**Interpretable deep learning architectures for improving drug response prediction performance: myth or reality?**

Yihui Li[1], David Earl Hostallero[1,2], and Amin Emad[1,2,3,*]

[1] Department of Electrical and Computer Engineering, McGill University, Montreal, QC, Canada

[2] Mila, Quebec AI Institute, Montreal, QC, Canada

[3] The Rosalind and Morris Goodman Cancer Institute, Montreal, QC, Canada

* Corresponding Author:

Amin Emad

755 McConnell Engineering Building

3480 University Street, Montreal, Quebec, Canada, H3A 0E9

Email: amin.emad@mcgill.ca

1

1  **Abstract**

2  Motivation: Recent advances in deep learning model development have enabled more accurate

3  prediction of drug response in cancer. However, the black-box nature of these models still

4  remains a hurdle in their adoption for precision cancer medicine. Recent efforts have focused on

5  making these models interpretable by incorporating signaling pathway information in model

6  architecture. While these models improve interpretability, it is unclear whether this higher

7  interpretability comes at the cost of less accurate predictions, or a prediction improvement can

8  also be obtained. Results: In this study, we comprehensively and systematically assessed four

9  state-of-the-art interpretable models developed for drug response prediction to answer this

10  question using three pathway collections. Our results showed that models that explicitly

11  incorporate pathway information in the form of a latent layer perform worse compared to

12  models that incorporate this information implicitly. Moreover, in most evaluation setups the best

13  performance is achieved using a simple black-box model. In addition, replacing the signaling

14  pathways with randomly generated pathways shows a comparable performance for the majority

15  of these interpretable models. Our results suggest that new interpretable models are necessary

16  to improve the drug response prediction performance. In addition, the current study provides

17  different baseline models and evaluation setups necessary for such new models to demonstrate

18  their superior prediction performance. Availability and Implementation: Implementation of all

19  methods are provided in https://github.com/Emad-COMBINE-lab/InterpretableAI_for_DRP.

20  Generated uniform datasets are in https://zenodo.org/record/7101665#.YzS79HbMKUk.

21  Contact: amin.emad@mcgill.ca

22  Supplementary Information: Online-only supplementary data is available at the journal's website.

23 **Introduction**

24 Machine learning models have found various applications in medicine, including drug

25 repositioning (Jarada, et al., 2020), drug discovery (Vamathevan, et al., 2019), gene prioritization

26 (Emad, et al., 2017; Zhang, et al., 2021), and drug response prediction (Adam, et al., 2020;

27 Ballester, et al., 2022; Costello, et al., 2014; Huang, et al., 2020). Models for drug response

28 prediction (DRP) are typically trained using various data modalities such as molecular 'omics'

29 profiles of samples (e.g., cancer cell lines or tumors), drug representations, and network

30 information (Adam, et al., 2020; Ballester, et al., 2022; Guvenc Paltun, et al., 2021). In recent

31 years, various models have been proposed using deep learning (DL) for drug response prediction

32 (Baptista, et al., 2021; Chen and Zhang, 2022; El Khili, et al., 2022; Hostallero, et al., 2022;

33 Hostallero, et al., 2021). In spite of their success in their perspective tasks, most DL models are

34 considered as "black-boxes" with inner operations that are difficult to interpret. This

35 characteristic of DL models is undesirable for applications in the biomedical field, as identifying

36 the set of biological features that contribute to the model prediction outputs and understanding

37 the relationship between these features are crucial when conducting further experimental

38 studies to validate these computational findings. To address these challenges, the concept of

39 interpretable artificial intelligence (Azodi, et al., 2020; Barredo Arrieta, et al., 2020; Malioutov, et

40 al., 2017) has been introduced to create models that can achieve both high performance and

41 interpretability.

42

43 In the context of DRP, model interpretability can be achieved in two ways: 1) using post-hoc

44 analysis to determine feature attributions and identify important features without explicitly

45    incorporating prior knowledge in model architecture, and 2) integrating prior knowledge (e.g.,

46    signaling pathways) to add meaningful structure to the model, which can then be interpreted (for

47    example using post-hoc feature importance methods). While we and others have successfully

48    used the former strategy in DRP (Hostallero, et al., 2022; Hostallero, et al., 2021) and other

49    applications (Caruana, et al., 2015; Che, et al., 2016), the latter strategy can potentially allow the

50    interpretability to go one step further to provide systems biology insights regarding the

51    mechanisms involved in response to drug treatments. Incorporating prior information such as

52    biological pathway and subsystem information allows the model embeddings to reflect

53    subsystem activities and state changes, which can then be computationally or experimentally

54    investigated to reveal different biological mechanisms that confer specific drug sensitivities

55    (Kuenzi, et al., 2020). In fact, post-hoc feature importance analysis can be incorporated in these

56    models to identify not only important input features, but also embeddings that reflect crucial

57    subsystems for cellular response to a particular drug.

58

59    The models that incorporate pathway information have generated valuable insights regarding

60    drugs' mechanisms of action and gene-pathway relationships, some of which have been validated

61    experimentally (Kuenzi, et al., 2020). However, there have been conflicting reports on their ability

62    in providing accurate drug response predictions (Deng, et al., 2020; Jin and Nam, 2021; Kuenzi,

63    et al., 2020; Snow, et al., 2021; Tang and Gottlieb, 2021; Zhang, et al., 2021). Ideally,

64    interpretability should not come at the expense of prediction performance, since a lower

65    prediction performance of interpretable models may reflect that the black-box models are better

66    capable at extracting patterns of the data and incorporating informative signals that are not being

67   utilized by the more interpretable models. For example, consider a hypothetical model that is

68   completely interpretable, but generates random drug response predictions that do not reflect

69   the measured drug responses of samples. No matter how interpretable this model may be, the

70   insights obtained from it is not going to reflect the biological and chemical mechanisms involved

71   in drug response.

72

73   Recognizing the intertwined relationship between interpretability and performance, the majority

74   of recent models that incorporate pathway information for better interpretability have also

75   sought and reported an improved prediction performance (Deng, et al., 2020; Jin and Nam, 2021;

76   Snow, et al., 2021; Tang and Gottlieb, 2021; Zhang, et al., 2021). On the other hand, some studies

77   have reported comparable or slightly worse model performance after incorporating pathway

78   information (Kuenzi, et al., 2020). However, it is rather difficult to gauge the (potential)

79   contribution of pathway information in DRP performance from the original studies, due to

80   differences between data used in each study, their evaluation setup, and in many cases a lack of

81   appropriate baseline models to act as control. To investigate these inconsistent findings in state-

82   of-the-art models, we conducted a study that comprehensively evaluates the effect of pathway

83   incorporation on performance of DRP models and aims to answer five main questions:

   1. Does the inclusion of biological pathway information improve model performance when

85      evaluated strictly and comprehensively?

   2. Which type of pathway incorporation strategy is best capable of improving the

87      performance?

88    3. Are interpretable models better suited for prediction of response of unseen cell lines or

89         unseen drugs?

90    4. Can the performance of the interpretable models be attributed to biological information

91         present in the pathway datasets, or a similar improvement can be also achieved through

92         the use of randomly generated pathways, reflecting a technical (instead of a biological)

93         origin for the performance?

94    5. What pathway database is most helpful in improving model performance?

95

96    To answer the proposed questions, we performed 189 experiments evaluating 21 computational

97    models with three pathway collections (Kanehisa and Goto, 2000), (Schaefer, et al., 2009),

98    (Fabregat, et al., 2017) and under three data splitting strategies. The models included four state-

99    of-the-art interpretable DL architectures that incorporate pathway information (Deng, et al.,

100   2020; Jin and Nam, 2021; Tang and Gottlieb, 2021; Zhang, et al., 2021) (henceforth pathway-

101   based models) and four variants of them, as well as thirteen baseline models that can evaluate

102   the performance of these models from different angles (discussed in Methods). We selected

103   these interpretable models since they use similar type of information for cancer cell lines (CCLs)

104   and drugs and utilize gene-pathway membership in their architectures, allowing us to compare

105   them fairly and comprehensively. Moreover, they represent two categories of strategies to

106   incorporate pathway information in DL architectures: methods that use a pathway layer

107   connecting genes to pathway nodes (*explicit* models) such as PathDNN (Deng, et al., 2020) and

108   ConsDeepSignaling (CDS) (Zhang, et al., 2021), and those that do not directly use a pathway layer

109   (*implicit* models) such as HiDRA (Jin and Nam, 2021), and PathDSP (Tang and Gottlieb, 2021).

110

111    Our baseline models included a traditional machine learning model (random forests), a black-box

112    fully connected neural network with a similar architecture to those of the interpretable models,

113    as well as "naive" predictors and "random-pathway" predictors, two important baselines that

114    have been largely overlooked in previous studies. The naive predictor uses the average drug

115    response of samples in the training set and reports that for each testing sample. This baseline is

116    particularly important in controlling for inflation of prediction performance due to distinct range

117    of log IC50 (natural log of the half maximal inhibitory concentration, a drug sensitivity measure)

118    of different drugs. In other words, it is possible to obtain a good approximation of drug response

119    by simply knowing the identity of the drug, resulting in artificially inflated performance metrics.

120    Each random-pathway predictor exactly matches the architecture and pipeline of an

121    interpretable model, but randomly assigns genes to pathways, while preserving the size of each

122    pathway. These baselines allow us to determine whether potential performance improvement

123    of an interpretable model is truly due to the added value of the biological information, or instead

124    is a technical artifact of modifying the model architecture.

125

126    Our analysis showed that overall, incorporating pathway information *does not* lead to improved

127    prediction performance, confirming the observations reported by Kuenzi et al. (Kuenzi, et al.,

128    2020) for their proposed model. In particular, in many cases a simple black-box multilayer

129    perceptron (MLP) achieves the best performance. Moreover, even in instances that performance

130    improvement compared to an MLP or a naive predictor was observed, a similar performance was

131    achieved using randomly generated pathways. This suggests that such improvements should not

132     be attributed to the biological information carried by pathway collections and is likely a technical

133     artifact. We also observed that the strategy used to include pathway information in the models

134     has a significant influence on the performance, and explicit models seem to perform worse

135     compared to implicit models. Finally, Reactome pathways seemed to provide slightly better

136     predictions compared to other pathway collections.

137

138     **Methods**

139     **Data preprocessing and uniform dataset formation**

140     To form uniform datasets for our analyses, we first evaluated different data modalities and

141     datasets used by each of the pathway-based models (Supplementary Table S1). In these studies,

142     gene expression (GEx), somatic mutation (Mut), and copy number variation (CNV) of samples

143     were used, while for drugs their targets (T) or their Morgan fingerprints (FP) capturing their

144     chemical structure were used. In order to maintain fairness and consistency of model

145     performance comparisons, for each choice of pathway collection we compiled a uniform dataset

146     that was used by all models evaluated in this study (three uniform datasets in total). These

147     datasets are freely available in https://zenodo.org/record/7101665#.YzS79HbMKUk.

148

149     We collected GEx, Mut, CNV, and drug sensitivity data (in the form of log IC50) of 959 cancer cell

150     lines (CCLs) from Genomics of Drug Sensitivity in Cancer (GDSC) (Yang, et al., 2013) database. We

151     obtained drug target information from STITCH (Szklarczyk, et al., 2016) and drug structural data

152     from PubChem (Kim, et al., 2021). Protein-protein interactions (PPI) that were used by one of the

153     models were obtained from the STRING database (Szklarczyk, et al., 2019) (only experimental
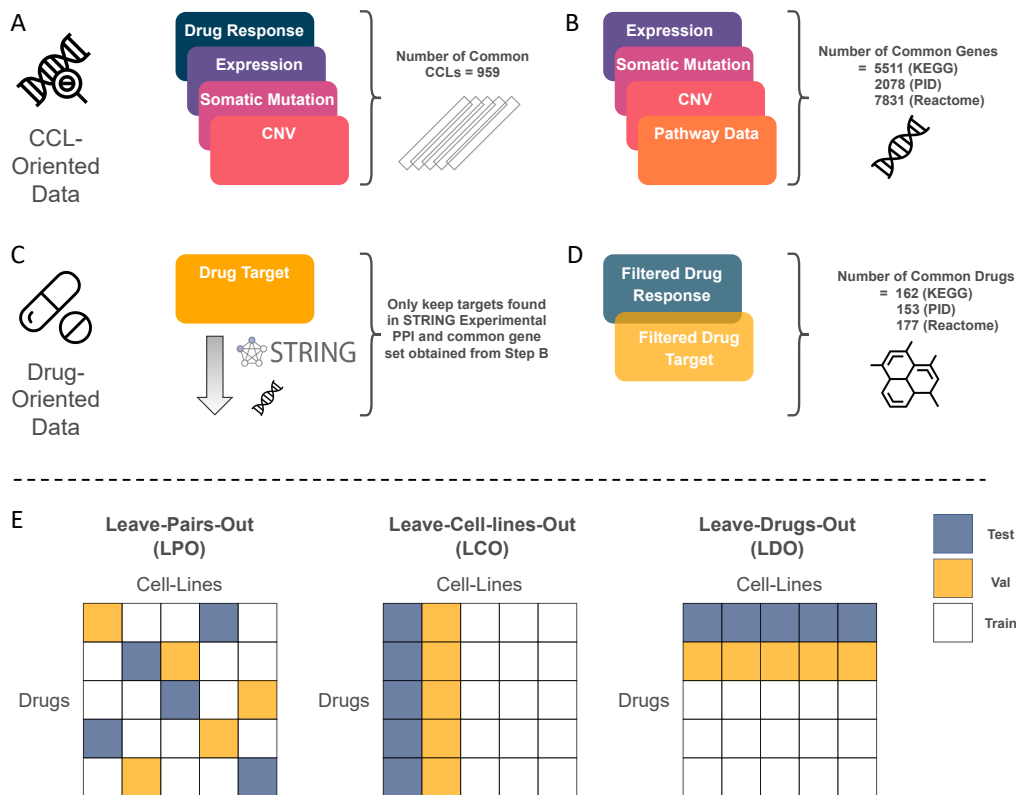
8

154    PPIs were used). Finally, gene-pathway membership information was obtained from KEGG (Kyoto

155    Encyclopedia of Genes and Genomes) (Kanehisa and Goto, 2000), PID (Pathway Interaction

156    Database) (Schaefer, et al., 2009), and Reactome (Fabregat, et al., 2017). Supplementary Table

157    S2 outlines the data used in this study and their sources. We obtained drug response data in the

158    form of log IC50 values and removed duplicate drugs that came from different experimental

159    batches. In such cases, we kept the drug whose response was measured across a larger number

160    of CCLs. We collected drug InChI (International Chemical Identifier) strings (Heller, et al., 2015)

161    from PubChem and used the RDKit (Landrum, 2006) software to generate 512-bit Morgan

162    fingerprints for these drugs. We obtained drug target data from the STITCH database, where we

163    only kept drug targets with confidence score larger than 800 (out of 1000) and coming from the

164    "experimental" and "database" channels.

165

166    **Table 1:** Summary of pathway-specific uniform datasets.

| Pathway Database | Num. CCLs | Num. Drugs | Num. Unique Drug Targets | Num. (Drug, CCL) Pairs | Num. Pathways | Num. Unique Genes |
|---|---|---|---|---|---|---|
| KEGG | 959 | 162 | 446 | 118,896 | 332 | 5511 |
| PID | 959 | 153 | 321 | 112,781 | 196 | 2078 |
| Reactome | 959 | 177 | 493 | 128,324 | 1608 | 7831 |

167

168    We performed log2(FPKM+1) normalization on the GEx data and removed genes whose

169    expression showed low variability across different CCLs (standard deviation < 0.1). We also

170    removed genes for which there were no somatic mutations, CNV, pathway information, drug

171    target data, and STRING Experimental PPI information. This formed our common gene set (Figure

172    1A and 1B). In parallel, drug targets that were not present in the common gene set above or in

173    the PPI network were excluded. Only drugs that had both log IC50 measurements and drug

174 targets were kept in the final uniform datasets (Figure 1C and 1D). The PPI network was involved

175 in the data preprocessing step as PathDSP (Tang and Gottlieb, 2021) incorporated it to perform

176 pathway enrichment analysis. Since we needed the uniform datasets to be usable by all models,

177 we included this step in the pre-processing procedure.



178

179 **Figure 1:** Construction of pathway-specific uniform datasets and data splitting approaches. (A)
180 cancer cell lines (CCLs) with available data for drug response, gene expression, somatic mutation,
181 and copy number variation (CNV) were selected. (B) Genes shared between different sources of
182 data were identified. Genes that were not present in any pathway were removed. (C) Drug target
183 genes that were not found in the common gene set obtained from Step B and the STRING
184 Experimental protein-protein interaction (PPI) network were removed. (D) Drugs and small
185 molecules that had measured log IC50 values and drug target information were selected. E)
186 Model input data was split into five folds, with the training, validation, and test set ratio of 3: 1:
187 1. Folds in the leave-pairs-out (LPO) validation scheme are formed by randomly selecting
188 mutually exclusive (CCL, drug) pairs, whereas in leave-cell-lines-out (LCO) and leave-drugs-out
189 (LDO) validation schemes, mutually exclusive cell lines and drugs are randomly selected,
190 respectively.

191

192   Figure 1 illustrates the process of constructing the pathway-specific uniform datasets. Since each

193   source of pathway collection contained different number of genes, the final dataset for each

194   collection was slightly different. Table 1 summarizes the number of CCLs, drugs, genes, and

195   pathways for each pathway collection in the uniform dataset, while Supplementary Table S3

196   provides details about CCLs and drugs.

197

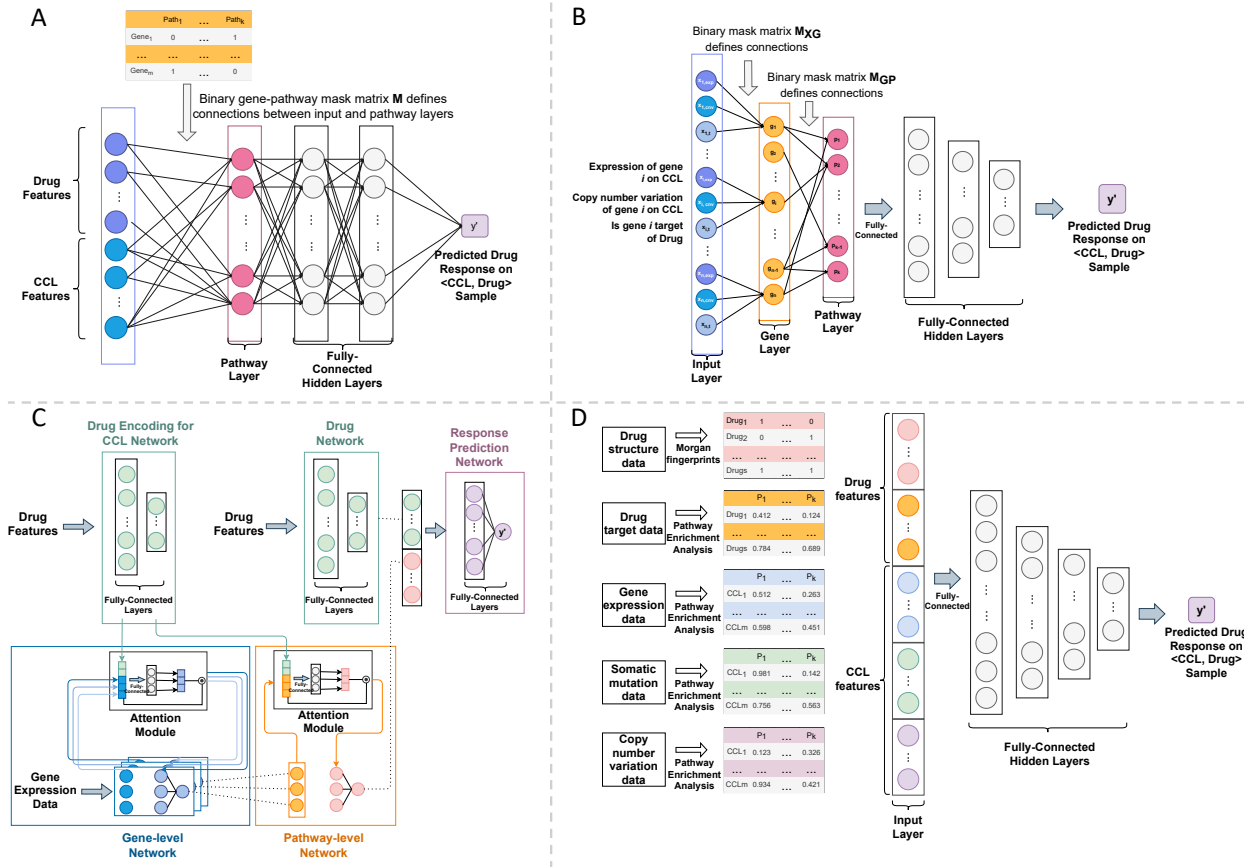198   **Model evaluation and data split**

199   We split our data randomly into five disjoint folds, where the training, validation, and test ratio

200   was 3: 1: 1. The validation set was used for hyperparameter tuning and the test set was used for

201   final model evaluation. The details of hyperparameter tuning, model training, and final

202   architectures are provided in Supplementary File S2. We adopted three data splitting methods

203   (validation schemes) to generate these folds: leave-pairs-out (LPO), leave-cell-lines-out (LCO),

204   and leave-drugs-out (LDO), as depicted in Figure 1E. These three strategists were adopted to

205   comprehensively assess the models for different drug response prediction tasks (for unseen (CCL,

206   drug) pairs, unseen CCLs, and unseen drugs, respectively), and to determine in which one of these

207   tasks (if any) pathway incorporation improves model prediction performance. To ensure fairness,

208   same folds were used for all models.

209

210   We evaluated the performance of each model using two main performance measures:

211   Spearman's correlation coefficient (SCC) and root mean squared error (RMSE), but various other

212   measures are also reported in supplementary tables. First, for a fixed CCL, the predicted values

213   across all drugs of the test set were compared with the measured log IC50 values to calculate a

214    CCL-specific performance measure (SCC or RMSE). Then, the mean and standard deviation of the

215    performance measure was calculated across all CCLs.

216

217    **Overview of interpretable models and their variants**

218    To study the effect of incorporating pathway information on drug response prediction, we

219    selected four pathway-based state-of-the-art models: PathDNN (Deng, et al., 2020),

220    ConsDeepSignaling (CDS) (Zhang, et al., 2021), HiDRA (Jin and Nam, 2021), and PathDSP (Tang

221    and Gottlieb, 2021). We selected these models since 1) they use similar types of information for

222    CCLs and drugs and utilize gene-pathway membership in their architectures (instead of other

223    types of prior information such as hierarchical relationships of gene ontologies), 2) they all

224    showed improved drug response prediction performance in their original studies compared to

225    their black-box counterparts, and 3) they represent two important categories of implicit and

226    explicit models (as discussed earlier). While other important models also exist (e.g., DrugCell

227    (Kuenzi, et al., 2020)), they did not satisfy the conditions above. For example, DrugCell (unlike

228    the models above) uses the hierarchical structure of gene ontologies and pathways, making it

229    rather difficult to compare against the models above in a fair manner, since it takes advantage of

230    more detailed information. Moreover, the original study of DrugCell showed that while including

231    prior information improved interpretability of their model, it did not improve the performance

232    of drug response prediction compared to its black-box counterpart. Due to the reasons above,

233    we decided to exclude it from this analysis.

234

**Figure 2:** An overview of pathway-based models considered in this study. A) PathDNN uses cancer cell line (CCL) gene expression profiles as CCL features and drug target information as drug features. The input features (genes) are connected to the pathway nodes through gene-pathway membership. The pathway layer is followed by a set of fully connected layers. B) ConsDeepSignaling (CDS) takes gene expression profile and copy number variation as CCL features and drug target information as drug features. Each node in the gene layer represents a gene and is connected to its corresponding input features in the input layer (through connection matrix $M_{XG}$). Connections between the gene and pathway layer are defined by gene-pathway membership (binary connection matrix $M_{GP}$). A set of fully connected layers follow the pathway layer. C) HiDRA has a hierarchical network architecture. It uses gene expression profiles as CCL features. Drug target information and structural data can be both used as drug features. The pathway information is incorporated using an attention module, where a small neural network is dedicated to each pathway. Pathway activation scores are calculated by the gene-level network and are concatenated with drug feature embeddings learned by the drug encoding network to generate the final input to the drug response prediction network. D) In PathDSP, drug target, gene expression, somatic mutation, and copy number variation data are processed using pathway enrichment analysis to from matrices of enrichment scores, which act as input features to the model, which is a set of fully connected layers.

255    Models with an *explicit* pathway layer (e.g., PathDNN and CDS) typically define a gene and a

256    pathway layer with connections between these layers reflecting gene-pathway membership

257    (Figure 2A-2B). The input layer of this category of models contains drug and cell line features at

258    gene level. As a result, only drug gene targets can be used with these models and Morgan

259    Fingerprint (and other structural data) is not usable without altering the model architecture. The

260    pathway layer is then connected to a group of fully connected layers to predict drug response for

261    a given sample. The inclusion of the pathway layer allows identification of important pathways

262    for a particular drug treatment or cancer type through post-hoc feature importance analysis.

263

264    Models that *implicitly* incorporate pathway information take various forms. For example, HiDRA

265    (Jin and Nam, 2021) uses a gene-level and pathway-level attention module to calculate pathway

266    importance scores, where a small-scale neural network is dedicated to each pathway by only

267    using features associated with the member genes of that specific pathway as inputs (Figure 2C).

268    On the other hand, PathDSP uses a classic fully connected feedforward architecture, but the input

269    features are pathway-enrichment scores rather than gene-level features (Figure 2D). See

270    Supplementary Files S1 and S2 for details regarding models' architectures and their training

271    procedure.

272

273    Each of the pathway-based models used different data modalities in their original study

274    (Supplementary Table S1). We tested all models using the three pathway collections discussed

275    earlier. For implicit models, we tried both drug targets (T) and Morgan fingerprints (FP); however,

276    for explicit models, only drug targets could be used due to their requirements that the drug

277     features must be at gene level. For CCL features, we used all data modalities used by the original

278     study. However, since gene expression data was used by all models (alone or in combination with

279     other omics data, Supplementary Table S1), we also implemented model variants that only

280     utilized GEx data. This ensured that one architecture is not given an unfair advantage due to

281     access to a larger number of modalities. Table 2 provides a summary of all variations of the

282     models considered in this study.

283

284     **Table 2:** List of evaluated models. △ = universal baseline, ◊ = model-specific random pathway
285     baseline, ★ = original pathway-based model, ■ = model variant, GEx = gene expression, CNV =
286     copy number variation, Mut = somatic mutation, T = drug target data, FP = Morgan fingerprint
287     (drug structural data)

| Model Name | Model Variant Name | Cell Line Features | Drug Features |
|---|---|---|---|
| Five-Layer MLP | MLP (GEx, FP) △ | GEx | FP |
| | MLP (GEx, T) △ | GEx | T |
| Naive Predictor | Naive Predictor △ | N/A | N/A |
| Random Forests | RF (GEx, FP) △ | GEx | FP |
| | RF (GEx, T) △ | GEx | T |
| PathDNN (Deng, et al., 2020) | PathDNN (GEx, T) ★ | GEx | T |
| | PathDNN_rand (GEx, T) ◊ | GEx | T |
| CDS (Zhang, et al., 2021) | CDS (GEx, CNV, T)★ | GEx, CNV | T |
| | CDS_rand (GEx, CNV, T) ◊ | GEx, CNV | T |
| | CDS (GEx, T) ■ | GEx | T |
| | CDS_rand (GEx, T) ◊ | GEx | T |
| HiDRA (Jin and Nam, 2021) | HiDRA (GEx, FP) ★ | GEx | FP |
| | HiDRA_rand (GEx, FP) ◊ | GEx | FP |
| | HiDRA (GEx, T) ■ | GEx | T |
| | HiDRA_rand (GEx, T) ◊ | GEx | T |
| PathDSP (Tang and Gottlieb, 2021) | PathDSP (GEx, CNV, MuT, FP, T) ★ | GEx, CNV, Mut | FP, T |
| | PathDSP_rand (GEx, CNV, MuT, FP, T) ◊ | GEx, CNV, Mut | FP, T |
| | PathDSP (GEx, FP) ■ | GEx | FP |
| | PathDSP_rand (GEx, FP) ◊ | GEx | FP |
| | PathDSP (GEx, T) ■ | GEx | T |
| | PathDSP_rand (GEx, T) ◊ | GEx | T |

288

289 **Baseline Models**

290 We used four types of baseline models to benchmark the pathway-based models and their

291 variants. First, we used a multilayer perceptron (MLP) with five layers as a universal baseline for

292 all models. This MLP represents a black-box feedforward neural network that is often used for

293 benchmarking of other deep learning architectures (including the pathway-based models). Since

294 all pathway-based models had a variant trained with GEx data, along with drug targets (or

295 Morgan fingerprints), we trained two MLP models, MLP (GEx, FP) and MLP (GEx, T), representing

296 the data input options above (Table 2).

297

298 The second type of baseline used in our study was a predictor that simply calculates the average

299 drug sensitivity measure of samples in the training set and reports their average for all samples

300 in the test set (henceforth referred as naive predictor). More specifically, the naive predictor does

301 not use any CCL or drug features, but instead simply relies on the identity of the CCL or the drug

302 (depending on the data splitting strategy). As shown in Supplementary Figure S1, in the LCO setup

303 and for a (CCL, drug) pair in the test set, the naive predictor reports the average response of all

304 CCLs in the training set to that drug. As a result, all CCLs in the test set will have the same response

305 value for a drug (i.e., only the drug identity determines the response). On the other hand, in the

306 LDO setup and for a (CCL, drug) pair, the average response of the CCL to all drugs in the training

307 set is reported as the prediction (i.e., only the identity of the CCL determines the response). In

308 the case of LPO, the averaging is done across all drugs and all CCLs corresponding to a (CCL, drug)

309 in the test set. The naive predictors reveal the performance of a model that does not learn the

310    relationship between the input features and output drug response and can control for inflation

311    in the performance metrics.

312

313    The third type of baselines correspond to model-specific baselines that have the exact same

314    architecture of a pathway-based model (with all their input data), but instead of gene-pathway

315    membership information from pathway databases use randomly generated pathways. This type

316    of baseline model (shown with a suffix of "_rand" in Table 2) allows us to determine if the

317    (potential) performance improvement of a pathway-based model is due to the added value of

318    biological information, or instead is a technical artifact. Let a pathway collection (e.g., KEGG)

319    contain $m$ pathways $P_i, i = 1, 2, \dots, m$, each with $N_i$ genes. Then, a randomly generated pathway

320    collection was produced by randomly assigning $N_i$ genes to pathway $P_i$. We evaluated the

321    performance of each pathway-based model with multiple randomly generated pathway

322    collections to determine the mean, standard deviation and histogram of the performance metrics

323    of these random pathway baselines.

324

325    Finally, the fourth type of baselines correspond to traditional machine learning algorithms,

326    namely random forests (RF). We trained two variations of RF, one with (GEx, T) as input and one

327    with (GEx, FP) as input.

328

329    **Cross-dataset analysis by predicting drug responses in CTRPv2 using models trained on GDSC**

330    In addition to the analysis performed using GDSC, we also assessed the generalizability of the

331    deep learning models by performing a cross-dataset analysis. Following the guidelines in a

332    previous study (Sharifi-Noghabi, et al., 2021), we trained the models using area under the dose

333    response curves (AUC) from GDSC dataset to predict AUC of drugs in CTRPv2 (Rees, et al., 2016).

334    For drugs in CTRPv2 dataset, we used their gene expression profile from the cancer cell line

335    encyclopedia (CCLE) (Barretina, et al., 2012). All models were trained using Reactome pathway

336    collection, gene expression and drug targets. Since in this dataset, the gene expressions were

337    quantified using transcript per million (TPM), we also used TPM values for the training set (GDSC).

338    Only common genes between GDSC and CCLE were included. The rest of the preprocessing steps

339    were as described earlier in the manuscript.

340

341    **Results**

342    **Models that incorporate KEGG pathway information implicitly outperformed explicit models**

343    Since KEGG was the most commonly used pathway collection in the original studies

344    (Supplementary Table S1), we used the uniform dataset that we formed for this collection to

345    comprehensively evaluate all models. We first focused on GEx data to represent CCLs since all

346    models used GEx modality in their original studies. We also used drug targets to represent

347    compounds since all models could take advantage of this data modality (Morgan fingerprints are

348    not compatible with PathDNN and CDS). Table 3 shows SCC and RMSE values for LCO, LDO, and

349    LPO data splitting strategies (see Supplementary Table S4 for other performance measures and

350    Supplementary Table S5 for statistical tests comparing these models).

351

352    PathDSP (GEx, T) outperformed all models in LCO and LPO data splitting schemes, while its

353    performance was close to MLP (GEx, T) baseline in the LDO scheme. Compared to the naive
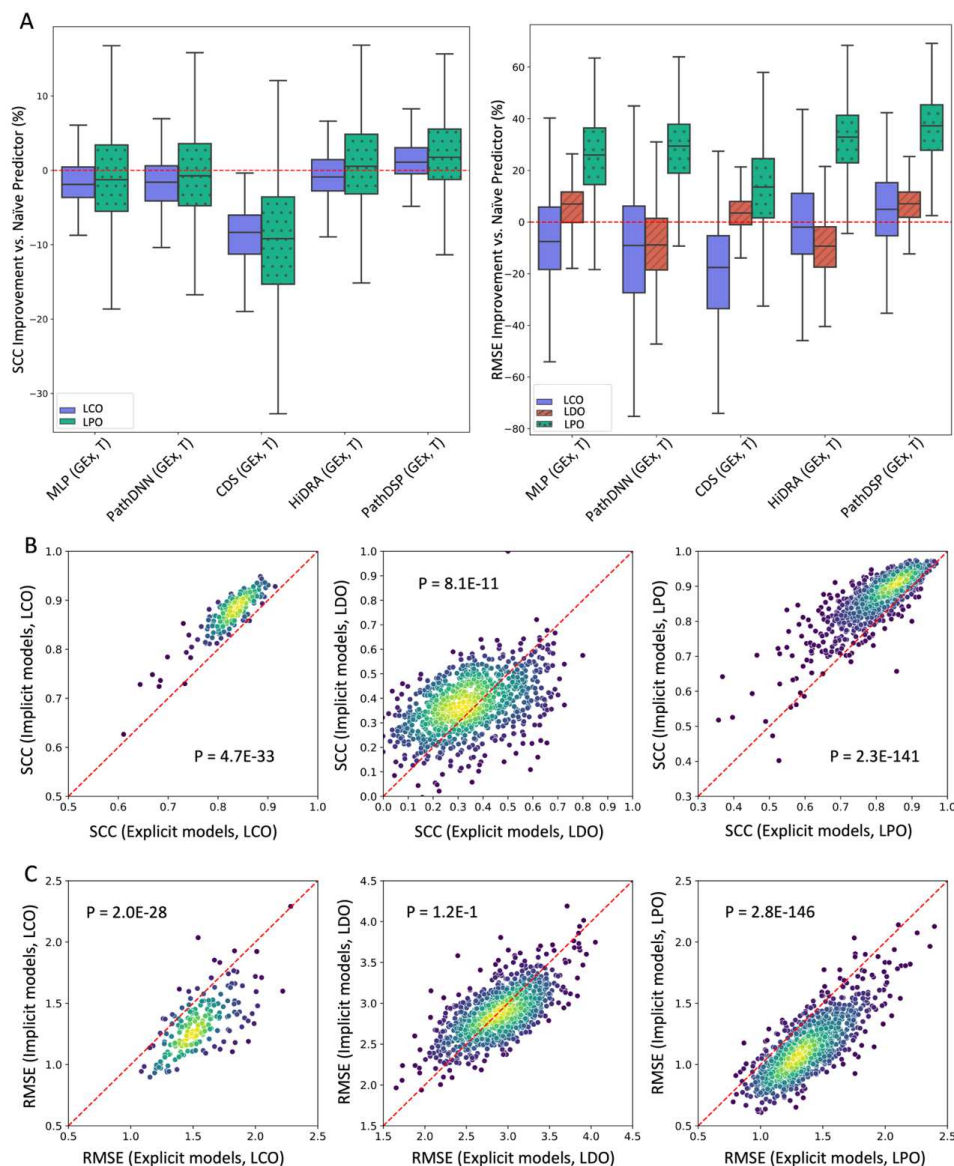
354  predictor, PathDSP (GEx, T) had a better performance in all evaluations, where the highest

355  difference was observed under the LPO validation scheme with 37% lower average RMSE. Overall,

356  the implicit models (HiDRA and PathDSP) outperformed the universal baselines (MLP and naive

357  predictor) for the majority of evaluations, while the explicit models (PathDNN and CDS) did not

358  outperform them in a considerable number of evaluations (Table 3 and Supplementary Table S4).

359

360  **Table 3:** Performance of pathway-based models using KEGG collection, with gene expression
361  (GEx) and drug targets (T) as inputs. The mean and standard deviations (std) are calculated across
362  cancer cell lines (CCLs). The best performing model is bold-faced, while worst performing model
363  is underlined. Models are ranked by their leave-cell-lines-out (LCO) RMSE. The following symbols
364  are used in this table: △ = universal baseline, ★ = original pathway-based model, ■ = model
365  variant, ↑ = higher value indicates better performance, ↓ = lower value indicates better
366  performance. ∗ The leave-drugs-out (LDO) Spearman's correlation coefficient (SCC) cannot be
367  calculated for the naive predictor since in this case it outputs the same value for all CCLs. For
368  performance of these models based on Pearson correlation coefficient, R-squared, mean squared
369  error (MSE), and concordance index, see Supplementary Table S4. Supplementary Figure S2
370  provides visualization of these values in the form of bar plots.

| Model Name | LCO | | LDO | | LPO | |
|---|---|---|---|---|---|---|
| | SCC ↑ (±std) | RMSE ↓ (±std) | SCC ↑ (±std) | RMSE ↓ (±std) | SCC ↑ (±std) | RMSE ↓ (±std) |
| PathDSP (GEx, T) ■ | **0.882 (±0.045)** | **1.283 (±0.230)** | 0.380 (±0.146) | **2.648 (±0.287)** | **0.876 (±0.074)** | **1.103 (±0.256)** |
| RF (GEx, T) △ | 0.867 (±0.042) | 1.329 (±0.205) | <u>0.342</u> (±0.175) | 2.955 (±0.474) | 0.824 (±0.096) | 1.349 (±0.301) |
| HiDRA (GEx, T) ■ | 0.864 (±0.048) | 1.368 (±0.253) | 0.356 (±0.156) | <u>3.110</u> (±0.400) | 0.863 (±0.078) | 1.174 (±0.260) |
| Naive Predictor △ | 0.871 (±0.045) | 1.373 (±0.292) | NA * | 2.826 (±0.274) | 0.855 (±0.087) | <u>1.742</u> (±0.280) |
| MLP (GEx, T) △ | 0.858 (±0.042) | 1.420 (±0.231) | **0.382 (±0.160)** | 2.686 (±0.366) | 0.845 (±0.086) | 1.294 (±0.286) |
| PathDNN (GEx, T)★ | 0.857 (±0.044) | 1.494 (±0.292) | <u>0.342</u> (±0.179) | 3.054 (±0.485) | 0.851 (±0.083) | 1.245 (±0.273) |
| CDS (GEx, T) ■ | <u>0.789</u> (±0.098) | <u>1.603</u> (±0.254) | 0.345 (±0.156) | 2.724 (±0.317) | <u>0.769</u> (±0.106) | 1.508 (±0.304) |

371 Next, we assessed the improvement provided by each deep learning method compared to the

372 corresponding naive predictor (Figure 3A). With regards to RMSE, all these models provided

373 improvement for the majority of CCLs in the LPO framework, which is expected since the

374 prediction task in LPO is significantly easier than LCO and LDO. However, in LDO and LCO, many

375 models could not provide a lower RMSE compared to the naive predictor. The improvement was

376 even less in terms of SCC (Figure 3A). However, PathDSP outperformed the naive predictor for

377 the majority of CCLs in all data splitting setups in terms of RMSE and SCC.

378

379 Next, we sought to directly compare the performance of explicit models against implicit models.

380 For this purpose, we calculated the average performance of the two implicit (PathDSP (GEx, T),

381 HiDRA (GEx T)) and the two explicit (PathDNN (GEx, T), CDS (GEx, T)) models for each CCL, and

382 used a two-sided Wilcoxon signed rank test to assess if one strategy outperforms the other

383 (Figure 3B). Based on SCC, the implicit strategy significantly outperformed the explicit strategy

384 that utilizes a pathway layer, for all three data splitting strategies. A similar pattern was observed

385 using RMSE, but for LDO strategy the difference was not statistically significant. These results

386 further confirm the observation that utilizing an explicit pathway layer does not seem to perform

387 well in prediction of drug response. Supplementary Figure S4 also shows similar scatter plots in

388 which the cancer types of cell lines are depicted, which does not suggest a cancer type-specific

389 pattern.

**Figure 3:** Performance of deep learning models in different data splitting setups. A) The improvement of each model versus naive predictor. Box plots show the distribution of performance improvement for cancer cell lines (CCLs). Each box shows the range between 25th and 75th percentiles, while whiskers show the range of the improvement (excluding outliers). See Supplementary Figure S3 in which the performance improvement of each datapoint (CCL) is also depicted. Spearman's correlation coefficient (SCC) for naive predictor cannot be calculated in leave-drugs-out (LDO). B-C) Performance of implicit pathway models versus explicit models that use a pathway layer. Each circle represents a CCL. The color of each circle represents the density of circles in its vicinity, where yellow indicates higher density and blue indicates lower density. P-values are calculated using a two-sided Wilcoxon signed-rank test. The average performance of explicit models (PathDNN and CDS) is shown on the x-axis, while the performance of implicit models (PathDSP and HiDRA) is shown on the y-axis. Panel B shows the performance in terms of SCC, while panel C shows it in terms of RMSE. See Supplementary Figure S4 in which the cancer types of CCLs are also depicted.

405

**Morgan fingerprints of compounds were more informative than drug targets for predicting**

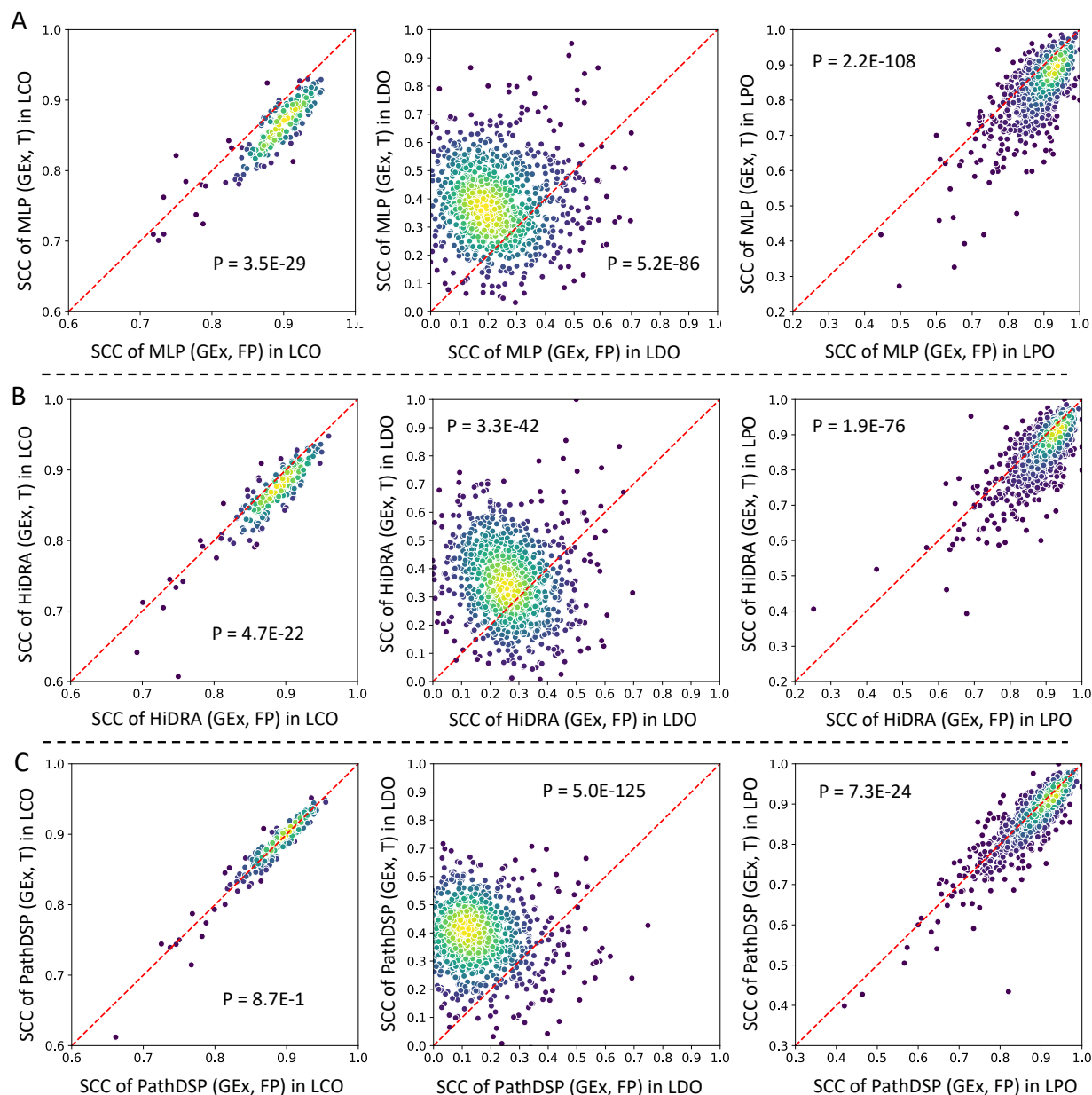**response of unseen cell lines**

408      Since three of the considered deep learning models (MLP, HiDRA, and PathDSP) can utilize both

409      drug targets (T) and Morgan fingerprints (FP) to represent drugs, we sought to determine which

410      compound representation is most informative for drug response prediction. As can be seen in

411      Figure 4, in all three models, using FP to represent compounds in most cases was superior in

412      terms of SCC in predicting unseen CCLs (LCO) or in predicting unseen CCL-drug pairs (LPO) (Two-

413      sided Wilcoxon signed-rank $P<0.05$, except for PathDSP LCO). On the other hand, in all three

414      models drug targets were more informative in predicting the response of unseen drugs (LDO).

415      However, one should note that none of the three models performed very well in the LDO data

416      splitting setup and more informative compound representations (e.g., transcriptomic changes in

417      response to compounds (El Khili, et al., 2022) or DL models that directly learn compound

418      representations (Zagidullin, et al., 2021)) may be necessary for such an application to allow

419      generalization to new compounds.

420

421      It is worth noting that although PathDSP (GEx, T) and HiDRA (GEx, T) both outperformed the MLP

422      (GEx, T) baseline in LCO and LPO evaluation, MLP (GEx, FP) baseline outperformed all other

423      models, independent of which CCL or drug representations they used in both LCO and LPO

424      (Supplementary Table S4). This is an important observation that shows that a simple MLP

425      baseline, when used with appropriate inputs could achieve comparable or better results

426      compared to various interpretable models. This observation is concordant with the observation

427    in (Kuenzi, et al., 2020), where the authors found that the interpretable version of their model

428    resulted in comparable performance to the matched black-box model. Interestingly, RF (GEx, FP)

429    provided the best performance in terms of SCC and RMSE in LCO, showing that sometimes

430    traditional machine learning methods can achieve similar or better results compared to deep

431    learning methods, an observation also made in (Chen and Zhang, 2022).

432

433    **Integrating multiple data modalities improves performance of PathDSP and CDS**

434    Among pathway-based methods that we considered in this study, two of them (CDS and PathDSP)

435    used multiple data modalities in their original study (Supplementary Table S1). Table 4 compares

436    the performance of these methods when data modalities chosen by the original study were used

437    as inputs against their performance when only GEx was used (see Supplementary Table S6 for

438    comparison of these models using two-sided Wilcoxon signed rank tests). The original PathDSP

439    model uses GEx, somatic mutation, and CNV as CCL features, as well as Morgan fingerprints and

440    drug targets as compound features, which for clarity we denote as PathDSP (GEx, CNV, MuT, FP,

441    T). PathDSP (GEx, CNV, MuT, FP, T) outperformed both PathDSP (GEx, T) and PathDSP (GEx, FP)

442    in 5 out of 6 evaluations (all except RMSE in LCO, Table 4). The original CDS model uses GEx and

443    CNV as CCL features and drug targets as compound features, which for clarity we denote as CDS

444    (GEx, CNV, T). The original CDS (GEx, CNV, T) model also outperforms CDS (GEx, T) in all

445    evaluations except for LCO approach. Overall, these results suggest that using multiple data

446    modalities can improve the performance of each model. However, it is important to remind that

447    MLP (GEx, FP) outperformed all models (including the multi-modality versions of PathDSP and

448    CDS) in LCO and LPO evaluations (Supplementary Table S4).

**Figure 4:** Performance of three models when using drug targets or Morgan fingerprints (FP) to represent drugs in terms of Spearman's correlation coefficient (SCC). Each circle represents a cancer cell line (CCL). The color of each circle represents the density of circles in its vicinity, where yellow indicates higher density and blue indicates lower density. P-values are calculated using a two-sided Wilcoxon signed-rank test. For all models, the mean SCC when using FP was higher in leave-pairs-out (LPO) and leave-cell-lines-out (LCO), and lower in leave-drugs-out (LDO) compared to when using drug targets (T). A) Performance of MLP (GEx, T) versus MLP (GEx, FP). B) Performance of HiDRA (GEx, T) versus HiDRA (GEx, FP). C) Performance of PathDSP (GEx, T) versus PathDSP (GEx, FP). Only models that could utilize both FP and T to represent drugs were used for this analysis.

**Table 4:** Performance of CDS and PathDSP using KEGG collection, with different input choices. The mean and standard deviations (std) are calculated across cancer cell lines (CCLs). For each method, the input choice that performs best is bold-faced. Since CDS can only use drug targets to represent compounds, the only considered baseline for it is CDS (GEx, T). The following symbols are used in this table: ★ = original pathway-based model, ■ = model variant, ↑ = higher value indicates better performance, ↓ = lower value indicates better performance. Supplementary Figure S5 provides visualization of these values in the form of bar plots.

| Model Name | LCO | | LDO | | LPO | |
|---|---|---|---|---|---|---|
| | SCC ↑ (±std) | RMSE ↓ (±std) | SCC ↑ (±std) | RMSE ↓ (±std) | SCC ↑ (±std) | RMSE ↓ (±std) |
| PathDSP (GEx, CNV, Mut, FP, T) ★ | **0.883** (±0.044) | 1.308 (±0.276) | **0.470** (±0.130) | **2.477** (±0.286) | **0.893** (±0.068) | **1.020** (±0.239) |
| PathDSP (GEx, T) ■ | 0.882 (±0.045) | **1.283** (±0.230) | 0.380 (±0.146) | 2.648 (±0.287) | 0.876 (±0.074) | 1.103 (±0.256) |
| PathDSP (GEx, FP) ■ | 0.882 (±0.043) | 1.297 (±0.227) | 0.139 (±0.146) | 2.944 (±0.327) | 0.887 (±0.068) | 1.051 (±0.243) |
| CDS (GEx, CNV, T) ★ | 0.777 (±0.049) | 1.625 (±0.189) | **0.378** (±0.164) | **2.606** (±0.320) | **0.776** (±0.083) | **1.478** (±0.287) |
| CDS (GEx, T) ■ | **0.789** (±0.098) | **1.603** (±0.254) | 0.345 (±0.156) | 2.724 (±0.317) | 0.769 (±0.106) | 1.508 (±0.304) |

**Randomly generated pathways provide comparable results to biological pathway collections for prediction of drug response in unseen cell lines**
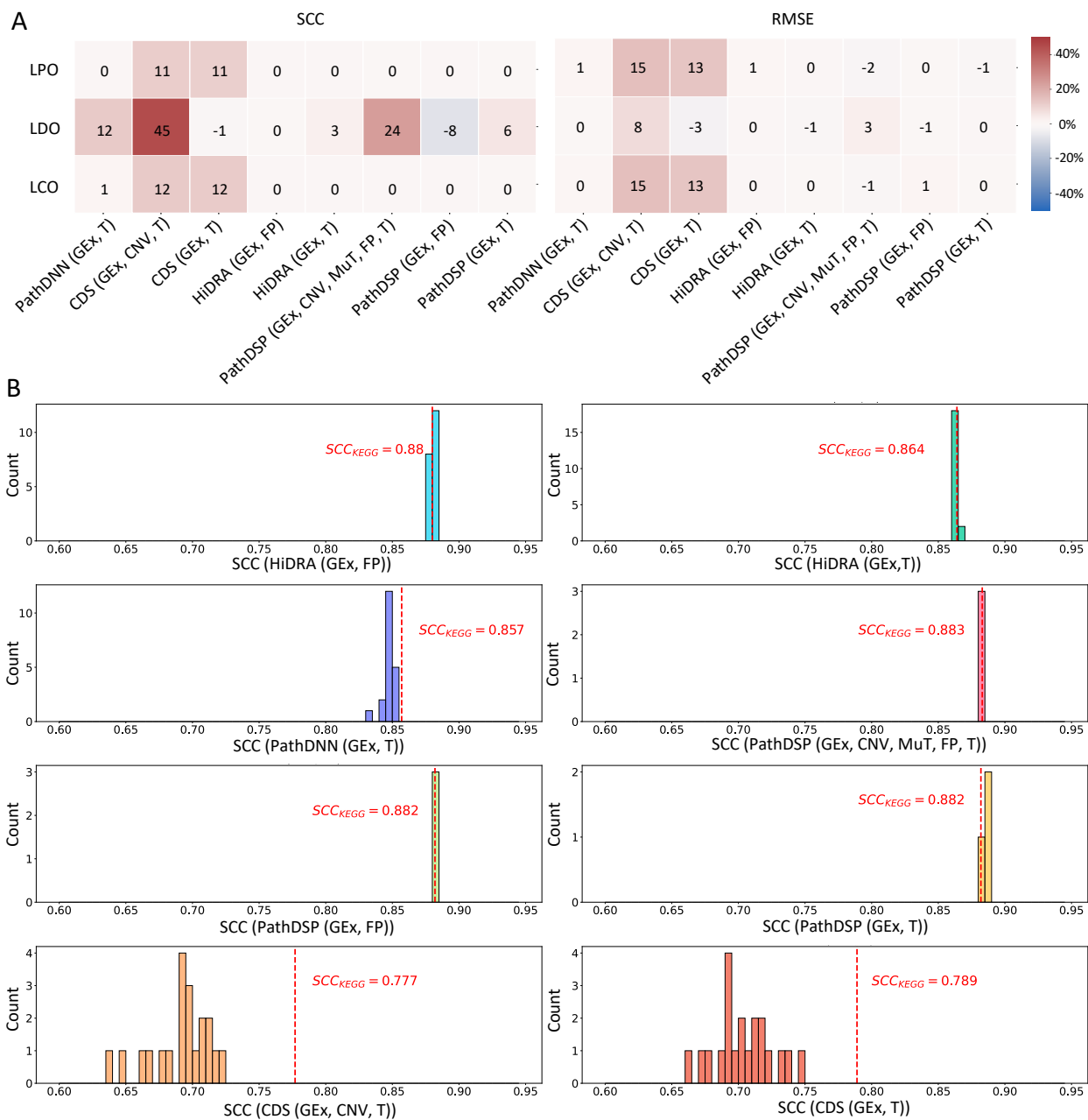
Next, we sought to determine whether the performance of pathway-based models can be attributed to the biological information in the pathways, or if randomly generated pathways can also result in a similar performance. For this purpose, we randomly assigned genes to pseudo-pathways while matching the size of the pathways in the KEGG collection. Figure 5A shows the percentage of improvement in the form of a heat map, where we compared the original pathway-based models (PathDNN, CDS, HiDRA, PathDSP) and their model variants with their corresponding random pathway baselines. Figure 5B and Supplementary Figures S6-S10 show

478     the distribution of SCC and RMSE of different models for each randomly generated pathway

479     collection in different validation schemes.

480

481     As can be seen in Figure 5, in the LCO and LPO evaluations, biological pathways provided almost

482     no improvement compared to their randomly generated counterparts (percentage of

483     improvement between -2% and 2%) for PathDSP, PathDNN, and HiDRA. CDS was the only

484     exception for which an improvement of up to 15% was obtained using biological pathways

485     (specifically for the original CDS (GEx, CNV, T) model). While it is difficult to conclusively

486     determine why CDS benefits from biological pathways, our conjecture is that this is due to its

487     unique architecture combined with the use of CNV input data. Moreover, it is important to note

488     that despite this improvement, CDS (GEx, CNV, T) and CDS (GEx, T) had much worse performance

489     compared to the other models (Figure 5B and Table 3).

490

491     In the LDO evaluation, the two models that used Morgan fingerprints to represent compounds,

492     PathDSP (GEx, FP) and HiDRA (GEx, FP) did not perform better with biological pathways compared

493     to randomly generated pathways. On the other hand, the majority of models that used drug

494     targets experienced an improvement compared to randomly generated pathways. We

495     investigated this behavior further by inspecting the number of drug targets in KEGG pathways

496     and the randomly generated pathways (Supplementary Figure S11). Comparing the number of

497     targets in each KEGG pathway with the randomly generated pathways of the same size showed

498     that in the majority of pathways (230 out of 332 pathways, approximately 70%), the number of

499     drug targets in the KEGG pathways were larger. Since drug targets are integrated with pathway

500    information to obtain drug embeddings, this difference in the number of drug targets results in

501    less informative and less distinguishable drug embeddings in the case of randomly generated

502    pathways. For example, in PathDNN where drug targets are represented as binary features, the

503    random pathway nodes are connected to many zero-valued drug features. Such nodes do not

504    participate much in capturing the similarities or differences of drugs, leading to embeddings that

505    are not as informative as their biological pathway counterparts in capturing patterns of similarity

506    and dissimilarity of drugs. This observation is also in line with a recent study that showed better

507    predictions could be obtained for compounds with diverse target classes (Kuenzi, et al., 2020).

508    The issue mentioned above is particularly important in the case of LDO, since unlike LCO and LPO

509    where all drugs in the test set have been seen by the model during training, the model must learn

510    drug similarity/dissimilarity patterns in order to make predictions for new drugs not observed

511    during training. This results in a deterioration of performance in random-pathway models

512    (parituclarly in LDO) compared to their biological counterparts observed in Figure 5.

**Figure 5:** Performance of pathway-based models using KEGG or randomly generated pathways. 1) Percentage of improvement (or deterioration) of different models when using KEGG compared to their mean performance when using randomly generated pathways. B) The histograms show the distribution of mean Spearman's correlation coefficient (SCC) of random pathway baselines using the leave-cell-lines-out (LCO) validation scheme. Vertical dashed red lines show SCC of the model when using KEGG pathways. Twenty random pathway baselines were constructed for each model, except for PathDSP models. Since PathDSP requires 1000 permutation tests for each type of input data, only three random pathway baselines were constructed due to its extremely high computational requirement.

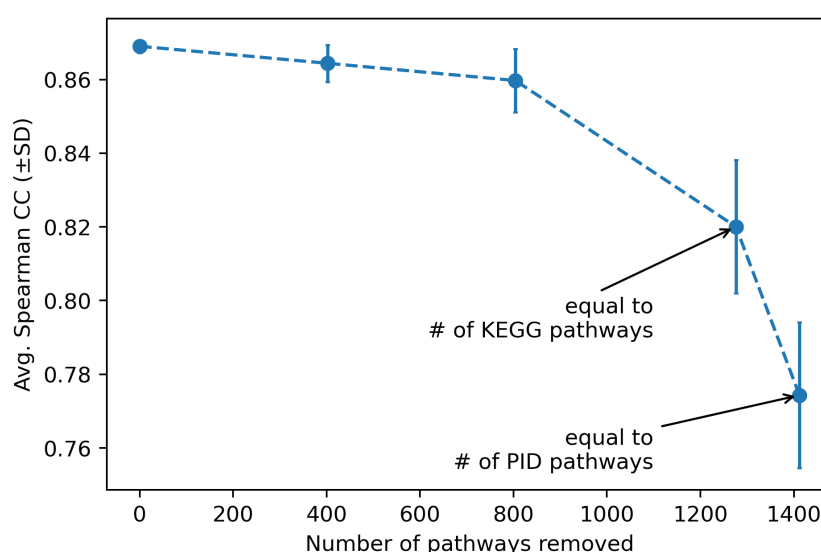524 **Effect of pathway collection choice on drug response prediction**

525 We next sought to investigate which pathway collection is more suitable for the drug response

526 prediction task. For this purpose, we compared the performance of each pathway-based model

527 (in their original architecture and using original input features) using each of these collections

528 (see Supplementary Tables S7 and S8 for the performance of all models and their variants using

529 PID and Reactome). To ensure a fair comparison, we only included (drug, CCL) pairs in the test

530 sets that were shared among all three uniform datasets. We used the LCO data splitting approach,

531 since the overlap among the test samples of the three uniform datasets was largest in this

532 strategy (21525 pairs versus 3277 in LDO and 851 in LPO).

533

534 **Table 5:** Performance of pathway-based models using different pathway collections. Models with
535 input data used in their original studies are used in this table. More specifically, the models
536 correspond to PathDNN (GEx, T), CDS (GEx, CNV, T), HiDRA (GEx, FP), and PathDSP (GEx, CNV,
537 MuT, FP, T). Mean and standard deviation are calculated across cell lines using the leave-cell-
538 lines-out (LCO) evaluation. Supplementary Figure S12 provides visualization of these values in the
539 form of bar plots and Supplementary Table S9 provides comparison of these models using
540 Wilcoxon signed rank tests.

| Pathway Collection | PathDNN | | CDS | | HiDRA | | PathDSP | |
|---|---|---|---|---|---|---|---|---|
| | SCC (±std) | RMSE (±std) | SCC (±std) | RMSE (±std) | SCC (±std) | RMSE (±std) | SCC (±std) | RMSE (±std) |
| Reactome | 0. 86 (±0.04) | 1.35 (±0.24) | 0.76 (±0.05) | 1.63 (±0.24) | 0.88 (±0.05) | 1.26 (±0.25) | 0.87 (±0.05) | 1.29 (±0.25) |
| PID | 0.85 (±0.05) | 1.42 (±0.26) | 0.78 (±0.06) | 1.63 (±0.27) | 0.88 (±0.05) | 1.28 (±0.29) | 0.88 (±0.05) | 1.29 (±0.25) |
| KEGG | 0.85 (±0.05) | 1.49 (±0.29) | 0.77 (±0.05) | 1.61 (±0.20) | 0.87 (±0.05) | 1.29 (±0.26) | 0.88 (±0.05) | 1.30 (±0.28) |

541

542 Table 5 and Supplementary Figure S12 show the mean and standard deviation of SCC and RMSE

543 of each model using all three pathway collections. Overall, we observed that the performance of

544    most models did not vary drastically based on the choice of pathway collection. However,

545    Reactome pathway provided slightly better results for the majority of the methods, being the top

546    performing option for 3 (out of 4 methods) based on RMSE. We hypothesized that this is due to

547    the larger number of pathway annotations included in this database for our use-case (1608

548    pathways in Reactome compared to 332 in KEGG and 196 in PID), resulting in a more

549    comprehensive representation of the input data.



550

**Figure 6:** Performance of PathDNN (GEx, T) with downsampled Reactome pathways. The y-axis
shows the mean (Avg.) and standard deviation (SD) of Spearman's correlation coefficient (SCC)
and the x-axis shows the number of pathways removed from the Reactome collection.

554

555    To test whether the large number of pathways in Reactome can explain its better performance,

556    we randomly downsampled the pathways in this collection. Figure 6 shows the SCC for PathDNN

557    (GEx, T) using different number of pathways removed (x-axis). We focused on PathDNN (GEx, T),

558    since it achieved its best performance when using Reactome pathway collection (compared to

559    PID or KEGG). For each value on the x-axis, downsampling was performed ten times and the

560    results were used to calculate the mean and standard deviation in the LCO setup (Supplementary

561   Table S10 provides details of each run). This figure shows that indeed, the number of pathways

562   in the Reactome collection plays a major role in its performance: as more pathways are removed,

563   the performance of PathDNN (GEx, T) deteriorates, with the lowest mean SCC value obtained

564   when only 196 (equal to the number of pathways in PID) have remained. This signifies that the

565   comprehensiveness of Reactome has enabled PathDNN to achieve better results. Interestingly,

566   the performance of this model with PID or KEGG was much better compared to the downsampled

567   version of Reactome with the same number of pathways (Table 5 and Figure 6). We attribute this

568   to the increasing probability of removing an important pathway during random downsampling of

569   Reactome, as well as the quality of the curated pathways in KEGG and PID.

570

571   **Cross-dataset performance of pathway-based models**

572   In addition to the analysis performed using GDSC reported earlier, we also assessed the

573   generalizability of the deep learning models to predict response of drugs in CTRPv2 (Rees, et al.,

574   2016). For this purpose, we trained the models on GDSC using drug AUC values and assessed their

575   performance on the prediction of AUC of drugs in CTRPv2. For consistency, all models were

576   trained using gene expression and drug targets. Supplementary Table S11 shows the results in

577   various data splitting and evaluation setups. Similar to our previous analyses on GDSC, in LCO and

578   LPO setup, implicit models performed better compared to explicit models and also outperformed

579   MLP. However, in the LDO setup MLP baseline achieved the best performance. The performance

580   of all models on CTRPv2 deteriorated compared to their performance on GDSC, highlighting the

581   challenging nature of this task.

582

583    **Discussion**

584    Recently, several deep learning methods have been proposed to enable a higher interpretability

585    of drug response prediction and to improve the prediction performance. In this study, we set out

586    to investigate four methods that try to achieve these goals by incorporating pathway information

587    under various validation schemes and answer five important questions discussed earlier. The

588    models were tested to predict drug response for unseen (CCL, drug) pairs, unseen CCLs, and

589    unseen drugs. We compared these methods against four types of baseline models, two of which

590    were usually overlooked in previous studies.

591

592    First, we observed that models that incorporate a dedicated explicit pathway layer and connect

593    gene nodes in a previous layer to pathways based on pathway membership perform worse

594    compared to models that implicitly (e.g., using attention mechanisms or pathway enrichment

595    scores) incorporate pathway information. In fact, in many occasions explicit models'

596    performance was inferior to a black-box simple MLP model with similar input. This suggests that

597    direct encoding of gene-pathway membership is not an effective strategy to incorporate pathway

598    information. The overly sparse connections between the gene and pathway layer may be the

599    cause for the unsatisfactory performance of these methods (due to a reduction in their capacity),

600    supported by the observation that their MLP counterparts (with fully connected layers) achieved

601    a better performance. Another limitation of explicit models is that they can only utilize gene-level

602    drug representations, limiting usable drug features to drug targets. Our analysis using methods

603    that could utilize both drug targets and Morgan fingerprints showed the latter to be superior in

604    prediction of response for unseen CCLs or unseen pairs. However, recent studies have suggested

605    that alternative drug representations such as transcriptomic changes in response to compounds

606    (El Khili, et al., 2022) or DL-based fingerprints (Zagidullin, et al., 2021) may improve performance

607    of drug response predictors.

608

609    Our analyses also showed that while implicit models generally performed better in predicting

610    unseen CCLs and unseen pairs, a comparable performance can be achieved when instead of

611    biological pathways, randomly generated pathways are used. Moreover, in these validation

612    setups a black-box MLP that used Morgan fingerprints for drug representation outperformed all

613    pathway-based models. Put together, these results suggest that to make the models

614    interpretable, these approaches inevitably make assumptions that cannot fully capture the

615    nuances of drugs' mechanisms of action in cancer cell lines, resulting in comparable or worse

616    performance compared to black-box models.

617

618    Our analyses also allowed us to assess the difficulty of drug response prediction in different

619    setups. While at first glance, Table 3 may suggest that predicting response of unseen drugs are

620    much more challenging than unseen CCLs, a more appropriate comparison can be made using

621    Figure 3A, where the models' performance improvement under each validation setup was

622    compared against a naive predictor. This figure shows that predicting response of unseen pairs

623    is much easier compared to prediction for unseen CCLs and unseen drugs. This is not surprising,

624    since this is a transductive setup and drugs and CCLs in the test set are present in the training set

625    (but not together). This setup is useful for imputation of missing drug response values but cannot

626    be used to predict response to new CCLs or new drugs. On the other hand, predicting response

627　of unseen drugs and unseen CCLs are much more difficult and most models cannot provide a

628　better prediction for the majority of CCLs in these two setups (Figure 3A). Moreover, the

629　performance of all models deteriorated when used to predict the drug response in a different

630　dataset (CTRPv2), revealing the challenging nature of this task.

631

632　The contrast between conclusions one may draw from Table 3 and Figure 3A (discussed above)

633　demonstrates the potential for obtaining inflated performance measures in the LCO framework,

634　in which for a given CCL in the test set, the log IC50 value of different drugs are to be predicted.

635　Supplementary Figure S13 shows the distributions of log IC50 values for each drug (across CCLs,

636　panel A) and each CCL (across drugs, panel B) in our dataset. While these values vary across both

637　drugs and CCLs, the identity of a drug plays a bigger role in determining its log IC50 value

638　compared to the identity of the CCL to which it was administered (i.e., log IC50 values are more

639　drug-specific than CCL-specific). Supplementary Figure S13C better clarifies this point by

640　depicting the histogram of false discovery rates (FDRs) obtained from comparing the local

641　distribution of log IC50 values per drug (purple) or per CCL (green) and the global distribution of

642　the log IC50 values using Mann-Whitney U tests. Although for the majority of drugs (93%) drug-

643　specific log IC50 values (across all CCLs, but for a specific drug) are significantly different (FDR <

644　0.05) from the global log IC50 values (across all drugs and CCLs), that is true for only 43% of CCLs.

645　This implies that by simply knowing the identity of a drug, a model can rank different drugs based

646　on their log IC50 values for an unseen CCL rather well. To overcome this issue and avoid reporting

647　unrealistically inflated metrics, one should focus on improvement compared to a naive predictor

648　(an approach that we adopted in this study), or should normalize log IC50 values of each drug

649     across CCLs to make them comparable to each other (an approach that we adopted in (Hostallero,

650     et al., 2022)).

651

652     We also compared the performance of different models using three pathway collections, PID,

653     KEGG, and Reactome. Even though there was not a major difference in the performance of

654     models when substituting one collection for the other, Reactome collection resulted in slightly

655     better performance. This seems to be due to the larger number of pathways in this collection

656     compared to the other two. However, since randomly generated pathway collections also

657     provided comparable performance based on the models considered in this study, it is not

658     possible to draw a conclusive determination regarding which pathway collection may be more

659     useful for the drug response prediction task.

660

661     This study focused on evaluating the effect of incorporating pathway information from the

662     perspective of model performance and we did not evaluate these models based on their level of

663     interpretability. A study that focuses on interpretability aspect of these models would be very

664     insightful and complementary to the current study. For example, one can take a closer look at

665     the feature attributions of these pathway-based models using explainers such as DeepLIFT (Deep

666     Learning Important FeaTures) (Shrikumar, et al., 2017), CXPlain (Schwab and Karlen, 2019), and

667     SHAP (Shapley Additive exPlanations) (Lundberg and Lee, 2017) to estimate feature importance

668     and identify genes or other biological features that have substantial influence on the model

669     predictions. Such analysis can be done for all pathway-based models to check if the most

670     important/predictive sub-networks or the top contributing genes extracted from each model

671   have any overlap. If such overlap exists between the pathway-based models, further studies can

672   be done by validating the findings with existing literature or conducting experiments under a lab

673   setting. Analysis on model interpretability will complement the insights obtained from model

674   performance evaluation and together provide a more holistic view for the effect of pathway

675   incorporation on drug response prediction.

676

677   In conclusion, we believe that while interpretability is a very crucial aim in precision medicine,

678   new models are necessary to enable a higher degree of interpretability while at the same time

679   improve the drug response prediction performance. In addition, it is not sufficient for these

680   models to show a better performance compared to their black-box counterparts, and they need

681   to also evaluate their models against randomly generated pathways (with similar pathway sizes

682   to the original collection) and naive predictors to control for different types of biases.

683

684   **Data and Code Availability:** Input data for the evaluated models is provided at

685   https://zenodo.org/record/7101665#.YzS79HbMKUk. The implementation of the models are

686   available at https://github.com/Emad-COMBINE-lab/InterpretableAI_for_DRP.

687

695

**696      Competing Interests:** None of the authors have any competing interests.

697

698    **References**

699    Adam, G*., et al.* Machine learning approaches to drug response prediction: challenges and recent progress.
700    *NPJ Precis Oncol* 2020;4:19.

701    Azodi, C.B., Tang, J. and Shiu, S.H. Opening the Black Box: Interpretable Machine Learning for Geneticists.
702    *Trends Genet* 2020;36(6):442-455.

703    Ballester, P.J*., et al.* Artificial intelligence for drug response prediction in disease models. *Brief Bioinform*
704    2022;23(1).

705    Baptista, D., Ferreira, P.G. and Rocha, M. Deep learning for drug response prediction in cancer. *Brief*
706    *Bioinform* 2021;22(1):360-379.

707    Barredo Arrieta, A*., et al.* Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and
708    challenges toward responsible AI. *Information Fusion* 2020;58:82-115.

709    Barretina, J*., et al.* The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug
710    sensitivity. *Nature* 2012;483(7391):603-7.

711    Caruana, R*., et al.* Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day
712    readmission. In, *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and*
713    *Data Mining*. 2015. p. 1721-1730.

714    Che, Z*., et al.* Interpretable Deep Models for ICU Outcome Prediction. In, *AMIA Annu Symp Proc*. 2016. p.
715    371-380.

716    Chen, Y. and Zhang, L. How much can deep learning improve prediction of the responses to drugs in cancer
717    cell lines? *Brief Bioinform* 2022;23(1).

718    Costello, J.C*., et al.* A community effort to assess and improve drug sensitivity prediction algorithms.
719    *Nature Biotechnology* 2014;32:1202-1212.

720    Deng, L*., et al.* Pathway-guided deep neural network toward interpretable and predictive modeling of
721    drug sensitivity. *Journal of Chemical Information and Modeling* 2020;60:4497-4505.

722    El Khili, M.R., Memon, S.A. and Emad, A. MARSY: A multitask deep learning framework for prediction of
723    drug combination synergy scores. *bioRxiv* 2022:bioRxiv 2022.06.07.495155.

724    Emad, A*., et al.* Knowledge-guided gene prioritization reveals new insights into the mechanisms of
725    chemoresistance. *Genome Biol* 2017;18(1):153.

726    Fabregat, A*., et al.* Reactome pathway analysis: A high-performance in-memory approach. *BMC*
727    *Bioinformatics* 2017;18.

728    Guvenc Paltun, B., Mamitsuka, H. and Kaski, S. Improving drug response prediction by integrating multiple
729    data sources: matrix factorization, kernel and network-based approaches. *Brief Bioinform* 2021;22(1):346-
730    359.

731    Heller, S.R*., et al.* InChI, the IUPAC International Chemical Identifier. *Journal of Cheminformatics* 2015;7.

732  Hostallero, D.E., Li, Y. and Emad, A. Looking at the BiG picture: incorporating bipartite graphs in drug
733  response prediction. *Bioinformatics* 2022;38:3609-3620.

734  Hostallero, D.E*., et al.* A Deep Learning Framework for Prediction of Clinical Drug Response of Cancer
735  Patients and Identification of Drug Sensitivity Biomarkers using Preclinical Samples. *bioRxiv* 2021:bioRxiv
736  2021.07.06.451273.

737  Huang, E.W*., et al.* Tissue-guided LASSO for prediction of clinical drug response using preclinical samples.
738  *PLoS Comput Biol* 2020;16(1):e1007607.

739  Jarada, T.N., Rokne, J.G. and Alhajj, R. A review of computational drug repositioning: strategies,
740  approaches, opportunities, challenges, and directions. *J Cheminform* 2020;12(1):46.

741  Jin, I. and Nam, H. HiDRA: Hierarchical Network for Drug Response Prediction with Attention. *Journal of*
742  *Chemical Information and Modeling* 2021;61:3858-3867.

743  Kanehisa, M. and Goto, S. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Research*
744  2000;28:27-30.

745  Kim, S*., et al.* PubChem in 2021: New data content and improved web interfaces. *Nucleic Acids Research*
746  2021;49:D1388-D1395.

747  Kuenzi, B.M*., et al.* Predicting Drug Response and Synergy Using a Deep Learning Model of Human Cancer
748  Cells. *Cancer Cell* 2020;38:672-684.e6.

749  Landrum, G. RDKit: Open-source Cheminformatics. In, *Http://Www.Rdkit.Org/*. 2006.

750  Lundberg, S.M. and Lee, S.-I. A unified approach to interpreting model predictions. *Advances in neural*
751  *information processing systems* 2017;30.

752  Malioutov, D.M*., et al.* Learning Interpretable Classification Rules with Boolean Compressed Sensing. In:
753  Cerquitelli, T., Quercia, D. and Pasquale, F., editors, *Transparent Data Mining for Big and Small Data*. Cham:
754  Springer International Publishing; 2017. p. 95-121.

755  Rees, M.G*., et al.* Correlating chemical sensitivity and basal gene expression reveals mechanism of action.
756  *Nat Chem Biol* 2016;12(2):109-16.

757  Schaefer, C.F*., et al.* PID: The pathway interaction database. *Nucleic Acids Research* 2009;37.

758  Schwab, P. and Karlen, W. CXPlain: Causal explanations for model interpretation under uncertainty.
759  *Advances in Neural Information Processing Systems* 2019;32.

760  Sharifi-Noghabi, H*., et al.* Drug sensitivity prediction from cell line-based pharmacogenomics data:
761  guidelines for developing machine learning models. *Brief Bioinform* 2021;22(6).

762  Shrikumar, A., Greenside, P. and Kundaje, A. Learning important features through propagating activation
763  differences. In, *34th International Conference on Machine Learning, ICML 2017*. 2017. p. 4844-4866.

764     Snow, O*., et al.* Interpretable Drug Response Prediction using a Knowledge-based Neural Network. In,
765     *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2021.
766     p. 3558-3568.

767     Szklarczyk, D*., et al.* STRING v11: Protein-protein association networks with increased coverage,
768     supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Research*
769     2019;47:D607-D613.

770     Szklarczyk, D*., et al.* STITCH 5: Augmenting protein-chemical interaction networks with tissue and affinity
771     data. *Nucleic Acids Research* 2016;44:D380-D384.

772     Tang, Y.C. and Gottlieb, A. Explainable drug sensitivity prediction through cancer pathway enrichment.
773     *Scientific Reports* 2021;11.

774     Vamathevan, J*., et al.* Applications of machine learning in drug discovery and development. *Nat Rev Drug*
775     *Discov* 2019;18(6):463-477.

776     Yang, W*., et al.* Genomics of Drug Sensitivity in Cancer (GDSC): A resource for therapeutic biomarker
777     discovery in cancer cells. *Nucleic Acids Research* 2013;41.

778     Zagidullin, B*., et al.* Comparative analysis of molecular fingerprints in prediction of drug combination
779     effects. *Brief Bioinform* 2021;22(6).

780     Zhang, H., Chen, Y. and Li, F. Predicting Anticancer Drug Response With Deep Learning Constrained by
781     Signaling Pathways. *Frontiers in Bioinformatics* 2021;1.

782     Zhang, H*., et al.* Benchmarking network-based gene prioritization methods for cerebral small vessel
783     disease. *Brief Bioinform* 2021;22(5).

784