

IMPACT OF MOLECULAR REPRESENTATIONS ON DEEP LEARNING MODEL COMPARISONS IN DRUG RESPONSE PREDICTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Deep learning (DL) plays a crucial role in tackling the complexity and heterogeneity of cancer, particularly in predicting drug response. However, the effectiveness of these models is often hindered by inconsistent benchmarks and disparate data sources. To address the gaps in comparisons, we introduce CoMParison workflow for Cross Validation (CMP-CV), an automated cross-validation framework that trains multiple models with user-specified parameters and evaluation metrics. The effectiveness of DL models in predicting drug responses is closely tied to the methods used to represent drugs at the molecular level. In this contribution, we benchmarked commonly leveraged drug representations (graph, molecular descriptors, molecular fingerprints, and SMILES) to learn and understand the predictive capabilities of the models. We compare the ability of different drug representations to encode different structural properties of the drugs by using prediction errors made by models in different drug descriptor domains. We find that, in terms of the average prediction error over the entire test set, molecular descriptors and Morgan fingerprints perform slightly better than the others. However, we also observe that the rankings of the model performance vary in different regions over the descriptor space studied in this work, emphasizing the importance of domain-based model comparison when selecting a model for a specific application. Our efforts are part of CANcer Distributed Learning Environment (CANDLE), enhancing the model comparison capabilities in cancer research and driving the development of more effective strategies for drug response prediction and optimization.

23 1 INTRODUCTION

24 Cancer research is currently exploring innovative techniques to enhance treatment outcomes through the use of analytical models called Drug Response Prediction (DRP) models Yancovitz et al. (2012); Fisher et al. (2013); Adam et al. (2020). These models utilize machine learning (ML) and deep learning (DL) algorithms to forecast tumor responses to various drug treatments without the need for specific biomarkers. However, accurately predicting drug responses using ML and DL models is a critical challenge Baptista et al. (2020); Adam et al. (2020); Zuo et al. (2021). Each study typically develops custom model implementation and validation strategies, making it difficult to assess model capabilities across drug representation methods, architectures, and datasets Partin et al. (2023). With the increasing complexity of models and the diversity of datasets, there is a pressing need for robust methodologies to compare these models Park et al. (2023). However, the current landscape lacks consistency and standardization in terms of model comparison techniques. Traditional approaches often rely on performance scores from original publications, which leads to incomparable and inconsistent results. This hinders elucidating the precise factors that drive predictive performance. Therefore, it is crucial to establish a standardized and comprehensive comparison workflow to address the urgent need to understand drug representation and its impact on drug response prediction error.

40 In light of these challenges, we recently implemented the [CoMParison workflow for Cross Validation \(CMP-CV\)](#) - an automated cross-validation framework that enables simultaneous training and evaluation of multiple DL models using standardized datasets, preprocessing, and performance

43 metrics. CMP-CV provides infrastructure for controlled experimentation by systematically varying
44 model hyperparameters and architectures. It also has built-in support for custom analytical func-
45 tions, which facilitates deeper analysis of model representations and uncertainties.

46 When applying DRP models in real-world applications, such as predicting drug efficacy or identify-
47 ing suitable cancer treatments, selecting the best model is crucial. While existing comparison meth-
48 ods utilize metrics like R2 (coefficient of determination), RMSE (Root Mean Squared Error), and
49 AUC (Area Under the ROC Curve) to assess overall model accuracy, they fail to reveal critical in-
50 formation about each model’s unique strengths and weaknesses. For instance, certain models might
51 excel in specific domains of the drug descriptor space but be less accurate in other regions. In this
52 work we analyze model performance within distinct domains of the drug descriptor space to iden-
53 tify the most effective models for specific drug candidates and determine if certain drug’s molecular
54 representations are superior to others. This type of analysis enables more informed decision-making
55 when selecting a model for practical applications.

56 A significant challenge in drug response prediction is the lack of consensus on a suitable molecular
57 representation, which is further complicated by the diversity of DRP models. Therefore, large-scale
58 model comparison is necessary, and CMP-CV serves as a robust framework for this purpose. Its
59 ability to accommodate user-defined Python functions to analyse model predictions allows for com-
60 prehensive benchmarking of models to determine the impact of various molecular representations on
61 prediction errors. The current application of CMP-CV focuses mainly on comparing Cancer Drug
62 Response Prediction (CDRP) models across diverse molecular descriptor spaces. This comprehen-
63 sive comparison not only provides a deeper understanding of drug representation and its impact on
64 drug response prediction errors but also highlights the relative strengths of various models on drug
65 properties in different domains.

66 2 RESULTS AND DISCUSSION

67 2.1 CMP-CV: DEEP LEARNING MODEL COMPARISON FRAMEWORK

68 The CANDLE/Supervisor framework (Wozniak et al., 2018) is a workflow application system de-
69 signed for HPC infrastructure. Supervisor consists of multiple exemplar workflows, including sim-
70 ple sweeps, automated hyperparameter optimization, and other data analysis workloads. It is capable
71 of calling into user-specified model codes via multiple techniques, including direct Python library in-
72 vocation, shell command lines, and Linux container invocation. Supervisor coordinates these model
73 executions via CANDLE “hyperparameters,” which extend the notion of model training hyperpa-
74 rameters to include a range of other control variables. The hyperparameter set is standardized by the
75 CANDLE Library (CANDLE Team, 2018).

76 The CMP-CV employs the Supervisor framework, which facilitates the integration of the container-
77 ized models described here along with their hyperparameters. Inside the workflow, depicted in
78 Figure 1, a list of hyperparameter combinations is specified in an external file, encoded in a JSON
79 format, and each training run is performed concurrently. In this manner, a very large HPC system
80 can be efficiently used. Supervisor monitors training progress and keeps resources busy, almost
81 eliminating the need for the workflow developer to consider concurrency. As each training run
82 completes, a comparison function is invoked across the error metrics produced during training.

83 The CMP-CV system’s unique integrated functionality offers a seamless process for analyzing pre-
84 diction results, delivering comparable output metrics, and facilitating the integration of custom ana-
85 lytical functions, thereby providing users with a tailored analytical experience. One key feature that
86 sets CMP-CV apart is its ability to accommodate user-defined Python functions, enabling users to
87 seamlessly integrate custom analytical functions into the workflow. We utilised this capability to
88 obtain drug response prediction errors in different regions in a drugs’ molecular descriptor space.
89 Our results highlight the importance of understanding where each model excels; this will enable us
90 and the rest of the community to better leverage their predictive power in future applications.

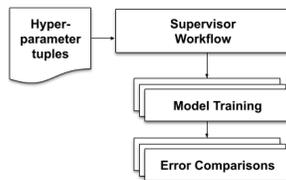


Figure 1: Architecture of the CMP-CV. The ‘Error Comparisons’ functionality contains python scripts to calculate the model errors corresponding to different regions in a drug’s molecular descriptor space.

91 2.2 OVERVIEW OF DRUG FEATURES AND REPRESENTATIONS

92 In the field of drug design and characterization, each drug is distinguished by a unique set of descrip-
 93 tors such as molecular structure, substructures, functionalities, physicochemical and biochemical
 94 properties, known targets, and clinical usage. These descriptors form the drug or molecular descrip-
 95 tor space. To apply machine learning techniques, it is necessary to create a numerical representation
 96 of these multifaceted descriptors. Investigating the effects of molecular representation on prediction
 97 accuracy provides valuable insights into current limitations of drug response modeling approaches.
 98 Our hypothesis is that the efficiency of a molecular representation depends on the model’s ability to
 99 predict outcomes across various domains of the molecular descriptor space.

100 For instance, a molecular representation that includes fine details about ring structure can ensure
 101 good performance of the model, regardless of the number of rings in the drug molecule. It is im-
 102 portant to mention that the model’s performance variation for molecules with different numbers of
 103 rings is not solely due to its molecular representation strength. Other aspects of the molecule, such
 104 as molecular weight or number of atoms/hydrogen bonds can also change. However, if a model con-
 105 sistentlly fails to achieve good performance in a particular domain of the descriptor space, it indicates
 106 that the model’s molecular representation is weak in that region.

107 2.3 CURATED EXISTING MACHINE LEARNING MODELS FOR COMPARISON AND 108 BENCHMARKING

109 In our effort to understand the relationship between molecular representations and drug response
 110 predictions, we conducted a thorough curation and analysis of existing CDRP models, such as
 111 GraphDRP, DeepTTC, and HiDRA Nguyen et al. (2022); Jiang et al. (2022); Jin & Nam (2021).
 112 By applying CMP-CV to a standardized CTRPv2¹ dataset, we were able to compare and cross-
 113 validate these models, yielding important metrics that highlight their relative performance across
 114 the molecular descriptor space. This approach to curation and comparison represents a significant
 115 step towards enhancing the field of drug response prediction models.

116 Based on our literature survey on CDRP models Baptista et al. (2020); Partin et al. (2023), we
 117 identify that the models primarily use four categories of molecular representations: graph structures,
 118 SMILES encodings, Morgan fingerprints, and molecular descriptors. In Table 1, we list the CDRP
 119 models that leverage these distinct molecular representations. Our work focuses on comparing these
 120 four types of representations to understand their strengths and limitations.

121 To ensure a fair comparison of different drug representations, we also developed a model with the
 122 ability to switch between different molecular representations while using the same cell line represen-
 123 tation. These models are hereafter referred to as **Graph**, **SMILES**, **Morgan** and **Descriptor**. More
 124 details about these models are given in the Appendix. Below is a brief description of the models
 125 from the literature.

126 **GraphDRP.** Nguyen et al. (2022) GraphDRP encodes drug molecules using graph convolutional
 127 layers followed by fully connected layers to arrive at a vector representation of length 128. The cell
 128 lines are initially represented using one hot encoding (735 dimensions). 1D convolutional operations
 129 followed by fully connected layers are used to convert the one hot encoded representation to a vector

¹CSA Benchmark Datasets

Table 1: Models categorized based on the kind of drug representation they use

Representation type	Models
Graph structure	SWnet (Zuo et al., 2021), DRPreter (Shin et al., 2022), GraphDRP Nguyen et al. (2022), DrugGCN(Kim et al., 2021)
SMILES encoding	DeepTTC Jiang et al. (2022), Paccmann Oskooei et al. (2019), tCNNS Liu et al. (2019)
Morgan fingerprints	DrugCell Kuenzi et al. (2020), HiDRA Jin & Nam (2021), DeepDSC Li et al. (2021), PathDSP Tang & Gottlieb (2021)
Molecular descriptors	CDRscan Chang et al. (2018), REFINED Bazgir et al. (2020), IGTD Zhu et al. (2021)

of 128 elements. The drug and cell line representations are concatenated and fed through another fully connected neural network to arrive at the final prediction.

DeepTTC. Jiang et al. (2022) In DeepTTC, the SMILES string is tokenized using Explainable Substructure Partition Fingerprints (ESPF) Huang et al. (2019). The SMILES string is decomposed into multiple substructures and each substructure is assigned a number based on a provided vocabulary of substructures. This sequence of numbers is converted to a one-hot encoded matrix, and then transformed using a weight matrix. To this representation, a positional encoding is added to create the initial representation of the drug. This representation is sent through transformer encoder layers that contain multihead attention to arrive at the final drug representation.

HiDRA. Jin & Nam (2021) HiDRA is an attention-based model that aggregates gene expression data to drug fingerprint features to create a pathway-level network between the drug and cell line. The overall architecture is composed of four networks encompassing a drug, gene, and pathway level network followed by the response prediction network. Morgan fingerprints are used for drug representations and genes were grouped to pathways through the KEGG Pathway database to create the cell line feature. 4592 unique genes were used to create these features.

ExtraTreesRegressor. Geurts et al. (2006); Pedregosa et al. (2011) For the comparison, we also use an ExtraTreesRegressor model. This model is based on an ensemble of decision trees and does not utilize DL techniques. The model takes a simple concatenation of drug features and gene expression values of the cell lines as input.

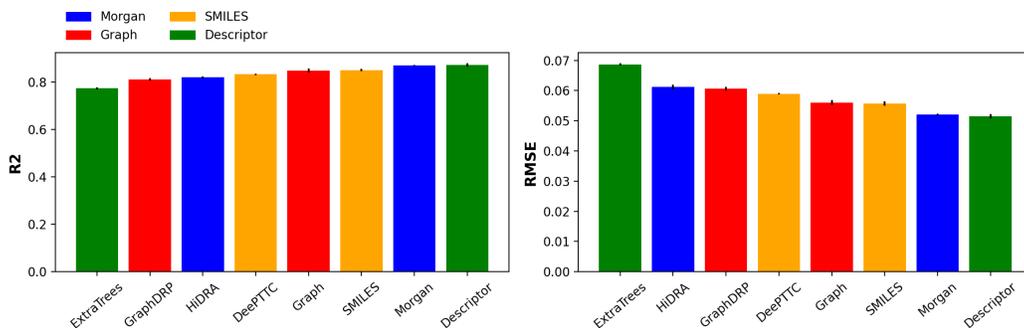


Figure 2: Comparative analysis of model prediction errors based on AUC^{DR} . Colors represent the type of representation used in each model.

2.4 MODEL COMPARISON

The CDRP models mentioned earlier were trained using the CTRPv2 dataset, which measures gene expression values in transcripts per million (TPM). These values were obtained from the CCLE DepMap² portal, while the response data were sourced from CTRP. As the dose-independent drug response metric, we use area under the dose response curve (AUC^{DR}). This AUC^{DR} is what the

²<https://depmap.org/portal/>

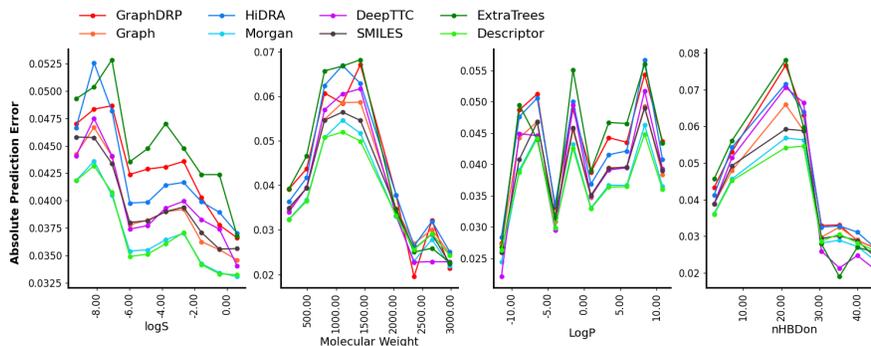


Figure 3: This figure presents a detailed analysis of AUC^{DR} prediction errors in the domains of important drug properties such as $\log S$, molecular weight, LogP , and $n\text{HBDon}$.

154 **CDRP models attempt to predict.** Further details on the dataset and model training are given in the
 155 **Methods section.**

156 The prediction accuracies for AUC^{DR} are displayed in Figure 2. Based on the R^2 results, it is
 157 observed that models utilizing molecular descriptors and Morgan fingerprints perform marginally
 158 better than the others. However, in this work, we aim to compare the performance of different
 159 models across various regions in the molecular descriptor space. To facilitate this comparison, we
 160 use Mordred Moriwaki et al. (2018) to generate molecular descriptors of the drugs. Descriptors
 161 that require three-dimensional coordinates were not taken into consideration. After obtaining the
 162 molecular descriptor values, they were divided into bins based on their ranges. These bins define
 163 the domains of the descriptors. Domain boundaries of continuous descriptors were found using
 164 NumPy³'s histogram function. Every unique value of a categorical descriptor was considered as a
 165 domain. A categorical descriptor is defined as one which consists of less than 20 unique integer
 166 values.

167 For instance, if a molecular descriptor value ranges from 5 to 95, to evaluate the performance of
 168 each model, we can group the molecules into intervals of 10 descriptor value units, such as 5-15,
 169 15-25, and so on. This approach allows us to analyze a model's predictions in different regions in the
 170 descriptor space. In Figure 3, we present the variations in the AUC^{DR} prediction error in the domains
 171 of solubility ($\log S$), molecular weight, LogP , and the number of hydrogen bond donors ($n\text{HBDon}$),
 172 which are crucial descriptors in drug design Di & Kerns (2016). The information presented in
 173 Figure 3 offers two main advantages: Firstly, it increases the awareness of the users of these models
 174 regarding the limitations of the models in terms of the properties of the drug molecules. Secondly,
 175 it provides model developers with valuable insights into the deficiencies of their models.

176 2.4.1 EXPLORING DESCRIPTOR DOMAINS OF MODEL APPLICABILITY

177 Drug response prediction errors in the domains of $\log S$, molecular weight, LogP , and $n\text{HBDon}$ can
 178 significantly impact the performance of drug response prediction models. By identifying the domain
 179 errors of different models, we can determine which molecular descriptors have not been adequately
 180 represented in the model. This information can be used to enhance the performance of models by
 181 improving their representation in these descriptor domains.

182 Based on Figure 3, none of the ML models appear to perform well when the $\log S$ of the drugs is
 183 less than -7, and their errors decrease as the drug solubility increases. The Descriptor and Morgan
 184 models can be expected to perform best when predicting highly soluble drug candidates. These
 185 results facilitate the domain-wise representations comparison. For instance, in the high solubility
 186 regime ($\log S > 0$), considering only the models with the same cell line representation, the goodness
 187 of the drug representation can be ranked as Morgan > Descriptor > Graph > SMILES.

188 In fact, one can construct a table showing the error-based model rankings for each domain as shown
 189 in Table (a), Figure 4. This resource empowers the systematic evaluation and determination of the

³NumPy

190 most efficacious models for drugs, characterized by distinct molecular attribute. For example, if we
 191 need to determine the best model for drugs with solubility varying in a wide range, the Descriptor
 192 model is the clear winner, followed by the Morgan model. For nHBDon however, the Descriptor
 193 model is more suitable when $2 > \text{nHBDon} < 8$ (see Appendix Table 3). For drugs with over 35
 194 hydrogen bond donors, DeepTTC is a superior model (Appendix Table 4). These tables system-
 195 atically categorize models based on their error rates within specific molecular descriptor domains,
 196 aiding in the seamless identification of the most adept models for predicting drug responses for drug
 197 candidates with particular molecular properties. Such information is useful for the robustness and
 198 reliability of drug response predictions.

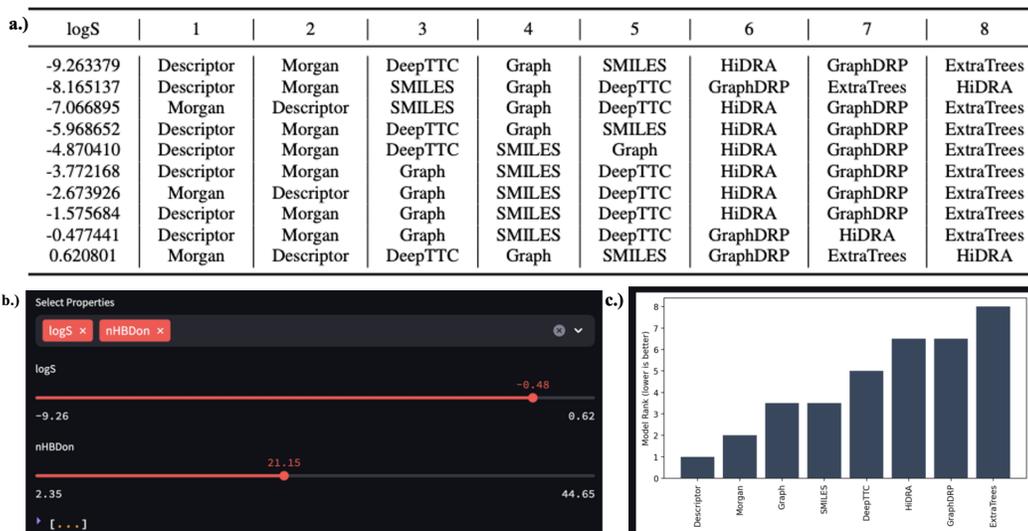


Figure 4: Table (a) systematically categorizes models based on their error rates within specific logS domains. Images (b) and (c) depict a web application that allows users to find model ranking based on multiple distinct molecular descriptor values. Values of more than 700 molecular descriptors can be changed (b) to obtain the corresponding model ranks (c).

199 We also designed a web application which allows a user to identify the models best suited for drug
 200 candidates described using multiple molecular descriptors. This interface allows the user to add as
 201 many as 786 molecular descriptors and adjust their values using the associated sliders. As shown in
 202 Figure 4 (b) and (c), once the descriptor values are chosen, a rank for each model is presented. These
 203 ranks are calculated by first looking up the model ranks corresponding to the chosen properties from
 204 tables similar to Figure 4, Table (a). If n property values are selected, we have n sets of model ranks.
 205 Each set contains ranks of m models considered in the comparison. Next, the average rank of each
 206 model is found which is considered as the final model rank. Models are ranked from 1 to m , where
 207 1 is the best rank and m is the worst rank.

208 2.4.2 IDENTIFYING MODEL REPRESENTATION DEFICIENCIES

209 When dealing with over 1000 molecular descriptors, it can be challenging to determine which ones
 210 are most important for understanding how drug representation affects model performance. A logical
 211 assumption is that if a particular descriptor has been accurately encoded by a representation, then
 212 domain errors associated with that descriptor will be minimal. Conversely, if a representation fails to
 213 capture the intricate details of a molecular descriptor, domain errors corresponding to that descriptor
 214 will be significant.

215 We can determine the maximum error of a model for a specific domain. For instance, HiDRA has
 216 a maximum error of approximately 0.0525 at $\log S = -8$ (Figure 3). These errors can be utilized to
 217 identify molecular descriptors that are not adequately encoded in the model’s representation. This
 218 particular insight into individual errors per model can act as a pivotal tool for discerning molecular

219 descriptors that remain inadequately encoded within the model’s architecture. Figure 5 displays the
 220 largest maximum and smallest minimum domain errors for each model, consisting of the top 5.

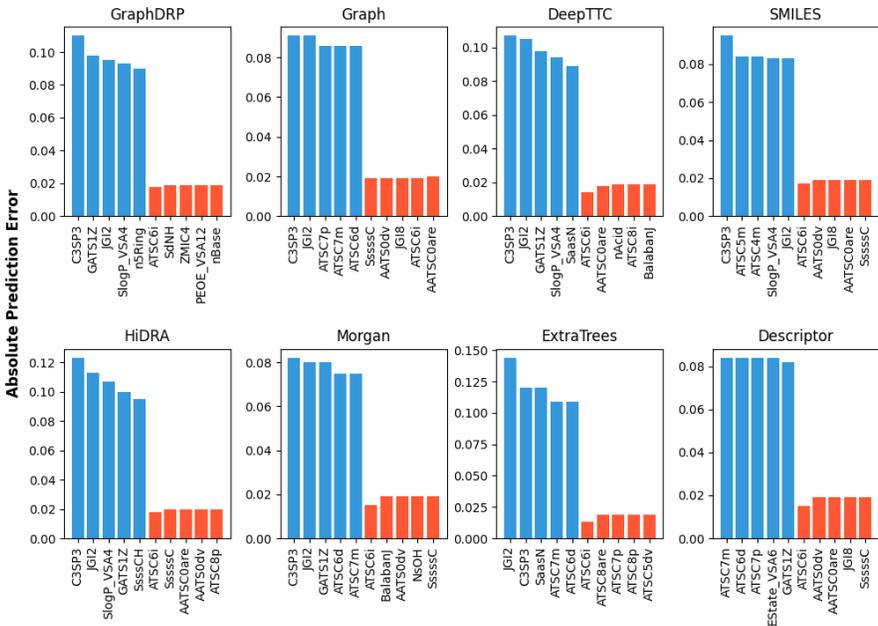


Figure 5: Descriptors that made maximum and minimum domain errors. Vertical axis is the number of drug response values.

221 We notice that GATS1Z, C3SP3, SlogP_VSA4 and JGI2 are among the descriptors having the largest
 222 domain errors for most of the models. GATS1Z is the geary coefficient of lag 1 weighted by atomic
 223 number, C3SP3 is SP3 carbon bound to 3 other carbons, SlogP_VSA4 is a MOE type descriptor
 224 based on Wildman-Crippen LogP and surface area contribution, and JGI2 is the mean topological
 225 charge index of order 2 Moriawaki et al. (2018).

Table 2: Drug response prediction errors associated with $AUC^{DR} < 0.75$ and $AUC^{DR} \geq 0.75$ cell-line – drug pairs.

	MAE	RMSE
$AUC^{DR} < 0.75$	0.06 ± 0.003	0.082 ± 0.005
$AUC^{DR} \geq 0.75$	0.032 ± 0.001	0.044 ± 0.001

226 Figure 6 further demonstrates the error oscillations for the aforementioned descriptors, unfolding
 227 domains with the most significant errors: $GATS1Z < 0.2$, $C3SP3 > 9$, $50 > SlogP_VSA4 <$
 228 55 , and $JGI2 < 0.04$. Such intricate data prove invaluable in decoding the root causes of subpar
 229 model performance and paves the path for consequential model enhancements. In fact, we notice
 230 that the prediction errors associated with $AUC^{DR} < 0.75$ drugs are significantly higher than those of
 231 $AUC^{DR} \geq 0.75$ drugs (see Table 2). [In the Appendix, we investigate whether the error from the](#)
 232 [above descriptors is due to a common molecular structure motif or a deficiency of the representa-](#)
 233 [tion.](#)

234 Investigating further, observing drug response values (AUC^{DR}) in domains $GATS1Z < 0.2$ and
 235 $GATS1Z > 1.5$ (refer to Figure 7) reveals certain AUC^{DR} values in the $GATS1Z < 0.2$ distribu-
 236 tion do not originate from a densely populated region in the complete distribution. This correlation
 237 highlights the association of $GATS1Z < 0.2$ drugs with diminished drug response values.

238 [In order to demonstrate how one can potentially use the information about domain errors to improve](#)
 239 [the model predictions, we pretrained the GraphDRP model to predict the molecular descriptors](#)
 240 [corresponding to largest error domains; GATS1Z, C3SP3, SlogP_VSA4, JGI2 and n5Ring. The](#)

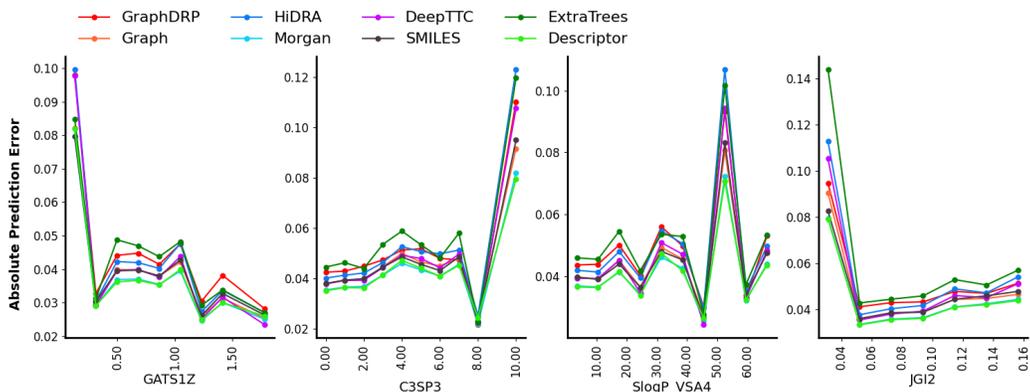


Figure 6: Visualization of error fluctuations within high-error descriptors domains. This plot is crucial for identifying and understanding the underlying causes of model performance

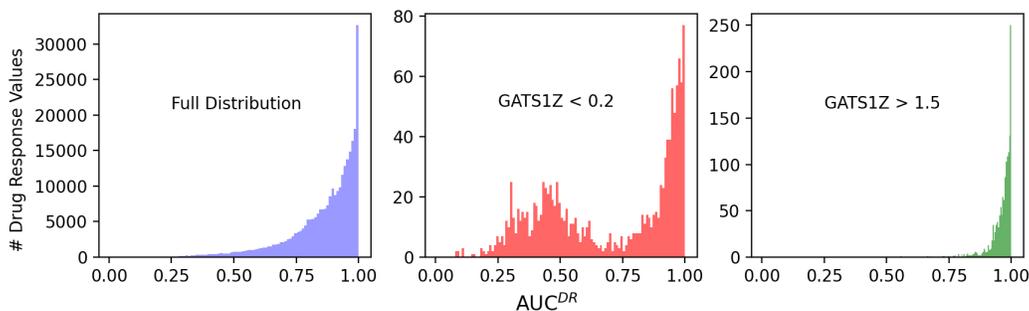


Figure 7: Examination of AUC^{DR} distributions in different GATS1Z regions in the dataset.

241 pretraining GraphDRP model was created by replacing the last linear layer with three layers; one
 242 with three outputs for GATS1Z, SlogP_VSA4 and JGI2, another two with 11 and 8 outputs for
 243 unique values of C3SP3 and n5Ring respectively. The model was trained for 100 epochs with
 244 early stopping. After training, the weights of this model were loaded to the original GraphDRP
 245 model and trained for 100 epochs. Using the predictions of this model we obtained the domain
 246 errors again. Comparison of the logS, Molecular Weight, logP and nHBDOn domain errors
 247 before and after pretraining are shown in Figure 8. We see significant error reductions in logS and LogP
 248 domains. We also observe a test set R2 improvement from 0.812 to 0.838 due to pretraining.

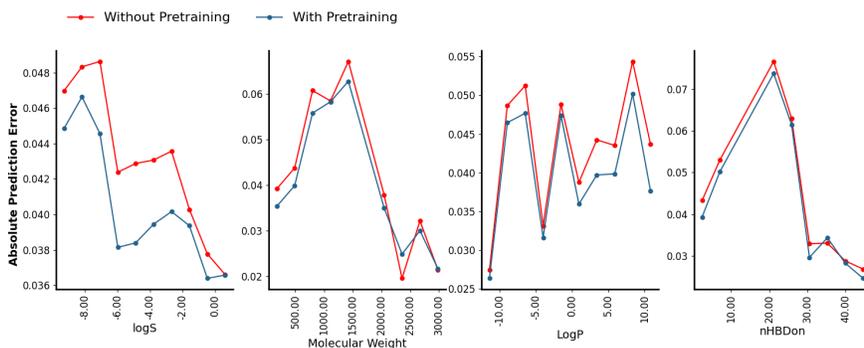


Figure 8: Reduction in GraphDRP error after pretraining.

249 3 METHODS

250 3.1 DATA AND MODEL TRAINING

251 The CTRPv2 dataset used in this work is from the CSA Benchmark Datasets curated as part of
252 the IMPROVE⁴ project. Cell line response data of this dataset were extracted from the Cancer
253 Therapeutics Response Portal version 2. After extracting multi-dose viability data, a unified dose
254 response fitting pipeline was used to calculate the dose-independent response metric, area under the
255 dose response curve (AUC^{DR}). Drug data have been retrieved from PubChem (Kim et al., 2023).
256 The CTRPv2 dataset has 720 cell lines and 494 unique drugs. The total number of drug response
257 values is 286665.

258 The full dataset was divided into ten random train, validation, and test folds using different random
259 seeds. This ensured that every drug-cancer cell combination was predicted at least once. The models
260 were trained using the train set, the validation set is used for saving the best models. Except for the
261 HiDRA model, others were trained for 100 epochs. As it takes about 30 minutes for a HiDRA epoch
262 to complete, it was trained for 20 epochs. The predictions made by each of the test sets are recorded.
263 These predictions are used to find the mean and the standard deviation of the prediction errors across
264 the ten runs.

265 4 CONCLUSIONS

266 Domain error is a significant factor that can impact the performance of drug response prediction
267 models. By utilizing our recently implemented CMP-CV framework and understanding the domain
268 errors of different CDRP models, we can identify the molecular descriptors that have not been en-
269 coded with sufficient detail by the model’s representation. This knowledge can be used to guide the
270 selection of models for specific applications. [We also introduce a web application which enables
271 users to find the CDRP models better suited for drugs having specific molecular properties.](#) We
272 found that the prediction accuracy for drugs with a low solubility, particularly below the threshold
273 $\log S < -7$, dramatically decreases regardless of molecular representation. Increased drug solubility
274 notably improves prediction accuracy with two models based on molecular descriptors and Morgan
275 fingerprints performing substantially better than other representation across the entire range for sol-
276 ubility. In addition, we can use the domain errors of models to improve the performance of models
277 by focusing on improving their representation in these descriptor domains. Our analysis revealed
278 that GATS1Z, C3SP3, SlogP_VSA4 and JGI2 are among the domains that might not be encoded
279 with adequate detail by any of the molecular representations that could help improve the model pre-
280 diction. By avoiding models with large errors in the domain of interest, we can obtain more reliable
281 predictions from the models. [We also show that using the descriptors corresponding to high-error
282 domains as pretraining targets has a potential to improve model predictions.](#)

283 In conclusion, molecular representation and feature domain exploration lays a robust foundation
284 for not only recognizing and comprehending the domains contributing to the largest errors but also
285 offers an opportunity for substantial model improvement.

⁴IMPROVE

286 REPRODUCIBILITY STATEMENT

287 We have provided the instructions to run the CMP-CV and the code to perform the data analysis
288 shown in the paper in the code.zip file.

289 REFERENCES

- 290 George Adam, Ladislav Rampásek, Zhaleh Safikhani, Petr Smirnov, Benjamin Haibe-Kains,
291 and Anna Goldenberg. Machine learning approaches to drug response prediction: chal-
292 lenges and recent progress. *npj Precision Oncology*, 4(1):1–10, June 2020. ISSN 2397-
293 768X. doi: 10.1038/s41698-020-0122-1. URL [https://www.nature.com/articles/
294 s41698-020-0122-1](https://www.nature.com/articles/s41698-020-0122-1). Number: 1 Publisher: Nature Publishing Group.
- 295 Timothy G. Armstrong, Justin M. Wozniak, Michael Wilde, and Ian T. Foster. Compiler techniques
296 for massively scalable implicit task parallelism. In *Proc. SC*, 2014.
- 297 Delora Baptista, Pedro G Ferreira, and Miguel Rocha. Deep learning for drug response prediction
298 in cancer. *Briefings in Bioinformatics*, 22(1):360–379, 01 2020. ISSN 1477-4054. doi: 10.1093/
299 bib/bbz171. URL <https://doi.org/10.1093/bib/bbz171>.
- 300 Omid Bazgir, Ruibo Zhang, Saugato Rahman Dhruba, Raziur Rahman, Souparno Ghosh, and
301 Ranadip Pal. Representation of features as images with neighborhood dependencies for com-
302 patibility with convolutional neural networks. *Nature Communications*, 11(1):4391, September
303 2020. ISSN 2041-1723. doi: 10.1038/s41467-020-18197-y. URL [https://www.nature.
304 com/articles/s41467-020-18197-y](https://www.nature.com/articles/s41467-020-18197-y). Number: 1 Publisher: Nature Publishing Group.
- 305 CANDLE Team. The `CandleLib` GitHub repository, 2018. [https://github.com/ECP-
306 CANDLE/candle_lib](https://github.com/ECP-CANDLE/candle_lib).
- 307 Yoosup Chang, Hyejin Park, Hyun-Jin Yang, Seungju Lee, Kwee-Yum Lee, Tae Soon Kim, Jongsun
308 Jung, and Jae-Min Shin. Cancer Drug Response Profile scan (CDRscan): A Deep Learning
309 Model That Predicts Drug Effectiveness from Cancer Genomic Signature. *Scientific Reports*,
310 8(1):8857, June 2018. ISSN 2045-2322. doi: 10.1038/s41598-018-27214-6. URL [https://
311 //doi.org/10.1038/s41598-018-27214-6](https://doi.org/10.1038/s41598-018-27214-6).
- 312 Li Di and Edward H. Kerns. Chapter 1 - introduction. In Li Di and Edward H.
313 Kerns (eds.), *Drug-Like Properties (Second Edition)*, pp. 1–3. Academic Press, Boston,
314 second edition edition, 2016. ISBN 978-0-12-801076-1. doi: [https://doi.org/10.1016/
315 B978-0-12-801076-1.00001-0](https://doi.org/10.1016/B978-0-12-801076-1.00001-0). URL [https://www.sciencedirect.com/science/
316 article/pii/B9780128010761000010](https://www.sciencedirect.com/science/article/pii/B9780128010761000010).
- 317 R. Fisher, L. Pusztai, and C. Swanton. Cancer heterogeneity: implications for targeted therapeutics.
318 *British Journal of Cancer*, 108(3):479–485, February 2013. ISSN 1532-1827. doi: 10.1038/
319 bjc.2012.581. URL <https://www.nature.com/articles/bjc2012581>. Number: 3
320 Publisher: Nature Publishing Group.
- 321 Message Passing Interface Forum. MPI: A message-passing interface standard, 1994. URL
322 citeseer.ist.psu.edu/forum94mpi.html.
- 323 Pierre Geurts, Damien Ernst, and Louis Wehenkel. Extremely randomized trees. *Machine Learning*,
324 63(1):3–42, April 2006. ISSN 1573-0565. doi: 10.1007/s10994-006-6226-1. URL [https://
325 //doi.org/10.1007/s10994-006-6226-1](https://doi.org/10.1007/s10994-006-6226-1).
- 326 Kexin Huang, Cao Xiao, Lucas Glass, and Jimeng Sun. Explainable substructure partition finger-
327 print for protein, drug, and more. *NeurIPS Learning Meaningful Representation of Life Workshop*,
328 2019.
- 329 Likun Jiang, Changzhi Jiang, Xinyu Yu, Rao Fu, Shuting Jin, and Xiangrong Liu. DeepTTA: a
330 transformer-based model for predicting cancer drug response. *Briefings in Bioinformatics*, 23(3):
331 bbac100, May 2022. ISSN 1477-4054. doi: 10.1093/bib/bbac100.

- 332 Iljung Jin and Hojung Nam. Hidra: Hierarchical network for drug response prediction with attention.
333 *Journal of Chemical Information and Modeling*, 61(8):3858–3867, 2021. doi: 10.1021/acs.jcim.
334 1c00706. URL <https://doi.org/10.1021/acs.jcim.1c00706>. PMID: 34342985.
- 335 Seonghun Kim, Seockhun Bae, Yinhua Piao, and Kyuri Jo. Graph convolutional network for drug
336 response prediction using gene expression data. *Mathematics*, 9(7):1–17, 2021. ISSN 22277390.
337 doi: 10.3390/math9070772.
- 338 Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Ben-
339 jamin A Shoemaker, Paul A Thiessen, Bo Yu, Leonid Zaslavsky, Jian Zhang, and Evan E Bolton.
340 PubChem 2023 update. *Nucleic Acids Research*, 51(D1):D1373–D1380, jan 2023. ISSN 0305-
341 1048. doi: 10.1093/nar/gkac956. URL [https://academic.oup.com/nar/article/
342 51/D1/D1373/6777787](https://academic.oup.com/nar/article/51/D1/D1373/6777787).
- 343 Brent M. Kuenzi, Jisoo Park, Samson H. Fong, Kyle S. Sanchez, John Lee, Jason F. Kreisberg,
344 Jianzhu Ma, and Trey Ideker. Predicting drug response and synergy using a deep learning model
345 of human cancer cells. *Cancer Cell*, 38(5):672–684.e6, 2020. ISSN 1535-6108. doi: [https://doi.
346 org/10.1016/j.ccell.2020.09.014](https://doi.org/10.1016/j.ccell.2020.09.014). URL [https://www.sciencedirect.com/science/
347 article/pii/S1535610820304888](https://www.sciencedirect.com/science/article/pii/S1535610820304888).
- 348 Min Li, Yake Wang, Ruiqing Zheng, Xinghua Shi, Yaohang Li, Fang-Xiang Wu, and Jianxin Wang.
349 Deepdsc: A deep learning method to predict drug sensitivity of cancer cell lines. *IEEE/ACM
350 Transactions on Computational Biology and Bioinformatics*, 18(2):575–582, 2021. doi: 10.1109/
351 TCBB.2019.2919581.
- 352 Pengfei Liu, Hongjian Li, Shuai Li, and Kwong-Sak Leung. Improving prediction of phenotypic
353 drug response on cancer cell lines using deep convolutional network. *BMC Bioinformatics*, 20
354 (1):408, July 2019. ISSN 1471-2105. doi: 10.1186/s12859-019-2910-6. URL [https://doi.
355 org/10.1186/s12859-019-2910-6](https://doi.org/10.1186/s12859-019-2910-6).
- 356 Leland McInnes, John Healy, and James Melville. UMAP: Uniform Manifold Approximation and
357 Projection for Dimension Reduction. feb 2018. URL [http://arxiv.org/abs/1802.
358 03426](http://arxiv.org/abs/1802.03426).
- 359 Hirotomo Moriwaki, Yu-Shi Tian, Norihito Kawashita, and Tatsuya Takagi. Mordred: a molec-
360 ular descriptor calculator. *Journal of Cheminformatics*, 10(1):4, February 2018. ISSN
361 1758-2946. doi: 10.1186/s13321-018-0258-y. URL [https://doi.org/10.1186/
362 s13321-018-0258-y](https://doi.org/10.1186/s13321-018-0258-y).
- 363 Tuan Nguyen, Giang T. T. Nguyen, Thin Nguyen, and Duc-Hau Le. Graph Convolutional Networks
364 for Drug Response Prediction. *IEEE/ACM transactions on computational biology and bioinform-
365 matics*, 19(1):146–154, 2022. ISSN 1557-9964. doi: 10.1109/TCBB.2021.3060430.
- 366 Ali Oskooei, Jannis Born, Matteo Manica, Vigneshwari Subramanian, Julio Sáez-Rodríguez, and
367 María Rodríguez Martínez. Pacmann: Prediction of anticancer compound sensitivity with multi-
368 modal attention-based neural networks, 2019.
- 369 Gihan Panapitiya, Michael Girard, Aaron Hollas, Jonathan Sepulveda, Vijayakumar Murugesan,
370 Wei Wang, and Emily Saldanha. Evaluation of deep learning architectures for aqueous solubility
371 prediction. *ACS Omega*, 7(18):15695–15710, 2022. doi: 10.1021/acsomega.2c00642. URL
372 <https://doi.org/10.1021/acsomega.2c00642>.
- 373 Aron Park, Yeeun Lee, and Seungyoon Nam. A performance evaluation of drug response prediction
374 models for individual drugs. *Scientific Reports*, 13:11911, July 2023. ISSN 2045-2322. doi: 10.
375 1038/s41598-023-39179-2. URL [https://www.ncbi.nlm.nih.gov/pmc/articles/
376 PMC10366128/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10366128/).
- 377 Alexander Partin, Thomas S. Brettin, Yitan Zhu, Oleksandr Narykov, Austin Clyde, Jamie Over-
378 beek, and Rick L. Stevens. Deep learning methods for drug response prediction in cancer: Pre-
379 dominant and emerging trends. *Frontiers in Medicine*, 10:1086097, February 2023. ISSN 2296-
380 858X. doi: 10.3389/fmed.2023.1086097. URL [https://www.ncbi.nlm.nih.gov/pmc/
381 articles/PMC9975164/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9975164/).

- 382 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
383 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
384 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
385 12:2825–2830, 2011.
- 386 Jihye Shin, Yinhua Piao, Dongmin Bang, Sun Kim, and Kyuri Jo. DRPreter: Interpretable Anti-
387 cancer Drug Response Prediction Using Knowledge-Guided Graph Neural Networks and Trans-
388 former. *International Journal of Molecular Sciences*, 23(22):13919, nov 2022. ISSN 1422-
389 0067. doi: 10.3390/ijms232213919. URL [https://www.mdpi.com/1422-0067/23/
390 22/13919](https://www.mdpi.com/1422-0067/23/22/13919).
- 391 Yi-Ching Tang and Assaf Gottlieb. Explainable drug sensitivity prediction through can-
392 cer pathway enrichment. *Scientific Reports*, 11(1):3128, February 2021. ISSN 2045-
393 2322. doi: 10.1038/s41598-021-82612-7. URL [https://www.nature.com/articles/
394 s41598-021-82612-7](https://www.nature.com/articles/s41598-021-82612-7). Number: 1 Publisher: Nature Publishing Group.
- 395 Justin M. Wozniak, Timothy G. Armstrong, Michael Wilde, Daniel S. Katz, Ewing Lusk, and Ian T.
396 Foster. Swift/T: Scalable data flow programming for distributed-memory task-parallel applica-
397 tions. In *Proc. CCGrid*, 2013.
- 398 Justin M. Wozniak, Timothy G. Armstrong, Ketan C. Maheshwari, Daniel S. Katz, Michael Wilde,
399 and Ian T. Foster. Interlanguage parallel scripting for distributed-memory scientific computing.
400 In *Proc. WORKS @ SC*, 2015.
- 401 Justin M. Wozniak, Rajeev Jain, Prasanna Balaprakash, Jonathan Ozik, Nicholson Collier, John
402 Bauer, Fangfang Xia, Thomas Brettin, Rick Stevens, Jamaludin Mohd-Yusof, Cristina Garcia
403 Cardona, Brian Van Essen, and Matthew Baughman. CANDLE/Supervisor: A workflow frame-
404 work for machine learning applied to cancer research. *BMC Bioinformatics*, 19(18):491, 2018.
405 ISSN 1471-2105. doi: 10.1186/s12859-018-2508-4. URL [https://doi.org/10.1186/
406 s12859-018-2508-4](https://doi.org/10.1186/s12859-018-2508-4).
- 407 Molly Yancovitz, Adam Litterman, Joanne Yoon, Elise Ng, Richard L. Shapiro, Russell S.
408 Berman, Anna C. Pavlick, Farbod Darvishian, Paul Christos, Madhu Mazumdar, Iman Os-
409 man, and David Polsky. Intra- and Inter-Tumor Heterogeneity of BRAFV600EMutations in
410 Primary and Metastatic Melanoma. *PLoS ONE*, 7(1):e29336, January 2012. ISSN 1932-
411 6203. doi: 10.1371/journal.pone.0029336. URL [https://www.ncbi.nlm.nih.gov/
412 pmc/articles/PMC3250426/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3250426/).
- 413 Yitan Zhu, Thomas Brettin, Fangfang Xia, Alexander Partin, Maulik Shukla, Hyunseung Yoo,
414 Yvonne A. Evrard, James H. Doroshov, and Rick L. Stevens. Converting tabular data into im-
415 ages for deep learning with convolutional neural networks. *Scientific Reports*, 11(1):11325, May
416 2021. ISSN 2045-2322. doi: 10.1038/s41598-021-90923-y. URL [https://www.nature.
417 com/articles/s41598-021-90923-y](https://www.nature.com/articles/s41598-021-90923-y). Number: 1 Publisher: Nature Publishing Group.
- 418 Zhaorui Zuo, Penglei Wang, Xiaowei Chen, Li Tian, Hui Ge, and Dahong Qian. SWnet:
419 a deep learning model for drug response prediction from cancer genomic signatures and
420 compound chemical structures. *BMC Bioinformatics*, 22(1):434, September 2021. ISSN
421 1471-2105. doi: 10.1186/s12859-021-04352-9. URL [https://doi.org/10.1186/
422 s12859-021-04352-9](https://doi.org/10.1186/s12859-021-04352-9).

423 A APPENDIX

424 A.1 MODEL RANKINGD FOR nHBDON DOMAIN.

425 How the model rankings change at different domains of nHBDon (number of hydrogen bond donors)
426 space.

Table 3: Rankings of models for each nHBDon descriptor domain.

nHBDon	1	2	3	4	5	6	7	8
2.350000	Descriptor	Morgan	DeepTTC	SMILES	Graph	HiDRA	GraphDRP	ExtraTrees
7.050000	Descriptor	Morgan	Graph	SMILES	DeepTTC	GraphDRP	HiDRA	ExtraTrees
11.750000	GraphDRP	Graph	HiDRA	Morgan	DeepTTC	SMILES	ExtraTrees	Descriptor
16.450000	GraphDRP	Graph	HiDRA	Morgan	DeepTTC	SMILES	ExtraTrees	Descriptor
21.150000	Descriptor	Morgan	SMILES	Graph	DeepTTC	HiDRA	GraphDRP	ExtraTrees
25.850000	Descriptor	Morgan	SMILES	Graph	ExtraTrees	GraphDRP	HiDRA	DeepTTC
30.550000	DeepTTC	ExtraTrees	Morgan	Descriptor	SMILES	Graph	HiDRA	GraphDRP
35.250000	ExtraTrees	DeepTTC	Morgan	SMILES	Descriptor	Graph	HiDRA	GraphDRP
39.950000	DeepTTC	ExtraTrees	Morgan	Descriptor	GraphDRP	SMILES	Graph	HiDRA
44.650000	DeepTTC	Morgan	SMILES	ExtraTrees	Descriptor	Graph	HiDRA	GraphDRP

Table 4: Average RMSE values for DeepTTC and Descriptor models in two regions in the nHBDon space.

nHBDon	DeepTTC	Descriptor
< 8	0.0587 ± 0.0004	0.0513 ± 0.0015
> 35	0.0336 ± 0.0044	0.0385 ± 0.0006

427 A.2 SUPERVISOR FRAMEWORK

428 SUPERVISOR FRAMEWORK SCALABILITY

429 Supervisor was designed as an Exascale Computing Project application, meaning it was designed
430 from the beginning for exascale computers. Supervisor is built around the Swift/T Wozniak et al.
431 (2013); Armstrong et al. (2014) workflow language and runtime. Swift/T is an MPI-based work-
432 flow system, so communication for task distribution and monitoring is performed over MPI Forum
433 (1994), the messaging layer provided by machine vendors for efficient use of large-scale computers.
434 Swift/T is scalable through two architectural innovations. First, the task distribution is coordinated
435 by a network of multiple task servers, which use an efficient work-stealing algorithm to distribute
436 work. Secondly, the control logic itself generated from the user workflow script is evaluated over
437 this distributed fabric, meaning that the workflow evaluation itself is also scalable.

438 SUPERVISOR USABILITY FOR DEEP LEARNING WORKFLOWS

439 Supervisor has many usability features for deep learning applications. These are based on features
440 of the Swift/T language and the Supervisor scripts that wrap the core workflow features with easier
441 to use scripts for launching workflows. For example, Swift/T contains multiple mechanisms for
442 calling back into user code through Python interfaces, command lines, and other languages like Tcl
443 and R Wozniak et al. (2015). Supervisor is launched with the `supervisor` tool, which accepts
444 a workflow name, site specification, and configuration file. The workflow name is essentially a
445 label to the workflow, such as `"CMP-CV"` for the present case. The site specification contains
446 settings for the computing system in question, such as program locations for Swift/T, Python, etc.
447 The configuration file contains any additional settings, such as scheduler items including walltime,
448 resources to allocate, parameters for a workflow control algorithm in use, etc.

449 Internally, Supervisor contains scripts to glue the workflow system to user models through the
450 "model shell." For the CMP-CV case, this script sets up the container for execution, handles the
451 hyperparameters, finds and runs the container with its standard command line, and collects results.
452 Everything here is logged into a per-model log for human examination and possible debugging later.

453 A.3 GRAPH, SMILES, MORGAN AND DESCRIPTOR MODELS

454 This section contains details on the **Graph**, **SMILES**, **Morgan** and **Descriptor** models introduced
455 in the Section 2.3. These four models use the same cell-line representation but different drug rep-
456 resentations. The cell-line representation is created by feeding the 1007 gene expression values to
457 a fully connected neural network. The drug representation of the **Graph** model is created using a

458 graph neural network Panapitiya et al. (2022) consisting of graph convolutional layers. In the **Mor-**
459 **gan** model, a drug molecule is represented using a Morgan fingerprint in the form of a bit vector of
460 size 1024. RDKit⁵ is used for this fingerprint generation. The drug representation of the **SMILES**
461 model is created as it is done in the DeepTTC⁶Jiang et al. (2022) model. The drug representation
462 of the **Descriptor** model is initialized using 783 molecular descriptors generated using the Mordred
463 package Moriwaki et al. (2018). These descriptors are fed into a fully connected neural network to
464 create the final drug representation.

465 A.4 UNRAVELING THE ROLE OF MOLECULAR STRUCTURE IN DRUG ERROR

466 By learning from the feature domain, we explored the potential relationships between the drug
467 structures and their corresponding features. In Figure 9, a visualization technique is employed to
468 embed the circular Morgan fingerprint representations of the drugs, utilizing UMAP (McInnes et al.,
469 2018). This method allows for the reduction of high-dimensional (2048-bit) fingerprint vectors into
470 an accessible two-dimensional representation. Subsequently, the desired descriptor values were
471 overlaid utilizing a color spectrum.

472 Upon close scrutiny of the four plotted graphs covering GATS1Z, SlogP_VSA4, C3SP3, and JGI2,
473 an identification of the regions of chemical space encompassed by the data is unveiled. This visu-
474 alization serves as a tool, highlighting the specific regions of space each descriptor predominantly
475 occupies, offering an insightful glance into the diverse chemical territories. From this figure we see
476 that there is close clustering for the $50 < S \log P_VSA4 < 55$ and $JGI2 < 0.04$ molecules, high-
477 lighted with red X's. This infers that the high error drugs in these descriptor domains exhibit similar
478 structural motifs that possibly contribute to the error. The opposite is also true where the descriptors
479 $GATS1Z < 0.2$ and $C3SP3 > 9$ show sparser data points. This points to these descriptors being less
480 correlated with certain similar structural motifs. This embedding offers yet another way to utilize
481 the results gathered above to draw conclusions about a model's weaknesses and strengths. A closer
482 look at examples of these structures can be found in Figure 10 and Figure 11.

⁵RDKit

⁶DeepTTC

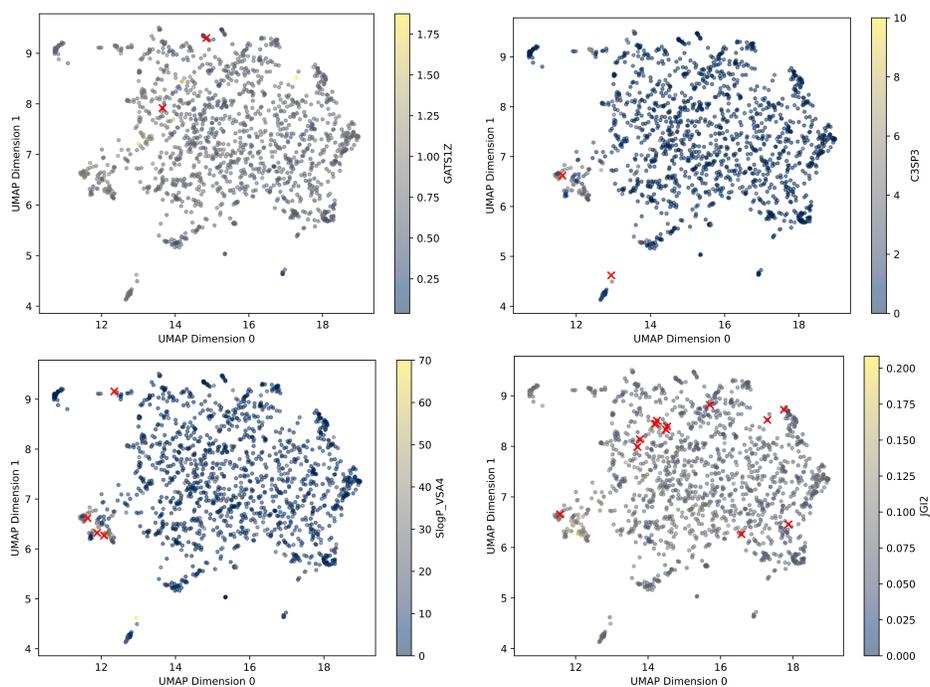


Figure 9: UMAP embedding of molecular fingerprints with overlay of molecular features of interest: GATS1Z, C3SP3, SlogP_VSA4, and JGI2. The color of each point corresponds with its associated value and the red X's highlight the molecules identified as having the highest error from Figure 5.

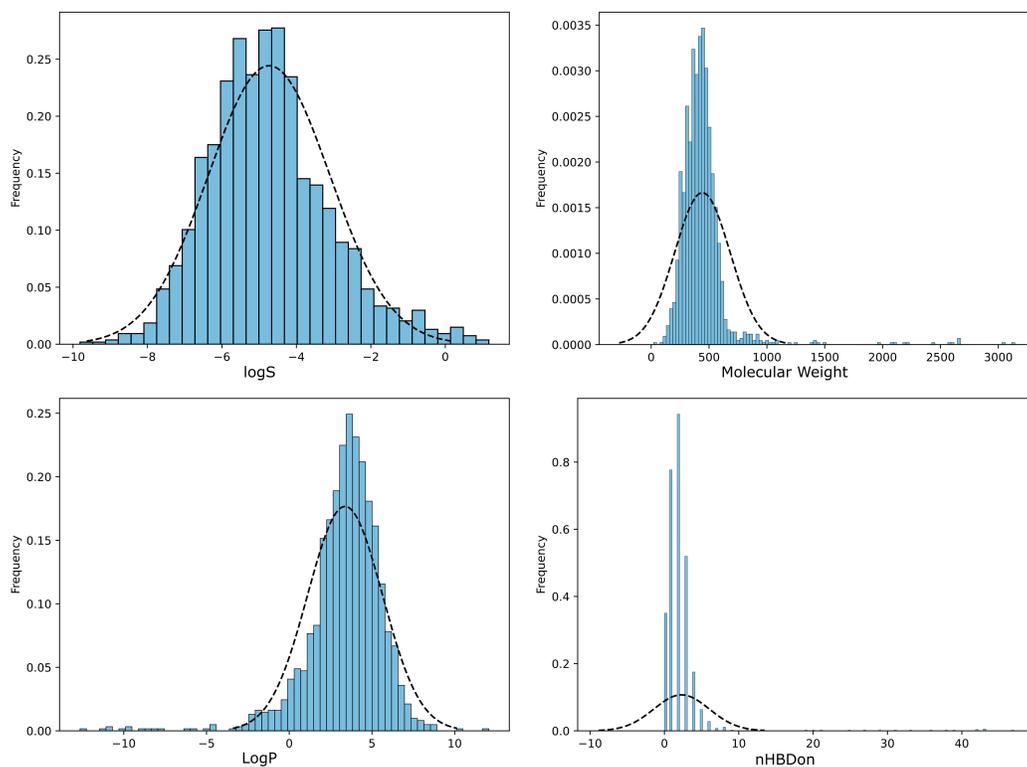


Figure 10: Distributions of drug like properties over the used dataset. Covers solubility (logS), Molecular Weight, the partition coefficient (LogP), and number of Hydrogen donors (nHBDon).

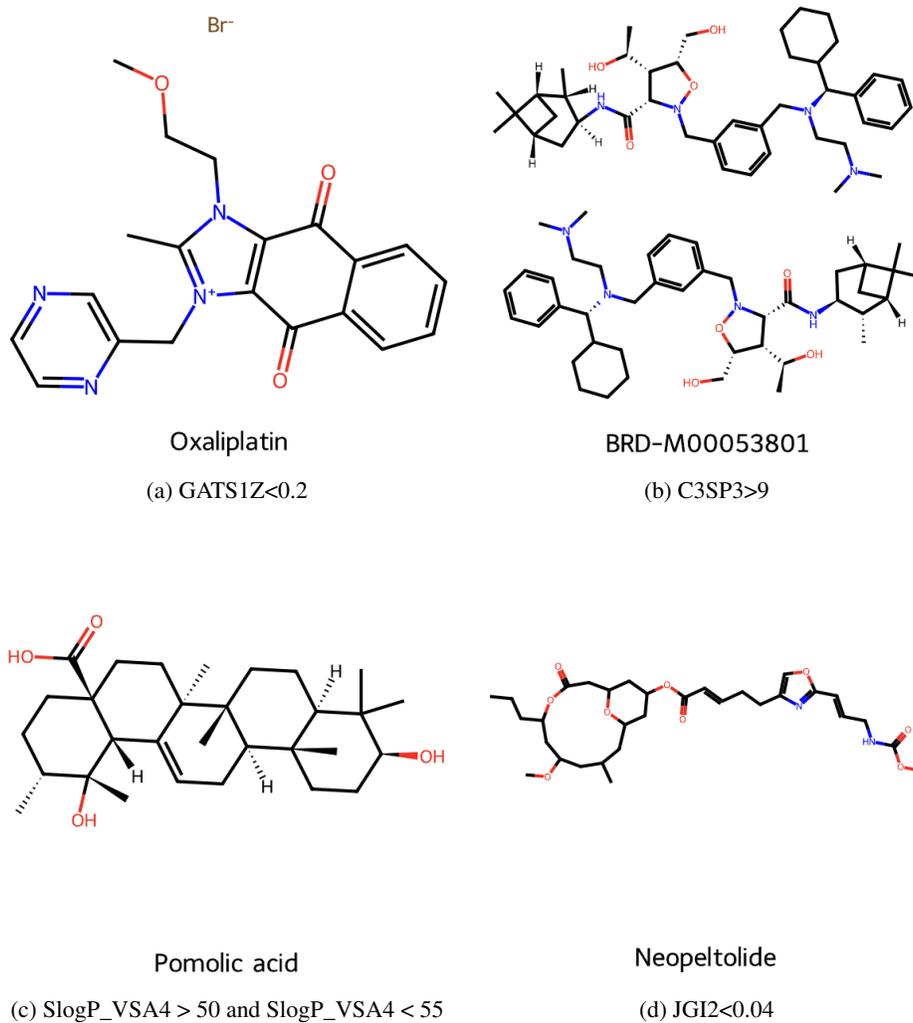


Figure 11: Example drug molecules in GATS1Z, C3SP3, SlogP_VSA4 and JGI2 domains.