

# Evolutionary couplings detect side-chain interactions

Adam J. Hockenberry<sup>a,1</sup> and Claus O. Wilke<sup>a</sup>

<sup>a</sup>Department of Integrative Biology, The University of Texas at Austin

This manuscript was compiled on October 18, 2018

1 **Patterns of amino acid covariation in large protein sequence align-**  
2 **ments can inform the prediction of *de novo* protein structures, bind-**  
3 **ing interfaces, and mutational effects. While algorithms that de-**  
4 **tect these so-called evolutionary couplings between residues have**  
5 **proven useful for practical applications, less is known about how**  
6 **and why these methods perform so well, and what insights into bio-**  
7 **logical processes can be gained from their application. Researchers**  
8 **frequently benchmark the performance of evolutionary coupling al-**  
9 **gorithms by comparing results with true structural contacts that are**  
10 **derived from solved protein structures. However, the method used to**  
11 **determine true structural contacts is not standardized and different**  
12 **definitions of structural contacts may have important consequences**  
13 **for comparing methods and understanding their overall utility. Here,**  
14 **we show that structural contacts between side-chain atoms are sig-**  
15 **nificantly more likely to be identified by evolutionary coupling analy-**  
16 **ses compared with backbone-based interactions. We use both simu-**  
17 **lation and empirical analyses to highlight that backbone-based defi-**  
18 **nitions of true residue-residue contacts may underestimate the accu-**  
19 **racy of evolutionary coupling algorithms by as much as 40%. These**  
20 **findings suggest that more advanced machine learning and neural**  
21 **network models developed to predict residue-residue contacts may**  
22 **be hindered by the use of mislabeled true positive training data.**

Protein evolution | Structural constraints | Contact prediction

1 **A** long-standing problem in biology is to predict the struc-  
2 ture of a protein based solely on its primary amino acid  
3 sequence (1–3). Despite advances in x-ray crystallography,  
4 NMR spectroscopy, and cryo-electron microscopy, the pace  
5 at which researchers are accumulating new genomes and gene  
6 sequences far outstrips the ability of traditional biophysical  
7 methods to describe these genomes at the level of 3D-structure  
8 (4–7). A variety of computational methods—such as homology  
9 modeling (8, 9)—have been developed to support traditional  
10 biophysical methods, but *de novo* structural determination  
11 from primary sequence information alone remains elusive for  
12 all but the smallest proteins.

13 In recent years, however, computational researchers have  
14 made substantial improvements to *de novo* structural deter-  
15 mination by leveraging co-evolutionary information contained  
16 within large sequence databases (10–13). Residues that co-  
17 evolve with one another across time may do so as a result of  
18 their spatial proximity with protein structures—i.e. mutations  
19 to an individual residue may be compensated for by subse-  
20 quent mutations to other directly interacting residues (14–16).  
21 By determining an ‘evolutionary coupling’ score for all pairs  
22 of amino acid residues within a structure—and assuming that  
23 the highest-scoring residue-residue pairs are in close spatial  
24 proximity—researchers can constrain the search space of pro-  
25 tein folding methods and accurately predict 3D-structures  
26 (10, 17, 18). Other applications have used evolutionary cou-  
27 pling scores to predict protein binding partners and interfaces

(19, 20), as well as to predict the effect of mutations on pro-  
tein stability and function (21). Many of these applications  
have been further improved through the use of machine  
learning (22–24) and deep neural networks that leverage evolu-  
tionary couplings along-side numerous other protein features  
(25–36).

Despite the progress that has been made in this field—  
spurred by the development of so-called Direct Coupling Anal-  
yses and related methods—there are a number of known limita-  
tions to current methods for computing evolutionary couplings  
(37–40). Perhaps most importantly is a requirement for vast  
numbers of sequence homologs (18). Additionally, the evolu-  
tionary relatedness of sequences and the heterogeneity of  
substitution rates across sites may impose further constraints  
on the overall identifiability of evolutionary couplings. Finally,  
the more distantly related a given homolog is to the target  
structure, the more likely it is that there will be actual struc-  
tural differences between molecules making the designation of  
a protein family based solely on sequence homology potentially  
problematic.

As researchers develop and refine algorithms to better pre-  
dict evolutionary couplings from large multiple sequence align-  
ments, a common work-flow is to benchmark methods against  
known protein structures to determine the accuracy of residue-  
residue contact predictions (41, 42). The large number of  
protein structures that have been solved at atomic resolution  
provides a training data set where intra-molecular contacts  
are *known* (7). However, even the most high-resolution crys-  
tal structures of proteins require researchers to extrapolate

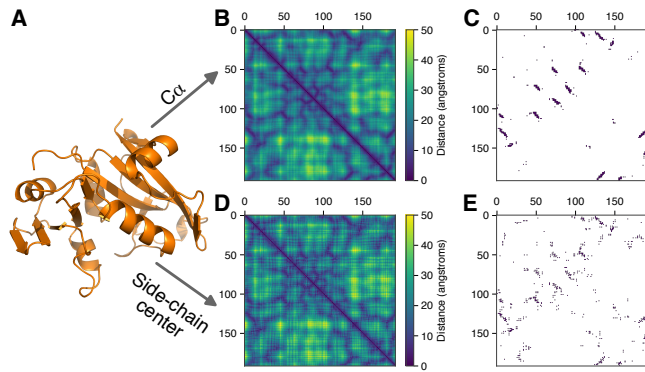
## Significance Statement

Evolutionary couplings between residues within a protein can provide valuable information about protein structures, protein-protein interactions, and the mutability of individual residues. However, the mechanistic factors that determine whether two residues will co-evolve remains unknown. We show that structural proximity by itself is not sufficient for co-evolution to occur between residues. Rather, evolutionary couplings between residues are specifically governed by interactions between side-chain atoms. By contrast, intramolecular contacts between atoms in the protein backbone display only a weak signature of evolutionary coupling. These findings highlight that different types of stabilizing contacts exist within protein structures and that these types have a differential impact on the evolution of protein structures.

Please provide details of author contributions here.

Please declare any conflict of interest here.

<sup>1</sup>To whom correspondence should be addressed. E-mail: adam.hockenberry@utexas.edu



**Fig. 1.** Constructing contact maps from protein structures. (A) An example structure (PDB:1AOE). (B) A symmetrical distance matrix between all pairs of amino acid residues measured from each residues  $C\alpha$  atom. (C) Medium- to long-range contacts ( $> 12$  residues apart) are identified using an  $8\text{\AA}$  cutoff (dark blue). (D) and (E) Same methodology as depicted in (B) and (C), using the geometric center of each residues side-chain as a reference point for measuring distances.

57 from the location of particular atoms and residues to classify  
58 residue-residue ‘bonds’ or ‘contacts’ (43–46). A commonly  
59 used heuristic is to determine that any amino acid residue  
60 that lies within some pre-defined physical distance—frequently  
61  $8\text{\AA}$ —of another amino acid residue is said to be in structural  
62 ‘contact’ (10).

63 Some current applications of evolutionary coupling analyses  
64 use  $C\alpha$  atoms as a reference point for determining residue-  
65 residue distances while others use  $C\beta$  or other complex meth-  
66 ods such as the minimum distance between heavy atoms be-  
67 tween two residues (18). Prior research has shown that the  
68 number of residue-residue contacts identified via side-chain  
69 centers is more closely related to evolutionary rates than simi-  
70 lar metrics derived from  $C\alpha$  atoms (47–49). However, the  
71 consequences of choosing different reference points to deter-  
72 mine the accuracy of modern evolutionary coupling approaches  
73 is unknown.

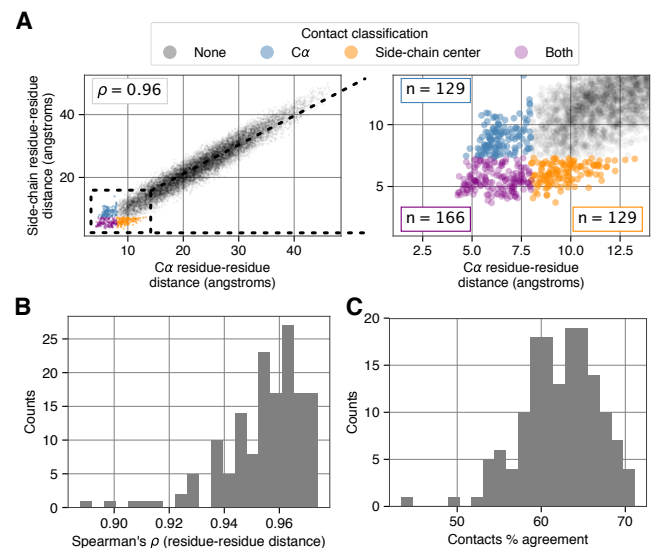
74 Currently, there are no accepted standards in the field for  
75 how to determine a network of residue-residue contacts for a  
76 given protein structure. Further, there has yet to be a system-  
77 atic analysis of whether co-evolutionary signatures are more-  
78 or less-closely related to different types of intra-molecular  
79 contacts that may exist within a protein structure. Here, we  
80 systematically test the the accuracy of several evolutionary  
81 coupling algorithms against true positive contacts defined via  
82  $C\alpha$ ,  $C\beta$ , or side-chain geometric centers. We find that residue-  
83 residue contacts defined according to the distances between  
84 side-chain centers are much more accurately predicted by evo-  
85 lutionary couplings. These results imply that the dominant  
86 epistatic effects resulting in co-evolutionary signatures arise  
87 from side-chain::side-chain interactions. Our findings highlight  
88 the importance of the choice of contact-definition and provide  
89 insight into the constraints governing the evolution of protein  
90 structures.

## 91 Results

92 **Structural contact definitions.** Putatively true interactions be-  
93 tween amino acid residues within a given protein are frequently  
94 derived from the distance between residues in known protein  
95 structures. Figure 1 depicts an example protein structure  
96 (PDB:1AOE) as well as a symmetric matrix depicting all

97 residue-residue contact distances (in angstroms,  $\text{\AA}$ ) defined  
98 according to the distance between the  $C\alpha$  atoms of individual  
99 residues. By convention, we define true contacts as residue-  
100 residue pairs that are less than  $8\text{\AA}$  apart. We further note  
101 that, for most applications, the most structurally *interesting*  
102 contacts are mid- to long-range, which we define here as amino  
103 acid pairs separated by a primary chain distance of at least 12  
104 residues (Fig. 1B,C). We only consider this subset of possible  
105 contacts for the remainder of this manuscript.

106 Many researchers have noted that the distance between  
107 amino acid residues need not be defined by  $C\alpha$  atom-based  
108 distances, and many applications rely on  $C\beta$  atoms (43–46).  
109 A logical question is whether using different reference points  
110 to define contacts matters in practice. To compare the con-  
111 sequences of choosing different reference points, we define all  
112 residue-residue contacts according to the  $8\text{\AA}$ ,  $C\alpha$  atom-based  
113 distance threshold for a given protein. Next, we use the same  
114 absolute number of contacts to determine a comparable dis-  
115 tance threshold (specific to each protein) to use for both  $C\beta$   
116 atom and side-chain center based distances such that an equal  
117 number of putatively true contacts are identified regardless  
118 of the distance metric employed (SI Fig. S1). Although the  
119 distance matrices look similar for an example protein when cal-  
120 culated via  $C\alpha$  atoms or side-chain centers (Fig. 1B compared  
121 to D), the resulting maps of residue-residue contacts show  
122 considerable heterogeneity (Fig. 1C compared to E). More  
123 quantitatively, the set of all residue-residue distances mea-  
124 sured by either  $C\alpha$  atoms,  $C\beta$  atoms, or side-chain centers are  
125 highly correlated with one another (Fig. 2A (left), SI Fig. S2).  
126 However, this strong *overall* correlation obscures important  
127 differences in *contact* definitions which we observe when focus-  
128 ing within the narrow region where direct amino acid residue  
129 contacts are defined (Fig. 2A (right), SI Fig. S2). For 1AOE,



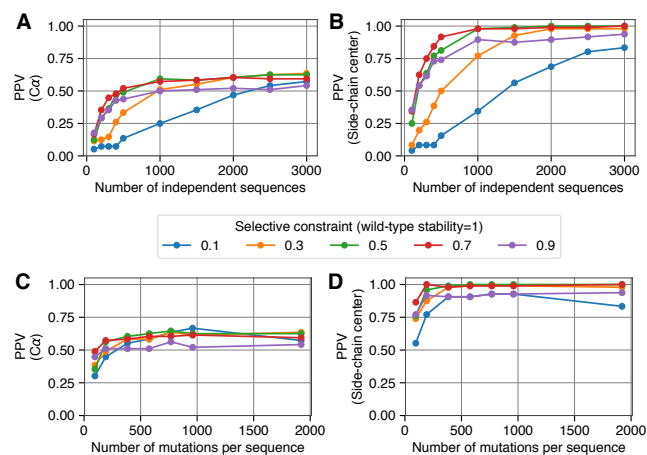
**Fig. 2.** Relationship between different contact identification methods. (A) Correlation between residue-residue distances in PDB:1AOE measured according to  $C\alpha$  atoms and side-chain centers (left). A zoomed in view (right) highlights variably defined residue-residue contacts indicated by the various colors. (B) Distribution of Spearman's correlation coefficient values ( $\rho$ ) between residue-residue distances for 150 different proteins. (C) Distribution of the percent agreement for contact definitions for the same set of proteins. (SI Fig.S2 shows a comparable comparison between  $C\beta$  and side-chain center-based distances.)

130 we identified a total of 295 contacts according to the 8Å  $C\alpha$   
 131 atom-based distance threshold. Of the shortest 295 contact  
 132 distances identified via side-chain centers—corresponding to  
 133 a distance threshold of 7.33Å for this protein—the percent of  
 134 residue-residue pairs that appear in both definitions is only  
 135 56% (78% for  $C\alpha$  compared to  $C\beta$  and 76% for  $C\beta$  compared  
 136 to side-chain centers).

137 To assess the generality of these findings, we applied this  
 138 analysis to a commonly used benchmark set of 150 proteins  
 139 (23, 34, 38). Across all of these proteins, we observed a median  
 140 correlation of 0.97 between residue-residue distances calculated  
 141 via  $C\alpha$  atoms and side-chain centers (Fig. 2B) and a median  
 142 overlap of 63% between contacts defined via  $C\alpha$  and side-chain  
 143 centers (Fig. 2C).  $C\alpha$ - and  $C\beta$ -based definitions, as well as  
 144  $C\beta$ - and side-chain center-based definitions, both had median  
 145 overlaps of 78% (SI Fig. S2). Together, these results highlight  
 146 that true contacts vary substantially according to the reference  
 147 point used to measure residue-residue distances.

148 **Simulation analyses.** While the previous analysis of empirical  
 149 structures shows that the choice of reference point has impor-  
 150 tant consequences for true contact identification, it is not clear  
 151 which of the different methods is more biologically "correct"  
 152 or practically meaningful. We thus turned our attention to a  
 153 simplified biophysical system to test the ability of evolution-  
 154 ary coupling analyses to recover intramolecular contacts. We  
 155 used the ROSETTA modeling software (50–52) to perform  
 156 all-atom evolutionary simulations of the evolutionary process  
 157 (53, 54) while selecting for the maintenance of protein stability  
 158 (expressed as a fraction of the initial PDB model stability).  
 159 We simulated thousands of independent evolutionary trajec-  
 160 tories, and used the resulting amino acid sequences from these  
 161 simulations to calculate evolutionary couplings. We used 3 sep-  
 162 arate algorithms to calculate evolutionary couplings, but the  
 163 main text results depict predictions using CCMpred. Within  
 164 this defined system, we are able to remove the constraints of  
 165 phylogenetic biases, limited data availability, homopolymeriza-  
 166 tion, and changes in evolutionary pressures over time between  
 167 species—all of which partially limit the power of algorithms to  
 168 detect true evolutionary couplings in real data (55).

169 We continued to use 1AOE as an example protein and varied  
 170 several parameters of our simulation to ensure robust results.  
 171 We defined true positive residue-residue contacts according  
 172 to the original PDB structure using residue-residue distances  
 173 calculated between different quantities for comparison ( $C\alpha$ ,  $C\beta$ ,  
 174 and side-chain center). To assess the accuracy of evolutionary  
 175 couplings, we determined the positive predictive value (PPV)  
 176 of the top  $L/2$  couplings—where  $L$  is the primary chain length  
 177 of the protein under investigation. For  $C\alpha$ -based contact  
 178 definitions, we found that the PPV increases rapidly according  
 179 to the number of independent sequences that we simulated  
 180 and consequently used as input for evolutionary coupling  
 181 analyses (Fig. 3A). In each case, we ran these simulations  
 182 until we accepted mutations totaling 10x the length of the  
 183 protein sequence. However, regardless of the selection strength  
 184 that we imposed on the sequence evolution (colored lines in  
 185 Fig. 3A), PPV values plateaued at a value  $\ll 1$  indicating  
 186 that evolutionary couplings were failing to accurately capture  
 187 protein contacts. By contrast, when we analyzed the same  
 188 evolutionary coupling values but used side-chain center-based  
 189 distances to define contacts we observed that PPV values  
 190 approached 1 (representing perfect prediction accuracy for



191 **Fig. 3.** Comparing simulation-derived evolutionary couplings to different contact  
 192 definitions. (A) For each of 5 separate selection strengths (colored lines), we ran  
 193 simulations until a number of mutations totaling 10 times the length of the protein were  
 194 accumulated per replicate. We varied the number of independent replicate sequences  
 195 (x-axis) that were used as input for evolutionary coupling analysis, and found that  
 196 couplings fail to fully recover  $C\alpha$  defined contacts for PDB:1AOE. (B) By contrast,  
 197 contacts defined via side-chain centers are near-perfectly recovered across a range  
 198 of simulation parameters. (C) and (D) Similar to parts (A) and (B), but along the x-axis  
 199 we now show results from simulations where a different number of accepted mutations  
 200 were accumulated per sequence. We fixed the number of replicate sequences that  
 201 were simulated—and used for evolutionary coupling analysis—at 3,000 for each of these  
 202 data points. (Results comparing  $C\beta$  and side-chain center-based contact definitions,  
 203 can be found in SI Fig. S3.)

204 this subset of couplings) in almost all cases.

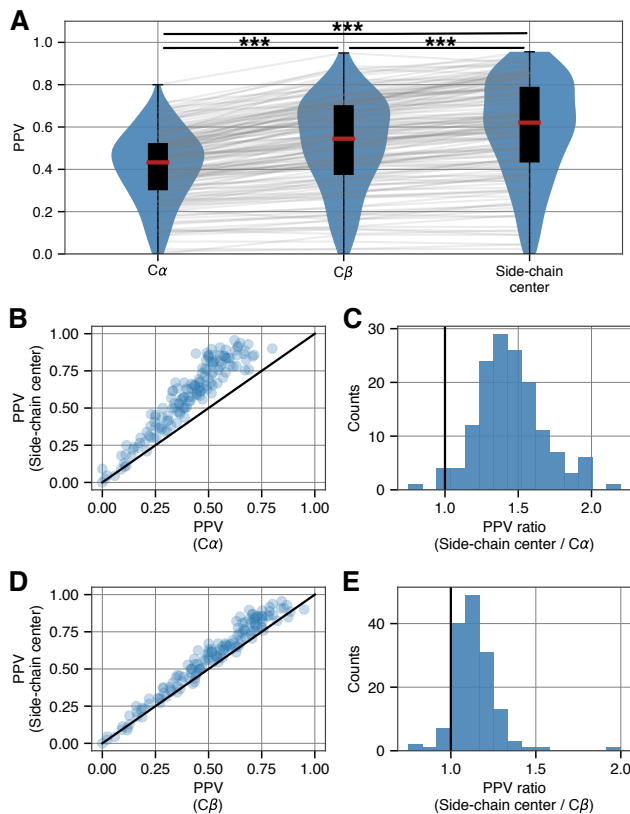
205 We additionally explored how the number of mutations  
 206 accumulated per sequence affected the ability of evolutionary  
 207 coupling algorithms to recover intramolecular contacts. We  
 208 fixed the number of replicate sequences at 3000, and observed  
 209 that PPV values showed minimal variation according to the  
 210 number of accepted mutations per sequence (Fig. 3C). As be-  
 211 fore, however, prediction accuracies were substantially higher  
 212 when we defined true contacts according to side-chain center  
 213 distances (Fig. 3D).

214 These simulation results highlight that—across numerous  
 215 parameter combinations—the top  $L/2$  evolutionary couplings  
 216 were all true positive intramolecular contacts as long as true  
 217 positives were defined according to side-chain centers and not  
 218  $C\alpha$  carbons. Additionally, we depict  $C\alpha$  and side-chain center  
 219 based methods here because they represent extreme ends of  
 220 the spectrum from backbone to side-chains.  $C\beta$ -based contact  
 221 definitions had intermediate accuracy, plateauing at higher  
 222 values than  $C\alpha$  but lower than side-chain center definitions  
 (SI Fig. S3).

223 **Empirical analyses.** To see how evolutionary couplings com-  
 224 pare to different definitions of true residue-residue contacts in  
 225 empirical data, we used PHMMER to identify sequence ho-  
 226 mologs for each of the 150 proteins (see Materials and Methods  
 227 for details). We assessed the relationship between evolutionary  
 228 couplings and structural contacts for all proteins by calculat-  
 229 ing the positive predictive value (PPV) of the highest  $L/2$   
 230 couplings.

231 As expected, the PPV between the top  $L/2$  evolutionary  
 232 couplings and  $C\alpha$ -based contacts varied substantially across  
 233 the 150 structures. This variation may result from a num-  
 234 ber of different effects, and we observed a clear and expected





**Fig. 4.** Accuracy of evolutionary couplings derived from empirical alignments. (A) For a diverse set of 150 proteins, the PPV of the top  $L/2$  evolutionary coupling scores—derived from empirical sequence alignments—is progressively higher when intramolecular contacts are defined according to  $C\alpha$  atoms,  $C\beta$  atoms, and side-chain centers. (\*\*\*) indicates  $p < 10^{-20}$ , Wilcoxon signed-rank test) (B) Scatter plot of PPVs for each protein according to  $C\alpha$  and side-chain center-based contact identification methods. (C) Histogram of the ratios from the data in (B) indicate that the median percent increase in accuracy is 43%. (D) and (E) As in (B) and (C), comparing  $C\beta$  and side-chain center-based contact identification methods. Results show a median percent increase in contact identification accuracy of 13%. (Results for other evolutionary coupling algorithm implementations can be found in SI Figs. S5 & S6.)

shown that PPVs are substantially higher when using side-chain-based distances to identify true positive intramolecular contacts compared with either  $C\alpha$  or  $C\beta$ -based distances. To look more specifically at why these differences were so pronounced, we decided to investigate the orientation of residue-residue pairs identified by the various criteria. At a simple level, any two residues can be in structural contact across a number of orientations of their respective side-chains: i) both residue's side-chains may point towards one another with the energetic interactions occurring through side-chain atoms, ii) one residue's side-chain may point towards the other residue while that residue's side-chain points away, or iii) both residue's side-chains may point away from one another with energetic interactions occurring between the respective amino acid backbones (Fig. 5A). As expected, when we look only at residue-residue pairs that are defined as contacts via different reference points, we see that side-chain based contact definitions strongly enrich for cases where both side-chains point towards one another in an example protein (Fig. 5B).

Across all 150 proteins in our dataset, we calculated the fraction of *all* residue-residue pairs (regardless of whether they are putative contacts, but subject to the same primary chain distance constraints applied throughout this manuscript) where both side-chains point towards one another and found it to be relatively small (Fig. 5C, “All pairs”). However, this fraction increases progressively when we limit our analysis to the subset of residue-residue pairs identified as true contacts for each protein according to  $C\alpha$ ,  $C\beta$  and side-chain centers—illustrating that the trend observed in (Fig. 5B) applies broadly. If instead we only look at the top-ranked evolutionary couplings (ignoring whether or not the residue-residue pairs are putatively true structural contacts), we observe that a large fraction of the strongest identified evolutionary couplings are between residues that point towards one another in the reference protein structure. Additionally, this fraction is highest for the most highly ranked evolutionary couplings and is substantially higher than the proportion identified by  $C\alpha$ -based distances.

To further illustrate this point, we turned to an alternative method for determining intramolecular contacts that we have not yet systematically explored: determining structural contacts based on the minimum distance between any two heavy atoms for each residue-residue pair. We implemented two versions of this algorithm, determining the minimum distance between: i) all heavy atoms within residues and ii) side-chain heavy atoms only. In each case, and to facilitate comparison between methods, we again selected the shortest  $X$  distances as contacts where  $X$  is the number of contacts identified for each protein via the 8Å distance threshold using  $C\alpha$ . For the set of 150 proteins, the resulting PPVs were significantly higher when contacts were defined *only* according to side-chain atoms as opposed to the complete set of backbone and side-chain atoms (SI Fig. S7). Furthermore, PPVs calculated via side-chain center distances were statistically indistinguishable from PPVs derived from the minimum distance between all heavy atoms within side-chains.

Taken together, our analysis of side-chain orientations and our analysis of contacts identified via minimum atomic distances both highlight that evolutionary couplings frequently occur between residues whose side-chains point towards one another.  $C\alpha$ -(and to a lesser extent  $C\beta$ -) based contact defini-

correlation between PPV values and the number of available homologous sequences used to determine evolutionary couplings (SI Fig. S4). Despite the variability in prediction accuracy between proteins, we observed systematic variability in the PPV according to which metric was used to identify true positive contacts (Fig. 4A). When compared with  $C\alpha$ -based distances, residue-residue distances measured according to  $C\beta$  atoms resulted in significantly higher PPVs, and side-chain-based contact distances resulted in even further improvements. Further, the magnitude was substantial: across all 150 proteins the median percent increase in PPV between  $C\alpha$ - and side-chain center-based contact identification methods was 43% (Fig. 4B,C). Even between the more similar  $C\beta$ - and side-chain center-based methods, the median percent increase in accuracy was 13% (Fig. 4D,E). Both comparisons were highly significant and persisted across the entire range of PPVs represented within our dataset (Fig. 4B,D). Additionally, these results were highly consistent across different evolutionary coupling algorithms (SI Figs. S5 & S6).

**Side-chain orientation and evolutionary couplings.** Using the exact same evolutionary couplings, the previous analyses have

305 tions classify a smaller number of contacts in this orientation,  
 306 and including backbone atoms in minimum-distance based  
 307 contact identification methods actually decreases the accuracy  
 308 of contact predictions based on evolutionary couplings.

## 309 Discussion

310 The co-evolutionary patterns of amino acid substitutions can  
 311 provide important information about protein structures. There  
 312 are a number of competing methods currently employed by  
 313 different researchers to detect evolutionary couplings between  
 314 residues, and the ability to recover true residue-residue con-  
 315 tacts has been the primary metric used to assess performance  
 316 of various methods. However, true structural contacts are  
 317 ill-defined and variability in contact definitions can prohibit  
 318 comparison between the efficacy of different methods, as well  
 319 as obscure the biological interpretation of evolutionary con-  
 320 straints. We show here that evolutionary couplings are signifi-  
 321 cantly more accurate at detecting true residue-residue contacts  
 322 based off of side-chain center distances. Critically, these find-  
 323 ings provide important biological insight protein evolution and  
 324 epistatic interactions between residues. Our model posits that  
 325 although different types of interactions between amino acid  
 326 residues may stabilize protein structures, evolutionary cou-  
 327 plings predominantly consist of residues whose contact occurs  
 328 via interactions between the side-chain atoms of both residues.

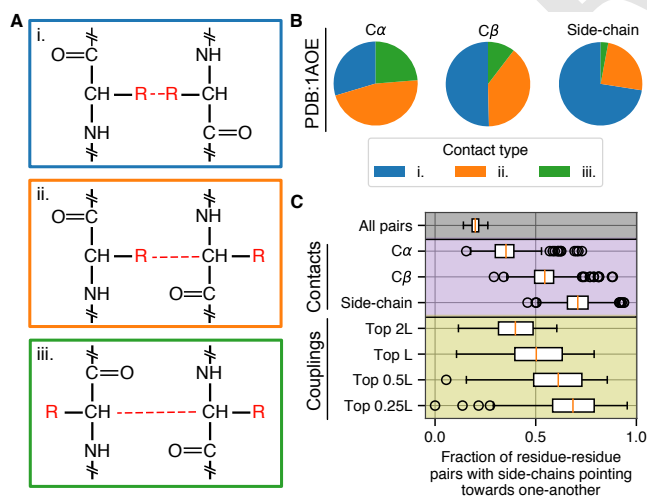
329 Evolutionary couplings are themselves important, but more  
 330 recently these values have been used as input into more ad-  
 331 vanced machine learning and neural network-based algorithms  
 332 that supplement evolutionary couplings with a variety of other  
 333 information to predict intramolecular contacts. However, our  
 334 analysis suggests that there may be biases in training data-  
 335 essential to supervised learning techniques—based on how in-

336 tramolecular contacts are defined; the definition of putatively  
 337 "true" structural contacts relies on the method used to cal-  
 338 culate residue-residue distances. We show that evolutionary  
 339 couplings more accurately predict side-chain center-based con-  
 340 tacts, and the strongest evolutionary couplings are consistently  
 341 enriched for residue-residue pairs where the side-chains are  
 342 oriented towards one-another. We speculate that accuracy  
 343 of supervised algorithms may be improved with more prop-  
 344 erly labeled training data that corresponds with these known  
 345 biophysical constraints. Alternatively, supervised learning al-  
 346 gorithms may be able to achieve even greater improvements in  
 347 accuracy by separating different *types* of residue-residue con-  
 348 tacts according to their atomic interactions, training separate  
 349 models to detect each type, and integrating the results.

350 A number of issues constrain the maximal accuracy that  
 351 can be expected from using evolutionary couplings alone to  
 352 predict contacts. Anischenko *et al.* (2017) illustrated that  
 353 many so-called false positive signals resulting from evolu-  
 354 tionary coupling analyses arise from repeat proteins, homo-  
 355 oligomerization, and structural variation within protein fam-  
 356 ilies (55). Here, we show that another source of false posi-  
 357 tive signals may simply be ill-defined true positive contacts.  
 358 Without changing anything about the way evolutionary cou-  
 359 plings are calculated, we show their accuracy at predicting  
 360 intramolecular contacts is progressively higher for different  
 361 contact-definitions. Further, the magnitude of this difference  
 362 is not trivial: across a diverse set of proteins we show that  
 363 side-chain center based contacts are predicted with a median  
 364 of 43% and 13% higher accuracy than comparable  $C\alpha$  and  
 365  $C\beta$  based contacts. Thus a substantial number of *false posi-*  
 366 *tive* predictions made by evolutionary coupling analyses may  
 367 simply be due to the false classification of true positives.

368 While improving contact identification methods is an impor-  
 369 tant practical result, our findings improve our understanding  
 370 of protein evolution by showing that side-chain interactions  
 371 are more important for governing epistasis between amino  
 372 acid residues within individual protein structures. Although  
 373 the overall structural geometry of a protein is dictated by the  
 374 shape of the protein backbone, consideration of side-chains  
 375 is critical for maintaining this geometry and determining the  
 376 co-evolutionary dynamics of substitutions. Our findings do not  
 377 suggest that intra-molecular contacts between the backbone  
 378 atoms of residues are not important for folding or stabilizing  
 379 protein structures. Rather, our results suggest that contacts  
 380 between backbone atoms are not likely to be detected by evo-  
 381 lutionary coupling analyses and imply that epistasis between  
 382 residues is largely governed by whether side-chain atoms are  
 383 in direct contact.

384 Direct coupling analyses and related methods have sig-  
 385 nificantly improved our ability to leverage the experiment  
 386 of natural sequence evolution for the purpose of predicting  
 387 important properties of proteins. While these methods contin-  
 388 ue to find novel applications, they are beginning to provide  
 389 mechanistic insight into the evolutionary process (56). Ulti-  
 390 mately, it may even be possible to incorporate more realistic  
 391 pair-wise interactions into models of sequence evolution and  
 392 inference, which are almost exclusively site-independent. Fur-  
 393 ther technical improvements, such as explicitly accounting  
 394 for the phylogenetic relatedness of sequences, may allow for  
 395 even more accurate inference of evolutionary couplings and  
 396 consequently insight into biological mechanisms.



**Fig. 5.** Different types of residue-residue interactions are possible. (A) Two interacting residues may interact via: each residue's side chain atoms (type i), the side-chain of one residue and the backbone of the other (type ii), or the backbone atoms of each residue (type iii). (B) For intramolecular contacts identified in PDB:1AOE, the relative proportion of different interaction types varies according to contact identification method. Residue-residue contacts defined via side-chain centers are enriched in type i interactions (blue). (C) For 150 proteins, the fraction of residue-residue pairs where the side-chains point towards one-another is highest in contacts defined via side-chain centers (purple). Further, the top ranked evolutionary couplings (regardless of whether they are defined as contacts) are progressively enriched in residue-residue pairs where the side-chains point towards one another (yellow).

## 397 Materials and Methods

398 **Dataset compilation and processing.** We downloaded the so-called  
399 PSICOV dataset of 150 proteins that have been extensively studied  
400 (23, 34, 38). We processed each starting “.PDB” file to select a single  
401 chain, ensure a consistent numbering of residues (1...*n*), test for  
402 unknown or non-standard residues, select the most likely state for all  
403 disordered sequence atoms, and remove all extraneous information  
404 (including “HETATM” lines). Next, to ensure that all residues  
405 were represented in full and repair those that were not, we used  
406 PYROSETTA to read in the “.PDB” files using the ‘pose\_from\_pdb’  
407 function and wrote the output as our final clean structure.

408 **Determining structural contacts and contact-types.** From each  
409 cleaned “.PDB” file, we calculated residue-residue distance matrices  
410 using custom python scripts (the euclidean distance from  
411 3-dimensional atomic coordinates). All residues contain a *C $\alpha$*  atom  
412 so this calculation was straightforward. For *C $\beta$*  calculations, we  
413 used the *C $\beta$*  atom of all residues except glycine, where we continued  
414 to use the *C $\alpha$*  atom. For side-chain center calculations, we calculated  
415 the geometric center of each residue based on the coordinates of  
416 all non-backbone heavy atoms. This calculation included *C $\beta$*   
417 atoms but excluded *C $\alpha$*  atoms for all amino acids except glycine,  
418 where we continued to use *C $\alpha$*  as the reference point.

419 To calculate minimum atomic distances between two residues,  
420 we calculated all pairwise euclidean distances between heavy atoms  
421 and selected the minimum distance. In extending this analysis  
422 to only consider side-chain atoms, we continued to consider *C $\beta$*   
423 atoms as part of the side-chain but not *C $\alpha$* . Again, we relaxed  
424 this restriction for glycine and included *C $\alpha$*  as a side-chain atom to  
425 permit calculations.

426 For all methods, contacts were assessed by first removing all  
427 residue-residue pairs where the two residues were shorter than 12  
428 amino acids apart in primary chain distance. Contacts were determined  
429 throughout this manuscript for each structure according to  
430 an 8Å cutoff between *C $\alpha$*  atoms. Since accuracy values are partially  
431 dependent on the number of true positives that are called,  
432 we maintained a constant number of true positive contact classifications  
433 throughout to facilitate comparison between methods. For each contact  
434 definition (*C $\beta$* , side-chain center, minimum atomic distances), we  
435 selected *n* residue-residue pairs with the shortest distances where  
436 *n* is the number of contacts defined according to the aforementioned  
437 *C $\alpha$* -based method.

438 To classify residue-residue pairs (*a* and *b*) via their side-chain  
439 orientations, we chose a residue (*a*) and drew two vectors: i) from  
440 the *C $\alpha$*  atom coordinate to the side-chain center for that residue and  
441 ii) from the *C $\alpha$*  atom coordinate to the *C $\alpha$*  atom coordinate for the  
442 other residue in question (*b*). If the angle between these two vectors  
443 was less than  $\pi/2$  radians, the side-chain of residue *a* was said to  
444 point towards residue *b*. To determine the residues classification as  
445 in Fig.5A, we next repeated the calculation using residue *b* as the  
446 reference and classified the residue-residue pair accordingly.

447 **Evolutionary coupling algorithms.** For each of the 150 proteins in our  
448 dataset, we followed a principled method to retrieve homologous  
449 sequences. We first extracted the primary amino acid sequence from  
450 from the “.PDB” file. We next used PHMMER to search through  
451 progressively larger databases in order to retrieve up to 10,000  
452 homologous sequences. To do so, we downloaded local versions  
453 of the rp15, rp35, rp55, and rp75 databases. We first searched  
454 the smallest, least redundant, database for each sequence using an  
455 E-value threshold of 0.0001. For any sequence with greater than  
456 10,000 hits we stopped and selected the top scoring 10,000 hits  
457 for further analysis. For sequences with fewer than 10,000 hits  
458 we moved to the next largest database and repeated the process.  
459 Finally for the small number of sequences for which we did not  
460 accumulate at least 1,000 sequences in the largest database (rp75),  
461 we used the online version of PHMMER to search the UniprotKB  
462 database and downloaded the maximum results.

463 For each protein, we next aligned the hits along-side the reference  
464 sequence using MAFFT (L-INS-i method with default parameters).  
465 Next, we cleaned these results to remove all columns that were  
466 gapped in the reference (“.PDB”) sequence. All other columns and  
467 sequences in the sequence alignments were retained regardless of  
468 gap coverage.

Using these alignments, we next calculated evolutionary couplings  
between residue-residue pairs. All results in the main manuscript  
are displayed using CCMpred with default parameters (0.8 local  
sequence re-weighting threshold, 0.2 pairwise regularization coefficients,  
average product correction). We additionally used the ‘plmc’ method  
from the EVcouplings framework with default parameters (no average  
product correction) and PSICOV (default parameters excepting: “-z  
50 -r 0.001”) to ensure the robustness of our findings.

Except where otherwise noted (Fig. 5), main text results (Fig. 3,  
Fig. 4) were calculated using the top *L/2* couplings for each protein  
where *L* is the length of the reference amino acid sequence. Positive  
Predictive Value (PPV) is calculated as the number of classified  
contacts among these top couplings divided by the total number of  
top couplings considered.

**Evolutionary simulations.** For the example protein used throughout  
the text (PDB:1AOE) we performed mutation accumulation simulations  
using PYROSETTA. We first read in the “.PDB” structure (with  
disulfide bonds turned off), and minimized it so as to optimize  
thermodynamic stability by rotamer selection and backbone movements.  
We next fixed the backbone, and implemented an expedited  
evolutionary algorithm to select amino acid point mutations (no  
insertions or deletions were allowed) according to their predicted  
impact on structural stability. At each step, we selected a random  
amino acid position, and attempted to mutate it randomly to one  
of the remaining 19 amino acids. We re-packed the structure within  
a 12Å radius of the mutation and determined whether or not to  
accept it based off of the resulting change in structural stability.  
Mutations which either did not alter or which increased stability  
(i.e. resulted in a decreased  $\Delta G$ ) were accepted. Mutations that  
decreased stability were accepted with a probability proportional  
to their  $\Delta\Delta G$  as in Teufel and Wilke (2017). At the end of the  
evolutionary process, the resulting amino acid sequence was stored  
for future analysis.

We performed thousands of independent replicates of this expedited  
evolutionary process where we altered the number of accepted  
mutations that we accumulated, the number of replicate evolutionary  
experiments that we performed, and the fraction of the initial  
wild-type stability value that we used for our selection criteria.  
Collections of the resulting sequences were analyzed via evolutionary  
coupling algorithms in the same manner as empirical sequences,  
with no need for sequence alignment.

**Data access.** All code and data are currently being compiled and  
edited. They will be made freely available at the following link:  
<https://github.com/adamhockenberry/side-chain-couplings> which  
will itself evolve throughout the submission process to reflect final  
manuscript analyses.

**ACKNOWLEDGMENTS.** The authors acknowledge valuable feedback  
and support from members of the Wilke lab.

1. Anfinsen CB (1973) Principles that Govern the Folding of Protein Chains. *Science* 181(4096):223–230.
2. Sadowski MI, Jones DT (2009) The sequence-structure relationship and protein function prediction. *Current Opinion in Structural Biology* 19(3):357–362.
3. Marks DS, Hopf TA, Sander C (2012) Protein structure prediction from sequence variation. *Nature Biotechnology* 30(11):1072–1080.
4. Liao M, Cao E, Julius D, Cheng Y (2013) Structure of the TRPV1 ion channel determined by electron cryo-microscopy. *Nature* 504(7478):107–112.
5. Amunts A, et al. (2014) Structure of the Yeast Mitochondrial Large Ribosomal Subunit. *Science* 343:1485–1489.
6. Punjani A, Rubinstein JL, Fleet DJ, Brubaker MA (2017) CryoSPARC: Algorithms for rapid unsupervised cryo-EM structure determination. *Nature Methods* 14(3):290–296.
7. Rose PW, et al. (2017) The RCSB protein data bank: Integrative view of protein, gene and 3D structural information. *Nucleic Acids Research* 45:D271–D281.
8. Söding J (2005) Protein homology detection by HMM-HMM comparison. *Bioinformatics* 21(7):951–960.
9. Biasini M, et al. (2014) SWISS-MODEL: Modelling protein tertiary and quaternary structure using evolutionary information. *Nucleic Acids Research* 42:252–258.
10. Marks DS, et al. (2011) Protein 3D structure computed from evolutionary sequence variation. *PLoS ONE* 6(12).
11. Morcos F, et al. (2011) Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proceedings of the National Academy of Sciences* 108(49):E1293–E1301.



- 541 12. Sulkowska JI, Morcos F, Weigt M, Hwa T, Onuchic JN (2012) Genomics-aided structure prediction. *Proceedings of the National Academy of Sciences* 109(26):10340–10345. 625
- 542 13. Kamisetty H, Ovchinnikov S, Baker D (2013) Assessing the utility of coevolution-based 626
- 543 residue-residue contact predictions in a sequence- and structure-rich era. *Proceedings of* 627
- 544 *the National Academy of Sciences* 110(39):15674–15679. 628
- 545 14. Gobel U, Sander C, Schneider R, Valencia A (1994) Correlated Mutations and Residue Con- 629
- 546 tacts in Proteins. *Proteins* 18:309–317. 630
- 547 15. Shindyalov IN, Kolchanov NA, Sander C (1994) Can three-dimensional contacts in protein 631
- 548 structures be predicted by analysis of correlated mutations? *Protein Engineering, Design* 632
- 549 *and Selection* 7(3):349–358. 633
- 550 16. Weigt M, White RA, Szurmant H, Hoch JA, Hwa T (2009) Identification of direct residue 634
- 551 contacts in protein-protein interaction by message passing. *Proceedings of the National* 635
- 552 *Academy of Sciences* 106(1):67–72. 636
- 553 17. Hopf TA, et al. (2012) Three-dimensional structures of membrane proteins from genomic 637
- 554 sequencing. *Cell* 149(7):1607–1621. 638
- 555 18. Ovchinnikov S, et al. (2017) Protein structure determination using metagenome sequence 639
- 556 data. *Science* 355(6322):294–298. 640
- 557 19. Hopf TA, et al. (2014) Sequence co-evolution gives 3D contacts and structures of protein 641
- 558 complexes. *eLife* 3:1–45. 642
- 559 20. Ovchinnikov S, Kamisetty H, Baker D (2014) Robust and accurate prediction of residue- 643
- 560 residue interactions across protein interfaces using evolutionary information. *eLife* pp. 1–21.
- 561 21. Hopf TA, et al. (2017) Mutation effects predicted from sequence co-variation. *Nature Biotech-* 644
- 562 *nology* 35(2):128–135.
- 563 22. Cheng J, Baldi P (2007) Improved residue contact prediction using support vector machines 645
- 564 and a large feature set. *BMC Bioinformatics* 8(1).
- 565 23. Jones DT, Singh T, Kosciolk T, Tetchner S (2015) MetaPSICOV: Combining coevolution 646
- 566 methods for accurate prediction of contacts and long range hydrogen bonding in proteins. 647
- 567 *Bioinformatics* 31(7):999–1006.
- 568 24. Michel M, Skwark MJ, Hurtado DM, Ekeberg M, Elofsson A (2017) Predicting accurate con- 648
- 569 tacts in thousands of Pfam domain families using PconsC3. *Bioinformatics* 33(18):2859– 649
- 570 2866.
- 571 25. Tegge AN, Wang Z, Eickholt J, Cheng J (2009) NNcon: Improved protein contact map predic- 650
- 572 tion using 2D-recursive neural networks. *Nucleic Acids Research* 37:515–518.
- 573 26. Lena PD, Nagata K, Baldi P (2012) Deep architectures for protein contact map prediction. 651
- 574 *Bioinformatics* 28(19):2449–2457.
- 575 27. Xiong D, Zeng J, Gong H (2017) A deep learning framework for improving long-range residue- 652
- 576 residue contact prediction using a hierarchical strategy. *Bioinformatics* 33(17):2675–2683.
- 577 28. Stahl K, Schneider M, Brock O (2017) EPSILON-CP: Using deep learning to combine infor- 653
- 578 mation from multiple sources for protein contact prediction. *BMC Bioinformatics* 18(1):1–11.
- 579 29. He B, Mortuza SM, Wang Y, Shen HB, Zhang Y (2017) NeBcon: Protein contact map predic- 654
- 580 tion using neural network training coupled with naive Bayes classifiers. *Bioinformatics* 655
- 581 33(15):2296–2306.
- 582 30. Wang S, Sun S, Li Z, Zhang R, Xu J (2017) *Accurate De Novo Prediction of Protein Contact* 656
- 583 *Map by Ultra-Deep Learning Model*. Vol. 13.
- 584 31. Riesselman AJ, Ingraham JB, Marks DS (2018) Deep generative models of genetic variation 657
- 585 capture mutation effects. *Nature Methods*.
- 586 32. Liu Y, Palmedo P, Ye Q, Berger B, Peng J (2018) Enhancing Evolutionary Couplings with 658
- 587 Deep Convolutional Neural Networks. *Cell Systems* 6(1):65–74.e3.
- 588 33. Wozniak PP, Pelc J, Skrzynecki M, Vriend G, Kotulska M (2018) Bio-knowledge based filters 659
- 589 improve residue-residue contact prediction accuracy. *Bioinformatics* (May):1–9.
- 590 34. Jones DT, Kandathil SM (2018) High precision in protein contact prediction using fully convo- 660
- 591 lutional neural networks and minimal sequence features. *Bioinformatics* (April):1–8.
- 592 35. Adhikari B, Hou J, Cheng J (2018) DNCON2: Improved protein contact prediction using two- 661
- 593 level deep convolutional neural networks. *Bioinformatics* 34(9):1466–1472.
- 594 36. Hanson J, Paliwal K, Litfin T, Yang Y, Zhou Y (2018) Accurate Prediction of Protein Contact 662
- 595 Maps by Coupling Residual Two-Dimensional Bidirectional Long Short-Term Memory with 663
- 596 Convolutional Neural Networks. *Bioinformatics* (June):bt481–bt481.
- 597 37. Burger L, Van Nimwegen E (2010) Disentangling direct from indirect co-evolution of residues 664
- 598 in protein alignments. *PLoS Computational Biology* 6(1).
- 599 38. Jones DT, Buchan DW, Cozzetto D, Pontil M (2012) PSICOV: Precise structural contact predic- 665
- 600 tion using sparse inverse covariance estimation on large multiple sequence alignments. 666
- 601 *Bioinformatics* 28(2):184–190.
- 602 39. Ekeberg M, Lövkqvist C, Lan Y, Weigt M, Aurell E (2013) Improved contact prediction in pro- 667
- 603 teins: Using pseudolikelihoods to infer Potts models. *Physical Review E* 87(1):1–16.
- 604 40. Seemayer S, Gruber M, Söding J (2014) CCMpred - Fast and precise prediction of protein 668
- 605 residue-residue contacts from correlated mutations. *Bioinformatics* 30(21):3128–3130.
- 606 41. Schaarschmidt J, Monastyrskyy B, Kryshchak A, Bonvin AM (2018) Assessment of con- 669
- 607 tact predictions in CASP12: Co-evolution and deep learning coming of age. *Proteins: Struc-* 670
- 608 *ture, Function and Bioinformatics* 86(October 2017):51–66.
- 609 42. Wang S, Sun S, Xu J (2018) Analysis of deep learning methods for blind protein contact 671
- 610 prediction in CASP12. *Proteins: Structure, Function and Bioinformatics* 86(August 2017):67– 672
- 611 77.
- 612 43. Seeliger D, de Groot BL (2007) Atomic contacts in protein structures. A detailed analysis of 673
- 613 atomic radii, packing, and overlaps. *Proteins* 68:595–601.
- 614 44. SathyaPriya R, Duarte JM, Stehr H, Filippis I, Lappe M (2009) Defining an essence of struc- 674
- 615 ture determining residue contacts in proteins. *PLoS Computational Biology* 5(12).
- 616 45. Duarte JM, SathyaPriya R, Stehr H, Filippis I, Lappe M (2010) Optimal contact definition for 675
- 617 reconstruction of Contact Maps. *BMC Bioinformatics* 11.
- 618 46. Yuan C, Chen H, Kihara D (2012) Effective inter-residue contact definitions for accurate pro- 676
- 619 tein fold recognition. *BMC Bioinformatics* 13(1).
- 620 47. Lin CP, et al. (2008) Deriving protein dynamical properties from weighted protein contact 677
- 621 number. *Proteins* 72(3):929–935.
- 622 48. Marcos ML, Echave J (2015) Too packed to change: side-chain packing and site-specific 678
- 623 substitution rates in protein evolution. *PeerJ* 3:e911.
- 624