaving more freedom for evolution
29 to proceed in incremental steps [17–19]. Interestingly, conservation signals be-
30 come more convincing if one turns away from individual– and towards patches
31 [20] or clusters of residues [21].

32 Water forms an essential part of protein-protein interfaces [9,22]. The occlusion
33 of bulk solvent is a common denominator not only of classical hotspots [23],
34 but also of the more recently identified interaction modules [14], which are
35 delimited by structural water. In fact, the removal of water from partially
36 solvated backbone hydrogen bonds has been argued to be a driving force of
37 binding [24,25].

38 Recently, Mihalek and coworkers [26] went one step further and classified inter-
39 face residues by the dynamics of surrounding water molecules. They asserted
40 that the important residues are the ones whose interactions are not disturbed
41 by water fluxes. These ‘dry’ residues (some of which may actually be in con-
42 tact with immobile, structural water molecules) were found to correlate better
43 with conservation than the overall geometric footprint and to feature some
44 characteristic properties of classical hotspots. The dryness results collated by
45 these authors on a variety of systems thus represent valuable information as
46 a measure of residue importance; we will constantly refer to them during this
47 work.

48 However, the method suffers from some drawbacks. It relies on molecular dy-
49 namics simulations which are computationally expensive and sensitive to setup

50 and parameterization. Furthermore, it cannot itself distinguish between inter-
51 face and noninterface residues. Mihalek and coworkers addressed this problem
52 by discarding residues that are also dry in the isolated partners, hereby further
53 increasing computational costs and neglecting the possibility of conformational
54 transitions upon binding.

55 All in all, the combination of the large size of protein-protein interfaces, the
56 relatively small areas that appear actually important and the lack of unam-
57 biguous ways to identify them, amounts to a difficult problem for which novel
58 approaches are highly desirable. We present a method based on the shelling
59 of the Voronoi interface of protein-protein complexes. The method quantifies
60 the depth of any given atom inside the interface, in a manner accounting for
61 both the geometry and the topology of the interface. The method is simulta-
62 neously accurate, computationally inexpensive, and elegant in that it does not
63 require parameterization. Voronoi shelling order features an excellent correla-
64 tion with the water shielding observed by Mihalek et al., without the need for
65 simulations or geometric footprinting. We analyze the relationship between
66 three quantities of interest (Voronoi shelling order, dryness and conservation)
67 on the same set of protein complexes. We illustrate the advantages as well as
68 potential improvements of the geometric measure with detailed examples and
69 elaborate on the more complex correlation with evolutionary information.

70 2 THEORY

71 2.1 Voronoi description of protein-protein interfaces

72 In this section, we briefly summarize the Voronoi model of protein-protein
73 interfaces, which is described in more detail in [27], together with a compre-
74 hensive bibliography. Given a collection of sample points equipped with the
75 Euclidean distance, the Voronoi diagram is the space partition which assigns
76 to every sample the convex polyhedron containing all points in space closer to
77 it than to any other sample. In 3D space, these Voronoi regions are bounded
78 by Voronoi facets (resp. edges, vertices) which consist of points equidistant
79 from two (resp. three, four) samples.

80 The Euclidean Voronoi diagram of atom centers in a molecule, first employed
81 by Richards [28] to investigate packing properties in proteins, is unable to
82 account for the fact that different atoms have different radii. A convenient
83 generalization thereof, which overcomes this limitation while retaining non-
84 curved bisectors, is the power diagram[29]. It replaces the Euclidean distance
85 with the ‘power distance’ of a point to a sphere centered at \mathbf{a} and of radius
86 r : $p(\mathbf{x}) = |\mathbf{a} - \mathbf{x}|^2 - r^2$. The power diagram is an extension of the Voronoi
87 diagram (to which it reverts for atoms of equal radii); hence, we continue to
88 refer to it as such in the text. Throughout the study, we compute it for atomic
89 spheres whose radii are the so-called group radii [30], expanded by the radius
90 of a probe water molecule $r_w = 1.4$ Å. This effectively models the solvent-
91 accessible surface (SAS) of the protein, as defined by Lee and Richards [31].
92 An example Voronoi diagram for a hypothetical two-dimensional molecule is
93 shown on Figure 1.

94 The Voronoi diagram has a dual (an associated and strictly equivalent struc-
95 ture) called the Delaunay triangulation; in practice, Voronoi diagrams are
96 calculated via their Delaunay triangulation rather than directly. The Delau-
97 nay triangulation consists of edges (resp. triangles, tetrahedra) that connect
98 the centers of two (resp. three, four) adjacent spheres whose corresponding
99 Voronoi regions share a facet (resp. an edge, a vertex).

100 When modeling molecules, a drawback of the Voronoi diagram is that atoms
101 located on the convex hull have unbounded Voronoi regions (all but the region
102 of atom a_2 , on Figure 1). An elegant way of solving this problem is to use
103 a restriction of the Delaunay triangulation called the α -complex [32]. For a
104 fixed value of α , each ball of center \mathbf{a}_i and radius r_i is replaced by a ball of
105 center \mathbf{a}_i and radius $\sqrt{r_i^2 + \alpha}$. Given these expanded balls, the construction
106 of the α -complex mimics that of the Delaunay triangulation, to the extent
107 that one focuses on the intersection of the restriction of each expanded ball
108 to its Voronoi region rather than the Voronoi region itself; see Figure 1 for an
109 illustration. Varying the value of α allows for the investigation of properties
110 at different scales. In particular, for very large values of α the α -complex
111 is identical to the Delaunay triangulation. In rare occurrences of desolvated
112 models, an additional filtering step may be necessary to discard all instances
113 of unphysically large facets at the rim of the interface [27]; we do not discuss
114 this issue further since this study involves solvated models only.

115 We now apply this methodology to model the interface between two proteins
116 A and B . Following [27], the AB interface consists of the Delaunay edges
117 found in the 0-complex – the α -complex for $\alpha = 0$, and whose endpoints
118 belong to A and B . Because of the duality between the Delaunay and Voronoi
119 representations, the interface can also be described using the Voronoi facets

120 dual to the aforementioned edges. The interface model can be extended to
121 accommodate interface water molecules W , defined as sharing at least one
122 edge with each partner in the 0-complex. This allows for the definition of
123 the following interfaces: AB between the protein partners; AW (resp. BW)
124 between partner A (resp. B) and interface water; $AW-BW$ as the union of the
125 interfaces AW and BW ; ABW as the union of the interfaces AB and $AW-BW$.
126 Like other methods mentioned above, our model correctly identifies any
127 atom losing solvent accessibility as an interface atom. Unlike these methods
128 however, it also detects interface atoms that do not lose solvent accessibility
129 – essentially buried backbone atoms, these represent a non-negligible 13% of
130 the interface [27].

131 2.2 Shelling the ABW interface

132 The next step of the algorithm attributes a Voronoi shelling order (VSO)
133 to each facet of the ABW interface. This represents the number of ‘jumps’
134 between adjacent facets that needs to be performed, from the currently con-
135 sidered location, to reach the rim of the interface (Figures 2a and 3a). The
136 Voronoi interface is thus partitioned into concentric shells of increasing selling
137 order.

138 The calculation of VSO values for all interface facets requires two passes.
139 During the first pass, boundary Voronoi facets located at the rim of the in-
140 terface are enumerated and given a VSO of one. Voronoi facets are bounded
141 by Voronoi edges, each of which is incident to exactly three Voronoi facets
142 in the Voronoi diagram; however, some of these facets may not belong to the
143 interface (their dual Delaunay edges are not in the 0-complex). This allows us

144 to detect rim Voronoi facets as the ones featuring at least one Voronoi edge
145 that is incident to one interface Voronoi facet only. The second pass explores
146 the interface breadth-first starting from the previously identified rim facets.
147 Given an interface Delaunay edge (of shelling order n), the algorithm checks
148 all incident Delaunay triangles, as each such triangle contributes zero, one or
149 two additional interface edges. If these have not already been shelled, they are
150 given a VSO of $n + 1$. To speed up the search operations, a temporary map
151 storing edges of VSO $n - 1$, n and $n + 1$ is used, since these are the only ones
152 that can be encountered at level n ; the contents of this map are copied over
153 to a permanent structure each time n increases.

154 The outcome of this process is the association of an integer VSO value to
155 each Delaunay edge (or equivalently, Voronoi facet) of the *ABW* interface.
156 However, our ultimate goal is to quantify the depth of any given atom in-
157 side the interface. This is done by tagging the atom with the minimum value
158 among the shelling orders of the Delaunay edges to which the atom contributes
159 (Figures 2b and 3b). The maximum or average values have also been consid-
160 ered as candidates, but their variation throughout the interface were found to
161 closely mimic that of the minimum. Finally, the shelling order of a residue,
162 defined as the average VSO value over its constituent atoms contributing to
163 the Voronoi interface, is employed when comparing to residue-based measures
164 such as conservation or dryness.

165 3 RESULTS

166 3.1 Voronoi shelling order, conservation and water dynamics

167 A recent simulation study examined the rate at which residues in protein-
168 protein interfaces exchange surrounding water molecules [26]. Residues that
169 were mostly shielded from mobile water molecules, defined as “dry” by Mi-
170 halek et al., turned out to be more conserved and were thus interpreted as
171 the active part of the interface. Our initial goal is to assess how well shelling
172 order is able to predict dryness on the set of homo- and heterodimer complexes
173 studied by Mihalek et al. [26]. As a yardstick, we compare to the previously
174 established correlation between conservation and dryness. Conservation is de-
175 termined from pFam [33] hidden Markov models [34] using a relative entropy
176 scheme [35]. In order to characterize all possible relationships, we also examine,
177 further down in the text, how good a predictor of shelling order conservation is.
178 We generate three ROC plots for each complex, describing the performance of
179 shelling order as predictor of dryness, of conservation as predictor of dryness
180 and of conservation as predictor of shelling order, respectively. A represen-
181 tative example set of ROC curves is shown in Figure 4. The area between
182 each ROC curve and the diagonal quantifies the predictive power of a score
183 (i.e. VSO, conservation) in terms of sensitivity and specificity. An area of 0.5
184 corresponds to a perfect prediction, which in the example of shelling order
185 predicting dryness means that the n dry residues in the interface perfectly
186 match the n residues with highest shelling order without any over-prediction.
187 By contrast, a ROC area of 0 corresponds to the performance of a pure random
188 classifier. See Section 5.4 for details.

189 The results are compiled in Tables 1 and 2 for heterodimers and homodimers,
190 respectively, and summarized in Figure 5. Evidently, Voronoi shelling order
191 is a very good predictor of dryness and outperforms conservation for 35 of
192 the 36 homodimers and 17 of the 18 heterodimers. VSO always performs
193 better than a purely random classifier, whereas conservation fails to do so
194 in seven cases (five homodimers and two heterodimers). The third columns of
195 Tables 1 and 2 quantify the ability of sequence conservation to predict Voronoi
196 shelling order. We define the n_{core} residues with highest VSO as ‘core’ and
197 the remainder as ‘rim’ and test the ability of conservation to discriminate
198 between the two. We adjust n_{core} for each complex so as to exactly match
199 the number of residues classified as dry. We thus tie ourselves to a threshold
200 chosen by Mihalek et al. [26] rather than optimizing our own. Nevertheless,
201 the connection from conservation to Voronoi shelling order appears as good
202 as it is to dryness. While the results differ in detail, the average ROC area
203 is 0.15 for heterodimers and 0.12 for homodimers, which compares well with
204 the respective figures of 0.14 and 0.13 for the prediction of water shielding.
205 However, both conservation-based predictions are outperformed by the much
206 closer correlation between shelling order and dryness, reflected by average
207 ROC areas of 0.31 and 0.34. This notable discrepancy indicates a more direct
208 link between the two latter properties, both of which are structure-based.

209 3.2 *Spatial distribution of conserved residues*

210 The analysis of the ROC curves provides insight into the location of highly
211 conserved residues across the interface shells: conservation becomes a mediocre
212 predictor for Voronoi shelling order when highly conserved residues are found

213 at low VSO (such residues are expected to be wet) and/or when poorly con-
214 served residues are found at high VSO (such residues are expected to be dry).
215 However, this simplified focus on extreme values can not fully capture the
216 spatial distribution of conservation. We therefore now address two comple-
217 mentary points, namely (i) the average residue conservation as a function of
218 VSO, and (ii) the cumulated conservation score over consecutive shells.

219 (i) Guharoy and Chakrabarti showed that residues at the interface core are,
220 on average, more conserved than those on the rim [36]. Their binary interface
221 model defined the rim as all residues that are not fully buried inside the com-
222 plex. Our more quantitative description helps to refine the prior conclusion.
223 We normalize conservation scores and Voronoi shelling order so that both span
224 the range 0 to 1 for each interface. We then compute the average conserva-
225 tion score as a function of VSO using a large moving window comprising 1/4
226 of all interface residues. Figures 6 and 7 show this running average for all
227 complexes. The relation between residue conservation on the one hand, and
228 depth within the interface on the other, is evidently not a simple one. The
229 non-averaged original values (gray lines) highlight the scattering of conserva-
230 tion across shells: highly conserved residues are found even at the very rim.
231 Only the extensive averaging reveals a clear correlation between increases in
232 shelling order and residue conservation. This observation is not sensitive to
233 the actual averaging window and the curves remain very similar for window
234 sizes between 1/8 and 1/2 of the interface (data not shown).

235 The overall correlation between shelling order and conservation can be quan-
236 tified in a single number by double integration over the running average. We
237 denote $c(x)$ the average conservation score at $VSO = x$ and reset the baseline

238 of this function to 0 by subtracting the minimum value m : $\bar{c}(x) = c(x) - m$.
239 We now define $A = \int_0^1 \bar{c}(t)dt$ to be the area under this running average and we
240 normalize $\bar{c}(x)$ to cover an area of 1: $f(x) = (c(x)-m)/A$. Function $f(x)$ can be
241 seen as a probability density function, with associated cumulated distribution
242 function $F(x) = \int_0^x f(t)dt$ (dash-dotted line in figures 6 and 7). One always has
243 $F(1) = 1$, but the speed at which F reaches 1 depends on whether conserved
244 residues are picked up early (in the outer shells) or late (inner shells). F thus
245 encodes the cumulative conservation score up to shelling order x . To provide
246 a concise measure of this property, we report $g(x) = \int_0^x F(t)dt$ (dotted line in
247 figures 6 and 7). The total area under F depends on the overall distribution of
248 conservation across shells. Lower values of $g(1)$ thus indicate that conserved
249 residues tend to cluster towards the *core* of the interface; values above 0.5 (the
250 double integral over a flat line) denote clustering near the *rim*. The deviation
251 $\Delta = g(1) - 0.5$ is reported in the lower right corner of each plot in figures
252 6 and 7. $g(1)$ falls below a value of 0.5 for 15 out of 18 heterodimers and
253 28 out of 36 homodimers. Conservation thus generally increases towards the
254 interface core. Nevertheless, apart from the few obvious exceptions, closer in-
255 spection also reveals some interesting systematic deviations: (i) Conservation
256 density often reaches its maximum before the innermost shell – the interface
257 center thus appears under less constraint than a surrounding outer core; (ii)
258 contrary to the overall trend, a pronounced secondary peak of conservation is
259 sometimes apparent at the very edge of the interface.

260 (ii) While the previous analysis focuses on the spatial distribution of conser-
261 vation per se, it is also worthwhile to compare the spatial distribution of con-
262 servation for two sets of residues: the interface residues and the dry residues.
263 The detailed analysis is described in section A.1 of the supplemental material.

264 Non-interface residues account for a proportion of the total conservation score
265 (over the whole protein) in the range 60% to 84% in heterodimers (average
266 76%), and 36% to 97% for homodimers (average 73%) –see the second column
267 of Tables A.1 and A.2 in the supplemental material. These results alone show
268 that the effect of the majority of conserved residues on the interface is at best
269 an indirect one –for example, through the imposition of a protein fold which in
270 turn dictates interface structure. Moreover, the comparison of the area under
271 the cumulated distribution function for interfacial and dry residues performed
272 in Section A.1 confirms that the rim amino-acids account for a non-negligible
273 part of the conservation. The good agreement with the scattered conservation
274 signals and conserved interface rims observed in figures 6 and 7 allows us to
275 rule out a purely statistical effect where a large number of moderately con-
276 served rim residues might end up having more weight than a small number
277 of highly conserved core amino-acids: highly conserved residues do occur on a
278 non fortuitous basis at the rim of protein-protein interfaces.

279 The in-depth examination of average and cumulated conservation thus con-
280 firms the general trend of higher conservation towards core shells but also
281 hints at a more complex fine structure. The very center of an interface often
282 appears more amenable to change than its immediate surroundings; further-
283 more, numerous interfaces seem to bear substantial evolutionary pressure on
284 their outer rims. From the inspection of examples, we speculate this latter
285 signal to be a signature of electrostatic steering [37] but the issue deserves
286 further scrutiny.

287 *3.3 Case-studies: best and worst case scenarios for shelling order*

288 To identify in more detail the incentives and shortcomings of using shelling
289 order for the description of interfaces and as a predictor of water dynamics, we
290 focus on three extreme cases of application, which are presented in Figure 8.

291 *The ideal case.* The interface of the homodimer complex 1E2D (left) features a
292 compact and planar core composed of a single patch of atoms with high shelling
293 orders (large panel), which the MD simulations of Mihalek and coworkers also
294 identify as dry (lower left-hand panel). Such compact interfaces with disk-like
295 topologies and no holes represent best case scenarios for the predictive power
296 of our model. Also conservation performs well for this complex. However, in
297 contrast to shelling order, the conservation score delimitates a patch which
298 extends far beyond the dry residues, resulting in a good sensitivity but a
299 poor selectivity. In fact, the most highly conserved residues are catalytic in
300 nature, and located at the entrance of a finger-like cavity which extends, from
301 the other side of the protein, in the direction of the interface (not visible
302 in the figure). The co-crystallized thymidine monophosphate and adenosine
303 diphosphate substrates [38] allowed Mihalek and coworkers to identify these
304 residues as catalytic and as such to exclude them from their analysis. However,
305 the detection of catalytic residues is not always as straightforward and the
306 influence of this and a variety of other factors hamper the use of conservation
307 measures for specific predictions.

308 *Stacks of water molecules.* The interface of the homodimer 1L5W is quite
309 extensive and highly non planar, consisting of two ‘prongs’ separated by a cleft.
310 Two high-VSO patches are found on either of the prongs. The *ABW* interface

311 is discontinuous in the region of the cleft, due to the presence of more than
312 one layer of solvent molecules sandwiched between the partners (Figure 9);
313 this resets the shelling order to low values in that area. On the other hand,
314 MD simulations find a much smaller patch of dry residues that extends inside
315 the cleft, which means that some of the aforementioned solvent molecules
316 are in fact structural in nature, and do not move during the simulation. A
317 remarkable example of this occurs for tryptophane 203 (located inside the
318 cleft), which is classified as dry by Mihalek and coworkers but is surrounded by
319 numerous water molecules on Figure 9. Here we are confronted with the main
320 advantage of MD simulations over our model: they are able to discriminate
321 structural water on the basis of residence times, whereas our static model
322 relies on the fact that buried interfacial water does not usually form multiple
323 layers. However, it is clear from Tables 1 and 2 that situations featuring water
324 molecules structured along more than one layer rarely occur; we discuss this
325 issue further in section 4. Within the interface, conservation fares better since
326 one of the prongs and the cleft region are fairly well conserved. However, the
327 most conserved regions lie at the protein core (not visible on the figure) and,
328 to a lesser extent, elsewhere on the protein surface.

329 *Discontinuities of the interface.* Figure 8 shows a graphical representation
330 of shelling, conservation and dryness for complex 1A59. 1A59 has an intri-
331 cate topology, consisting of two monomers of predominantly globular nature
332 linked by long ‘tails’ wrapped around the partner. Dry residues appear both
333 on the globular part and on the first segment of the tail (Figure 8). Voronoi
334 shelling order very accurately predicts the latter patch of dry residues, but
335 over-predicts the entire tail as being dry or active, too. More interestingly, it
336 also misses the lower part of the dry patch on the globular side of the protein.

337 A careful inspection of the interface reveals two holes in the AB interface which
338 reset the shelling order there, preventing the shelling order from peaking in
339 this region (Figure 10). The fact that such holes are visible in the AB interface
340 hints at a sizable packing issue: minute defects do not usually result in such
341 discontinuities of the AB interface[27]. Indeed, the gaps between the atoms of
342 the two monomers ¹ span the range 5.2-6.2 Å and 5.9-6.3 Å, respectively, and
343 could accommodate a water molecule each. Since the crystal structure does
344 not contain structural water, we cannot ascertain whether this is the case and
345 our fast solvation procedure proved unable to fill the holes – even though it
346 did successfully place isolated water molecules in three other locations. By
347 comparison, conservation correlates with dryness on the globular part of the
348 interface, but also features widespread conserved patches covering most of the
349 protein surface.

¹ Hole 1: residues 209 to 213 (chain A) and 583 to 587 (chain B); hole 2: residues 206 to 210 (chain A) and 586 to 590 (chain B).

350 4 DISCUSSION AND CONCLUSION

351 4.1 *A quantitative interface definition*

352 Among the various definitions of what exactly constitutes a protein-protein
353 interface, the planar facets obtained from a Voronoi tessellation [39,40] ar-
354 guably present the closest ties to the literal meaning of the term ‘interface’.
355 Indeed, such facets stem from pairs of directly interacting atoms, and the
356 definition of the interaction area is simpler than that required by analytical
357 interface models [41]. The Voronoi model shows excellent correlation with clas-
358 sically defined curvature and solvent accessible area but captures the interface
359 more fully than methods based on solvent accessibility [27] —see also [42] for
360 a review on the use of Voronoi diagrams in protein structure and interface
361 analysis. By contrast, the widely used geometric footprint (based on residue
362 contacts) yields an ambiguous interaction layer biased towards large residues
363 and subject to an arbitrary distance cut-off [3].

364 Here, we go beyond the binary classification of whether or not a given atom
365 is part of the interface and furthermore quantify how many facets separate it
366 from the edge of the interface. The idea is related to the concept of residue or
367 atom depth [43,44] which shows some correlation with thermodynamic prop-
368 erties [43] and residue conservation [45] in globular proteins. Previous studies
369 have defined atomic depth as the simple Euclidean distance to the closest sol-
370 vent molecule. By contrast, Voronoi shelling order partitions the interface into
371 concentric shells, accounting for both the geometry and topology of the inter-
372 face and appears closer to physical reality. Yet other previous studies have dis-
373 sected protein interfaces into “inner” and “outer” or “core” and “rim” residues

374 (for example, [46–48,36]). Although a number of general trends emerge, con-
375 clusions from these works are hindered by distinct definitions of the interface
376 combined with different classifications for core and rim. Voronoi shelling order
377 provides a more quantitative, parameter-free and unambiguous alternative to
378 the ad-hoc classifications previously employed.

379 *4.2 Shelling order and water dynamics*

380 The shelling of the Voronoi interface yields an accurate quantification for the
381 concept of burial depth. Shelling order quantifies the number of atomic shells
382 a water molecule must pass on the shortest path to a given position (facet) in
383 the interface. This description is particularly valuable for highly curved inter-
384 faces (1A59, 1L5W...) which the Euclidean distance cannot correctly measure.
385 We have here revealed a clear correlation between Voronoi shelling order and
386 the ‘dryness’ of a residue, that is, its shielding from itinerant bulk solvent
387 molecules. While one could expect some ties between the two measures, the
388 extent of the agreement over a representative set of complexes is intriguing.
389 After all, dryness was derived from exhaustive molecular dynamics simulations
390 which consider hundreds of additional parameters and details that are totally
391 ignored by our model. On the contrary, Voronoi shelling order is a purely
392 geometric property, calculated from a static set of atomic positions without
393 any further parameter. In particular, we do *not* consider: electrostatic charges,
394 polarity, hydrogen bonds, or any kind of fluctuations – all of which are ex-
395 pected to influence water dynamics. This suggests that the seemingly complex
396 dynamic exchange of bulk solvent with interfacial water primarily depends on
397 a simple path length and could tentatively be approximated by an analytical

398 model of diffusion along a gradient.

399 *4.3 Complementarity of conservation and Voronoi shelling order*

400 Evolutionary conservation alone cannot usually be employed to predict the
401 active part of an interface, let alone the interface itself. Hence the neces-
402 sity to cross-correlate it with some other measure (like geometric footprint or
403 change in solvent accessibility) before using it for such purposes. By compar-
404 ison, Voronoi shelling order simultaneously offers an unambiguous definition
405 of the protein-protein interface and a more fine-grained classification within
406 this interface.

407 Furthermore, the quantification of evolutionary signals is not trivial. pFam
408 sequence alignments are considered high quality but are not guaranteed to be
409 homogeneously distributed between protein families, hereby introducing bias.
410 Moreover, some protein stretches cannot be aligned at all, and needed to be
411 excluded from our analysis of conservation. We quantify conservation with an
412 entropy-based measure that has been shown to outperform other conservation
413 scores [35]; alternative means can be employed but the actual method of choice
414 seems to have limited effect on the correlation with dryness[26].

415 Bearing in mind the interference from many other factors, sequence conserva-
416 tion can, nevertheless, provide independent testimony of an area's importance.
417 It confirms the notion of water shielding as an indicator of binding activity
418 and it supports the functional relevance of shelling order. In fact, conservation
419 and VSO are best used in conjunction rather than as competitors. We find a
420 general correlation between shelling order and conservation but, in contrast to

421 a simple classification into rim and core, our continuous measure also resolves
422 interesting deviations from this trend. Such deviations hint at catalytic sites,
423 defects in solvation and packing, but may also indicate binding contributions
424 that do not directly rely on water shielding.

425 4.4 Methodological improvements

426 As previously discussed, discrepancies between dryness and shelling order arise
427 for cases where structural (slow moving) water molecules form more than
428 one layer inside a cavity. This is due to the fact that in our current model,
429 interfacial water molecules must make simultaneous contact with both protein
430 partners; any additional layer of water molecules not fulfilling this criterion
431 will be considered as bulk and lead to the splitting of the *ABW* interface.
432 However, ‘trapped’ water molecules are known to stabilize turns and bends
433 through hydrogen bonding with main-chain atoms in otherwise unstructured
434 regions [49], and cannot be ignored. Their behavior is so different from that of
435 bulk water that it is debatable whether they should be considered as delimiters
436 for the interface, even when stacked in more than one layer – dryness results
437 from MD simulations tend to show that they shouldn’t.

438 The most straightforward approach to alleviate discrepancies between dryness
439 and shelling order in these difficult cases would be to optimize the threshold
440 separating ‘dry’ from ‘wet’, instead of using Mihalek’s choice [26]. Our model
441 could also be extended so as to declare as interface water all solvent molecules
442 W_i found on a path $AW_1 \dots W_k B$ joining both partners. Using $k = 2$ or $k = 3$
443 could allow to infer similar properties for water molecules organized in layers,
444 as in complex 1L5W. Nevertheless, the current interface model, despite using

445 $k = 1$, demonstrates that it is legitimate to infer dryness/activity from a
446 purely geometric perspective. This effectively replaces a costly MD simulation
447 by a very fast computation on a structure taken directly from the PDB.

448 Another worthwhile methodological improvement would address rare cases
449 where discontinuities in the interface appear due to packing or solvation de-
450 fects. An example thereof is the previously discussed 1A59 interface (Fig-
451 ure 10). Regardless of the quality of the structure or the equilibration proce-
452 dure, such cases could be accommodated by using a water probe radius larger
453 than 1.4 Å, or by devising an adaptive scheme for the value of α ($\alpha > 0$)
454 employed to construct the α -complex. In any case, these extensions should be
455 investigated in conjunction with the threshold used to define dryness.

456 *4.5 Conclusion*

457 In this paper, we present a novel method to explore protein-protein interfaces.
458 The interface is defined using the Voronoi diagram of interacting atom pairs;
459 unlike geometric footprinting methods, all atoms involved in the interface are
460 identified with little to no over-prediction and without resorting to a distance
461 threshold. We have shelled the Voronoi interface from the rim to the core, thus
462 associating an interface depth to each atom. This Voronoi shelling order (VSO)
463 correlates very well with the protection of residues from itinerant water fluxes,
464 as computed by Mihalek and coworkers [26] which, in turn, can be considered
465 a measure of residue activity. The calculation of shelling orders, however, is
466 about five orders of magnitude faster than a typical MD simulation. Moreover,
467 the rather accurate prediction from a simplistic and purely geometric model
468 hints at the possibility to approximate the complex dynamics of interfacial

469 water by simple analytic diffusion models. Comparison with evolutionary sig-
470 nals confirms the functional relevance of ‘dry’ residues and, likewise, reveals
471 a general increase of conservation towards inner interface shells. Systematic
472 deviations from this trend may inform about distinct binding mechanisms,
473 catalytic activities but also modeling errors. Our accurate and continuous
474 scale of burial depths could also be used to delimitate patches on an interface.
475 Hence, it appears as a worthy candidate for the theoretical study of collective
476 effects in protein-protein interfaces [13], which are progressively replacing the
477 traditional ‘hotspot’ view.

478 5 METHODS

479 5.1 *Complex preparation*

480 The coordinates for the homo- and heterodimer complexes listed in Tables 1
481 and 2 originate from the PDB database. Crystallographic water molecules
482 were removed in order to exclude bias from different structure qualities. Miss-
483 ing atoms, including polar hydrogens, were added and briefly minimized. The
484 structure was surrounded by a 9 Å layer of water molecules from an equili-
485 brated TIP3P box. The water was briefly minimized by 3 rounds of conjugate-
486 gradient optimization of 40 steps each with, initially (round 1), frozen and
487 later (rounds 2 and 3) harmonically restrained protein coordinates. Keeping
488 this restraint, the water was then further relaxed by 100 2-fs steps of molecu-
489 lar dynamics at 100 K, followed by 40 steps conjugate gradient minimization.
490 Optimizations and simulations were performed using the CHARMM19 force
491 field [50] and an electrostatic cutoff of 12 Å with force shifting [51] inside the
492 X-PLOR package. This structure preparation protocol is automated by the
493 `pdb2explor.py` program which is part of the open source Biskit package [52].
494 The final structure was stripped of its hydrogen atoms and used as input for
495 the Voronoi interface calculations (see below).

496 To test the legitimacy of this economical solvation procedure, a more thorough
497 approach was employed on complex 1M0S. After an initial re-optimization of
498 the crystal structure (retaining crystal water), the complex was placed in-
499 side a triclinic box, solvated with SPC water molecules from an equilibrated
500 box and neutralized by 8 Na^+ ions. The solvent molecules were then relaxed
501 around the fixed solute by a steepest-descent optimization followed by 100 ps

502 of molecular dynamics (MD) simulation with position restraints on the so-
503 lute. The entire system was then simulated for 5 ns without restraints, with a
504 300 K Maxwellian distribution of initial velocities. MD simulations employed
505 the particle-mesh Ewald treatment of long-range electrostatics and periodic
506 boundary conditions, as well as couplings to heat (300 K, 1 ps) and pressure
507 (1 bar, 1 ps) baths; they were performed with GROMACS 3.3.2 [53] using
508 the OPLS all-atom force field [54]. The final equilibrated box had dimensions
509 76x92x69 Å and comprised 13460 water molecules. Convergence of the protein
510 structure was reached after 2 ns of simulation, at a mean RMSD of 1.90 Å
511 from the crystal structure.

512 Section A.2 of Supplemental Material compares the Voronoi interfaces of com-
513 plex 1M0S using these two equilibration procedures. The very similar results,
514 both in terms of interface topology and the identification of interfacial wa-
515 ter, justify the economical solvation method and indicate the robustness of
516 our model against minor changes both in protein conformation and hydration
517 patterns.

518 *5.2 Calculation of shelling orders*

519 The program Intervor, responsible for the actual computation and shelling of
520 the Voronoi interface, is based on the CGAL computational geometry library
521 [55]; an online version of Intervor is available [56]. On an Intel Pentium IV 3
522 GHz CPU, an Intervor run for a typical complex takes less than 5 seconds. We
523 also provide a wrapper (Biskit.Intervor) for integrating the stand-alone pro-
524 gram in Biskit workflows. Residue shelling orders were calculated by averaging
525 over a residue's interface atoms.

526 *5.3 Dryness and conservation*

527 Dryness results were those discussed in [26] and were kindly provided to us by
528 O. Lichtarge and coworkers.

Multiple sequence alignments were obtained from the pFam database [33] of HMMER profiles [34] using the HMMER software version 2.3.1. Protein family profiles matching a given sequence were identified with hmmpfam using a conservative E-value and bit score cutoff of 1e-8 and 60, respectively. The sequence was then aligned to the matching profile with the hmalign program. Following [35], the conservation of each alignment position was quantified by the Kullback-Leibler divergence (relative entropy) between the HMM emission probabilities p and the background distribution of amino acids in SwissProt q :

$$s = \sum_{i=1}^{20} p_i \log \frac{p_i}{q_i}.$$

529 The complete procedure is automated in the Hmmer.py module of Biskit.
530 Before further analysis, residues outside the interface (average $VSO = 0$) or
531 lacking conservation scores were removed and conservation scores were inde-
532 pendently normalized to the maximum of each monomer face.

533 *5.4 ROC curves*

534 Receiver Operating Characteristics (ROC) curves[57] are an efficient way of
535 representing the accuracy of a binary classifier. A binary classifier maps in-
536 stances of an object into two categories, positive or negative, based on each
537 instance's position relative to a threshold. The quality of the classifier is then
538 assessed by how well the prediction relates to the actual value of the instance.

539 Four cases are possible: true positive (both the outcome from a prediction and
540 the actual value are positive), false positive (the prediction is positive while
541 the actual value is negative), true negative (prediction and value are both
542 negative) and false negative (prediction is negative while value is positive).
543 From this contingency table, the notions of selectivity and sensitivity can be
544 defined as

$$Sensitivity = \frac{True\ Positive}{True\ Positive + False\ Negative}$$

and

$$Specificity = \frac{True\ Negative}{True\ Negative + False\ Positive}$$

545 A ROC curve is the 2D plot of sensitivity versus specificity, where each point
546 corresponds to a different threshold value. A perfect predictor, which features
547 neither false positive nor false negative occurrences, should pass through the
548 point (1,1) for the optimal threshold value. Therefore, the closer the ROC
549 plot is to the upper right corner, the higher the overall accuracy of the test
550 [58]. A purely random classifier, with equal chances of making correct or er-
551 roneous predictions, has a linear ROC curve connecting points (0,1) and (1,0)
552 – the first diagonal. How much better than random a predictor is can hence
553 be quantified by calculating the area between its ROC curve and the diag-
554 onal, which varies from -0.5 (worst-case classifier) to 0.5 (perfect classifier)
555 through 0 (pure random classifier). ROC curve and ROC area calculations
556 were performed with the Biskit.ROCalyzer module.

557 By way of example, figure 4 shows typical ROC curves for shelling order and
558 conservation as predictors for dryness, in the specific case of the 1HE1 complex.
559 For this system, shelling order is systematically better than conservation at

560 predicting dryness, regardless of the threshold chosen to discriminate between
561 positive and negative predictions in each case. This translates into a larger
562 area between the diagonal (representing a random prediction) and the shelling
563 order ROC plot, than between the diagonal and the conservation ROC plot.

564 5.5 Miscellaneous

565 The Biskit python package [52] was also used for various other scripting tasks
566 and the collation of results. All parts of Biskit are open source and available at
567 <http://biskit.sf.net>. Pymol [59], Ipe [60] and CGAL-Ipelets [61] were employed
568 for the rendering of figures.

569 **Acknowledgments.** *We would like to express our gratitude to Olivier Lichtarge*
570 *and Tuan Anh Tran for providing us with their detailed dryness results. The*
571 *automatic generation of conservation profiles was implemented by Johan Leck-*
572 *ner. B. Bowvier acknowledges funding from the INRIA Cooperative project*
573 *ReflexP. R. Grünberg is supported by the Human Frontiers Science Program.*

574 References

- 575 [1] A.-C. Gavin, P. Aloy, P. Grandi, R. Krause, M. Boesche, M. Marzioch,
576 C. Rau, L. J. Jensen, S. Bastuck, B. Dmpelfeld, A. Edelmann, M.-A. Heurtier,
577 V. Hoffman, C. Hoefert, K. Klein, M. Hudak, A.-M. Michon, M. Schelder,
578 M. Schirle, M. Remor, T. Rudi, S. Hooper, A. Bauer, T. Bouwmeester,
579 G. Casari, G. Drewes, G. Neubauer, J. M. Rick, B. Kuster, P. Bork, R. B.
580 Russell, G. Superti-Furga, Proteome survey reveals modularity of the yeast cell
581 machinery., Nature 440 (7084) (2006) 631–636.

- 582 URL <http://dx.doi.org/10.1038/nature04532>
- 583 [2] J. J. Gray, High-resolution protein-protein docking., *Curr Opin Struct Biol*
584 16 (2) (2006) 183–193.
585 URL <http://dx.doi.org/10.1016/j.sbi.2006.03.003>
- 586 [3] R. Grünberg, J. Leckner, M. Nilges, Complementarity of structure ensembles
587 in protein-protein binding., *Structure* 12 (12) (2004) 2125–2136.
588 URL <http://dx.doi.org/10.1016/j.str.2004.09.014>
- 589 [4] R. Grünberg, M. Nilges, J. Leckner, Flexibility and conformational entropy in
590 protein-protein binding., *Structure* 14 (4) (2006) 683–693.
591 URL <http://dx.doi.org/10.1016/j.str.2006.01.014>
- 592 [5] C. Chotia, J. Janin, Principles of protein-protein recognition, *Nature* 256 (1975)
593 705–708.
- 594 [6] J.-L. K. Kouadio, J. R. Horn, G. Pal, A. A. Kossiakoff, Shotgun alanine scanning
595 shows that growth hormone can bind productively to its receptor through a
596 drastically minimized interface., *J Biol Chem* 280 (27) (2005) 25524–25532.
597 URL <http://dx.doi.org/10.1074/jbc.M502167200>
- 598 [7] S. Jones, J. Thornton, Analysis of protein-protein interaction sites using surface
599 patches, *J. Mol. Biol.* 272.
- 600 [8] B. Ma, T. Elkayam, H. Wolfson, R. Nussinov, Protein-protein interactions:
601 Structurally conserved residues distinguish between binding sites and exposed
602 protein surfaces, *Proceedings of the National Academy of Sciences* 100 (10)
603 (2003) 5772–5777.
604 URL <http://www.pnas.org/cgi/content/abstract/100/10/5772>
- 605 [9] L. Lo Conte, C. Chothia, J. Janin, The atomic structure of protein-protein
606 recognition sites, *Journal of Molecular Biology* 285 (1999) 2177–2198.

607 URL [http://www.sciencedirect.com/science/article/
608 B6WK7-45R884M-RY/2/6f4a9866e2495a34273695de046893dc](http://www.sciencedirect.com/science/article/B6WK7-45R884M-RY/2/6f4a9866e2495a34273695de046893dc)

609 [10] S. A. Teichmann, Principles of protein-protein interactions, *Bioinformatics*
610 18 (suppl. 2) (2002) S249–.

611 URL [http://bioinformatics.oxfordjournals.org/cgi/content/
612 abstract/18/suppl_2/S249](http://bioinformatics.oxfordjournals.org/cgi/content/abstract/18/suppl_2/S249)

613 [11] T. Clackson, J. Wells, A hot spot of binding energy in a hormone-receptor
614 interface, *Science* 267 (5196) (1995) 383–386.

615 URL <http://www.sciencemag.org/cgi/content/abstract/267/5196/383>

616 [12] I. S. Moreira, P. A. Fernandes, M. J. Ramos, Hot spots – a review of the
617 protein-protein interface determinant amino-acid residues, *Proteins: Structure,
618 Function, and Bioinformatics* 68 (4).

619 URL <http://dx.doi.org/10.1002/prot.21396>

620 [13] D. Reichmann, O. Rahat, S. Albeck, R. Meged, O. Dym, G. Schreiber, The
621 modular architecture of protein-protein binding interfaces, *PNAS* 102 (1) (2005)
622 57–62.

623 URL <http://www.pnas.org/cgi/content/abstract/102/1/57>

624 [14] D. Reichmann, O. Rahat, M. Cohen, H. Neuvirth, G. Schreiber, The molecular
625 architecture of protein-protein binding sites, *Current Opinion in Structural
626 Biology* 17 (2007) 67–76.

627 [15] D. R. Caffrey, S. Somaroo, J. D. Hughes, J. Mintseris, E. S. Huang, Are protein-
628 protein interfaces more conserved in sequence than the rest of the protein
629 surface?, *Protein Sci* 13 (1) (2004) 190–202.

630 URL <http://www.proteinscience.org/cgi/content/abstract/13/1/190>

631 [16] B. E. Shakhnovich, N. V. Dokholyan, C. DeLisi, E. I. Shakhnovich, Functional
632 fingerprints of folds: Evidence for correlated structure-function evolution,

- 633 Journal of Molecular Biology 326 (1) (2003) 1–9.
634 URL [http://www.sciencedirect.com/science/article/
635 B6WK7-47RC32Y-3/2/780a30e8b2d84bdefba0ea80b0a6b6b2](http://www.sciencedirect.com/science/article/B6WK7-47RC32Y-3/2/780a30e8b2d84bdefba0ea80b0a6b6b2)
- 636 [17] A. Valencia, Automatic annotation of protein function, Current Opinion in
637 Structural Biology 15 (3) (2005) 267–274.
638 URL [http://www.sciencedirect.com/science/article/
639 B6VS6-4G7X9HS-3/2/90d887674957c53860854c3254a94748](http://www.sciencedirect.com/science/article/B6VS6-4G7X9HS-3/2/90d887674957c53860854c3254a94748)
- 640 [18] P. Aloy, H. Ceulemans, A. Stark, R. B. Russell, The relationship between
641 sequence and interaction divergence in proteins, Journal of Molecular Biology
642 332 (5) (2003) 989–998.
643 URL [http://www.sciencedirect.com/science/article/
644 B6WK7-49HMCGT-4/2/6d13750d3d40cc10b07ef096ef54961b](http://www.sciencedirect.com/science/article/B6WK7-49HMCGT-4/2/6d13750d3d40cc10b07ef096ef54961b)
- 645 [19] R. P. Bahadur, P. Chakrabarti, F. Rodier, J. Janin, A dissection of specific and
646 non-specific protein-protein interfaces, Journal of Molecular Biology 336 (4)
647 (2004) 943–955.
648 URL [http://www.sciencedirect.com/science/article/
649 B6WK7-4BRTKR9-C/2/ce716dee1132d2605f8e1c96d1154253](http://www.sciencedirect.com/science/article/B6WK7-4BRTKR9-C/2/ce716dee1132d2605f8e1c96d1154253)
- 650 [20] O. Lichtarge, H. R. Bourne, F. E. Cohen, An evolutionary trace method defines
651 binding surfaces common to protein families, Journal of Molecular Biology
652 257 (2) (1996) 342–358.
653 URL [http://www.sciencedirect.com/science/article/
654 B6WK7-45PV59P-1T/2/08c0824641e25a3fe8dba3e81635a653](http://www.sciencedirect.com/science/article/B6WK7-45PV59P-1T/2/08c0824641e25a3fe8dba3e81635a653)
- 655 [21] O. Rahat, A. Yitzhaky, G. Schreiber, Cluster conservation as a novel tool for
656 studying protein-protein interactions evolution., Proteins.
657 URL <http://dx.doi.org/10.1002/prot.21749>
- 658 [22] F. Rodier, R. Bahadur, P. Chakrabarti, J. Janin, Hydration of protein - protein

- 659 interfaces, *Proteins* 60 (1) (2005) 36–45.
- 660 [23] A. A. Bogan, K. S. Thorn, Anatomy of hot spots in protein interfaces, *Journal*
661 *of Molecular Biology* 280 (1998) 1–9.
- 662 URL [http://www.sciencedirect.com/science/article/
663 B6WK7-45S49GB-9C/2/b3d9c6f299c1eec3933d2774dffaf67d](http://www.sciencedirect.com/science/article/B6WK7-45S49GB-9C/2/b3d9c6f299c1eec3933d2774dffaf67d)
- 664 [24] A. Fernandez, R. S. Berry, Extent of Hydrogen-Bond Protection in Folded
665 Proteins: A Constraint on Packing Architectures, *Biophys. J.* 83 (5) (2002)
666 2475–2481.
- 667 URL <http://www.biophysj.org/cgi/content/abstract/83/5/2475>
- 668 [25] A. Fernandez, H. A. Scheraga, Insufficiently dehydrated hydrogen bonds as
669 determinants of protein interactions, *Proceedings of the National Academy of*
670 *Sciences* 100 (1) (2003) 113–118.
- 671 URL <http://www.pnas.org/cgi/content/abstract/100/1/113>
- 672 [26] I. Mihalek, I. Res, O. Lichtarge, On itinerant water molecules and detectability
673 of protein-protein interfaces through comparative analysis of homologues,
674 *Journal of Molecular Biology* 369 (2) (2007) 584–595.
- 675 [27] F. Cazals, F. Proust, R. P. Bahadur, J. Janin, Revisiting the Voronoi description
676 of protein-protein interfaces, *Protein Sci* 15 (9) (2006) 2082–2092.
- 677 URL <http://www.proteinscience.org/cgi/content/abstract/15/9/2082>
- 678 [28] F. M. Richards, The interpretation of protein structures: Total volume, group
679 volume distributions and packing density, *Journal of Molecular Biology* 82
680 (1974) 1–14.
- 681 [29] F. Aurenhammer, Power diagrams: properties, algorithms and applications,
682 *SIAM J. Comput.* 16 (1987) 78–96.
- 683 [30] C. Chotia, The nature of accessible and buried surfaces in proteins, *J. Mol. Bio.*
684 105 (1976) 1–12.

- 685 [31] B. Lee, F. M. Richards, The interpretation of protein structures: Estimation of
686 static accessibility, *Journal of Molecular Biology* 55 (3) (1971) 379–380.
687 URL [http://www.sciencedirect.com/science/article/
688 B6WK7-4DNGV9F-34/2/659a5209b127663bc3133c0b801a29a5](http://www.sciencedirect.com/science/article/B6WK7-4DNGV9F-34/2/659a5209b127663bc3133c0b801a29a5)
- 689 [32] H. Edelsbrunner, E. P. Mücke, Three-dimensional alpha shapes, *ACM Trans.*
690 *Graph.* 13 (1) (1994) 43–72.
- 691 [33] R. D. Finn, J. Mistry, B. Schuster-Bockler, S. Griffiths-Jones, V. Hollich,
692 T. Lassmann, S. Moxon, M. Marshall, A. Khanna, R. Durbin, S. R. Eddy,
693 E. L. L. Sonnhammer, A. Bateman, Pfam: clans, web tools and services, *Nucl.*
694 *Acids Res.* 34 (suppl. 1) (2006) D247–251.
695 URL [http://nar.oxfordjournals.org/cgi/content/abstract/34/suppl_
696 1/D247](http://nar.oxfordjournals.org/cgi/content/abstract/34/suppl_1/D247)
- 697 [34] R. Durbin, S. Eddy, A. Krogh, G. Mitchison, *Biological sequence analysis:*
698 *probabilistic models of proteins and nucleic acids*, Cambridge University Press,
699 1998, Ch. The theory behind profile HMMs.
- 700 [35] K. Wang, R. Samudrala, Incorporating background frequency improves entropy-
701 based residue conservation measures., *BMC Bioinformatics* 7 (2006) 385.
702 URL <http://dx.doi.org/10.1186/1471-2105-7-385>
- 703 [36] M. Guharoy, P. Chakrabarti, Conservation and relative importance of residues
704 across protein-protein interfaces., *Proc Natl Acad Sci U S A* 102 (43) (2005)
705 15447–15452.
706 URL <http://dx.doi.org/10.1073/pnas.0505425102>
- 707 [37] T. Selzer, G. Schreiber, Predicting the rate enhancement of protein complex
708 formation from the electrostatic energy of interaction., *J Mol Biol* 287 (2) (1999)
709 409–419.
710 URL <http://dx.doi.org/10.1006/jmbi.1999.2615>

- 711 [38] N. Ostermann, I. Schlichting, R. Brundiers, M. Konrad, J. Reinstein, T. Veit,
712 R. S. Goody, A. Lavie, Insights into the phosphoryltransfer mechanism of human
713 thymidylate kinase gained from crystal structures of enzyme complexes along
714 the reaction coordinate, *Structure* 8 (6) 629–642.
715 URL [http://www.sciencedirect.com/science/article/
716 B6VSR-40H578S-B/2/72acba30cd1e54cc5111cf865a8374b2](http://www.sciencedirect.com/science/article/B6VSR-40H578S-B/2/72acba30cd1e54cc5111cf865a8374b2)
- 717 [39] A. Varshney, F. P. Brooks, D. C. Richardson, W. V. Wright, D. Manocha,
718 Defining, computing, and visualizing molecular interfaces, in: *IEEE*
719 *Visualization*, Atlanta, USA, 1995, pp. 36–43.
- 720 [40] Y.-E. A. Ban, , H. Edelsbrunner, J. Rudolph, Interface surfaces for protein-
721 protein complexes, in: *RECOMB*, San Diego, 2004, pp. 205–212.
- 722 [41] R. R. Gabdoulline, R. C. Wade, D. Walther, Molsurfer: A macromolecular
723 interface navigator., *Nucleic Acids Res* 31 (13) (2003) 3349–3351.
- 724 [42] A. Poupon, Voronoi and voronoi-related tessellations in studies of protein
725 structure and interaction., *Curr Opin Struct Biol* 14 (2) (2004) 233–241.
726 URL <http://dx.doi.org/10.1016/j.sbi.2004.03.010>
- 727 [43] S. Chakravarty, R. Varadarajan, Residue depth: a novel parameter for the
728 analysis of protein structure and stability., *Structure* 7 (7) (1999) 723–732.
- 729 [44] A. Pintar, O. Carugo, S. Pongor, Atom depth in protein structure and function.,
730 *Trends Biochem Sci* 28 (11) (2003) 593–597.
- 731 [45] A. Pintar, O. Carugo, S. Pongor, Atom depth as a descriptor of the protein
732 interior., *Biophys J* 84 (4) (2003) 2553–2561.
- 733 [46] P. Chakrabarti, J. Janin, Dissecting protein-protein recognition sites, *Proteins*
734 47.
- 735 [47] I. M. A. Nooren, J. M. Thornton, Structural characterisation and functional

736 significance of transient protein-protein interactions., J Mol Biol 325 (5) (2003)
737 991–1018.

738 [48] A. Bordner, R. Abagyan, Statistical analysis and prediction of protein-protein
739 interfaces, Proteins 60 (3) (2005) 353–66.

740 [49] J. G. S. Sheldon Park, Statistical and molecular dynamics studies of buried
741 waters in globular proteins, Proteins: Structure, Function, and Bioinformatics
742 60 (2005) 450–463.

743 URL <http://dx.doi.org/10.1002/prot.20511>

744 [50] B. Brooks, R. Bruccoleri, Olafson B.D., D. States, S. Swaminathan, M. Karplus,
745 CHARMM: a program for macromolecular energy, minimization and dynamics
746 calculations., J Comp Chem 4 (1983) 187–217.

747 [51] P. Steinbach, R. Loncharich, B. Brooks, The effects of environment and
748 hydration on protein dynamics: A simulation study of myoglobin., Chem Phys
749 158 (1991) 383–94.

750 [52] R. Grünberg, M. Nilges, J. Leckner, Biskit–A software platform for structural
751 bioinformatics, Bioinformatics 23 (6) (2007) 769–770.

752 URL [http:
753 //bioinformatics.oxfordjournals.org/cgi/content/abstract/23/6/769](http://bioinformatics.oxfordjournals.org/cgi/content/abstract/23/6/769)

754 [53] D. V. D. Spoel, E. Lindahl, B. Hess, G. Groenhof, A. E. Mark, H. J. C.
755 Berendsen, Gromacs: fast, flexible, and free., J Comput Chem 26 (16) (2005)
756 1701–1718.

757 URL <http://dx.doi.org/10.1002/jcc.20291>

758 [54] W. Damm, A. Frontera, J. Tirado-Rives, W. L. Jorgensen, Opls all-atom force
759 field for carbohydrates, Journal of Computational Chemistry 18 (1997) 1955–
760 1970.

761 URL [http://dx.doi.org/10.1002/\(SICI\)1096-987X\(199712\)18:16<1955::](http://dx.doi.org/10.1002/(SICI)1096-987X(199712)18:16<1955::)
762 [AID-JCC1>3.0.CO;2-L](http://dx.doi.org/10.1002/(SICI)1096-987X(199712)18:16<1955::AID-JCC1>3.0.CO;2-L)

763 [55] CGAL, Computational Geometry Algorithms Library, <http://www.cgal.org>.

764 [56] <http://cgal.inria.fr/Intervor>.

765 [57] D. M. Green, J. M. Swets, Signal detection theory and psychophysics, John
766 Wiley and Sons Inc., New York, 1966.

767 [58] M. Zweig, G. Campbell, Receiver-operating characteristic (ROC) plots: a
768 fundamental evaluation tool in clinical medicine [published erratum appears
769 in Clin Chem 1993 Aug;39(8):1589], Clin Chem 39 (4) (1993) 561–577.

770 URL <http://www.clinchem.org/cgi/content/abstract/39/4/561>

771 [59] W. DeLano, The Pymol molecular graphics system, <http://www.pymol.org>
772 (2002).

773 [60] O. Cheong, The Ipe extensible drawing editor, [http://tclab.kaist.ac.kr/](http://tclab.kaist.ac.kr/ipe/)
774 [ipe/](http://tclab.kaist.ac.kr/ipe/) (1993-2007).

775 [61] <http://cgal-ipelets.gforge.inria.fr/>.

776 **6 FIGURES LEGENDS**

777 **Legend of Fig. 1.** Voronoi diagram (light solid lines) for a hypothetical
778 molecule consisting of four atoms (a_1 to a_4), and restriction of the balls to
779 their Voronoi regions. The α -complex ($\alpha = 0$) consists of the four vertices a_1
780 to a_4 , of the three edges a_1a_2 , a_1a_3 , a_2a_3 , and of the triangle $a_1a_2a_3$ formed
781 between them.

782 **Legend of Fig. 2.** (a) Shelling of the Voronoi interface of a dimer complex,
783 seen from the top. Solid dots represent protein atoms' centers, hollow dots
784 water atoms' centers; for clarity, all atomic radii have been taken equal and
785 the corresponding spheres omitted. The Voronoi facets composing the protein-
786 protein interface are colored according to their shelling order: one (light gray,
787 at the rim), two (middle gray), three (dark gray). (b) Two-dimensional illustration
788 of the Voronoi interface shelling of a dimer complex. Red and blue circles
789 represent the atoms of each partner, the green circle a water molecule. Inter-
790 face Delaunay edges, which connect atoms on different partners, are shown as
791 solid black (AB interface) or green ($AW - BW$ interface) lines; the Voronoi
792 facets are shown as dashes. Black numerals denote the shelling order of each
793 Delaunay edge/Voronoi facet, from which the atomic shelling orders (red, blue
794 and green numerals) can be derived (refer to text for details). On this simple
795 illustration, the high curvature of the $AW - BW$ interface due to the water
796 molecule accounts for the high shelling order of the blue atoms.

797 **Legend of Fig. 3.** (a) Voronoi interface of the 2DOR homodimer complex,
798 superimposed on the solvent accessible surface representation of one of the
799 monomers (gray); for clarity, the second monomer is not shown. The facet

800 shelling order varies from 1 (blue) to 6 (red). (b) Solvent accessible surface of
801 one monomer of the 2DOR complex, showing the shelling order of interface
802 atoms (color-coded as in panel b).

803 **Legend of Fig. 4.** ROC plots evaluating shelling order (solid line) and con-
804 servation (dashed line) as predictors for dryness. Each point on a ROC plot
805 corresponds to a different threshold value for the prediction. The plot for a
806 perfect predictor should pass through (1, 1); that of a random predictor (on
807 average) is the diagonal (dotted line). The area between the ROC curve and
808 the diagonal measure the performance of the predictor compared to random.

809 **Legend of Fig. 5.** Performance of shelling order (circles, solid line) and
810 conservation (squares, dashed line) as predictors of dryness, for all studied
811 heterodimer (left panel) and homodimer (right panel) complexes. Scores are
812 measured as the area between the corresponding ROC curve and the diago-
813 nal; complexes are sorted by decreasing shelling order score. Negative values
814 (hatched area) denote a performance that is no better (on average) than that
815 of a purely random classifier.

816 **Legend of Fig. 6.** Spatial distribution of conservation across heterodimer
817 interfaces. The conservation score for each interface residue, normalized to
818 the maximum score, is plotted against its normalized shelling order. Black —:
819 running average with a large window size (1/4 of all interface residues); Gray
820 —: all data points; - · -: Integral over running average; · · ·: Double integral
821 over running average; Δ : deviation of the double integral from 0.5 – values
822 below zero indicate conservation bias towards high shelling order (the core).

823 **Legend of Fig. 7.** Spatial distribution of conservation across homodimer
824 interfaces. See figure 6 and text for a detailed description.

825 **Legend of Fig. 8.** Projection of shelling order (large panels), dryness (lower
826 left-hand panel) and conservation (lower right-hand panel) on the molecular
827 surface of homocomplexes 1E2D (left), 1L5W (center) and 1A59 (right); one
828 of the monomers was removed for clarity. Cold (resp. hot) colors represent
829 low (resp. high) values; gray areas denote residues for which conservation
830 information was unavailable.

831 **Legend of Fig. 9.** View of the cleft region of the 1L5W interface, showing
832 the two protein partners as solid and mesh surfaces, respectively. Colors code
833 for shelling order, which is low inside the cleft due to the presence of numerous
834 water molecules which fragment the interface.

835 **Legend of Fig. 10.** Boundary of the AB interface of complex 1A59 (red
836 line), interfacial water (gray spheres), and $AW - BW$ interface (grey and green
837 Voronoi polygons). The holes pointed out by arrows prevent the shelling order
838 from peaking in the middle of the interface patch –compare to the bottom left
839 panel of complex 1A59 on Fig. 8.

PDB Id.	VSO \rightarrow dry	Conservation \rightarrow dry	Conservation \rightarrow VSO
1HE1	0.42	0.28	0.02
1CXZ	0.39	0.24	0.19
1CEE	0.39	0.12	0.11
1C1Y	0.36	0.17	0.05
1RRP	0.34	0.22	0.21
1FIN	0.34	0.10	0.18
1E96	0.34	-0.02	0.15
1ZBD	0.33	0.09	0.19
1FOE	0.33	0.19	0.27
1A0O	0.32	0.23	0.12
2TRC	0.32	-0.08	0.11
1GOT	0.32	0.13	0.23
1WQ1	0.31	0.19	0.08
1IBR	0.30	0.01	-0.14
1A2K	0.26	0.15	0.28
1LFD	0.25	0.26	0.15
1AGR	0.19	0.10	0.25
1YCS	0.16	0.16	0.29
avg.	0.31	0.14	0.15

Table 1

Heterodimers. Performance of shelling order (VSO) as a predictor for dryness, of conservation as a predictor for dryness, and of conservation as a predictor for shelling order, for each of the considered heterodimer complexes.

PDB Id.	VSO \rightarrow dry	Conservation \rightarrow dry	Conservation \rightarrow VSO
2BIF	0.45	0.09	0.02
1E5Q	0.45	0.15	0.31
1E2D	0.45	0.37	0.38
1H7T	0.45	0.12	0.17
1TB5	0.43	0.14	0.02
2DOR	0.42	0.19	0.13
1QIN	0.42	0.14	0.14
1E98	0.42	0.40	0.45
1J79	0.40	-0.09	-0.08
1NYW	0.40	-0.09	0.04
1BTO	0.38	0.27	0.12
1Y6R	0.38	0.17	0.03
1KER	0.37	0.14	0.08
1EK4	0.37	0.15	0.21
1LBX	0.37	0.21	0.11
1L9W	0.36	0.29	0.27
1AI2	0.36	0.18	-0.05
1W1U	0.35	0.07	-0.03
1DQX	0.33	0.10	-0.09
1E7Y	0.32	0.24	-0.06
1HKV	0.32	0.09	0.04
1M0S	0.32	0.07	0.34
1KC3	0.32	0.35	0.32
1M4N	0.31	0.17	0.14
1A59	0.31	0.15	0.19
1DQR	0.31	0.09	0.08
1AN9	0.30	0.11	0.06
1M7P	0.29	0.01	0.08
1TC2	0.29	-0.01	0.17
1AD3	0.28	-0.03	0.16
1ALN	0.27	0.14	0.04
1H16	0.27	-0.06	-0.02
1M9N	0.26	0.09	0.20
1L5W	0.24	0.18	0.25
1CG0	0.22	0.12	0.05
1LXY	0.21	0.10	0.11
avg.	0.34	0.13	0.12

Table 2

Homodimers. Performance of shelling order (VSO) as a predictor for dryness, of conservation as a predictor for dryness, and of conservation as a predictor for shelling order, for each of the considered homodimer complexes.

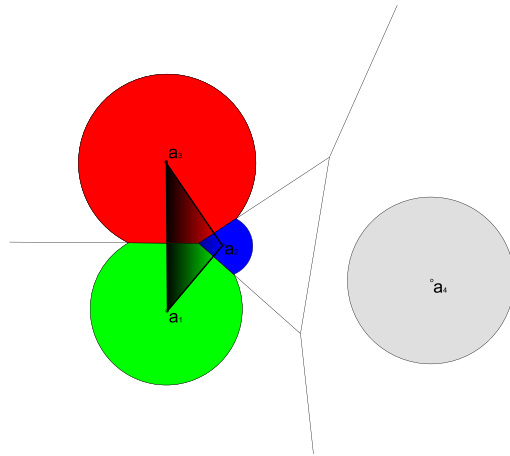


Fig. 1.

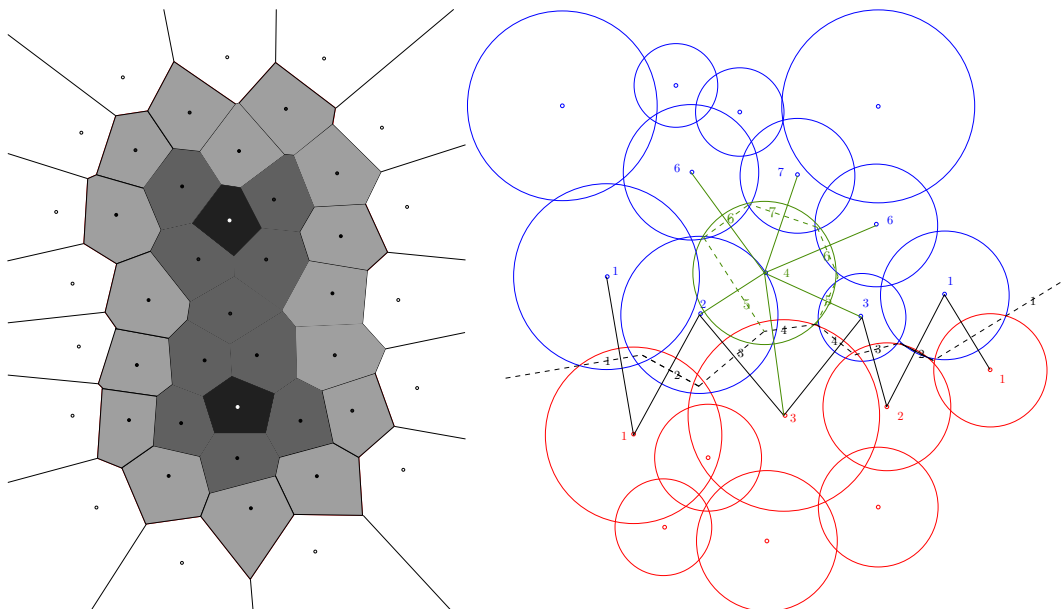


Fig. 2. (a) and(b)

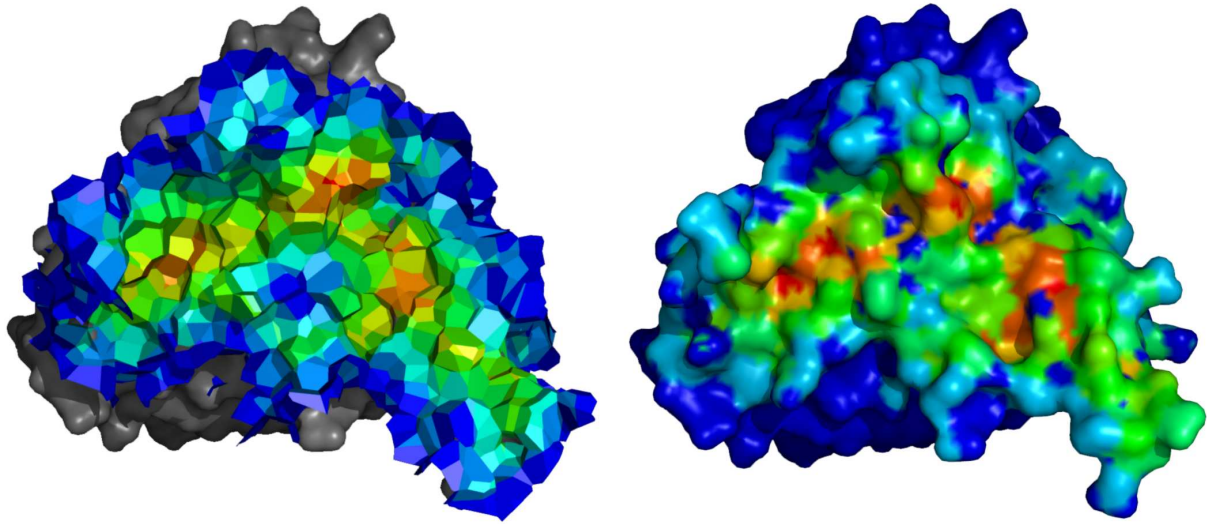


Fig. 3. (a) and(b)

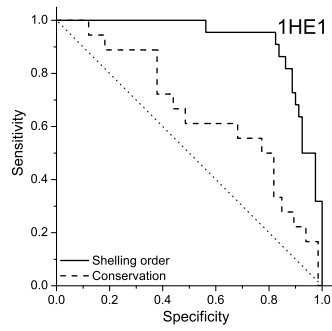


Fig. 4.

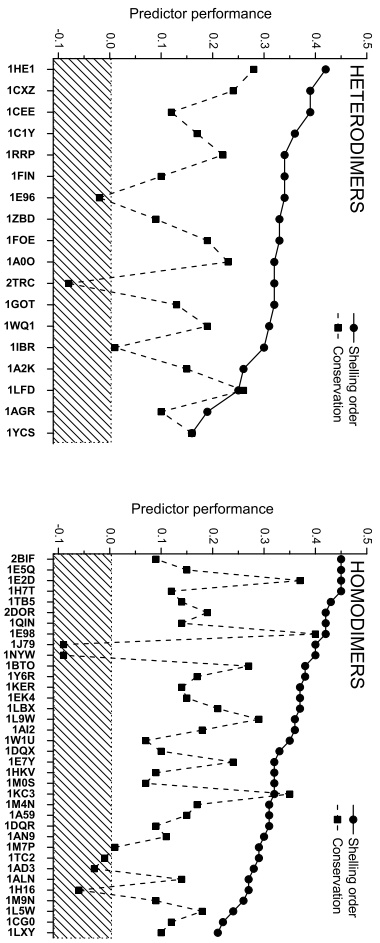


Fig. 5.

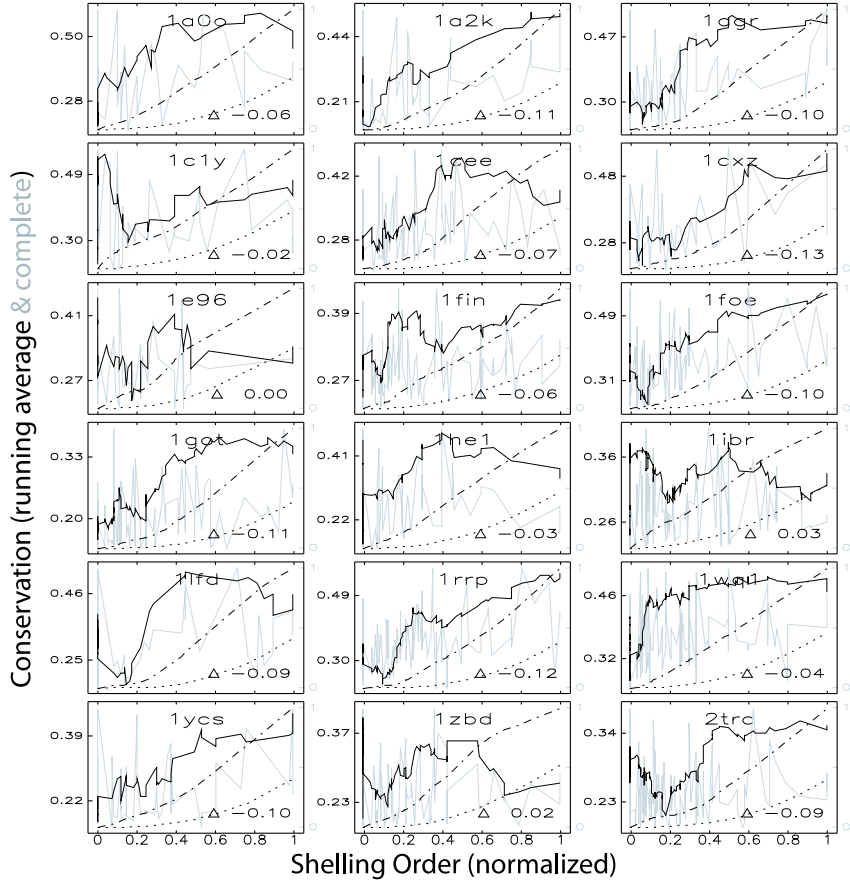


Fig. 6.

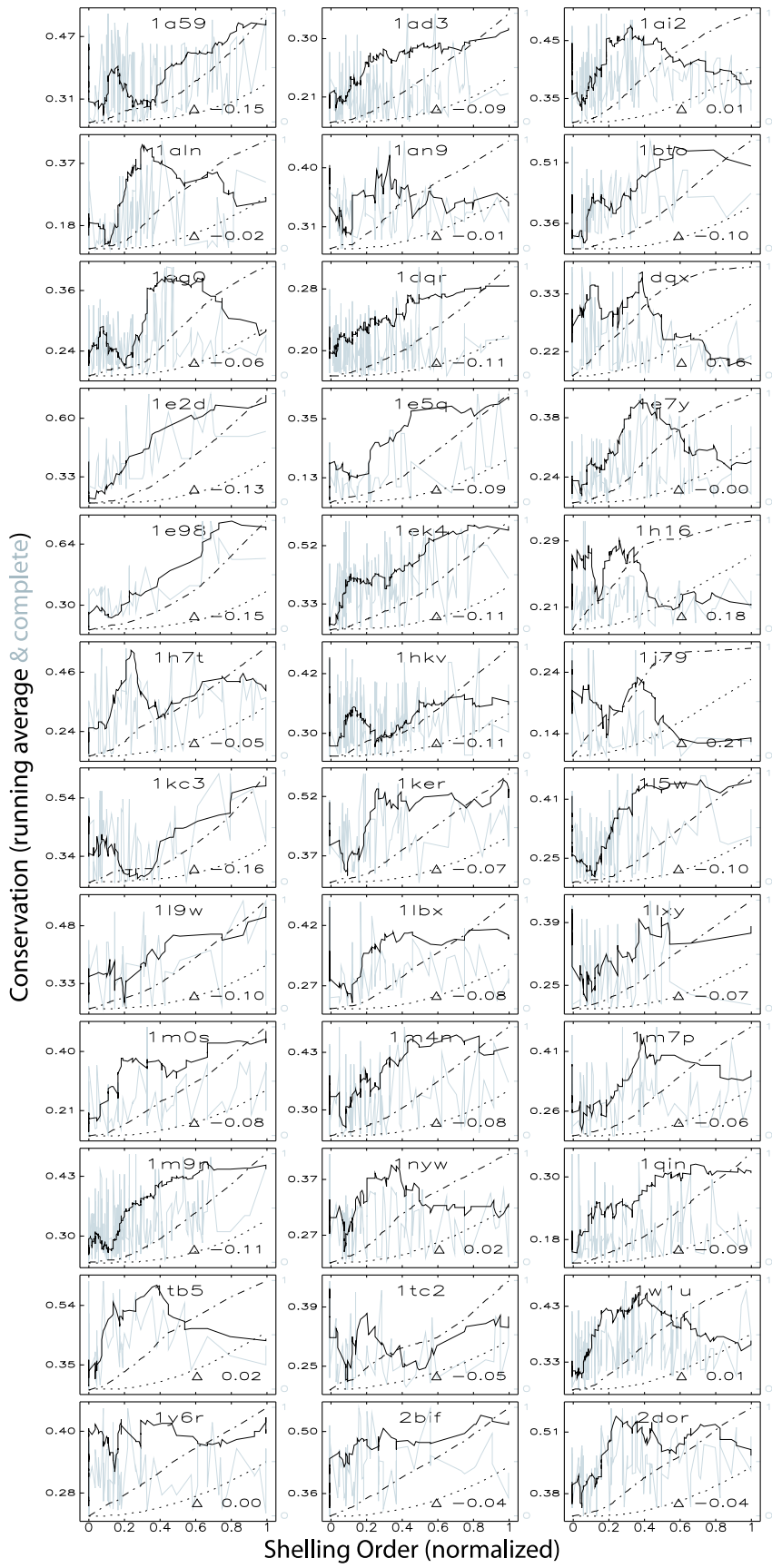


Fig. 7.

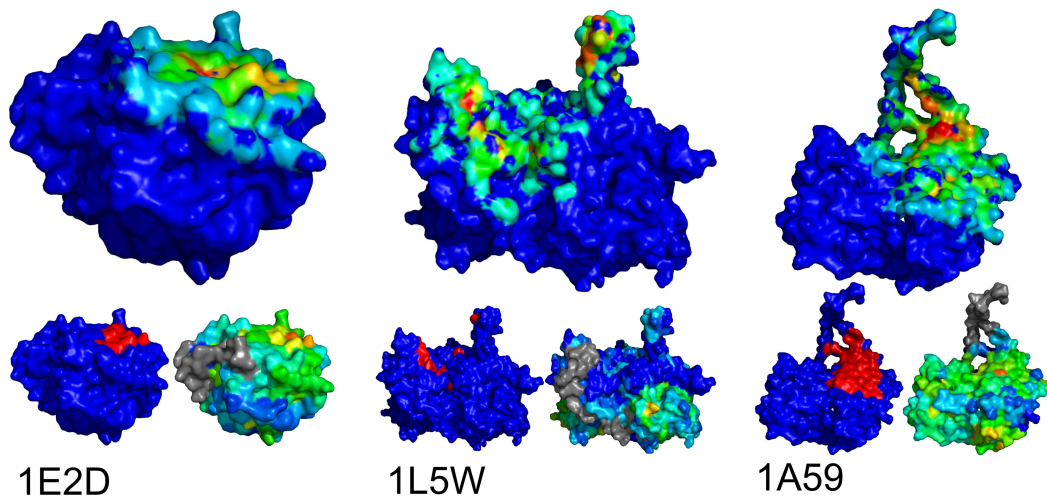


Fig. 8.

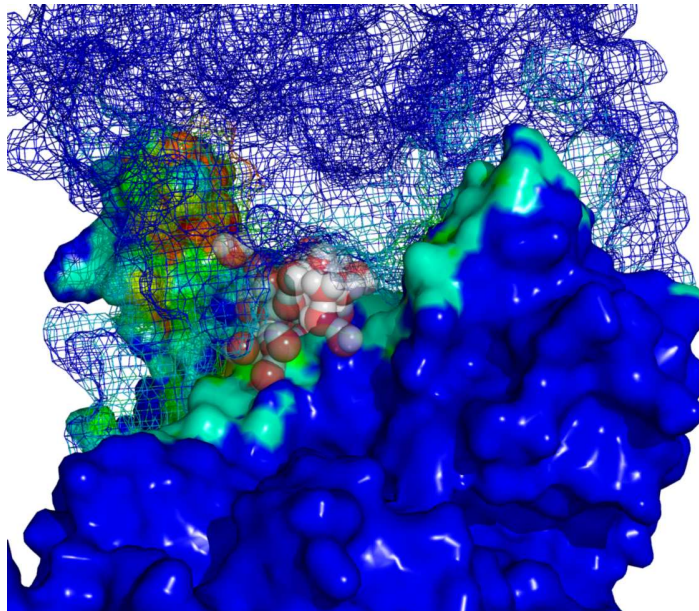


Fig. 9.

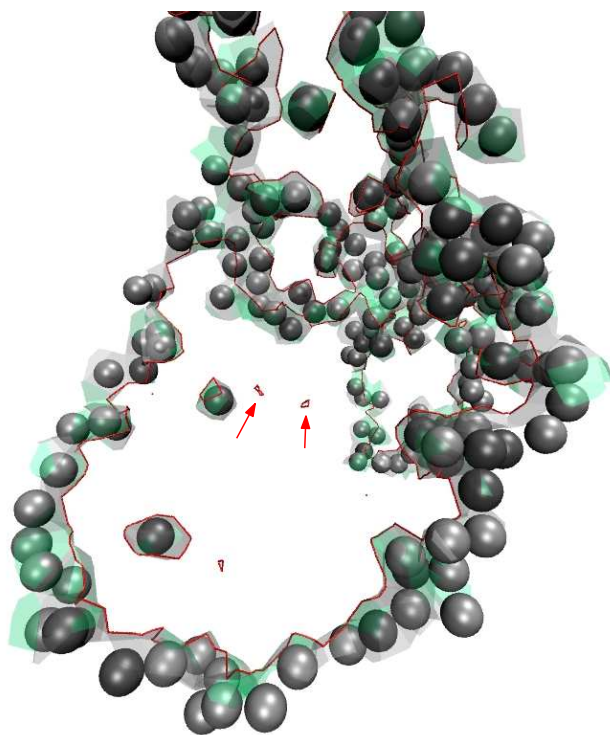


Fig. 10.

842 A Supplemental Material

843 A.1 *Distribution of conserved residues: interface residues versus dry residues*

844 As outlined in section 3.2, we compare the spatial distribution of conservation
845 in the entire set of interface residues with that of the dry residues.

846 We first consider all interface residues. To study the cumulated conservation
847 score over consecutive shells, we compute the proportion of the interface con-
848 servation score which is contained in the subset of residues whose average
849 VSO is lower than some value. Normalizing over shelling orders and varying
850 the threshold yields a curve that rises from $(0, 0)$ (no residues selected, zero
851 cumulative conservation) to $(1, 1)$ (all residues selected, 100% cumulative con-
852 servation). The area under this curve provides information about the variation
853 of conservation with shelling order, since numerous highly conserved residues
854 with low (high) shelling order will cause the curve to rise early (late) and
855 result in large (small) areas.

856 Next, we focus on the dry residues and construct references with which to
857 compare the previously computed areas, that quantify the relevance of rim
858 residue conservation in each case. Denoting n_{dry} the number of dry residues of
859 a given complex as reported in [26], we sort the interface residues by decreasing
860 shelling order and assume the first n_{dry} only to be conserved –those with
861 highest shelling orders. Let m and M be the minimum and maximum shelling
862 orders in this subset, respectively (note that M is also the highest VSO found
863 in the entire complex), and let $x = m/M$. The step function which is null from
864 0 to x , and equal to 1 from x to 1, maximizes the area $1 - x$ under the curve

865 relative to the conservation of the subset of n_{dry} residues.

866 As seen from Figure A.1, the rim residues account for a non-negligible part
867 of the conservation: the area under the corresponding curve was found to be
868 greater than the reference in all but two homodimer complexes, for which
869 both measures were roughly equal. This could, in part, be due to a purely
870 statistical effect: a large number of moderately conserved rim residues might
871 end up having more weight than a small number of highly conserved core
872 amino-acids. However, the peak in average conservation observed at the rim
873 of many complexes (Section 3.2 (i)) proves that highly conserved residues
874 occur on a non fortuitous basis at the rim of protein-protein interfaces – most
875 likely as anchors for important electrostatic interactions that dictate complex
876 formation and activity.

877 *A.2 Validation of the sample preparation procedure*

878 The procedure employed for the rehydration and equilibration of each of the
879 complexes (Section 5) has deliberately been kept short, and can be run in
880 minutes on a desktop computer. In this paragraph, we ascertain whether the
881 placement and equilibration of the water molecules added using this fast pro-
882 tocol are of sufficient quality for the current application. Of particular interest
883 are the interfacial water molecules. When in simultaneous contact with both
884 protein partners, they form the $AW - BW$ interface (Figure 2b and 10); but
885 several layers of water inside a larger pocket will create holes in the interface,
886 possibly splitting it into several connected components. The implications for
887 shelling orders are crucial: in the first case, the water molecules will not af-
888 fect the SO, while in the second scenario a boundary is created and the SO

889 consequently reset to 1.

890 The complex 1M0S, which features a large pocket filled with crystal water
891 molecules, was used for the test. A rigorous equilibration procedure, retaining
892 the crystal water molecules and involving a 5 ns molecular dynamics simulation
893 with state-of-the-art algorithms and parameters (Section 5), provided us with
894 a reference structure. Both this structure and the one from the fast procedure
895 were used as input to Intervor. Figure A.2 shows the tessellation of the AB
896 interface and the interfacial water molecules for both cases. Due to minor
897 conformational transitions that have occurred during the 5 ns MD simulation,
898 the two interfaces are not superposable. However, they retain the same shape
899 and number of connected components. In both cases, the central cavity is filled
900 with interfacial water that participates to the ABW interface. Both interfaces
901 feature boundaries of comparable lengths and topologies.

902 This difficult test case provides justification for our sample preparation method-
903 ology. It also represents a tribute to the robustness of our model, which de-
904 livers stable results upon variation of the solvation of the complex within a
905 reasonable range.

907 **Legend of Fig. A.1.** Area under the normalized cumulative conservation vs.
908 shelling order curve (circles, solid line) and reference area (squares, dashed
909 line), for all studied heterodimer (left panel) and homodimer (right panel)
910 complexes – see text for details. Areas larger than the reference denote com-
911 plexes for which rim residues are significantly conserved.

912 **Legend of Fig. A.2.** The AB interface (colored Voronoi facets) and the
913 interfacial water molecules W (grey spheres) for two distinct rehydration and
914 equilibration procedures – a fast (a) and a more exhaustive one (b); see text
915 for details. Boundaries of the AB and $AW - BW$ interfaces are shown as red
916 and green sticks, respectively.

PDB Id.	Proportion of conservation score for noninterface residues	of Area under curve, interface residues	Reference
1YCS	0.76	0.57	0.53
1RRP	0.61	0.66	0.57
1E96	0.83	0.65	0.52
1CXZ	0.78	0.61	0.52
1LFD	0.80	0.51	0.16
1WQ1	0.64	0.67	0.66
1FOE	0.77	0.68	0.67
1AGR	0.77	0.64	0.64
1IBR	0.77	0.70	0.66
1FIN	0.75	0.61	0.59
1HE1	0.61	0.70	0.60
1A2K	0.70	0.71	0.66
1A0O	0.71	0.64	0.48
1ZBD	0.79	0.72	0.66
1GOT	0.83	0.60	0.51
2TRC	0.71	0.71	0.66
1CEE	0.62	0.61	0.42
1C1Y	0.77	0.66	0.47

Table A.1

Relationship of shelling order and conservation for the heterodimer set: proportion of total conservation provided by noninterface residues, area under the normalized cumulative conservation vs. VSO curve (see text), area under the corresponding 'reference' curve (see text).

PDB Id.	Proportion of conservation score for noninterface residues	of Area under curve, interface residues	Reference
1A59	0.72	0.63	0.60
1H16	0.89	0.70	0.60
1M0S	0.73	0.49	0.42
1E5Q	0.97	0.55	0.32
1H7T	0.83	0.59	0.32
1E7Y	0.86	0.62	0.53
1ALN	0.64	0.60	0.62
1CG0	0.71	0.66	0.66
1E2D	0.81	0.64	0.55
1W1U	0.84	0.66	0.62
1KER	0.86	0.59	0.55
1EK4	0.74	0.63	0.64
1BTO	0.74	0.70	0.56
1QIN	0.36	0.62	0.45
1TB5	0.84	0.62	0.43
1M4N	0.76	0.64	0.52
2BIF	0.86	0.65	0.56
1M9N	0.57	0.70	0.68
1M7P	0.74	0.62	0.51
1E98	0.83	0.55	0.49
1L5W	0.95	0.70	0.62
1AD3	0.74	0.68	0.65
1J79	0.85	0.69	0.47
1AI2	0.62	0.68	0.61
1L9W	0.90	0.58	0.53
1LXY	0.87	0.66	0.51
1NYW	0.64	0.65	0.52
1KC3	0.87	0.66	0.58
1Y6R	0.72	0.68	0.66
1LBX	0.76	0.65	0.26
2DOR	0.72	0.59	0.43
1DQR	0.64	0.67	0.62
1AN9	0.85	0.64	0.56
1TC2	0.79	0.67	0.61
1HKV	0.72	0.63	0.54
1DQX	0.57	0.62	0.45

Table A.2

Relationship of shelling order and conservation for the homodimer set: proportion of total conservation provided by noninterface residues, area under the normalized cumulative conservation vs. VSO curve (see text), area under the corresponding 'reference' curve (see text).

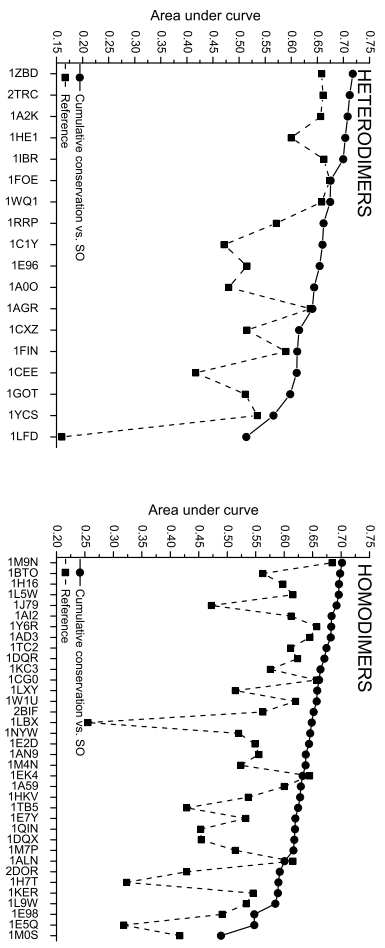


Fig. A.1.

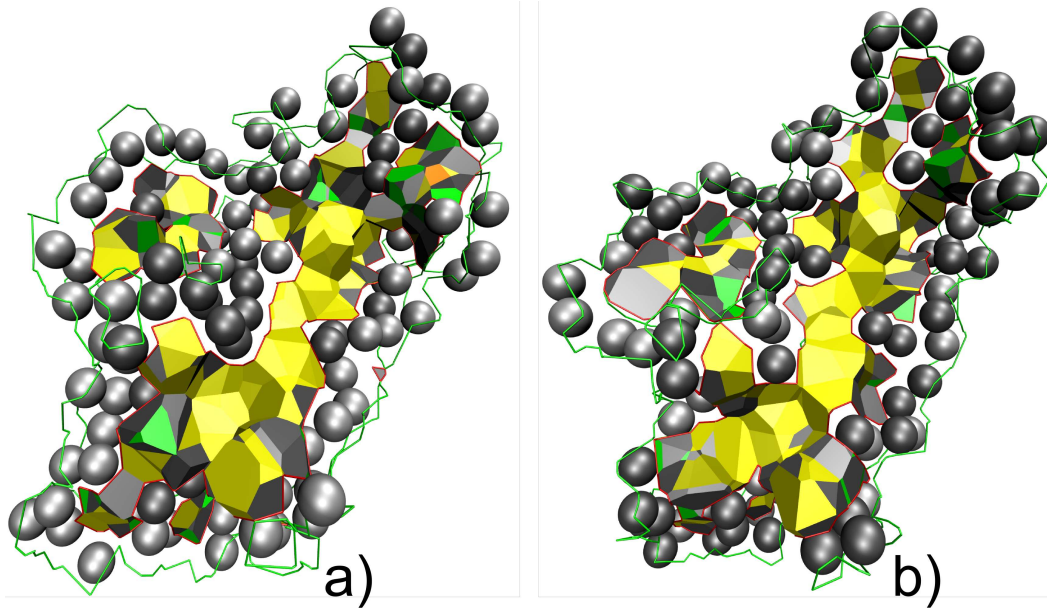


Fig. A.2.