

pyDockDNA: a new web server for energy-based protein-DNA docking and scoring

1 Luis Ángel Rodríguez-Lumbreras^{1,2}, Brian Jiménez-García^{1,3}, Silvia Giménez-Santamarina^{1,4}
2 and Juan Fernández-Recio^{1,2*}

3 ¹Barcelona Supercomputing Center, Barcelona, Spain

4 ²Instituto de Ciencias de la Vid y del Vino (ICVV), CSIC-UR-Gobierno de La Rioja, Logroño,
5 Spain.

6 ³Zymvol Biomodeling SL, Barcelona, Spain

7 ⁴ICMol, Universitat de València, Paterna, Spain

8

9 * **Correspondence:**

10 Juan Fernández-Recio

11 juan.fernandezrecio@icvv.es

12 **Keywords:** Structural modeling, Ab initio docking, protein-DNA interaction, Scoring function,
13 Nucleotide parameters, Docking benchmark.

14 **Abstract**

15 Proteins and nucleic acids are essential biological macromolecules for cell life. Indeed, interactions
16 between proteins and DNA regulate many biological processes such as protein synthesis, signal
17 transduction, DNA storage, or DNA replication and repair. Despite their importance, less than 4% of
18 total structures deposited in the Protein Data Bank (PDB) correspond to protein-DNA complexes, and
19 very few computational methods are available to model their structure. We present here the
20 pyDockDNA web server, which can successfully model a protein-DNA complex with a reasonable
21 predictive success rate as benchmarked in a standard dataset of proteins in complex with DNA in B-
22 DNA conformation. The server implements the pyDockDNA program, as a module of pyDock suite,
23 thus including third-party programs, modules, and previously developed tools, as well as new modules
24 and parameters to handle the DNA properly. The user is asked to enter PDB files for protein and DNA
25 input structures (or suitable models) and select the chains to be docked. The server calculations are
26 mainly divided into three steps: sampling by FTDOCK, scoring with new energy-based parameters and
27 the possibility of applying external restraints. The user can select different options for these steps. The
28 final output screen shows a 3D representation of the top 10 models and a table sorting the model
29 according to the scoring function selected previously. All these output files can be downloaded,
30 including the top 100 models predicted by pyDockDNA. The server can be freely accessed for
31 academic use (<https://model3dbio.csic.es/pydockdna>).

32

33 **1 INTRODUCTION**

34 Proteins and nucleic acids are fundamental biological macromolecules whose functions and
35 interactions are vital to regulating cell's life. Their interactions regulate many biological processes such

36 as protein synthesis, signal transduction, DNA storage, and DNA replication and repair, among others.
37 Learning how protein and DNA interact is fundamental to fully elucidate many central biological
38 processes and disease mechanisms, and can also support the discovery of novel therapeutic targets.
39 Although 192,025 structures have been experimentally determined and deposited in the June 2022
40 release of Protein Data Bank (PDB), only 10,480 of them correspond to protein-nucleic acid complexes
41 (this includes 6,732 protein-DNA complexes). Thus, the number of protein-DNA structures
42 experimentally determined is clearly much smaller than the number of protein-DNA complexes that
43 are expected to be formed in cells. This gap is partially explained by the difficulty of the experimental
44 determination process, i.e. a very time-consuming process in the best scenarios or impossible in many
45 cases due to limitations on the experimental techniques. For this reason, a computational approach on
46 modelling protein-DNA interactions could be of enormous help.

47 Even though theoretical models of macromolecular structures are usually less accurate than direct
48 experimental measurements, they can yield sufficient information to build a working hypothesis,
49 complementing experimental approaches in elucidating protein-DNA interactions and guiding further
50 experimental analyses to identify essential amino acids or nucleotide residues. From a computational
51 point of view, there are two main approaches to model the structure of a protein-DNA complex:
52 template-based modelling and *ab initio* docking. Template-based modelling aims to model a complex
53 based on the structure of a homologous complex. The popularity of template-based methods has
54 increased in the past years, especially for modelling protein-protein complexes, thanks to the
55 development and support of many structural databases of protein interactions that can provide the
56 required templates, such as 3D Complex (Levy et al., 2006), Dockground (Kundrotas et al., 2018), or
57 Interactome3D (Mosca et al., 2013). However, the quality of template-based predictions clearly
58 depends on the availability of suitable templates, not particularly high in the case of protein-DNA
59 interactions (see for instance PDIDb (Norambuena and Melo, 2010), which makes these methods of
60 very limited applicability. On the other hand, *ab initio* docking methods aim at predicting the three-
61 dimensional structures of macromolecular complexes, starting from the atomic coordinates of their
62 components. *Ab initio* docking methods do not depend on a priori in external information which makes
63 them more useful in the actual protein-DNA context.

64 The methodology for prediction and modelling of protein-protein complexes is very well established
65 despite there are still many challenges to be addressed. Numerous protein-protein docking methods
66 have been developed and assessed as shown in the Critical Assessment of PRediction of Interactions
67 (CAPRI) community-wide experiment. During the past editions of the CAPRI experiment (Janin et al.,
68 2003), targets other than protein-protein complexes were proposed: protein-RNA complex (Lensink
69 and Wodak, 2010) (T33, T34), protein-peptide (T60-64) or protein-heparin (T57) among others.
70 However, protein-DNA docking received limited attention from the CAPRI community and developers
71 of computational methods. Macromolecular docking protocols that accept protein and DNA
72 coordinates as input include FTDock (Gabb et al., 1997), GRAMM-X (Tovchigrechko and Vakser,
73 2006), HEX (Macindoe et al., 2010), PatchDock (Schneidman-Duhovny et al., 2005; Macindoe et al.,
74 2010) and NPDock (Tuszynska et al., 2015), HDock (Yan et al., 2017), ClusPro (Comeau et al., 2004)
75 and HADDOCK (Van Zundert et al., 2016) servers. From this list of tools, only NPDock and HDock
76 were originally developed for protein-nucleic acid docking; the rest were developed as protein-protein
77 docking tools that also accept nucleic acids coordinates, but they lack an intrinsic scoring function
78 dedicated to assessing protein-DNA interactions. These protocols usually report high predictive rates
79 in bound conditions, i.e. when the co-crystallized partners in a known complex structure are separated
80 and re-docked. However, despite bound docking is useful for testing and development purposes, it does
81 not represent realistic conditions and therefore it is of limited practical value for biology. Therefore, it
82 is important to have available datasets to test protein-DNA docking tools in unbound conditions.

83 Compared to protein-protein docking, where the most recent release of the Weng's group Protein-
84 Protein Docking Benchmark 5.5 (Vreven et al., 2015) has 257 entries, and to protein-RNA docking,
85 where there are different reported benchmarks (Barik et al., 2012; Pérez-Cano et al., 2012; Huang and
86 Zou, 2013; Nithin et al., 2017), for protein-DNA docking there is only one available benchmark, which
87 contains 47 complexes (van Dijk and Bonvin, 2008). Using this benchmark, protein-DNA docking
88 protocols report moderate success rates in unbound conditions. For instance, on a subset of 23 cases
89 from this benchmark, HDock success rate for top 10 models (i.e. at least one near-native structure
90 within the top 10 models) is less than 10%, while success rate for top 100 is slightly over 30% (Yan et
91 al., 2017). NPDock reports a maximum success rate (i.e. at least one near-native conformation found
92 in the entire prediction set) of 7/47 (15%) (Tuszynska et al., 2015). Protein-DNA docking with
93 HADDOCK reported an excellent performance (van Dijk and Bonvin, 2010) when using restraints
94 from the real interface. This represents a very promising approach, but in a realistic scenario, lack of
95 knowledge on the actual complex interface might limit its application. A more recent coarse-version
96 of HADDOCK protein-DNA docking shows similar accuracy with ~6-fold speed increase over
97 atomistic calculations (Honorato et al., 2019). The need of new computational tools to address
98 unbound protein-DNA docking is clear. We present here a new web server that implements the
99 pyDockDNA protein-DNA docking and scoring protocol, as a new module of pyDock version 4
100 (upcoming publication). The original pyDock docking and scoring approach (Cheng et al., 2007),
101 which showed excellent performance for the prediction of protein-protein docking (Lensink et al.,
102 2019; Rosell et al., 2020), has been rewritten in Python 3 and extended for its application to protein-
103 DNA docking, with new functionalities to handle the nucleic acid structures and upgraded atomic
104 solvation parameters for a more accurate scoring of protein-DNA interactions.

105

106 2 MATERIALS AND METHODS

107 Data Sets: protein-DNA docking benchmark and external case studies

108 In order to test the new pyDockDNA docking protocol, we used a previously developed protein-DNA
109 docking benchmark (version 1.2) (van Dijk and Bonvin, 2008). The benchmark contains bound and
110 unbound x-ray crystallography and NMR structures for 47 protein-DNA complexes, in which DNA is
111 in B-DNA conformation. These are classified as 'easy', 'intermediate' or 'difficult' cases, based on the
112 interface RMSD values between the bound and unbound components of the complex.

113 An additional set of case studies was compiled following the criteria selection of the above described
114 protein-DNA docking benchmark. This test set is composed of ten protein-DNA complexes, where
115 both bound and unbound structures are available for each reference complex, and the sequences are
116 different from those in the first protein-DNA docking benchmark. Protein-DNA complex and unbound
117 structures were compiled from the Protein-DNA Interface Database (PDIdb) (Norambuena and Melo,
118 2010) and the Protein Data Bank (PDB) (Berman et al., 2000). Only complexes that meet the following
119 conditions were considered: i) DNA sequence length larger than eight base pairs, and ii) proteins
120 without mutations in the core of the complex interface. To find the protein unbound structures of the
121 protein-DNA complexes selected, all the PDB entries containing only protein structures were retrieved,
122 including structures solved by NMR. Crystallographic structures with a resolution worse than 3.0 Å
123 were not considered. To avoid redundancy, entries with sequence similarity larger than or equal to 90%
124 were discarded. PDBeFOLD (Krissinel and Henrick, 2004) was used to find correspondences between
125 bound and unbound protein structures. This tool performs structural alignments between two (pairwise
126 alignment) or more (multi-alignment) molecules using their 3-dimensional structures. The alignment

127 is based on the Secondary Structure Matching algorithm (Krissinel and Henrick, 2004). Alignments
128 with a Q-score higher than 8.0, high P-score and sequence similarity around 90-100% were accepted
129 as the corresponding unbound. Then, both bound and unbound structures for each case, were post-
130 processed according to the protocol followed in a previously developed protein-DNA docking
131 benchmark, for instance by checking consistency between unbound and bound coordinates in chain
132 IDs, residue numbers and atom names (van Dijk and Bonvin, 2008). The unbound DNA models were
133 generated by using the software 3DNA (Lu and Olson, 2003; Lu and Olson, 2008), in canonical B-
134 DNA conformation (fiber model 4).

135 This additional test set (Table 1) is freely available at the "Help" section of the server
136 (https://model3dbio.csic.es/pydockdna/info/faq_and_help#extended_benchmark).

137 [INSERT HERE TABLE 1]

138 **Sampling**

139 In this first step, the input files with the coordinates in PDB format for the structures (or models) of a
140 protein and a DNA molecule (which can be B-DNA or any other conformation) are checked for
141 potential format errors, missing side-chains in the protein are rebuilt with SCWRL 3.0 (Bower et al.,
142 1997), and the electrostatics Amber94 force field (Cornell et al., 1995) is loaded, assigning the charges
143 to the atoms. Then, rigid-body docking poses between the protein and the DNA, represented as 3D
144 grids, are generated with a faster and parallelized version of the original FTDock (v2.0) software (Gabb
145 et al., 1997) in which the number of cells in the grid is optimized for maximum computing efficiency
146 (Jiménez-García et al., 2013). The molecule with the longest maximal distance between any pair of
147 atoms is considered the receptor, that is, the fixed molecule, and the other one is the ligand or mobile
148 molecule. By default, the program uses 0.7 Å grid cell size, 1.3 Å surface thickness, 12° rotation
149 sampling, and keeps the best 3 poses for each rotation. For each target, a total of 10,000 docking poses
150 were generated.

151 **Scoring**

152 Finally, the protein-DNA docking poses are ranked using a scoring function composed of electrostatics,
153 desolvation and van der Waals energy. This new pyDockDNA scoring function is adapted from the
154 previously pyDock scoring function for protein-protein docking (Grosdidier et al., 2007; Jiménez-
155 García et al., 2013), which now includes atom types for nucleotides from Amber94 force field (Cornell
156 et al., 1995) in order to calculate for the modelled protein-DNA complexes. The nucleotide AMBER
157 atom types have been mapped to the previously defined atom types in pyDock within a new parameter
158 set (*nuc.dat*).

159 **Implementation of pyDockDNA web server**

160 The program pyDockDNA is built as a module of the new pyDock 4.0 version (upcoming publication),
161 thus include the same third-party programs, modules and tools from previous versions of pyDock as
162 well as new functionalities to handle the nucleic acid structures properly. The user can select the chains
163 to be docked, the energetic scoring function, and even include external information (from available
164 experimental data or using predictive methods such as the DBSI server [REF:
165 <https://doi.org/10.1093/bioinformatics/btw315>]) as residue-nucleotide distance restraints to rescore
166 docking models as previously described for pyDockRST (Chelliah et al., 2006). The output will be a
167 set of docking models represented in different formats: i) the 3D structure of the best-scoring 10
168 docking models in terms of scoring can be visualized in the output screen, ii) the PDB files for the best-

169 scoring 100 models can be directly downloaded, and iii) the rotation/translation vectors are provided
170 to generate up to a total of 10,000 docking poses. A summary of the docking results can be visualized
171 as a plot with the distribution of the different energy values obtained for all docking poses (Figure 1).

172 [INSERT HERE FIGURE 1]

173 **Clustering of protein-DNA docking models in benchmarking**

174 When testing this software (see Results) we have run several docking executions in parallel, using
175 different initial random rotations for the input structures, and the best-scoring 100 resulting models for
176 each individual run were merged into a single pool. To avoid redundancy in the final set, all docking
177 orientations were clustered by pyProCT analysis software (Gil and Guallar, 2014), which implements
178 the GROMOS clustering algorithm (Daura et al., 1999). Distance matrix is built with pyRMSD with
179 the option "QCP OMP CALCULATOR" to compute the ligand root-mean-square deviation (L-RMSD)
180 values for all pairs of docking orientations after their receptors were superimposed
181 (<https://github.com/victor-gil-sepulveda/pyRMSD/>). The L-RMSD cut-off value 4.0 Å was used to
182 define the clusters. For each defined cluster of models, the orientation with the lowest docking score is
183 defined as the cluster representative.

184 **Docking performance**

185 We have evaluate the predicted performance of pyDockDNA in different conditions as the success
186 rates for the obtained top N docking models, which is the % of benchmark cases in which a near-native
187 (acceptable) solution is found within the top N docking models. A near-native solution is defined as a
188 docking orientation model with L-RMSD ≤ 10 Å with respect to the reference structure.

189

190 **3 RESULTS AND DISCUSSION**

191 **Performance of pyDockDNA evaluated on the protein-DNA docking benchmark**

192 The pyDockDNA web server has been tested on the 47 cases of a previously reported protein-DNA
193 docking benchmark (see Methods). It is known that using different randomly rotated input structures
194 can slightly affect docking predictions in FFT-based docking protocols as in FTDOCK, because this
195 can modify the mapping of the atom positions on the 3D grids (Garzon et al., 2009; Pallara et al., 2016).
196 To check for convergence, we applied pyDockDNA to 10 different random rotations of the initial input
197 structures for each benchmark case and computed the predictive success rates for the results obtained
198 from each randomly rotated input structures. The results indicate even more differences in the
199 predictive values than previously reported for protein-protein docking (Table S1). For instance, the
200 success rates for the top 10 models ranged from 12.8% to 21.3%. Therefore, for a more robust
201 evaluation, we merged the results of all 10 docking executions and clustered the obtained docking
202 models to remove similar orientations (see Methods). Figure 2 shows the predictive success rates of
203 the cluster representatives resulting from merging these 10 docking runs. The predictive success for
204 the default pyDock scoring function (including parameters for nucleotide atoms, see Methods) are
205 better than those obtained for the individual docking runs, which means that increasing sampling
206 variability when using different random initial rotations, followed by redundancy removal with
207 clustering, have improved the docking results.

208 [INSERT HERE FIGURE 2]

209 We further analyzed whether a scoring function previously developed for protein-protein docking was
210 really optimal for protein-DNA docking, since for the latter, electrostatics energy term is expected to
211 have a larger contribution to binding energy due to the higher overall charge of DNA molecules.
212 Moreover, desolvation atomic parameters were previously derived for protein-protein docking in
213 pyDock, but they were not specifically optimized here for nucleotide atoms. To analyze the role of
214 desolvation in protein-DNA scoring, we rescored the generated docking models with the pyDockDNA
215 scoring function but excluding desolvation energy. This greatly improved the success rates, as the curve
216 *pyDockDNA (no desolv)* shows in Figure 2. This indeed indicates that desolvation is not really needed
217 for the scoring of the protein-DNA docking models generated by FFT-based sampling, perhaps because
218 the parameters have not been yet optimized for nucleotide atoms, or because electrostatics is more
219 relevant in protein-DNA interactions than in protein-protein complexes, as above discussed. We tested
220 other solvation parameters for protein-DNA reported in the literature (Kagawa et al., 1989), but the
221 docking results did not improve (we are currently working on the optimization of these parameters in
222 search of a better desolvation for protein-DNA).

223 In addition, we have also tried other combinations of energy terms, for instance, increasing the factor
224 for van der Waals to 1.0 (we previously found that geometrical complementarity was very important
225 in protein-RNA; (Pérez-Cano et al., 2016), or removing desolvation and van der Waals terms from the
226 scoring function to test the relevance of electrostatics scoring alone, but none of these new combined
227 scoring functions improved the prediction rates (Figure S1).

228 In a rigid-body docking approach as pyDock, it is known that protein flexibility upon binding is perhaps
229 the most determinant factor for docking success. To further analyze whether the docking performance
230 of pyDockDNA is affected by the flexibility of the protein or the DNA input molecules during the
231 complex formation, we have grouped the docking results on the protein-DNA docking benchmark
232 according to the flexibility of the protein or the DNA, that is, based on the RMSD between the unbound
233 molecules and the corresponding ones in the complex. Regarding protein flexibility, in order to make
234 groups of similar size, we defined these three categories: low (unbound-bound RMSD $< 1 \text{ \AA}$), medium
235 ($1 \text{ \AA} \leq \text{unbound-bound RMSD} < 3 \text{ \AA}$) and high (unbound-bound RMSD $\geq 3 \text{ \AA}$) flexible cases. As for
236 DNA flexibility, we defined these three categories: low (unbound-bound RMSD $< 3 \text{ \AA}$), medium (3 \AA
237 $\leq \text{unbound-bound RMSD} < 5 \text{ \AA}$) and high (unbound-bound RMSD $\geq 5 \text{ \AA}$) flexible cases. The results
238 are shown in Figure 3. We can observe that the docking predictive performance does not get worse
239 when protein flexibility is higher (actually, for pyDockDNA with no desolvation, success rates increase
240 when protein flexibility is medium or high). However, we can see that the docking performance for
241 highly flexible DNA molecules is dramatically low. We should note that in this benchmark, proteins
242 in general show smaller unbound-bound RMSD values (average 2.6 \AA) than DNA (average 4.2 \AA). In
243 addition, due to the different RMSD cut-off values used for proteins and for DNA, the unbound-bound
244 RMSD values for high flexible proteins (average 4.8 \AA) are much smaller than those for DNA (average
245 7.8 \AA), which could explain the much worse predictive rates in the group of highly flexible DNA.

246 [INSERT HERE FIGURE 3]

247 **Application to external case studies**

248 For further testing, we have applied pyDockDNA to a set of ten additional protein-DNA cases (Table
249 1) where the structures for the complex and the unbound protein were available at PDB, and the
250 unbound DNA was modelled in canonical B-DNA conformation (see Methods).

251 For each case study, we have performed a single pyDockDNA execution on the randomly rotated
252 unbound protein and DNA structures, a realistic scenario, since the pyDockDNA server only provides
253 results for a docking execution (randomly rotated input structures should be provided to the server in
254 independent executions for a more thorough docking study similar to the benchmark performance
255 analysis above shown). Overall, we obtained predictive success rates of 10% and 30% (for the top 10
256 and 100 models, respectively) when using pyDockDNA scoring function, and 10% and 60% (for the
257 top 10 and 100 models, respectively), when using pyDockDNA without desolvation. Given the small
258 number of cases of these additional set, these values are within the expected range according to the
259 larger docking benchmark set.

260 The most successful case is the complex between the DNA binding domain of Early B-cell Factor 1
261 (Ebf1) bound to a 22bp DNA (PDB 3MLO), where a near-native docking model (L-RMSD 3.33 Å
262 with respect to the reference) is found with rank 5 when using pyDockDNA (no desolvation) scoring
263 function (Figure 4A). When using pyDockDNA (including desolvation) scoring function, this docking
264 model is ranked 6, so it is still within top 10 models. This case has low-flexible protein but high-flexible
265 DNA.

266 [INSERT HERE FIGURE 4]

267 Another case is the complex between the catabolite gene activator protein and a 11bp DNA (PDB
268 1O3R), where we found an almost acceptable docking model (L-RMSD 10.76 Å with respect to the
269 reference) with rank 5, when using pyDockDNA either including solvation or not (Figure 4B). This
270 case has also low-flexible protein but medium-flexible DNA. If this case had been considered
271 acceptable, the success rates for the top 10 would have been 20%. However, these percentage values
272 are perhaps not very meaningful considering the low number of cases in this external test set.
273 Interestingly, when using van der Waals term with weighing factor 1.0 (instead of the default factor in
274 pyDock and pyDockDNA, that is 0.1), we find near-native solutions in 3 more cases, in addition to
275 3MLO: i) 5JLT (L-RMSD 7.08 Å) with rank 1 when using desolvation; ii) 2NTC (L-RMSD 7.25 Å)
276 with rank 3 not using desolvation, and iii) 2PI0 (L-RMSD 6.63 Å) with rank 3 and 2, when using
277 desolvation or not using it, respectively. Therefore, for half of these external case studies, we found
278 near-native docking models within the top 10 models with pyDockDNA, using different variants of the
279 scoring function.

280 In summary, we present here the pyDockDNA web server to model protein-DNA complexes, which
281 implements a docking method based on pyDock, with new parameters for DNA. We have evaluated
282 the performance on unbound proteins and modelled DNA molecules in canonical B-DNA
283 conformation, using a known protein-DNA docking benchmark. The results show near 40% success
284 rate for the top 10 models when using the pyDockDNA (no desolvation) scoring function, after merging
285 the results from 10 docking executions using different randomly rotated initial structures, and
286 clustering the models to remove redundant ones. The method has been applied to external case studies,
287 with similar predictive performance.

288

289 4 AUTHOR CONTRIBUTIONS

290 Conflict of Interest

291 The authors declare that the research was conducted in the absence of any commercial or financial
292 relationships that could be construed as a potential conflict of interest.

293 L.A.R.-L. wrote the first draft, performed the analysis, optimized the energy-based parameters,
 294 implemented the final version of the server and updated the module on the standalone version of
 295 pyDock 4.0. B.J.-G. implemented the first version of the server and that of the new module on the
 296 standalone version of the pyDock software, and reviewed the draft. S.-G.-S. compiled the external case
 297 studies and validated the software. J.F.-R. devised the idea, optimized the energy-based parameters,
 298 analyzed the results, and wrote the final manuscript.

299 5 FUNDING

300 This work was supported by grant number PID2019-110167RB-I00 / AEI / 10.13039/501100011033.
 301 B.J.-G. is employed by Zymvol Biomodeling on a project which received funding from the European
 302 Union's Horizon 2020 research and innovation programme under Marie Skłodowska-Curie grant
 303 agreement No. 801342 (Tecniospring INDUSTRY) and the Government of Catalonia's Agency for
 304 Business Competitiveness (ACCIÓ).

305

306 6 ACKNOWLEDGMENTS

307 7 REFERENCES

- 308 Barik, A., Nithin, C., Manasa, P., and Bahadur, R.P. (2012). A protein-RNA docking benchmark (I):
 309 Nonredundant cases. *Proteins: Structure, Function and Bioinformatics* 80(7), 1866-1871. doi:
 310 10.1002/prot.24083.
- 311 Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., et al. (2000). The
 312 Protein Data Bank. *Nucleic Acids Research* 28(1), 235-242. doi: 10.1093/nar/28.1.235.
- 313 Bower, M.J., Cohen, F.E., and Dunbrack, R.L. (1997). Prediction of protein side-chain rotamers from
 314 a backbone-dependent rotamer library: A new homology modeling tool. *Journal of Molecular*
 315 *Biology* 267(5), 1268-1282. doi: 10.1006/jmbi.1997.0926.
- 316 Chelliah, V., Blundell, T.L., and Fernandez-Recio, J. (2006). Efficient restraints for protein-protein
 317 docking by comparison of observed amino acid substitution patterns with those predicted
 318 from local environment. *J Mol Biol* 357(5), 1669-1682. doi: 10.1016/j.jmb.2006.01.001.
- 319 Cheng, T.M.-K., Blundell, T.L., and Fernandez-Recio, J. (2007). pyDock: Electrostatics and
 320 desolvation for effective scoring of rigid-body protein-protein docking. *Proteins: Structure,*
 321 *Function, and Bioinformatics* 68(2), 503-515. doi: 10.1002/prot.21419.
- 322 Comeau, S.R., Gatchell, D.W., Vajda, S., and Camacho, C.J. (2004). ClusPro: A fully automated
 323 algorithm for protein-protein docking. *Nucleic Acids Research* 32(WEB SERVER ISS.), 96-
 324 99. doi: 10.1093/nar/gkh354.
- 325 Cornell, W.D., Cieplak, P., Bayly, C.I., Gould, I.R., Merz, K.M., Ferguson, D.M., et al. (1995). A
 326 Second Generation Force Field for the Simulation of Proteins, Nucleic Acids, and Organic
 327 Molecules. *Journal of the American Chemical Society* 117(19), 5179-5197. doi:
 328 10.1021/ja00124a002.
- 329 Daura, X., Gademann, K., Jaun, B., Seebach, D., van Gunsteren, W.F., and Mark, A.E. (1999).
 330 Peptide Folding: When Simulation Meets Experiment. *Angewandte Chemie International*
 331 *Edition* 38(1/2), 236-240. doi: 10.1002/(SICI)1521-3773(19990115)38:1/2<236::AID-
 332 ANIE236>3.3.CO;2-D.

- 333 Gabb, H.a., Jackson, R.M., and Sternberg, M.J. (1997). Modelling protein docking using shape
334 complementarity, electrostatics and biochemical information. *Journal of molecular biology*
335 272(1), 106-120. doi: 10.1006/jmbi.1997.1203.
- 336 Garzon, J.I., Lopez-Blanco, J.R., Pons, C., Kovacs, J., Abagyan, R., Fernandez-Recio, J., et al.
337 (2009). FRODOCK: a new approach for fast rotational protein-protein docking.
338 *Bioinformatics* 25(19), 2544-2551. doi: 10.1093/bioinformatics/btp447.
- 339 Gil, V.A., and Guallar, V. (2014). PyProCT: Automated cluster analysis for structural bioinformatics.
340 *Journal of Chemical Theory and Computation* 10(8), 3236-3243. doi: 10.1021/ct500306s.
- 341 Grosdidier, S., Pons, C., Solernou, A., and Fernández-Recio, J. (2007). Prediction and scoring of
342 docking poses with pyDock. *Proteins* 69(4), 852-858. doi: 10.1002/prot.21796.
- 343 Honorato, R.V., Roel-Touris, J., and Bonvin, A. (2019). MARTINI-Based Protein-DNA Coarse-
344 Grained HADDOCKing. *Front Mol Biosci* 6, 102. doi: 10.3389/fmolb.2019.00102.
- 345 Huang, S.Y., and Zou, X. (2013). A nonredundant structure dataset for benchmarking protein-RNA
346 computational docking. *Journal of Computational Chemistry* 34(4), 311-318. doi:
347 10.1002/jcc.23149.
- 348 Janin, J., Henrick, K., Moult, J., Eyck, L.T., Sternberg, M.J.E., Vajda, S., et al. (2003). CAPRI: A
349 critical assessment of PRedicted interactions. *Proteins: Structure, Function and Genetics*
350 52(1), 2-9. doi: 10.1002/prot.10381.
- 351 Jiménez-García, B., Pons, C., and Fernández-Recio, J. (2013). pyDockWEB: A web server for rigid-
352 body protein-protein docking using electrostatics and desolvation scoring. *Bioinformatics*
353 29(13), 1698-1699. doi: 10.1093/bioinformatics/btt262.
- 354 Kagawa, T.F., Stoddard, D., Zhou, G., and Ho, P.S. (1989). Quantitative Analysis of DNA Secondary
355 Structure from Solvent-Accessible Surfaces: The B- to Z-DNA Transition as a Model.
356 *Biochemistry* 28(16), 6642-6651. doi: 10.1021/bi00442a017.
- 357 Krissinel, E., and Henrick, K. (2004). Secondary-structure matching (SSM), a new tool for fast
358 protein structure alignment in three dimensions. *Acta Crystallographica Section D:*
359 *Biological Crystallography* 60(12 I), 2256-2268. doi: 10.1107/S09074444904026460.
- 360 Kundrotas, P.J., Anishchenko, I., Dauzhenka, T., Kotthoff, I., Mnevets, D., Copeland, M.M., et al.
361 (2018). Dockground: A comprehensive data resource for modeling of protein complexes.
362 *Protein Sci* 27(1), 172-181. doi: 10.1002/pro.3295.
- 363 Lensink, M.F., Brysbaert, G., Nadzirin, N., Velankar, S., Chaleil, R.A.G., Gerguri, T., et al. (2019).
364 Blind prediction of homo- and hetero-protein complexes: The CASP13-CAPRI experiment.
365 *Proteins* 87(12), 1200-1221. doi: 10.1002/prot.25838.
- 366 Lensink, M.F., and Wodak, S.J. (2010). Docking and scoring protein interactions: CAPRI 2009.
367 *Proteins: Structure, Function and Bioinformatics* 78(15), 3073-3084. doi:
368 10.1002/prot.22818.
- 369 Levy, E.D., Pereira-Leal, J.B., Chothia, C., and Teichmann, S.A. (2006). 3D complex: a structural
370 classification of protein complexes. *PLoS Comput Biol* 2(11), e155. doi:
371 10.1371/journal.pcbi.0020155.
- 372 Lu, X.J., and Olson, W.K. (2003). 3DNA: A software package for the analysis, rebuilding and
373 visualization of three-dimensional nucleic acid structures. *Nucleic Acids Research* 31(17),
374 5108-5121. doi: 10.1093/nar/gkg680.

- 375 Lu, X.J., and Olson, W.K. (2008). 3DNA: A versatile, integrated software system for the analysis,
376 rebuilding and visualization of three-dimensional nucleic-acid structures. *Nature Protocols*
377 3(7), 1213-1227. doi: 10.1038/nprot.2008.104.
- 378 Macindoe, G., Mavridis, L., Venkatraman, V., Devignes, M.D., and Ritchie, D.W. (2010).
379 HexServer: An FFT-based protein docking server powered by graphics processors. *Nucleic*
380 *Acids Research* 38(SUPPL. 2), 445-449. doi: 10.1093/nar/gkq311.
- 381 Mosca, R., Céol, A., and Aloy, P. (2013). Interactome3D: adding structural details to protein
382 networks. *Nature Methods* 10(1), 47-53. doi: 10.1038/nmeth.2289.
- 383 Nithin, C., Mukherjee, S., and Bahadur, R.P. (2017). A non-redundant protein-RNA docking
384 benchmark version 2.0. *Proteins* 85(2), 256-267. doi: 10.1002/prot.25211.
- 385 Norambuena, T., and Melo, F. (2010). The Protein-DNA Interface database. *BMC bioinformatics* 11,
386 262-262. doi: 10.1186/1471-2105-11-262.
- 387 Pallara, C., Rueda, M., Abagyan, R., and Fernández-Recio, J. (2016). Conformational heterogeneity
388 of unbound proteins enhances recognition in protein-protein encounters. *Journal of Chemical*
389 *Theory and Computation* 12(7), 3236-3249. doi: 10.1021/acs.jctc.6b00204.
- 390 Pérez-Cano, L., Jiménez-García, B., and Fernández-Recio, J. (2012). A protein-RNA docking
391 benchmark (II): Extended set from experimental and homology modeling data. *Proteins:*
392 *Structure, Function and Bioinformatics* 80(7), 1872-1882. doi: 10.1002/prot.24075.
- 393 Pérez-Cano, L., Romero-Durana, M., and Fernández-Recio, J. (2016). Structural and energy
394 determinants in protein-RNA docking. *Methods*, 1-8. doi: 10.1016/j.ymeth.2016.11.001.
- 395 Rosell, M., Rodríguez-Lumbreras, L.A., Romero-Durana, M., Jimenez-Garcia, B., Diaz, L., and
396 Fernandez-Recio, J. (2020). Integrative modeling of protein-protein interactions with pyDock
397 for the new docking challenges. *Proteins* 88(8), 999-1008. doi: 10.1002/prot.25858.
- 398 Schneidman-Duhovny, D., Inbar, Y., Nussinov, R., and Wolfson, H.J. (2005). PatchDock and
399 SymmDock: Servers for rigid and symmetric docking. *Nucleic Acids Research* 33(SUPPL. 2),
400 363-367. doi: 10.1093/nar/gki481.
- 401 Tovchigrechko, A., and Vakser, I.A. (2006). GRAMM-X public web server for protein-protein
402 docking. *Nucleic Acids Research* 34(WEB. SERV. ISS.), 310-314. doi: 10.1093/nar/gkl206.
- 403 Tuszynska, I., Magnus, M., Jonak, K., Dawson, W., and Bujnicki, J.M. (2015). NPDock: A web
404 server for protein-nucleic acid docking. *Nucleic Acids Research* 43(W1), W425-W430. doi:
405 10.1093/nar/gkv493.
- 406 van Dijk, M., and Bonvin, A.M.J.J. (2008). A protein-DNA docking benchmark. *Nucleic Acids*
407 *Research* 36(14), e88-e88. doi: 10.1093/nar/gkn386.
- 408 van Dijk, M., and Bonvin, A.M.J.J. (2010). Pushing the limits of what is achievable in protein-DNA
409 docking: Benchmarking HADDOCK's performance. *Nucleic Acids Research* 38(17), 5634-
410 5647. doi: 10.1093/nar/gkq222.
- 411 Van Zundert, G.C.P., Rodrigues, J.P.G.L.M., Trellet, M., Schmitz, C., Kastiris, P.L., Karaca, E., et
412 al. (2016). The HADDOCK2.2 Web Server: User-Friendly Integrative Modeling of
413 Biomolecular Complexes. *Journal of Molecular Biology* 428(4), 720-725. doi:
414 10.1016/j.jmb.2015.09.014.
- 415 Vreven, T., Moal, I.H., Vangone, A., Pierce, B.G., Kastiris, P.L., Torchala, M., et al. (2015).
416 Updates to the Integrated Protein-Protein Interaction Benchmarks: Docking Benchmark

417 Version 5 and Affinity Benchmark Version 2. *Journal of Molecular Biology* 427(19), 3031-
418 3041. doi: 10.1016/j.jmb.2015.07.016.

419 Yan, Y., Zhang, D., Zhou, P., Li, B., and Huang, S.Y. (2017). HDOCK: A web server for protein-
420 protein and protein-DNA/RNA docking based on a hybrid strategy. *Nucleic Acids Research*
421 45(W1), W365-W373. doi: 10.1093/nar/gkx407.

422

423

424

425

426

427 **FIGURE LEGENDS**

428 **Figure 1.** Schematic representation of the pyDockDNA web server main functionalities.

429

430 **Figure 2.** Predictive performance for the top $N=1, 5, 10, 100$ models of pyDockDNA (with and without
431 desolvation) on the protein-DNA docking benchmark.

432

433 **Figure 3.** Predictive performance for the top 10 models of pyDockDNA (with and without desolvation)
434 on the protein-DNA docking benchmark when cases are grouped according to (A) protein flexibility
435 (low: $\text{RMSD} < 1 \text{ \AA}$; medium: $1 \text{ \AA} \leq \text{RMSD} < 3 \text{ \AA}$; high: $\text{RMSD} \geq 3 \text{ \AA}$), and (B) DNA flexibility (low:
436 $\text{RMSD} < 3 \text{ \AA}$; medium: $3 \text{ \AA} \leq \text{RMSD} < 5 \text{ \AA}$; high: $\text{RMSD} \geq 5 \text{ \AA}$). See more details about flexibility
437 definition in main text.

438

439 **Figure 4.** Application of pyDockDNA to case studies. (A) Near-native model (in yellow) obtained by
440 pyDockDNA docking between a modelled 22bp DNA (receptor) and Ebf1 (ligand). This model was
441 ranked 5 with pyDockDNA (no desolvation) scoring function and has L-RMSD 3.33 \AA with respect to
442 the reference (PDB 3MLO; in red). (B) Reasonable model (in yellow) obtained by pyDockDNA
443 docking between the catabolite gene activator protein (receptor) and a modelled 11bp DNA (ligand).
444 This model was ranked 5 with pyDockDNA (either with desolvation or with no desolvation) scoring
445 function and has L-RMSD 10.76 \AA with respect to the reference (PDB 1O3R; in red).

446

447

448

449

450

451

452

453

454

455 TABLES

456 Table 1: List of case studies

PDB complex	Protein	PDB unbound protein	RMSD unbound-bound protein	DNA	RMSD unbound-bound DNA
5JLT	phage T4 MotA DNA-binding domain	1KAF	0.83 ^a	22bp dsDNA	1.89
2X6V	TBX5	2X6V	0.55	11bp DNA	2.03
3POV	SOX	3FHD	1.46	19bp DNA	2.26
4UUV	ETV4 DNA-binding ETS domain	5ILU	1.24	10bp DNA	2.81
2NTC	sv40 large T antigen	2FUF	1.13 ^a	21-nt PEN element of the SV40 DNA origin	2.96
2ITL	sv40 large T antigen	4NBP	5.37 ^a	24-nt PEN element of the SV40 DNA origin	3.84
3MFK	Protein C-Ets1	1GVJ	5.61 ^a	stromelysin-1 promoter DNA	4.34
2PI0	IRF-3	3QU6	0.76 ^a	PRDIII-I region of human interferon-B promoter strand 1	4.46
1O3R	catabolite gene activator protein	4R8H	0.65	11bp DNA	4.77
3MLO	Ebf1	3LYR	0.71 ^a	22bp DNA	5.11

457 ^a In cases with more than one protein-DNA interface in the x-ray structure, the average value is
458 provided.

459

460

461

462

463