

Exposing the limitations of molecular machine learning with activity cliffs

Derek van Tilborg,^{1,2} Alisa Alenicheva,³ and Francesca Grisoni^{1,2*}

¹Molecular Machine Learning group, Eindhoven University of Technology, Institute for Complex Molecular Systems and Dept. Biomedical Engineering, Groene Loper 7, 5612AZ Eindhoven, Netherlands.

²Centre for Living Technologies, Alliance TU/e, WUR, UU, UMC Utrecht, Princetonlaan 6, 3584 CB Utrecht, The Netherlands.

³JetBrains Research. Saint Petersburg, Russia.

*Corresponding author: f.grisoni@tue.nl

Abstract

Machine learning has become a crucial tool in drug discovery and chemistry at large, e.g., to predict molecular properties, such as bioactivity, with high accuracy. However, activity cliffs – pairs of molecules that are highly similar in their structure but exhibit large differences in potency – have been underinvestigated for their effect on model performance. Not only are these edge cases informative for molecule discovery and optimization, but models that are well-equipped to accurately predict the potency of activity cliffs have an increased potential for prospective applications. Our work aims to fill the current knowledge gap on best-practice machine learning methods in the presence of activity cliffs. We benchmarked a total of 720 machine and deep learning models on curated bioactivity data from 30 macromolecular targets for their performance on activity cliff compounds. While all methods struggled in the presence of activity cliffs, machine learning approaches based on molecular descriptors outperformed more complex deep learning methods. Our findings highlight large case-by-case differences in performance, advocating for (a) the inclusion of dedicated “activity-cliff-centered” metrics during model development and evaluation, and (b) the development of novel algorithms to better predict the properties of activity cliffs. To this end, the methods, metrics, and results of this study have been encapsulated into an open-access benchmarking platform named MoleculeACE (Activity Cliff Estimation, available on GitHub at: <https://github.com/molML/MoleculeACE>). MoleculeACE is designed to steer the community towards addressing the pressing but overlooked limitation of molecular machine learning models posed by activity cliffs.

Introduction

In the last decade, artificial intelligence (AI) in the form of machine learning has permeated many domains of science. The chemical sciences have particularly benefited from the AI *renaissance*^{1–3}. In multiple applications, machine learning has performed *on par* or even outperformed existing approaches, e.g., for computer-assisted synthesis planning^{4–6}, protein structure prediction^{7,8}, and de novo molecular design^{9–11}. Most AI breakthroughs in chemistry have been driven by *deep learning* – based on neural networks with multiple processing layers^{12–14}. However, there is currently no consensus on whether deep learning models outperform simpler machine learning approaches when it comes to molecular property prediction^{15–17}. The identification of current gaps in machine and deep learning approaches would allow the development of more reliable and widely applicable models to accelerate molecule discovery.

Molecular property prediction has the principle of similarity at its heart¹⁸ – postulating that similar compounds are likely to have similar properties. Notably, one particular exception to this principle holds great insights into the underlying

structure-activity (or structure-property) relationships¹⁹. Such an exception is constituted by *activity cliffs*²⁰ – pairs of structurally similar molecules that exhibit a large difference in their biological activity. Activity cliffs may cause machine learning models to remarkably mispredict the activity of certain molecules, even with an overall high model predictivity. Although generally constituting a source of “disappointment”²⁰, activity cliffs also encode valuable information for many applications¹⁹ (e.g., hit-to-lead optimization^{21,22}, structural alert development²³), since the large change in activity is induced by small structural changes^{24,25}. Activity cliffs are particularly relevant in the context of virtual screening, with the number of highly similar molecules in commonly used commercial libraries varying between 10,000 and 170,000 (Supp. Table S1). While numerous studies have focused on defining activity cliffs^{19,24,26,27}, their detrimental effect on machine learning models has been disproportionately underinvestigated²⁵. Arguably, models that can provide better predictions on activity cliffs are overall better, as they capture the underlying “structure-activity landscape”²⁰ more accurately. Finally, although (macromolecular) structure-based

approaches can aid in identifying discontinuities in the activity landscape²⁸, ligand-based methods are routinely employed “out-of-the-box” for virtual screening, without incorporating considerations on activity cliffs.

Stemming from these considerations, the presented work has a threefold goal: (1) benchmark the performance of several machine- and deep-learning methods on activity cliffs, (2) quantify the effect of activity cliffs on the overall performance of machine learning, and (3) identify promising approaches and future directions in the field of molecular machine learning. To this end, we compared sixteen “classical” machine learning methods – based on human-engineered features (“molecular descriptors”²⁹) – with six deep learning approaches based on molecular strings or graphs, to predict the biological activity of more than 35,000 molecules over 30 macromolecular targets. Our results highlight a generally poor performance of machine learning approaches on activity cliff compounds (particularly evident for deep learning), thereby further underscoring the relevance of assessing structure-activity “discontinuities” during model training and selection.

To further steer the community's efforts toward the underinvestigated topic of activity cliffs, the results of our study were encapsulated in a dedicated benchmarking platform, called MoleculeACE (“Activity Cliff Estimation”). MoleculeACE complements existing benchmarks and datasets for molecular property prediction^{30–33} by providing a novel framework specifically focused on identifying activity cliffs and quantifying the corresponding model performance. MoleculeACE positions itself in a broader movement within the machine learning community^{34–36}, and aims to survey the landscape of existing AI approaches systematically for molecular property prediction³⁷.

Results and discussion

Study set-up

Datasets and activity cliff definition.

To ensure a comprehensive analysis of model performance, we collected and curated data on 30 macromolecular targets from ChEMBL³⁸ v29 (Table 1). To rule out the presence of significant sources of error as much as possible, we curated molecules by considering best practice^{39–41}. In particular, we checked for (a) the presence of duplicates, salts, and mixtures, (b) the consistency of structural annotations (*i.e.*, molecular validity and “sanity”, charge standardization, and stereochemistry definition), and (c) the reliability of the reported experimental values in terms of annotated validity, standard deviation of

multiple entries, and presence of outliers (see Materials and Methods). The curated collection contains a total of 48,707 molecules (of which 35,632 were unique) and mimics typical drug discovery datasets, as it: (a) includes several target families relevant for drug discovery (*e.g.*, kinases, nuclear receptors, G-protein-coupled receptors, transferases, and proteases), and (b) spans different training scenarios, from small (*e.g.*, 615 molecules for Janus Kinase 1 [JAK1]) to large (*e.g.*, 3657 molecules, dopamine D3 receptor [DRD3]) datasets (Table 1).

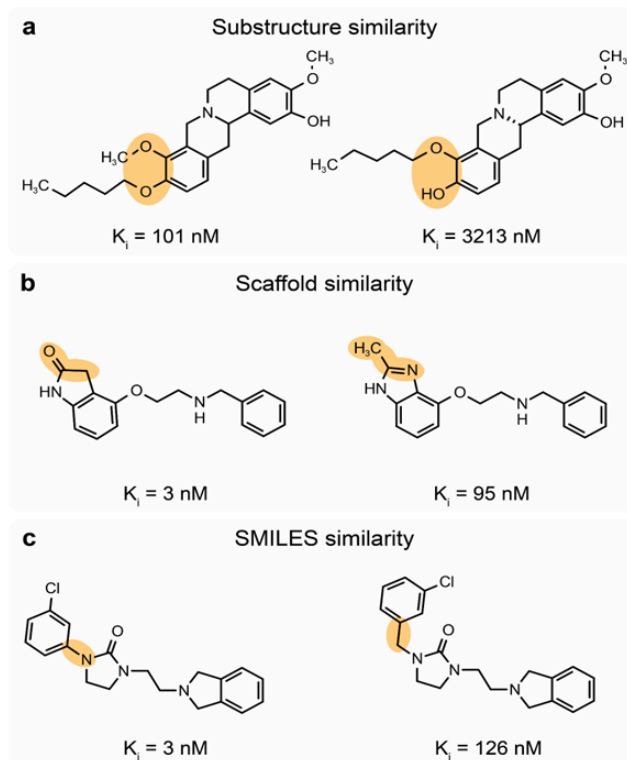


Fig. 1 | Selected examples of activity cliffs (on Dopamine D3 receptor, D3R). **a**, General substructure similarity (Tanimoto coefficient on ECFP). **b**, Scaffold similarity that quantifies the similarity between molecular cores or scaffold decorations (Tanimoto coefficient on scaffold ECFP). **c**, SMILES similarity which detects string insertions, deletions, and translocations (scaled Levenshtein distance).

For each macromolecular target, activity cliffs were identified by considering pairwise structural similarity and differences in potency. We quantified molecular similarity between any pairs of molecules belonging to the same dataset with three distinct approaches:

1. **Substructure similarity.** We computed the Tanimoto coefficient⁴² on extended connectivity fingerprints⁴³ (ECFPs), to capture the presence of shared radial, atom-centered substructures among pairs of molecules. This approach captures “global” differences between molecules by considering the entire set of substructures they contain (Fig. 1a).
2. **Scaffold similarity,** determined by computing ECFPs on atomic scaffolds⁴⁴ and calculating the

respective Tanimoto similarity coefficient. The scaffold similarity allows identifying pairs of compounds that have minor differences in their molecular cores or differ based on their scaffold decoration (Fig. 1b).

3. *Similarity of SMILES strings*, captured by using the Levenshtein distance⁴⁵. This metric detects character insertions, deletions, and translocations (Fig. 1c).

Although there is no widely accepted definition of activity cliffs^{19,46,47}, these three definitions were chosen to cover different types of structural differences relevant to medicinal chemistry. Pairs of molecules that had a computed similarity larger than 90% with at least one of the three methods were considered as ‘highly similar’ in structure. Such pairs of compounds were then checked for their difference in reported potency. In agreement with previous studies²¹, a ten-fold (10x) or larger difference in bioactivity (*i.e.*, on reported K_i or EC_{50} values) was used to identify activity cliff pairs. Compounds that formed at least one activity cliff pair were labeled as “activity cliff compounds”. Although widespread in their usage, we did not consider matched molecular pairs^{47,48}, as they almost doubled the number of cliff compounds compared to our initial approach, while covering 86.6% of cliff compounds identified by our approach. The percentage of activity cliff compounds identified with our approach varied from 7% (JAK1) to 52% (OX2R, Table 1).

Data splitting strategy

The nature of activity cliffs complicates data splitting into training and test sets. Having high structural similarity, but vastly differing bioactivities, makes it infeasible to evenly distribute activity cliff molecules across sets by both their structure and activity. Besides, multiple molecules are often involved in the same activity cliff series: across all datasets, molecules have on average 2.7 ± 0.9 activity cliff ‘partners’ identified by our approach (Supp. Table S3). In this work, we set out to ensure (a) equal representation of the number of activity cliff compounds in the train and test set (to avoid an over/underestimation of their effect on the performance), and (b) preserving structural similarity between training and test molecules, as previously suggested⁴⁹.

To this end, for each dataset, molecules were clustered based on substructure similarity, using spectral clustering⁵⁰ on extended connectivity fingerprints (ECFPs)⁴³. For each cluster, molecules were split into a training (80%) and test set (20%) by stratified random sampling, using their activity cliff label (see Materials and Methods). This method ensured that, even in the case where all activity cliff ‘partners’ end up in the test set ($9.1 \pm 5.3\%$ of activity

cliff molecules on average), highly similar molecules (in terms of substructure [0.80 ± 0.03], scaffold [0.93 ± 0.02], and SMILES [0.95 ± 0.01] similarity) are still present in the training set (Supp. Table S3).

To rule out any potential bias in favor of ECFPs, we set out to compare the similarities of different molecular descriptors in the training and test set, for each macromolecular target (see Materials and Methods). An FDR-adjusted Mann-Whitney-U test ($\alpha = 0.05$) revealed no statistical difference between the distributions of the two sets across all descriptors and all targets. This indicates that the train-test similarity is preserved also when using different molecular descriptors.

“Classical” machine learning strategies

In this work, we considered four “classical” machine learning algorithms that are commonly used for structure-activity relationships prediction (Fig. 2), namely:

1. *K-nearest neighbor* (KNN)⁵¹, a non-parametric approach that uses the k most similar training molecules to predict the response of a new molecule (as the average of the response values).
2. *Random forest* (RF)⁵², based on an ensemble of t distinct decision trees, each trained on various subsamples of the training set (built by bootstrapping). The molecule’s response is predicted as average over the t predictions.
3. *Gradient boosting machine* (GBM)⁵³. Like RF, this algorithm uses multiple decision trees. However, each next decision tree is optimized to predict the prediction residuals of the previous tree.
4. *Support vector regression* (SVM)⁵⁴, which maps data into higher dimensions via a kernel function (a radial basis function in this work) to fit an optimal hyperplane to the training data.

Each algorithm was combined with four types of molecular descriptors²⁹ (Fig. 2), *i.e.*, human-engineered numerical features designed to capture pre-determined chemical information. We explored molecular descriptors with several levels of complexity: (1) extended connectivity fingerprints⁴³ (ECFPs), encoding atom-centered radial substructures⁴³ in the form of a binary array; (2) Molecular ACCess System keys⁵⁵ (MACCs), which encode the presence of predefined substructures in a binary array; (3) 11 physicochemical properties relevant for drug-likeness⁵⁶ (see Materials and Methods), (4) Weighted Holistic Invariant Molecular (WHIM) descriptors⁵⁷, capturing three-dimensional molecular size, shape, symmetry, and atom distribution. Although this selection is not comprehensive (owing to the high number of existing

molecular descriptors²⁹), we consider it a good overview of existing descriptors.

Graph-based deep learning

Molecular graphs are a mathematical representation of molecular topology, with nodes and edges representing atoms and chemical bonds, respectively (Fig. 2b). Neural networks that can learn directly from graphs are becoming increasingly popular for molecular property prediction^{14,58–61}. In this work, we explored four neural network architectures that can directly operate on molecular graphs (Fig. 2d), namely:

1. *Message passing neural network* (MPNN)⁶². For every node in the molecular graph, information (the

'message') from neighboring nodes is aggregated by transforming it with a learnable function.

2. *Graph attention network* (GAT)⁶³. Instead of a message passed across edges, this algorithm also learns attention coefficients that determine the importance of features.

3. *Graph convolutional networks* (GCN),⁶⁴ which aggregate information from neighboring nodes using a fixed convolution.

4. *Attentive fingerprint* (AFP),⁵⁹ which uses attention mechanisms at both the atom and molecule level, allowing it to better capture subtle substructure patterns.

Table 1 | Dataset overview, with response type (inhibition [inhibitory constant, K_i] or agonism [half maximal effective concentration, EC_{50}]), the number of total and test set molecules (n and n_{TEST} , respectively), along with the percentage of total and test activity cliffs (%cliff and %cliff_{test}). An extensive description of the datasets can be found in Supp. Table S2.

Target name	Type	n (n_{TEST})	%cliff (%cliff _{TEST})
Androgen Receptor (AR)	K_i	659 (134)	24 (23)
Cannabinoid receptor 1 (CB1)	EC_{50}	1031 (208)	36 (36)
Coagulation factor X (FX)	K_i	3097 (621)	44 (43)
Delta opioid receptor (DOR)	K_i	2598 (521)	37 (37)
Dopamine D3 receptor (D3R)	K_i	3657 (734)	39 (40)
Dopamine D4 receptor (D4R)	K_i	1859 (374)	38 (38)
Dopamine transporter (DAT)	K_i	1052 (213)	25 (25)
Dual specificity protein kinase CLK4	K_i	731 (149)	9 (9)
Farnesoid X receptor (FXR)	EC_{50}	631 (128)	39 (39)
Ghrelin receptor (GHSR)	EC_{50}	682 (139)	48 (49)
Glucocorticoid receptor (GR)	K_i	750 (152)	31 (31)
Glycogen synthase kinase-3 beta (GSK3)	K_i	856 (173)	18 (18)
Histamine H1 receptor (HRH1)	K_i	973 (197)	23 (23)
Histamine H3 receptor (HRH3)	K_i	2862 (574)	38 (38)
Janus kinase 1 (JAK1)	K_i	615 (126)	7 (8)
Janus kinase 2 (JAK2)	K_i	976 (197)	12 (13)
Kappa opioid receptor (KOR) agonism	EC_{50}	955 (193)	42 (42)
Kappa opioid receptor (KOR) inhibition	K_i	2602 (521)	36 (36)
μ -opioid receptor (MOR)	K_i	3142 (630)	35 (35)
Orexin receptor 2 (OX2R)	K_i	1471 (297)	52 (52)
Peroxisome proliferator-activated receptor alpha (PPAR α)	EC_{50}	1721 (344)	41 (41)
Peroxisome proliferator-activated receptor gamma (PPAR γ)	EC_{50}	2349 (470)	38 (38)
Peroxisome proliferator-activated receptor delta (PPAR δ)	EC_{50}	1125 (225)	42 (42)
PI3-kinase p110-alpha subunit (PIK3CA)	K_i	960 (193)	37 (36)
Serine/threonine-protein kinase PIM1	K_i	1456 (294)	33 (33)
Serotonin 1a receptor (5-HT1A)	K_i	3317 (666)	35 (35)
Serotonin transporter (SERT)	K_i	1704 (342)	35 (35)
Sigma opioid receptor (SOR)	K_i	1328 (267)	35 (35)
Thrombin (F2)	K_i	2754 (553)	36 (36)
Tyrosine-protein kinase ABL1	K_i	794 (161)	32 (32)

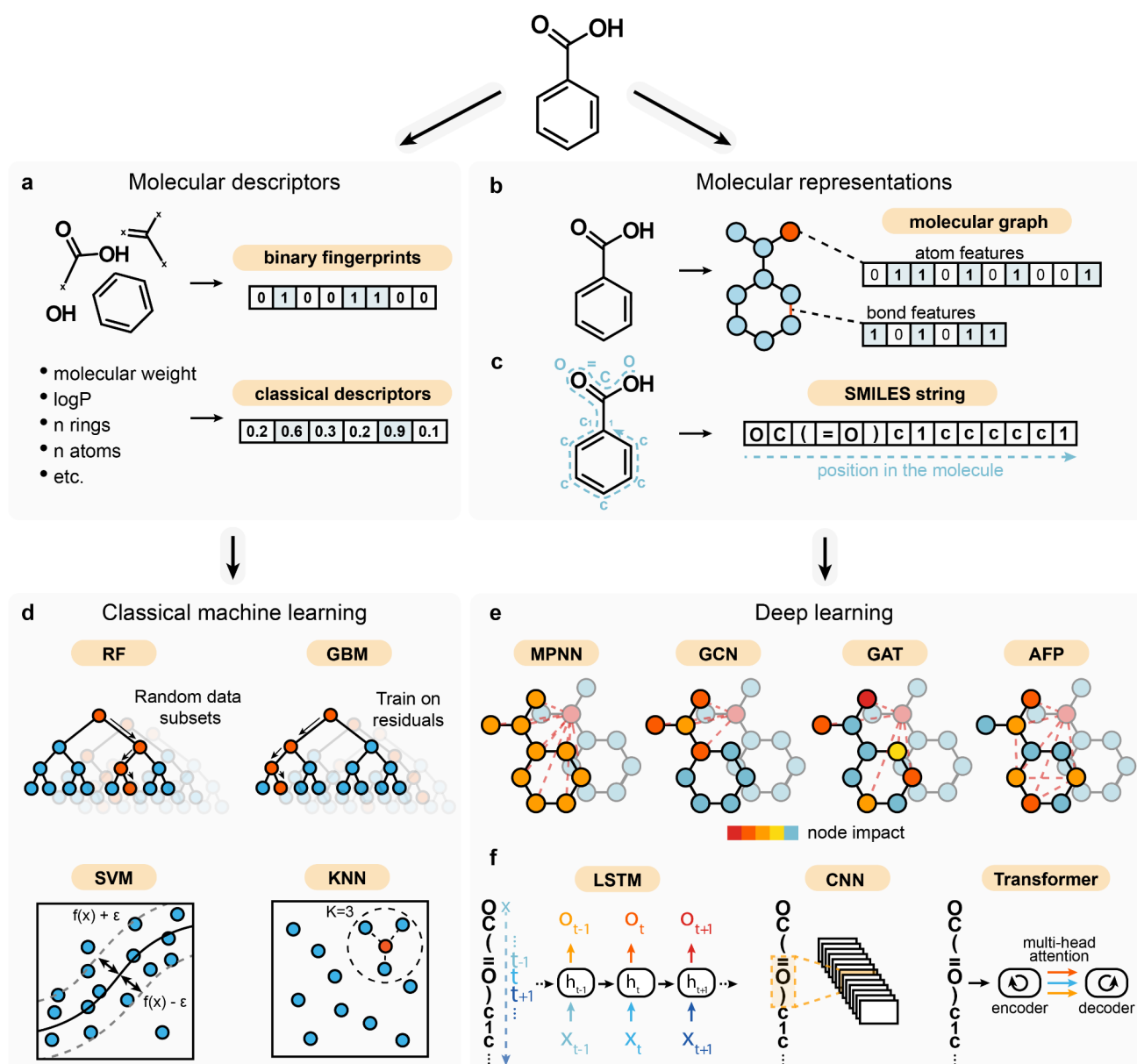


Fig. 2 | Machine learning strategies. **a**, Simplified representation of molecular descriptors, which capture pre-defined molecular features. Both binary fingerprints and classical molecular descriptors are used in this work. **b**, Molecular graph, in which atoms are represented as nodes (with corresponding node features) and bonds are represented as edges (with corresponding edge features, if any). **c**, SMILES strings, which capture two-dimensional information (atom and bond type, and molecular topology) into a string. **d**, Selected “classical” machine learning algorithms that are trained on molecular descriptors: random forest (RF), gradient boosting (GBM), support vector regression (SVM), and k -nearest neighbor (KNN). **e**, Deep learning methods. Four graph neural networks that can learn from molecular graphs were used: message-passing neural network (MPNN), graph convolutional network (GCN), graph attention network (GAT), and attentive fingerprint (AFP). Node colors indicate the impact of other nodes during feature aggregation (indicated by dashed lines). Three SMILES-based methods that can learn from sequential data were used: long short-term memory networks (LSTM), 1D convolutional neural networks (CNN), and transformers.

SMILES-based deep learning methods

As an additional representation, we employed the Simplified Molecular Input Line Entry Systems (SMILES) strings⁶⁵, which have recently become particularly popular for de novo molecular design^{9–11}, and capture two-dimensional molecular information in a textual format (Fig. 2c). Here, we explored three types of neural networks suitable to learn from SMILES strings:

1. *Convolutional neural networks* (CNN)⁶⁶. This neural network architecture uses a learnable convolutional

filter to aggregate information from neighboring positions in a SMILES string with a sliding window approach.

2. *Long short-term memory* (LSTM)⁶⁷ networks. LSTMs – a type of recurrent neural network – can learn from string sequences by keeping track of long-range dependences. As in a previous study⁶⁸, LSTM models were pre-trained on SMILES obtained by merging all training sets with no repetitions (36,281 molecules), using next-character prediction, before applying transfer learning for bioactivity prediction.

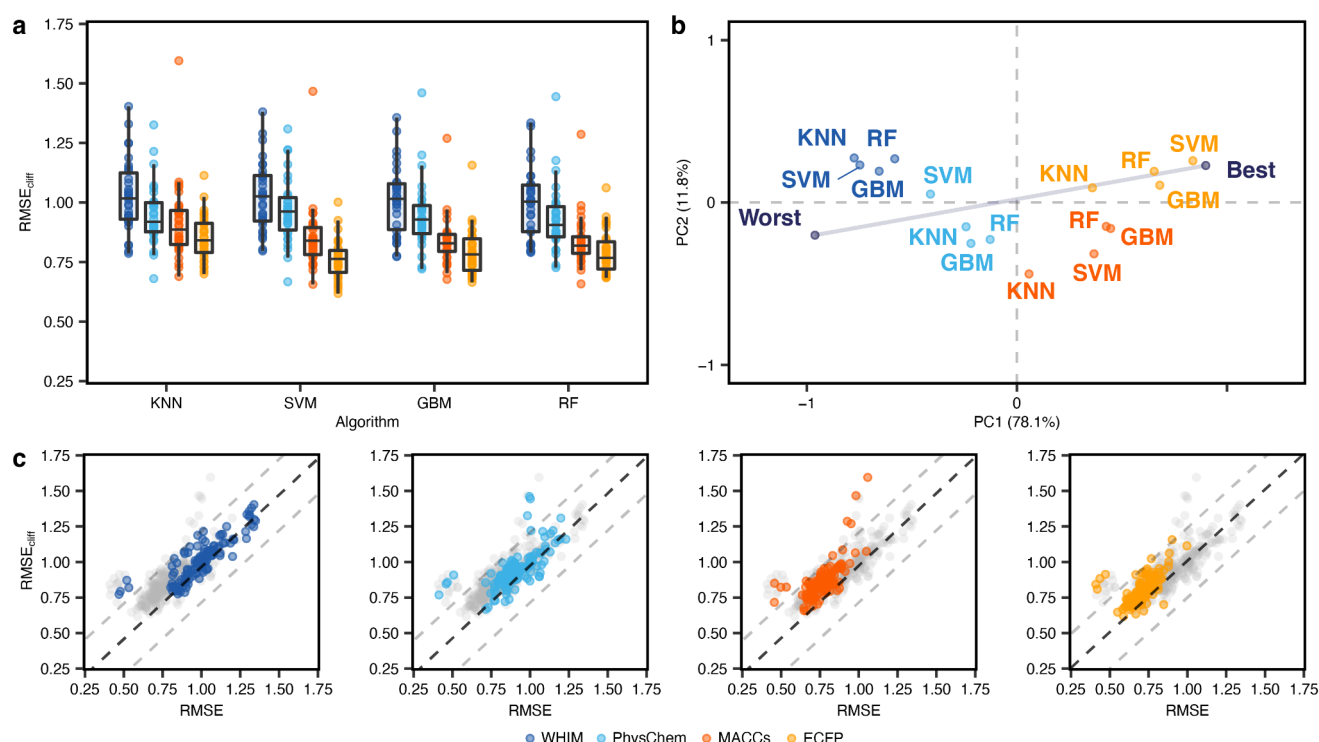


Fig. 3 | Performance of classical machine learning methods. **a**, RMSE on activity cliff compounds using different machine learning algorithms and molecular descriptors (indicated by colors). **b**, Global ranking of all methods using PCA (first two principal components, PC1 and PC2), scaled between best and worst performance. Percentages represent the variance explained by each principal component. **c**, Comparison between the error on activity cliff compounds ($RMSE_{cliff}$) and the error on all compounds (RMSE) for all methods. Black dashed lines indicate $RMSE = RMSE_{cliff}$, while grey dashed lines indicate a difference of ± 0.5 log units between $RMSE_{cliff}$ and RMSE.

3. *Transformer model.* Transformers process the whole sequence at once in a graph-like manner, using positional embedding to capture positional information⁶⁹. Transformers implement the so-called attention⁶⁹, which enables the model to learn which portions of the sequence are more relevant for a given task. The pre-trained ChemBERTa⁷⁰ architecture (10M compounds) was used in combination with transfer learning for bioactivity prediction.

In agreement with previous studies^{71,72} and thanks to the non-univocal character of SMILES strings, we used 10-fold SMILES augmentation to artificially increase the number of training samples for all approaches.

Model performance with activity cliffs

Classical machine learning methods

First, we evaluated the ability of classical machine learning approaches to predict bioactivity (expressed as pEC_{50} or pK_i) in the presence of activity cliffs. The performance was quantified using the root mean squared error on test set molecules (RMSE – the lower, the better; Eq. 1) and activity cliff molecules in the test set ($RMSE_{cliff}$ – the lower, the better; Eq. 2). Overall, large differences in predictive performance on activity cliff compounds can be observed among datasets, with $RMSE_{cliff}$ values ranging from 0.62 to 1.60 log units (Fig. 3a). This effect was also observed

on the overall performance of test set molecules, with RMSE values ranging from 0.41 to 1.35 log units (Supp. Fig. S1a), in line with previous works^{73–75}. Differences in performance relate mostly to the chosen molecular descriptor rather than the machine learning algorithm ($p < 0.05$, Wilcoxon rank-sum test with Benjamini-Hochberg correction, Supp. Fig. S4), with ECFPs yielding the lowest average prediction error on average. Non-binary descriptors (WHIM and physicochemical properties) performed considerably worse overall than binary fingerprints (ECFPs and MACCs), with a higher variation among datasets.

To provide a global assessment of methods across the analyzed datasets, we performed a principal component analysis (PCA) on the obtained $RMSE_{cliff}$ values (Fig. 3b and Supp. Fig. S2a). PCA is a multivariate analysis technique used for data visualization and dimensionality reduction, which linearly combines the original variables into new orthogonal variables (principal components), sorted by the variance they explain. To enhance the interpretability, rows capturing the best and worst $RMSE_{cliff}$ for each dataset were added in order to stretch the PCA results along the direction of the best and worst results^{76,77}. The higher the deviation from the best-worst line, the higher the variability of a method's performance based on the dataset. This analysis confirms the higher impact of molecular

descriptors than the chosen machine learning algorithm on the model performance^{78,79}. SVM coupled with ECFPs resulted in the best method on activity cliffs on average, in agreement with a previous study⁸⁰. However, no statistical difference was found between SVM, GBM, or RF coupled with ECFPs (Wilcoxon rank-sum test, Supp. Fig. S4). In the case of our results, however, the superior performance of ECFPs is somewhat surprising, given that they were used for the definition of activity cliffs (criteria 1 and 2, Fig. 1a), which was expected to introduce an unfavorable bias.

To further investigate the relevance of considering activity cliffs for model assessment, we compared $RMSE_{cliff}$ with the overall error on the test set molecules (Fig. 3c and Supp. Fig. S3). As expected, activity cliff compounds tend to yield higher prediction errors, regardless of the considered approach⁸¹. Although in most of the cases RMSE and $RMSE_{cliff}$ are highly correlated ($\rho = 0.81$ on average), the model performance on activity cliff compounds might be overestimated when considering RMSE alone, up to 0.54 log units. For instance, SVM coupled with ECFP descriptors – resulting in the best performance on average – ranged greatly in its ability to handle activity cliffs. While the mean difference between RMSE and $RMSE_{cliff}$ for this method was only 0.094 log units, large differences were observed in certain data sets (e.g., up to 0.39 log units for the JAK1 receptor). This underscores that strategies with a low overall prediction error might not necessarily be the best ones at handling activity cliffs, thereby hampering their potential for prospective applications.

Deep learning methods

In contrast to classical machine learning algorithms, neural networks allow to bypass human-constructed molecular descriptors and can learn directly from ‘unstructured’ representations of chemical structures. Deep learning approaches trained on either graphs or SMILES strings were compared with (a) a multi-layer perceptron (MLP) based on ECFPs, and (b) the best performing “classical” machine learning method (SVM with ECFP fingerprints), both serving as a reference point (Fig. 4a).

Transfer learning⁸² – applying a models’ previously learned knowledge to a new, related problem by further training – was applied to the LSTM and transformer models in agreement with previous studies^{68,70,83,84}. In a preliminary analysis, we explored transfer learning approaches for graph neural networks using self-supervision (context prediction⁸⁵, infomax⁸⁶, edge prediction⁸⁷, and masking⁸⁵). Since, in line with a recent study⁸⁸, no approach yielded a

notable increase in predictive performance, we did not consider transfer learning further for graph neural networks. When comparing the performance of all tested deep learning methods, we found large differences in predictive performance across datasets – similar to classical machine learning approaches – with $RMSE_{cliff}$ values ranging from 0.68 to 1.44 log units (Fig. 4a). Among the graph-based neural networks, MPNN models resulted in the lowest error on activity cliff compounds on average, although without statistical significance (Wilcoxon rank-sum test with Benjamini-Hochberg correction, Supp. Fig. S4). SMILES-based methods outperformed graph-based methods on average, with LSTM models outperforming all other deep learning methods, including the SMILES-based CNN and transformer models. For CNNs, we did not implement any transfer learning strategy, which could explain their poor(er) performance compared to the other SMILES-based methods. Notably, despite transformers were pretrained on a larger corpus of SMILES strings (10M compounds⁷⁰), they did not perform better than LSTMs, which were pretrained on 36,281 molecules only.

When inspecting the PCA performed on the obtained $RMSE_{cliff}$ values for each target (Fig. 4b), the multi-layer perceptron coupled with ECFPs outperformed all other neural networks based on SMILES or graphs. This is surprising to a certain extent, considering that ECFPs and SMILES are constructed from a molecular graph. This aspect further underscores a current gap in learning efficient features from ‘raw’ molecular representations in the small data regimes typical of drug discovery. Compared to most classical machine learning approach, deep neural networks seem to fall short at picking up subtle structural differences (and the corresponding property change) that give rise to activity cliffs. Similar results were obtained when comparing graph networks for (a) feature attribution with activity cliffs⁸⁹, and (b) bioactivity prediction³⁰. A recent analysis on physicochemical-property cliffs highlights an opposite trend, with deep learning methods performing better than simpler machine learning approaches⁹⁰ – potentially due to the higher number of training samples (approx. 20,000 molecules).

Interestingly, no deep learning method was stable across datasets, as shown by the large deviation from the worst-best line (Fig. 4b and Supp. Fig. 2b). This highlights the need of evaluating the usage of such methods on a case-by-case basis.

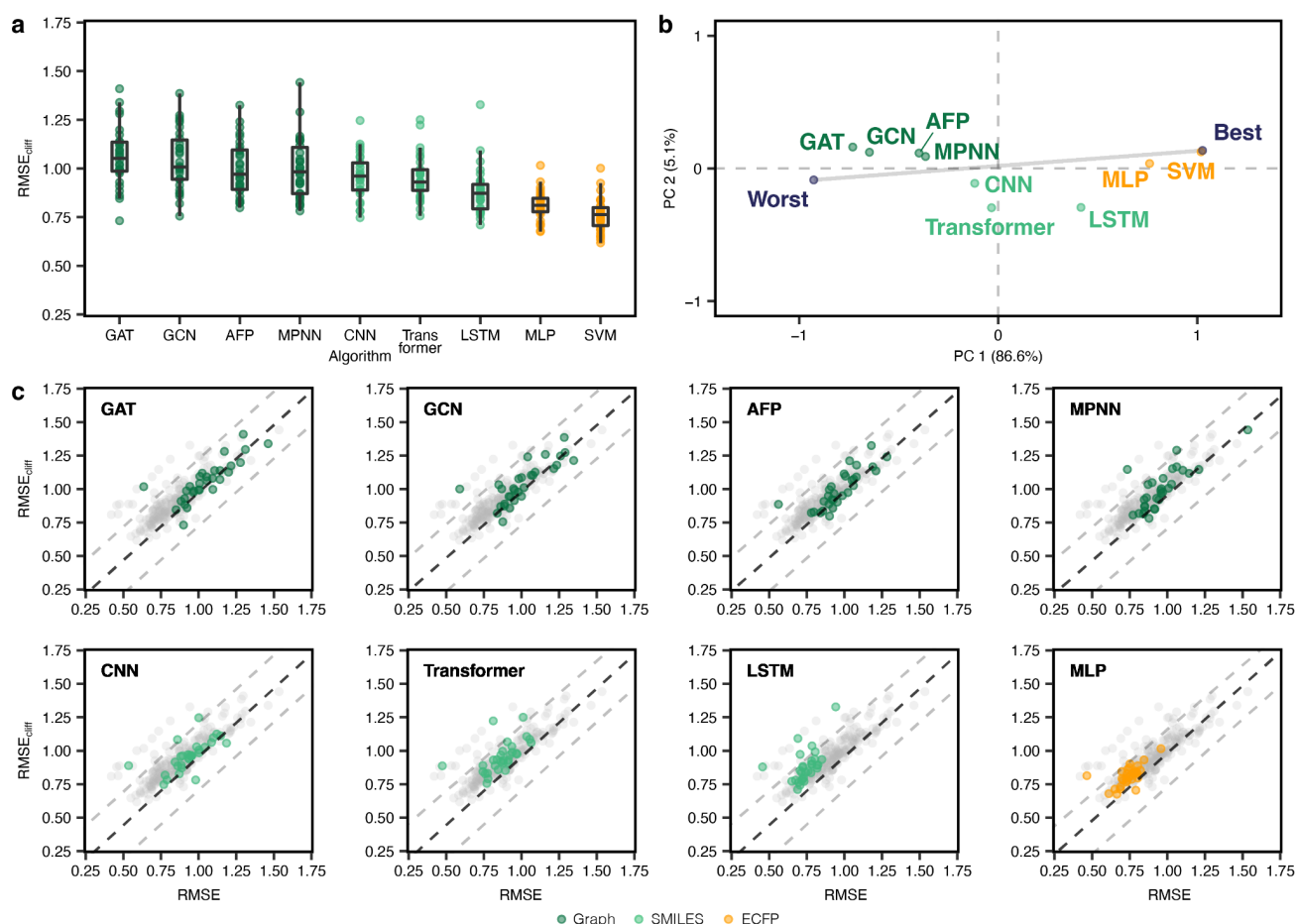


Fig. 4 | Performance of deep learning methods. **a**, RMSE on activity cliff compounds on different deep learning strategies. SVM is reported as a reference. **b**, Global ranking of all methods using PCA (first two principal components, PC1 and PC2), scaled between best and worst performance. Percentages indicate the explained variance by each principal component. **c**, Prediction error on activity cliff compounds (RMSE_{cliff}) compared to all compounds (RMSE) for all methods.

Failure modes of machine learning on activity cliffs

The systematic training and assessment of 720 machine learning models allowed us to investigate the potential “failure modes” of machine learning approaches on activity cliffs. All methods tend to struggle in the presence of activity cliffs (Fig. 3 and Fig. 4). Our first analysis addressed the variation of RMSE_{cliff} across methods and datasets, in search of causes of poor performance. Although small-data regimes are well-known to affect the performance of machine- and deep-learning methods, no correlation was found between the number of molecules in the training set and the prediction error on activity cliffs (Supp. Fig. S5). Furthermore, no relationship between the percentage of activity cliff compounds in the data and model performance was found, except for differences between RMSE and RMSE_{cliff}. This relates to the fact that, the higher the percentage of activity cliffs, the more the RMSE_{cliff} values (computed on a subset of molecules, Eq. 2) will approach RMSE values (Supp. Fig. 6). At the same time, the drug target family did not seem to affect RMSE_{cliff} either (Supp. Fig. S7), further highlighting the difficulties in

forecasting the performance of machine learning on activity cliffs a priori.

We then compared the overall prediction error (RMSE on test set molecules) with the performance on activity cliffs (RMSE_{cliff} on test set molecules). While RMSE and RMSE_{cliff} tend to correlate to a high degree ($p > 0.70$ for 25 datasets out of 30, Fig. 5a), we observed large case-by-case variations. In most cases, the difference between RMSE_{cliff} and RMSE is similar among methods (Fig. 5a and b). This implies that, when choosing a method for its overall error on test set molecules, the performance on activity cliff compounds will be implicitly accounted for. However, for some targets (e.g., CLK4), methods with comparable RMSE scores can exhibit large differences in RMSE_{cliff} scores (Fig. 5c). This indicates that, in these specific cases, choosing a model based on RMSE only might lead to poor prospective performance, e.g., for hit-to-lead optimization or virtual screening in the presence of congeneric compounds (Supp. Table S1). These “islands” of poor performance on activity cliffs were observed across the whole spectrum of machine learning strategies, independently of the reported average performance.

To better elucidate the “drivers of failure” on activity cliffs, we investigated the effect of the training set size on (a) the difference between predictivity on the entire test set and on activity cliffs only ($RMSE_{cliff} - RMSE$) and (b) the correlation between the overall performance ($RMSE$) and the performance on activity cliffs ($RMSE_{cliff}$). The absolute difference between $RMSE$ and $RMSE_{cliff}$ does not correlate with the

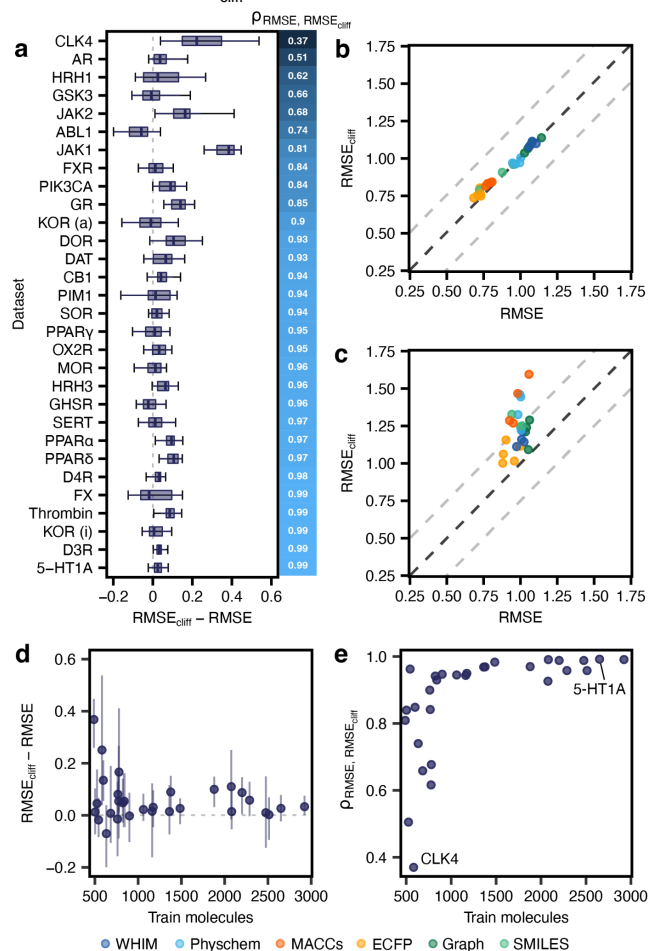


Fig. 5 | Comparing overall model performance and performance on activity cliff compounds. a, Method-wide differences between overall $RMSE$ and $RMSE_{cliff}$ for all targets ordered by Pearson correlation (ρ) between $RMSE$ and $RMSE_{cliff}$. Error bars indicate lowest and highest $RMSE_{cliff}$. b, Comparison between $RMSE$ and $RMSE_{cliff}$ of all methods on 5-HT1A. c, Comparison between $RMSE$ and $RMSE_{cliff}$ of all methods on CLK4. d, Effect of the number of training molecules on the difference between $RMSE$ and $RMSE_{cliff}$. Error bars indicate lowest and highest $RMSE_{cliff}$. e, Relationship between the number of training molecules and the Pearson correlation (ρ) of $RMSE$ and $RMSE_{cliff}$.

number of training molecules ($\rho = -0.15$, Fig 5d). However, the number of training molecules is an important factor in determining the correlations between $RMSE$ and $RMSE_{cliff}$ (Fig. 5e). Datasets containing a sufficient number of training molecules (e.g., larger than 1000) showed a high correlation between $RMSE$ and $RMSE_{cliff}$ ($\rho > 0.80$). In other words, if the number of training molecules increases, the “relative difficulty” of predicting bioactivity on activity cliff molecules decreases. This implies that,

with a sufficient number of training molecules, optimizing $RMSE$ alone will allow to implicitly optimize $RMSE_{cliff}$, too. However, the problem of determining the targets on which $RMSE_{cliff}$ will be suboptimal remains, especially in small data regimes, further underscoring the relevance of implementing activity-cliff-related evaluation approaches.

Bringing it all together: the MoleculeACE benchmark and future applications

Our results and systematic analysis expose current limitations of molecular machine learning, and motivate the use of dedicated metrics and tools for assessing the model performance on activity cliffs, especially in low data regimes. Hence, we collected the modeling and assessment strategies of this study into a dedicated, “activity-cliff-centered” benchmark tool, called *MoleculeACE* (available at URL: <https://github.com/molML/MoleculeACE>).

MoleculeACE integrates standardized data processing for molecular bioactivity data, a comprehensive approach to quantifying activity cliffs, and the tailored performance evaluation strategies presented in this work. Thanks to its modular character, *MoleculeACE* will allow researchers to:

1. *Systematically benchmark a model’s performance* on activity cliffs compounds, in comparison with well-established machine- and deep learning methods.
2. *Evaluate the deck of chosen models on a new dataset* not included in our benchmark, thanks to the data collection and curation pipeline.
3. *Further expand the definition of activity cliffs^{91–93}*, based on specific use cases⁹⁴. It is possible to use custom thresholds for potency differences and structural similarity (e.g., matched molecular pairs, which are already supported) in determining cliff compounds. As this work relies on public bioactivity data, which might be affected by undetectable experimental noise^{81,95} (despite the best data curation efforts), we hope in the future to also see applications of *MoleculeACE* on more homogeneous data, e.g., in terms of used in vitro assay and assay conditions.

We hope that *MoleculeACE*, along with the results of this benchmark study, will incentivize machine learning researchers to consider the crucial topic of activity cliffs in model evaluation and development pipelines. We envision that *MoleculeACE* will serve as a platform for the wider community to develop models that can more accurately capture complex structure-activity landscapes — and ultimately boost the capabilities of machine learning for molecule discovery.

Conclusions and outlook

While machine learning is increasingly often employed for early drug discovery, the topic of activity cliffs has found only limited attention by the scientific community. As shown by our results, not only do machine learning strategies struggle with activity cliffs compared to their overall performance but deep learning methods are particularly challenged by the presence of such compounds. Approaches based on human-engineered molecular descriptors resulted to outperform deep learning based on graphs or SMILES, with no machine learning strategy being consistently better at handling activity cliffs compared to their absolute performance. Our results corroborate previous evidence showing that deep learning methods do not necessarily hold up against simpler machine learning methods (yet) for drug discovery purposes^{15–17}. We envision the development of deep learning strategies that are (a) more efficient in low-data scenarios and (b) better suited to capture structure-activity “discontinuities” to be key for future prospective applications. Structure-based deep learning approaches^{28,96,97} (considering the structure of the macromolecular target in addition to ligand information) might be key to address the issue of activity cliff detection. However, to date, there is no consensus on the benefit of including structural information into machine learning for bioactivity prediction⁹⁸, potentially due to undesirable bias in existing databases^{98–100}.

In the framework of our study design, the model's performance on activity cliff compounds resulted to be highly dataset-dependent, especially for deep learning methods in low-data scenarios. Although the overall prediction error often approximates the performance on activity cliffs, “islands” of poor performance on activity cliffs exist when different strategies are compared on the same data set. These results highlight the importance of evaluating machine learning models for their performance on activity cliffs, especially when prospective applications are envisioned (e.g., virtual screening).

To facilitate such an “activity cliff-centered” model evaluation and development, we developed MoleculeACE. By estimating a model's performance in the presence of activity cliffs alongside regular performance, MoleculeACE has the goal of incentivizing researchers in molecular machine learning to consider the long-standing issue of activity cliffs fully. Models that can accurately predict the effects of subtle structural changes on molecular properties will ultimately give rise to more effective hit-to-lead optimization and the identification of activity cliffs during lead optimization. We envision these

improvements as key to propelling the potential of deep learning in drug discovery and beyond.

Materials and Methods

Data curation

Data collection and preparation. For each macromolecular target, compound bioactivity values were collected from ChEMBL³⁸ v29 via the ‘ChEMBL webresource’ client (*Homo sapiens*). Molecules in the form of canonical SMILES strings were sanitized using RDKit¹⁰¹ v. 2020.09.5¹⁰¹ with default settings and neutralized if charged. Compounds with failed sanitization, annotated in the form of salts, and/or with doubtful data validity (as in the “data_validity_comment” entry of ChEMBL) were removed. For each unique SMILES string, experimental bioactivity data (*i.e.*, K_i or EC_{50} values [nM]) were collected. Dixon's Q test¹⁰² was used to detect the presence of outliers among multiple annotations of a given molecule ($\alpha = 0.05$). The average K_i or EC_{50} value for each molecule was computed and subsequently converted into pEC_{50}/pK_i values (as the negative logarithm of molar concentrations). If the standard deviation of the multiple annotations used to compute the average was above 1 log unit, the corresponding molecule was removed. To rule out errors due to inconsistent annotation of stereochemistry, pairs of compounds having different canonical SMILES but identical ECFPs were removed.

Molecular descriptors calculation. Molecular descriptors were computed from canonicalized SMILES strings using RDkit v. 2020.09.5¹⁰¹. (a) Extended Connectivity Fingerprints (ECFP)⁴³ were computed with a length of 1024 bits and a radius of 2 bonds. (b) MACCS keys⁵⁵ with a length of 166, were computed with default settings. (c) Weighted Holistic Invariant Molecular (WHIM) descriptors⁵⁷ (114 descriptors), were computed on the minimum energy conformers generated with experimental-torsion knowledge distance geometry¹⁰³ and MMFF94¹⁰⁴ force field optimization. (d) “Physico-chemical descriptors” included 11 properties of drug-likeness, *i.e.*, molecular weight, predicted octanol-water partitioning coefficient¹⁰⁵, molar refractivity, topological polar surface area, formal charge, and the number of hydrogen bond donors, hydrogen bond acceptors, rotatable bonds, atoms, rings, and heavy atoms. Real-valued descriptors were standardized by gaussian normalization using the training data mean and standard deviation values.

Detection of activity cliffs. Pairs of structurally similar molecules were detected with three approaches: (a)

substructure similarity, computed via the Tanimoto coefficient on ECFP; (b) *scaffold similarity*, calculated on the ECFP of molecular graph frameworks⁴⁴ (Tanimoto coefficient); (c) (canonical) *SMILES similarity*, computed using the Levenshtein distance^{45,106} (scaled and subsequently converted into '1-distance'). Pairs of compounds having a computed similarity equal to or larger than 0.9 according to at least one of these metrics were checked for the fold-difference in their respective bioactivity (in nM units). Pairs of highly similar compounds showing more than a 10-fold difference in their respective bioactivity were considered activity cliffs.

Train/test splitting. For each target, molecules were clustered by their molecular structure (described as ECFP) into five clusters using spectral clustering⁵⁰ implemented with *sklearn* v. 1.0.2¹⁰⁷ (using a Gaussian kernel and a precomputed affinity matrix of Tanimoto distances). For each cluster, 80% of molecules were assigned to the training data, and 20% to the test data, by stratified splitting (using their belonging to at least one activity cliff pair ['yes'/'no'] as a label).

Descriptor similarity between training and test sets. Similarity among molecular descriptors in the training set was calculated as the mean distance of each molecule in the train set to its five nearest neighbors in the train set. Similarity between molecular descriptors of each molecule in the test set was calculated for the five nearest neighbors in the train set. Graph representations were not considered, as computing graph distances is non-trivial and ECFPs are directly related to molecular graphs. A Mann-Whitney U test, corrected for a false discovery rate of 0.05, was performed using SciPy v. 1.8.1¹⁰⁸.

Molecular graph featurization. For all methods, atom features were encoded as follows. (a) One-hot-encoded properties: atom type, orbital hybridization, atomic vertex degree, aromaticity, and ring membership. (b) Numerically-encoded properties: atomic weight, partial charge (Gasteiger-Marsili¹⁰⁹), number of valence electrons, and number of bound hydrogens. Atomic weight and partial charge were scale-transformed via a sigmoidal function. For MPNN and AFP architectures, bond features were included, i.e., with bond type and conjugation (one-hot encoded).

Model implementation

Hyperparameter optimization. Hyperparameter optimization was performed with Bayesian optimization using a Gaussian process

(method-based specifics are mentioned below). For all models, a maximum of 50 hyperparameter combinations were evaluated, using five-fold cross validation.

Classical machine learning algorithms. KNN, SVM, GBM, and RF regression models were implemented using *sklearn* v. 1.0.2¹⁰⁷. For each approach, the model hyperparameters were optimized as follows: (a) KNN, optimization of the number of neighbors (k), $k = [3, 5, 11, 21]$; (b) SVM, optimization of kernel coefficient (γ) and regularization parameter (C), $\gamma = [1 \times 10^{-6}, 1 \times 10^{-5}, 1 \times 10^{-4}, 1 \times 10^{-3}, 1 \times 10^{-2}, \text{ or } 1 \times 10^{-1}]$ and $C = [1, 10, 100, 1000, 10000]$; (c) GBM, optimization of number of boosting stages (n_b) and maximal model depth (m_d), $n_b = [100, 200, 400]$ and $m_d = [5, 6, 7]$; (d) RF, number of decision trees (t), $t = [100, 250, 500, 1000]$.

Graph neural networks. All regression models were implemented using the PyTorch Geometric package v. 2.0.4¹¹⁰. In MPNN, GCN, and GAT models, global pooling was implemented with a graph multiset transformer¹¹¹ using eight attention heads, followed by a fully connected prediction head. For all models, we optimized the learning rate (lr), $lr = [5 \times 10^{-4}, 5 \times 10^{-5}, \text{ or } 5 \times 10^{-6}]$. The following hyperparameters were optimized:

(a) GCN, hidden atom features (h_a), number of convolutional layers (n_c), hidden multiset transformer nodes (h_t), hidden predictor features (h_p), $h_a = [32, 64, 128, 256, 512]$, $n_c = [1, 2, 3, 4, 5]$, $h_t = [64, 128, 256, 512]$, $h_p = [128, 256, 512]$;

(b) GAT, the hyperparameter search space used for GCN models and the use of GATv1⁶³ or GATv2¹¹² convolutions;

(c) MPNN, hidden atom features (h_a), hidden edge features (h_e), number of message passing steps (s_m), hidden multiset transformer nodes (h_t), hidden predictor features (h_p), $h_a = [32, 64, 128, 256]$, $h_e = [32, 64, 128, 256]$, $s_m = [1, 2, 3, 4, 5]$, $h_t = [64, 128, 256, 512]$, $h_p = [128, 256, 512]$;

(d) AFP, number of attentive layers (n_a), timesteps (n_t), number of hidden predictor features (h_p), $n_a = [1, 2, 3, 4, 5]$, $n_t = [1, 2, 3, 4, 5]$, $h_p = [32, 64, 128, 256]$.

All models were trained for 300 epochs, using early-stopping with a patience of 10 epochs.

Feed-forward neural network. A multi-layer perceptron was implemented using Pytorch v. 1.11.0¹¹³. It was optimized for: (a) learning rate ($lr = [5 \times 10^{-4}, 5 \times 10^{-5}, 64, 5 \times 10^{-6}]$), (b) number of hidden features ($n_h = [256, 512, 1024]$) and (c) number of layers ($n_l = [1, 2, 3, 4, 5]$). Models were trained for 500 epochs using early-stopping with a patience of 10 epochs.

SMILES-based models. SMILES strings were encoded as one-hot vectors. SMILES strings longer than 200 characters were truncated (0.71% on average). 10-fold data augmentation was applied to all SMILES-based methods using a maximum of 9 extra non-canonical SMILES strings for every SMILES string in the data set. Non-canonical SMILES strings were generated using RDKit¹⁰¹.

(a) LSTM models were pre-trained on SMILES obtained by merging all training sets with no repetitions (36,281 molecules), using next-character prediction as in a recent study⁶⁸. The network was composed of four layers comprising 5,820,515 parameters (layer 1: batch normalization; layer 2: LSTM with 1024 units; layer 3: LSTM with 256 units; layer 4: batch normalization). We used the Adam optimizer with a learning rate of 10^{-4} for 100 epochs. Regression models were then obtained by transfer learning (with weight freezing for layer no. 2) for 100 epochs with a regression head.

(b) 1D CNNs were adapted from a recent study⁶⁶. We used a single 1D convolutional layer with a stepsize equal to 1, followed by a fully connected layer, with training for 500 epochs. It was optimized for the learning rate (lr), number of hidden features in the fully connected layer (n_h), and convolution kernel size (n_k), $lr = [5 \times 10^{-4}, 5 \times 10^{-5}, 64, 5 \times 10^{-6}]$, $n_h = [128, 256, 512, 1024]$, $n_k = [4, 8, 10]$.

(c) Transformer models and the corresponding SMILES tokenization were based on the ChemBERTa⁷⁰ architecture. We used the pre-trained ChemBERTa model weights based on 10M compounds from PubChem¹¹⁴. We fine-tuned the model by freezing its weights and replacing the final pooling layer by a regression head with 1 fully connected layer, and trained for 100 epochs. We used the Adam optimizer with a learning rate of 5×10^{-4} . For all methods, we used early-stopping with a patience of 10 epochs.

Performance evaluation. The overall model performance was quantified via the root mean square error (RMSE) computed on the bioactivity values (*i.e.*, pK_i or pEC_{50}), as follows (Eq. 1):

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}} \quad (1)$$

where \hat{y}_i is the predicted bioactivity of the i -th compound and y_i is the corresponding experimental value, while n represents the number of considered molecules.

The performance on activity cliffs compounds was quantified by computing the root mean square error ($RMSE_{cliff}$) only on compounds that belonged to at least one activity cliff pair, as follows (Eq.2):

$$RMSE_{cliff} = \sqrt{\frac{\sum_{j=1}^{n_c} (\hat{y}_j - y_j)^2}{n_c}} \quad (2)$$

where \hat{y}_j is the predicted bioactivity of the j -th activity cliff compound, y_j is the corresponding experimental value, and n_c represents the total number of activity cliff compounds considered. R^2 and Q^2 metrics were not considered to avoid the introduction of undesired biases related to the different range of the training/test set responses across datasets^{115,116}.

Author contributions

Conceptualization: FG, DvT. **Data curation:** DvT, FG. **Formal analysis:** DvT, AA. **Methodology:** DvT, AA, FG. **Software:** DvT, AA. **Writing - original draft:** DvT. **Writing - review & editing:** all authors. All authors have given approval to the final version of the manuscript.

Abbreviations

AFP: Attentive fingerprint
 CNN: Convolutional neural network
 ECFP: Extended connectivity fingerprints
 GAT: Graph attention network
 GBM: Gradient boosting machine
 GCN: Graph convolutional network
 KNN: K-nearest neighbor
 LSTM: Long short-term memory network
 MACCs: Molecular ACCess system
 MLP: Multilayer perceptron
 MPNN: Message passing neural network
 RF: Random forest
 RMSE: Root mean square error
 SMILES: Simplified molecular input line entry system
 SVM: Support vector machine
 WHIM: Weighted holistic invariant molecular descriptors

Acknowledgments

The authors thank the group of Prof. Luc Brunsveld for valuable discussion and feedback, Dr. Jiménez-Luna and Luke Rossen for comments on the code, and Rıza Özçelik for feedback on the manuscript. FG acknowledges support from the Irène Curie Fellowship and from the Centre for Living Technologies.

Conflict of interest

None to declare.

References

- Almeida, A. F. de, Moreira, R. & Rodrigues, T. Synthetic organic chemistry driven by artificial intelligence. *Nature Reviews Chemistry* vol. 3 589–604

- (2019).
- Baskin, I. I., Winkler, D. & Tetko, I. V. A renaissance of neural networks in drug discovery. *Expert Opin. Drug Discov.* **11**, 785–795 (2016).
 - Chen, H., Engkvist, O., Wang, Y., Olivecrona, M. & Blaschke, T. The rise of deep learning in drug discovery. *Drug Discov. Today* **23**, 1241–1250 (2018).
 - Segler, M. H. S., Preuss, M. & Waller, M. P. Planning chemical syntheses with deep neural networks and symbolic AI. *Nature* **555**, 604–610 (2018).
 - Schwaller, P. *et al.* Molecular Transformer: A Model for Uncertainty-Calibrated Chemical Reaction Prediction. *ACS Central Science* vol. 5 1572–1583 (2019).
 - Coley, C. W., Green, W. H. & Jensen, K. F. Machine Learning in Computer-Aided Synthesis Planning. *Acc. Chem. Res.* **51**, 1281–1289 (2018).
 - Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
 - Baek, M. *et al.* Accurate prediction of protein structures and interactions using a three-track neural network. *Science* **373**, 871–876 (2021).
 - Segler, M. H. S., Kogej, T., Tyrchan, C. & Waller, M. P. Generating Focused Molecule Libraries for Drug Discovery with Recurrent Neural Networks. *ACS Central Science* vol. 4 120–131 (2018).
 - Merk, D., Friedrich, L., Grisoni, F. & Schneider, G. De Novo Design of Bioactive Small Molecules by Artificial Intelligence. *Mol. Inform.* **37**, (2018).
 - Yuan, W. *et al.* Chemical Space Mimicry for Drug Discovery. *J. Chem. Inf. Model.* **57**, 875–882 (2017).
 - Schmidhuber, J. Deep learning in neural networks: An overview. *Neural Networks* vol. 61 85–117 (2015).
 - LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* vol. 521 436–444 (2015).
 - Atz, K., Grisoni, F. & Schneider, G. Geometric deep learning on molecular representations. *Nature Machine Intelligence* **3**, 1023–1032 (2021).
 - Jiang, D. *et al.* Could graph neural networks learn better molecular representation for drug discovery? A comparison study of descriptor-based and graph-based models. *J. Cheminform.* **13**, 12 (2021).
 - Yang, K. *et al.* Analyzing Learned Molecular Representations for Property Prediction. *J. Chem. Inf. Model.* **59**, 3370–3388 (2019).
 - Valsecchi, C. *et al.* Predicting molecular activity on nuclear receptors by multitask neural networks. *J. Chemom.* (2020).
 - Johnson, M. A. & Maggiora, G. M. *Concepts and applications of molecular similarity*. (Wiley, 1990).
 - Stumpfe, D., Hu, H. & Bajorath, J. Advances in exploring activity cliffs. *J. Comput. Aided Mol. Des.* **34**, 929–942 (2020).
 - Maggiora, G. M. On outliers and activity cliffs—why QSAR often disappoints. *J. Chem. Inf. Model.* **46**, 1535 (2006).
 - Stumpfe, D. & Bajorath, J. Exploring activity cliffs in medicinal chemistry. *J. Med. Chem.* **55**, 2932–2942 (2012).
 - Dimova, D., Heikamp, K., Stumpfe, D. & Bajorath, J. Do medicinal chemists learn from activity cliffs? A systematic evaluation of cliff progression in evolving compound data sets. *J. Med. Chem.* **56**, 3339–3345 (2013).
 - Wedlake, A. J. *et al.* Structural Alerts and Random Forest Models in a Consensus Approach for Receptor Binding Molecular Initiating Events. *Chem. Res. Toxicol.* **33**, 388–401 (2020).
 - Hu, Y. & Bajorath, J. Extending the activity cliff concept: structural categorization of activity cliffs and systematic identification of different types of cliffs in the ChEMBL database. *J. Chem. Inf. Model.* **52**, 1806–1811 (2012).
 - Cruz-Monteagudo, M. *et al.* Activity cliffs in drug discovery: Dr Jekyll or Mr Hyde? *Drug Discov. Today* **19**, 1069–1080 (2014).
 - Stumpfe, D., Hu, Y., Dimova, D. & Bajorath, J. Recent progress in understanding activity cliffs and their utility in medicinal chemistry. *J. Med. Chem.* **57**, 18–28 (2014).
 - Bajorath, J. *et al.* Navigating structure–activity landscapes. *Drug Discov. Today* **14**, 698–705 (2009).
 - Husby, J., Bottegoni, G., Kufareva, I., Abagyan, R. & Cavalli, A. Structure-based predictions of activity cliffs. *J. Chem. Inf. Model.* **55**, 1062–1076 (2015).
 - Consonni, V. & Todeschini, R. *Molecular Descriptors for Chemoinformatics: Volume I: Alphabetical Listing / Volume II: Appendices, References*. (John Wiley & Sons, 2009).
 - Wu, Z. *et al.* MoleculeNet: a benchmark for molecular machine learning. *Chem. Sci.* **9**, 513–530 (2018).
 - Feinberg, E. N. *et al.* PotentialNet for Molecular Property Prediction. *ACS Cent Sci* **4**, 1520–1530 (2018).
 - Hu, W. *et al.* Open graph benchmark: Datasets for machine learning on graphs. *Adv. Neural Inf. Process. Syst.* **33**, 22118–22133 (2020).
 - Stanley, M. *et al.* FS-Mol: A Few-Shot Learning Dataset of Molecules. (2021).
 - Deng, J. *et al.* ImageNet: A large-scale hierarchical image database. in *2009 IEEE Conference on Computer Vision and Pattern Recognition* 248–255 (2009).
 - Wang, A. *et al.* GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. *arXiv [cs.CL]* (2018).
 - Brown, N., Fiscato, M., Segler, M. H. S. & Vaucher, A. C. GuacaMol: Benchmarking Models for de Novo Molecular Design. *J. Chem. Inf. Model.* **59**, 1096–1108 (2019).
 - Raji, I. D., Bender, E. M., Paullada, A., Denton, E. & Hanna, A. AI and the Everything in the Whole Wide World Benchmark. *arXiv [cs.LG]* (2021).
 - Gaulton, A. *et al.* ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic Acids Res.* **40**, D1100–7 (2012).
 - Tiikkainen, P., Bellis, L., Light, Y. & Franke, L. Estimating error rates in bioactivity databases. *J. Chem. Inf. Model.* **53**, 2499–2505 (2013).
 - Mansouri, K., Grulke, C. M., Richard, A. M., Judson, R. S. & Williams, A. J. An automated curation procedure for addressing chemical errors and inconsistencies in public datasets used in QSAR modelling. *SAR QSAR Environ. Res.* **27**, 939–965 (2016).
 - Fourches, D., Muratov, E. & Tropsha, A. Trust, but verify: on the importance of chemical structure curation in cheminformatics and QSAR modeling research. *J. Chem. Inf. Model.* **50**, 1189–1204 (2010).
 - Cereto-Massagué, A. *et al.* Molecular fingerprint similarity search in virtual screening. *Methods* **71**, 58–63 (2015).
 - Rogers, D. & Hahn, M. Extended-connectivity fingerprints. *J. Chem. Inf. Model.* **50**, 742–754 (2010).
 - Bemis, G. W. & Murcko, M. A. The Properties of Known Drugs. 1. Molecular Frameworks. *Journal of Medicinal Chemistry* vol. 39 2887–2893 (1996).
 - Yujian, L. & Bo, L. A normalized Levenshtein distance metric. *IEEE Trans. Pattern Anal. Mach. Intell.* **29**, 1091–1095 (2007).
 - Hussain, J. & Rea, C. Computationally efficient algorithm to identify matched molecular pairs (MMPs)

- in large data sets. *J. Chem. Inf. Model.* **50**, 339–348 (2010).
47. Hu, X., Hu, Y., Vogt, M., Stumpfe, D. & Bajorath, J. MMP-Cliffs: systematic identification of activity cliffs on the basis of matched molecular pairs. *J. Chem. Inf. Model.* **52**, 1138–1145 (2012).
 48. Bajorath, J. Representation and identification of activity cliffs. *Expert Opin. Drug Discov.* **12**, 879–883 (2017).
 49. Puzyn, T., Mostrag-Szlichtyng, A., Gajewicz, A., Skrzyński, M. & Worth, A. P. Investigating the influence of data splitting on the predictive ability of QSAR/QSPR models. *Struct. Chem.* **22**, 795–804 (2011).
 50. Stella & Shi. Multiclass spectral clustering. *Proc. IEEE Int. Conf. Comput. Vis.* (2008).
 51. Fix, E. & Hodges, J. L. Discriminatory Analysis. Nonparametric Discrimination: Consistency Properties. *Int. Stat. Rev.* **57**, 238–247 (1989).
 52. Breiman, L. Bagging predictors. *Mach. Learn.* **24**, 123–140 (1996).
 53. Friedman, J. H. Greedy Function Approximation: A Gradient Boosting Machine. *Ann. Stat.* **29**, 1189–1232 (2001).
 54. Cristianini, N. & Shawe-Taylor, J. *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. (Cambridge University Press, 2000).
 55. Durant, J. L., Leland, B. A., Henry, D. R. & Nourse, J. G. Reoptimization of MDL keys for use in drug discovery. *J. Chem. Inf. Comput. Sci.* **42**, 1273–1280 (2002).
 56. Walters, W. P. & Murcko, M. A. Prediction of 'drug-likeness'. *Adv. Drug Deliv. Rev.* **54**, 255–271 (2002).
 57. Todeschini, R. & Gramatica, P. New 3D molecular descriptors: The WHIM theory and QSAR applications. in *3D QSAR in Drug Design* 355–380 (Kluwer Academic Publishers, 2005).
 58. Stokes, J. M. *et al.* A Deep Learning Approach to Antibiotic Discovery. *Cell* **181**, 475–483 (2020).
 59. Xiong, Z. *et al.* Pushing the Boundaries of Molecular Representation for Drug Discovery with the Graph Attention Mechanism. *J. Med. Chem.* **63**, 8749–8760 (2020).
 60. Chen, C., Ye, W., Zuo, Y., Zheng, C. & Ong, S. P. Graph Networks as a Universal Machine Learning Framework for Molecules and Crystals. *Chem. Mater.* **31**, 3564–3572 (2019).
 61. Kearnes, S., McCloskey, K., Berndl, M., Pande, V. & Riley, P. Molecular graph convolutions: moving beyond fingerprints. *J. Comput. Aided Mol. Des.* **30**, 595–608 (2016).
 62. Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O. & Dahl, G. E. Neural Message Passing for Quantum Chemistry. in *Proceedings of the 34th International Conference on Machine Learning* vol. 70 1263–1272 (PMLR, 2017).
 63. Velickovic, P. *et al.* Graph attention networks. *Stat* **1050**, 20 (2017).
 64. Kipf, T. N. & Welling, M. Semi-Supervised Classification with Graph Convolutional Networks. *arXiv [cs.LG]* (2016).
 65. Weininger, D. SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *J. Chem. Inf. Comput. Sci.* **28**, 31–36 (1988).
 66. Kimber, T. B., Gagnebin, M. & Volkamer, A. Maxsmi: Maximizing molecular property prediction performance with confidence estimation using SMILES augmentation and deep learning. *Artificial Intelligence in the Life Sciences* **1**, 100014 (2021).
 67. Hochreiter, S. & Schmidhuber, J. Long Short-Term Memory. *Neural Computation* vol. 9 1735–1780 (1997).
 68. Moret, M., Grisoni, F., Katzberger, P. & Schneider, G. Perplexity-based molecule ranking and bias estimation of chemical language models. *ChemRxiv* (2021).
 69. Vaswani, Shazeer & Parmar. Attention is all you need. *Adv. Neural Inf. Process. Syst.*
 70. Chithrananda, S., Grand, G. & Ramsundar, B. ChemBERTa: Large-Scale Self-Supervised Pretraining for Molecular Property Prediction. *arXiv [cs.LG]* (2020).
 71. Arús-Pous, J. *et al.* Randomized SMILES strings improve the quality of molecular generative models. *J. Cheminform.* **11**, 71 (2019).
 72. Bjerrum, E. J. SMILES Enumeration as Data Augmentation for Neural Network Modeling of Molecules. *arXiv [cs.LG]* (2017).
 73. Sheridan, R. P. Three useful dimensions for domain applicability in QSAR models using random forest. *J. Chem. Inf. Model.* **52**, 814–823 (2012).
 74. Guha, R., Dutta, D., Jurs, P. C. & Chen, T. Local lazy regression: making use of the neighborhood to improve QSAR predictions. *J. Chem. Inf. Model.* **46**, 1836–1847 (2006).
 75. Subramanian, G., Ramsundar, B., Pande, V. & Denny, R. A. Computational Modeling of β -Secretase 1 (BACE-1) Inhibitors Using Ligand Based Approaches. *J. Chem. Inf. Model.* **56**, 1936–1949 (2016).
 76. Todeschini, R., Ballabio, D., Cassotti, M. & Consonni, V. N3 and BNN: Two New Similarity Based Classification Methods in Comparison with Other Classifiers. *J. Chem. Inf. Model.* **55**, 2365–2374 (2015).
 77. Grisoni, F., Merk, D., Byrne, R. & Schneider, G. Scaffold-Hopping from Synthetic Drugs by Holistic Molecular Representation. *Sci. Rep.* **8**, 16469 (2018).
 78. Tetko, I. V. *et al.* Critical Assessment of QSAR Models of Environmental Toxicity against *Tetrahymena pyriformis*: Focusing on Applicability Domain and Overfitting by Variable Selection. *Journal of Chemical Information and Modeling* vol. 48 1733–1746 (2008).
 79. Zhu, H. *et al.* Combinatorial QSAR modeling of chemical toxicants tested against *Tetrahymena pyriformis*. *J. Chem. Inf. Model.* **48**, 766–784 (2008).
 80. de la Vega de León, A. & Bajorath, J. Prediction of compound potency changes in matched molecular pairs using support vector regression. *J. Chem. Inf. Model.* **54**, 2654–2663 (2014).
 81. Sheridan, R. P. *et al.* Experimental Error, Kurtosis, Activity Cliffs, and Methodology: What Limits the Predictivity of Quantitative Structure-Activity Relationship Models? *J. Chem. Inf. Model.* **60**, 1969–1982 (2020).
 82. Cai, C. *et al.* Transfer Learning for Drug Discovery. *J. Med. Chem.* **63**, 8683–8694 (2020).
 83. Gupta, A. *et al.* Generative Recurrent Networks for De Novo Drug Design. *Mol. Inform.* **37**, (2018).
 84. Awale, M., Sirockin, F., Stiefl, N. & Reymond, J.-L. Drug Analogs from Fragment-Based Long Short-Term Memory Generative Neural Networks. *J. Chem. Inf. Model.* **59**, 1347–1356 (2019).
 85. Hu, W. *et al.* Strategies for pre-training Graph Neural Networks. (2019).
 86. Veličković, P. *et al.* Deep Graph Infomax. (2018).
 87. Hamilton, W. L., Ying, R. & Leskovec, J. Inductive Representation Learning on Large Graphs. (2017).
 88. Wang, H. *et al.* Evaluating Self-Supervised Learning for Molecular Graph Embeddings. *arXiv [cs.LG]* (2022).
 89. Jiménez-Luna, J., Skalic, M. & Weskamp, N. Benchmarking Molecular Feature Attribution Methods

- with Activity Cliffs. *Journal of Chemical Information and Modeling* vol. 62 274–283 (2022).
90. Kwapien, K. *et al.* Implications of additivity and nonadditivity for machine learning and deep learning models in drug design. *Research Square* (2022).
 91. Stumpfe, D., Hu, H. & Bajorath, J. Introducing a new category of activity cliffs with chemical modifications at multiple sites and rationalizing contributions of individual substitutions. *Bioorg. Med. Chem.* **27**, 3605–3612 (2019).
 92. Hu, H. & Bajorath, J. Introducing a new category of activity cliffs combining different compound similarity criteria. *RSC Med Chem* **11**, 132–141 (2020).
 93. Guha, R. & Van Drie, J. H. Structure–activity landscape index: identifying and quantifying activity cliffs. *Journal of chemical information and modeling* vol. 48 646–658 (2008).
 94. Stumpfe, D., Hu, H. & Bajorath, J. Advances in exploring activity cliffs. *J. Comput. Aided Mol. Des.* **34**, 929–942 (2020).
 95. Gogishvili, D., Nittinger, E., Margreitter, C. & Tyrchan, C. Nonadditivity in public and inhouse data: implications for drug design. *J. Cheminform.* **13**, 47 (2021).
 96. Wallach, I., Dzamba, M. & Heifets, A. AtomNet: A Deep Convolutional Neural Network for Bioactivity Prediction in Structure-based Drug Discovery. *arXiv [cs.LG]* (2015).
 97. Gentile, F. *et al.* Deep Docking: A Deep Learning Platform for Augmentation of Structure Based Drug Discovery. *ACS Cent Sci* **6**, 939–949 (2020).
 98. Volkov, M. *et al.* On the Frustration to Predict Binding Affinities from Protein–Ligand Structures with Deep Neural Networks. *J. Med. Chem.* **65**, 7946–7958 (2022).
 99. 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.* **59**, 947–961 (2019).
 100. Chen, L. *et al.* Hidden bias in the DUD-E dataset leads to misleading performance of deep learning in structure-based virtual screening. *PLoS One* **14**, e0220113 (2019).
 101. RDkit. RDKit: Open-source cheminformatics; <http://www.rdkit.org>.
 102. Rorabacher, D. B. Statistical treatment for rejection of deviant values: critical values of Dixon's 'Q' parameter and related subrange ratios at the 95% confidence level. *Analytical Chemistry* vol. 63 139–146 (1991).
 103. Riniker, S. & Landrum, G. A. Better Informed Distance Geometry: Using What We Know To Improve Conformation Generation. *J. Chem. Inf. Model.* **55**, 2562–2574 (2015).
 104. Halgren, T. A. Merck molecular force field. I. Basis, form, scope, parameterization, and performance of MMFF94. *Computational Chemistry* **17**, 490–519 (1996).
 105. Wildman, S. A. & Crippen, G. M. Prediction of Physicochemical Parameters by Atomic Contributions. *J. Chem. Inf. Comput. Sci.* **39**, 868–873 (1999).
 106. python-Levenshtein. <https://pypi.org/project/python-Levenshtein/>.
 107. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *the Journal of machine Learning research* **12**, 2825–2830 (2011).
 108. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* **17**, 261–272 (2020).
 109. Gasteiger, J. & Marsili, M. Iterative partial equalization of orbital electronegativity—a rapid access to atomic charges. *Tetrahedron* **36**, 3219–3228 (1980).
 110. Fey, M. & Lenssen, J. E. Fast Graph Representation Learning with PyTorch Geometric. *arXiv [cs.LG]* (2019).
 111. Baek, J., Kang, M. & Hwang, S. J. Accurate Learning of Graph Representations with Graph Multiset Pooling. *arXiv [cs.LG]* (2021).
 112. Brody, S., Alon, U. & Yahav, E. How Attentive are Graph Attention Networks? *arXiv [cs.LG]* (2021).
 113. Paszke, A. *et al.* PyTorch: An imperative style, high-performance deep learning library. *Adv. Neural Inf. Process. Syst.* **32**, (2019).
 114. Kim, S. *et al.* PubChem 2019 update: improved access to chemical data. *Nucleic Acids Res.* **47**, D1102–D1109 (2019).
 115. Alexander, D. L. J., Tropsha, A. & Winkler, D. A. Beware of R2: Simple, Unambiguous Assessment of the Prediction Accuracy of QSAR and QSPR Models. *J. Chem. Inf. Model.* **55**, 1316–1322 (2015).
 116. Consonni, V., Todeschini, R., Ballabio, D. & Grisoni, F. On the Misleading Use of Q2 F3 for QSAR Model Comparison. *Mol. Inform.* **38**, e1800029 (2019).

Table S1 | Presence of highly similar compounds in commercially available libraries. For each library, the number of molecules having *at least* one highly similar neighbor (determined as having a Tanimoto similarity on ECFPs larger than 90%) was reported, along with the corresponding percentage. ECFPs (1024 bits, radius = 2) were computed with RDKit on canonical SMILES strings, after filtering out duplicates and salts, within KNIME 4.3.3.

Provider and library	no. molecules	no. similar	perc. similar
<i>Asinex</i> ^a	572,393	76,794	13.42%
<i>Specs</i> (10, 20, 50 mg) ^b	199,965	13,369	6.69%
<i>Enamine Advanced</i> ^c	604,507	173,383	28.77%
<i>Enamine Premium</i> ^c	40,694	18,936	46.53%

^a"All screening compounds", downloaded from [https://www.asinex.com/screening-libraries-\(all-libraries\)](https://www.asinex.com/screening-libraries-(all-libraries)) on February 2022.

^bDownloaded from <https://enamine.net/compound-collections/screening-collection/> on February 2022.

^cProvided by Specs (<https://www.specs.net/>) on March 2021.

Supporting Information

Table S2 | Dataset overview, with receptor class, ChEMBL ID, response type (inhibition [K_i] or agonism [EC_{50}]), number of compounds in the train and test set, along with the number of activity cliff compounds in the train and test set.

Target name	Receptor Class	ChEMBL ID	Type	n train/test	n cliff train/test
Androgen Receptor (AR)	NR	CHEMBL1871	K_i	525/134	126/31
Cannabinoid receptor 1 (CB1)	GPCR	CHEMBL218	EC_{50}	823/208	292/75
Coagulation factor X (FX)	Protease	CHEMBL244	K_i	2476/621	1080/270
Delta opioid receptor (DOR)	GPCR	CHEMBL236	K_i	2077/521	772/193
Dopamine D3 receptor (D3R)	GPCR	CHEMBL234	K_i	2923/734	1150/291
Dopamine D4 receptor (D4R)	GPCR	CHEMBL219	K_i	1485/374	572/143
Dopamine transporter (DAT)	Other	CHEMBL238	K_i	839/213	209/54
Dual specificity protein kinase CLK4	Kinase	CHEMBL4203	K_i	582/149	51/13
Farnesoid X receptor (FXR)	NR	CHEMBL2047	EC_{50}	503/128	195/50
Ghrelin receptor (GHSR)	GPCR	CHEMBL4616	EC_{50}	543/139	262/68
Glucocorticoid receptor (GR)	NR	CHEMBL2034	K_i	598/152	183/47
Glycogen synthase kinase-3 beta (GSK3)	Kinase	CHEMBL262	K_i	683/173	127/31
Histamine H1 receptor (HRH1)	GPCR	CHEMBL231	K_i	776/197	178/46
Histamine H3 receptor (HRH3)	GPCR	CHEMBL264	K_i	2288/574	865/219
Janus kinase 1 (JAK1)	Kinase	CHEMBL2835	K_i	489/126	36/10
Janus kinase 2 (JAK2)	Kinase	CHEMBL2971	K_i	779/197	95/25
Kappa opioid receptor (KOR) agonism	GPCR	CHEMBL237	EC_{50}	762/193	319/81
Kappa opioid receptor (KOR) inhibition	GPCR	CHEMBL237	K_i	2081/521	753/188
μ -opioid receptor (MOR)	GPCR	CHEMBL233	K_i	2512/630	889/222
Orexin receptor 2 (OX2R)	GPCR	CHEMBL4792	K_i	1174/297	610/153
Peroxisome proliferator-activated receptor alpha (PPAR α)	NR	CHEMBL239	EC_{50}	1377/344	568/141
Peroxisome proliferator-activated receptor delta (PPAR δ)	NR	CHEMBL3979	EC_{50}	900/225	373/94
Peroxisome proliferator-activated receptor gamma (PPAR γ)	NR	CHEMBL235	EC_{50}	1879/470	703/178
PI3-kinase p110-alpha subunit (PIK3CA)	Transferase	CHEMBL4005	K_i	767/193	281/70
Serine/threonine-protein kinase PIM1	Kinase	CHEMBL2147	K_i	1162/294	387/98
Serotonin 1a receptor (5-HT1A)	GPCR	CHEMBL214	K_i	2651/666	917/230
Serotonin transporter (SERT)	Other	CHEMBL228	K_i	1362/342	479/120
Sigma opioid receptor (SOR)	Other	CHEMBL287	K_i	1061/267	371/93
Thrombin	Protease	CHEMBL204	K_i	2201/553	790/199
Tyrosine-protein kinase ABL1	Kinase	CHEMBL1862	K_i	633/161	202/51

Table S3 | Dataset overview, with number of compounds in the train and test set, the number of activity cliff compounds, the mean number of activity cliff 'partners' per activity cliff compound, the number of activity cliff compounds with all activity cliff 'partners' in the test set, and the mean maximal substructure/scaffold/SMILES similarity of all test activity cliff compounds to the train set.

Target name	n train/test	n cliff train/test	Mean cliff partners	All cliff partners in test	Mean max. similarity test cliff to train
AR	525/134	126/31	1.66	3	0.78 / 0.89 / 0.96
CB1	823/208	292/75	2.25	3	0.81 / 0.92 / 0.96
FX	2476/621	1080/270	3.25	29	0.83 / 0.94 / 0.96
DOR	2077/521	772/193	2.91	23	0.83 / 0.92 / 0.96
D3R	2923/734	1150/291	2.73	24	0.81 / 0.95 / 0.95
D4R	1485/374	572/143	2.68	12	0.79 / 0.95 / 0.95
DAT	839/213	209/54	1.73	11	0.75 / 0.90 / 0.92
CLK4	582/149	51/13	1.25	0	0.67 / 0.93 / 0.92
FXR	503/128	195/50	2.96	2	0.81 / 0.94 / 0.97
GHSR	543/139	262/68	5.51	0	0.82 / 0.94 / 0.96
GR	598/152	183/47	2.64	9	0.80 / 0.92 / 0.95
GSK3	683/173	127/31	1.59	3	0.78 / 0.92 / 0.93
HRH1	776/197	178/46	1.75	9	0.77 / 0.90 / 0.93
HRH3	2288/574	865/219	2.82	17	0.81 / 0.96 / 0.95
JAK1	489/126	36/10	1.43	0	0.82 / 0.91 / 0.96
JAK2	779/197	95/25	2.3	1	0.77 / 0.95 / 0.94
KOR (a)	762/193	319/81	4.43	4	0.84 / 0.92 / 0.96
KOR (i)	2081/521	753/188	2.99	15	0.83 / 0.92 / 0.96
MOR	2512/630	889/222	3.71	11	0.79 / 0.95 / 0.96
OX2R	1174/297	610/153	2.37	15	0.84 / 0.92 / 0.96
PPAR α	1377/344	568/141	2.84	6	0.83 / 0.93 / 0.96
PPAR δ	900/225	373/94	2.5	18	0.83 / 0.93 / 0.96
PPAR γ	1879/470	703/178	2.18	11	0.80 / 0.95 / 0.95
PIK3CA	767/193	281/70	3.74	6	0.79 / 0.94 / 0.95
PIM1	1162/294	387/98	2.31	21	0.82 / 0.94 / 0.95
5-HT1A	2651/666	917/230	2.16	15	0.80 / 0.94 / 0.94
SERT	1362/342	479/120	2.27	14	0.80 / 0.94 / 0.94
SOR	1061/267	371/93	2.32	17	0.83 / 0.94 / 0.96
Thrombin	2201/553	790/199	3.23	6	0.80 / 0.96 / 0.96
ABL1	633/161	202/51	3.16	18	0.82 / 0.92 / 0.96

Supporting Information

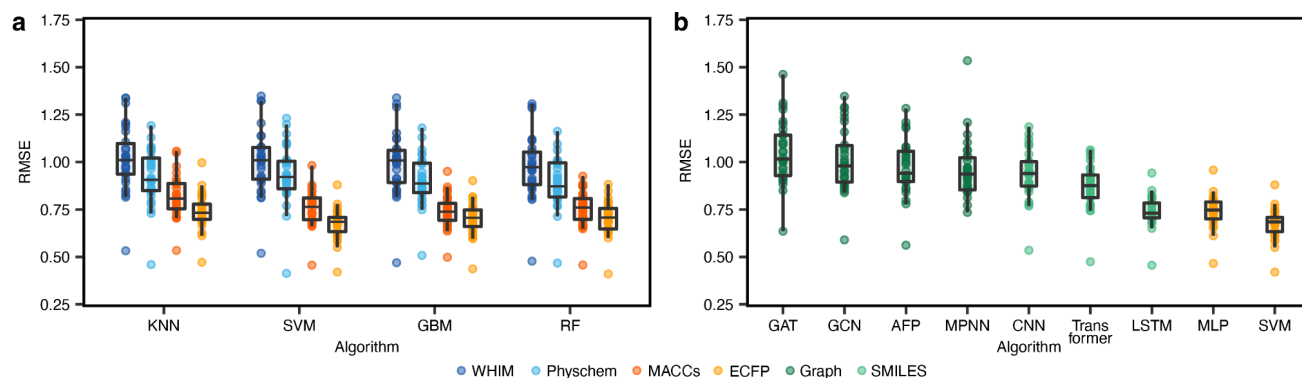


Fig. S1 | Overall performance of machine learning methods on all targets. a, RMSE using different 'classical' machine learning algorithms and molecular descriptors. b, RMSE using deep learning methods and unstructured molecular representations.

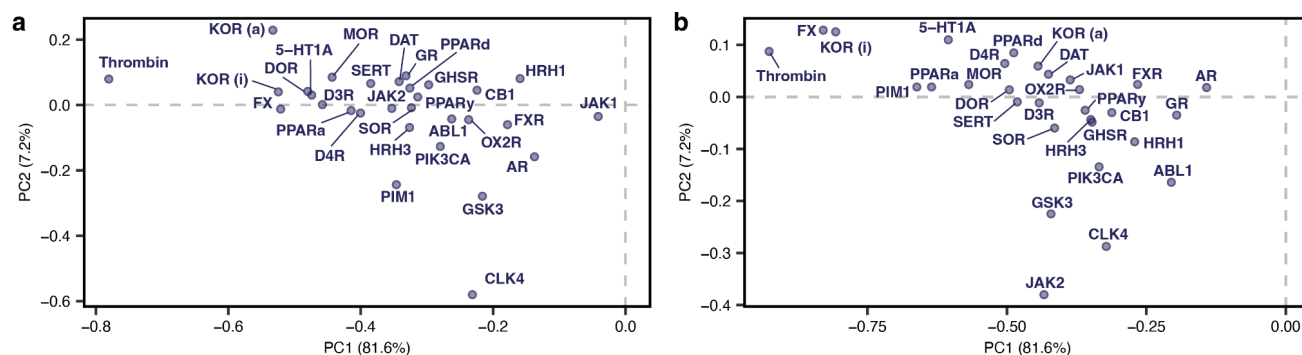


Fig. S2 | PCA loadings of all methods. a, Effects of individual data sets (loadings for PC1 and PC2) on the PCA of 'classical' machine learning methods (see Fig. 3b). b, Effects of individual data sets (loadings for PC1 and PC2) on the PCA of deep learning methods (see Fig. 4b).

Supporting Information

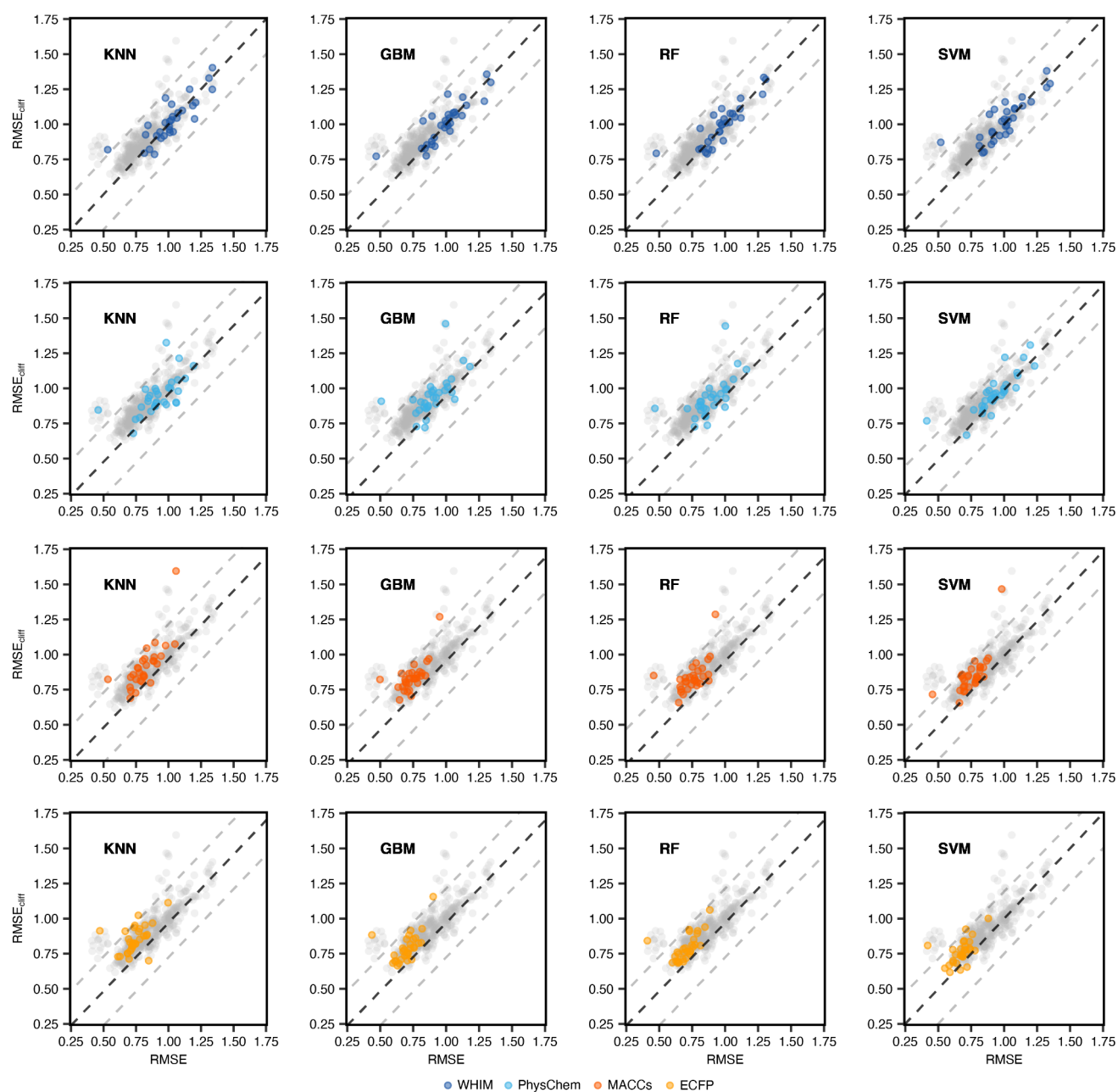


Fig. S3 | Relative prediction error of activity cliff compounds. Prediction error on activity cliff compounds compared to all compounds for all classical machine learning algorithms and molecular descriptor combinations.



Fig. S4 | Statistical differences between the $RMSE_{cliff}$ values obtained by different machine learning strategies. Asterisks indicate statistically significant differences ($p < 0.05$) between pairs of methods, as obtained by the Wilcoxon rank-sum test (adjusted for false discovery rate using a Benjamini-Hochberg procedure).

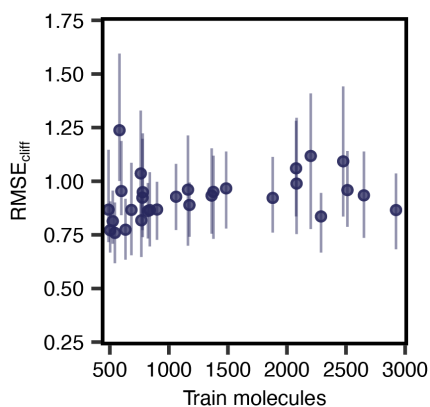


Fig. S5 | The effect of the number of training molecules on $RMSE_{cliff}$. Error bars indicate the lowest and highest $RMSE_{cliff}$.

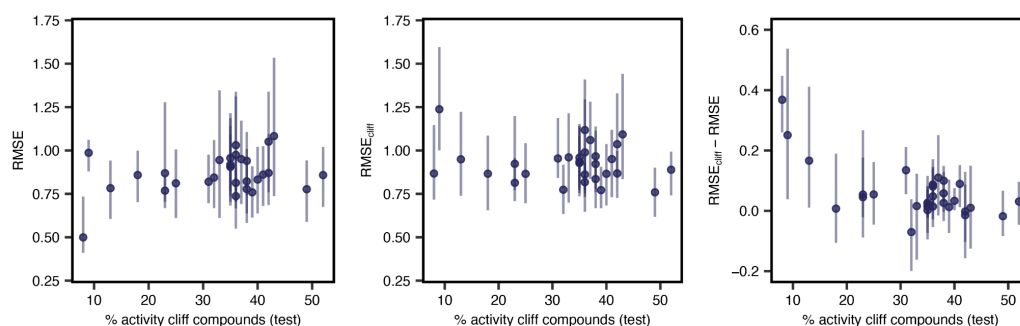


Fig. S6 | The effect of the fraction of activity cliff compounds and model performance. **a**, Overall model performance (RMSE). Pearson correlation (ρ) = 0.21. **b**, Performance on activity cliff compounds ($\text{RMSE}_{\text{cliff}}$, ρ =-0.12). **c**, Difference between $\text{RMSE}_{\text{cliff}}$ and RMSE (ρ =-0.54). Error bars indicate the lowest and highest RMSEs.

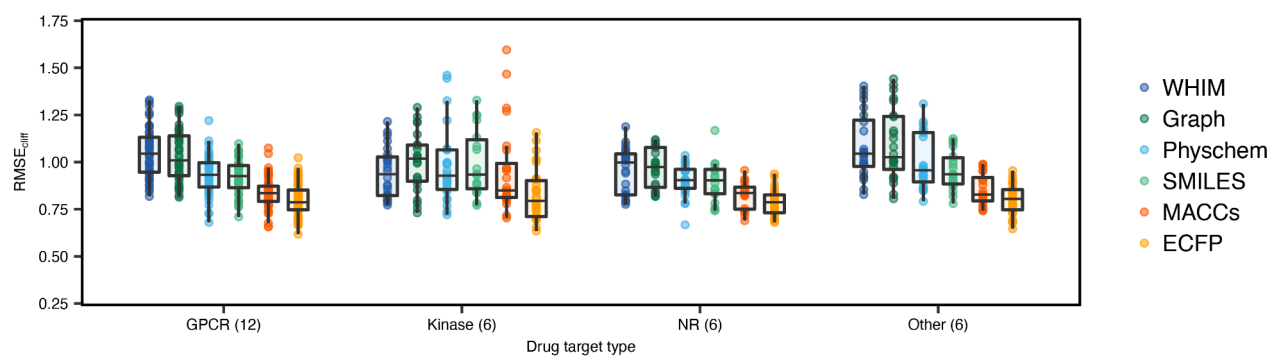


Fig. S7 | The effect of different drug target classes on $\text{RMSE}_{\text{cliff}}$. All machine learning strategies are grouped by molecular descriptor/representation.