

# Digital Discovery

Accepted Manuscript

This article can be cited before page numbers have been issued, to do this please use: J. Poziemski, A. Yurkevych and P. Siedlecki, *Digital Discovery*, 2025, DOI: 10.1039/D5DD00452G.



This is an Accepted Manuscript, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this Accepted Manuscript with the edited and formatted Advance Article as soon as it is available.

You can find more information about Accepted Manuscripts in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this Accepted Manuscript or any consequences arising from the use of any information it contains.

# Assessment of molecular dynamics time series descriptors in protein-ligand affinity prediction.

Jakub Poziemski <sup>1</sup>, Artur Yurkevych <sup>2</sup>, Paweł Siedlecki <sup>1\*</sup>

<sup>1</sup> Institute of Biochemistry and Biophysics, Polish Academy of Sciences, Warsaw, Poland

<sup>2</sup> Institute of Chemistry, University of Silesia in Katowice, Katowice, Poland

\* Corresponding author: pawel@ibb.waw.pl

## Abstract

The advancements of computational methods in drug discovery, particularly through the use of machine learning (ML) and deep learning (DL), have significantly enhanced the precision of binding affinity predictions. However, accurate prediction of binding affinity remains a challenge due to the complex, non-linear character of molecular interactions. Generalizability continues to limit the current models, with performance discrepancies noted between training datasets and external test conditions. This study explores the integration of molecular dynamics (MD) simulations with ML to assess its predictive performance and limitations. In particular MD simulations offer a dynamic perspective by depicting the temporal interactions within protein-ligand complexes, potentially supplementing additional information for affinity and specificity estimates. By generating and analyzing over 800 unique protein-ligand MD simulations, we evaluate the utility of MD-derived descriptors based on time series in enhancing predictive accuracies. The findings suggest specific and generalizable features derived from MD data and propose approaches to augment the current *in silico* affinity prediction methods.



## Keywords:

View Article Online  
DOI: 10.1039/D5DD00452G

molecular dynamics, time series descriptors, protein ligand affinity, machine learning, simulation dataset

## Introduction

Computer-aided drug discovery (CADD) techniques have made an impact on the pharmaceutical industry by enhancing the efficiency of the drug development process, reducing time, cost, and labor [1]. Despite these advancements, accurate prediction of binding affinity continues to pose a considerable challenge, often bottlenecked by the inherent complexities of molecular interactions [2,3]. Progress in machine learning (ML) and deep learning (DL) has shown promise in overcoming some of these hurdles [4–6], revealing in large datasets intricate, non-linear properties and relationships in protein-ligand complexes. Current state-of-the-art methods achieve Pearson correlation coefficient (PCC) around 0.7–0.85 on the CASF2016 benchmark [7,8], which is a significant improvement compared to the classical scoring functions. Despite this achievement, challenges remain, particularly with the generalizability of these models. While definitely useful, traditional static computational approaches like molecular docking only provide a snapshot view of a molecular complex, without its temporal dynamics [9,10].

Molecular dynamics (MD) simulations introduce a vital temporal dimension to protein-ligand complex studies. Such simulations allow for more detailed observations of how drug molecules interact with biological targets over time [11–13]. Over the last years integration of MD simulations with ML and DL has been applied with varying success in different drug discovery tasks and specific campaigns. In the case of affinity prediction, Ash and Fourches in 2017 [14] analyzed 87 ERK2-docked ligand complexes by computing chemical descriptors derived from 20ns molecular dynamics (MD) trajectories. They showed that models trained on MD derived



descriptors were able to distinguish the most active ERK2 inhibitors from the moderate/weak actives and inactive. They claimed that the descriptors extracted from MD trajectories are highly informative and, having little correlation with classical 2D/3D descriptors, could augment chemical libraries screening tasks, candidate design and lead prioritization. A similar conclusion was presented by [15], who performed molecular docking of 43 compounds associated with Caspase-8, with consecutive 10-ns MD simulations of top scoring complex for each ligand. They investigated 770 2D and 115 3D descriptors together with 4 descriptors extracted from MD simulations: solvent accessible surface area (SASA), radius of gyration (Rg), potential energy and total energy, in the form of mean and standard deviation (8 descriptors in total). They reported that ML models trained on MD data had the most balanced accuracies and AUC values, compared to the 2D and 3D descriptor models, and that models using a combination of 3D and MD descriptors had the best performance.

A counter experiment was performed by [16] using the BCR-ABL tyrosine-kinase and 15ns MD simulations of Imatinib and a large series of its derivatives. In conclusion the authors stated that incorporating MD based matrices could not improve the binding affinity prediction ability of either deep NN nor random forest (RF) QSAR models. However, their models did show reduced prediction error, indicating the negative effect of simulation noise becoming stronger as the number of snapshots increased. Another approach to compare different ML models trained on descriptors obtained from MD trajectories was presented by [10]. Using three different targets and a maximum of 433 complexes predicted by docking per target, the results for MD augmented approaches were greatly dependent by target. The paper concludes the use of MD does not generally improve screening results and may only be justified in certain cases. Given that the models were trained with descriptors generated from every frame, this may have been a challenge for simple ML models due to the large number of frames and small number of examples. In addition, the low MM/PBSA and Glide scores suggest that the analyzed collections were rather difficult. In conclusion, current research indicates the complexity of leveraging molecular dynamics (MD) data, suggesting it is target specific and can depend on the noise to signal ratio, e.g. number of frames, the length of MD simulation,



etc. It is difficult however to draw definite conclusions as only a handful of targets have been tested so far.

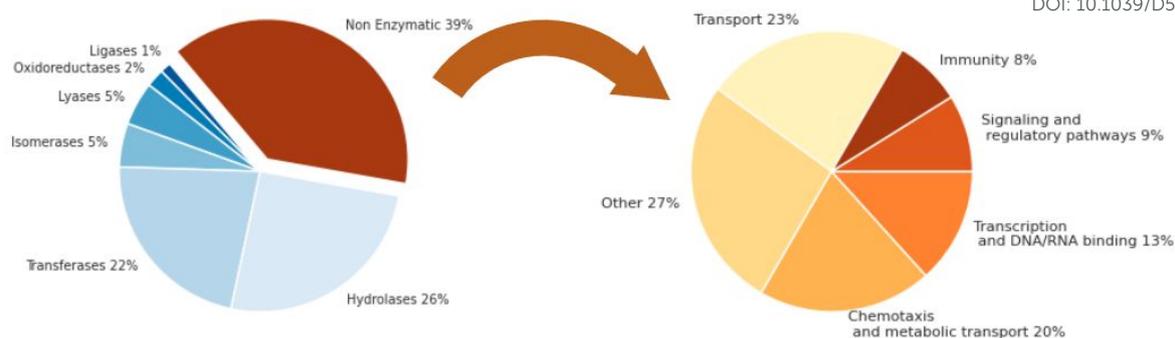
In this study we have generated the largest set of MD simulations to date, encompassing a broad array of protein-ligand complexes. We developed a comprehensive set of descriptors that utilize various aspects of MD-derived data, and implemented a rigorous feature selection mechanism to tailor the number of features to the analytical methods employed. By training ML models on MD-derived time series descriptors from over 800 unique protein-ligand complexes, we seek answers to the following questions: 1) what are the complex specific and simulation specific features influencing models' affinity prediction outcomes? 2) whether time series representations of MD are beneficial and how do they generalize? 3) Can the MD-derived descriptors augment and/or replace current crystallographic derived descriptors? Based on our findings we present general guidelines and assessments on how such an approach influences the *in silico* affinity prediction on a large scale.

## Methods

### Dataset compilation

The Molecular Dynamics Dataset (MDD) comprises 231 unique targets (Figure 1) from 862 protein-ligand complexes, sourced from the PDBBind collection v2020 [19]. Only complexes with well-defined active sites and ligands with unambiguous affinity values ( $-\log K_i$ ,  $K_d$  and  $IC_{50}$ ) were considered. More details on the filtering procedure, functional composition and similarity assessments of MDD targets and ligands are available in Supplementary materials.





**Figure 1: Functional characterization of MDD targets.** The “Non-enzymatic” part of MDD ( $\frac{1}{3}$  of the targets) are described on the right chart by 5 distinct biological processes (GO annotation). Nearly 40% of all MDD targets are non-enzymatic proteins with other functions.

## MD simulation procedure

MD preparation protocol, identical for all MDD complexes, is detailed in the supplementary materials. Briefly, MD simulations were executed using the GROMACS [17], utilizing a cubic box with periodic boundary conditions and a TIP3P water model. An initial minimization cycle, followed by temperature equilibration in the NVT ensemble and pressure equilibration in the NPT ensemble was carried out. Production simulations were conducted over a 200 ns timeframe, with a timestep of 100 ps.

## Representation

Given the relatively small size of the MDD dataset compared to e.g. PDBBind, we limit the number of descriptors to minimize overfitting, data sparsity and avoid other unfavorable phenomena caused by the curse of dimensionality [18].

The crystallographic (static) representation is composed of 63 descriptors: 24 pocket descriptors calculated with RDKit and MDAnalysis [19,20]), 11 interaction descriptors calculated with ProLIF [21]) and 28 ligand descriptors calculated with RDKit and SciPy [19,22] (for details see “List of descriptors” section in the Supplementary material).



In the case of MD simulation data, for each complex we calculate 51 descriptors: 24 pocket descriptors, 11 interaction descriptors, 9 ligand descriptors (19 ligand property descriptors are omitted as they do not change during simulation), and additional 7 motion descriptors not present in a static representations. These 51 descriptors are calculated for each frame of the MD simulation. For each descriptor sequence we extract its unique time series value with the use of a multistep procedure:

1. For each descriptor we calculate all 788 time series descriptors (`ts_descriptors`) supported by `tsfresh` (see Figure 2) .
2. For each `ts_descriptor`, its p-value is determined in terms of statistical significance against experimentally determined affinity (sourced from PDBBind), using a univariate test with FDR set to 0.001.
3. Next, a trimming procedure is used. For each of the 75 features types [23], the `ts_descriptor` with the lowest p-value is selected. Some `ts_descriptors` are parametric in nature; for such cases we use different thresholds to generate several versions of that `ts_descriptor` and select the one with the lowest p-value. At this stage, there could be a maximum of 75 `ts_descriptors` per a single descriptor.
4. To avoid caveats in training, trimming of correlations was applied. All descriptors and `ts_descriptors` were tested for correlations with each other. Correlated descriptors were dropped if PCC was  $\geq 0.8$ . Among a group of correlated `ts_descriptors`, the one with the highest cardinality was left.
5. In the final filtering step, for each descriptor, the `ts_descriptor` with the highest mutual information score between itself and the experimentally determined affinity value was selected. As a result, each descriptor is described with at most one `ts_descriptor` (Figure 2, bottom left).

In total, each crystallographic complex is described by 63 descriptors. For the MD data representation the same complex is described by a maximum of 51 `ts_descriptors` and 19 ligand property (static) descriptors, for which time series cannot be generated. Both

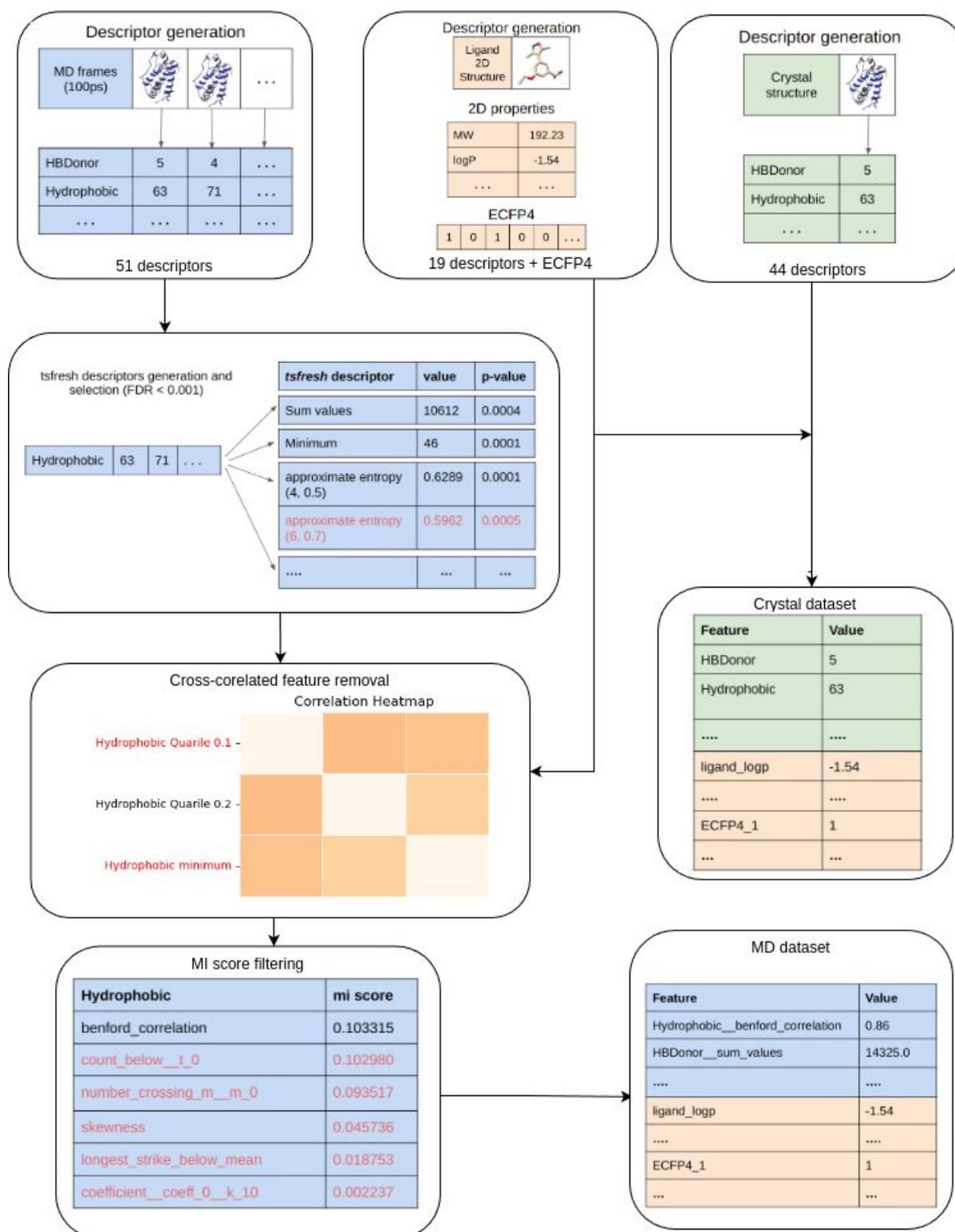


representations are augmented with the ECFP4 fingerprint (1024bits) to directly incorporate ligands molecular connectivity information.

[View Article Online](#)

DOI: 10.1039/D5DD000452G





**Figure 2. Flowchart of descriptors and ts\_descriptors selection procedure.** Color code: light green - crystallographic (static) descriptors, light blue - multiframe MD descriptors and ts\_descriptors, light orange - ligand descriptors. See section “Representation” for more details.



## Data splits, training and testing

Two ways of splitting the target data were used: random split and target split at a ratio of 4:1 (80% training collection, 20% test collection). In the random split, complexes were randomly allocated to the training and test sets. In target split, UniprotIDs split the complexes therefore there are no identical training and testing examples. The target split can therefore be used to approximate the generalization potential.

Ligands were split using DeepChem [24]. Scaffold Splitter divides molecules into groups based on their Bemis-Murcko scaffolds; the smallest scaffold groups form the test set. Although not without its caveats, this type of ligand split is more challenging than random splits as it tests more thoroughly the generalizability to new or less abundant areas of chemical space.

Three different ML models were trained using the MDD dataset; Random Forest and SVM with [25] and XGBoost with [26]. Model parameters were selected using 5-fold cross-validation (CV). The models were fitted to the training sets and evaluated on the test sets. Throughout this work boxplots represent results obtained on the test set, with mean (triangle) and median (horizontal line inside the boxplot) values. More details on the Random Forest models and on the Descriptor model (XGB) are present in the Supplementary material section

## Results

### Baseline performance

To test the hypothesis that time series descriptors improve affinity prediction, we first assess the difficulty of predicting the affinity of MDD complexes using published models with publicly available code and training procedures [27]. Selected models were trained on the PDBBind dataset (v2020) with 862 MDD complexes excluded (Table 1). We compare these results against those obtained for CASF2016 dataset presented in the literature (Table 2).



Both Pearson Correlation Coefficient (PCC) and Root Mean Square Error (RMSE) values (Table 1) render the MDD as a more difficult dataset compared to CASF2016 (Table 2). All tested models show a rather consistent drop in performance, independent of their complexity. The observed decrease in performance may be multifaceted: CASF2016 is a relatively small dataset compared to MDD (285 vs 862 complexes) therefore it may be easier to optimize or overtrain the models. Also, excluding MDD complexes from the training data may have influenced the availability of information necessary for higher affinity prediction performance. Our descriptor model (see Representation section in Materials and Methods) shows the same consistent performance drop as seen with other models. Interestingly, our models' affinity prediction performance is on par with some of the best, highly sophisticated methods. This result highlights that a carefully selected set of descriptors and relatively simple ML model can show a level of performance comparable to specialized neural networks.

Model name	Description	PCC	RMSE	Training size	Code and feature reference
OnionNet2	CNN trained on contact descriptors	0.75	1.26	13 546	[6]
PLEC-NN	Extended Connectivity FP & Neural Network	0.74	1.54	11 203	[28]
Descriptors model (XGB)	XGBoost trained on 63 descriptors + ECFP4	0.72	1.29	12 349	this work
NN-Score	Feed-forward neural network	0.67	1.40	4647	[29]
RF-Score v3	Random Forest with spatial distance count	0.65	1.45	4647	[4]
Vina	Hybrid empirical scoring function	0.49	-	-	[30]



**Table 1. Performance of selected affinity prediction methods on MDD.** All models were trained on PDBBind v2020 complexes with MDD complexes excluded. Both PCC (Pearson Correlation Coefficient) and RMSE (Room Mean Square Error) values were calculated for MDD complexes. The Training size column shows the number of unique protein-ligand complexes used for training the ML/DL models; according to their original implementation either with general set or refined set; Vina uses a classical hand-designed scoring function with optimized parameters.

Model name	Description	PCC CASF   CoreSet	RMSE CASF   CoreSet	Training size	Original reference
OnionNet2	CNN trained on contact descriptors	0.86   0.82	1.16   1.36	-	[6]
TopBP	Topological Descriptors with GBT	0.86   0.81	1.19   1.95	3 767	[31]
SS-GNN	Graph Neural Network	0.85   0.82	1.18   1.35	15 394	[32]
Descriptors model (XGB)	XGBoost	0.85   0.81	1.20   1.37	12 866	this work
EBA-AY	Ensemble Attention Based	0.86   0.79	1.20   1.44	10324	[33]
OPRC-GBT	Ollivier persistent Ricci curvature	0.84   0.79	1.25   2.01	3 772	[34]
DCML	Dowker complex based	0.84   0.78	1.25   1.43	3 772	[35]



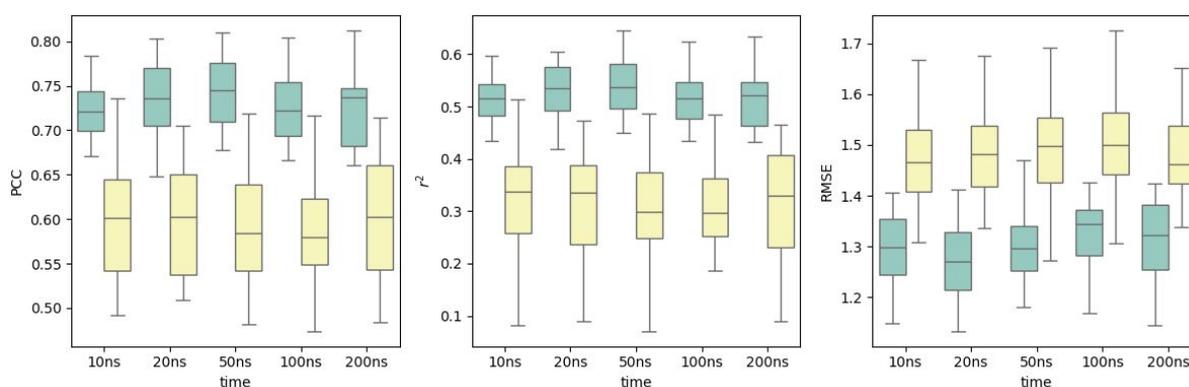
Model name	Description	PCC CASF   CoreSet	RMSE CASF   CoreSet	Training size	Original reference
OnionNet2	CNN trained on contact descriptors	0.86   0.82	1.16   1.36	-	[6]
TopBP	Topological Descriptors with GBT	0.86   0.81	1.19   1.95	3 767	[31]
SS-GNN	Graph Neural Network	0.85   0.82	1.18   1.35	15 394	[32]
Descriptors model (XGB)	XGBoost	0.85   0.81	1.20   1.37	12 866	this work
EBA-AY	Ensemble Attention Based	0.86   0.79	1.20   1.44	10324	[33]
	machine learning				
CAPLA	Sequence based cross-attention with 1D CNN	0.84   0.77	1.2   1.36	11 906	[36]
PLANET	Graph Neural Network	0.82   N/A	1.24   N/A	15 616	[37]
K <sub>DEEP</sub>	Convolution Neural Network	0.82   N/A	1.27   N/A	3 767	[38]
PLEC-NN	Extended Connectivity FP & Neural Network	0.82   0.77	1.25   1.43	12 906	[28]
OnionNet	Convolutional Neural Network	0.82   0.78	1.27   1.50	11 906	[39]
RF-Score v3	Random Forest	0.80   0.74	1.39   1.51	3 767	[4]
Pafnucy	Convolution Neural Network	0.78   0.70	1.42   1.62	11 906	[5]
Vina	Hybrid empirical scoring function	0.60   0.56	1.75   1.86	-	[40]

**Table 2: Published affinity prediction performance obtained with models of increasing complexity, tested with CASF\_2016 and CoreSet\_2013 benchmarks.** PCC (Pearson Correlation Coefficient), RMSE (Room Mean Square Error), training size (number of unique protein-ligand complexes used for training ML/DL models). The use of simple models on crystallographic data can yield comparable results to the use of complex neural networks based models.



## Simulation length

To determine the optimal length of MD simulations for information extraction, we tested 5 timescales from 10ns up to 200ns. summarizes the results of testing different MD lengths with the RF models. The results (Figure 3) show varying correlations, both with random and target splits. For random splits 50ns runs have slightly higher correlation values (PCC and  $r^2$ ) but show a lower RMSE at 20ns. For the target split, 20ns show best performance estimates for all three measures (PCC,  $r^2$  and RMSE values). Taken together the 20ns simulation length should provide good performance balance for both random and target splits.



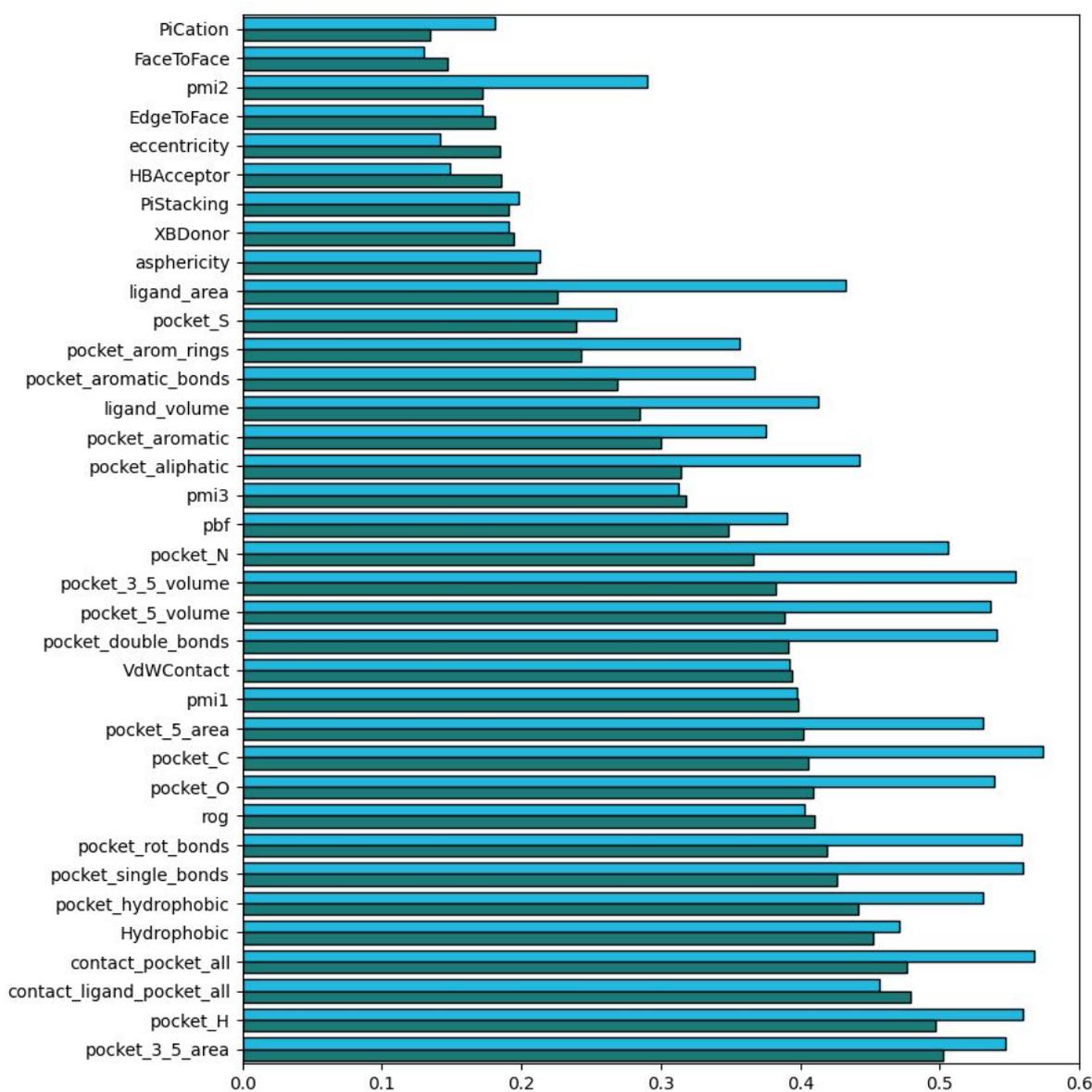
**Figure 3: Dependence of MD simulation length on model performance.** RF models trained with trajectories of different lengths (from 10ns to 200ns) tested on two types of data splits: random (green) and target (yellow). 20ns trajectories show good overall performance for both random and target splits.

We assessed what is the difference in affinity correlations between individual time series descriptors (`ts_descriptors`) derived from 20ns compared to 200ns MD simulations. We filtered all cross-correlated `ts_descriptors` from the two simulation lengths, and compared the shared 36 `ts_descriptors`. The results show a significant gain in around 40% of tested `ts_descriptors` (14 out of 36, with  $\Delta > 0.1$ ) in favor of 20ns simulation length. Comparable correlations are registered for 22 `ts_descriptors` ( $\Delta < 0.1$ ). Interestingly there was no descriptor that had more than 0.1 correlation difference in favor of the 200ns simulation.



In the assessed timescale, the obtained results indicate that short MD simulations may capture useful steric conformation changes and that longer MD simulations may contain more random movements or noise, which may lower the individual correlation of some of the ts\_descriptors. Similar results have been presented in works concerning single targets [10,14,15]. Taken together our results show that time series descriptors from longer MD do not provide advantage over short MD runs, possibly due to higher noise content and increased noise fitting during model training.

View Article Online  
DOI: 10.1039/D5DD000452G



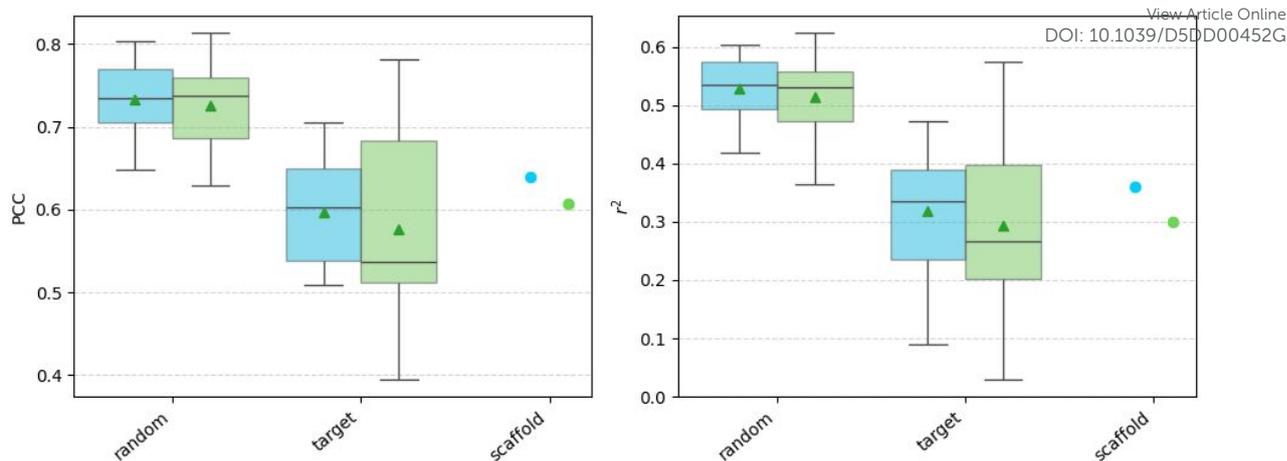
**Figure 4:** Correlations between affinity prediction and time series descriptors describing the MDD protein-ligand complexes. ts\_descriptors derived from 20ns (light blue) and 200ns (dark blue) molecular dynamics simulations.

We also tested if the frequency of saving the trajectory frames had an influence on model performance. We observed no significant improvement of results when training on an extensive number of frames (see Supplementary figure S3). From a practical point, in similar setups we recommend using a less frequent recording (larger time interval), which can significantly reduce storage requirements without decreasing prediction quality.

## MD representation

To estimate the impact of time series descriptors on affinity prediction we compared the overall performance of two types of models; trained with crystallographic descriptors, and with MD derived ts\_descriptors, with respect to different target and ligand data splits. The results of these experiments are summarized on Figure 5. All analyses refer to the RF models trained on 20ns molecular dynamics runs. For randomly split data both models achieve comparable results with respect to PCC (0.73 vs 0.73),  $r^2$ (0.53 vs 0.51) and SD (0.06 vs 0.06). However, in the case of a more challenging target split, an advantage of the model trained on time series descriptors can be seen (PCC mean: 0.6 vs 0.58, median 0.6 vs 0.53), along with a smaller SD (0.07 vs 0.10). This slight advantage can also be seen with respect to the Bemis-Murcko scaffold split. Both target and scaffold splits are more challenging tests as they try to minimize data leakage events. For scaffold split, the MD-based model also achieved better results (PCC 0.63 VS 0.60,  $r^2$  0.35 vs 0.29). Taken together, the performance gain and the lower SD of models trained on time series descriptors would suggest better generalizability potential, also with respect to uncharted chemical space.





**Figure 5: Affinity prediction performance of models trained on crystallographic-only data and models augmented with MD data.** RF models trained on static descriptors (crystallographic data: green), and time series descriptors (MD simulations: blue). Triangle: mean; horizontal line inside box: median. Scaffold split for a given ligands dataset is deterministic in nature, therefore only a single measurement point is visible.

Next we compared each crystallographic descriptor with its time series counterpart to assess their correlations with affinity. Figure 6 shows the results obtained for the 6 different groups of descriptors used. In the group of interaction descriptors, all time series descriptors have increased correlations compared to their crystallographic counterparts (with the exception of the Anionic term). For both types of descriptors, the hydrophobic and vdW interactions show the highest affinity correlation, importantly the time series correlation at least doubles compared to crystallographic descriptors. In the case of a single descriptor a correlation around 0.5 PCC would be considered fairly high. The rest of the interaction descriptor time series are higher than crystallographic, nevertheless overall their correlation values are low for both types of data.

With respect to pocket descriptors, both property and geometric, the results show a substantial number of them with PCC close or above 0.5. Overall, pocket property descriptors show the highest difference between static and dynamic treatment. Out of 8 best correlating time series descriptors only one static descriptor (pocket\_hydrophobic) has a comparable correlation. The



results show that simple pocket composition features, such as atom or bond types count, which change with respect to ligand movement, correlate with high PCC values with small molecule affinity.

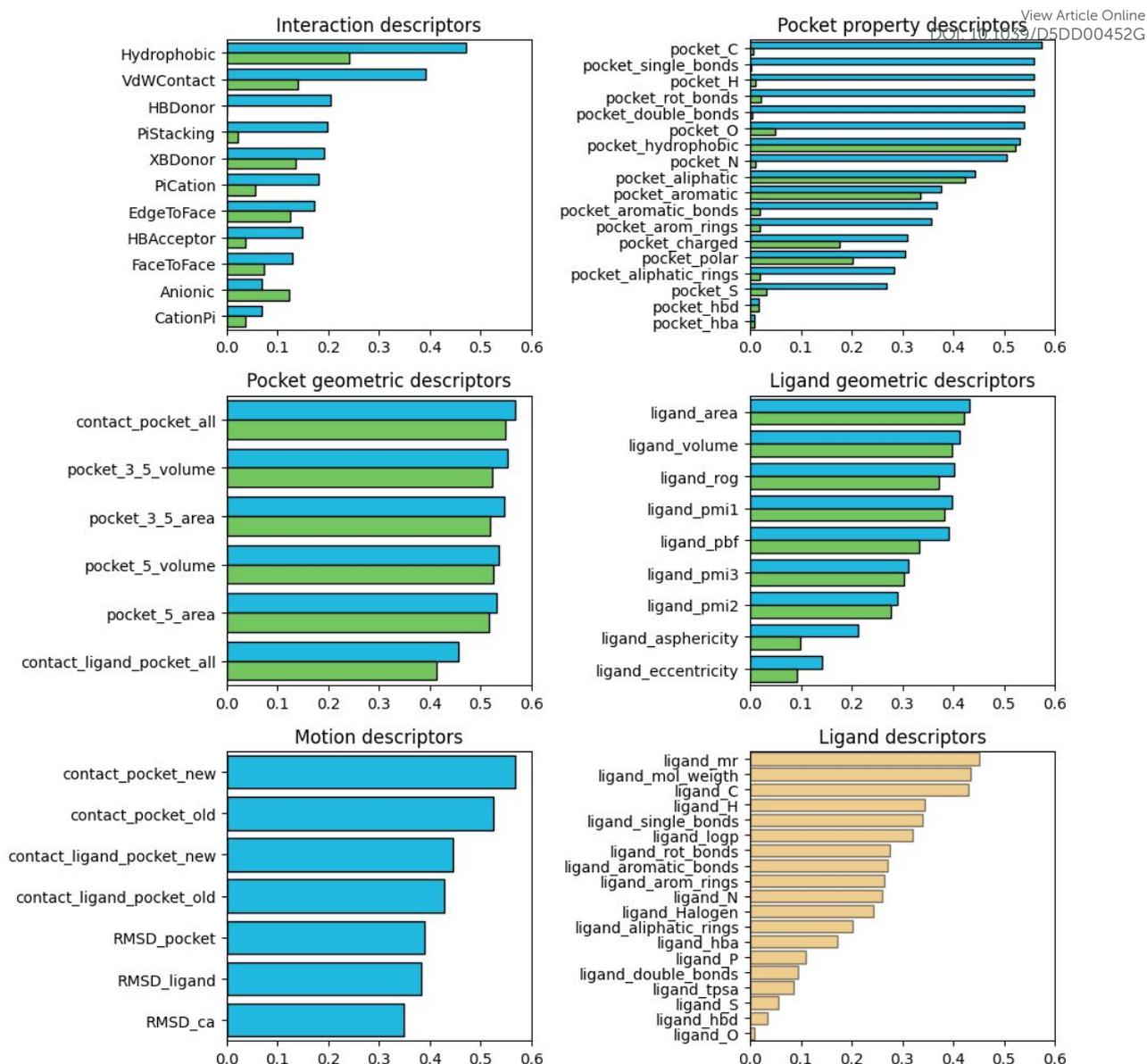
In the case of geometric descriptors correlation values are comparable between time series and static treatment. Similar results are obtained for ligand geometric descriptors. Here molecular eccentricity and molecular asphericity descriptor correlations gain the most from a time series representation, however their PCC values are rather low.

Surprisingly, motion descriptors which we thought would substantially contribute to the affinity prediction performance show only minor gains compared to static, crystallographic descriptors. Only two of them (pocket\_contact\_new and pocket\_contacts\_old) have PCC around 0.5 or above. Results suggest that pocket internal contacts (both preserved and new) are more correlated to affinity than e.g. pocket-ligand interactions and changes in RMSD of both molecular entities.

Taken together, although individual performance of the ts\_descriptors is favorable, it does not seem to add up to the final model performance (Figure 5). This effect is probably due to non-linear cross-correlations present between ts\_descriptors. This effect is further enhanced by the noise in the data. Common correlation methods such as Pearson's, Spearman's or Kendall are able to find linear and monotonic relationships between variables, assuming little noise. Other methods measuring potential non-linear correlations test only the independence of variables without quantifying the strength of the relationship, and are also very sensitive to noise and outliers. Therefore such nonlinearity is currently difficult to filter beforehand [41].

View Article Online  
DOI: 10.1039/D5DD000452G





**Figure 6: Pearson correlation coefficient of six types of descriptors expressed by the absolute value of the descriptor in relation to affinity.** Green represents results obtained by descriptors calculated from crystallographic structures, blue with their corresponding  $ts\_descriptors$  as calculated from the MD simulations (20ns), used in the model.

## Screening and pose selection.

We tested the impact of models trained with molecular dynamics (MD) data and with static crystallographic structures on a screening campaign using one of the targets derived from the DUD-E dataset [42]. The glucocorticoid receptor (GCR, PDB ID: 3BQD) was selected along



with a randomly drawn subset of 100 ligands; comprising 10 active molecules and 90 decoys

[View Article Online](#)  
DOI: 10.1039/D5DD000452G

This composition mimics a virtual screening scenario with a high imbalance between actives and inactives. Each GCR-ligand complex was docked using Vina with default parameters, except exhaustiveness was set to 32. Next, the best scoring conformation was subject to a 20ns MD simulation using the same parametrization procedure as for all other complexes (see Materials and Methods section). Table 3 summarizes the results of GCR screening.

	Static model	MD-derived model
<b>Screening task (GCR target)</b>		
EF10%	1.92	2.82
ROC AUC	0.615	0.675
<b>Pose selection (7 targets structures)</b>		
Top-1	2/7	4/7
Top-3	3/7	5/7
Mean std	0.22	0.28

**Table 3: Screening results for the static and MD-derived model (GCR target) and pose selection.**

The MD-derived model showed better performance in early enrichment (EF10%) compared to the static model (2.82 vs 1.92). This result was accompanied by a higher ROC AUC score (0.675 vs. 0.615). Overall the MD-derived model improved the ability to discriminate between active and inactive compounds, indicating a potential benefit of using the dynamic information in the training process.

To evaluate the ability of the two models to select the correct ligand binding pose, we conducted an experiment using a subset of targets derived from the CASF-2016 dataset. We selected 3 targets: FAX (1lpg, 1z6e, 2xbv), CDK2 (1pxn, 4eor) and CAH (2weg, 3dd0). For each of the 7 target-ligand complexes a docking run was performed following the same procedure as described for the screening experiment (i.e. Vina, default parameters,



exhaustiveness=32). To minimize bias, ligand structures used for docking were generated from SMILES strings using OpenBabels lowest energy parameter. Each docking run resulted in nine putative conformations for every protein target. The conformation derived directly from the crystallographic structure was included, resulting in 10 poses per target. Finally, each of the obtained poses were subject to a 20ns MD simulation using the same parametrization procedure, to obtain MD-derived descriptors (see Materials and Methods section). Results of the ability of the two types of models to correctly identify the native crystallographic structure (i.e. assign the highest affinity value) among the docking results are shown in Table X. The model based on MD trajectories performs better than the model based on crystallographic structures only. However, it is also noticeable that the standard deviation is higher for the MD model (0.28 vs 0.22), suggesting a greater sensitivity of this approach to the choice of initial conformation. This may reflect the complex nature of the molecular dynamics process and indicates that accounting for conformational variation affects the stability of the prediction. Ultimately, these results suggest that models based on MD simulation data may be an effective alternative to approaches based only on static white structures.

## Ablation studies

To understand the contributions of static and time series descriptors on model behavior we performed ablation studies, including SHAP analysis [43]. Table 4 provides a summary of ablation studies where only a certain group of descriptors or ts\_descriptors are used for training. Overall, for the target splits the time series descriptors based models perform better compared to their static models counterpart, however we note the difference is not major. The highest performance was obtained when training exclusively with pocket property ts\_descriptors, closely followed by pocket geometry ts\_descriptors. These results are in line with Figure 6 where the time series representation of pocket properties provided substantially more information useful for affinity correlation than a static representation.

Interestingly, for random splits, models trained only on static pocket property descriptors perform better than models trained on their time series counterparts. The situation changes



with target splits; the static models perform significantly worse compared to the MD-derived models, from which they were superior. One explanation of this result is the interdependence of pocket and ligand descriptors. Since pocket descriptors are calculated with respect to the ligand position, ligand information is implicitly contained, even more with a dynamic representation of the pocket.

The highest difference between crystal and MD derived models is obtained when training with interaction descriptors (0.450 vs 0.257). This result confirms our initial hypothesis that the interaction and motion descriptors when represented as time series provide novel and useful information in the context of affinity prediction. However these ts\_descriptors, when combined with other types fail to make a substantial difference in performance. This might indicate the need for a more thorough representation of motion and interactions present in ligand-receptor complexes, compared to the setup tested in this work.

An interesting result was achieved by the models trained only with ECFP4 and ligand properties. These models, having no information about the target, show elevated performance in affinity prediction, suggesting they mostly learn biases and random relationships in the data rather than predict affinity as a function of both target and ligand complex. Similar conclusions in the context of protein-ligand affinity predictions have been noted in other work as well [44–46].

View Article Online  
DOI: 10.1039/D5DD000452G



	Target split				Random Split			
	PCC		r <sup>2</sup>		PCC		r <sup>2</sup>	
Descriptors	MD	Crystal	MD	Crystal	MD	Crystal	MD	Crystal
All descriptors	<b>0.60</b> (0.07)	0.58 (0.10)	<b>0.32</b> (0.10)	0.29 (0.14)	<b>0.73</b> (0.04)	0.73 (0.05)	<b>0.53</b> (0.06)	0.51 (0.06)
Pocket property	<b>0.54</b> (0.07)	0.52 (0.07)	<b>0.26</b> (0.10)	0.23 (0.10)	0.68 (0.04)	<b>0.71</b> (0.04)	0.45 (0.05)	<b>0.50</b> (0.060)
Motion	0.53 (0.07)	-	0.24 (0.10)	-	0.58 (0.05)	-	0.33 (0.06)	-
Pocket geometric	<b>0.53</b> (0.07)	0.53 (0.08)	<b>0.23</b> (0.09)	0.22 (0.12)	<b>0.58</b> (0.04)	0.58 (0.05)	<b>0.33</b> (0.05)	0.33 (0.07)
Ligand property + ECFP4	0.52 (0.09)		0.21 (0.13)		0.69 (0.04)		0.46 (0.04)	
Interaction	<b>0.45</b> (0.08)	0.26 (0.14)	<b>0.15</b> (0.11)	-0.01 (0.10)	<b>0.51</b> (0.05)	0.46 (0.06)	<b>0.25</b> (0.05)	0.20 (0.05)
Ligand geometric	<b>0.44</b> (0.11)	0.42 (0.11)	<b>0.12</b> (0.17)	0.10 (0.15)	<b>0.51</b> (0.06)	0.49 (0.06)	<b>0.25</b> (0.06)	0.23 (0.06)

**Table 4. Ablation studies of the crystallographic and MD-derived models.** The 'Descriptors' column defines the sole group of descriptors on which the model has been trained. Standard deviation values are presented in brackets.

SHAP analysis presents the 20 most important features, together with their utilization in the form of counts, over all 20 models of each type; static- and time series based (Table S7). Each model is using a slightly different set of descriptors (or ts\_descriptors) depending on data splits it was trained on. SHAP analysis of a single static- and time series based model is presented on Figure S4 in the supplementary materials as reference.

One striking difference between the static and time series models is the heavy reliance on ligand descriptors by the former. Out of 20 most influential descriptors, nearly half of them (9/20, 45%) are ligand based. The same comparison with time series models shows they rely only on 25% (5/20) of ligand based ts\_descriptors. Even more, with static models 5/9 of the most influential ligand descriptors are simple 2D physicochemical features. With time series



models just 1/5. These results may explain the poorer performance of static models observed with the scaffold splits.

SHAP shows the MD-derived models rely mostly on a different set of features compared to models trained on static representations. There are three motion descriptors important for the time series models. The pocket contacts (both old and new) and ligand RMSD are especially interesting as they are exclusive to the MD-derived representation.

Taken together, the SHAP results show that the time series models use different sets of features, rely on simulation exclusive motion descriptors, and use less simple ligand features, rendering them possibly less prone to small molecule bias present in datasets.

## Conclusion

In this work, we introduce the MDD dataset, the largest publicly available collection of 200 ns molecular dynamics simulations of 862 protein-ligand complexes. This resource enables systematic investigation of time-series descriptor extraction strategies, feature-selection protocols, and provides a scalable foundation for extending both the size of training sets and the duration of simulations, thereby supporting the development and benchmarking of next-generation MD aware structure-based models in drug discovery.

The novelty of our approach lies in a feature-centric representation of MD data, where selected protein-ligand interaction properties are tracked over time and summarized using time-series descriptors (ts\_descriptors). Unlike snapshot-based augmentation strategies commonly employed in prior studies, this framework captures ensemble-level interaction tendencies through temporal statistics, enabling efficient learning from MD trajectories using tabular representations.

One of the central findings of this study is that longer MD simulations do not necessarily improve predictive performance. Across more than 800 simulations, models trained on descriptors derived from 20 ns trajectories consistently matched or outperformed those based on 200 ns simulations. This effect is likely attributable to increased variational noise introduced at longer timescales, which is difficult for machine learning models to filter. We further



investigated whether the optimal MD simulation length depends on macroscopic protein-ligand properties, such as protein or ligand size, or intrinsic flexibility. Despite extensive analysis across the MDD dataset, we did not observe any consistent relationships between the best-performing simulation length and these features (see supplementary Figure S6). Nevertheless, we note that more subtle dependencies may emerge when considering specific protein families, enzymatic classes, or systems involving pronounced allosteric motions, which remain important directions for future investigation.

Another important finding of this study is even a rather generic MD protocol and a relatively small number of complexes can be used with success to achieve predictive accuracy on par or better than highly complex models based on neural networks, with a much larger number of parameters. It is therefore expected that increasing the number of short MDs will further improve prediction performance. At scale this conclusion brings hope to the inclusion of short MD simulations into protocols concerning diverse chemical library screening and hit prioritization, and is also consistent with some previous works done on single targets [10,14,15].

The choice of traditional machine learning methods over deep learning was motivated by the tabular, physically interpretable nature of the descriptor set, as well as the need for transparency and robustness in the presence of correlated interaction features. While deep learning models have demonstrated strong performance in affinity prediction, the performance gap remains limited for engineered descriptors, and classical approaches provide a strong, interpretable baseline. However, the MDD dataset provides MD trajectories that can be used to construct more sophisticated representations, including graph-based or sequence-aware encodings that are better suited for deep learning architectures, for further comparison and benchmarking.

Comparisons between models trained on static crystallographic descriptors and those derived from MD trajectories reveal that the two approaches consider different descriptors to be most relevant. Interestingly, we note a number of time series derived descriptors with significantly better correlations compared to their static counterparts (Figure 6). Their summarized



influence however did transfer only slightly to improved affinity prediction performance. Given that extensive cross-correlation filtering was performed, this would suggest that non-linear correlations may decrease the overall performance. Ablation analysis and SHAP studies (see supplementary Figure S3) further confirm these findings, showing the two models employ different types of descriptors. In the case of more challenging target splits, an advantage of MD-derived models can also be observed, highlighting potential generalizability advantages of the time series descriptors.

In both the screening and pose selection tasks, MD-based models showed better performance than models based solely on static crystallographic structures. However, the higher standard deviation in the results of the MD models suggests a greater sensitivity to the selection of the initial conformation, which may affect the reproducibility and stability of the prediction. This phenomenon may be due to the greater complexity of dynamic representations and should be further investigated, taking into account different classes of molecular targets and simulation conditions. Importantly, while the descriptors used in this study are derived from molecular dynamics trajectories, they represent statistical summaries of interaction properties sampled over time, rather than an explicit encoding of discrete dynamic events. The time series descriptors (ts\_descriptors) capture ensemble-level tendencies, such as the persistence and variability of classical protein-ligand features, such as van der Waals contacts, hydrophobic interactions, hydrogen bonding, etc., aggregated across the simulation window. Our analysis shows the model's performance emerges from multivariate combinations of interaction features rather than individual interaction terms.

Consequently, the presented representation should be viewed as an approximation of interaction dynamics, designed with physical interpretability, robustness, and scalability for machine-learning applications. Capturing true dynamic processes, such as interaction lifetimes, state transitions, or allosteric shifts would require alternative time-resolved or event-based representations, which remain an important direction for future work. Within this context, the MDD dataset and the descriptor framework introduced here provide a reproducible and



extensible baseline upon which more mechanistically explicit modeling approaches can be developed and benchmarked.

In summary, we demonstrate that short MD simulations combined with time-series descriptor representations and classical machine learning models can achieve predictive performance on par with, and in some cases exceeding, models based on static structures, while offering improved generalization and scalability. Although the developed MDD dataset remains limited in size relative to static structures, our findings already show strong potential of MD aware models and lay foundations for further development of larger and more diverse protein-ligand MD collections.

## Acknowledgements

This work was sponsored by grant 2020/39/B/ST4/02747 obtained by PS from the Polish National Science Center. Computational resources were partially provided by the POL-OPENSREEN HE ERIC project.

## Data and code availability

Descriptor generation scripts are available from

Github: [https://github.com/JPoziemski/md\\_for\\_affinity\\_prediction](https://github.com/JPoziemski/md_for_affinity_prediction) and

Zenodo: <https://doi.org/10.5281/zenodo.18805105>

Trajectories of molecular dynamics are deposited at

Zenodo: <https://doi.org/10.5281/zenodo.18805105>

PDBBind 2020 R1 dataset was downloaded from: <https://www.pdbbind-plus.org.cn/>

DUD-E dataset was downloaded from: <https://dude.docking.org>



## Authors contribution

J.P. was responsible for methodology design, software development, investigation and data curation. A.Y. contributed to data interpretation and investigation. P.S. led the conceptualization of the study, methodology formulation, funding acquisition, data curation and drafted the original manuscript. All authors have approved the final manuscript.

## References

1. Kairys V, Baranauskiene L, Kazlauskiene M, Matulis D, Kazlauskas E. Binding affinity in drug design: experimental and computational techniques. *Expert Opin Drug Discov.* 2019;14: 755–768.
2. Mobley DL, Gilson MK. Predicting Binding Free Energies: Frontiers and Benchmarks. *Annu Rev Biophys.* 2017;46: 531–558.
3. Parks CD, Gaieb Z, Chiu M, Yang H, Shao C, Walters WP, et al. D3R grand challenge 4: blind prediction of protein-ligand poses, affinity rankings, and relative binding free energies. *J Comput Aided Mol Des.* 2020;34: 99–119.
4. Ballester PJ, Mitchell JBO. A machine learning approach to predicting protein-ligand binding affinity with applications to molecular docking. *Bioinformatics.* 2010;26: 1169–1175.
5. Stepniewska-Dziubinska MM, Zielenkiewicz P, Siedlecki P. Development and evaluation of a deep learning model for protein-ligand binding affinity prediction. *Bioinformatics.* 2018. doi:[10.1093/bioinformatics/bty374](https://doi.org/10.1093/bioinformatics/bty374)
6. Wang Z, Zheng L, Liu Y, Qu Y, Li Y-Q, Zhao M, et al. OnionNet-2: A Convolutional Neural Network Model for Predicting Protein-Ligand Binding Affinity Based on Residue-Atom Contacting Shells. *Front Chem.* 2021;9: 753002.
7. Su M, Yang Q, Du Y, Feng G, Liu Z, Li Y, et al. Comparative Assessment of Scoring Functions: The CASF-2016 Update. *J Chem Inf Model.* 2019;59: 895–913.
8. Shen C, Hu Y, Wang Z, Zhang X, Zhong H, Wang G, et al. Can machine learning consistently improve the scoring power of classical scoring functions? Insights into the role of machine learning in scoring functions. *Brief Bioinform.* 2021;22: 497–514.
9. Ganesan A, Coote ML, Barakat K. Molecular dynamics-driven drug discovery: leaping forward with confidence. *Drug Discov Today.* 2017;22: 249–269.
10. Gu S, Shen C, Yu J, Zhao H, Liu H, Liu L, et al. Can molecular dynamics simulations improve predictions of protein-ligand binding affinity with machine learning? *Brief Bioinform.* 2023;24. doi:[10.1093/bib/bbad008](https://doi.org/10.1093/bib/bbad008)
11. Gioia D, Bertazzo M, Recanatini M, Masetti M, Cavalli A. Dynamic Docking: A Paradigm



Shift in Computational Drug Discovery. *Molecules*. 2017;22.  
doi:[10.3390/molecules22112029](https://doi.org/10.3390/molecules22112029)

View Article Online  
DOI: 10.1039/D5DD00452G

12. Śledź P, Caflisch A. Protein structure-based drug design: from docking to molecular dynamics. *Curr Opin Struct Biol*. 2018;48: 93–102.
13. Guterres H, Im W. Improving Protein-Ligand Docking Results with High-Throughput Molecular Dynamics Simulations. *J Chem Inf Model*. 2020;60: 2189–2198.
14. Ash J, Fourches D. Characterizing the Chemical Space of ERK2 Kinase Inhibitors Using Descriptors Computed from Molecular Dynamics Trajectories. *J Chem Inf Model*. 2017;57: 1286–1299.
15. Jamal S, Grover A, Grover S. Machine Learning From Molecular Dynamics Trajectories to Predict Caspase-8 Inhibitors Against Alzheimer's Disease. *Front Pharmacol*. 2019;10: 780.
16. Kyaw Zin PP, Borrel A, Fourches D. Benchmarking 2D/3D/MD-QSAR Models for Imatinib Derivatives: How Far Can We Predict? *J Chem Inf Model*. 2020;60: 3342–3360.
17. Abraham MJ, Murtola T, Schulz R, Páll S, Smith JC, Hess B, et al. GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX*. 2015;1-2: 19–25.
18. Altman N, Krzywinski M. The curse(s) of dimensionality. *Nat Methods*. 2018;15: 399–400.
19. RDKit: Open-source cheminformatics. Available: <http://www.rdkit.org>
20. Michaud-Agrawal N, Denning EJ, Woolf TB, Beckstein O. MDAAnalysis: a toolkit for the analysis of molecular dynamics simulations. *J Comput Chem*. 2011;32: 2319–2327.
21. Bouysset C, Fiorucci S. ProLIF: a library to encode molecular interactions as fingerprints. *J Cheminform*. 2021;13: 72.
22. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nat Methods*. 2020;17: 261–272.
23. Overview on extracted features — tsfresh 0.20.2.post0.dev4+g3da2360 documentation. [cited 23 Jun 2024]. Available: [https://tsfresh.readthedocs.io/en/latest/text/list\\_of\\_features.html](https://tsfresh.readthedocs.io/en/latest/text/list_of_features.html)
24. Ramsundar B, Eastman P, Walters P, Pande V, Leswing K, Wu Z. *Deep Learning for the Life Sciences*. O'Reilly Media; 2019.
25. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine Learning in Python. *J Mach Learn Res*. 2011;12: 2825–2830.
26. Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System. *arXiv [cs.LG]*. 2016. Available: <http://arxiv.org/abs/1603.02754>
27. Menardi G, Torelli N. Training and assessing classification rules with imbalanced data. *Data Min Knowl Discov*. 2014;28: 92–122.
28. Wójcikowski M, Kukielka M, Stepniewska-Dziubinska MM, Siedlecki P. Development of a protein-ligand extended connectivity (PLEC) fingerprint and its application for binding



affinity predictions. *Bioinformatics*. 2019;35: 1334–1341.

View Article Online  
DOI: 10.1039/D5DD00452G

29. Durrant JD, McCammon JA. NNScore 2.0: a neural-network receptor-ligand scoring function. *J Chem Inf Model*. 2011;51: 2897–2903.
30. Trott O, Olson AJ. AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *J Comput Chem*. 2010;31: 455–461.
31. Cang Z, Mu L, Wei G-W. Representability of algebraic topology for biomolecules in machine learning based scoring and virtual screening. *PLoS Comput Biol*. 2018;14: e1005929.
32. Zhang S, Jin Y, Liu T, Wang Q, Zhang Z, Zhao S, et al. SS-GNN: A Simple-Structured Graph Neural Network for Affinity Prediction. *ACS Omega*. 2023;8: 22496–22507.
33. Mohamed Abdul Cader J, Newton MAH, Rahman J, Mohamed Abdul Cader AJ, Sattar A. Ensembling methods for protein-ligand binding affinity prediction. *Sci Rep*. 2024;14: 24447.
34. Wee J, Xia K. Ollivier Persistent Ricci Curvature-Based Machine Learning for the Protein-Ligand Binding Affinity Prediction. *J Chem Inf Model*. 2021;61: 1617–1626.
35. Liu X, Feng H, Wu J, Xia K. Dowker complex based machine learning (DCML) models for protein-ligand binding affinity prediction. *PLoS Comput Biol*. 2022;18: e1009943.
36. Jin Z, Wu T, Chen T, Pan D, Wang X, Xie J, et al. CAPLA: improved prediction of protein-ligand binding affinity by a deep learning approach based on a cross-attention mechanism. *Bioinformatics*. 2023;39. doi:[10.1093/bioinformatics/btad049](https://doi.org/10.1093/bioinformatics/btad049)
37. Zhang X, Gao H, Wang H, Chen Z, Zhang Z, Chen X, et al. PLANET: A Multi-objective Graph Neural Network Model for Protein-Ligand Binding Affinity Prediction. *J Chem Inf Model*. 2023. doi:[10.1021/acs.jcim.3c00253](https://doi.org/10.1021/acs.jcim.3c00253)
38. Jiménez J, Škalič M, Martínez-Rosell G, De Fabritiis G. KDEEP: Protein-Ligand Absolute Binding Affinity Prediction via 3D-Convolutional Neural Networks. *J Chem Inf Model*. 2018;58: 287–296.
39. Zheng L, Fan J, Mu Y. OnionNet: a multiple-layer inter-molecular contact based convolutional neural network for protein-ligand binding affinity prediction. *arXiv [physics.bio-ph]*. 2019. Available: <http://arxiv.org/abs/1906.02418>
40. Shen C, Zhang X, Hsieh C-Y, Deng Y, Wang D, Xu L, et al. A generalized protein-ligand scoring framework with balanced scoring, docking, ranking and screening powers. *Chem Sci*. 2023;14: 8129–8146.
41. Chatterjee S. A New Coefficient of Correlation. *J Am Stat Assoc*. 2021;116: 2009–2022.
42. Mysinger MM, Carchia M, Irwin JJ, Shoichet BK. Directory of useful decoys, enhanced (DUD-E): better ligands and decoys for better benchmarking. *J Med Chem*. 2012;55: 6582–6594.
43. Lundberg SM, Lee S-I. A Unified Approach to Interpreting Model Predictions. In: Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, et al., editors. *Advances in Neural Information Processing Systems*. Curran Associates, Inc.; 2017. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/8a20a8621978632d76c43dfd](https://proceedings.neurips.cc/paper_files/paper/2017/file/8a20a8621978632d76c43dfd)



[28b67767-Paper.pdf](#)

View Article Online  
DOI: 10.1039/D5DD00452G

44. Sieg J, Flachsenberg F, Rarey M. In Need of Bias Control: Evaluating Chemical Data for Machine Learning in Structure-Based Virtual Screening. *J Chem Inf Model.* 2019;59: 947–961.
45. Volkov M, Turk J-A, Drizard N, Martin N, Hoffmann B, Gaston-Mathé Y, et al. On the Frustration to Predict Binding Affinities from Protein-Ligand Structures with Deep Neural Networks. *J Med Chem.* 2022;65: 7946–7958.
46. Libouban P-Y, Aci-Sèche S, Gómez-Tamayo JC, Tresadern G, Bonnet P. The Impact of Data on Structure-Based Binding Affinity Predictions Using Deep Neural Networks. *Int J Mol Sci.* 2023;24. doi:[10.3390/ijms242216120](https://doi.org/10.3390/ijms242216120)



# Assessment of molecular dynamics time series descriptors in protein-ligand affinity prediction.

Jakub Poziemski <sup>1</sup>, Artur Yurkevych <sup>2</sup>, Paweł Siedlecki <sup>1\*</sup>

<sup>1</sup> Institute of Biochemistry and Biophysics, Polish Academy of Sciences, Warsaw, Poland

<sup>2</sup> Institute of Chemistry, University of Silesia in Katowice, Katowice, Poland

\* Corresponding author: [pawel@ibb.waw.pl](mailto:pawel@ibb.waw.pl)

## Data and code availability

Descriptor generation scripts are available from

Github: [https://github.com/JPoziemski/md\\_for\\_affinity\\_prediction](https://github.com/JPoziemski/md_for_affinity_prediction) and

Zenodo: <https://doi.org/10.5281/zenodo.18805105>

Trajectories of molecular dynamics are deposited at

Zenodo: <https://doi.org/10.5281/zenodo.18805105>

PDBBind 2020 R1 dataset was downloaded from: <https://www.pdbbind-plus.org.cn/>

DUD-E dataset was downloaded from: <https://dude.docking.org>

