

A geometric deep learning approach to predict binding conformations of bioactive molecules

Oscar Méndez-Lucio, ^{*,a} Mazen Ahmad,^a Ehecatl Antonio del Rio-Chanona,^b
Jörg Kurt Wegner ^{*,a}

^a Generative AI Team, High Dimensional Biology and Discovery Data Sciences, Janssen Research & Development, Janssen Pharmaceutica N.V., Turnhoutseweg 30, Beerse B-2340, Belgium

^b Centre for Process Systems Engineering (CPSE), Department of Chemical Engineering, Imperial College London, UK

* E-mail: omendezl@its.jnj.com, jwegner@its.jnj.com

Abstract

Understanding the interactions formed between a ligand and its molecular target is key to guide the optimization of molecules. Different experimental and computational methods have been key to understand better these intermolecular interactions. Herein, we report a method based on geometric deep learning that is capable of predicting the binding conformations of ligands to protein targets. Concretely, the model learns a statistical potential based on distance likelihood which is tailor-made for each ligand-target pair. This potential can be coupled with global optimization algorithms to reproduce experimental binding conformations of ligands. We show that the potential based on distance likelihood described in this paper performs similar or better than well-established scoring functions for docking and screening tasks. Overall, this method represents an example of how artificial intelligence can be used to improve structure-based drug design.

28 Introduction

29 There is no doubt that drug design is a challenging task. One of the difficulties
30 arises from the fact that only a small portion of the large chemical space (circa
31 10^{60} drug-like molecules^{1,2}) will bind to a specific biological target resulting in a
32 therapeutical effect. In this context, knowing up-front the biological target and its
33 three-dimensional structure seems to be associated with higher success rates³.
34 To a very large extent, this success results from the use of experimental and
35 computational methods that can help understand the key interactions between
36 a ligand and its molecular target to guide the optimization of molecules. In fact,
37 it is known that these intermolecular interactions are a key factor driving drug
38 potency and selectivity⁴. Experimental methods such as X-ray diffraction, NMR
39 crystallography and more recently Cryo-EM have been of paramount
40 importance for drug discovery projects to explore and understand these
41 intermolecular interactions^{3,5}. In a similar way, computational methods have
42 also played an important role since they allow to virtually study compounds that
43 have not been synthesized yet. In particular, molecular docking has been
44 recently used to virtually screen ultra-large compound libraries^{6,7}, although other
45 methods such as molecular dynamics are also commonly used for drug
46 discovery.

47 In the recent years, the explosion of experimental structural data has also
48 allowed the application of machine learning and artificial intelligence to study
49 ligand-target interactions. For example, machine learning has been successfully
50 applied to identify regions of a protein where a ligand can directly bind⁸⁻¹⁰.
51 Additionally, a wide a variety of methods have been developed to predict
52 binding affinity from the three-dimensional structure of a ligand-target
53 complex^{11,12}. Many of these methods make use of engineered descriptors that
54 capture the main ligand-target interactions which can be fed into a predictive
55 algorithm¹³⁻¹⁶, while others directly use convolutional neural networks (CNNs)¹⁷⁻
56 ²⁰ or graph convolutional neural networks (GNNs)^{21,22} for the prediction task.

57 Despite the need for more computationally efficient methods for structure-based
58 design, there are few efforts to accelerate or improve the structure prediction of
59 a bound ligand by using artificial intelligence or machine learning. Most of

60 current artificial intelligence methods applied to structure-based drug discovery
61 rely on the 3D structure of a ligand-target complex previously obtained either by
62 experimental or computational approaches. Herein, we report DeepDock, a
63 method based on geometric deep learning that is capable of predicting the
64 binding conformation of ligands to protein targets. For this, the method learns a
65 statistical potential based on distance likelihood which is tailor-made for each
66 ligand-target. Statistical potentials have been used to sample small molecule
67 conformations in an efficient manner^{23–26}. In particular, the work of Klebe and
68 Mietzner²⁶ pave the road for using statistical potentials on torsion angles to
69 generate molecular conformations. Nonetheless, learning these potentials using
70 deep learning confers some advantages such as taking larger portions of the
71 molecule into account or inferring the potential for a combination of atoms not
72 included in the training set. Similar advantages have been observed in deep
73 learning potentials recently used for protein structure prediction²⁷. In this work
74 we show that the proposed potential based on distance likelihood performs
75 similar or better than well-established scoring function for docking and
76 screening tasks. In addition, it can be coupled with global optimization
77 algorithms to reproduce experimental binding conformations of ligands.

78 **Results**

79 **Learning a customized potential based on distance likelihood**

80 Contrary to most computational methods that predict the binding conformation
81 of a ligand (e.g., docking), our geometric deep learning approach learns a
82 potential that is specific for each ligand-target complex and which global
83 minimum corresponds to the optimal binding conformation. To learn this
84 potential we trained DeepDock using experimental three-dimensional data of
85 ligands bound to protein targets (e.g., X-ray crystallography), extracted from the
86 PDBbind database²⁸. DeepDock is a neural network responsible for two main
87 tasks: feature extraction from the input data and identify key ligand-target
88 interactions, as shown in Fig. 1.

89 In a first step, the neural network extracts relevant representations of the input
90 data, namely ligand and target structures. Our approach directly uses the
91 molecular surface of the binding site in the form of a polygon mesh. In this

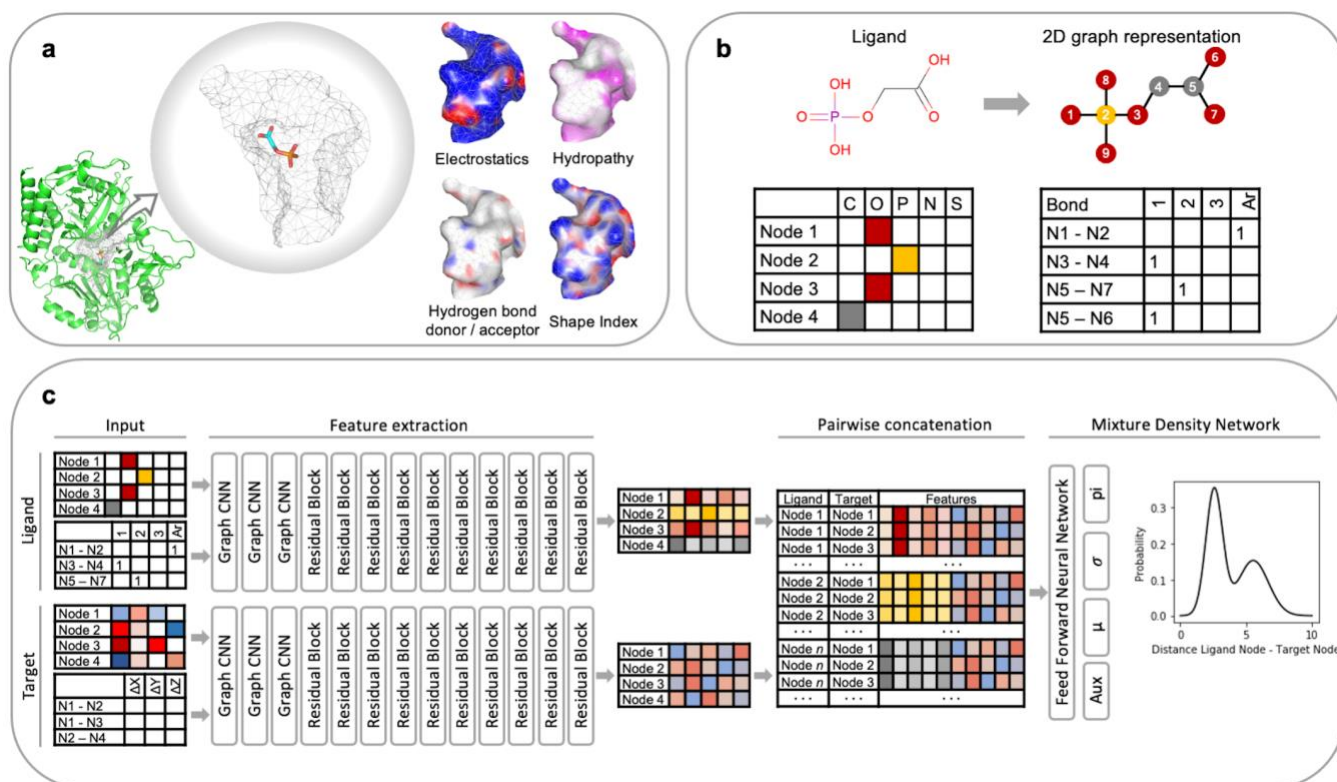
92 mesh, a collection of nodes, edges and faces defines the shape of the
93 molecular surface as a polygon (Fig. 1a). Moreover, the nodes also contain
94 features encoding chemical and topological information at that specific point of
95 the molecular surface, whereas edge features encode the connectivity between
96 nodes. In a similar way, ligands are represented as a two-dimensional
97 undirected graph, where atoms are designated by nodes and bonds are
98 represented by edges (Fig. 1b). In this case, node and edge features encode
99 the atom and bond types, respectively. Both, the target mesh and the ligand
100 graph, are processed by independent residual graph convolutional neural
101 networks (GNNs). Through this procedure, the processed node features not
102 only contain information of an individual atom or point in the molecular surface,
103 but also have information about the other nodes around them. In other words,
104 the processed atom features encode all the atomic environment around a
105 specific atom, whereas the target features encode a patch of the molecular
106 surface around a specific point. A more detailed description of the feature
107 extraction can be found in the Methods section.

108 In a following step, the processed node features from the target and ligand were
109 combined in order to model the interaction of the ligand with the target (Fig. 1c).
110 For this, we concatenate all node features in a pairwise manner meaning each
111 ligand atom will be paired with each node in the molecular surface of the target.
112 In a final step, these concatenated features are processed by mixture density
113 network (MDN)²⁹. This network is composed by a feed forward neural network
114 that predict a set of means, standard deviations and mixing coefficients needed
115 to parametrize a mixture density model for each ligand-target node-pair. The
116 mixture model represents the conditional probability density function of distance
117 for any given ligand-target node pair $P(d_{ij}|v_i^l, v_j^t)$. In other words, using this
118 probability density function we can estimate the likelihood of finding ligand node
119 i separated from a target node j by any distance d_{ij} . Using an MDN is essential
120 since it allows to learn the distribution of distance data i.e., the distribution of all
121 values d_{ij} separating ligand node i from target node j observed in the training
122 set. On the contrary, a simple feed forward neural trained by minimizing the
123 error (e.g., using RMSE or MAE as loss function) only approximates the
124 average distance $\overline{d_{ij}}$ that separates ligand node i from a target node j in the

125 training data. This would be inadequate to model the multi-valued nature of the
 126 data used to train the model since d_{ij} can take an infinite number of valid
 127 values, but some of them are more likely to be observed than others.

128 Finally, the probability density functions of all pairwise combinations of ligand
 129 atoms and points in the molecular surface are aggregated into a statistical
 130 potential. This is simply done by adding up all the independent negative log
 131 likelihood values calculated for each ligand-target pair. This results in an energy
 132 function that can be minimized, and whose minimum correspond to the
 133 conformation of ligand in which all atoms are separated from all points in the
 134 target surface by the most likely distance.

135



136 Fig. 1 Deep learning model used to learn a potential to predict binding
 137 conformations. **a**, The protein target is represented as a polygon mesh of the
 138 molecular surface with four properties encoded in each node, namely
 139 electrostatics, hydrophathy, hydrogen bond donor / acceptor and shape index. **b**,
 140 The ligand is represented as a graph where each node corresponds to one
 141 atom and each edge to one bond. **c**, Ligand and target representations are
 142 processed by a neural network that extracts features using graph convolutions,
 143 which then are pairwise concatenated and used as input of a mixture density
 144 network. As a result, the model predicts a set of probability distributions that are
 145 assembled into a potential.

146 **Potential based on distance likelihood can be used as an accurate scoring**
147 **function**

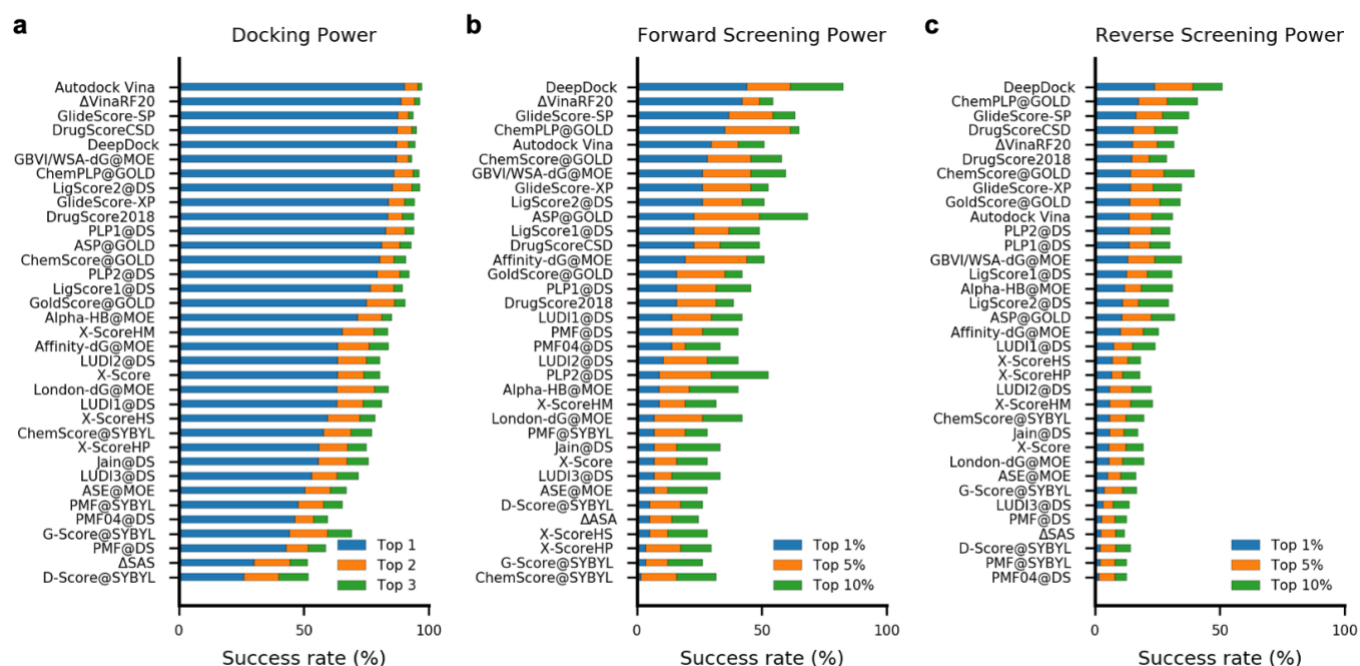
148 We used the CASF-2016 benchmark^{30,31} in order to evaluate if our approach is
149 suitable to be used as an accurate scoring function for an optimization
150 algorithm. The CASF-2016 benchmark is composed of 285 protein-ligand
151 complexes carefully selected to contain diverse proteins in terms of amino acid
152 sequence and unique ligands with a wide binding affinity range. This particular
153 benchmark is designed to assess scoring functions in four demanding tasks,
154 namely scoring power, ranking power, docking power and screening power.
155 Since DeepDock is not specifically trained to predict binding affinities, only the
156 docking and screening power tasks are relevant in this study.

157 The evaluation of docking power measures the ability of a scoring function to
158 identify native ligand binding poses among a set of decoys. For this, the CASF-
159 2016 benchmark provides a set of ~100 decoy conformations for each of the
160 285 ligand-protein complexes, with an RMSD ranging from 0 to 10 Å from the
161 native binding pose. The scoring function under evaluation is used to rank all
162 decoys expecting those with similar conformations to the native ligand-binding
163 pose (i.e. RMSD < 2 Å) to be among the top-ranked. Fig. 2a shows the results
164 of our approach compared to results obtained for other 34 frequently used
165 scoring functions evaluated in the same benchmark by Su et al.³¹. In 87% of the
166 cases, the top ranked decoy using our approach was within an RMSD < 2 Å
167 from the native ligand-binding pose and this amount increased to 94.7% if the 3
168 top-ranked decoys are considered. Based on these results, DeepDock is
169 ranked among the top 5 best performing scoring function in this benchmark,
170 and not far from the best performing scoring function, Autodock Vina, in which
171 the conformation of the best ranked decoy was similar to the native binding
172 pose in 90% of the cases. In addition, it is important to mention that our
173 approach presented a Spearman's rank correlation of 0.83 between the
174 computed score and the decoy RMSD from the native binding pose
175 (Supplementary Fig. 1). In other words, this value indicates that the more similar
176 the decoy conformation is to the native binding pose, the higher the score
177 computed by the scoring function. This correlation has been used as an
178 indicator of the efficiency of a scoring function since it is believed that scoring

179 functions with a high rank correlation can improve conformation sampling to find
180 the native pose³¹.

181 The evaluation of screening power in CASF-2016 is designed to measure the
182 ability of a scoring function to identify true binders of a specific target from a
183 pool of random compounds. The CASF-2016 benchmark is composed of 57
184 protein targets each with a set of 5 true ligands (i.e., a total of 285 compounds)
185 covering a wide range of binding affinity (at least 100-fold). The benchmark
186 provides 100 precomputed binding conformations for each of the 285 ligands in
187 each of the protein target (i.e., 28,500 conformations per target and 1,624,500
188 in total). First, the scoring function is used to assess all conformations in each
189 target, then all compounds are ranked based on the score of their best
190 conformation, and it is expected that true binders are among the top ranked
191 compounds. The ability of distinguishing true binders from random molecules is
192 evaluated using an enhancement factor (EF)^{30,31}. DeepDock presented a mean
193 EF of 16.41 (90% CI, 12.67 - 19.91) when the top 1% ranked compounds are
194 considered, which is the highest compared to other scoring functions previously
195 evaluated in this benchmark (mean EF < 12) as shown in Supplementary Fig. 2.
196 In addition, Fig. 2b shows the success rate of identifying the most potent true
197 binder among the top 1%, 5% or 10% ranked compounds using different
198 scoring functions previously evaluated³¹. Our approach showed the best
199 performance by finding the most potent ligand among the top 1% ranked
200 compounds for 25 protein targets (43.9%), among the top 5% for 35 targets
201 (61.4%) and among the top 10% for 47 targets (82.5%). Other scoring functions
202 among the best performers are $\Delta_{Vina}RF_{20}$, GlideScore-SP, ChemPLP@GOLD
203 and Autodock Vina which were able to rank the most potent ligand among the
204 top 10% ranked compounds for just 37 targets or less (< 65%).

205 The above framework can also be used to evaluate the reverse screening
206 power of a scoring function, that is the ability of identify the real target of a
207 molecule among a set of random targets. Fig. 2c shows the reverse screening
208 power of our approach compared to the performance of other scoring functions
209 previously evaluated³¹. Our approach identified the true target among the 1%
210 top ranked targets for 68 ligands (23.9%), among the 5% for 112 ligands
211 (39.3%) and among the 10% for 145 ligands (50.9%).



213 Fig. 2 Results of the distance likelihood potential in the CASF-2016 benchmark
 214 compared to other scoring functions reported by Su et al.³¹. **a**, Success rate of
 215 detecting real binding pose of a ligand (with an RMSD < 2 Å) among the top 1,
 216 2, and 3 ranked poses during the docking power evaluation task. **b**, Success
 217 rate of detecting the highest affinity ligand for a given target (among the top 1%,
 218 5%, and 10% candidates) during the forward screening task. **c**, Success rate of
 219 detecting the best target protein for a given ligand (among the top 1%, 5%, and
 220 10% possible targets) during the reverse screening task.

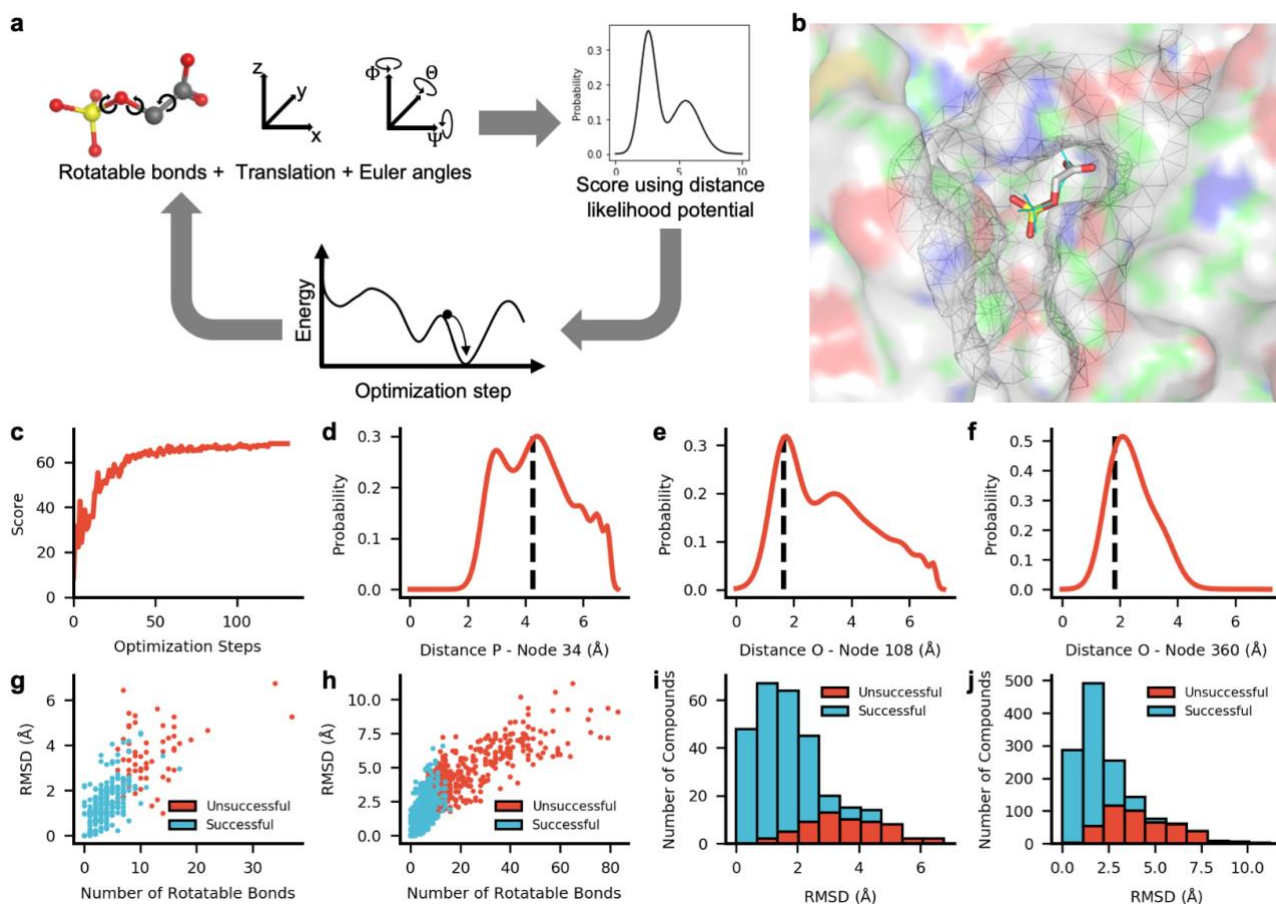
221

222 Potential based on distance likelihood can reproduce experimental 223 binding conformations

224 An advantage of this deep learning approach is that it can be easily combined
 225 with optimization algorithms in order to find the ligand conformation associated
 226 with the global minimum of the potential. In other words, it can find the ligand
 227 conformation with highest likelihood of binding. For this, the optimization
 228 algorithm carefully rotates each rotatable bond in the molecule, and at the same
 229 time it translates and rotates the whole ligand using a transformation matrix until
 230 it finds the conformation that best fits the binding pocket (Fig. 3a). In this case
 231 we used differential evolution³² as the optimization algorithm, but others such as
 232 particle swarm optimization (PSO), simulated annealing (SA), or even gradient-
 233 based algorithms can be adapted to DeepDock. For example, gradient descent

234 has recently been used to minimize a potential learnt by neural networks in
235 order to predict protein structures²⁷.

236 To assess if an optimization algorithm can minimize a potential based on
237 distance likelihood, we tried to reproduce real binding poses of different ligand-
238 target when starting from a random conformation and position. For a more
239 realistic test, this was done using the 285 ligand-target pairs in the CASF-2016
240 coresets plus 1,367 ligand-target pairs used as the validation set. It is worth
241 mentioning that none of these complexes were included in the training set.
242 DeepDock was able to find conformations corresponding to a minimum for 225
243 (87%) of the compounds in the CASF-2016 coresets and for 917 (67%) of the
244 molecules in the validation set. Interestingly, the optimization failed for most of
245 the compounds with more than 10 rotatable bonds (Fig. 3g-h). The effect of the
246 number of rotatable bonds on optimization has been noticed before and is
247 linked to the inefficiency of the optimization algorithms when dealing with a
248 large number of degrees of freedom³³. In general, all compounds for which the
249 optimization finished correctly presented a conformation very similar to the real
250 binding pose, that is, a median (IQR) RMSD of 1.33 (0.81 to 1.99) Å for the
251 CASF-2016 molecules and a median (IQR) RMSD of 1.47 (1.00 to 2.11) Å for
252 molecules in the validation set (Fig. 3i-j). These similarities are also evident
253 from the high correlation between the scores produced by the predicted and the
254 real binding pose ($R^2 = 0.81$ for CASF-2016 coresets and $R^2 = 0.82$ for the
255 validation set). It is important to mention that no correlation was found between
256 the compound binding affinity and the score of the predicted or real binding
257 pose (Supplementary Fig. 3-4).



259 Fig. 3 Use of distance likelihood potential to predict ligand binding
 260 conformations. **a**, Representation of the optimization process where a ligand
 261 conformation is represented by a vector of the values of all rotatable bonds in
 262 the molecule, the displacement across the three dimensions of the Euclidean
 263 space and the three Euler angles that represent the rotation of the molecule.
 264 This conformation is scored using the distance likelihood potential and then
 265 optimized using differential evolution to produce a new conformation, which
 266 follows the same procedure until the optimization has successfully finished. **b**,
 267 Example of the predicted binding conformation of 2-phosphoglycolic acid to rat
 268 PEPCK (PDB ID: 2RKA). Experimental binding conformation is depicted in
 269 cyan lines and the polygon mesh in gray lines. **c**, Optimization process of 2-
 270 phosphoglycolic acid to rat PEPCK. **d-f**, Examples of predicted distance
 271 probability distributions between ligand atoms and target nodes for 2RKA. The
 272 dashed line indicates the distance of between ligand atoms and target node for
 273 the predicted binding conformation. **g-h**, Scatter plots for 285 compounds in
 274 CASF-2016 (**g**) and 1,367 compounds in the validation set (**h**) showing that
 275 RMSD between predicted and experimental binding conformations is lower for
 276 compounds with less rotatable bonds. The optimization using differential
 277 evolution successfully finished for most compounds bearing less than 10
 278 rotatable bonds. **g-h**, Distributions of RMSD between predicted and
 279 experimental binding conformations in CASF-2016 (**g**) and in the validation set
 280 (**h**). Color code refers to compounds for which the optimization successfully
 281 finished or not.

282 **Conclusions**

283 In this work, we report a method that exploits geometric deep learning to predict
284 ligand binding conformations. Contrary to docking methods where a one-fits-all
285 scoring function is used, here a deep neural network learns a potential that is
286 specific for each protein-ligand complex, which is then used to find the optimal
287 binding conformation. In a first instance, the deep neural network learns the
288 parameters of a mixture model that is employed as a probability density
289 function. This probability density function is used to determine the most likely
290 distance separating a ligand atom from a specific point in the molecular surface
291 of the binding site. The potential is determined as the combination of the
292 negative log likelihood of all pairwise combination of ligand atoms and points in
293 the molecular surface. The optimal conformation is the one that minimizes the
294 potential, that is, the ligand conformation in which every atom is separated from
295 the target surface by the most likely distance. We demonstrate that this
296 potential can be used as an accurate scoring function for molecular docking and
297 virtual screening. In fact, this potential performs equally or better than many of
298 the most widely used scoring functions in the CASF-2016 benchmark. It is
299 important to mention that this benchmark only contains a small number of
300 compounds compared to real screening libraries (usually composed by millions
301 of molecules). Finally, we also show that this potential can be minimized using
302 global optimization methods, such as differential evolution³², in order to find the
303 most likely binding conformation of a ligand. More in concrete, we reproduce the
304 binding conformation of 868 ligands (177 from CASF-2016 and 691 from the
305 test set) with an RMSD $< 2 \text{ \AA}$ from the experimental structure using the method
306 described in this paper. Overall, we have presented evidence that geometric
307 deep learning can be used to predict the binding conformation of ligands to their
308 biological target. Although the results presented in this work were mainly
309 focused on small molecules, similar approaches can be used to predict binding
310 conformations of larger molecules such as peptides or even protein-protein
311 interactions. We anticipate that further developments in geometric deep
312 learning will help to significantly improve and speed up structure-based virtual
313 screening.

314

315 **Methods**

316 **Data set:** The model reported in this study was trained using the general set of
317 the PDBbind database (v.2019)²⁸, which contains a collection of 17,679 protein-
318 ligand structures with their respective potency (e.g., IC₅₀, K_d, etc). From these,
319 we removed those complexes that are included in the CASF-2016 benchmark
320 and those that failed during the pre-processing step leaving a total 16,367
321 protein-ligand complexes which were randomly divided in a training set
322 containing 15,000 complexes and a test set with 1,367. Each of these
323 complexes was processed in order to be used as an input for the model.

324 The chemical structures of ligands were represented as undirected graphs $\mathcal{G}^l =$
325 $(\mathcal{V}^l, \mathcal{E}^l)$ where nodes $v_i^l \in \mathcal{V}^l$ represent atoms in the molecule and edges $e_{i,j}^l \in$
326 \mathcal{E}^l represent bonds. In this case each node v_i^l is represented by a one-hot
327 vector that indicates the atom type among 28 possibilities (Be, B, C, N, O, F,
328 Mg, Si, P, S, Cl, V, Fe, Co, Cu, Zn, As, Se, Br, Ru, Rh, Sb, I, Re, Os, Ir, Pt, Hg).
329 Similarly, each edge $e_{i,j}^l$ is represented by a one-hot vector that indicates the
330 bond type, either single, double, triple or aromatic. It is important to mention that
331 no information regarding the three-dimensional conformation of the ligand was
332 used for training the model.

333 The protein targets were processed using a pipeline based on the one
334 previously described by Gainza et al.³⁴. As in MaSIF, protein surfaces were
335 triangulated using MSMS³⁵ with a density of $3.0 \text{ nodes}/\text{\AA}^2$ and a probe radius of
336 1.5 \AA . The resulting meshes were down sampled to a resolution of 1 \AA and
337 processed using pymesh. The resulting mesh $\mathcal{G}^t = (\mathcal{V}^t, \mathcal{E}^t)$ is composed of a
338 fixed set of nodes $v_i^t \in \mathcal{V}^t$ and edges $e_{i,j}^t \in \mathcal{E}^t$. Each node v_i^t is represented by
339 vector of four features calculated using MaSIF, namely, Poisson-Boltzmann
340 continuum electrostatics, free electrons and proton donors, hydrophathy and
341 shape index. In a similar way, each edge $e_{i,j}^t$ is represented by a vector defining
342 the relative cartesian coordinates of the linked nodes i.e., $r_{i,j} = (p_i - p_j)$ where
343 $p_i \in \mathbb{R}^3$ represents the coordinates of v_i^t in a three-dimensional Euclidean
344 space. It is worth mentioning that only nodes defining the binding site (i.e.,
345 within 10 \AA or less from any ligand atom) were used to train the model.

346 **Model:** The model construction can be divided into three stages: feature
 347 extraction, feature concatenation and a mixed density network. In the first
 348 stage, features are extracted by two independent residual graph convolutional
 349 neural networks (GNNs), one for the ligand and the other for the target. Despite
 350 being independent, both residual GNN have the same architecture. First, the
 351 node and edge features are projected to a 128-dimensional embedding using a
 352 linear layer as in Eqs. (1) and (2). Then we used a sequence of three GNNs
 353 to update each node and edge based on their neighbouring nodes and the type
 354 of edges connecting them. The GNN first updates each edge in the graph by
 355 applying a multi-layer perceptron (MLP) on the concatenation of the edge
 356 features and the features of the two connecting nodes as shown in Eq. (3).
 357 The updated edge features $e_{i,j}^\ell$ are used to update the node features as shown
 358 in Eq. (4). The updated edge and node features ($e_{i,j}^\ell$ and v_i^ℓ , respectively)
 359 contain information of the central atom but also of the neighbouring atoms
 360 around it and can be used as input of another convolution round (Eqs. (3) and
 361 (4)). In this case, we used three convolutions, i.e. up to $\ell = 3$.

$$e_{i,j}^0 = Linear(e_{i,j}) \quad (1)$$

$$v_i^0 = Linear(v_i) \quad (2)$$

$$e_{i,j}^\ell = MLP([v_i^{\ell-1}, v_j^{\ell-1}, e_{i,j}^{\ell-1}]) \quad (3)$$

$$v_i^\ell = MLP\left(\left[v_i^{\ell-1}, \frac{1}{\|j\|} \sum_j MLP([v_j^{\ell-1}, e_{i,j}^\ell])\right]\right) \quad (4)$$

362 After the initial processing by the GNNs, the node and edge features were
 363 processed by 10 residual GNN blocks. Each residual block starts by projecting
 364 the node and edge features (v_i^{h-1} and $e_{i,j}^{h-1}$, respectively) to a 32-dimensional
 365 vector using an MLP as shown in Eqs. (5) and (6). The resulting vectors are
 366 used as inputs to a GNN (Eq. (7)) resulting in aggregated node and edge
 367 features v_i' and $e_{i,j}'$, respectively, which are projected back to 128-dimensional
 368 vectors v_i^{up} and $e_{i,j}^{up}$ as shown in Eqs. (8) and (9). Finally, the resulting
 369 vectors (v_i^{up} and $e_{i,j}^{up}$) are added to the input vectors (v_i^{h-1} and $e_{i,j}^{h-1}$) to create

370 a skip connection and later modified by an activation function (Eqs. (10) and (
 371 11)).

$$e_{i,j}^{down} = MLP(e_{i,j}^{h-1}) \quad (5)$$

$$v_i^{down} = MLP(v_i^{h-1}) \quad (6)$$

$$v'_i, e'_{i,j} = GNN(v_i^{down}, v_j^{down}, e_{i,j}^{down}) \quad (7)$$

$$e_{i,j}^{up} = Dropout\left(BatchNorm\left(Linear(e'_{i,j}) \right) \right) \quad (8)$$

$$v_i^{up} = Dropout\left(BatchNorm\left(Linear(v'_i) \right) \right) \quad (9)$$

$$e_{i,j}^h = ELU(e_{i,j}^{h-1} + e_{i,j}^{up}) \quad (10)$$

$$v_i^h = ELU(v_i^{h-1} + v_i^{up}) \quad (11)$$

372 The extracted node features by the GNNs and residual GNNs for both target \bar{v}_r^t
 373 and ligand \bar{v}_s^l are then pairwise concatenated and used as input of a mixture
 374 density network (MND)²⁹. The MND uses an MLP to create a hidden
 375 representation $h_{r,s}$ that combines the concatenated target and ligand node
 376 information as shown in Eq. (12). The hidden representation is used to
 377 compute the outputs of the MND, which consist of the means ($\mu_{r,s}$), standard
 378 deviations ($\sigma_{r,s}$) and mixing coefficients ($\alpha_{r,s}$) that are necessary to parametrize
 379 a mixture of gaussians (Eqs. (13)- (15)). In this particular case, the mixture
 380 model uses 10 gaussians to simulate the probability density distribution of the
 381 distance between the ligand and target nodes (\bar{v}_s^l and \bar{v}_r^t , respectively).

$$h_{r,s} = Dropout\left(MLP([\bar{v}_r^t, \bar{v}_s^l]) \right) \quad (12)$$

$$\mu_{r,s} = ELU\left(Linear(h_{r,s}) \right) + 1 \quad (13)$$

$$\sigma_{r,s} = ELU\left(Linear(h_{r,s}) \right) + 1 \quad (14)$$

$$\alpha_{r,s} = Softmax\left(Linear(h_{r,s}) \right) \quad (15)$$

382 In addition, the extracted ligand node features \bar{v}_s^l were used to predict auxiliary
 383 tasks, namely atom type and bond type with connecting neighbouring nodes.
 384 These auxiliary tasks help to learn molecular structures which accelerates
 385 training. All MLPs used are composed by a linear layer followed by batch
 386 normalization and an ELU activation function. The dropout rate used was of 0.1
 387 in all experiments.

388 **Training:** We employed the Adam optimizer with a learning rate of 0.002 to
 389 update model weights. The model was trained to minimize the loss function
 390 shown in Eq. (16) where \mathcal{L}_{MDN} represents the loss of the mixture density
 391 network whereas \mathcal{L}_{atoms} and \mathcal{L}_{bonds} are the cross-entropy cost functions of
 392 predicting atom and bond types, respectively, that were used as auxiliary tasks.
 393 In particular, \mathcal{L}_{MDN} minimizes the negative log-likelihood of $d_{r,s}$, which
 394 represents the distance separating the target node v_r^t from the ligand node v_s^l ,
 395 computed using the mixture model formed by $k = 10$ gaussians and
 396 parametrised by $\alpha_{r,s}$, $\mu_{r,s}$ and $\sigma_{r,s}$ that were predicted by the model (Eq. (17)).
 397 The model was trained for 150 epochs using a batch size of 16 protein-ligand
 398 complexes. Contributions of ligand-target node pairs separated by a $d_{r,s} > 7 \text{ \AA}$
 399 were masked since we considered that those atoms cannot form relevant
 400 interactions.

$$\mathcal{L}_{total} = \mathcal{L}_{MDN} + \mathcal{L}_{atoms} + \mathcal{L}_{bonds} \quad (16)$$

$$\mathcal{L}_{MDN} = -\log P(d_{r,s} | v_r^t, v_s^l) = -\log \sum_{k=1}^K \alpha_{r,s,k} \mathcal{N}(d_{r,s} | \mu_{r,s,k}, \sigma_{r,s,k}) \quad (17)$$

$$U_{(x)} = -\sum_{r=1}^R \sum_{s=1}^S \log P(d_{r,s} | v_r^t, v_s^l) \quad (18)$$

401 The loss function shown in Eq. (17) can be easily used to define a potential
 402 $U_{(x)}$ which is tailored for a particular target-ligand complex (Eq. (18)). It is
 403 possible to use this potential to score the 3D structure of a target-ligand
 404 complex by computing the distances $d_{r,s}$ separating each target node v_r^t from
 405 each ligand node v_s^l in that specific conformation, calculating the negative log
 406 likelihood $-\log P(d_{r,s} | v_r^t, v_s^l)$ for each target-ligand node pair, and summing

407 across all possible pairs. The lower the value of $U_{(x)}$, the more likely is to find
408 the target-ligand complex in that specific conformation.

409 **Benchmark:** This approach was evaluated using the CASF-2016
410 benchmark^{30,31}, which contains 285 protein-ligand complexes carefully curated.
411 Structures from this benchmark were preprocessed in the same way as the
412 training set. There are four different tasks in this benchmark in order to evaluate
413 the scoring power, ranking power, docking power and screening power of a
414 scoring function. Only the docking and screening power are relevant for this
415 evaluation. The former evaluates the ability of the scoring function to identify the
416 real binding conformation among generated decoy conformations of the same
417 ligand. The latter evaluates if the scoring function can identify true binders for a
418 particular target using the enhancement factor (Eq. (19)) as metric. Results can
419 be directly compared with other scoring functions previously evaluated by Su et
420 al.³¹. The complete protocol and scripts are fully described in the original
421 publication³⁰.

$$EF_{1\%} = \frac{\text{Number of true binders detected among 1\%}}{(\text{Total true binders}) \times 1\%} \quad (19)$$

422 **Prediction of binding conformations:** We represent the ligand conformation
423 as vector of the Euler angels, the relative position of the ligand in the Euclidean
424 space and the dihedral angles of all rotatable bonds in the molecule. We
425 employed differential evolution³² to find the ligand conformation that minimize
426 the potential $U_{(x)}$ learnt by the model for that specific complex, that is, the
427 resulting ligand conformation will be the most likely to interact with target
428 binding site according to the model. We run the global optimization for a
429 maximum of 500 iterations using a population size of 150, the mutation constant
430 was randomly changed each generation from a (0.5, 1) interval and with a
431 recombination constant of 0.8. The values of dihedrals of rotatable bonds and
432 Euler angles were restricted to be between $-\pi$ and π . In interest of
433 reproducibility, the calculation was seeded but this is not a requirement.

434

435

436 **Author contributions**

437 O.M.L. conceived the idea, wrote the code, performed the experiments and
438 wrote the paper. M.A., A.E.D.C. and J.K.W. helped with the preparation of the
439 manuscript and with insightful discussions. A.E.D.C. helped to improve code.

440 **Acknowledgments**

441 Authors thank Dries Van Rompaey, Jonas Verhoeven and Natalia Dyubankova
442 for supporting this project.

443 **Competing interests**

444 O.M.L., M.A. and J.K.W. are employees of Janssen Pharmaceutica N.V.

445 **Data availability**

446 The data that support the findings of this study will be available after the final
447 publication of this manuscript.

448 **Code availability**

449 The code used to generate results shown in this study will be available after the
450 final publication of this manuscript.

451 **References**

- 452 1. Hert, J., Irwin, J. J., Laggner, C., Keiser, M. J. & Shoichet, B. K.
453 Quantifying biogenic bias in screening libraries. *Nat. Chem. Biol.* **5**, 479–
454 483 (2009).
- 455 2. Dobson, C. M. Chemical space and biology. *Nature* **432**, 824–828 (2004).
- 456 3. Congreve, M., Murray, C. W. & Blundell, T. L. Keynote review: Structural
457 biology and drug discovery. *Drug Discov. Today* **10**, 895–907 (2005).
- 458 4. Klebe, G. Protein--Ligand Interactions as the Basis for Drug Action. in
459 *Drug Design: Methodology, Concepts, and Mode-of-Action* (ed. Klebe, G.)
460 61–88 (Springer Berlin Heidelberg, 2013). doi:10.1007/978-3-642-17907-
461 5_4.
- 462 5. Renaud, J. P. *et al.* Cryo-EM in drug discovery: Achievements, limitations
463 and prospects. *Nat. Rev. Drug Discov.* **17**, 471–492 (2018).
- 464 6. Gorgulla, C. *et al.* An open-source drug discovery platform enables ultra-
465 large virtual screens. *Nature* **580**, 663–668 (2020).
- 466 7. Lyu, J. *et al.* Ultra-large library docking for discovering new chemotypes.
467 *Nature* **566**, 224–229 (2019).
- 468 8. Krivák, R. & Hoksza, D. P2Rank: machine learning based tool for rapid
469 and accurate prediction of ligand binding sites from protein structure. *J.*
470 *Cheminform.* **10**, 1–12 (2018).
- 471 9. Pu, L., Govindaraj, R. G., Lemoine, J. M., Wu, H. C. & Brylinski, M.
472 Deepdrug3D: Classification of ligand-binding pockets in proteins with a
473 convolutional neural network. *PLoS Comput. Biol.* **15**, 1–23 (2019).
- 474 10. Jiménez, J., Doerr, S., Martínez-Rosell, G., Rose, A. S. & De Fabritiis, G.
475 DeepSite: Protein-binding site predictor using 3D-convolutional neural
476 networks. *Bioinformatics* **33**, 3036–3042 (2017).
- 477 11. Ain, Q. U., Aleksandrova, A., Roessler, F. D. & Ballester, P. J. Machine-
478 learning scoring functions to improve structure-based binding affinity

- 479 prediction and virtual screening. *Wiley Interdiscip. Rev. Comput. Mol. Sci.*
480 **5**, 405–424 (2015).
- 481 12. Li, H., Sze, K. H., Lu, G. & Ballester, P. J. Machine-learning scoring
482 functions for structure-based virtual screening. *Wiley Interdiscip. Rev.*
483 *Comput. Mol. Sci.* **11**, (2021).
- 484 13. Sánchez-Cruz, N., Medina-Franco, J. L., Mestres, J. & Barril, X. Extended
485 connectivity interaction features: improving binding affinity prediction
486 through chemical description. *Bioinformatics* **btaa982**, (2020).
- 487 14. Wójcikowski, M., Ballester, P. J. & Siedlecki, P. Performance of machine-
488 learning scoring functions in structure-based virtual screening. *Sci. Rep.*
489 **7**, 1–10 (2017).
- 490 15. Wójcikowski, M., Kukielka, M., Stepniewska-Dziubinska, M. M. &
491 Siedlecki, P. Development of a protein-ligand extended connectivity
492 (PLEC) fingerprint and its application for binding affinity predictions.
493 *Bioinformatics* **35**, 1334–1341 (2019).
- 494 16. Ballester, P. J. & Mitchell, J. B. O. A machine learning approach to
495 predicting protein-ligand binding affinity with applications to molecular
496 docking. *Bioinformatics* **26**, 1169–1175 (2010).
- 497 17. Ragoza, M., Hochuli, J., Idrobo, E., Sunseri, J. & Koes, D. R. Protein-
498 Ligand Scoring with Convolutional Neural Networks. *J. Chem. Inf. Model.*
499 **57**, 942–957 (2017).
- 500 18. Stepniewska-Dziubinska, M. M., Zielenkiewicz, P. & Siedlecki, P.
501 Development and evaluation of a deep learning model for protein–ligand
502 binding affinity prediction. *Bioinformatics* **34**, 3666–3674 (2018).
- 503 19. Hassan-Harrirou, H., Zhang, C. & Lemmin, T. RosENet: Improving
504 Binding Affinity Prediction by Leveraging Molecular Mechanics Energies
505 with an Ensemble of 3D Convolutional Neural Networks. *J. Chem. Inf.*
506 *Model.* **60**, 2791–2802 (2020).
- 507 20. Jiménez, J., Škalič, M., Martínez-Rosell, G. & De Fabritiis, G. KDEEP:

- 508 Protein-Ligand Absolute Binding Affinity Prediction via 3D-Convolutional
509 Neural Networks. *J. Chem. Inf. Model.* **58**, 287–296 (2018).
- 510 21. Feinberg, E. N. *et al.* PotentialNet for Molecular Property Prediction. *ACS*
511 *Cent. Sci.* **4**, 1520–1530 (2018).
- 512 22. Lim, J. *et al.* Predicting Drug-Target Interaction Using a Novel Graph
513 Neural Network with 3D Structure-Embedded Graph Representation. *J.*
514 *Chem. Inf. Model.* **59**, 3981–3988 (2019).
- 515 23. Gasteiger, J., Rudolph, C. & Sadowski, J. Automatic generation of 3D-
516 atomic coordinates for organic molecules. *Tetrahedron Comput.*
517 *Methodol.* **3**, 537–547 (1990).
- 518 24. Velec, H. F. G., Gohlke, H. & Klebe, G. DrugScoreCSDKnowledge-Based
519 Scoring Function Derived from Small Molecule Crystal Data with Superior
520 Recognition Rate of Near-Native Ligand Poses and Better Affinity
521 Prediction. *J. Med. Chem.* **48**, 6296–6303 (2005).
- 522 25. Fan, H. *et al.* Statistical Potential for Modeling and Ranking of Protein–
523 Ligand Interactions. *J. Chem. Inf. Model.* **51**, 3078–3092 (2011).
- 524 26. Klebe, G. & Mietzner, T. A fast and efficient method to generate
525 biologically relevant conformations. *J. Comput. Aided. Mol. Des.* **8**, 583–
526 606 (1994).
- 527 27. Senior, A. W. *et al.* Improved protein structure prediction using potentials
528 from deep learning. *Nature* **577**, 706–710 (2020).
- 529 28. Liu, Z. *et al.* PDB-wide collection of binding data: Current status of the
530 PDBbind database. *Bioinformatics* **31**, 405–412 (2015).
- 531 29. Bishop, C. M. Mixture Density Networks. *Tech. Report. Ast. Univ.*
532 *Birmingham.* (1994).
- 533 30. Li, Y. *et al.* Assessing protein-ligand interaction scoring functions with the
534 CASF-2013 benchmark. *Nat. Protoc.* **13**, 666–680 (2018).
- 535 31. Su, M. *et al.* Comparative Assessment of Scoring Functions: The CASF-

- 536 2016 Update. *J. Chem. Inf. Model.* **59**, 895–913 (2019).
- 537 32. Storn, R. & Price, K. Differential Evolution – A Simple and Efficient
538 Heuristic for global Optimization over Continuous Spaces. *J. Glob. Optim.*
539 **11**, 341–359 (1997).
- 540 33. Li, H., Leung, K. S., Ballester, P. J. & Wong, M. H. Istar: A web platform
541 for large-scale protein-ligand docking. *PLoS One* **9**, (2014).
- 542 34. Gainza, P. *et al.* Deciphering interaction fingerprints from protein
543 molecular surfaces using geometric deep learning. *Nat. Methods* **17**, 184–
544 192 (2020).
- 545 35. Sanner, M. F., Olson, A. J. & Spehner, J. C. Reduced surface: An efficient
546 way to compute molecular surfaces. *Biopolymers* **38**, 305–320 (1996).
- 547