



Ensembles from Ordered and Disordered Proteins Reveal Similar Structural Constraints during Evolution

Q1
4
5
6
7
Q4 Q3
8 **Julia Marchetti¹, Alexander Miguel Monzon^{1,2}, Silvio C.E. Tosatto²,
9 Gustavo Parisi¹ and María Silvina Fornasari¹**

10 **1 - Departamento de Ciencia y Tecnología, CONICET, Universidad Nacional de Quilmes, Roque Sáenz Peña 352, B1876BXD,**
11 **Bernal, Provincia de Buenos Aires, Argentina**

12 **2 - Department of Biomedical Sciences, University of Padua, Padua, Italy**

13
14
15 **Correspondence to Gustavo Parisi: gusparisi@gmail.com**

16 **<https://doi.org/10.1016/j.jmb.2019.01.031>**

17 **Edited by Monika Fuxreiter**

18 Abstract

20 The conformations accessible to proteins are determined by the inter-residue interactions between amino acid
21 residues. During evolution, structural constraints that are required for protein function providing biologically
22 relevant information can exist. Here, we studied the proportion of sites evolving under structural constraints
23 in two very different types of ensembles, those coming from ordered and disordered proteins. Using a
24 structurally constrained model of protein evolution, we found that both types of ensembles show comparable,
25 near 40%, number of positions evolving under structural constraints. Among these sites, ~68% are in
26 disordered regions and ~57% of them show long-range inter-residue contacts. Also, we found that disordered
27 ensembles are redundant in reference to their structurally constrained evolutionary information and could be
28 described on average with ~11 conformers. Despite the different complexity of the studied ensembles and
29 proteins, the similar constraints reveal a comparable level of selective pressure to maintain their biological
30 functions. These results highlight the importance of the evolutionary information to recover meaningful
31 biological information to further characterize conformational ensembles.

© 2019 Published by Elsevier Ltd.

35 Introduction

36 The protein native state is described by a collection
37 of the different conformers which a given sequence
38 could adopt. This collection is also called a confor-
39 mational ensemble and is an essential concept to
40 understand protein biology [1,2]. The existence of
41 conformational ensembles is known since the crys-
42 tallization of hemoglobin with its two conformational
43 states T and R (deoxy and oxygenated forms) in the
44 early 1960. The growth of Protein Data Bank (PDB)
45 redundancy, refinement and development of tech-
46 niques such as NMR, small-angle X-ray scattering,
47 and single-molecule spectroscopy over the last years
48 have allowed the experimental characterization of a
49 large number of protein ensembles [2,3]. Structural
50 differences between conformers could result from
51 the relative movements of large domains as rigid
52 bodies [4], secondary and tertiary element rear-
53 rangements [5], and loop movements [6]. Apparently,

most globular proteins have very few conformers 54
describing their native state to achieve their functions 55
[7]. Proteins with low flexibility at the backbone 56
level, called rigids, have only one conformer in their 57
ensembles [7] like the cellulase from *Clostridium* 58
cellulolyticum [8]. Hemoglobin, as mentioned previ- 59
ously, is the paradigm for proteins with two con- 60
formers [9], while the dimeric catabolite activator 61
protein [10] and the human glucokinase have three 62
[11]. Complex proteins composed of several different 63
chains, like mitochondrial ATP synthase, could have 64
at least seven conformers [12]. As protein flexibility 65
increases, the number of conformers in the ense- 66
mble increases as well, giving rise to very complex 67
ensembles as in the case of intrinsically disordered 68
proteins (IDPs) or regions (IDRs). IDPs are character- 69
ized by the lack of tertiary structure under physiological 70
conditions [13,14]. IDP ensembles are composed by 71
a large number of interconverting conformers given 72
their low free-energy barriers among them [15]. Far 73

74 from being random polymers or random-coiled en-
75 sembles, it is becoming evident that IDP ensembles
76 are not fully disordered, showing transient short
77 and long-range structural organization [16]. Order-
78 disorder transitions are frequently observed in IDPs or
79 IDRs, sometimes associated with ligand binding [17]
80 but in other cases just reflecting the heterogeneous
81 composition of the ensembles [7,18].

82 Here, we studied the level of structural constraints
83 in IDPs ensembles compared with those found in
84 globular proteins. Structural constraints could be
85 studied using direct methods such as the measure-
86 ments of contacts between residues in a given
87 conformer and some derived parameters such as
88 the contact density (mean number of residue-residue
89 contacts per residue) or their interaction networks [19].
90 However, inter-residue contacts could be artifacts
91 or simply be irrelevant in very complex ensembles
92 such as those found in IDPs, making it difficult to
93 detect biologically relevant conformers [20]. For these
94 reasons, in this work, we evaluated the amount of
95 structural constraints using an evolutionary approach.
96 It is a well-established concept that conservation of
97 protein structures during evolution constrains se-
98 quence divergence modulating in this way the amino
99 acid substitution pattern of certain positions [21,22].
100 These structural constraints are evidenced in se-
101 quence alignments as differentially conserved posi-
102 tions, showing a given physicochemical bias or
103 subject to coevolutionary processes due to their
104 relative importance to maintain protein fold and
105 dynamics (i.e., conservation of given interactions to
106 increase stability, sustain protein movements). This
107 structurally constrained (SC) substitution pattern has
108 been exploited to improve models of molecular
109 evolution [23–25], explain rate heterogeneity [26],
110 make functional predictions [27], and compare the
111 substitution process in ordered and disordered
112 proteins [28] and in the inference of given tertiary
113 folds [29] to mention just a few examples of their many
114 applications. Furthermore, evolutionary information
115 could be used to predict native contacts and structural
116 models of globular domains [30–32]. More recently,
117 these methods were adapted to successfully predict
118 globular states in disordered proteins and to show the
119 evolutionary constraints in protein interfaces between
120 disordered and ordered proteins again showing the
121 importance of SC information during evolution [33,34].

122 Substitution patterns observed in sequence align-
123 ments can be described by evolutionary models
124 [35]. Alternative models, making different assump-
125 tions about the amino acid substitution pattern,
126 can be compared using maximum likelihood (ML)
127 estimations to decide which assumptions better
128 describe the evolutionary process in a given family.
129 In particular, in this work, a model of protein
130 evolution using protein structure to derive an SC
131 site-specific substitution pattern was used [24].
132 As this model is structure-specific, each protein

conformation represents different evolutionary models. 133
Using ML estimations, we then compared how the SC 134
substitution pattern outperforms models of evolution 135
lacking structural information (e.g., JTT [36], Dayhoff 136
[37], WAG [38]) in its ability to explain the observed 137
site-specific substitution pattern in a set of homologous 138
proteins for each studied protein. Interestingly, con- 139
sidering all conformers in the ensembles of globular 140
and IDP proteins, we found that the number of SC 141
positions is similar for both kinds of proteins. 142

Results 143

Description of the data sets 144

145 In the last years, an emerging picture evidences that 146
increasing structural differences between conformers, 147
connected by very different dynamical behaviors, 148
produces a continuum in protein space [39]. One 149
extreme feature of this continuum is the presence of 150
rigids proteins with almost no backbone differences 151
among their conformers and just displaying only 152
conformational diversity at the residue level [7]. 153
Increasing conformational diversity at the backbone 154
level could evidence the presence of disorder, where 155
the appearance of short-time dynamical behavior 156
allows for the sampling of a large conformational 157
space [40]. Figure 1 shows different types of 158
ensembles as protein conformational diversity in- 159
creases. In one extreme of the distribution (left-side 160
panel in Fig. 1), typical globular or ordered proteins 161
are shown. These proteins generally show large 162
proportions of secondary structure where their spatial 163
arrangement defines a single tertiary structure and 164
hydrophobic core. The higher density of inter-residue 165
interactions of this core constrains evolutionary rates 166
when compared to exposed residues [41] and also 167
contains enough information to define a global tertiary 168
arrangement [42]. As mentioned before, ordered 169
proteins could also contain different conformers 170
to achieve their biological functions (Fig. 1, middle 171
panel), giving place to additional restrictions in the 172
protein substitution pattern [43]. Middle-panel exam- 173
ples of Fig. 1 also display proteins with ordered or 174
globular regions as well as with very flexible regions 175
showing different dynamical behavior and possibly 176
originating disordered regions of different lengths. 177
Right panel in Fig. 1 shows a typical ensemble of IDPs 178
showing a collection of conformers determined by 179
NMR. These ensembles show highly flexible chains 180
and eventually small and transient segments of 181
secondary or tertiary structure [44]. Consequently, 182
IDPs have a large degree of conformational entropy 183
that can be limited by inter-residue interactions 184
originating a complex mixture of conformers in the 185
ensemble [15,20]. As described in Materials and 186
Methods, two hand-curated data sets were analyzed.

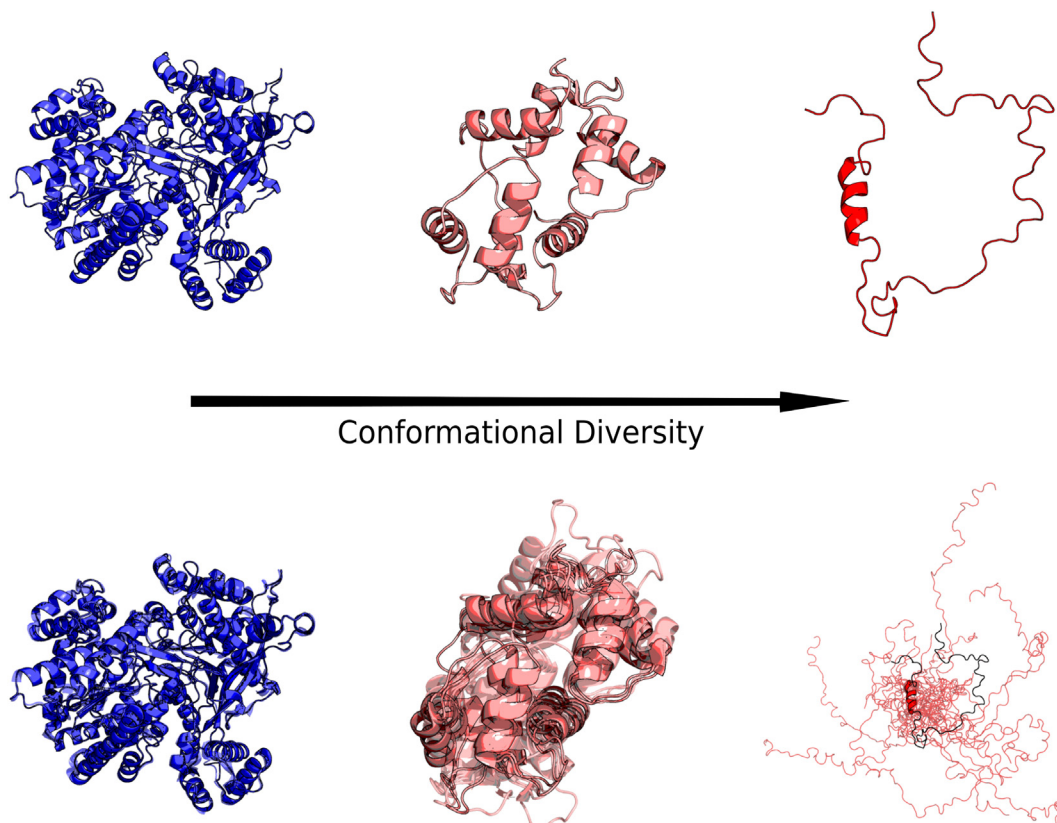


Fig. 1. Different protein ensembles as a function of flexibility increment. Top panel shows a given conformer, while the bottom panel shows all the available conformers in the ensemble. Left, maltodextrin phosphorylase, (PDB codes = 1AHP_A, 1AHP_B, 1L5V_B) showed as a rigid protein with 6.53% disordered and taken as a representative of ordered proteins. Calmodulin (PDB codes = 2FOT_A, 1LIN_A, 1NIW_E, 3G43_A, 2BE6_A, 1CDL_A, 3GP2_A, 4L79_B, 1CLL_A) shows 10.64% of disorder. Thylakoid soluble phosphoprotein, (PDB ID = 2FFT_A) is a typical IDP ensemble with 100% of estimated disorder. The percentages of disorder were estimated with ESpritz.

187 The ordered data set is composed of 183 proteins
 188 with known crystallographic structure containing non-
 189 missing residues, and a disordered data set contains
 190 93 NMR ensembles of different proteins. Disorder
 191 has been estimated in both data sets using ESpritz
 192 and Mobi 2.0 for the disordered and ordered data sets,
 193 respectively (see [Materials and Methods](#)). As is it
 194 shown in [Fig. 2](#), ordered proteins show a low predicted
 195 content of disordered residues, while the disordered
 196 data set shows a distribution of disordered residues.
 197 The median of these distribution is 58% of disordered
 198 positions (minimum 40% and up to 98%). It is then
 199 expected that the disordered data set contains small
 200 globular regions and more than the half of the protein
 201 in a disordered state. Sequence alignments for each
 202 protein in each data set were extracted from HSSP
 203 database (see [Materials and Methods](#)), and to avoid
 204 high occurrence of indels, sequences above 30%
 205 identity with the protein with known structure were only
 206 considered. Additional information about protein
 207 alignments could be found in [Fig. S1](#).

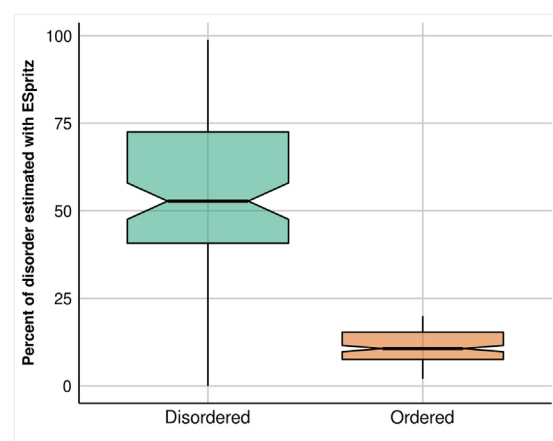


Fig. 2. Estimation of disorder content using NMR-ESpritz in the disordered set and ESpritz in the ordered set. It is shown that the ordered set has a low proportion of disorder well below the reported error in the estimation [45].

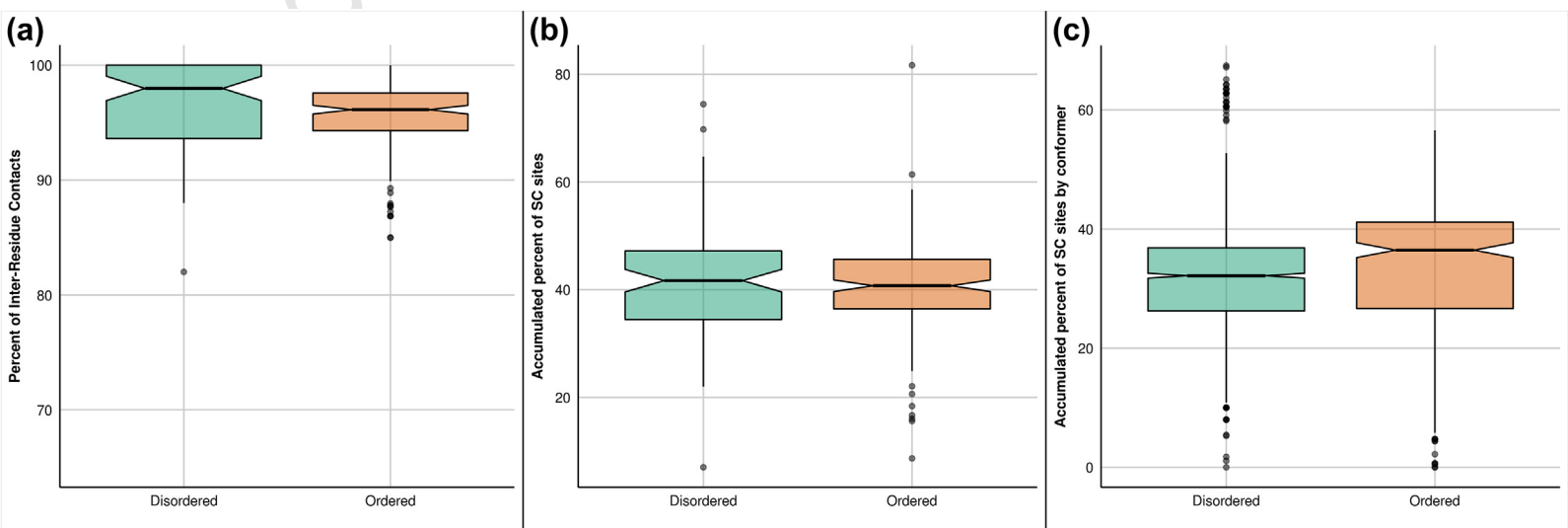


Fig. 3. (a) Percentage of inter-residue contacts for the disordered and ordered data sets (average median of 96.1%). (b) Distribution of the accumulated number of SCs for both data sets showing 41.6% and 40.5% of the positions. The distributions are statistically similar using a Kolmogorov–Smirnov test with p value = 0.39 and Mann–Whitney–Wilcoxon test with p value = 0.45. (c) Distribution of SCs per conformer per protein showing a median of 32.1% and 36.1% of their sites constrained.

208 Physical contacts versus structural constraints 209 during evolution

210 To assess the structural constraints in ordered and
211 disordered ensembles, we quantified the inter-residue
212 interactions accumulating the contact information
213 for each site through all the available conformers in
214 each corresponding ensemble (Fig. S2, panel A).
215 Accumulation is a reasonable idea sustained by the
216 particular contributions each conformer makes to the
217 biological function [2]. As a result, we obtained that
218 the great majority of residues are involved in inter-
219 residues contacts as it is shown in Fig. 3a. Permanent
220 secondary and tertiary contacts in ordered proteins
221 define their levels of structural constraints, while the
222 contribution of transient contacts along the entire
223 ensemble of IDPs produces almost the same amount
224 of accumulated inter-residues contacts (third quartile
225 is 100% and 97% for IDPs and ordered sets,
226 respectively). According to this result, the vast majority
227 of positions in IDPs are constrained by structural
228 restrictions as well as those for ordered proteins.
229 However, it is well established that the pattern of amino
230 acid substitutions in IDPs is different from the one
231 observed in ordered proteins. IDPs show also a highly
232 conserved composition of amino acids [46] instead of
233 the well-defined site-specific substitution pattern
234 observed in ordered proteins [47]. In addition, IDPs and
235 IDRs show higher evolutionary rates as well as higher
236 rates of insertions and deletions compared with
237 their ordered counterpart [13,44,48]. To elucidate the
238 influence of such high levels of structural constraints
239 (Fig. 3a), we turned to study the substitution pattern
240 observed in the homologous family of each protein
241 in both data sets. Using ML comparisons (Fig. S2,
242 panel B), we assessed if the observed substitution
243 pattern is better explained by an evolutionary model
244 containing structural information (like SCPE, see
245 Materials and Methods) or by other models not
246 containing this information (JTT, Dayhoff and WAG
247 models, see Materials and Methods). For every
248 position showing a SCPE site-specific substitution
249 matrix that outperforms each one of the other three
250 models, it is inferred as a site evolving under structural
251 constraints. Considering the different nature of ordered
252 and disordered ensembles, unexpectedly, we found
253 that the percentages of SCs are almost the same
254 in both types of ensembles (41.6% and 40.5% for
255 disordered and ordered data sets; Fig. 3b) and much
256 lower than estimations made using the accumulated
257 account of inter-residue contacts. Interestingly, the
258 individual conformers show slightly less percentages
259 of SC sites (Fig. 3c) showing 32.1% and 36.1% in
260 average for the disordered and ordered data sets.

261 SC sites

262 SC sites are then sites that at least have one physical
263 inter-residue contact in at least one conformer but also,

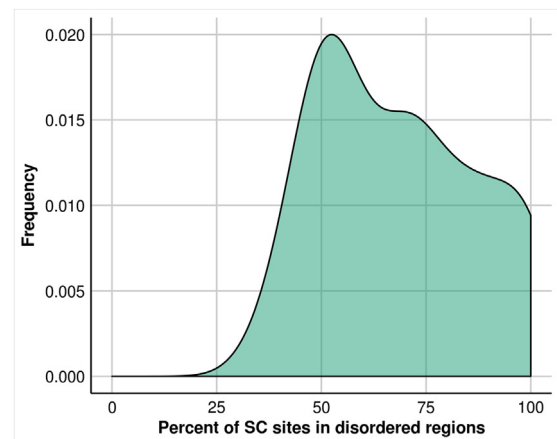


Fig. 4. Distribution of the accumulated number of SCs along all the ensembles. On average, 68.3% of the SC sites belong to predicted disordered regions.

and more importantly, modulates sequence diver- 264
gence in that specific position. To further investigate 265
these structural constraints, we studied the distribution 266
of SC sites. We found that ~68% of the SCs are located 267
in the disordered regions of the proteins belonging 268
to the disordered data set (Fig. 4). As we mentioned 269
before, disordered proteins could have permanent 270
or transient globular regions that could increase the 271
structural constraints of the protein as a whole. 272
However, the number of SC sites in the globular or 273
ordered regions of the disordered proteins is ~32%. 274
These results indicate that globular regions of disor- 275
dered proteins are less constrained than the corre- 276
sponding one observed in the ordered data set (see 277
Fig. 3b). Also, following our definition of inter-residue 278
contacts (see Materials and Methods), all estimated 279
contacts are tertiary and in ~57% the SCs are classified 280
as long-range inter-residue contacts (see Fig. 5). This 281

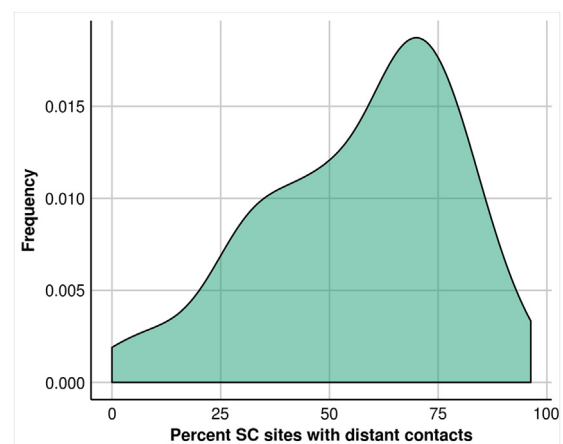


Fig. 5. Distribution of the accumulated number of SCs along all the ensembles, with long-distance contacts (at least five residues away). In average, 56.8% of the SC sites have long-range inter-residue contacts.

282 finding can explain how SC sites could appear in
 283 disordered regions. As we can see in Fig. 6, disordered
 284 proteins could have large conformational diversity.
 285 However, among the representative conformers of
 286 the ensembles, we can find some of them collapsing
 287 over the globular part of the protein or just adopting
 288 close conformations increasing in this way the number
 289 of contacts per site. As it is shown in Fig. 7, 51% of the
 290 positions have contacts that are present in the 100%
 291 of the conformers of the ensemble. However, there
 292 is still a tail in the distribution showing that single
 293 conformers could have SC sites; in other words, single

conformers could have inter-residue contacts that
 modulate the substitution pattern of those positions. 294 295

Ensemble redundancy 296

How many conformers are required to fully describe 297
 evolutionary structural constraints contained in se- 298
 quence alignments? When we calculated the mini- 299
 mum number of conformers per ensemble to reach 300
 the accumulated SC percentage per protein, we found 301
 that on average ~11 conformers are required for the 302
 proteins in the disordered data set (see Fig. 8), while 303

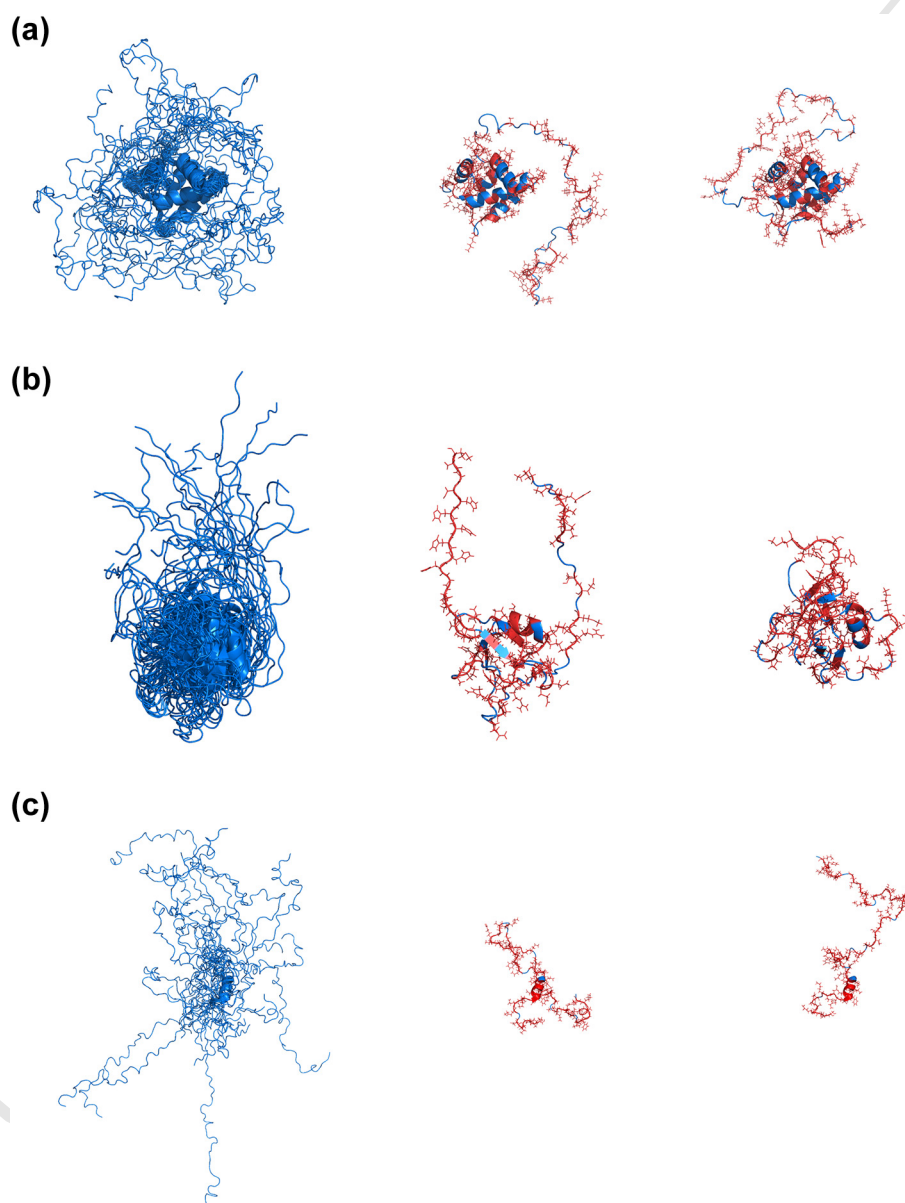


Fig. 6. Examples showing SC sites distribution in different conformers. The three panels (top, middle, and bottom) contain disordered proteins showing in the left the available ensemble, while in the middle and in the right, different conformers are shown. Proteins are shown. Cartoon representation was used. iSC sites are shown in red sticks, and the rest in blue. 2JRF_A, 2ADZ_A and 5MRG_A are the corresponding PDB codes for the top, middle, and bottom panels.

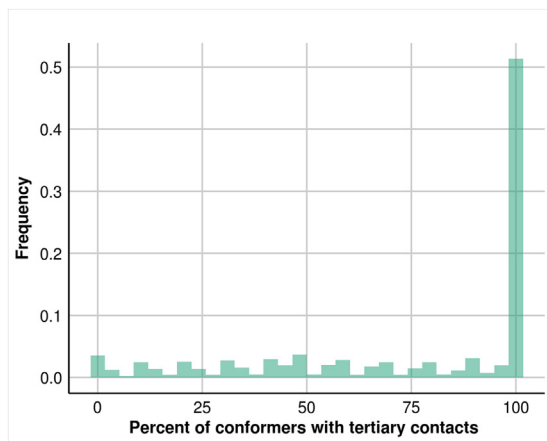


Fig. 7. Approximately ~51% of SC sites present contacts in 100% of the conformers, and only ~3% of SC sites present contacts in 50% of the conformers.

304 in the ordered data set, it is ~1.5. The value for the
 305 ordered data set is consistent with the available
 306 experimental evidence. Most ordered proteins show
 307 low conformational diversity, and then are called
 308 “rigid” [7], or could show very few conformers, mostly
 309 two, referring to the bound and unbound forms of the
 310 protein [49–51]. Due to the complexity of disordered
 311 ensembles, the number of conformers is difficult if
 312 not impossible to estimate. However, our measure of
 313 the number of conformers required to explain the
 314 evolutionary SC information in sequence alignments
 315 could offer a proxy to the number of conformers. Since
 316 the average of conformers in the NMR ensembles in
 317 our data set is ~20, our results indicate that they are
 Q8 mostly redundant.

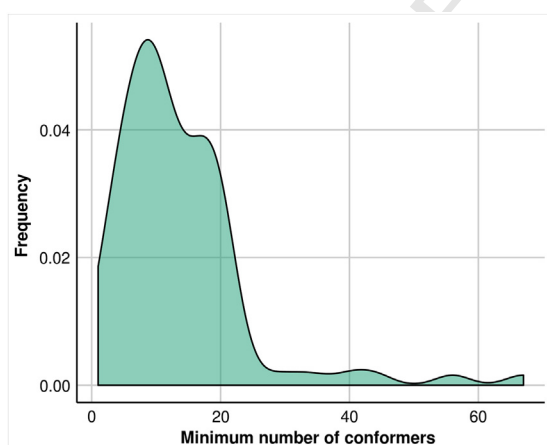


Fig. 8. Distribution of the minimum number conformers to reach the accumulated percentage of SC sites per protein for the 93 disordered proteins corresponding to the set obtained with Mobi 2.0 and ESpritz (NMR). Minimum = 1, average ~11, and maximum ~64.

Discussion

319

Two main findings emerge from the present work. 320
 First, the number of positions having inter-residue 321
 contacts accumulated along all available conformers 322
 in each ensemble approaches almost 100% of 323
 the positions (Fig. 3a). However, as we have shown, 324
 the average percentage of positions evolving under 325
 structural constraints is much lower, ~40% (Fig. 3b). 326
 Part of this reduction is expected, given that not all 327
 intramolecular non-covalent contacts could be equally 328
 relevant, for example, in structure stabilization [52]. 329
 Inaccurate models and atomic coordinate uncertain- 330
 tainties could also play a role to explain the observed 331
 difference between the amount of physical contacts 332
 and the observed evolutionary derived structural 333
 constraints [53–55]. In addition, the reduction could 334
 be also attributed to the lack of structure/conformer- 335
 specific information contained in sequence align- 336
 ments. This effect operates over SCPE substitution 337
 matrices, which are site and conformer specific but are 338
 evaluated using sequence alignments from corre- 339
 sponding homologous families. Thus, evolutionary 340
 information contained in those alignments reflects 341
 constraints of several sorts, such as structural 342
 divergence [41] or dynamical adaptations [56,57], 343
 which could certainly modify the contact pattern in the 344
 homologous proteins. It is then expected that this 345
 ~40% of SCs on average obtained for both ensem- 346
 bles does not capture subtle inter-residue contacts 347
 originated in functional adaptations for individual 348
 proteins. In line with this observation, it has been 349
 recently shown that the use of sequence alignments 350
 recovers the most conserved pattern of inter-residues 351
 contacts when co-evolutionary and evolutionary 352
 coupling methods are used [57]. The other important 353
 result is related with the comparable structural 354
 constraints on sequence divergence in ordered and 355
 disordered proteins (Fig. 3b). Our results suggest 356
 that individual contributions of each conformer in 357
 the disordered ensemble are required to sustain 358
 biological function as is well established for ordered 359
 proteins, and more recently suggested for disordered 360
 ones [2,13,48]. These small contributions from each 361
 disordered conformer give overall the same propor- 362
 tion of structural constraints as found in ordered 363
 proteins, possibly with different weights according to 364
 their biological role. 365

Interestingly, the number of conformers in the IDPs 366
 ensembles to reach the corresponding level of global 367
 constraints per protein is ~11 (Fig. 8). This means 368
 that IDP ensembles are redundant in terms of 369
 conformations and that possibly the number of 370
 biologically relevant conformers in IDP ensembles 371
 would not be so large as expected due to their 372
 high flexibility. These results are in agreement with 373
 the idea that different members of the ensemble 374
 could be directly involved in protein function, but 375
 also, they could be important as a local minimum 376

377 representatives in the interconversion of biologically
378 relevant conformations [58].

379 Our results highlight the importance of the evolu-
380 tionary analysis in the discrimination of inter-residue
381 contacts to detect meaningful biological information
382 as well as the estimation of the number of conformers
383 and structural constraints in such complex ensem-
384 bles as those belonging to IDPs.

385 Materials and Methods

386 Data set collection

387 Globular or ordered protein ensembles were
388 obtained from the CoDNas database [59]. Considering
389 the presence of missing residues as a primary indicator
390 of IDRs in proteins [60], we selected 183 proteins
391 having no missing residues in any of their available
392 conformers. These selected protein ensembles have
393 at least five conformers in the database to assure a
394 good estimation of the conformational variability [61].
395 Only the pair of conformers showing the maximum
396 RMSD along all the ensemble was considered in this
397 set. To obtain the IDPs data set, we predicted and
398 estimated disorder in all the available NMR protein
399 structures in PDB (available May 2018) using NMR-
400 ESpritz [45] and Mobi 2.0 [62]. After a hand-curated
401 revision considering length and protein biology, we
402 finally obtained 93 protein NMR ensembles with more
403 than 40% of disordered positions. Ordered set of
404 proteins showed negligible levels of disorder predicted
405 with ESpritz X-ray (see Figs. 3 and S3).

406 SC substitution pattern estimation

407 In Fig. S2, we resumed the workflow to analyze SCs
408 and physical contacts. For each conformer and each
409 protein in both data sets (for the disordered data set,
410 we considered all the NMR available conformers, and
411 for the ordered data set, we used those corresponding
412 for the maximum RMSD according to CoDNas), the
413 SCPE model of protein evolution was run [24]. SCPE
414 derives site-specific substitution matrices using evolu-
415 tionary simulations under neutral conditions for protein
416 fold conservation [47,63] (please see Fig. S4). Briefly,
417 it uses energetic calculations to evaluate the structural
418 perturbation introduced by non-synonymous substitu-
419 tions in the simulation process. Using ML estimations,
420 it is possible to compare SCPE matrices with models
421 lacking structural information such as JTT [36], Dayh-
422 off [64], and WAG [38]. Site-specific ML calculations
423 were performed with the HYPHY package [65]. The
424 alignments used for the ML analysis were obtained
425 from HSSP [66] database. Neighbor-joining distance
426 phylogenetic trees were obtained with the Phylip [67]
427 package. To define whether a site was SC, Akaike
428 information criteria (AIC) coefficient was used [68] and
429 a ranking for the estimated models was made using

Δ AIC [69] in which models having Δ AIC ≤ 2 have a
substantial support, those where Δ AIC is between
4 and 7 have an intermediate support, and those
with Δ AIC > 10 have no support. Tertiary contacts
were estimated considering the distance between two
non-contiguous residues having the van der Waals
spheres of each residue side chain heavy atoms
below 1.0 Å. Long-range inter-residues contacts were
estimated using same definition but considering ± 5
residues of a given residue.

Acknowledgments

G.P. and M.S.F. are CONICET researchers, and
J.M. and A.M.M. are PhD and postdoctoral fellows of
the same institution.

This work was supported by Universidad Nacional
de Quilmes (PUNQ 1004/11) (G.P.), Agencia de
Ciencia y Tecnología (PICT-2014-3430) (G.P.), and
COST Action (BM1405) NGP-net (S.C.E.T.). This
project has received funding from the European
Union's Horizon 2020 research and innovation
program under the Marie Skłodowska-Curie grant
agreement No 778247 (IDPfun). The funders had no
role in study design, data collection and analysis,
decision to publish, or preparation of the manuscript. Q9

Appendix A. Supplementary data

Supplementary data to this article can be found
online at <https://doi.org/10.1016/j.jmb.2019.01.031>.

Received 5 November 2018; 457

Received in revised form 23 January 2019; 458

Accepted 24 January 2019 459

Available online xxxx 460

461

Keywords: 462

protein evolution; 463

protein ensemble; 464

conformational diversity; 465

disordered proteins 466

467

Abbreviations used: Q6

PDB, Protein Data Bank; IDP, intrinsically disordered 469

protein; IDR, intrinsically disordered region; SC, structu- 470

rally constrained site; ML, maximum likelihood; AIC, 471

Akaike information criteria. 472

References

- [1] C.J. Tsai, S. Kumar, B. Ma, R. Nussinov, Folding funnels, 475
binding funnels, and protein function, *Protein Sci.* 8 (1999) 476
1181–1190. 477

- 478 [2] G. Wei, W. Xi, R. Nussinov, B. Ma, Protein ensembles: how
479 does nature harness thermodynamic fluctuations for life? The
480 diverse functional roles of conformational ensembles in the
Q12 cell, *Chem. Rev.* (2016) (acs.chemrev.5b00562).
482 [3] C. Marino-Buslje, A.M. Monzon, D.J. Zea, S. Fornasari, G.
483 Parisi, On the dynamical incompleteness of the Protein Data
484 Bank, *Brief. Bioinform.* (2017) 1–4.
485 [4] M. Gerstein, N. Echols, Exploring the range of protein
486 flexibility, from a structural proteomics perspective, *Curr.*
487 *Opin. Chem. Biol.* 8 (2004) 14–19.
488 [5] M. Gerstein, A database of macromolecular motions, *Nucleic*
489 *Acids Res.* 26 (1998) 4280–4290.
490 [6] Y. Gu, D.-W. Li, R. Brüschweiler, Decoding the mobility
491 and time scales of protein loops, *J. Chem. Theory Comput.*
492 11 (2015) 1308–1314.
493 [7] A.M. Monzon, D.J. Zea, M.S. Fornasari, T.E. Saldaño, S.
494 Fernandez-Alberti, S.C.E. Tosatto, G. Parisi, Conformational
495 diversity analysis reveals three functional mechanisms in
496 proteins, *PLoS Comput. Biol.* 13 (2017) 1–29.
497 [8] G. Parsiegla, C. Reverbel-Leroy, C. Tardif, J.P. Belaich, H.
498 Driguez, R. Haser, Crystal structures of the cellulase Cel48F
499 in complex with inhibitors and substrates give insights into its
500 processive action, *Biochemistry* 39 (2000) 11238–11246.
501 [9] M.F. Perutz, W. Bolton, R. Diamond, H. Muirhead, H.C.
502 Watson, Structure of haemoglobin. An X-ray examination of
503 reduced horse haemoglobin, *Nature* 203 (1964) 687–690.
504 [10] N. Popovych, S. Sun, R.H. Ebright, C.G. Kalodimos,
505 Dynamically driven protein allostery, *Nat. Struct. Mol. Biol.*
506 13 (2006) 831–838.
507 [11] K. Kamata, M. Mitsuya, T. Nishimura, J.-I. Eiki, Y. Nagata,
508 Structural basis for allosteric regulation of the monomeric
509 allosteric enzyme human glucokinase, *System* 12 (2004)
510 429–438.
511 [12] A. Zhou, A. Rohou, D.G. Schep, J.V. Bason, M.G. Montgomery,
512 J.E. Walker, N. Grigorieff, J.L. Rubinstein, Structure and
513 conformational states of the bovine mitochondrial ATP syn-
514 thase by cryo-EM, *elife* 4 (2015), e10180.
515 [13] J. Siltberg-Liberles, J.a. Grahnen, D.a. Liberles, The evolu-
516 tion of protein structures and structural ensembles under
517 functional constraint, *Genes* 2 (2011) 748–762.
518 [14] P. Tompa, Intrinsically unstructured proteins, *Trends Biochem.*
519 *Sci.* 27 (2002) 527–533.
520 [15] M. Varadi, S. Kosol, P. Lebrun, E. Valentini, M. Blackledge,
521 A.K. Dunker, I.C. Felli, J.D. Forman-Kay, R.W. Kriwacki, R.
522 Pierattelli, J. Sussman, D.I. Svergun, V.N. Uversky, M.
523 Vendruscolo, D. Wishart, P.E. Wright, P. Tompa, pE-DB: a
524 database of structural ensembles of intrinsically disordered
525 and of unfolded proteins, *Nucleic Acids Res.* 42 (2014)
526 D326–D335.
527 [16] P. Tompa, Unstructural biology coming of age, *Curr. Opin.*
528 *Struct. Biol.* 21 (2011) 419–425.
529 [17] D. Zea, A.M. Monzon, C. Gonzalez, M.S. Fornasari, S.C.E.
530 Tosatto, G. Parisi, Disorder transitions and conformational
531 diversity cooperatively modulate biological function in pro-
532 teins, *Protein Sci.* 25 (2016) 1138–1146.
533 [18] S. DeForte, V.N. Uversky, Resolving the ambiguity: making
534 sense of intrinsic disorder when PDB structures disagree,
535 *Protein Sci.* 25 (2016) 676–688.
536 [19] D. Piovesan, G. Minervini, S.C. Tosatto, The RING 2.0 web
537 server for high quality residue interaction networks, *Nucleic*
538 *Acids Res.* 44 (W1) (2016) W367–W374.
539 [20] P. Sormanni, D. Piovesan, G.T. Heller, M. Bonomi, P. Kucik,
540 C. Camilloni, M. Fuxreiter, Z. Dosztanyi, R.V. Pappu, M.M.
541 Babu, S. Longhi, P. Tompa, A.K. Dunker, V.N. Uversky, S.C.E.
Tosatto, M. Vendruscolo, Simultaneous quantification of
protein order and disorder, *Nat. Chem. Biol.* 13 (2017) 543
339–342. 544
[21] J. Overington, M.S. Johnson, A. Sali, T.L. Blundell, Tertiary
structural constraints on protein evolutionary diversity: 545
templates, key residues and structure prediction, *Proc. Biol.* 546
Sci. 241 (1990) 132–145. 548
[22] C.L. Worth, S. Gong, T.L. Blundell, Structural and functional
constraints in the evolution of protein families, *Nat. Rev. Mol.* 549
Cell Biol. 10 (2009) 709–720. 551
[23] U. Bastolla, H.E. Roman, M. Vendruscolo, Neutral evolution
of model proteins: diffusion in sequence space and over- 552
dispersion, *J. Theor. Biol.* 200 (1999) 49–64. 554
[24] G. Parisi, J. Echave, Structural constraints and emergence
of sequence patterns in protein evolution, *Mol. Biol. Evol.* 18 555
(2001) 750–756. 557
[25] C.L. Kleinman, N. Rodrigue, N. Lartillot, H. Philippe, 558
Statistical potentials for improved structurally constrained
evolutionary models, *Mol. Biol. Evol.* 27 (2010) 1546–1560. 559
[26] J. Echave, S.J. Spielman, C.O. Wilke, Causes of evolutionary 560
rate variation among protein sites, *Nat. Rev. Genet.* 17 (2016) 561
109–121. 563
[27] A.L. Simon, E.A. Stone, A. Sidow, Inference of functional 564
regions in proteins by quantification of evolutionary con-
straints, *Proc. Natl. Acad. Sci. U. S. A.* (2001). Q13 565
[28] C.J. Brown, A.K. Johnson, G.W. Daughdrill, Comparing 567
models of evolution for ordered and disordered proteins, *Mol.* 568
Biol. Evol. 27 (2010) 609–621. 569
[29] J. Surkont, J.B. Pereira-Leal, Evolutionary patterns in coiled- 570
coils, *Genome Biol. Evol.* 7 (2015) 545–556. 571
[30] T.A. Hopf, C.P.I. Schärfe, J.P.G.L.M. Rodrigues, A.G. Green, O. 572
Kohlbacher, C. Sander, A.M.J.J. Bonvin, D.S. Marks, Sequence
co-evolution gives 3D contacts and structures of protein 573
complexes, *elife* 3 (2014), <https://doi.org/10.7554/eLife.03430>. 575
[31] F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D.S. Marks, C. 576
Sander, R. Zecchina, J.N. Onuchic, T. Hwa, M. Weigt, Direct-
coupling analysis of residue coevolution captures native 577
contacts across many protein families, *Proc. Natl. Acad. Sci.* 579
U. S. A. 108 (2011) E1293–E1301. 580
[32] S. Ovchinnikov, H. Kamisetty, D. Baker, Robust and accurate 581
prediction of residue–residue interactions across protein
interfaces using evolutionary information, *elife* 3 (2014) 582
<https://doi.org/10.7554/elife.02030>. 584
[33] A. Toth-Petroczy, P. Palmedo, J. Ingraham, T.A. Hopf, B. 585
Berger, C. Sander, D.S. Marks, Structured states of disor-
dered proteins from genomic sequences, *Cell* 167 (2016) 587
158–170.e12. 588
[34] R. Pancsa, F. Zsolyomi, P. Tompa, Co-evolution of intrinsically 589
disordered proteins with folded partners witnessed by evolu-
tionary couplings, *Int. J. Mol. Sci.* 19 (2018), [https://doi.org/10.](https://doi.org/10.3390/ijms19113315) 591
[3390/ijms19113315](https://doi.org/10.3390/ijms19113315). 592
[35] J.L. Thorne, Models of protein sequence evolution and their 593
applications, *Curr. Opin. Genet. Dev.* 10 (2000) 602–605. 594
[36] D.T. Jones, W.R. Taylor, J.M. Thornton, The rapid generation 595
of mutation data matrices from protein sequences, *Comput.* 596
Appl. Biosci. 8 (1992) 275–282. 597
[37] M.O. Dayhoff, A model of evolutionary change in proteins, 598
Atlas Protein Seq. Struct. 5 (1972) 89–99. 599
[38] S. Whelan, N. Goldman, A general empirical model of protein 600
evolution derived from multiple protein families using a
maximum-likelihood approach, *Mol. Biol. Evol.* 18 (2001) 602
691–699. 603
[39] H.N. Motlagh, J.O. Wrabl, J. Li, V.J. Hilser, The ensemble 604
nature of allostery, *Nature* 508 (2014) 331–339. 605

- 606 [40] R.B. Berlow, H. Jane Dyson, P.E. Wright, Functional
607 advantages of dynamic protein disorder, *FEBS Lett.* 589
608 (2015) 2433–2440.
- 609 [41] K. Illergård, D.H. Ardell, A. Elofsson, Structure is three to ten
610 times more conserved than sequence—a study of structural
611 response in protein cores, *Proteins* 77 (2009) 499–508.
- 612 [42] A.F. Pereira de Araujo, J.N. Onuchic, A sequence-
613 compatible amount of native burial information is sufficient
614 for determining the structure of small globular proteins, *Proc.*
615 *Natl. Acad. Sci. U. S. A.* 106 (2009) 19001–19004.
- 616 [43] D. Javier Zea, A. Miguel Monzon, M.S. Fornasari, C. Marino-
617 Buslje, G. Parisi, Protein conformational diversity correlates
618 with evolutionary rate, *Mol. Biol. Evol.* 30 (2013) 1500–1503.
- 619 [44] J.D. Forman-Kay, T. Mittag, From sequence and forces to
620 structure, function, and evolution of intrinsically disordered
621 proteins, *Structure* 21 (2013) 1492–1499.
- 622 [45] I. Walsh, A.J.M. Martin, T. Di Domenico, S.C.E. Tosatto,
623 ESpritz: accurate and fast prediction of protein disorder,
624 *Bioinformatics* 28 (2011) 503–509.
- 625 [46] H.A. Moesa, S. Wakabayashi, K. Nakai, A. Patil, Chemical
626 composition is maintained in poorly conserved intrinsically
627 disordered regions and suggests a means for their classifica-
628 tion, *Mol. BioSyst.* 8 (2012) 3262–3273.
- 629 [47] M.S. Fornasari, G. Parisi, J. Echave, Site-specific amino acid
630 replacement matrices from structurally constrained protein
631 evolution simulations, *Mol. Biol. Evol.* 19 (2002) 352–356.
- 632 [48] C.J. Brown, A.K. Johnson, A. Keith Dunker, G.W. Daughdrill,
633 Evolution and disorder, *Curr. Opin. Struct. Biol.* 21 (2011)
634 441–446.
- 635 [49] A. Gutteridge, J. Thornton, Conformational changes ob-
636 served in enzyme crystal structures upon substrate binding,
637 *J. Mol. Biol.* 346 (2005) 21–28.
- 638 [50] A. Gutteridge, J. Thornton, Conformational change in substrate
639 binding, catalysis and product release: an open and shut case,
640 *FEBS Lett.* 567 (2004) 67–73.
- 641 [51] T. Amemiya, R. Koike, A. Kidera, M. Ota, PSCDB: a
642 database for protein structural change upon ligand binding,
643 *Nucleic Acids Res.* 40 (2012) D554–D558.
- 644 [52] R. Sathyapriya, J.M. Duarte, H. Stehr, I. Filippis, M. Lappe,
645 Defining an essence of structure determining residue
646 contacts in proteins, *PLoS Comput. Biol.* 5 (2009), e1000584.
- 647 [53] S.O. Garbuzynskiy, B.S. Melnik, M.Y. Lobanov, A.V.
648 Finkelstein, O.V. Galzitskaya, Comparison of X-ray and
649 NMR structures: is there a systematic difference in residue
650 contacts between X-ray- and NMR-resolved protein struc-
651 tures? *Proteins: Struct. Funct. Bioinf.* 60 (2005) 139–147.
- 652 [54] C.A.E.M. Spronk, J.P. Linge, C.W. Hilbers, G.W. Vuister,
653 Improving the quality of protein structures derived by NMR
654 spectroscopy, *J. Biomol. NMR* 22 (2002) 281–289.
- 655 [55] C.A.E.M. Spronk, S.B. Nabuurs, A.M.J.J. Bonvin, E. Krieger, 655
656 G.W. Vuister, G. Vriend, The precision of NMR structure
657 ensembles revisited, *J. Biomol. NMR* 25 (2003) 225–234.
- 658 [56] A.M. Monzon, D.J. Zea, C. Marino-Buslje, G. Parisi, 658
659 Homology modeling in a dynamical world, *Protein Sci.* 659
660 (2017), <https://doi.org/10.1002/pro.3274>.
- 661 [57] D.J. Zea, A.M. Monzon, G. Parisi, C. Marino-Buslje, How 661
662 is structural divergence related to evolutionary information? 662
663 *Mol. Phylogenet. Evol.* 127 (2018) 859–866.
- 664 [58] K.K. Turoverov, I.M. Kuznetsova, V.N. Uversky, The protein 664
665 kingdom extended: ordered and intrinsically disordered 665
666 proteins, their folding, supramolecular complex formation, 666
667 and aggregation, *Prog. Biophys. Mol. Biol.* 102 (2010) 73–84.
- 668 [59] A.M. Monzon, C.O. Rohr, M.S. Fornasari, G. Parisi, 668
669 CoDNaS 2.0: a comprehensive database of protein confor- 669
670 mational diversity in the native state, *Database* 2016 (2016), 670
671 baw038.
- 672 [60] A.K. Dunker, J.D. Lawson, C.J. Brown, R.M. Williams, P. 672
673 Romero, J.S. Oh, C.J. Oldfield, A.M. Campen, C.M. Ratliff, 673
674 K.W. Hipps, J. Ausio, M.S. Nissen, R. Reeves, C. Kang, C.R. 674
675 Kissinger, R.W. Bailey, M.D. Griswold, W. Chiu, E.C. Garner, 675
676 Z. Obradovic, Intrinsically disordered protein, *J. Mol. Graph.* 676
677 *Model.* 19 (2001) 26–59.
- 678 [61] R.B. Best, K. Lindorff-Larsen, M.A. DePristo, M. Vendruscolo, 678
679 Relation between native ensembles and experimental struc- 679
680 tures of proteins, *Proc. Natl. Acad. Sci. U. S. A.* 103 (2006) 680
681 10901–10906.
- 682 [62] D. Piovesan, S.C.E. Tosatto, Mobi 2.0: an improved method 682
683 to define intrinsic disorder, mobility and linear binding regions 683
684 in protein structures, *Bioinformatics* 34 (2017) 122–123.
- 685 [63] G. Parisi, J. Echave, Generality of the structurally con- 685
686 strained protein evolution model: assessment on represen- 686
687 tatives of the four main fold classes, *Gene* 345 (2005) 45–53.
- 688 [64] W.M. Fitch, M.O. Dayhoff, Atlas of protein sequence and 688
689 structure, 1972, *Syst. Zool.* 22 (1973) 196.
- 690 [65] S.L. Pond, S.D. Frost, S.V. Muse, HyPhy: hypothesis testing 690
691 using phylogenies, *Bioinformatics* 21 (2005) 676–679.
- 692 [66] W.G. Touw, C. Baakman, J. Black, T.A.H. te Beek, E. 692
693 Krieger, R.P. Joosten, G. Vriend, A series of PDB-related 693
694 databanks for everyday needs, *Nucleic Acids Res.* 43 (2015) 694
695 D364–D368.
- 696 [67] J. Feisenstein, PHYLIP: Phylogeny Inference Package 696
697 Version 3.2 Manual, 1989.
- 698 [68] H. Akaike, A new look at the statistical model identification, 698
699 *IEEE Trans. Autom. Contr.* 19 (1974) 716–723.
- 700 [69] F.S. Guthery, K.P. Burnham, D.R. Anderson, Model selection 700
701 and multimodel inference: a practical information-theoretic 701
702 approach, *J. Wildl. Manag.* 67 (2003) 655.
- 703