

Cite this: *Chem. Sci.*, 2019, 10, 10911

All publication charges for this article have been paid for by the Royal Society of Chemistry

# DeltaDelta neural networks for lead optimization of small molecule potency†

José Jiménez-Luna,<sup>a</sup> Laura Pérez-Benito,<sup>bc</sup> Gerard Martínez-Rosell,<sup>f</sup> Simone Sciabola,<sup>d</sup> Rubben Torella,<sup>e</sup> Gary Tresadern<sup>bc</sup> and Gianni De Fabritiis<sup>id</sup> \*<sup>afg</sup>

The capability to rank different potential drug molecules against a protein target for potency has always been a fundamental challenge in computational chemistry due to its importance in drug design. While several simulation-based methodologies exist, they are hard to use prospectively and thus predicting potency in lead optimization campaigns remains an open challenge. Here we present the first machine learning approach specifically tailored for ranking congeneric series based on deep 3D-convolutional neural networks. Furthermore we prove its effectiveness by blindly testing it on datasets provided by Janssen, Pfizer and Biogen totalling over 3246 ligands and 13 targets as well as several well-known openly available sets, representing one the largest evaluations ever performed. We also performed online learning simulations of lead optimization using the approach in a predictive manner obtaining significant advantage over experimental choice. We believe that the evaluation performed in this study is strong evidence of the usefulness of a modern deep learning model in lead optimization pipelines against more expensive simulation-based alternatives.

Received 13th September 2019

Accepted 15th October 2019

DOI: 10.1039/c9sc04606b

rsc.li/chemical-science

## 1 Introduction

In the lead optimization phase of drug discovery, the chemical structure of a molecule is typically modified by a medicinal chemistry team with the intent of improving its potency, selectivity, and many other pharmacokinetic and toxicological parameters.<sup>1–3</sup> These modifications result in congeneric series, a set of ligands with few atom changes between them, usually around a unique or small number of different scaffolds for which there are experimental structures of the complex with the target protein. Series range from few hundred to thousands of compounds and require considerable human, time and financial resources for synthesis and assays. It is therefore of great value to have *in silico* predictive tools to accelerate this process. Series typically feature very small potency differences, which in turn is a challenge for predictors, as having what could be

considered a low error in other scenarios (e.g. below 1 kcal mol<sup>−1</sup>) is not a guarantee for successful ranking.

It is therefore common to focus on relative binding free energy (RBFE) simulation methods,<sup>4–13</sup> where the difference in affinity between two ligands is computed using a thermodynamic cycle that alchemically perturbs only the small region associated with the changing atoms. RBFE methods have shown good results in several studies, with accuracy close to 1 kcal mol<sup>−1</sup> and reasonable correlations. Despite this, these methods suffer from several issues, such as system preparation, treatment of waters, force-field selection, protein flexibility and computational cost, making their prospective application difficult in practice.<sup>14</sup> On the other side, many empirical,<sup>15,16</sup> knowledge-based<sup>17,18</sup> and machine learning<sup>19–24</sup> scoring functions have been designed for the task of predicting absolute binding affinities. They mostly tackle the problem in a regression setup, where the binding affinity is to be predicted using a set of protein–ligand descriptors, modelling the interaction among both. The fact that they model absolute affinities and are trained on very chemically diverse bodies of data, such as iterations of the PDBbind<sup>25</sup> database, limits their applicability when predicting small structural differences between two ligands, such in the congeneric series case. While other machine learning approaches have been presented for this task,<sup>26–28</sup> here we propose a modern 3D-convolutional-neural-network-based continuous learning approach for relative binding affinity prediction in congeneric series and show strong predictive power using multiple blind benchmarks as well as public datasets at negligible computational costs. This study

<sup>a</sup>Computational Science Laboratory, Parc de Recerca Biomèdica de Barcelona, Universitat Pompeu Fabra, C Dr Aiguader 88, Barcelona, 08003, Spain. E-mail: gianni.defabritiis@upf.edu

<sup>b</sup>Laboratori de Medicina Computacional, Unitat de Bioestadística, Facultat de Medicina, Universitat Autònoma de Barcelona, Spain

<sup>c</sup>Janssen Research and Development, Turnhoutseweg 30, 2340 Beerse, Belgium

<sup>d</sup>Biogen Chemistry and Molecular Therapeutics, 115 Broadway Street, Cambridge, MA 02142, USA

<sup>e</sup>Pfizer I&I, 610 Main Street, Cambridge, MA 02139, USA

<sup>f</sup>Acellera, Carrer del Dr Trueta, 183, 08005 Barcelona, Spain

<sup>g</sup>Institució Catalana de Recerca i Estudis Avançats (ICREA), Passeig Lluís Companys 23, 08010 Barcelona, Spain

† Electronic supplementary information (ESI) available. See DOI: 10.1039/c9sc04606b

serves as a very large evaluation of a modern machine-learning pipeline for lead optimization in a real-life drug discovery scenario, thanks to the joint collaboration with several pharmaceutical companies.

## 2 Materials and methods

### 2.1 Datasets studied

The BindingDB protein–ligand validation sets<sup>29</sup> were used to pretrain our models. For testing, we also extracted well-known publicly-available literature test sets<sup>30</sup> used for benchmarking RBFE calculations. Furthermore we include a recent freely-available BRD4 bromodomain dataset.<sup>31</sup> In regards to internal pharmaceutical data, we tested on five different congeneric series from Janssen R&D. Three chemical series (sets 1, 2 and 3) were phosphodiesterase 2 (PDE2) inhibitors with bioactivity *versus* PDE2, PDE3, and PDE10 (ref. 32 and 33) (publication number WO2018083103A1), the fourth series were proto-oncogene tyrosine kinase (ROS1) inhibitors (publication number WO2015144799A1) and the final beta-secretase 1 (BACE1) inhibitors.<sup>34</sup> We tested six congeneric series with Pfizer, three of which target a kinase, and the remaining an enzyme, a phosphodiesterase (PDE) and an activator of transcription. The sizes of these vary from 93 molecules up to 362, for a total of 955 tested compounds. Lastly, Biogen tested the proposed procedure on two different series, composed of 196 and 220 analogues targeting a tyrosine-protein kinase and a receptor-associated kinase, respectively. The size of the sets presented here allow, to the best of our knowledge the largest evaluation yet of a modern machine learning pipeline in lead optimization.

Out of the total 645 available congeneric series available in BindingDB, 495 with IC<sub>50</sub> affinity values were extracted and processed for further evaluation, as it was the unit with most data available, containing a diverse set of targets. The majority of these sets encompass a single protein–ligand crystal structure, the rest of the ligands modelled against the reference using the Surflex docking software.<sup>35</sup> We then assign each protein structure in the database to a family cluster using a 90% sequence similarity threshold, as per PDB conventions.<sup>36</sup> For each series in the same protein cluster we use a maximum common substructure (MCS) protocol as available in rdkit<sup>37</sup> to remove identical ligands. This procedure ensures that the same ligand is not repeated against similar targets, avoiding potential overfitting problems and overoptimistic evaluations.<sup>38</sup> Affinity values were log-converted to avoid target scaling issues (pIC<sub>50</sub> = −log<sub>10</sub> IC<sub>50</sub>). Ligands that could not be read by rdkit were removed. Histograms of the number of ligands and their affinity range per series are provided in ESI Fig. 7,† with the average available number of ligands per series being 8.84. In the Schrödinger and BRD4 sets, since only Δ*G* (per kcal mol<sup>−1</sup>) information was available, we converted affinity values to the pIC<sub>50</sub> range assuming non-competitive binding. Descriptive information on these series is provided in ESI Table 1.† Compounds provided by Janssen were docked using a common scaffold structure *via* the Glide software. These congeneric series range from 48 up to a 900 different compounds with varying affinity ranges (ESI Table 2†).

### 2.2 Descriptor calculation

We have recently reported a machine learning approach that can learn based on 3D features of the binding site interactions.<sup>20</sup> A similar encoding was used here that represents the protein–ligand binding by voxelizing both using a 24 Å pocket centered box, with a density of 1 Å<sup>3</sup> per voxel. The contribution of each atom to each voxel is inversely proportional to their Euclidean distance *r* and the van der Waals radius *r*<sub>vdw</sub> of the first. We use several *channels* for both protein and ligand, in the sense that the atomic contribution to each voxel depends on their type. The contribution of each atom to each voxel is assigned according to a pair correlation function defined by:

$$n(r) = 1 - \exp\left(-\left(\frac{r_{\text{vdw}}}{r}\right)^{12}\right). \quad (1)$$

We define several *channels* for both protein and ligand, in the sense that the atomic contribution to each voxel depends on their type. For the protein we define eight pharmacophoric-like descriptors, as detailed in ESI Table 3.† For the ligands we use a simpler representation based on atom types contained in the set {C,N,O,F,P,S,Cl,Br,I,H}, for a total of 18 stacked channels. We note that there is no particular reasoning behind this choice of descriptors other than they showed promising practical performance in previous studies.

### 2.3 Neural network architecture

Regular feed-forward neural networks do not scale well when the input is high dimensional (as in images, or in this case atomic interactions). CNNs are specifically designed for handling lattices, where local spatial information needs to be preserved. While a feed-forward network would ignore such interactions, a convolutional one arranges its neurons spatially, and only connects locally to the output of the previous layer. The latter type of architectures have become the de-facto workhorse in computer vision problems,<sup>39–41</sup> providing state-of-the-art performance. Due to this success, many applications in bio-informatics and computational chemistry followed.<sup>42–51</sup>

The neural network we propose has a novel zero-symmetric architecture whose main building blocks are 3D-convolution operations. In this work we focus on predicting relative affinities for close analogues in lead optimization, therefore, our approach is to build a network whose input is a pair of ligand binding voxelized representations belonging to the same series. A two-input convolutional neural network was designed, with fixed weights on both legs. The inputs are forwarded through several convolution and pooling operations and then flattened into a 192-dimensional latent representation. The symmetry property of relative binding affinity requires that inverting the order of the input ligands should change the sign of the predicted value, and we embed such symmetry in the network by computing the difference between latent representations. A final linear layer with no bias is then applied to the result of this difference, ensuring zero-symmetry by design and producing the desired predicted difference in affinity. In contrast, calculating relative affinities from an absolute predictor inevitably



leads to the concatenation of errors from two separate predictions.

A schema of our architectural choice is provided in Fig. 1 is provided in the ESI.† It features two convolution operations with a kernel size of 3 in each leg, followed by a max-pooling operation, and finally another convolution operation with the same kernel size for both before flattening and performing the latent difference between analogues. The ReLU activation function was used for all layers in the network except for the last, which does not feature one. We include a dropout layer in the end to control for overfitting. Xavier initialization was used for the weights. Training is performed using the Adam stochastic gradient descent optimizer<sup>52</sup> with standard hyperparameters ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-4}$ ) using a batch size of 32 samples for 50 epochs. Data augmentation is performed during training by randomly rotating input coordinates to mitigate the lack of rotational equivariance in CNNs. Furthermore, given a set of relative binding predictions, its absolute counterparts can always be retrieved given a single experimentally determined absolute reference, such as the one provided by a lead. If more than one is available, absolute affinities can be computed towards each, in practice providing a predictive absolute affinity distribution, whose average can then be interpreted as a maximum *a posteriori* (MAP) estimate of the absolute affinity and its standard deviation as a measure of its uncertainty, given the current model state.

## 2.4 Continuous learning

In the proposed continuous-learning approach, we explicitly use the fact that congeneric series are sequentially generated in a lead optimization campaign, and follow an incremental training and testing procedure. For each congeneric series at a given time the affinity of previously tested ligands is known experimentally: differences for these are taken as training data, while for test data we predict differences between unknown and known ones. While this approach is less ambitious than having a predictor for relative affinity with no experimentally tested data (such as a physical-based model), its applicability is general, since it is the common scenario that medicinal chemists face in lead optimization campaigns. The training for the BindingDB sets starts with a reference structure in each series, for which we take the crystal structure ligand if available or the structure with the lowest average maximum common substructure (MCS) distance to the rest. Ligands from the rest of

the series are then sequentially added in a random order. It is well known that either a random<sup>53</sup> or scaffold-based training test split produce overoptimistic results when testing machine-learning algorithms on activity benchmarks. Since the industrial datasets in our study include a compound creation time-stamp, we also evaluate a more realistic temporal split,<sup>54</sup> where at each training step we consider the first  $n$  tested ligands and the differences of the posterior ones against the first are taken. The performance of the machine learned models is reported as the root mean squared error (RMSE) and either Pearson's correlation coefficient  $R$  or Spearman's  $\rho$  between experimental and predicted affinity differences. We note that in all blind tests a single model was provided, and no explicit attempt to optimize hyperparameters in each set was made.

## 3 Results

We first present results concerning our validation on the 495 protein-series datasets from the BindingDB, where the proposed model achieves an average correlation coefficient above 0.4 and an RMSE below 1.25 (pIC<sub>50</sub> units) even when only one binding-energy difference is taken per congeneric series (ESI Fig. 1†). This suggests that the method works reasonably well in the very low-data scenario, such as the beginning of a lead optimization campaign. A noticeable performance boost is seen as more differences are included in training, with a correlation coefficient above 0.62 and an RMSE below 1.05 when another four different ligands from the same congeneric series are known in advance, with performance plateauing beyond five additional training ligands. A comparison against an absolute affinity model is also provided (*i.e.* one of the legs of the architecture), where as expected it can be appreciated that it performs considerably worse than its relative counterpart.

Now we present results on the Wang *et al.*<sup>30</sup> and BRD4 inhibitor datasets.<sup>31</sup> In this and the rest of cases, we pretrained a model with all difference pairs available in the BindingDB database, which provides a prior for further fine-tuning. We then mixed new available data as training in each sequential iteration of each set with the rest of the BindingDB database for only 3 epochs, significantly reducing computational overhead. A FEP baseline provided by Wang *et al.*<sup>30</sup> is used for comparison. The model efficiently interpolates differences for unseen ligands, achieving considerably high correlation coefficients and low errors in all series with as few as 3–4 additional ligands and associated activity pairs, surpassing in many cases the much more expensive FEP baseline (ESI Fig. 2†). For instance, for the MCL1 target, after testing 3 ligands, the correlation coefficient is above 0.8, surpassing the FEP baseline, and the RMSE is below 1.2 (pIC<sub>50</sub> units).

The same evaluation procedure was taken for the compounds available in the Janssen PDE sets (Fig. 2 and ESI Fig. 3†) for both a random and a temporal split, where a baseline against Glide score<sup>55</sup> is also added. Excellent performance was seen on a random split given enough training data, and as expected, although the temporal split performance is lower, it is still sufficiently high to be used in a real-life prospective lead optimization scenario. For instance, for the first PDE2 activity

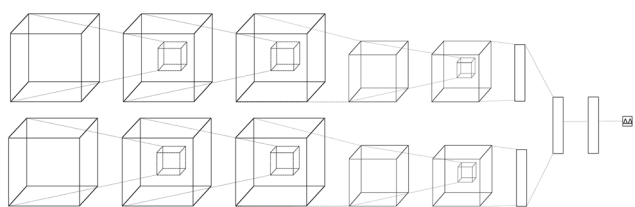


Fig. 1 Architecture of the proposed model. A two-legged neural network with tied weights was constructed, and a pair of protein–ligand voxelization is feed-forwarded through it to later perform a latent space difference.



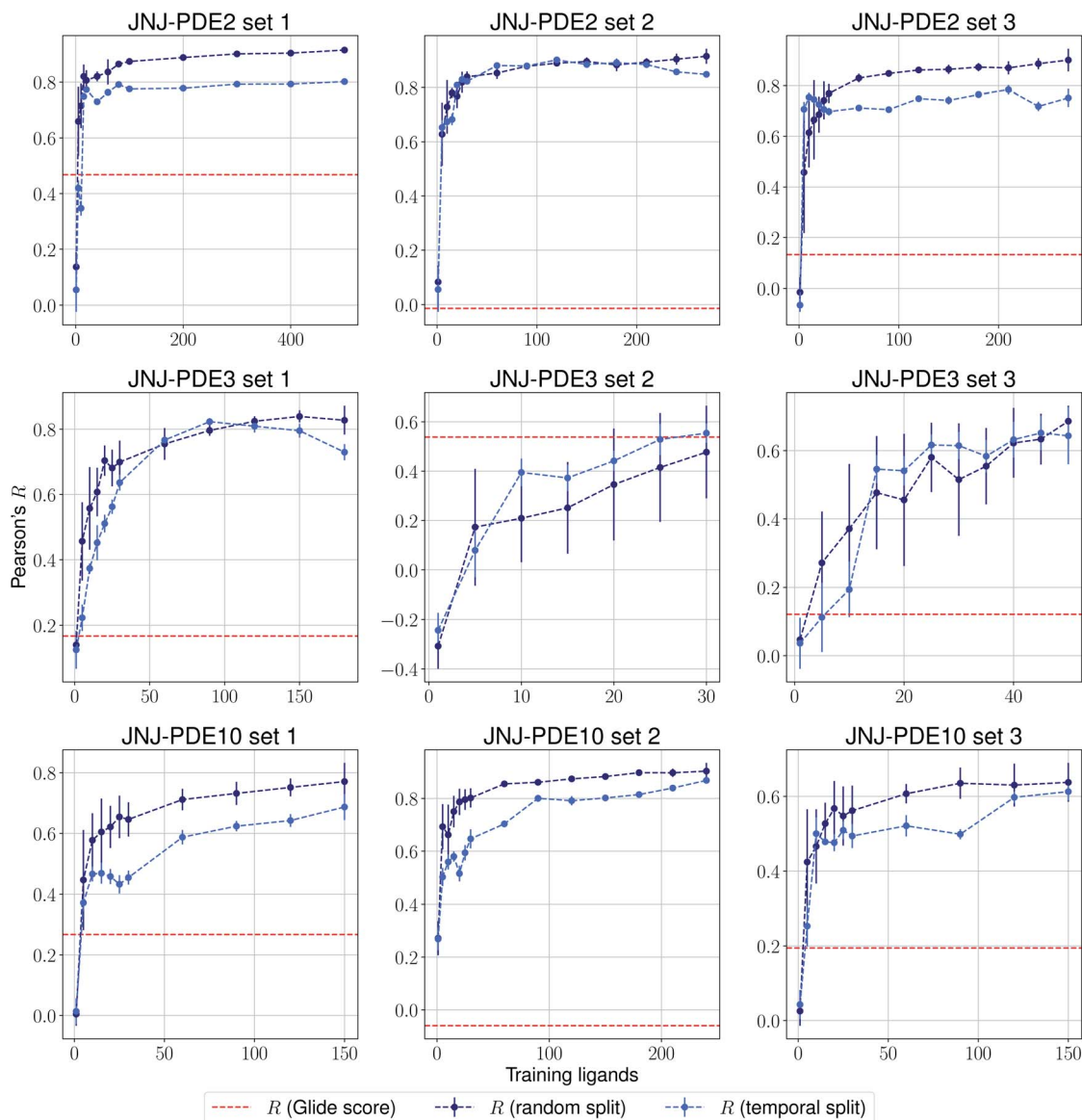


Fig. 2 Average Pearson's correlation coefficient  $R$  ( $\pm 1$  standard deviation) based on 25 independent runs on different sets for the Janssen PDE2, PDE3 and PDE10 targets.

set after 20 ligands sorted by time, the Pearson's correlation coefficient  $R$  and RMSE were 0.77 and 1.35 (in  $\text{pIC}_{50}$  units) respectively. Results for the ROS1 and BACE sets, show a similar trend and insights (Fig. 3 and ESI Fig. 4<sup>†</sup>). Interestingly, performance slightly decreases for the BACE target late in the evaluation, suggesting that the network may have found an unexplored activity cliff or that the last tested ligands are harder to predict than others in the series. Furthermore, we also provide a type of split where only differences among the most chemically close ligands are predicted, based on ECFP4 fingerprint similarity, as available in rdkit. That is, in each training step we predict from the remaining untested pool of ligands those that are closest to the ones in our training set, with the intention of resembling a real-life lead optimization RBFE scenario, typically applied to close analogues. Split-based

results on fingerprint similarity for the first PDE2 set (ESI Fig. 5<sup>†</sup>), show that after 20 ligands sorted by chemical similarity the  $R$  and RMSE were 0.83 and 1.12 (in  $\text{pIC}_{50}$  units). These suggest better performance in this scenario than the proposed temporal split, and closer to the random one.

We then present the results provided by Pfizer using a temporal split in Table 1, where specific target names cannot be disclosed. We compare such results with several baselines such as molecular weight,  $\text{clog } P$ , a MM-GBSA pipeline<sup>56,57</sup> and deep-learning absolute affinity predictor  $K_{\text{DEEP}}$ ,<sup>20</sup> trained on the v.2016 iteration of the PDBbind database. The model proposed here performs considerably better than the rest when given only 10% of the training data, again highlighting the importance of incrementally training these on the congeneric series of interest. An exception, however, is found in the Kinase #3 series,





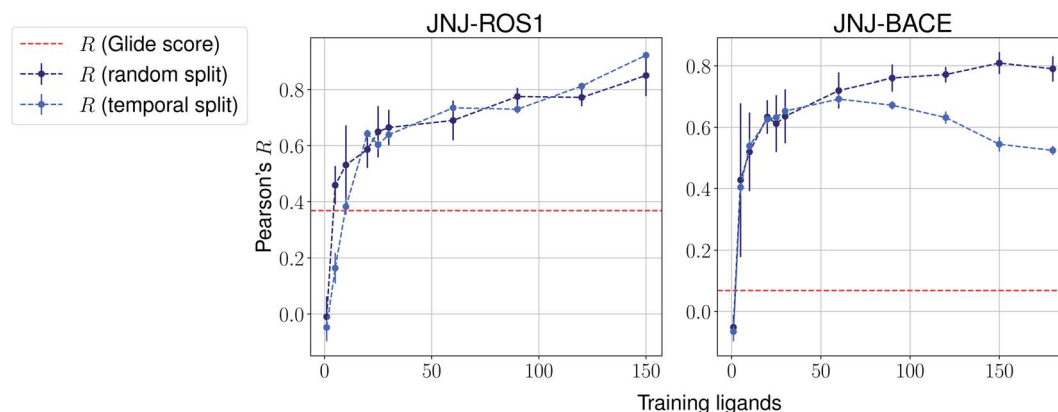


Fig. 3 Average Pearson's correlation coefficient  $R$  ( $\pm 1$  standard deviation) based on several independent runs on two sets for the Janssen ROS1 and BACE targets.

for which no significant improvement is observed when providing extra training data. This particular last case appears to be particularly hard to predict, as all tested methods perform poorly. We provide results using a temporal split for the last two congeneric series provided by Biogen, for which we also compare against several baselines: (a) Glide score, (b) an MM-GBSA pipeline, and (c) a standard QSAR approach using MACCS, ECFP4 and rdkit descriptors with a random forest model (Fig. 4 and ESI Fig. 6†). Our model reveals similar conclusions, significantly outperforming all baselines. Curiously, it can also be seen that the proposed method does not perform significantly worse than the aforementioned baselines in the second target when no training data is used. When some is used, such as only 5 analogues, our proposed machine-learning model significantly outperforms all baselines.

All the tests regarding internal pharmaceutical data were carried out blindly by providing fully-containerized software to our collaborators, who executed the application and reported corresponding results. Furthermore only one pretrained model was provided without any opportunity to overfit to each specific test set.

One aim of our study was to test whether machine-learning driven relative affinity predictions could efficiently identify key

high potency compounds in a close to real-life lead optimization scenario, by retrospectively comparing them to the experimental order of synthesis. With some of the large industrial datasets it was possible to test this and we used the most active compound as a surrogate interesting lead molecule. The model is trained on the first experimentally tested compounds, and then is incrementally trained by choosing from the remaining ones based on an upper confidence bound (UCB)-like criterion,<sup>58</sup> defined as:

$$\text{UCB} = \mu(x) + \beta\sigma(x), \quad (2)$$

where  $\mu$  and  $\sigma$  are the average and standard deviation predicted absolute affinities provided by the model for ligand  $x$  and  $\beta$  is a user-chosen factor controlling the balance between exploitation and exploration, that we fix in our study to  $\beta = 1.64$ .

We stop the procedure once the model retrieves the analogue with the highest associated affinity, and compare this with its original synthesis experimental order in its corresponding series. We present results for this simulation-based benchmark in Table 2. In 4 out of 5 sets our proposed model is able to reach the compound with the highest affinity faster than its experimental order or by random selection. Surprisingly, in all ten independent runs of the second set for the PDE2 target, the

Table 1 Spearman's  $\rho$  performance results between experimental and predicted absolute affinities provided by Pfizer I&I, where other empirical, simulation, and machine-learning based affinity prediction methods are compared on several congeneric series. Performance is poor for most tested model except for the sequential approach proposed here, with Pearson correlations averaging over 0.5 with as few as 10% used analogues from the congeneric series at hand

Target	# ligands	Mol. weight ( $\rho$ )	clog $P^a$ ( $\rho$ )	MM-GBSA ( $\rho$ )	$K_{\text{DEEP}}$ ( $\rho$ )	This work (10% training, $\rho$ )	This work (20% training, $\rho$ )	This work (30% training, $\rho$ )
Kinase #1	362	0.19	0.06	0.56	0.42	0.49	0.64	0.73
Kinase #2	106	0.1	0.28	0.25	0.25	0.25	0.41	0.51
Kinase #3	95	0	0.04	0.25	-0.27	0.3	0.3	0.31
Enzyme	93	0.43	0.24	0.01	0.49	0.43	0.26	0.59
Phosphodiesterase	100	0.37	0.36	0.67	0	0.49	0.64	0.73
Activator of transcriptions	199	0.13	0.08	0.66	0.29	0.72	0.84	0.94
Weighted avg.		0.19	0.14	0.47	0.25	<b>0.49</b>	<b>0.59</b>	<b>0.69</b>
Simple avg.		0.2	0.18	0.4	0.18	<b>0.45</b>	<b>0.52</b>	<b>0.64</b>

<sup>a</sup> Calculated log  $P$  as available in rdkit.



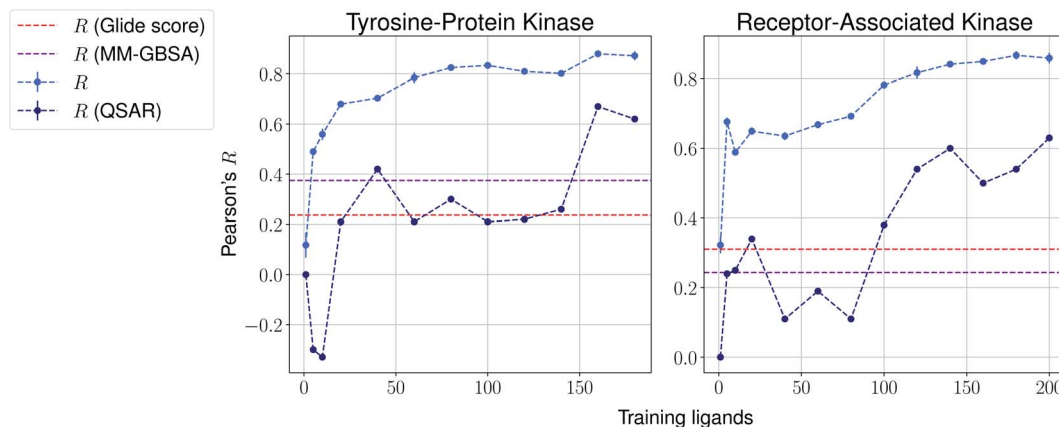


Fig. 4 Results over 5 runs on Biogen's tyrosine-protein kinase and receptor-associated kinase using a temporal split, and MM-GBSA and QSAR random forest pipelines as baselines.

compound with the highest affinity was found after only a single synthesis epoch. Furthermore, one would expect the average affinity in the training set to increase at each synthesis epoch (as the model is tasked to pick compounds with increasingly higher UCB). This is the case for 4 out of 5 sets again (Fig. 5), with the exception of the ROS1 target, which shows a non-monotonic trend, albeit its model reaches the compound with highest affinity before its experimental order. In all tested cases, the average training pool affinity for the ligands selected by the model is higher than experimental choice. Overall results are very promising and suggest that the proposed method could be applied successfully in a prospective scenario. Particularly, in the first PDE2 set, we were able to reach potent compounds synthesizing up to six times less molecules than the baseline method used by the medicinal chemistry team.

## 4 Discussion

In this work we have designed and tested a deep-learning based model for the task of predicting relative binding affinity predictions in congeneric series. This work provides evidence that the method is able to efficiently rank compounds as shown by an evaluation on both publicly available and industrial data and can be of use by computational and medicinal chemists in early drug-discovery projects by providing informed choices of future compounds to synthesize, as suggested by our

simulation-based benchmark. The accuracy of the method heavily depends on the amount of available data but can be trained and applied in minutes on a single GPU, offering a substantial improvement in performance compared with physics-based RBFE calculations which can take days for a small number of analogues. While the results presented here are encouraging, it is important to note that they remain retrospective: a proper prospective validation of the model, which would entail chemists synthesizing compounds according to the decisions taken by the trained model, remains a topic of future study. In the long term, however, we expect that improving molecular simulations accuracy<sup>59,60</sup> by the integration of physics and machine learning approaches would produce a more convenient approach for engineering drug discovery. In the meantime, methods such as the one proposed here provide accurate performance at a fraction of the computational cost of other approaches.

## Code & data availability

All the models here were developed using the PyTorch package for tensor computation and neural network training.<sup>61</sup> BindingDB, Wang *et al.* and Mobley *et al.* set results are available upon reasonable request. Python code for generating the proposed featurization is available within the open-source HTMD software.<sup>62</sup> The code of the network architecture in

**Table 2** Simulation-based benchmark results over 10 independent runs for the different datasets. We show the amount of molecules the model is allowed to pick at each synthesis epoch, the experimental order of the compound with the highest affinity in the series, the average synthesis epoch our model found said molecule, the total necessary sampled ligands the proposed model has chosen before the target compound, and the sampling advantages over the experimental and random orders

Target	Set	# ligands	Chosen per synthesis epoch	Experimental order	Found at synthesis epoch	Total sampled ligands	Advantage over experimental choice	Advantage over random choice
PDE2	1	900	10	766	12.2	132	<b>634</b>	<b>318</b>
PDE2	2	303	10	61	1	20	<b>41</b>	<b>131.5</b>
PDE2	3	278	10	253	5.9	69	<b>184</b>	<b>70</b>
ROS1	—	165	10	73	3.1	41	<b>32</b>	<b>41.5</b>
BACE	—	229	10	190	20.8	218	−28	−103.5



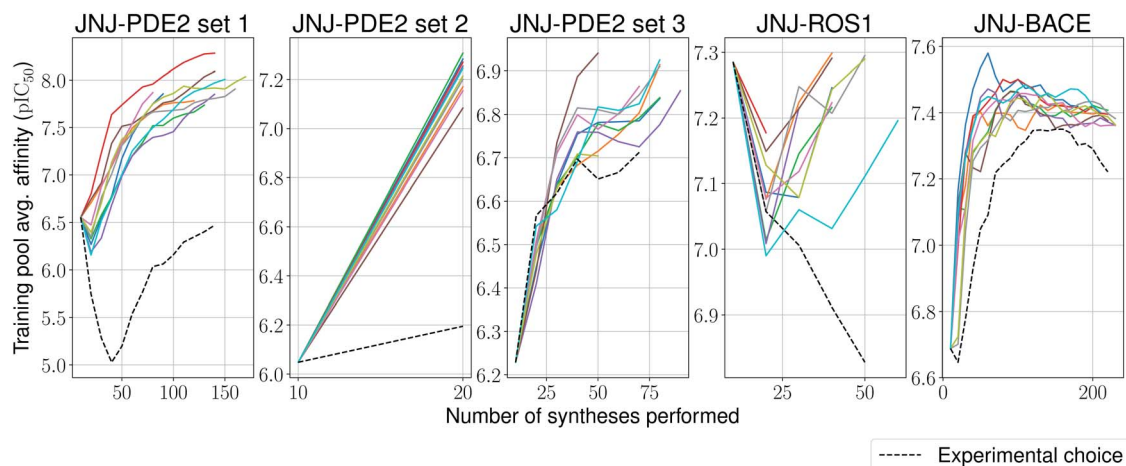


Fig. 5 Average model-picked training set affinity per number of compounds synthesized for the Janssen PDE2, ROS1 and BACE sets, as well as a baseline based on the actual experimental choice order of compounds.

a PyTorch implementation is provided in the ESI.† An implementation of this application is available through the <http://PlayMolecule.org> repository of applications, where users can freely submit their protein in PDB format and two sets of the same congeneric series, for training and validation respectively in SDF format. Depending on the size of these last two, training and prediction time may vary, as the order of data for training increases by  $\frac{n(n-1)}{2}$ , and for testing  $nm$  factors, where  $n$  and  $m$  are the number of training and testing instances respectively. At the moment, predictions are limited to a default total of a 1000 molecules per congeneric series, with runtimes averaging and hour on a modern GeForce 1080Ti GPU. Larger experiments can be arranged for users willing to run more computationally demanding experiments.

## Conflicts of interest

G. D. F is a founder and current CEO of Acellera Ltd. J. J. receives funding from Acellera Ltd. and G. M. R. is an employee.

## Acknowledgements

The authors thank Acellera Ltd. for funding. G. D. F. acknowledges support from MINECO (BIO2014-53095-P), MICINN (PTQ-17-09079) and FEDER. This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 675451 (CompBioMed project).

## Notes and references

- C. A. Nicolaou and N. Brown, *Drug Discovery Today: Technol.*, 2013, **10**, e427–e435.
- I. Kola and J. Landis, *Nat. Rev. Drug Discovery*, 2004, **3**(8), 711.
- S. Ekins, J. D. Honeycutt and J. T. Metz, *Drug discovery today*, 2010, **15**, 451–460.
- L. Wang, B. J. Berne and R. A. Friesner, *Proc. Natl. Acad. Sci. U. S. A.*, 2012, **109**, 1937–1942.
- E. B. Lenselink, J. Louvel, A. F. Forti, J. P. D. van Veldhoven, H. de Vries, T. Mulder-Krieger, F. M. McRobb, A. Negri, J. Goose, R. Abel, H. W. T. van Vlijmen, L. Wang, E. Harder, W. Sherman, A. P. IJzerman and T. Beuming, *ACS Omega*, 2016, **1**, 293–304.
- S. Wan, A. P. Bhati, S. Skerratt, K. Omoto, V. Shanmugasundaram, S. K. Bagal and P. V. Coveney, *J. Chem. Inf. Model.*, 2017, **57**, 897–909.
- D. A. Goldfeld, R. Murphy, B. Kim, L. Wang, T. Beuming, R. Abel and R. A. Friesner, *J. Phys. Chem. B*, 2015, **119**, 824–835.
- L. Pérez-Benito, H. Keränen, H. van Vlijmen and G. Tresadern, *Sci. Rep.*, 2018, **8**, 4883.
- M. Ciordia, L. Pérez-Benito, F. Delgado, A. A. Trabanco and G. Tresadern, *J. Chem. Inf. Model.*, 2016, **56**, 1856–1871.
- C. Schindler, F. Rippmann and D. Kuhn, *J. Comput.-Aided Mol. Des.*, 2017, **32**, 1–8.
- H. Keränen, L. Pérez-Benito, M. Ciordia, F. Delgado, T. B. Steinbrecher, D. Oehlrich, H. W. T. van Vlijmen, A. A. Trabanco and G. Tresadern, *J. Chem. Theory Comput.*, 2017, **13**, 1439–1453.
- G. Heinzelmann, N. M. Henriksen and M. K. Gilson, *J. Chem. Theory Comput.*, 2017, **13**, 3260–3275.
- M. Aldeghi, A. Heifetz, M. J. Bodkin, S. Knapp and P. C. Biggin, *Chem. Sci.*, 2016, **7**, 207–218.
- Z. Cournia, B. Allen and W. Sherman, *J. Chem. Inf. Model.*, 2017, **57**, 2911–2937.
- Y. Cao and L. Li, *Bioinformatics*, 2014, **30**, 1674–1680.
- M. P. Brenner, L. J. Colwell, et al., *Proc. Natl. Acad. Sci. U. S. A.*, 2016, **113**, 13564–13569.
- T. A. Halgren, R. B. Murphy, R. A. Friesner, H. S. Beard, L. L. Frye, W. T. Pollard and J. L. Banks, *J. Med. Chem.*, 2004, **47**, 1750–1759.
- O. Trott and A. Olson, *J. Comput. Chem.*, 2010, **31**, 455–461.
- P. J. Ballester and J. B. O. Mitchell, *Bioinformatics*, 2010, **26**, 1169–1175.

- 20 J. Jiménez, M. Škalič, G. Martínez-Rosell and G. De Fabritiis, *J. Chem. Inf. Model.*, 2018, **58**, 287–296.
- 21 E. N. Feinberg, D. Sur, Z. Wu, B. E. Husic, H. Mai, Y. Li, S. Sun, J. Yang, B. Ramsundar and V. S. Pande, *ACS Cent. Sci.*, 2018, **4**, 1520–1530.
- 22 M. Ragoza, J. Hochuli, E. Idrobo, J. Sunseri and D. R. Koes, *J. Chem. Inf. Model.*, 2017, **57**, 942–957.
- 23 D. D. Nguyen, Z. Cang, K. Wu, M. Wang, Y. Cao and G.-W. Wei, *J. Comput.-Aided Mol. Des.*, 2019, **33**, 71–82.
- 24 Z. Gaieb, C. D. Parks, M. Chiu, H. Yang, C. Shao, W. P. Walters, M. H. Lambert, N. Nevins, S. D. Bembenek, M. K. Ameriks, *et al.*, *J. Comput.-Aided Mol. Des.*, 2019, **33**, 1–18.
- 25 R. Wang, X. Fang, Y. Lu and S. Wang, *J. Med. Chem.*, 2004, **47**, 2977–2980.
- 26 W. Zhan, D. Li, J. Che, L. Zhang, B. Yang, Y. Hu, T. Liu and X. Dong, *Eur. J. Med. Chem.*, 2014, **75**, 11–20.
- 27 A. Amini, P. J. Shrimpton, S. H. Muggleton and M. J. Sternberg, *Proteins: Struct., Funct., Bioinf.*, 2007, **69**, 823–831.
- 28 D. Zilian and C. A. Sotriffer, *J. Chem. Inf. Model.*, 2013, **53**, 1923–1933.
- 29 T. Liu, Y. Lin, X. Wen, R. N. Jorissen and M. K. Gilson, *Nucleic Acids Res.*, 2006, **35**, D198–D201.
- 30 L. Wang, Y. Wu, Y. Deng, B. Kim, L. Pierce, G. Krilov, D. Lupyan, S. Robinson, M. K. Dahlgren, J. Greenwood, D. L. Romero, C. Masse, J. L. Knight, T. Steinbrecher, T. Beuming, W. Damm, E. Harder, W. Sherman, M. Brewer, R. Wester, M. Murcko, L. Frye, R. Farid, T. Lin, D. L. Mobley, W. L. Jorgensen, B. J. Berne, R. A. Friesner and R. Abel, *J. Am. Chem. Soc.*, 2015, **137**, 2695–2703.
- 31 D. L. Mobley and M. K. Gilson, *Annu. Rev. Biophys.*, 2017, **46**, 531–558.
- 32 F. J. Rombouts, G. Tresadern, P. Buijnsters, X. Langlois, F. Tovar, T. B. Steinbrecher, G. Vanhoof, M. Somers, J.-I. Andrés and A. A. Trabanco, *ACS Med. Chem. Lett.*, 2015, **6**, 282–286.
- 33 P. Buijnsters, M. De Angelis, X. Langlois, F. J. R. Rombouts, W. Sanderson, G. Tresadern, A. Ritchie, A. A. Trabanco, G. VanHoof, Y. V. Roosbroeck and J.-I. Andrés, *ACS Med. Chem. Lett.*, 2014, **5**, 1049–1053.
- 34 F. J. R. Rombouts, G. Tresadern, O. Delgado, C. Martínez-Lamenca, M. Van Gool, A. García-Molina, S. A. Alonso de Diego, D. Oehlrich, H. Prokopcova, J. M. Alonso, N. Austin, H. Borghys, S. Van Brandt, M. Surkyn, M. De Cleyn, A. Vos, R. Alexander, G. Macdonald, D. Moechars, H. Gijzen and A. A. Trabanco, *J. Med. Chem.*, 2015, **58**, 8216–8235.
- 35 R. Spitzer and A. N. Jain, *J. Comput.-Aided Mol. Des.*, 2012, **26**, 687–699.
- 36 C. Camacho, G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer and T. L. Madden, *BMC Bioinf.*, 2009, **10**, 1.
- 37 G. Landrum, *Online*, <http://www.rdkit.org>, 2006.
- 38 C. Kramer and P. Gedeck, *J. Chem. Inf. Model.*, 2010, **50**, 1961–1969.
- 39 A. Krizhevsky, I. Sutskever and G. E. Hinton, *Advances in Neural Information Processing Systems*, 2012, pp. 1–9.
- 40 O. Vinyals, A. Toshev, S. Bengio and D. Erhan, *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3156–3164.
- 41 A. Radford, L. Metz and S. Chintala, *ICLR*, 2016, pp. 1–16.
- 42 M. Skalic, J. Jiménez, D. Sabbadin and G. De Fabritiis, *J. Chem. Inf. Model.*, 2019, **59**(3), 1205–1214.
- 43 B. Ramsundar, B. Liu, Z. Wu, A. Verras, M. Tudor, R. P. Sheridan and V. Pande, *J. Chem. Inf. Model.*, 2017, **57**, 2068–2076.
- 44 Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing and V. Pande, *Chem. Sci.*, 2018, **9**, 513–530.
- 45 J. Jiménez, S. Doerr, G. Martínez-Rosell, A. S. Rose and G. De Fabritiis, *Bioinformatics*, 2017, **33**, 3036–3042.
- 46 M. Skalic, A. Varela-Rial, J. Jiménez, G. Martínez-Rosell and G. De Fabritiis, *Bioinformatics*, 2018, **35**, 243–250.
- 47 M. Skalic, A. Varela-Rial, J. Jiménez, G. Martínez-Rosell and G. De Fabritiis, *Bioinformatics*, 2019, **35**(2), 243–250.
- 48 C. Wehmeyer and F. Noé, *J. Chem. Phys.*, 2018, **148**, 241703.
- 49 G. Derevyanko, S. Grudin, Y. Bengio and G. Lamoureux, *Bioinformatics*, 2018, **34**, 4046–4053.
- 50 M. H. Segler, M. Preuss and M. P. Waller, arXiv preprint arXiv:1708.04202, 2017.
- 51 R. Gómez-Bombarelli, J. N. Wei, D. Duvenaud, J. M. Hernández-Lobato, B. Sánchez-Lengeling, D. Sheberla, J. Aguilera-Iparraguirre, T. D. Hirzel, R. P. Adams and A. Aspuru-Guzik, *ACS Cent. Sci.*, 2018, **4**, 268–276.
- 52 D. P. Kingma and J. L. Ba, *International Conference on Learning Representations 2015*, 2015, pp. 1–15.
- 53 C. Kramer and P. Gedeck, *J. Chem. Inf. Model.*, 2010, **50**, 1961–1969.
- 54 R. P. Sheridan, *J. Chem. Inf. Model.*, 2013, **53**, 783–790.
- 55 R. A. Friesner, J. L. Banks, R. B. Murphy, T. A. Halgren, J. J. Klicic, D. T. Mainz, M. P. Repasky, E. H. Knoll, M. Shelley, J. K. Perry, D. E. Shaw, P. Francis and P. S. Shenkin, *J. Med. Chem.*, 2004, **47**, 1739–1749.
- 56 S. Genheden and U. Ryde, *Expert Opin. Drug Discovery*, 2015, **10**, 449–461.
- 57 T. Hou, J. Wang, Y. Li and W. Wang, *J. Chem. Inf. Model.*, 2010, **51**, 69–82.
- 58 N. Srinivas, A. Krause, S. M. Kakade and M. Seeger, arXiv preprint arXiv:0912.3995, 2009.
- 59 K. T. Schütt, H. E. Sauceda, P.-J. Kindermans, A. Tkatchenko and K.-R. Müller, *J. Chem. Phys.*, 2018, **148**, 241722.
- 60 J. S. Smith, O. Isayev and A. E. Roitberg, *Chem. Sci.*, 2017, **8**, 3192–3203.
- 61 A. Paszke, G. Chanan, Z. Lin, S. Gross, E. Yang, L. Antiga and Z. Devito, *Advances in Neural Information Processing Systems*, 2017, vol. 30, pp. 1–4.
- 62 S. Doerr, M. J. Harvey, F. Noé and G. De Fabritiis, *J. Chem. Theory Comput.*, 2016, **12**, 1845–1852.

