

Docking-informed machine learning for kinome-wide affinity prediction

Jordy Schifferstein^{1,2}, Andrius Bernatavicius³, Antonius P.A. Janssen^{*1,2}

¹ Department of Molecular Physiology, Leiden Institute of Chemistry, Leiden University, The Netherlands

² Oncode Institute, The Netherlands

³ Leiden Institute of Advanced Computer Science, Leiden University, the Netherlands

* E-mail: a.p.a.janssen@lic.leidenuniv.nl

Abstract

Kinase inhibitors are an important class of anti-cancer drugs, with 80 inhibitors clinically approved, and >100 in active clinical testing. Most bind competitively in the ATP-binding site, leading to challenges with selectivity for a specific kinase, resulting in risks for toxicity and general off-target effects. Assessing the binding of an inhibitor for the entire kinome is experimentally possible but expensive. A reliable and interpretable computational prediction of kinase selectivity would greatly benefit the inhibitor discovery and optimisation process. Here, we use machine learning on docked poses to address this need. To this end, we aggregated all known inhibitor-kinase affinities and generated the complete accompanying 3D interactome by docking all inhibitors to the respective high quality X-ray structures. We then used this resource to train a neural network as a kinase-specific scoring function, which achieved an overall performance (R^2) of 0.63-0.74 on unseen inhibitors across the kinome. The entire pipeline from molecule to 3D-based affinity prediction has been fully automated and wrapped in a freely available package. This has a graphical user interface which is tightly integrated with PyMOL to allow immediate adoption in the medicinal chemistry practice.

Introduction

Protein kinases are one of the main protein families targeted by anti-cancer drugs, with 80 approved drugs and around 150 in clinical testing.¹ However, current FDA-approved kinase inhibitors are designed to target only a few percent of the entire protein family.² The so-far untargeted kinases, thus, offer great opportunities for the development of novel molecular therapies.

The chances of success for any drug greatly depend on two parameters: affinity of the drug for the intended target protein, and selectivity over the rest of the protein family. Off-target activity is often the main cause of (pre-)clinical toxicity, and side-effects in general.³ This issue is particularly pressing for kinase inhibitors, as these in most cases target the ATP-binding site of the protein, which is highly conserved across this large protein family.⁴ This leads to many kinase inhibitors potentially binding to many family members, sometimes inhibiting as much as 70% of all kinases.⁵ Determining the specificity of an inhibitor over all ± 500 kinases is experimentally feasible, but is prohibitively expensive in terms of time, material and funds to perform on a routine basis.

In recent years, various computational methods of predicting kinase inhibitor selectivity have thus been developed.⁶⁻⁸ Approaches vary from 'classical' protein structure-based techniques such as molecular docking to machine learning approaches such as Quantitative Structure Activity Relationship (QSAR) studies. The revolution of artificial intelligence (AI) has not gone unnoticed in this field, and *e.g.* AlphaFold⁹ will have a tremendous impact in the coming years, giving direct access to structures for all proteins. Structure-based methods typically rely on either classical physics-based scoring functions to 'score' a generated protein-ligand complex. More recently, machine learning-based scoring functions such as RFScore have reached state-of-the-art performance.¹⁰⁻¹² These scoring functions were trained on experimental datasets such as the PDBbind, offering a relatively broad set of protein-inhibitor complexes and their bioactivity data.¹³

We set out to develop a fully automated docking-based affinity prediction for kinases. As it is generally accepted that pose finding for most docking algorithms is very good¹⁴, we envisioned that a large docking-based protein-inhibitor dataset for which biochemical data is known should also function as a basis for training a scoring function. We demonstrated this approach by generating protein-inhibitor complexes for all kinase inhibitors in the Papyrus dataset¹⁵, a large aggregation of literature binding data, for kinases of which a high-quality experimental protein structure is available in the KLIFS database, a kinase specific mirror of the PDB.^{18,19} We used two docking algorithms: Autodock VinaGPU¹⁶ and DiffDock¹⁷. This generated database is then used to train a multi-layered Neural Network as scoring function, that shows excellent performance on bio-activity predictions for unseen inhibitors. The automated workflow has been wrapped in an easily installable Docker container²⁰ with a convenient PyMOL Graphical User Interface (GUI) plugin, allowing broad access to the methodology.

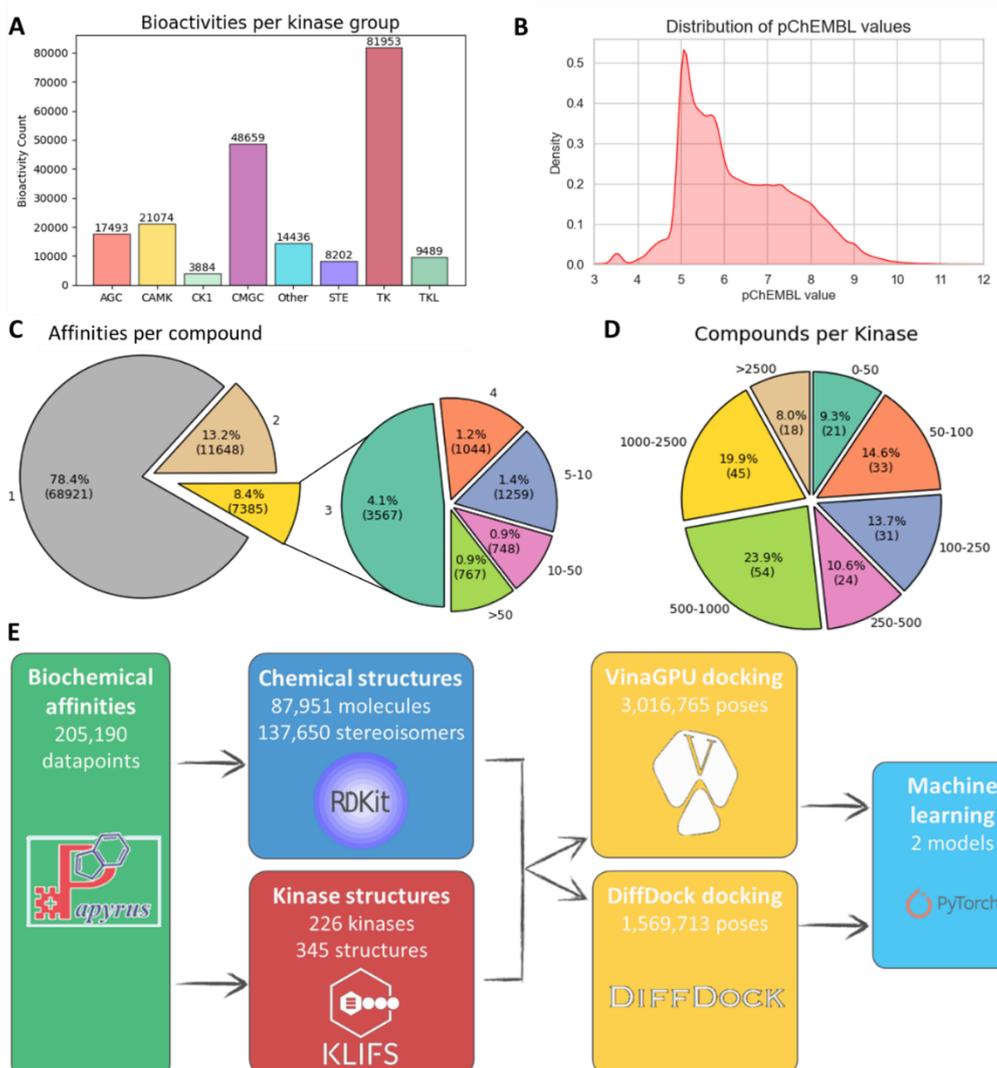
52 **Results**

53 **Extracting literature biochemical and structural data**

54 To generate our desired docking-based training dataset, we first needed to select all kinases of which we have a
 55 high-quality experimental structure. As a source of well curated and annotated kinase protein structures
 56 available in the Protein Data Bank we used the KLIFS database. These structures were filtered based on their
 57 resolution ($\leq 2.5 \text{ \AA}$) and KLIFS quality metric (≥ 8). We selected the best of each of the four possible combinations
 58 of DFG-in/out and α -C helix in/out as annotated in the KLIFS database. In total, this led to 345 protein structures
 59 for 226 unique kinases.

60 Next, we extracted all reported inhibitory activities for these kinases in the Leiden Papyrus dataset, a
 61 curated resource combining data from resources such as ChEMBL, PubChem and others. We chose to
 62 indiscriminately use pIC_{50} , pK_i and pK_d values, collectively from hereon: pChEMBL values. We filtered the
 63 compounds to entail only the more drug-like small molecules using quite lax criteria ($MW \leq 750$, $NumHBD \leq 10$,
 64 $NumHBA \leq 15$, $Rotatable \text{ Bonds} \leq 15$), which should have reasonable chance to dock well and form a
 65 representative training set for real world medicinal chemistry applications. An overview of the resulting
 66 physicochemical properties and chemical diversity is plotted in Supplementary Figure 1.

67 This procedure led to a completed dataset of in total 205,190 affinity values for 87,951 unique compounds
 68 against 226 unique kinases. A summative view of the workflow and complete resulting dataset is depicted in
 69 Figure 1 and Supplementary Figure 1.



70 **Figure 1: Kinase activity dataset** | A) Distribution of kinase inhibition values reported per kinase group; B) Distribution of the
 71 reported inhibitory values; C) Number of pChEMBL values for unique kinases reported per kinase inhibitor, *i.e.*, against how
 72 many kinases was a compound tested; D) Number of reported inhibitors per kinase, *i.e.*, how many compounds were tested
 73 for a kinase; E) Overview of the workflow of the work in this paper.
 74

75 Large scale Molecular Docking using Open Source software

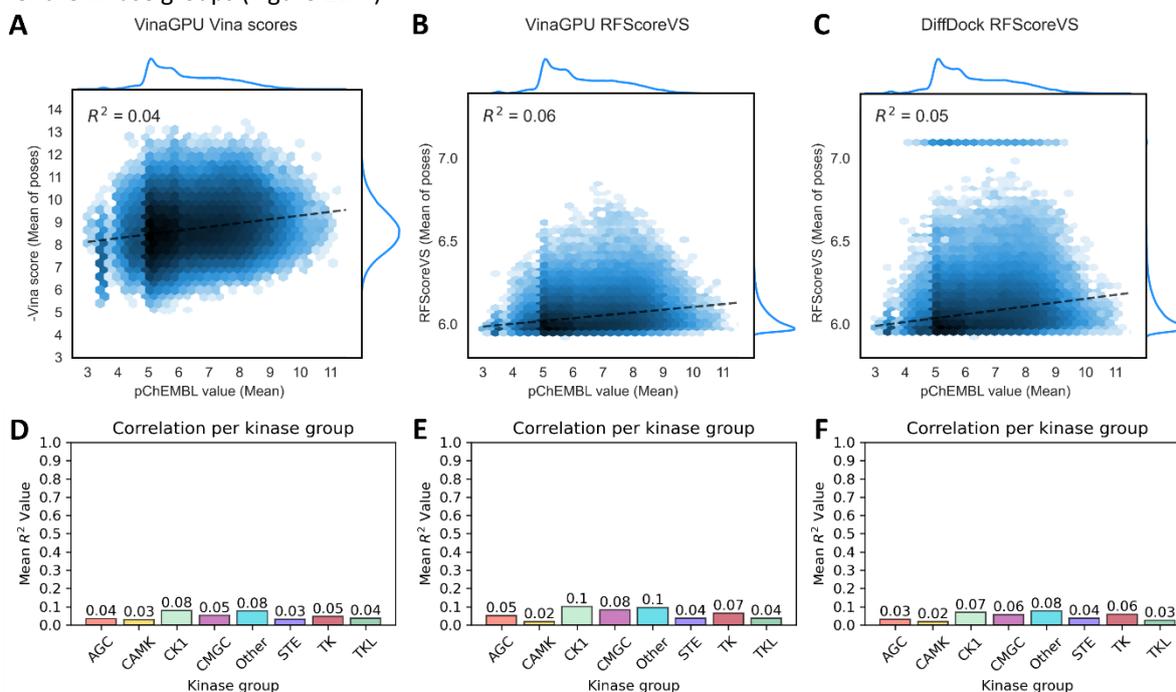
76 We set up an automated docking pipeline to generate a set of docking poses for all inhibitor-protein pairs in the
77 created dataset (Figure 1E). To this end, inhibitors were prepared for docking using an RDKit²¹ pipeline, which
78 enumerates potential stereo- and double bond isomers, and generates a 3D conformer. For each protein
79 structure, a binding site was defined using PyVOL to guide the VinaGPU docking algorithm.²² All prepared isomers
80 were consecutively docked in the known targets of these inhibitors using two docking algorithms: Autodock
81 VinaGPU and the diffusion-based DiffDock algorithm (version of December 2022).

82 For all compound-protein structure pairs, a maximum of 5 poses were generated. The poses were filtered
83 for excessive atomic overlap based on a tailored clash-score (see Methods and Supplementary Figure 2) to get
84 rid of unphysical poses generated, a problem especially prevalent in DiffDock generated poses. For the inhibitor-
85 kinase pairs in our dataset for which an experimental pose has been determined (only $\pm 0.2\%$ of the 205,000),
86 the root mean squared deviation (RMSD) was calculated for both docking algorithms. Median \pm absolute
87 deviation for DiffDock and VinaGPU were 1.296 ± 0.587 and 5.659 ± 4.177 , respectively.

88 The results of this large-scale docking project were aggregated and have been made available in an SQLite
89 database that holds all activities, compounds, isomers, protein information, kinase structure information and all
90 poses for both docking tools. A simplified schema of this database with statistics per table is depicted in
91 Supplementary Figure 3A. The database includes all .mol-formatted poses in a compressed format, as well as the
92 .pdb files for all KLIFS structures used. The database was designed to be readily usable for machine learning
93 applications. Additionally, a KNIME-based user interface has been built to browse and query the generated
94 docking complexes (Supplementary Figure 3B). The full database and accompanying application are freely
95 available on Zenodo and GitHub (*vide infra*).

96 Baseline performance of readily available docking scores

97 The performance of two readily available docking scores was assessed to establish a baseline for bioactivity
98 prediction. To this end, we assessed the predicted binding affinity by the Vina score, and used RFScoreVS²³
99 to rescore all poses generated by VinaGPU and DiffDock. The results are aggregated in Figure 2. Unsurprisingly,
100 neither of the scoring algorithms showed any productive correlation with the Papyrus pChEMBL values, either
101 when looking at the entire dataset (Figure 2A-C) or when aggregating the per-kinase correlation coefficient (R^2)
102 over the kinase groups (Figure 2D-F).

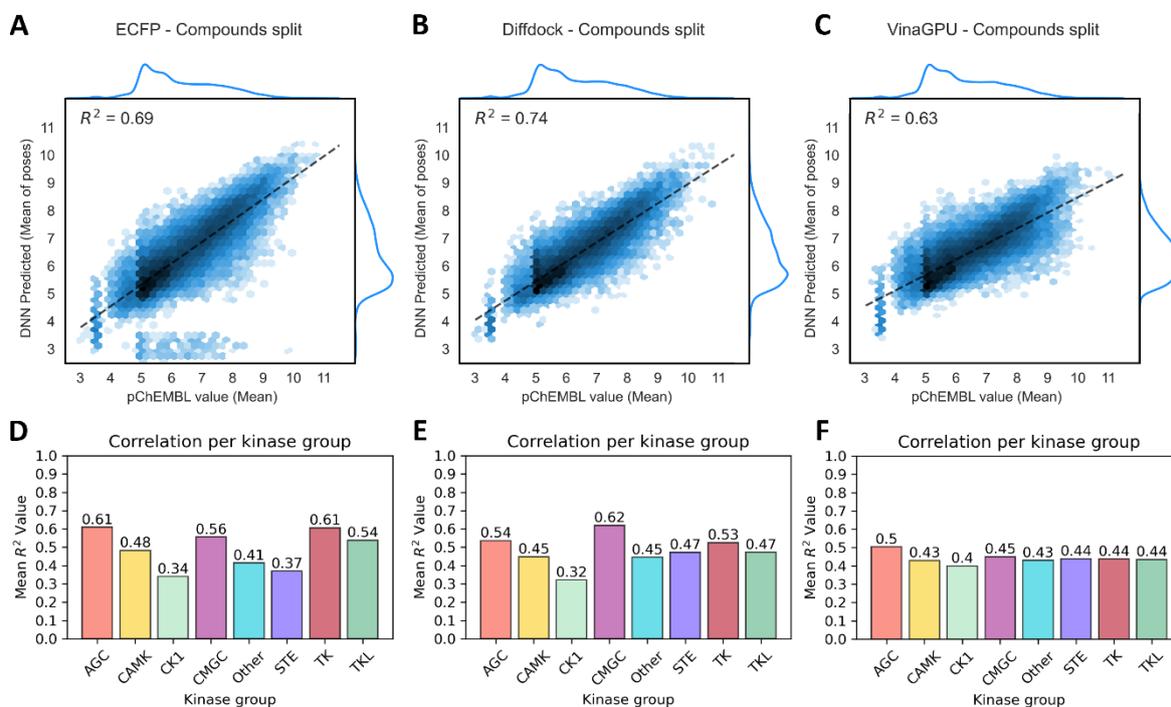


103 **Figure 2: Correlation of Vina and RFScoreVS scoring functions with Papyrus pChEMBL data** | Predicted affinity values vs.
104 literature values displayed as logarithmic hexbin plots, as based on the -Vina score for all VinaGPU poses (A), RFScoreVS for
105 all VinaGPU poses (B), RFScoreVS for all DiffDock poses (C), R^2 calculated per kinase and aggregated per kinase group for Vina
106 scores (D), RFScoreVS for VinaGPU poses (E) and RFScoreVS for DiffDock poses (F).
107

108 Kinome-wide activity predictions learned from docked poses

109 We then set out to train a more performant kinase specific scoring function on this unprecedentedly large structural
110 dataset. First, the database was used to generate protein-ligand extended connectivity (PLEC)²⁴ fingerprints for
111 the first three poses of every protein structure-inhibitor pair. All PLEC fingerprints were used as input for one
112 single 3-layer Deep Neural Network tasked with predicting the affinity value based on a given fingerprint. This
113 was done separately for the two docking algorithms, to compare their relative performance in this task. The
114 generated models, which function as kinase-specific scoring functions, were trained on either a random 80:20
115 split of protein-inhibitor activity pairs, an 80:20 compound-based split (completely unseen compounds) or an
116 80:20 split based on kinases (completely unseen kinases as test set). These latter splits are intended to assess
117 the generalisation capabilities of the models towards newly designed inhibitors or unseen kinase targets,
118 respectively. As a non-docking 2D comparison, in parallel we trained the same DNN on only the ECFP4
119 fingerprints of the inhibitors, to assess the added value of using docked poses as input. In this case we trained
120 one model per kinase for all kinases that had at least 100 unique inhibitors known (172 out of 226 kinases in the
121 dataset). The performance results of these models are shown in Figure 3 and Supplementary Figure 5 and
122 specified per kinase in Supplementary Tables 1-3 (Supplementary Materials).

123 Regardless of the underlying docking algorithm, the performance of the DNNs trained on the compound
124 splits ($R^2 = 0.63 - 0.74$) vastly outperformed both the original Vina scoring ($R^2 = 0.04$) as well as the rescoring
125 using RFScoreVS ($R^2 = 0.05 - 0.06$). For the DiffDock model, for 86 out of 214 kinases (40%) the R^2 of the
126 compound split was higher than 0.6, yielding predictions of sufficient quality to be genuinely informative in drug
127 discovery projects. This value is comparable to the ECFP models, where 84 out of the 172 were ≥ 0.6 . Of note,
128 the ECFP models were only trained for kinases with ≥ 100 compounds, which leads to fewer kinases covered
129 overall. The DiffDock model can extrapolate to some extent to the lower coverage kinases that are lacking in de
130 ECFP models, with an $R^2 \geq 0.6$ for 18% of these (8 out of 42). The VinaGPU model shows somewhat lower overall
131 performance, with 65 out of all 220 models having an $R^2 \geq 0.6$, and 5 of the low-coverage kinases. This
132 corresponds to the overall higher RMSD as observed in the docking procedure, pointing to the lower quality of
133 the underlying training data.



134 **Figure 3: Model performance** | Predicted affinity values vs. literature values for the compounds-split test set displayed as
135 logarithmic hexbin plots, as based on predictions for ECFP models (A), the DNN trained on DiffDock poses (B) and on the
136 VinaGPU poses (C). Panels D, E and F show the average performance per kinase group for ECFP, DiffDock and VinaGPU models,
137 respectively.
138

139 Comparing the DiffDock and VinaGPU-based models shows some intriguing results. There is only a low
140 correlation between the performances per kinase (Supplementary Figure 6). This can partially be attributed to
141 the smaller number of successful docking poses DiffDock generated but could also be due to the intrinsic
142 differences between the pose finding tools.

143 The different splits clearly showed that the random splits performed best overall, although only slightly
144 outcompeting the compound split. This is to be expected as for 78% of the compounds there is only 1 activity in
145 the dataset, meaning that the random and compound splits have highly similar difficulties in practice. However,
146 for unseen kinases the performance drops significantly (Supplementary Figure 5). This seems to indicate that the
147 model strongly relies on the kinase structure underlying the complexes and suggests that when appending new
148 kinases or KLIFS structures to the dataset, retraining of the model is warranted.

149 **KinaseDocker² release for direct local application**

150 Encouraged by the strong performance across the kinome we decided to wrap our workflow and models in a
151 user-friendly application that allows predictions to be generated by a medicinal chemist in real-world
152 applications. Because the model inherently generates docking poses on which the affinity prediction is based,
153 interpretation of the reliability of the output can be done on a per-compound and per-kinase basis. With this
154 interpretability endpoint in mind, we chose the open-source molecular viewer PyMOL as the basis for the
155 program. We wrote a plugin that allows the input of a (list of) SMILES strings, the selection of a (list of) kinases
156 and the choice of docking engine. The docking and consecutive bioactivity prediction by the neural network is
157 handled by a Docker image that requires minimal installation by the user. The output data is written to files as
158 well as presented in a table on screen. From this table generated complexes can be loaded and inspected in the
159 PyMOL session. Programmatic access is available if larger scale runs are desired. The whole codebase has been
160 designed to be modular, allowing the future implementation of different model architectures or structure
161 encodings. All code and Docker images are openly available, see section Code Availability.

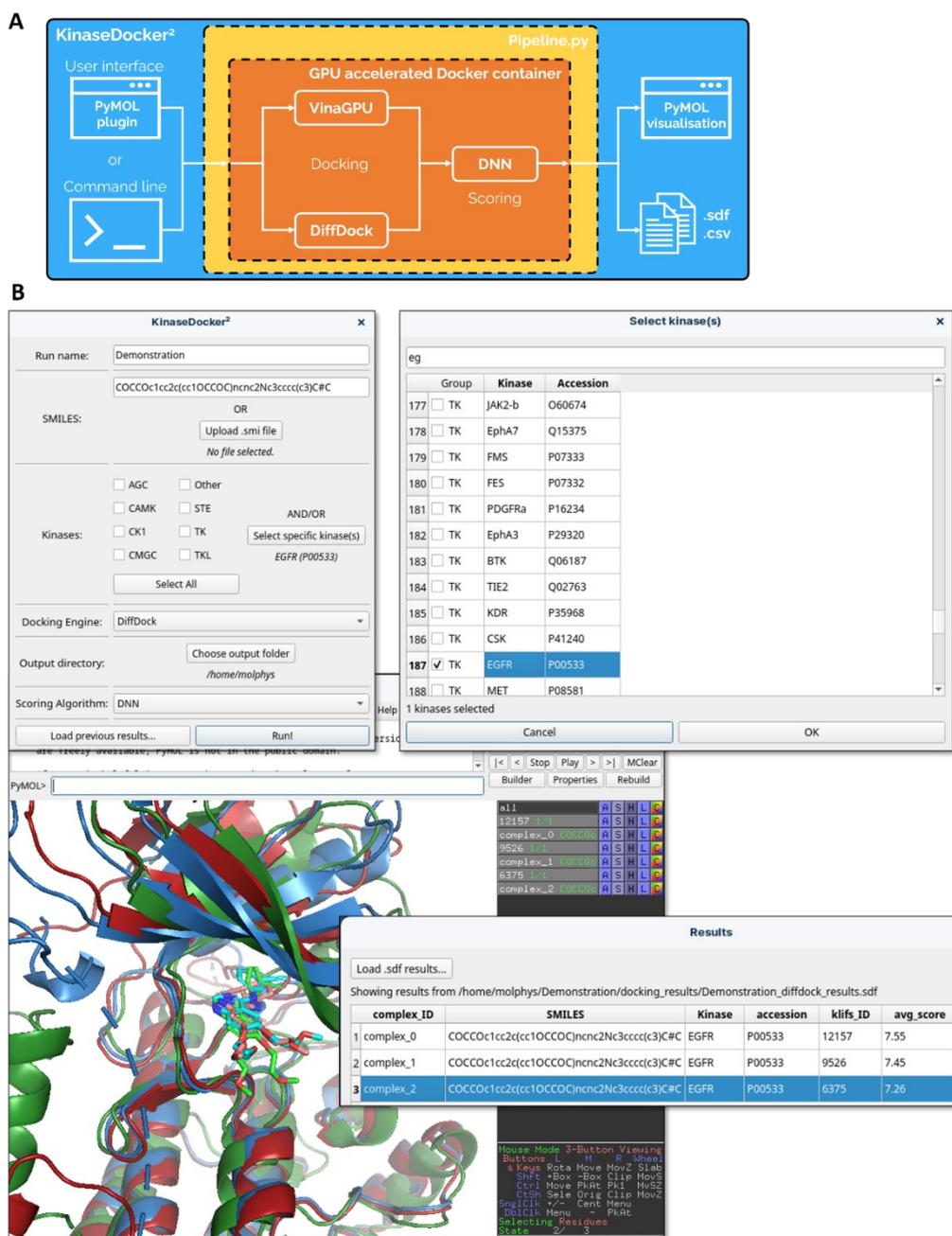


Figure 4: A user-friendly application: KinaseDock² | (A) Schematic overview of the software setup; (B) screenshots of the Graphical User Interface of KinaseDock².

162
163
164
165

166 Discussion

167 The homogeneity of the sources of biochemical data in the Papyrus dataset (and nearly every other publicly
168 available dataset) inherently means that there is considerable noise in the data. Realistically, R^2 values of around
169 0.8 are as high as can be achieved when taking experimental error into account.²⁵ This means that the DiffDock
170 model for 42 kinases ($\pm 20\%$) already reached this maximum. For these, no significant improvement on this
171 metric can be expected regardless of the methodological improvements or addition of further data. Adding more
172 (diverse) compounds would for these kinases merely expand the chemical space where the model is applicable.
173 For the kinases with lower performing predictions, the addition of more data and/or more structures could still
174 increase performance.

175 Training (machine learning-based) scoring functions on structural data has been a successful strategy for
176 years, enabled by datasets such as the PDBbind, as demonstrated by, for example, the RFScore series.^{12,23,26–28}
177 Utilising the accuracy of pose finding in docking algorithms to synthesize an orders of magnitude larger training
178 dataset has, to the best of our knowledge, not been attempted before. Here we clearly showed that the approach
179 in the basis works and outperforms current state-of-the-art in this kinase-specific use-case. There are many
180 possibilities for future improvement over the current machine learning implementation. The docking
181 performance of our VinaGPU workflow was not very high, with an average RMSD $> 5 \text{ \AA}$. More manual curation
182 of the dataset could reduce the amount of flawed docking poses, arguably positively impacting the quality of the
183 training data.

184 From a machine learning perspective, the current choice of encoding the poses (3D) using PLEC
185 fingerprints (1D) and utilising a basic DNN architecture is inherently lossy. Implementing geometric deep learning
186 models directly on the 3D data could positively impact performance if it can make better use of the available
187 information. Additionally, the attention mechanism of the Transformer architecture could be used to highlight
188 important regions in the complex for the generated prediction, yielding better interpretability and guidance for
189 compound optimisation.

190 There are more domain-focused improvements that could improve the performance too. The current
191 implementation uses every KLIFS structure available for a certain kinase when docking a compound, regardless
192 of inhibitor type (type I, II, III).^{29,30} Previous work has shown that ML models can differentiate Type I and II
193 inhibitors based on structure to a reasonable extent.³¹ By only considering the poses of a molecule in their
194 preferred activity state (DFG-in or -out), when available, the predictions should theoretically be improved.

195 To broaden the scope of kinases for which predictions can be made, structural data on the proteins is
196 currently the main bottleneck. Of the 636 kinases, 226 ($\pm 35\%$) have crystal structures that meet our criteria. Of
197 these, only about 26% (59) have both DFG-in and -out(-like) structures available. A strategy to enrich this dataset
198 could be through homology modelling. Considering the high sequence and structural similarity in the kinase
199 domains, for many if not most kinases a reliable homology model in both DFG-states should be feasible to obtain.
200 Adding these to the dataset would not only considerably extend the applicability of the model to the entire
201 kinome, it would also grow the size of the available biochemical training data with $>100,000$ datapoints for which
202 currently no high quality experimental structure is available.

203 Conclusion

204 Kinase inhibitors are an essential part of anti-cancer therapy. Developing new kinase inhibitors suitable for clinical
205 use requires these to be as specific as possible, targeting primarily the intended kinase. Due to the high homology
206 in kinase domains, this is not a trivial requirement. Computational tools to aid in the development of these
207 inhibitors by predicting inhibition across the kinome can be of great value. Current state-of-the-art struggles to
208 perform well across the protein family, in part due to the lack of suitable data. Here, we generate a large dataset
209 of predicted binding poses, each corresponding to an experimental binding affinity in the Papyrus dataset, where
210 a high-quality kinase domain structure of the target is available in the KLIFS database. We showed that this
211 dataset forms a strong basis on which to train machine learning models that can predict binding affinities of
212 compounds for a wide variety of targets. We trained a Deep Neural Network on a 1D protein-ligand interaction
213 fingerprint representation (PLEC) and showed that this vastly outperforms readily available (re-)scoring functions
214 like Vina score and RFScoreVS. Encouraged by these results, we developed a user-friendly interface to bring the
215 automated docking procedure and scoring function as a freely available tool called KinaseDocker² to the bench
216 chemist. Simultaneously, we ensured the modularity of the code, so that exchanging the protein-ligand complex

217 encoding or the predictive model for more advanced approaches is feasible. Finally, we set-up an interface for
218 the database of docking poses to expose the data encapsulated in this to the general (bio)chemist.

219 We expect that the scoring functions trained here are useful as is, but also that, together with the dataset
220 generated here, they form a starting point to further tackle the kinase selectivity question, enabling the reliable
221 prediction of affinities across the kinome to aid in bringing new and safe anti-cancer drugs to patients.

222 **Code & Data availability**

223 The 3D structure database generated as part of this work is available as an .sqlite database on Zenodo
224 (10.5281/zenodo.10894122), together with the KNIME workflow that provides a simple user interface to
225 search it. Code to reproduce the work described in this paper is available on GitHub
226 (<https://github.com/APAJanssen/KinaseDocker2-Paper-code>). The PyMOL plugin is available on its own GitHub
227 (<https://github.com/APAJanssen/KinaseDocker2>), which contains instructions on how to set up the Docker
228 environment. The Docker image is available on Docker Hub
229 (<https://hub.docker.com/repository/docker/apajanssen/kinasedocker2>).

230 **Acknowledgements**

231 The authors thank Roelof van der Kleij for helping with the usage of the university's computational resources.
232 Dr. Olivier Béquignon is acknowledged for fruitful discussions regarding the high-throughput docking. A.P.A.J.
233 expresses his sincere gratitude to prof.dr. Gerard J.P. van Westen and prof.dr. Mario van der Stelt for their
234 mentorship and critical feedback on the manuscript.

235 **Author Contributions**

236 APAJ conceived the project. JS, AB and APAJ performed the work described herein. JS and APAJ wrote the paper.

237 **Funding**

238 Research reported in this publication was supported by the Onco Institute and Onco Accelerator, a Dutch
239 National Growth Fund project under grant number NGFOP2201.

240

References

- 242 (1) Roskoski, R. Properties of FDA-Approved Small Molecule Protein Kinase Inhibitors: A 2024 Update.
243 *Pharmacol. Res.* **2024**, *200*, 107059. <https://doi.org/10.1016/j.phrs.2024.107059>.
- 244 (2) Fedorov, O.; Müller, S.; Knapp, S. The (Un)Targeted Cancer Kinome. *Nat. Chem. Biol.* **2010**, *6* (3), 166–
245 169. <https://doi.org/10.1038/nchembio.297>.
- 246 (3) van Esbroeck, A. C. M.; Janssen, A. P. A.; Cognetta, A. B.; Ogasawara, D.; Shpak, G.; van der Kroeg, M.;
247 Kantae, V.; Baggelaar, M. P.; de Vrij, F. M. S.; Deng, H.; Allarà, M.; Fezza, F.; Lin, Z.; van der Wel, T.;
248 Soethoudt, M.; Mock, E. D.; den Dulk, H.; Baak, I. L.; Florea, B. I.; Hendriks, G.; De Petrocellis, L.;
249 Overkleeft, H. S.; Hankemeier, T.; De Zeeuw, C. I.; Di Marzo, V.; Maccarrone, M.; Cravatt, B. F.; Kushner,
250 S. A.; van der Stelt, M. Activity-Based Protein Profiling Reveals off-Target Proteins of the FAAH Inhibitor
251 BIA 10-2474. *Science* **2017**, *356* (6342), 1084–1087. <https://doi.org/10.1126/science.aaf7497>.
- 252 (4) Manning, G.; Whyte, D. B.; Martinez, R.; Hunter, T.; Sudarsanam, S. The Protein Kinase Complement of
253 the Human Genome. *Science* **2002**, *298* (5600), 1912–1934. <https://doi.org/10.1126/science.1075762>.
- 254 (5) Zarrinkar, P. P.; Gunawardane, R. N.; Cramer, M. D.; Gardner, M. F.; Brigham, D.; Belli, B.; Karaman, M.
255 W.; Pratz, K. W.; Pallares, G.; Chao, Q.; Sprankle, K. G.; Patel, H. K.; Levis, M.; Armstrong, R. C.; James, J.;
256 Bhagwat, S. S. AC220 Is a Uniquely Potent and Selective Inhibitor of FLT3 for the Treatment of Acute
257 Myeloid Leukemia (AML). *Blood* **2009**, *114* (14), 2984–2992. <https://doi.org/10.1182/blood-2009-05-222034>.
- 259 (6) Sorgenfrei, F. A.; Fulle, S.; Merget, B. Kinome-Wide Profiling Prediction of Small Molecules.
260 *ChemMedChem* **2018**, *13* (6), 495–499. <https://doi.org/10.1002/cmdc.201700180>.
- 261 (7) Merget, B.; Turk, S.; Eid, S.; Rippmann, F.; Fulle, S. Profiling Prediction of Kinase Inhibitors: Toward the
262 Virtual Assay. *J. Med. Chem.* **2017**, *60* (1), 474–485. <https://doi.org/10.1021/acs.jmedchem.6b01611>.
- 263 (8) Cichonska, A.; Ravikumar, B.; Parri, E.; Timonen, S.; Pahikkala, T.; Airola, A.; Wennerberg, K.; Rousu, J.;
264 Aittokallio, T. Computational-Experimental Approach to Drug-Target Interaction Mapping: A Case Study
265 on Kinase Inhibitors. *PLOS Comput. Biol.* **2017**, *13* (8), e1005678.
266 <https://doi.org/10.1371/journal.pcbi.1005678>.
- 267 (9) Jumper, J.; Evans, R.; Pritzel, A.; Green, T.; Figurnov, M.; Ronneberger, O.; Tunyasuvunakool, K.; Bates, R.;
268 Židek, A.; Potapenko, A.; Bridgland, A.; Meyer, C.; Kohl, S. A. A.; Ballard, A. J.; Cowie, A.; Romera-Paredes,
269 B.; Nikolov, S.; Jain, R.; Adler, J.; Back, T.; Petersen, S.; Reiman, D.; Clancy, E.; Zielinski, M.; Steinegger, M.;
270 Pacholska, M.; Berghammer, T.; Bodenstein, S.; Silver, D.; Vinyals, O.; Senior, A. W.; Kavukcuoglu, K.; Kohli,
271 P.; Hassabis, D. Highly Accurate Protein Structure Prediction with AlphaFold. *Nature* **2021**, *596* (7873),
272 583–589. <https://doi.org/10.1038/s41586-021-03819-2>.
- 273 (10) Ain, Q. U.; Aleksandrova, A.; Roessler, F. D.; Ballester, P. J. Machine-Learning Scoring Functions to Improve
274 Structure-Based Binding Affinity Prediction and Virtual Screening. *WIREs Comput. Mol. Sci.* **2015**, *5* (6),
275 405–424. <https://doi.org/10.1002/wcms.1225>.
- 276 (11) Li, H.; Sze, K.; Lu, G.; Ballester, P. J. Machine-learning Scoring Functions for Structure-based Virtual
277 Screening. *WIREs Comput. Mol. Sci.* **2020**. <https://doi.org/10.1002/wcms.1478>.
- 278 (12) Ballester, P. J.; Mitchell, J. B. O. A Machine Learning Approach to Predicting Protein–Ligand Binding
279 Affinity with Applications to Molecular Docking. *Bioinformatics* **2010**, *26* (9), 1169–1175.
280 <https://doi.org/10.1093/bioinformatics/btq112>.
- 281 (13) Liu, Z.; Li, Y.; Han, L.; Li, J.; Liu, J.; Zhao, Z.; Nie, W.; Liu, Y.; Wang, R. PDB-Wide Collection of Binding Data:
282 Current Status of the PDBbind Database. *Bioinformatics* **2015**, *31* (3), 405–412.
283 <https://doi.org/10.1093/bioinformatics/btu626>.
- 284 (14) Leach, A. R.; Shoichet, B. K.; Peishoff, C. E. Prediction of Protein–Ligand Interactions. Docking and
285 Scoring: Successes and Gaps. *J. Med. Chem.* **2006**, *49* (20), 5851–5855.
286 <https://doi.org/10.1021/jm060999m>.
- 287 (15) Béquignon, O. J. M.; Bongers, B. J.; Jaspers, W.; IJzerman, A. P.; van der Water, B.; van Westen, G. J. P.
288 Papyrus: A Large-Scale Curated Dataset Aimed at Bioactivity Predictions. *J. Cheminformatics* **2023**, *15*
289 (1), 3. <https://doi.org/10.1186/s13321-022-00672-x>.
- 290 (16) Ding, J.; Tang, S.; Mei, Z.; Wang, L.; Huang, Q.; Hu, H.; Ling, M.; Wu, J. Vina-GPU 2.0: Further Accelerating
291 AutoDock Vina and Its Derivatives with Graphics Processing Units. *J. Chem. Inf. Model.* **2023**, *63* (7),
292 1982–1998. <https://doi.org/10.1021/acs.jcim.2c01504>.
- 293 (17) Corso, G.; Stärk, H.; Jing, B.; Barzilay, R.; Jaakkola, T. DiffDock: Diffusion Steps, Twists, and Turns for
294 Molecular Docking. arXiv October 4, 2022. <http://arxiv.org/abs/2210.01776> (accessed 2022-10-24).
- 295 (18) Kooistra, A. J.; Kanev, G. K.; Van Linden, O. P. J.; Leurs, R.; De Esch, I. J. P.; De Graaf, C. KLIFS: A Structural
296 Kinase-Ligand Interaction Database. *Nucleic Acids Res.* **2015**. <https://doi.org/10.1093/nar/gkv1082>.

- 297 (19) Kanev, G. K.; de Graaf, C.; Westerman, B. A.; de Esch, I. J. P.; Kooistra, A. J. KLIFS: An Overhaul after the
298 First 5 Years of Supporting Kinase Research. *Nucleic Acids Res.* **2021**, *49* (D1), D562–D569.
299 <https://doi.org/10.1093/nar/gkaa895>.
- 300 (20) Merkel, D. Docker: Lightweight Linux Containers for Consistent Development and Deployment. *Linux J.*
301 **2014**, *2014* (239), 2:2.
- 302 (21) Landrum, G. RDKit: Open-Source Cheminformatics; [Http://Www.Rdkit.Org](http://www.rdkit.org). <http://www.rdkit.org>.
- 303 (22) Smith, R. H. B.; Dar, A. C.; Schlessinger, A. PyVOL: A PyMOL Plugin for Visualization, Comparison, and
304 Volume Calculation of Drug-Binding Sites. *bioRxiv* October 24, 2019, p 816702.
305 <https://doi.org/10.1101/816702>.
- 306 (23) Wójcikowski, M.; Ballester, P. J.; Siedlecki, P. Performance of Machine-Learning Scoring Functions in
307 Structure-Based Virtual Screening. *Sci. Rep.* **2017**, *7*, 46710. <https://doi.org/10.1038/srep46710>.
- 308 (24) Wójcikowski, M.; Kukięka, M.; Stepniewska-Dziubinska, M. M.; Siedlecki, P. Development of a Protein-
309 Ligand Extended Connectivity (PLEC) Fingerprint and Its Application for Binding Affinity Predictions.
310 *Bioinforma. Oxf. Engl.* **2019**, *35* (8), 1334–1341. <https://doi.org/10.1093/bioinformatics/bty757>.
- 311 (25) Landrum, G. A.; Riniker, S. Combining IC50 or Ki Values from Different Sources Is a Source of Significant
312 Noise. *J. Chem. Inf. Model.* **2024**. <https://doi.org/10.1021/acs.jcim.4c00049>.
- 313 (26) Ballester, P. J.; Schreyer, A.; Blundell, T. L. Does a More Precise Chemical Description of Protein–Ligand
314 Complexes Lead to More Accurate Prediction of Binding Affinity? *J. Chem. Inf. Model.* **2014**, *54* (3), 944–
315 955. <https://doi.org/10.1021/ci500091r>.
- 316 (27) Li, H.; Leung, K.-S.; Wong, M.-H.; Ballester, P. J. Improving AutoDock Vina Using Random Forest: The
317 Growing Accuracy of Binding Affinity Prediction by the Effective Exploitation of Larger Data Sets. *Mol.*
318 *Inform.* **2015**, *34* (2–3), 115–126. <https://doi.org/10.1002/minf.201400132>.
- 319 (28) Li, H.; Leung, K.-S.; Wong, M.-H.; Ballester, P. J. Correcting the Impact of Docking Pose Generation Error
320 on Binding Affinity Prediction. *BMC Bioinformatics* **2016**, *17* (11), 308. [https://doi.org/10.1186/s12859-](https://doi.org/10.1186/s12859-016-1169-4)
321 [016-1169-4](https://doi.org/10.1186/s12859-016-1169-4).
- 322 (29) Dar, A. C.; Shokat, K. M. The Evolution of Protein Kinase Inhibitors from Antagonists to Agonists of
323 Cellular Signaling. *Annu. Rev. Biochem.* **2011**, *80* (1), 769–795. [https://doi.org/10.1146/annurev-](https://doi.org/10.1146/annurev-biochem-090308-173656)
324 [biochem-090308-173656](https://doi.org/10.1146/annurev-biochem-090308-173656).
- 325 (30) Gavrín, L. K.; Saiah, E. Approaches to Discover Non-ATP Site Kinase Inhibitors. *MedChemComm* **2012**, *4*
326 (1), 41–51. <https://doi.org/10.1039/C2MD20180A>.
- 327 (31) Abdelbaky, I.; Tayara, H.; Chong, K. T. Prediction of Kinase Inhibitors Binding Modes with Machine
328 Learning and Reduced Descriptor Sets. *Sci. Rep.* **2021**, *11* (1), 1–13.
329 <https://doi.org/10.1038/s41598-020-80758-4>.
- 330 (32) Consortium, T. U. UniProt: The Universal Protein Knowledgebase. *Nucleic Acids Res.* **2017**, *45* (D1), D158–
331 D169. <https://doi.org/10.1093/nar/gkw1099>.
- 332 (33) O’Boyle, N. M.; Banck, M.; James, C. A.; Morley, C.; Vandermeersch, T.; Hutchison, G. R. Open Babel: An
333 Open Chemical Toolbox. *J. Cheminformatics* **2011**, *3* (1), 33. <https://doi.org/10.1186/1758-2946-3-33>.
- 334 (34) Wójcikowski, M.; Zielenkiewicz, P.; Siedlecki, P. Open Drug Discovery Toolkit (ODDT): A New Open-Source
335 Player in the Drug Discovery Field. *J. Cheminformatics* **2015**, *7* (1), 26. [https://doi.org/10.1186/s13321-](https://doi.org/10.1186/s13321-015-0078-2)
336 [015-0078-2](https://doi.org/10.1186/s13321-015-0078-2).
- 337 (35) *CalcLigRMSD* - *GitHub*. GitHub. <https://github.com/rdkit/rdkit/tree/master/Contrib/CalcLigRMSD>
338 (accessed 2023-04-27).
- 339 (36) Trott, O.; Olson, A. J. AutoDock Vina: Improving the Speed and Accuracy of Docking with a New Scoring
340 Function, Efficient Optimization, and Multithreading. *J. Comput. Chem.* **2009**, *31* (2), 455–461.
341 <https://doi.org/10.1002/jcc.21334>.
- 342 (37) Andrius Bernatavicius. Highly Parallel Molecular Docking Pipeline Using Vina-GPU (Dockerized) +
343 AutoDock Vina CPU, 2024. <https://github.com/andriusbern/vinaGPU> (accessed 2024-02-27).
- 344 (38) *PeriodicTable.com*. <https://periodictable.com/> (accessed 2023-02-26).
- 345 (39) Pettersen, E. F.; Goddard, T. D.; Huang, C. C.; Meng, E. C.; Couch, G. S.; Croll, T. I.; Morris, J. H.; Ferrin, T.
346 E. UCSF ChimeraX: Structure Visualization for Researchers, Educators, and Developers. *Protein Sci. Publ.*
347 *Protein Soc.* **2021**, *30* (1), 70–82. <https://doi.org/10.1002/pro.3943>.
- 348

349 Methods

350 Biochemical data

351 Data was retrieved from Papyrus v5.5¹⁵, filtering on the Uniprot Protein Class ‘Kinase’ and data quality ‘High’. The
352 data was matched to the KLIFS¹⁹ dataset based on Uniprot³² accessions. Mutations were disregarded and
353 averages for unique compound – Uniprot pairs were used as activity value (pChEMBL). Included bioactivities
354 were filtered based on the drug-likeness of the measured compounds. Filters used were MW between 250 and
355 750 Da, rotatable bonds ≤ 15 , number of hydrogen bond donors ≤ 10 and number of hydrogen bond acceptors \leq
356 15, calculated using RDKit.²¹

358 Structural data

359 Kinase structures and annotations were retrieved from KLIFS in October 2022. The structures were filtered on
360 resolution (≤ 2.5 Å) and missing residues (≤ 5) after which the highest quality (KLIFS metric) structure was selected
361 based on DFG-in/out and α C-helix states as annotated in KLIFS, if available. The .mol2 files were downloaded and
362 converted to PDB files using OpenBabel³³. PDB structures thus generated were used as is for DiffDock or further
363 converted to .pdbqt format using the Open Drug Discovery Toolkit³⁴ for use with Autodock VinaGPU.

365 Docking benchmark set

366 Ligands from the KLIFS database were extracted and used as a benchmark dataset for the two docking algorithms
367 used. RMSD was determined using the CalcLigRMSD extension for RDKit.³⁵

369 Pocket definition

370 Pockets for Autodock Vina docking were automatically generated using PyVOL²² using default settings with
371 manual curation to encompass the entire ATP-binding pocket. The largest pocket detected in most cases
372 represented the ATP-binding site, to which a 5 Å padding was added for the docking box. DiffDock was executed
373 without restraints on binding site location (i.e., blind docking).

375 Ligand preparation

376 SMILES strings from the Papyrus dataset were transformed into 2D structures using default settings and
377 enantiomers and cis/trans isomers were enumerated using RDKit.²¹ These RDKit objects were converted to
378 .pdbqt format for VinaGPU docking using the Open Drug Discovery Toolkit.³⁴ The RDKit objects were written to
379 .csv files in canonical SMILES format with stereo information to use as DiffDock input.

381 Docking

382 Two docking procedures were employed: DiffDock¹⁷ and AutoDock VinaGPU^{16,36}, both installed through Docker²⁰.
383 Generated VinaGPU poses were converted to mol-format using RDKit.

384 AutoDock VinaGPU

385 A Docker image of AutoDock VinaGPU^{16,37} was used, running on commercial RTX4070 or RTX3070 GPUs. For each
386 protein, the corresponding KLIFS structures with predefined binding site boxes were iterated and all compounds
387 with known activities docked. The AutoDock VinaGPU implementation differs slightly from the well characterized
388 CPU version in its docking settings, where the *exhaustiveness* parameter is now replaced by *search_depth* and
389 *thread*. A small parameter optimisation was performed to benchmark the performance of VinaGPU on this
390 dataset, resulting in the final settings *search_depth* = 10, *threads* = 8192 which resulted in balanced performance
391 vs. run time (data not shown). Output .pdbqt formatted poses were converted to .mol format using OpenBabel
392 and aggregated in a tabular format for inclusion in the database.

393 DiffDock

394 The original DiffDock Github release of October 2021 was used. Compounds were provided in canonical SMILES
395 format with explicit stereochemistry. ESM embeddings were generated using the provided scripts and default
396 settings:

```
--repr_layers 33 --include_per_tok
```

397 For inference, the release inference.py script was used with minor changes relating to the output data structure.
398 The author recommended settings for high throughput inference were used:

```
--inference_steps 20 --samples_per_complex 5 --batch_size 10 --actual_steps 18  
--no_final_step_noise
```

399 Output .sdf formatted poses were expanded to .mol format and aggregated in a tabular format for inclusion in
400 the database.

401

402 Clash-score filtering

403 The filter criterium ($\text{clash} < 10$) was based on the Vina output, where after fitting a normal distribution on the
404 clash scores a 3σ upper limit was calculated to be 10.02, which was visually inspected to be sensible and used
405 for both docking algorithms. The clash-score was calculated per atom using the formula:

$$406 \max \left[0, 1 - \frac{d}{(r_1 + r_2)} \right]$$

407 where d is the Euclidian distance between the atoms, and r_1 and r_2 are the Van der Waals radii of the respective
408 atom types.³⁸ This per-atom contribution was calculated based on selections made using the PyMOL API. In brief,
409 KLIFS .pdb and docking pose .mol-files were loaded in PyMOL. A selection around 4 Å of the ligand was made,
410 and for all resulting atoms pairs the clashing contribution was determined. All contributions were summed to
411 yield the pose clash-score.

412

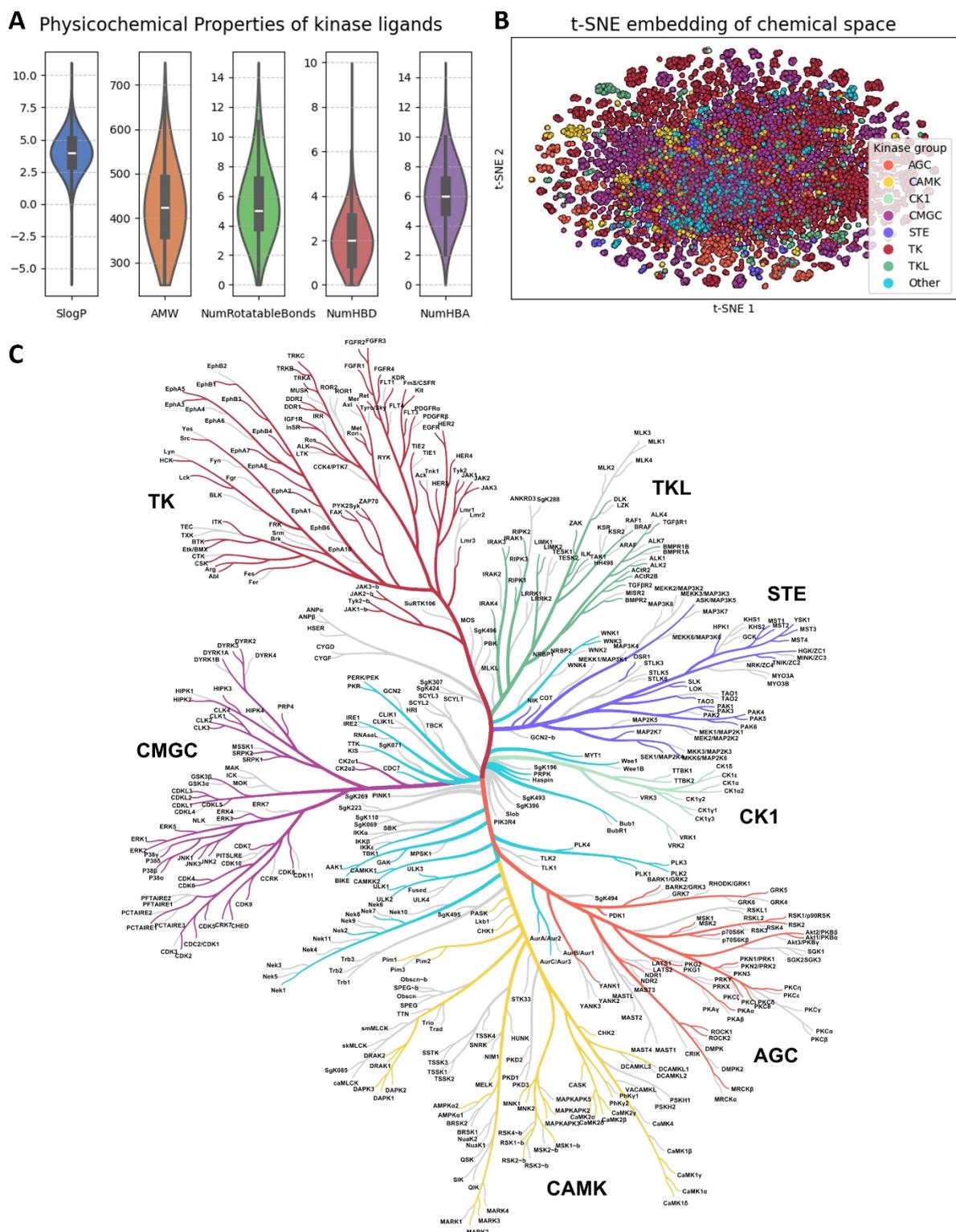
413 Machine learning

414 All machine learning algorithms were implemented in PyTorch 2.0. Splits were curated to ensure that the test set
415 pChEMBL distribution is similar to the train set distribution. All DNNs were 3-layer fully connected NNs with the
416 input layer either 2048 bits (ECFP) or 65536 bits (PLEC) to 4000, the hidden layer 4000 inputs to 1000 outputs
417 and the output layer using the 1000 inputs to 1 output value. All layers use ReLU activation functions and the
418 input and hidden layers use a dropout rate of 25% during training. The learning rate is fixed at 10^{-5} , batch size
419 128 with 100 epochs as fixed termination. After every epoch the performance on the test set is evaluated and
420 the best model is stored. Typically, 50-70 epochs are required to reach a plateau.

421

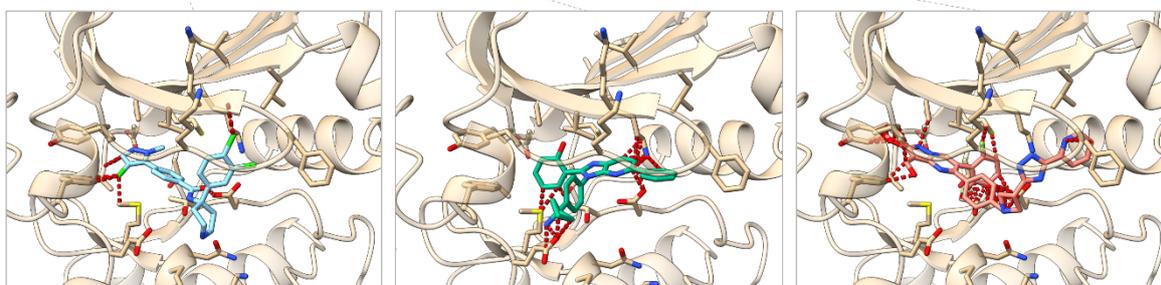
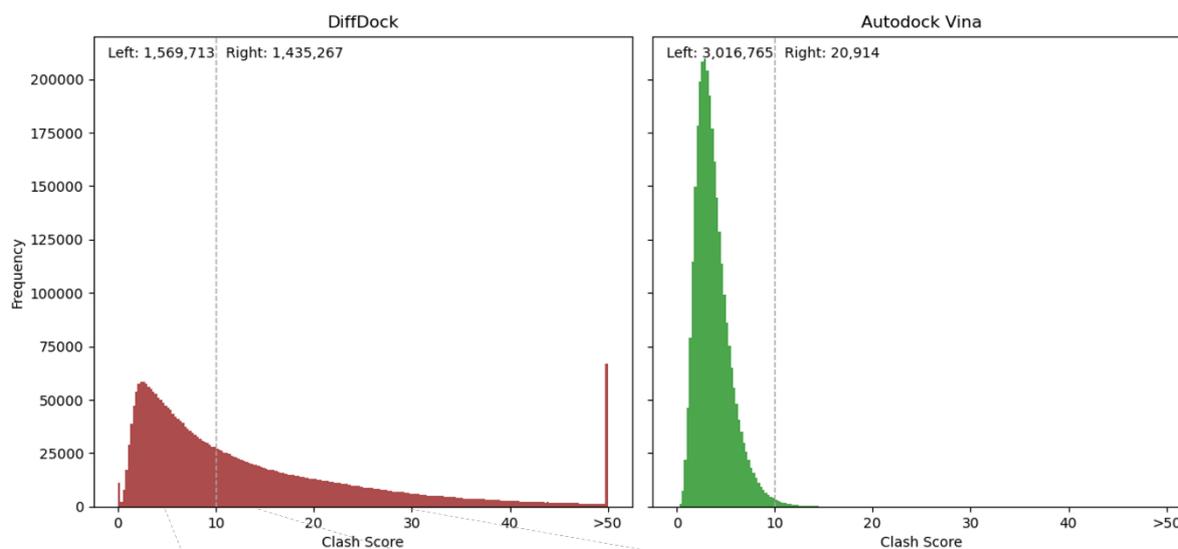
422 Prediction aggregation

423 For any given kinase-compound combination, the top 3 poses for all available KLIFS structures were scored by
424 the DNN. To get to a final activity prediction, we tested several aggregation strategies. Taking the mean value of
425 all options (aggregating the various KLIFS, all available stereoisomers, and the top 3 poses) yielded consistently
426 the highest R^2 (Supplementary Figure 8). As expected, using only the top 1 pose (according to Vina or DiffDock
427 ranking) performed slightly better than taking only the 2nd or 3rd ranked pose, showing that on average the built-
428 in scoring mechanism of both algorithms is able to prioritize the most relevant poses. However, averaging either
429 the top 2 or top 3 poses consistently improved the performance.



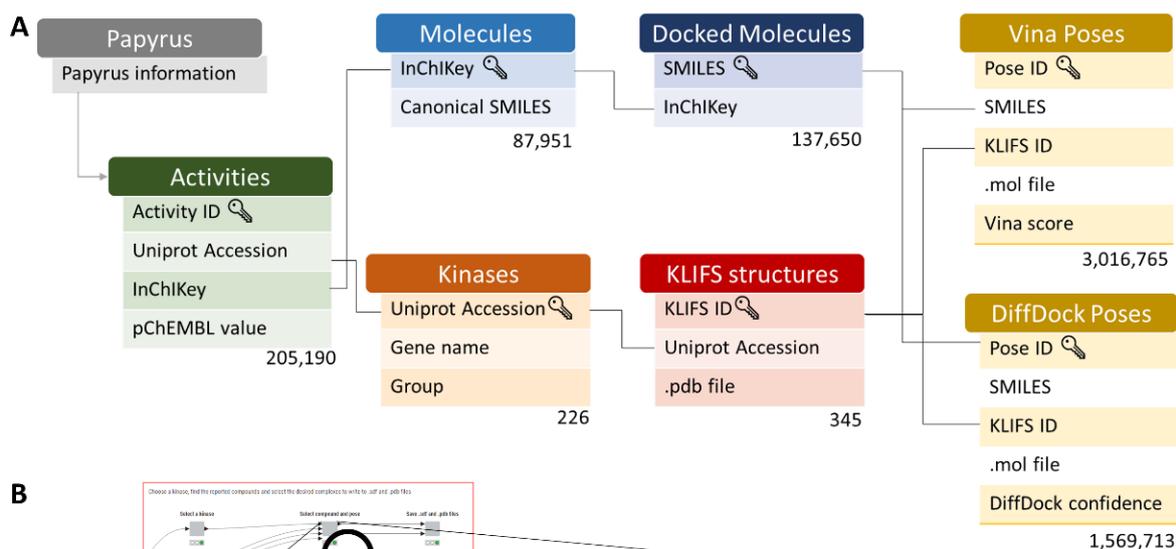
431
 432 **Supplementary Figure 1: Chemical and kinase diversity** | A) Violin plot of physicochemical properties of the kinase inhibitors;
 433 B) t-SNE embedding of the chemical space by ECFP4 fingerprints (2048 bits), coloured by majority kinase group target; C)
 434 View of included kinases coloured in phylogenetic tree.
 435

Comparison of Clash Scores



436
437
438
439
440

Supplementary Figure 2: Clashing score filtering | Histograms of the calculated clash scores for all DiffDock poses (left) and VinaGPU poses (right), illustrating the cut-off value of 10. Bottom inset shows 3 illustrations of poses with clash scores of 5, 15 and 30. Red dashed lines indicate atomic clashes. Insets were generated using UCSF ChimeraX³⁹.



B

Select the compound(s) of interest

InChIKey	pchembl_value_Mean	accession	Kinase	SMILES depiction
<input type="checkbox"/> IRUKIBDQFXZTM-UHFFFAOYSA-N	9.22	O15530	PDK1	
<input type="checkbox"/> YXZBQCGWYRRFP-UHFFFAOYSA-N	9.22	O15530	PDK1	
<input type="checkbox"/> RULJNHRYRQDQO-UHFFFAOYSA-N	9.05	O15530	PDK1	
<input checked="" type="checkbox"/> AFLQZVQTUSNSA-UHFFFAOYSA-N	9	O15530	PDK1	
<input type="checkbox"/> BEAUSMFBOFFBSA-UHFFFAOYSA-N	9	O15530	PDK1	

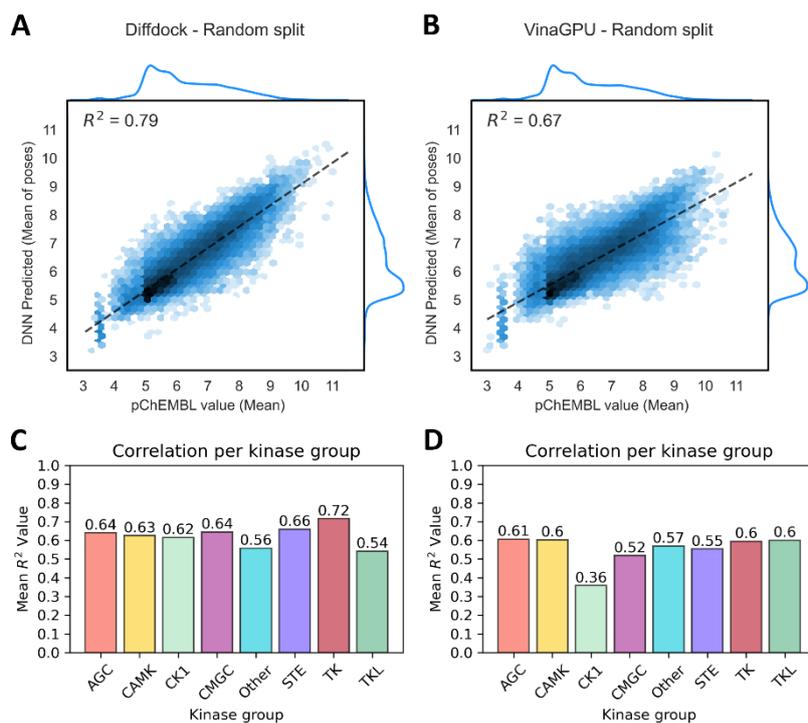
Select the pose(s) of interest

accession	Kinase	Kinasegroup	Fullname	klifs_id	inchikey	pchembl_value_Mean
<input checked="" type="checkbox"/> O15530	PDK1	AGC	3-phosphonoolido dependent protein kinase 1	3172	AFLQZVQTUSNSA-UHFFFAOYSA-N	9
<input checked="" type="checkbox"/> O15530	PDK1	AGC	5-phosphonoolido dependent protein kinase 1	3100	AFLQZVQTUSNSA-UHFFFAOYSA-N	9
<input checked="" type="checkbox"/> O15530	PDK1	AGC	5-phosphonoolido dependent protein kinase 1	3222	AFLQZVQTUSNSA-UHFFFAOYSA-N	9

Showing 1 to 3 of 3 entries (filtered from 6,897 total entries)

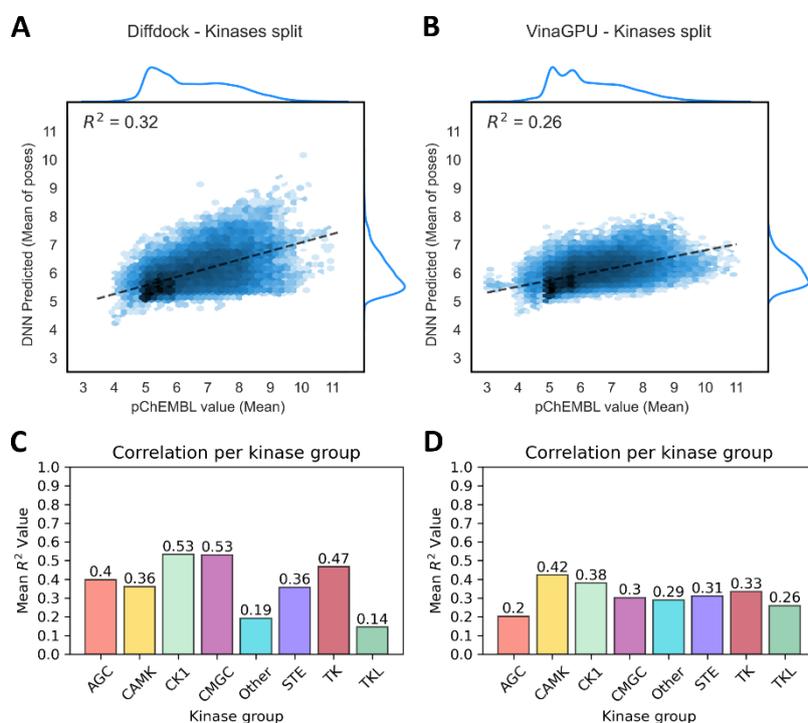
441
442
443
444
445

Supplementary Figure 3 Machine Learning-ready database of kinase-inhibitor complexes | A) Schematic and abbreviated database schema with statistics per table; B) Screenshots of the KNIME-based GUI that enables users to search and download data locally from the database.



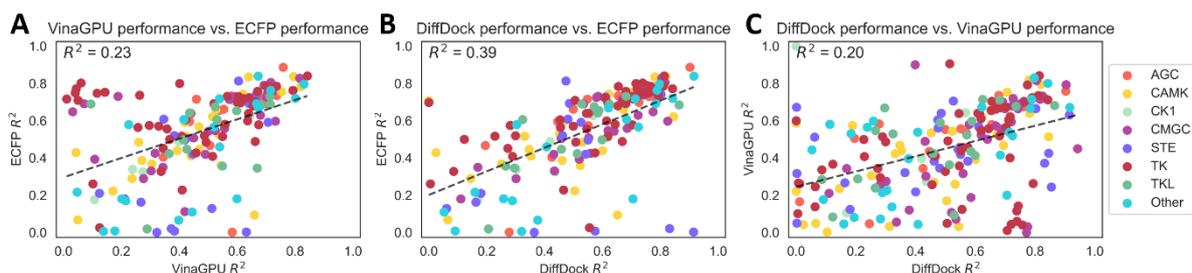
446
447
448
449
450
451

Supplementary Figure 4: Model performance on the random split | Predicted affinity values vs. literature values for the random-split test set displayed as logarithmic hexbin plots, as based on predictions of the DNN trained on DiffDock poses (A) and on the VinaGPU poses (B). Panels C and D show the average performance per kinase group for DiffDock and VinaGPU models, respectively.

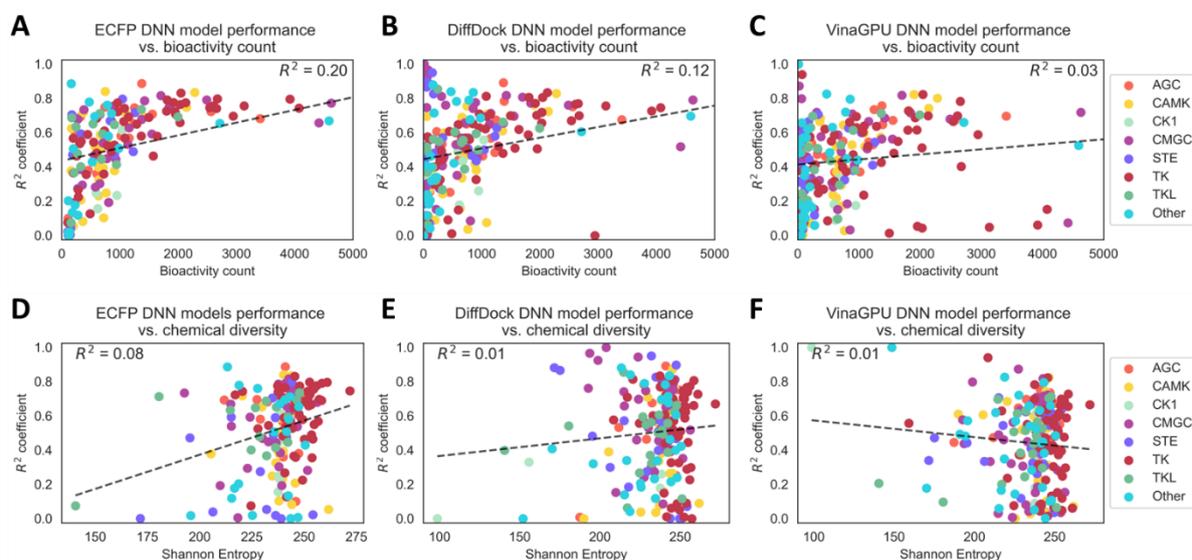


452
453
454
455
456

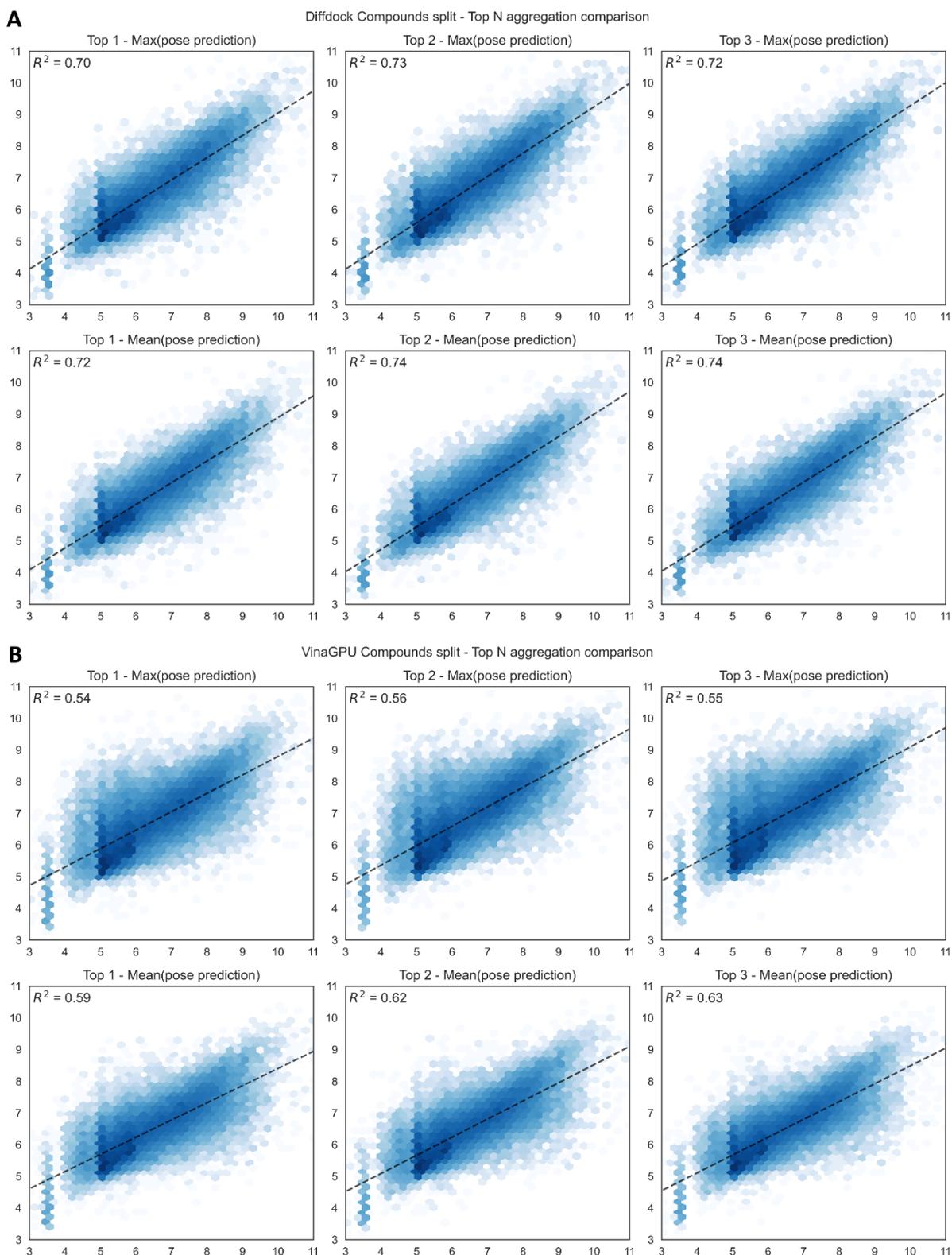
Supplementary Figure 5: Model performance on the kinases split | Predicted affinity values vs. literature values for the kinase-split test set displayed as logarithmic hexbin plots, as based on predictions of the DNN trained on DiffDock poses (A) and on the VinaGPU poses (B). Panels C and D show the average performance per kinase group for DiffDock and VinaGPU models, respectively.



Supplementary Figure 6: Performance correlation between models | Assessment of the correlation between the per kinase performance for VinaGPU and ECFP (A), DiffDock and ECFP (B) and DiffDock and VinaGPU (C) models. Kinases are coloured by kinase group.



Supplementary Figure 7: Correlation between performance and bioactivity count or chemical diversity | Assessment of the correlation between the per kinase performance for ECFP and bioactivity count (A), DiffDock and bioactivity count (B) and VinaGPU and bioactivity count (C), ECFP and chemical diversity (D), DiffDock and chemical diversity (E) and VinaGPU and chemical diversity (F). Chemical diversity is calculated as the Shannon entropy (higher = more diverse) of the ECFP fingerprint for all compounds included in the kinase's bioactivity data. Kinases are coloured by kinase group.



471
472
473
474
475

Supplementary Figure 8: Top N pose aggregation strategies | Predicted affinity values vs. literature values for the compounds-split test set displayed as logarithmic hexbin plots, as based on predictions of the DNN trained on DiffDock poses (A) and on the VinaGPU poses (B). Each sub-plot shows the aggregation of either the Top 1, 2 or 3 poses, using either the maximum or the mean of all predictions for all poses for that compound-kinase pair.