
A Foundational Multi-Modal Knowledge Graph for Pancreatic Cancer Drug Effects Prediction

Jingwen Hui

Shu Chien-Gene Lay Department of Bioengineering
University of California San Diego
La Jolla, CA 92093
jwhui@ucsd.edu

Shengchao Liu*

Independent
shengchao1224@gmail.com

Xiaohua Huang*

Shu Chien-Gene Lay Department of Bioengineering
University of California San Diego
La Jolla, CA 92093
x2huang@ucsd.edu

Anima Anandkumar*

Computing + Mathematical Sciences
California Institute of Technology
Pasadena, CA 91125
anima@caltech.edu

Abstract

1 AI-assisted drug discovery has revolutionized healthcare by accelerating virtual
2 screening methods as compared to traditional processes. Many advanced AI mod-
3 els have been developed to predict and generate drug candidates, with potential
4 applications across various diseases. However, challenges still remain in apply-
5 ing AI models in clinical settings. These include the lack of heterogeneity and
6 insufficient consideration of patient-specific treatment plans. To mitigate these
7 challenges, we propose PanRX, a cell-line-specific pancreatic cancer drug effect
8 model using multi-modal knowledge graphs. It aims at achieving a personalized
9 drug discovery framework by incorporating rich genetic and chemical information.
10 We first construct a multi-modal knowledge graph dataset PanCan-DrugsGenes. It
11 extracts textual genetic information from NCBI, mutation status from the Genomics
12 of Drug Sensitivity in Cancer (GDSC) dataset, textual descriptions of drugs from
13 PubChem, and chemical geometry from the PM6 dataset. Then, PanRX utilizes a
14 geometric model to learn chemical conformation, a language model to learn textual
15 description, and a graph neural network to fuse all information and predict the
16 target drug effects. We verify the effectiveness of PanRX by achieving the general-
17 ization performance with very low MSE (< 0.0000) and MAE (0.0009). This
18 work emphasizes the potential of merging knowledge graphs and deep learning in
19 the fields of genomics and medicine, enriching the intersection of human biological
20 expertise and AI in drug discovery and design tasks.

21 1 Introduction

22 Personalized medicine is becoming the pillar of modern medicine because it considers varying
23 phenotypes of diseases as a result of differences in genetic information. Genetic information,
24 including gene and copy number alterations (CNA) mutation status, gene regulatory networks, gene
25 function, and gene location are critical information that collectively determine the wellbeing of
26 individuals [1, 2]. Considering these factors is essential to understanding cancer mechanisms and
27 therapeutic requirements in addition to the chemical/medicinal properties of various drugs. Databases

*Equal advising author

28 like GDSC [3] have collected drug effect experiments and genomic association tests on over 1,000
29 cancer cell lines. However, predicting drug effects given genetic information is still limited by data
30 insufficiency and a lack of insights into the genetic circuit. Nevertheless, deep learning (DL) tools
31 can integrate multi-modal information from all human knowledge and are expected to generalize to
32 unseen tasks. They can be adapted to uncover hidden relationships of the genetic circuit to generalize
33 partial data and offer robust predictions efficiently compared to human methods.

34 Recent breakthroughs in deep learning have accelerated the drug discovery process due to its ability for
35 pattern recognition. Among these, conventional (transductive) knowledge graphs (KG) gained wide
36 attraction in computational biology due to their network-oriented structure [4, 5, 6]. Transductive KG
37 works with a fixed set of nodes and edges. In the biological context, transductive KG is constructed
38 only with known entities (e.g. genes, proteins, diseases) and their interactions (e.g. genetic pathways,
39 gene-protein interactions, gene-disease associations). Existing DL pipelines for drug discovery
40 predominantly use SMILES, 2D molecular graphs, and conformation data of chemical compounds
41 [7]. However, the complexity of the interactions within the living systems due to biophysical
42 conditions necessitates a more advanced method of drug effect predictions. Recent advancements
43 in geometric learning further popularized drug effect prediction, catalyzing state-of-the-art models
44 to capture these complex biological interactions. For instance, B. Kuenzi et al predicted the drug
45 response of human cancer cells using chemical structure data, genomics data, drug sensitivity data,
46 and protein activity data [8]. However, the choice of 2D structure instead of 3D topology and a
47 lack of comprehensive gene descriptions/networks limit the model's ability to fully capture complex
48 interactions of drugs and intricate biological pathways. In addition, drug response is often heavily
49 influenced by one's genetic makeup. The gene regulatory network poses a great challenge for
50 researchers in this field because of the hidden interactions and a lack of biological understanding of
51 this subject. To address these challenges, this research implements multi-modal data frameworks that
52 are descriptive (geometric/textual information) and relationship-focused.

53 **1.1 Preliminaries: Drug Discovery Process**

54 Drug discovery involves five main phases; 1) the pre-discovery stage where disease mechanisms are
55 explored; 2) the discovery stage where scientists search for appropriate small-molecule therapeutics
56 that interfere with the disease mechanism; 3) the preclinical stage where drug candidates are tested
57 for their efficacy on various in-vitro or in-vivo models; 4) the clinical stage for human testing; 5)
58 the post-market reviewing and approval of this drug [9]. Our ML pipeline aims to optimize the
59 pre-clinical phase by predicting in-vitro drug effects on various cell lines.

60 **Traditional Preclinical Development**

61 Extensive in vitro tests are performed during preclinical development. These experiments test
62 potential drug efficacy before proceeding to in-vivo studies and clinical trials. IC₅₀, AUC, and
63 Z-score are some of the most informative metrics that shed light on different aspects of the drug's
64 behavior. The IC₅₀ value measures the drug's potency. It shows the concentration needed to inhibit a
65 biological or biochemical function by 50% [10]. The AUC is derived from a dose-response curve that
66 represents the effect of various drug concentrations on each cell line. This value summarizes the drug
67 efficacy across a range of concentrations [11]. The Z-score compares an IC₅₀ value with those from
68 other cell lines, showcasing the effectiveness of a drug in comparison to the average response, and
69 highlighting whether the drug is more or less potent in a specific cell line relative to others.

70 Traditionally, screening these values costs hundreds of millions of dollars and requires years of
71 clinical testing [12]. On average, out of 10,000 molecules screened, only one may eventually lead to
72 the market. Due to the immense efforts of validating one drug molecule, there is limited flexibility in
73 the traditional pipeline. Such inefficiencies pose an economic barrier to drug development for serious
74 diseases like cancer. Consequently, companies place a lower priority on these endeavors and often
75 shift their focus toward more cost-effective avenues, potentially delaying the development of critical
76 treatments for diseases like cancer.

77 **AI-Assisted Preclinical Development**

78 Artificial Intelligence (AI) is revolutionizing the drug discovery process. The power of AI can be
79 hugely manifested in medicine due to its ability to recognize patterns in vast amounts of data with
80 varying modalities, personalize treatment plans, and predict patient outcomes [13]. For instance,
81 supervised learning is heavily used for the prediction of molecular properties, pharmacokinetics,

82 chemical synthesis, etc [14, 15]. However, the biggest challenge lies in the complexity of data. This
83 arises from the diversity of data types, the complexity of chemical interactions in biological systems,
84 and the difficulty of accurately encoding this information into an ML pipeline.

85 1.2 Our contributions.

86 We introduce PanRX, a deep-learning model designed to predict drug effects by integrating multi-
87 modal drug and genetic information through a transductive KG. We leverage the extensive training
88 capabilities of language models and geometric models with the relational structure of KGs to capture
89 complex interactions between drug molecules and cell lines. On the drug side, we incorporate
90 3D geometric information such as atomic numbers, bonding details, and 3D coordinates of atomic
91 positions as well as SMILES string representation of molecular topology. Also, textual information
92 (including drug summaries, pharmacodynamics, indications, and mechanism of action) and numeric
93 data (including charge, enthalpy, and free energy) further enrich the model. On the genetic side,
94 PanRX integrates gene-related textual data including gene summaries, expression patterns, cellular
95 locations, interactions, and binary mutation status for genes and copy number alterations. This
96 multi-modal fusion of 3D geometry and comprehensive genomic descriptions within a KG framework
97 emphasizes extensive biological relationships, providing a nuanced understanding of drug-cell line
98 interactions. Therefore, PanRX bridges the gap between current models and real-world biological
99 complexities, offering improved predictive accuracy and possibility in drug discovery.

100 To verify the effectiveness of PanRX, we created PanCan-DrugsGenes, a dataset specifically designed
101 for pancreatic cancer research. PanCan-DrugsGenes is a Lightning Memory-Mapped Database
102 (LMDB) with 204 pancreatic cancer drugs and 142 genes that are highly correlated with the disease.
103 Each drug or gene has multi-modal data that conveys rich information. Each KG is constructed from
104 the outputs of this dataset. For evaluation of PanRX, we utilize IC50, AUC, and Z-Score values to
105 measure the confidence of predicted drug effects.

106 To the best of our knowledge, PanRX is the first to combine 3D drug chemical structures with
107 genetic information to enhance personalized drug design. In addition, PanRX emphasizes complex
108 interactions within biological networks, labeling each connection to provide detailed information
109 about the nature of these relationships. This structure prepares for future downstream applications
110 such as entity/relationship predictions. Extensive experiments on clinical datasets were conducted.
111 The results show powerful drug effect predictive capabilities for pancreatic cancer.

112 2 Related Work

113 **Multi-Modal Modeling on Small Molecules.** In existing ML for the drug discovery community,
114 there are multiple modalities describing molecules, and they can be roughly divided into two venues:
115 internal chemical structure and external functional description [16]. For internal chemical structures,
116 GraphMVP [17] initiates the molecule pretraining by utilizing the 2D topology and 3D geometry.
117 Follow-up works like MoleculeSDE [18] extend this line by proposing a more advanced geometric
118 pretraining algorithm, and MoleculeJAE [19] utilizes the molecular dynamics for pretraining. On
119 the other hand, MolT5 [20] and MoleculeSTM [21] are the first two works to align both the internal
120 chemical structures and the external functional description for molecule design and editing.

121 **Multi-Modal Modeling on Genomes.** Genomic information is organized in various ways depending
122 on the research objective. FASTA formats are used to represent sequences (DNA, RNA, protein)
123 and are predominantly used for gene function prediction or drug-target interactions [22]. FASTA
124 provides efficient storing of biological sequence strings that allows for straightforward retrieval and
125 analysis. Feature tables gained popularity due to their ability to map genomic regions to biological
126 functions, assisting the understanding of gene regulation and expression under various conditions.
127 Sparse matrices are used to handle high-dimensional gene expression data used during drug response
128 prediction tasks [23]. Sparse matrices are also computationally efficient due to the majority of
129 elements being zero.

130 This research team utilizes textual descriptions of the genome (genes and CNAs) because text
131 descriptions can encapsulate information that is represented by any other types of genomic data
132 structures. Textual data can effectively represent nuanced information such as genomic interactions,
133 regulatory mechanisms, and expression conditions. By integrating these descriptions that are often lost

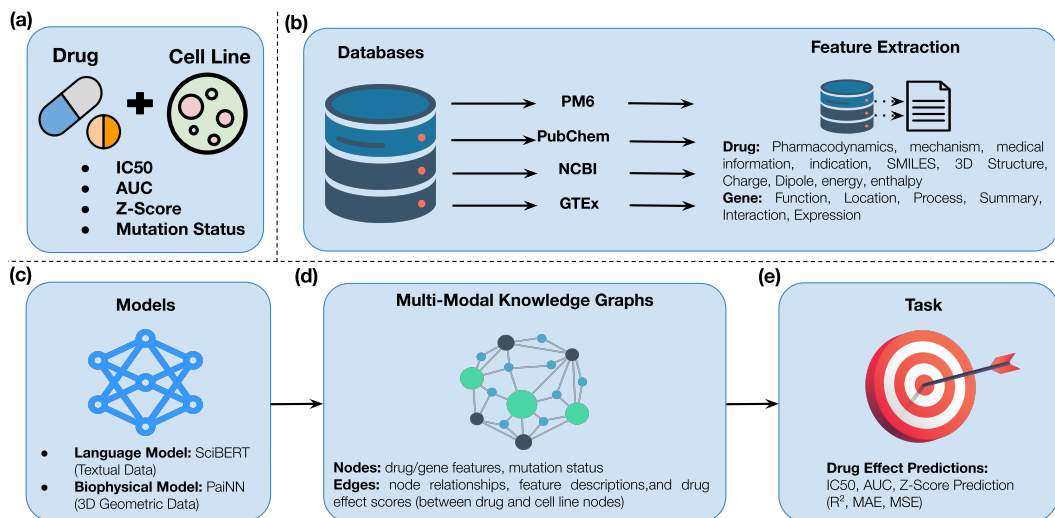


Figure 1: An illustration of the pancreatic cancer drug prediction (PanRX) pipeline. (a) The bioassays (IC₅₀, AUC, Z-Score) of each drug on diseased cell lines is extracted from GDSC. In addition, the mutation status of key genes and CNAs are recorded. (b) Drug and genetic features were extracted from open-access databases like PM6, PubChem, NCBI, and PTE_x. (c) Textual and geometric data are encoded by SciBERT and PaiNN respectively. (d) Multi-modal knowledge graphs are constructed for each pair of drug-cell line interactions. (e) The bioassays (IC₅₀, AUC, Z-Score) are predicted.

134 in purely numeric data modalities, we enhance the application of genomic information in personalized
 135 medicine.

136 **Multi-Modal Modeling on Knowledge Graph.** In this work, we are merging the multi-modal infor-
 137 mation of small molecules and genomes using a knowledge graph. Existing knowledge graph papers
 138 use either single-modal information on small molecules [24, 25, 26] or single-modal information on
 139 genomes [27, 28]. More recent works have started to merge the information of different entities, such
 140 as small molecules and proteins [29]. However, as illustrated in recent benchmark works [30], the
 141 geometric information has been more informative for molecule representation.

142 3 Methods

143 In drug development, performing precise drug effect predictions on specific cancer cell lines is crucial
 144 yet challenging. Genetic variations such as gene and copy number alteration (CNA) mutation status
 145 pose extreme obstacles to anticipating drug performance in different biological contexts. Traditional
 146 approaches rely on time-consuming procedures, significantly delaying the drug delivery pipeline.

147 Predicting drug effects in pancreatic cancer is a multi-dimensional barrier involving integrating
 148 various levels of knowledge such as molecular properties, genetic information, and relevant bioassays.
 149 Unique genetic profiles like the variation of copy number alternations (CNA) and mutation statuses
 150 remain the root cause of this challenge. Traditional methods rely on repetitive in vitro tests and linear
 151 models which lack insights against complex biological interactions that define pancreatic cancer drug
 152 responses.

153 This paper aims to build a predictive pipeline that captures the interaction between drugs and
 154 pancreatic cancer cell lines. We focus on genetic variations across cell lines by leveraging cell-line
 155 genetic characteristics and chemical properties. Specifically, we utilize numerical data to represent
 156 drug characteristics, including charge, dipole moment, energy, and enthalpy, alongside bioassay
 157 measurements such as IC₅₀, AUC, and Z-score. Additionally, we incorporate binary data to capture
 158 the mutation status of genes and copy number alterations (CNA) within each cell line. By constructing
 159 a multi-modal knowledge graph with other chemical and genetic features, the problem is now framed
 160 as an edge prediction task a multi-modal knowledge graph where the aim is to predict drug-cell line
 161 interaction outcomes (IC₅₀, AUC, Z-Score).

162 Our approach involves the construction of knowledge graphs where nodes represent drugs, genes,
163 and their inherent features to emphasize variability. This task requires multi-modal data including
164 textual descriptions, chemical structures, and numerical data.

165 3.1 Problem Formulation

166 Our model can be described by $p(x, y, z | KG)$, where x , y , and z represent the IC50, AUC, and Z-
167 Score values, respectively, and KG refers to the multi-modal knowledge graphs we constructed. We
168 use graph neural networks (GNN) to fuse the multi-modal information discussed above. Leveraging
169 graph structure, we encode drug molecules (geometric data) and cell line genetic information (textual
170 data) as nodes and relationships between these entities as edges. This multi-modal approach enables
171 the model to learn from diverse data types of the biological network that are critical to predicting
172 drug responses. In this project, we consider two architects of the GNN to perform such fusion: graph
173 convolutional networks (GCNs) and graph attention networks (GATs).

174 **Geometric Modeling on Conformation.** In our study, we employ a geometric model to represent
175 3D drug conformations to further investigate their interactions with biological targets. The chemical
176 geometry is described using the PaiNN model $f_g = \text{PaiNN}(a, r)$, where a is atom types and r is
177 atom coordinates in the molecule [31]. By leveraging a geometrically pre-trained PaiNN [18], our
178 method assures all drug spatial configurations are precisely represented. This approach facilitates
179 improved accuracy interaction and effect prediction.

180 **Language Model on Functional Description.** On the other hand, we leverage the BERT architecture
181 by utilizing SciBERT to effectively capture the existing single-modal data. Specifically, drug
182 textual descriptions include drug summaries, pharmacodynamics, drug mechanisms, and indications.
183 Genetic textual descriptions include Gene functions, locations, processes, summary, interactions, and
184 expression. The model has a representation of $f_t = \text{SciBERT}(w_t)$, where w_t is the textual description
185 for each chemical/gene feature node. These embeddings capture intricate details contained within
186 textual descriptions, improving the overall model with an extra layer of biological context.

187 3.2 Multi-modal Fusion

188 **Graph convolution network (GCN).** Our proposed GCN model consists of six convolutional layers
189 (GCNConv), each followed by an Exponential Linear Unit (ELU) activation to capture non-linear
190 relationships. To prevent overfitting, a dropout rate of 0.3 is applied after each layer. These layers
191 aggregate node features on a global level based on node connectivity, allowing the convolutional layers
192 to propagate node feature embeddings across the graph, with each layer adjusting the embeddings
193 by considering adjacent nodes. A layer normalization is applied after the final convolutional layer
194 to ensure stable learning dynamics. The model is trained using the Adam optimizer with a learning
195 rate of 0.00001 and a weight decay of 0.0001. Performance is evaluated using Mean Squared Error
196 (MSE), Mean Absolute Error (MAE), and the coefficient of determination (R^2).

197 **Graph attention network (GAT).** The GAT setup builds upon the GCN architecture by integrating
198 attention mechanisms to dynamically assess the significance of neighboring nodes during the aggre-
199 gation phase. The model comprises six GAT layers, each employing multi-head attention to explore
200 various facets of a node’s vicinity. The initial five layers utilize four attention heads each, while the
201 final layer uses a single head to achieve dimensionality reduction of the embedding. This approach
202 permits selective integration of neighboring information, fostering a context-sensitive aggregation
203 compared to the GCN. The GAT model utilizes the ELU activation function and a dropout rate of
204 0.3. The optimization is performed using the Adam optimizer with a learning rate of 0.00001 and a
205 weight decay of 0.0001. The evaluation metrics are consistent with those used in the GCN model,
206 including MSE, MAE, and R^2 .

207 4 Experiments

208 4.1 Data Acquisition and Feature Extraction

209 Our primary data source was the GDSC dataset [3]. To evaluate the effectiveness of multi-modal
210 knowledge graphs, we compiled a set of 5014 pancreatic cancer drug-cell line pairs from GDSC.
211 Descriptive features of drugs and cell lines are pulled from public databases. For drugs, we extracted

212 3D geometry, compound charge, dipole moment, energy, and enthalpy from PM6 [32], and pharma-
 213 codynamics, drug mechanisms, medical information, and drug indication from PubChem. For each
 214 cell line, GDSC provides the mutation status of key genes and copy number alterations (CNA). We
 215 extract descriptive information about each gene from NCBI, Biopython, and GTEx such as function,
 216 location, process, summary, interactions, and expression level.

217 GDSC dataset provides comprehensive drug response information on various cell lines. We build
 218 a knowledge graph for each drug-cell pair and extract supporting descriptive features of the drug
 219 and related genes from public databases like PubChem, PM6, and NCBI. By feeding them into a
 220 deep learning pipeline, we assess the accuracy of drug effect predictions and graph relationship
 221 predictions. We selected graph attention networks (GAT) and graph convolution networks (GCN) as
 222 our deep learning architects because they can process graph-structured data and capture local and
 223 global dependencies. Both architects are designed for edge attribute predictions within the graph data,
 224 emphasizing the edge attributes that connect nodes 0 and 1 (IC50, AUC, Z-score).

225 4.2 Data Pre-Processing and Integration

226 We built an automated pipeline that streamlines the extraction process of relevant drug and gene
 227 information systematically. Firstly, this pipeline extracts data from our primary sources and stores
 228 them in the appropriate format in a Lightning Memory-Mapped Database (LMDB) to ensure efficient
 229 access to large-scale genomic and pharmacological data of various modalities. Then, we implemented
 230 a custom heterogeneous data loader in which each output instance contains a specific drug and
 231 the complete set of genes required for constructing a corresponding knowledge graph representing
 232 the interaction between a drug and a cell line. Data of various modalities are encoded using the
 233 models described in section 4.3 before they are assembled into nodes and organized into a knowledge
 234 graph structure. Appropriate edges are included to represent relationships that enhance the logical
 235 connectivity of each graph. It is critical to note that the connection between node 0 (drug) and node 1
 236 (cell line) is characterized by three edges each representing a different bioassay measurement - IC50,
 237 AUC, Z-Score.

238 4.3 Experimental Setup

239 The objective of this experiment is to predict edge attributes between the drug node (node 0) and the
 240 cell line node (node 1). These edge attributes represent bioassays, namely IC50, AUC, and Z-score
 241 values. Our models were trained at 500 epochs with a batch size of 400 to ensure sufficient stability.
 242 To evaluate model performances, we split the dataset into training, validation, and testing sets of
 243 80/10/10 and 50/25/25 ratios. Training, validation, and testing results (MSE, MAE, and R^2) are
 244 reported. On average, training out models (GAT and GCN) with the current parameters requires
 245 approximately 10 hours.

246 **Hardware** All experiments were run on the Nvidia 4090 GPU with 24GB of GDDR6X RAM and
 247 16,384 CUDA cores. We utilized the Windows operating system with Python version 3.12.1 for all
 248 computations.

249 4.4 Results

Table 1: Regression Analysis of Predicted Bioassay Values (80/10/10 split)

GAT (500 epochs)				GCN (500 epochs)			
Test Set	R^2	MSE	MAE	Test Set	R^2	MSE	MAE
Training	0.3327	0.0482	0.0184	Training	0.3253	0.0348	0.0152
Validation	0.3330	0.0002	0.0090	Validation	0.3321	0.0009	0.0171
Testing	0.3331	0.0002	0.0089	Testing	0.3325	0.0009	0.0169

250 As depicted in Table 1 and Table 2, both the GAT and GCN architectures performed consistently across
 251 the 80/10/10 and 50/25/25 splits, with R^2 at approximately 0.333 across the training, validation, and
 252 testing sets, highlighting its stable predictive performance. Also, the MSE was notably low across both
 253 architects and splits. For instance, in the 80/10/10 split, the GAT model achieved an MSE of 0.0002

Table 2: Regression Analysis of Predicted Bioassay Values (50/25/25 split)

GAT (500 epochs)				GCN (500 epochs)			
Test Set	R ²	MSE	MAE	Test Set	R ²	MSE	MAE
Training	0.3332	0.0049	0.0070	Training	0.3332	0.0052	0.0073
Validation	0.3333	0.0000	0.0010	Validation	0.3333	0.0000	0.0028
Testing	0.3333	0.0000	0.0009	Testing	0.3333	0.0000	0.0025

254 during validation and testing, while the GCN exhibited a slightly higher MSE of 0.0009. Additionally,
 255 in the 50/25/25 split, the MSE of both the GAT and GCN showed 0.0000 across these phases.

256 On the other hand, the GAT consistently achieved lower values relative to GCN in terms of MAE.
 257 In the 80/10/10 split, the GAT had an MAE of 0.0090 during validation and 0.0089 during testing
 258 compared to the GCN’s higher values of 0.0171 and 0.0169. Similarly, in the 50/25/25 split, the GAT
 259 maintained MAE values of 0.0010 and 0.0009 for validation and testing, while the GCN resulted in
 260 slightly higher values of 0.0028 and 0.0025.

261 We conducted experiments using multiple training batch configurations, and the results demonstrated
 262 an exceptional degree of similarity across the different settings. These results demonstrate the
 263 outstanding performance of both models, with GAT exhibiting a slightly enhanced performance in
 264 error minimization.

265 5 Conclusion

266 In this paper, we proposed PanRX, a novel multi-modal deep-learning pipeline for prediction of
 267 pancreatic cancer drug effects. Unlike previous work that used 2D molecular representations, we
 268 combined 3D molecular topology with textual information to comprehensively capture each pair of
 269 drugs, cell line information, and relationships. Specifically, embeddings were created via the PaiNN
 270 model for geometric data, the SciBERT model for textual data, and direct encoding for numerical
 271 data. We then constructed a multi-modal knowledge graph for 5014 drug-cell line pairs and employed
 272 a GAT and GCN model to predict IC50, AUC, and Z-scores. Preliminary experiments show that
 273 PanRX achieves high accuracies on both GNN architects, showcasing the importance of heterogeneity
 274 of modalities with the addition of 3D molecular conformation.

275 However, our model requires further consideration before advancing to the clinical level. Firstly,
 276 GDSC performed in vitro experiments which drastically differ from in vivo physiological conditions
 277 within an organism. In addition, mutation status is not necessarily binary due to the continuous
 278 spectrum of functional consequences. Lastly, PanRX is currently only trained and experimented on
 279 pancreatic cancer data which cannot be generalized to other diseases. To overcome these challenges,
 280 we plan to further our research with the following points: 1) Train and experiment our model with
 281 other cancer types found in the GDSC dataset to validate the effectiveness of multi-model knowledge
 282 graphs. 2) Transition to a large language model (LLM) due to its improved efficiency for handling
 283 large texts, as it can represent the same information conveyed by knowledge graphs in pure text form.
 284 3) Verify the model with clinicians on the patient level.

285 Acknowledgement

286 The authors would like to thank Divin Yan from Fudan University and Haorui Li from Huazhong
 287 University of Science and Technology (HUST) for their interesting discussions, suggestions, and data
 288 extraction help during this research.

289 References

290 [1] S. L. Carey-Smith, R. S. Kotecha, L. C. Cheung, and S. Malinge. Insights into the clinical,
 291 biological and therapeutic impact of copy number alteration in cancer. *International Journal of*
 292 *Molecular Sciences*, 25(13):6815, 2024. doi: 10.3390/ijms25136815.

- 293 [2] P. Unger Avila, T. Padvitski, A. C. Leote, H. Chen, J. Saez-Rodriguez, M. Kann, and A. Beyer.
294 Gene regulatory networks in disease and ageing. *Nature Reviews Nephrology*, 20:616–633,
295 2024. doi: 10.1038/s41581-024-00849-7.
- 296 [3] W. Yang, J. Soares, P. Greninger, E. J. Edelman, H. Lightfoot, S. Forbes, N. Bindal, D. Beare,
297 J. A. Smith, I. R. Thompson, et al. Genomics of drug sensitivity in cancer (gdsc): a resource for
298 therapeutic biomarker discovery in cancer cells. *Nucleic Acids Research*, 41(D1):D955–D961,
299 2013. doi: 10.1093/nar/gks1111.
- 300 [4] H. Askr, E. Elgeldawi, H. Aboul Ella, et al. Deep learning in drug discovery: an integrative
301 review and future challenges. *Artificial Intelligence Review*, 56:5975–6037, 2023. doi: 10.1007/
302 s10462-022-10306-1.
- 303 [5] M. Sun, S. Zhao, C. Gilvary, O. Elemento, J. Zhou, and F. Wang. Graph convolutional networks
304 for computational drug development and discovery. *Briefings in Bioinformatics*, 21(3):919–935,
305 2020. doi: 10.1093/bib/bbz042.
- 306 [6] B. Tang, Z. Pan, K. Yin, and A. Khateeb. Recent advances of deep learning in bioinformatics
307 and computational biology. *Frontiers in Genetics*, 10, 2019. doi: 10.3389/fgene.2019.00214.
- 308 [7] J. Kim, S. Park, D. Min, and W. Kim. Comprehensive survey of recent drug discovery using
309 deep learning. *Int. J. Mol. Sci.*, 22(18):9983, 2021. doi: 10.3390/ijms22189983.
- 310 [8] B. M. Kuenzi, J. Park, S. H. Fong, K. S. Sanchez, J. Lee, J. F. Kreisberg, and T. Ideker. Predicting
311 drug response and synergy using a deep learning model of human cancer cells. *Cancer Cell*, 38
312 (5):672–684.e6, 2020. doi: 10.1016/j.ccell.2020.09.014.
- 313 [9] N. Singh, P. Vayer, K. Tsaioun, B. O. Voutireix, S. Tanwar, and J.-L. Poyet. Drug discovery
314 and development: introduction to the general public and patient groups. *Front. Drug Discov.*, 3:
315 Article 101419, 2023. doi: 10.3389/fdds.2023.1021419. URL [https://doi.org/10.3389/
316 fdds.2023.1021419](https://doi.org/10.3389/fdds.2023.1021419).
- 317 [10] S. Aykul and E. Martinez-Hackert. Determination of half-maximal inhibitory concentration
318 using biosensor-based protein interaction analysis. *Analytical Biochemistry*, 508:97–103, 2016.
319 doi: 10.1016/j.ab.2016.06.025.
- 320 [11] J. D. Scheff, R. R. Almon, D. C. DuBois, W. J. Jusko, and I. P. Androulakis. Assessment of
321 pharmacologic area under the curve when baselines are variable. *Pharmaceutical Research*, 28
322 (5):1081–1089, 2011. doi: 10.1007/s11095-010-0363-8.
- 323 [12] J. P. Hughes, S. Rees, S. B. Kalindjian, and K. L. Philpott. Principles of early drug discovery.
324 *British Journal of Pharmacology*, 162(6):1239–1249, 2011. doi: 10.1111/j.1476-5381.2010.
325 01127.x.
- 326 [13] G. Obaido, I. D. Mienye, I. D. Emmanuel, A. Ogunleye, O. F. Egbelowo, P. Miene, and
327 K. Aruleba. Supervised machine learning in drug discovery and development: Algorithms,
328 applications, challenges, and prospects. *Machine Learning with Applications*, 17:100576, 2024.
329 doi: 10.1016/j.mlwa.2024.100576.
- 330 [14] Z. A. Rollins, A. C. Cheng, and E. Metwally. Molprop: Molecular property predic-
331 tion with multimodal language and graph fusion. *J Cheminform*, 16:56, 2024. doi:
332 10.1186/s13321-024-00846-9. URL <https://doi.org/10.1186/s13321-024-00846-9>.
- 333 [15] M. Galushka, C. Swain, F. Browne, et al. Prediction of chemical compounds properties
334 using a deep learning model. *Neural Comput Appl*, 33:13345–13366, 2021. doi: 10.1007/
335 s00521-021-05961-4. URL <https://doi.org/10.1007/s00521-021-05961-4>.
- 336 [16] S. Liu. Ai for molecule discovery with multi-modal knowledge. 2023.
- 337 [17] S. Liu, H. Wang, W. Liu, J. Lasenby, H. Guo, and J. Tang. Pre-training molecular graph represen-
338 tation with 3d geometry. 2022. URL <https://openreview.net/forum?id=xQe1pOKPam>.
- 339 [18] S. Liu, W. Du, Z. M. Ma, H. Guo, and J. Tang. A group symmetric stochastic differential
340 equation model for molecule multi-modal pretraining. In *International Conference on Machine
341 Learning*, pages 21497–21526. PMLR, 2023.

- 342 [19] W. Du, J. Chen, X. Zhang, Z. Ma, and S. Liu. Molecule joint auto-encoding: trajectory
343 pretraining with 2d and 3d diffusion. In *Proceedings of the 37th International Conference on*
344 *Neural Information Processing Systems*, pages 55077–55096, 2023.
- 345 [20] C. Edwards, T. Lai, K. Ros, G. Honke, K. Cho, and H. Ji. Translation between molecules and
346 natural language. *arXiv preprint arXiv:2204.11817*, 2022.
- 347 [21] S. Liu, W. Nie, C. Wang, J. Lu, Z. Qiao, L. Liu, J. Tang, C. Xiao, and A. Anandkumar. Multi-
348 modal molecule structure–text model for text-based retrieval and editing. *Nature Machine*
349 *Intelligence*, 5(12):1447–1457, 2023.
- 350 [22] W. R. Pearson. *Using the FASTA Program to Search Protein and DNA Sequence Databases*,
351 volume 24 of *Methods in Molecular Biology*, page 307. Humana Press, Totowa, NJ, 1994.
352 ISBN 978-0-89603-246-0. doi: 10.1385/0-89603-246-9:307.
- 353 [23] A. Serra, P. Coretto, M. Frattello, and R. Tagliaferri. Robust and sparse correlation matrix
354 estimation for the analysis of high-dimensional genomics data. *Bioinformatics*, 34(4):625–634,
355 2018. doi: 10.1093/bioinformatics/btx642.
- 356 [24] A. Kroll, S. Ranjan, M. K. M. Engqvist, et al. A general model to predict small molecule
357 substrates of enzymes based on machine and deep learning. *Nature Communications*,
358 14:2787, 2023. doi: 10.1038/s41467-023-38347-2. URL [https://doi.org/10.1038/
359 s41467-023-38347-2](https://doi.org/10.1038/s41467-023-38347-2).
- 360 [25] A. Rossi, D. Barbosa, D. Firmani, A. Matinata, and P. Merialdo. Knowledge graph embedding
361 for link prediction: a comparative analysis. *ACM Transactions on Knowledge Discovery from*
362 *Data (TKDD)*, 15(2):14, 2020. doi: 10.1145/3424672. URL [https://doi.org/10.1145/
363 3424672](https://doi.org/10.1145/3424672).
- 364 [26] M. Wang, L. Qiu, and X. Wang. A survey on knowledge graph embeddings for link prediction.
365 *Symmetry*, 13(3):485, 2021. doi: 10.3390/sym13030485. URL [https://doi.org/10.3390/
366 sym13030485](https://doi.org/10.3390/sym13030485).
- 367 [27] J. Vilela, M. Asif, A. R. Marques, J. X. Santos, C. Rasga, A. Vicente, and H. Martiniano.
368 Biomedical knowledge graph embeddings for personalized medicine: Predicting disease-gene
369 associations. *Expert Systems*, 2022. doi: 10.1111/exsy.13181.
- 370 [28] Z. Gao, P. Ding, and R. Xu. Kg-predict: A knowledge graph computational framework for drug
371 repurposing. *Journal of Biomedical Informatics*, 132:104133, 2022. doi: 10.1016/j.jbi.2022.
372 104133.
- 373 [29] S. Liu, M. Qu, Z. Zhang, H. Cai, and J. Tang. Structured multi-task learning for molecular
374 property prediction. In *International conference on artificial intelligence and statistics*, pages
375 8906–8920. PMLR, 2022.
- 376 [30] S. Liu, Y. Li, Z. Li, Z. Zheng, C. Duan, Z. M. Ma, O. Yaghi, A. Anandkumar, C. Borgs,
377 J. Chayes, et al. Symmetry-informed geometric representation for molecules, proteins, and
378 crystalline materials. *Advances in neural information processing systems*, 36, 2024.
- 379 [31] K. Schütt, O. Unke, and M. Gastegger. Equivariant message passing for the prediction of
380 tensorial properties and molecular spectra. In *International Conference on Machine Learning*,
381 pages 9377–9388. PMLR, 2021.
- 382 [32] M. Nakata, T. Shimazaki, M. Hashimoto, and T. Maeda. Pubchemqc pm6: Data sets of 221
383 million molecules with optimized molecular geometries and electronic properties. *J. Chem. Inf.*
384 *Model.*, 60:5891–5899, 2020. doi: 10.1021/acs.jcim.0c00740. URL [https://doi.org/10.
385 1021/acs.jcim.0c00740](https://doi.org/10.1021/acs.jcim.0c00740).