

Building A Unified Model for Drug Synergy Analysis Powered by Large Language Models

Tianyu Liu^{1,2}, Tinyi Chu², Xiao Luo³, Hongyu Zhao^{1,2*}

¹Interdepartmental Program in Computational Biology & Bioinformatics, Yale University, New Haven, 06511, CT, USA.

²Department of Biostatistics, Yale University, New Haven, 06511, CT, USA.

³Department of Computer Science, University of California, Los Angeles, Los Angeles, 90095, CA, USA.

*Corresponding author(s). E-mail(s): hongyu.zhao@yale.edu;

Contributing authors: tianyu.liu@yale.edu; tinyi.chu@yale.edu; xiaoluo@cs.ucla.edu;

Abstract

Drug synergy prediction is a challenging and important task in the treatment of complex diseases including cancer. In this manuscript, we present a novel unified Model, known as BAITSAO, for tasks related to drug synergy prediction with a unified pipeline to handle different datasets. We construct the training datasets for BAITSAO based on the context-enriched embeddings from Large Language Models for the initial representation of drugs and cell lines. After demonstrating the relevance of these embeddings, we pre-train BAITSAO with a large-scale drug synergy database under a multi-task learning framework with rigorous selections of tasks. We demonstrate the superiority of the model architecture and the pre-trained strategies of BAITSAO over other methods through comprehensive benchmark analysis. Moreover, we investigate the sensitivity of BAITSAO and illustrate its unique functions including new drug discoveries, drug combinations-gene interaction, and multi-drug synergy predictions.

Keywords: Scaling Laws, Large Language Model, Multi-Task Learning, Transfer Learning, Drug Synergy Prediction, Cancer Genomics

1 Introduction

Treating patients with a combination of drugs has become common for various diseases, including HIV [1] and cancers [2, 3]. One key aspect of drug combinations is the synergistic effect, which means that the joint effect of multiple drugs is larger than the sum of individual drug effects [4]. Other definitions have also been used to define synergistic effects, such as [5]. Effective drug combination can reduce drug resistance of monotherapy [6] with relatively lower doses of individual drugs [7]. Since drugs can change gene expressions when applied to different systems, e.g., cell lines, their effects can be studied through the genomics lens [8, 9]. Currently, researchers use high-throughput combinatorial screening to identify drug combinations with synergistic effects for specific cell lines [10]. However, such experimental screening is laborious and time-consuming due to the very number of potential drug combinations, and it is even more challenging to assess the synergistic effect of combinations with three or more drugs [11]. Therefore, it is important to develop computational methods based on extensive experimental datasets in the public domain as well as diverse types of prior biological knowledge to predict the presence and strength of synergistic effects for candidate drug combinations. Accurate prediction methods can facilitate drug discovery [12] and clinical development [13].

Given its importance, it is no surprise that many machine learning methods, especially deep learning methods, have been proposed to predict drug synergy. These methods differ in model architecture, training strategies, and datasets used to build the models. DeepSynergy [14] is among the earliest tools by building a neural network for both regression and classification, with follow-up work such as TreeComb [15, 16] and MatchMarker [17]. Existing drug synergy prediction methods can be broadly classified into two groups. The first group of methods, such as MARSY [18], focus on predicting specific synergy scores, whereas the second group of methods, such as DeepDDs [19], transfer the continuous synergy score into a binary one via thresholding to infer drug combination synergy. However, most existing methods do not incorporate the extensive synergy information from public databases [20] in their predictions. [21] utilized a transfer learning approach and pre-trained the model based on large-scale databases while incorporating different types of features (e.g., gene expression, molecular structure). However, it did not consider datasets [14] with only partial information and treated drugs with the same molecular formula but different names as distinct ones. Therefore, the generalization ability of this model is limited by its input data format. Moreover, since public databases are updated constantly, it is important to track the versions of training datasets.

Large Language Models, as a type of Foundation Models (FMs) [22], have greatly improved the performance of deep learning on various tasks in Natural Language Processing (NLP) [23]. Such models have received broad attention from both industry and academia [24]. Researchers have proposed to use FMs to predict drug synergy via LLMs by transferring the drug synergy prediction problem into a Question-Answer problem [25, 26]. By incorporating prior information of single drug and single cell line from LLMs, it has become possible to predict drug synergy of unknown drug combinations in unknown cell lines. Text information may be less noisy than the features

(e.g. gene expression levels) that have been used in this task. However, such QA setting limits the task to a classification problem, which introduces the potential bias of pre-defined thresholds. Moreover, these two LLMs are not open-source so it is difficult for researchers to evaluate their performances. Open-source is important for the development of science [27]. Moreover, there is a lack of exploration on the utilization of the information in LLMs for more difficult drug synergy prediction problems, e.g., the effects of multiple drug combinations or model explainability.

Here we present a scalable unified model for drug synergy prediction called BAITSAO¹. BAITSAO utilizes the information from LLMs as input and was pre-trained based on large-scale known synergistic effect information of paired drug combinations and cell lines. The information of drug combinations and cell lines is necessary to predict synergy scores. We show that the embeddings of these features from LLMs can be effective input for drug synergy prediction as well as the effects of drugs on gene expressions. We further demonstrate the capability of building an effective predictor for synergy prediction under both the classification and regression settings through multi-task learning (MTL) [28]. Finally, we pre-train BAITSAO to predict synergistic effects for unseen drug combinations based on the zero-shot learning framework and the fine-tuning framework. The scalability of BAITSAO allows us to consider multiple drugs and incorporate extra meta information.

2 Results

Overview of BAITSAO. We highlight two major contributions of BAITSAO as a unified model. We first provide a new unified pipeline for pre-processing the information from both drugs and cell lines for machine learning in a tabular format, and generate training datasets from these embeddings for multiple tasks. We show that these embeddings contain functional information for prediction. We then utilize the unified training datasets for different synergistic effect prediction tasks under the multi-task learning framework. We demonstrate the superiority of the model architecture and the contribution of pre-training through comprehensive experiments. BAITSAO can be easily transferred to perform novel downstream tasks related to drug synergy analysis. We illustrate the landscape of BAITSAO in Figures 1 (a) and (b) and summarize the differences between BAITSAO and other synergy prediction methods in Figure 1 (c). The major functions of BAITSAO are shown in Figure 1 (d).

Drug embeddings from LLMs reflect functional similarity and responses at the cell level. In this section, we discuss the information offered by drug embeddings and cell-line embeddings. We generate the description for the drugs and cell lines from our training datasets based on designed prompts from LLMs, and then use the embedding module from GPT 3.5 [29] to generate the embeddings of such descriptions, where the embeddings become the features of drugs or cell lines. We utilized GPT 3.5 rather than GPT 4 [30] because the layer for generating embeddings is from GPT 3 [29] series and the querying time from GPT 4 with similar quality required much more time [31], and efficiency is very important in LLM deployment [32, 33]. Moreover, the performance difference between embeddings from GPT 3.5 and GPT

¹BAITSAO means collections of herbs and drugs in Chinese culture.

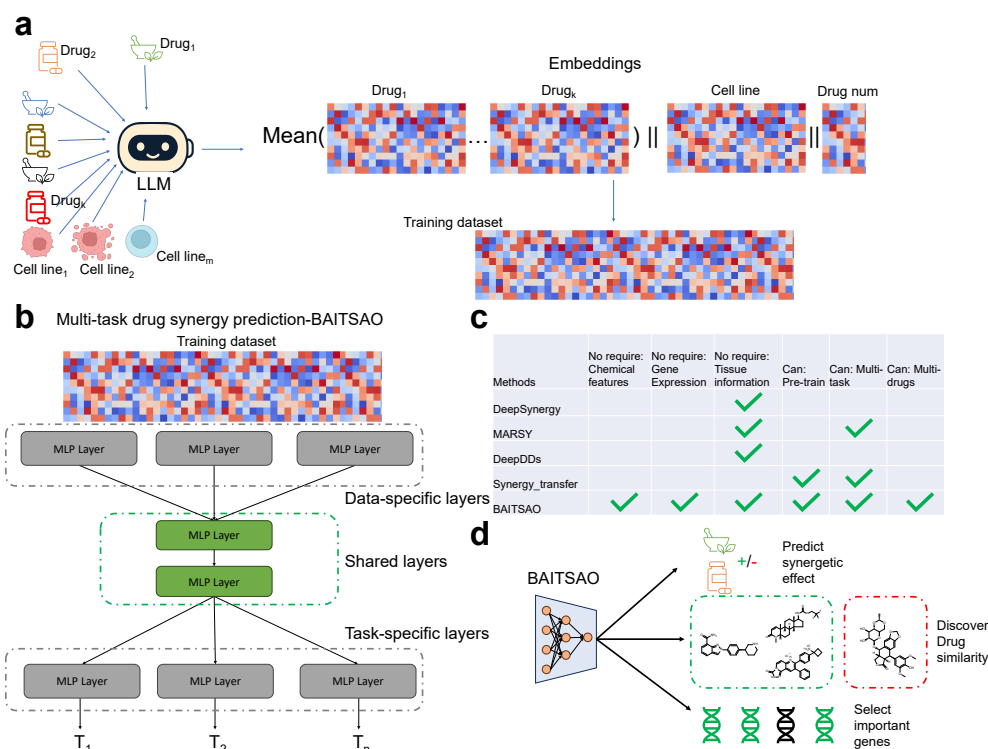


Fig. 1 An overview of BAITSAO as a FM under the pre-training and fine-tuning/zero-shot learning pipeline. (a) The pre-processing steps we used to transfer the meta information into embeddings to construct training datasets. (b) The model architecture of BAITSAO under a multi-task learning framework. (c) Comparisons of different methods for drug synergy analysis. (d) Different functions of BAITSAO.

4 or from GPT 3.5 and Claude 3.5 [34] is not significant based on our experiments, shown in Supplementary Figure 1 (a) (Wilcoxon rank-sum test, p-value=0.86 for GPT 3.5 vs. GPT 4, and p-value=0.44 for GPT 3.5 vs. Claude 3.5). In the same figure, we also found that embeddings from GPT 3.5 are better than embeddings from Gemini [35] (p-value=0.0039), and thus our current selection is well-designed. We visualize the drug embeddings and the cell-line embeddings based on Uniform Manifold Approximation and Projection (UMAP) [36] shown in Supplementary Figures 2 (a) and (b). We investigated the quality of the embeddings by considering both the quality of the description and the quality of the functions of embeddings.

For the first aspect, we recorded the outputs as descriptions from GPT 3.5 based on our prompts and compared the content with information from DrugBank [37] and NCBI [38]. Here we used drugs and cell lines from DeepSynergy, which contains 39 drugs and 38 cell lines. The descriptions summarized the functional information of drugs and cell lines. Based on our experiments, only one drug (MK-8669) has a mismatched generated description, while 13 drugs cannot be matched with indication

information if we search them in DrugBank. All descriptions are included in Supplementary file 1. We plot the Cosine Similarity (CS) for all drugs' embeddings in Figure 2 (a). We also randomly selected 10 drugs from this dataset and plot the CS for the embeddings of the same drug under 10 different descriptions by running GPT 3.5 multiple times in Supplementary Figure 3. These two figures show that the similarity from different drugs is generally lower than that from the same drug, suggesting that we can get informative embeddings from LLMs.

To perform a comprehensive analysis of our generated drug embeddings from LLMs, we downloaded the descriptions of drugs, including indication, summary, and background, from the DrugBank. We embedded these descriptions based on the same GPT-3.5 embeddings layer and computed the CS between embeddings from DrugBank descriptions and the LLM-generated descriptions. We found that embeddings from LLMs have a strong average similarity with all three descriptions from DrugBank (CS=0.87 for indication, CS=0.90 for summary, and CS=0.90 for background), and thus the generated drug embeddings preserved the important functional and chemical properties of the original drug. Furthermore, we visualize the CS based on the embeddings from drug indication (Supplementary Figure 4 (a)), drug summary (Supplementary Figure 4 (b)) and drug background (Supplementary Figure 4 (c)). We further computed the Pearson Correlation Coefficient (PCC) between the similarity matrix from DrugBank descriptions and LLM descriptions, which could be used to evaluate the ability of embeddings used by BAITSAO in preserving the drug-drug similarity. The PCCs are annotated under each figure, and all of the PCCs are high ($PCC \geq 0.76$) and significant ($p\text{-value} < 0.05$). Therefore, we demonstrated the ability of LLMs to generate meaningful descriptions as well as embeddings by comparing the generated information with known database, and further enhanced the reliability of the pipeline.

Furthermore, we performed Mann-Whitney U test [39] to compare the PCCs among the drugs from the MK class and the PCCs between the drugs from the MK class and other classes, and the test statistics showed a significant difference ($p\text{-value} = 9.9e-12$). Therefore, in Figure 2 (b), we used drug MK-4541 as one example and there is no clinical information for this drug in the DrugBank, to infer its function based on our embeddings. By excluding the drug MK-8669 due to mismatched information, drug MK-2206, and drug MK-4827 have the highest similarity with MK-4541. Since MK-2206 and MK-4827 have similar functions (e.g., treating Breast-cancer-related and Prostate-cancer-related diseases), we may infer that MK-4541 may have a similar effect. Among these drugs, EPTOPOSIDE has the lowest similarity and it also has different clinical trial information, suggesting that correlation between embedding similarity and function similarity. Therefore, our drug embeddings may help the inference of clinical functions of drugs based on the embeddings' similarity.

To investigate whether the embeddings can be used to predict drug response for the cell-level task, we utilized CPA [9] and single-cell RNA sequencing (scRNA-seq) [40] datasets with different perturbations (defined by different drugs or drug combinations) to evaluate whether our drug embeddings can facilitate the gene expression prediction task. With drug embeddings, we can use CPA to predict gene expression response

to unseen drugs. Cells with unknown perturbation results are also known as out-of-distribution (OOD) samples. The original implementation of CPA utilized the drug embeddings from Rdkit [41, 42] to encode the molecular structure of the selected drug into the embedding space. However, such methods could not handle drugs not in the Simplified Molecular-input Line-entry System (SMILES) [43], which limits the generalization of CPA. Here we considered replacing the original embeddings in CPA with the embeddings from GPT 3.5, enlarging the accessibility for drug embeddings. We compared three different embedding settings for two datasets (CPA example [9] and Openproblems [44]), which contain the gene expression profiles under the control case and drug-based perturbations. The results are shown in Supplementary Figures 5 (a)-(d), where stacking the embeddings from SMILES and GPT 3.5 achieved the best performance under both datasets. For the CPA dataset, both using the embeddings from GPT 3.5 and the setting of embeddings stacking can enhance the prediction performance significantly, compared with the mode of only using SMILES (Wilcoxon rank-sum test, p-values<0.05). For the Openproblems dataset, the contribution of such embeddings stacking for prediction is especially significant (p-values<0.05). Therefore, the embeddings from LLMs can improve the gene expression prediction for perturbed scRNA-seq data.

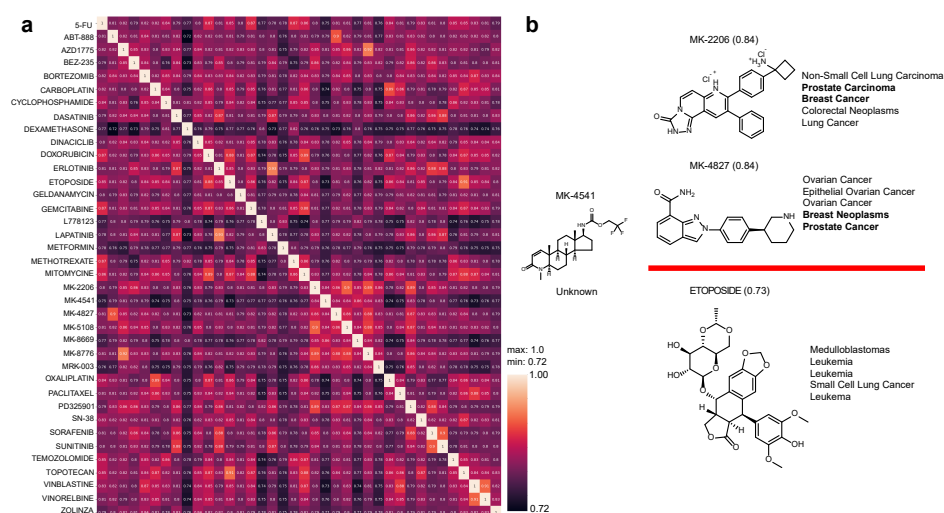


Fig. 2 Investigation of drug embeddings. (a) The heatmap for the similarity of embeddings across all the drugs. (b) Exploration of drug similarity related to MK-4541. The drugs above the red line represent the two most similar drugs, while the drugs below the red line represent the most different drug. We list five types of clinical trial information ranked by the phases. Source data are provided as a Source Data file.

Since our experiments demonstrate that drug embeddings and cell embeddings can summarize the functional information and drug embeddings can also interact with

cell-level gene expressions, we believe that these embeddings allow us to construct the training dataset to predict drug synergy effect in different cell lines.

Demonstration of powerful embeddings and architecture by evaluation without pre-training. In this section, we show the strength of LLM embeddings and select the choice of network structure for model pre-training based on two different drug synergy prediction tasks: classification and regression. For each task, we selected two datasets and two metrics for evaluation. For regression, we included the Pearson Correlation Coefficient (PCC) and Mean Squared Error (MSE) for model evaluation based on datasets D1 [14] and D2 [18]. For classification, we included ROC-AUC (ROCAUC) and Accuracy (ACC) for model evaluation based on datasets D1 and D3 [19]. These metrics and datasets were widely used in the related work [14, 18, 19, 45, 46]. First, we tested if the model performance will be affected by prompt engineering of LLMs, and we compared the raw embeddings with embeddings generated by drug descriptions from MetaPrompt [47] and Chain-of-Thought (COT) [48]. According to Supplementary Figure 1 (b), the differences between the default mode and these two prompt engineering methods are not significant (p-value=0.63 for raw mode vs. MetaPrompt, and p-value=0.43 for raw mode vs. COT). Therefore, our embeddings have enough information as inputs for synergetic effects. Second, we validated the contribution of BAITSAO's architecture shown in Supplementary Figure 1 (c). We compared the performances between BAITSAO and DeepSynergy with LLM embeddings as inputs. The difference is significant and thus our optimization of model architecture also contributed to the prediction task (p-value=0.002). Finally, we selected seven other methods (DeepSynergy, MARSY, TreeComb, SVM [39, 49], TabNet [50], BERT [51] and Lasso [39, 52]) for benchmarking the regression task and seven methods (DeepSynergy, DeepDDs, TreeComb, SVC, TabNet, BERT, and Lasso) for benchmarking the classification task. We utilized the best hyper-parameters of these methods for every dataset, with details of hyper-parameter tuning summarized in the Methods section. Our results based on five-fold cross-validation [14] are summarized in Figure 3. This figure shows that BAITSAO ranked the best in three out of four metrics. Moreover, BAITSAO was also the most stable among the top deep-learning-based methods (including MARSY, DeepSynergy, DeepDDs, and TabNet). The performance of BERT was worse than BAITSAO in three out of four metrics, thus using embeddings as input is better than using the combination of description in general. For the evaluation based on MSE, BAITSAO performed well on the D1 dataset. Our experiments showed that embeddings from LLMs with a suitable model architecture can formalize a better training-testing framework compared with data from the classical feature space. The details of our dataset information, model construction, and training process are summarized in the Methods section.

Explainability of BAITSAO for drug-gene interaction and drug-cell line interaction with multi-modal learning. We interpret contributions of different features for the prediction task with the help of SHAP [53]. Here we integrated known gene expression profiles of different cell lines in D3 to our input datasets and performed the same training process for the drug synergy prediction task. We then utilized SHAP to study the importance of different genes and the results could be treated as the relevance between the gene expression levels (as a new modality) and the possibility

239 to produce synergistic effect for drug combinations. We followed the default setting of
240 SHAP to fix the number of genes for explainability at 20. We also performed statistical
241 analysis based on the outputs of BAITSAO to discover the drug combination with the
242 largest range of synergistic targets. The details of our approach are provided in the
243 Methods section.

244 By collecting gene expression profiles of cell lines [54], we studied the explainability
245 of BAITSAO for DEXAMETHASONE (drug)-DINACICLIB (drug) across different
246 cell lines. In Figure 3 (b), we visualize the importance of different genes. The gene
247 VIM was top-ranked by the average importance, and VIM is known as important
248 for various cancers from pan-cancer analysis [55]. Furthermore, we conducted three
249 experiments to further investigate the contributions of the selected genes.

250 We first separated the samples into two groups based on the existence of the
251 synergistic effect and performed DEG analysis using DESeq2 [56, 57] between two
252 groups of cell lines. We present the adjusted p-values using Benjamini-Hochberg for
253 the selected genes in Figure 3 (b). Genes SPON2, HMCN1, and BMP4 listed in this
254 figure were significant DEGs. Our selected genes had significant overlap with DEGs
255 (Fisher's exact test p-value=0.0062), and the gene BMP4 is a validated targets for each
256 drug according to biology experiments [58, 59]. Moreover, these genes had relatively
257 lower expression levels with synergistic effect, which matched the distribution of their
258 SHAP values (enriched in the negative values).

259 We list the ranks of the selected genes based on their variances in Figure 3 (b).
260 The top five ranked genes had relatively greater variances, suggesting that the genes
261 we selected characterize the heterogeneity from both cell lines and the drug synergistic
262 effect. We further performed enrichment analysis based on Gene Ontology (GO) [60]
263 for biological pathways and Molecular Signature Database (MsigDB) [61] for cancer-
264 specific signals based on this set of genes, with results shown in Supplementary Figures
265 6 (a) and (c). These enriched pathways represent important biological processes and
266 cancer-specific signals. These results suggest that our method may uncover the het-
267 erogeneity in the drug synergy prediction process. We summarize our results for the
268 single cell line with the same drug combination in Appendix A. The plots for impor-
269 tant genes across different cell lines can be found in Supplementary Figure 16. The
270 test statistics used in this section are given in Supplementary file 2.

271 We further investigated the drug combination that showed synergistic effects on
272 the largest number of cell lines. We first plot the probabilities of all drug-cell line com-
273 binations to be classified as samples with synergistic effects in Figure 3 (c). This figure
274 shows that the distribution of such probability is different under different synergistic
275 labels. We performed the Rank-sum test [62] for these two sets of probabilities and
276 their difference is significant (p-value<2.22e-308 with two-side mode and no adjust-
277 ment is needed). Therefore, our model can uncover the relationship between input
278 features and the synergistic effect. Moreover, we ranked the drug combinations based
279 on the number of cell lines predicted to have synergistic effects in descending order.
280 We computed the Pearson correlation coefficient [62] between the count value based
281 on predicted labels and observed labels, summarized in Figure 3 (d). Based on this
282 figure, the count values based on the predicted labels had a strong positive correlation
283 with those based on the real labels, thus our model can also be used to identify the

drug combinations with the most synergistic targets given a set of cell lines and cancer types used in our experiments. We also highlight the drug combination L778123 and MK-8669 that has the largest number of targeted cell lines with synergistic effect in Figure 3 (d). The p-value is computed with two-side mode and no adjustment is needed. Therefore, BAITSAO can capture the variance of different drug combinations across cell lines, offering a novel option for selecting effective drug combinations.

Statistics of pre-training datasets. Here we summarize the statistics and properties of our pre-training datasets for BAITSAO. We collected information from DrugComb [20], which is known as the largest database containing synergistic effect information for drug pairs with different cell lines. We downloaded the updated version of DrugComb and removed the missing value or single-drug information. The major statistics of DrugComb are summarized in Figure 4, whose (a) represents the total number of *drug-cell line* combinations by tissue types and Figure 4 (b) represents the total number of *cell lines* by the type of tissues from DrugComb. Most of the drugs presented here were analyzed using cells from skin, lymphoid, and/or lung. These tissues are important for maintaining normal physiological activity in the human body. In total, DrugComb collects more than 700,000 available combinations. As shown in Figure 4 (c), the distribution of the synergy scores is not balanced, with a large number of combinations having low synergy scores. We further plot the Half Maximal Inhibitory Concentration (IC₅₀) for all drugs in Figure 4 (d) with a similar distribution to the synergy score. We illustrate the non-linear relationship between single drug IC₅₀ and synergy score in Figure 4 (e). Therefore, fitting non-linear models like neural networks may help the synergy prediction task. Finally, Figure 4 (f) shows the overlap of combinations by tissues, where most tissues have low overlap, and thus the pre-training dataset has information from diverse tissues. We plot the embeddings for drugs and cell lines in the pre-training dataset colored by clusters from Leiden [63] in Supplementary Figures 7 (a) and (b). The items in the same Leiden cluster can be treated in similar context of embeddings with functional information, so we can visualize the functional similarity of different drugs and cell lines through embeddings. Since our pre-training dataset was published in June 2021 and GPT 3.5 collected data for pre-training until Sep 2021, considering the time needed for pre-training a LLM, we believe that the data from DrugComb were precluded in GPT 3.5.

Pre-trained BAITSAO contributes to drug synergistic effect prediction under the multi-task condition. Here we investigated and pre-trained BAITSAO based on the optimal model structure. Specifically, we extended the model structure with a multi-task learning framework. By pre-training BAITSAO with large-scale synergy datasets, BAITSAO is able to predict both single drug inhibition and drug synergistic effect. For drug pairs, we expect to predict both drugs' inhibition, thus we have a total of four tasks inspired by the pre-training datasets, including the regression task for synergy prediction, the classification task for synergy prediction, and regression tasks for single-drug inhibition of each drug in the drug pairs. For the regression task of synergy prediction, we only considered predicting the synergy score under the Loewe setting because we show that the synergy scores computed based on other methods are positively correlated with the Loewe score [64] in Supplementary Figure 8, and literature [14, 25] suggests using the threshold for generating a classification task from

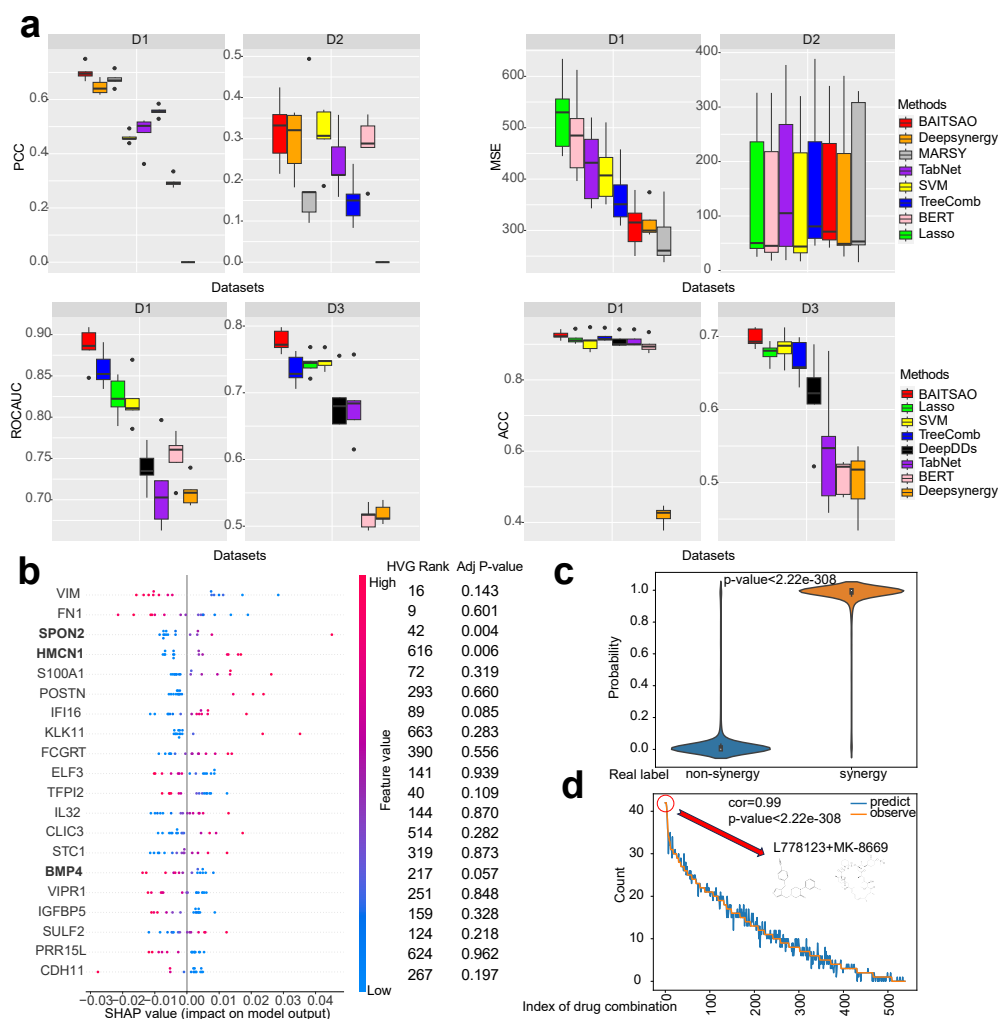


Fig. 3 Results of evaluations for model structure, reliability and explainability. (a) Evaluations for BAITSAO and other methods. Each panel represents one metric with two datasets. The ranks of methods are averaged by datasets. Data are presented as boxplots (n=5 per group; center line, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers). The explanations of datasets D1-D3 are summarized in the Methods section. (b) The explainability of BAITSAO for the combination DEXAMETHASONE (drug)-DINACICLIB (drug) for different cell lines. We also list the ranks based on variance (HVG rank) and the adjusted p-value based on DESeq2 analysis results for each gene. The genes with adjusted p-value for multiple comparisons smaller or close to 0.05 are boldfaced. (c) The violin plot (n=6299 for non-synergy group; n=6116 for synergy group; center point, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers) for the outputs of BAITSAO (Probability) across the synergistic labels. We also present the two-side p-value in this figure. This panel supports the reliability of selected features from SHAP. (d) The rank-based plot between the number of drugs-cell line combinations with synergy (Count) and the index of drug combination (Index of drug combination). The index is ranked by the value of Count. We present the Pearson correlation (corr) and two-side p-value in this figure. Source data are provided as a Source Data file.

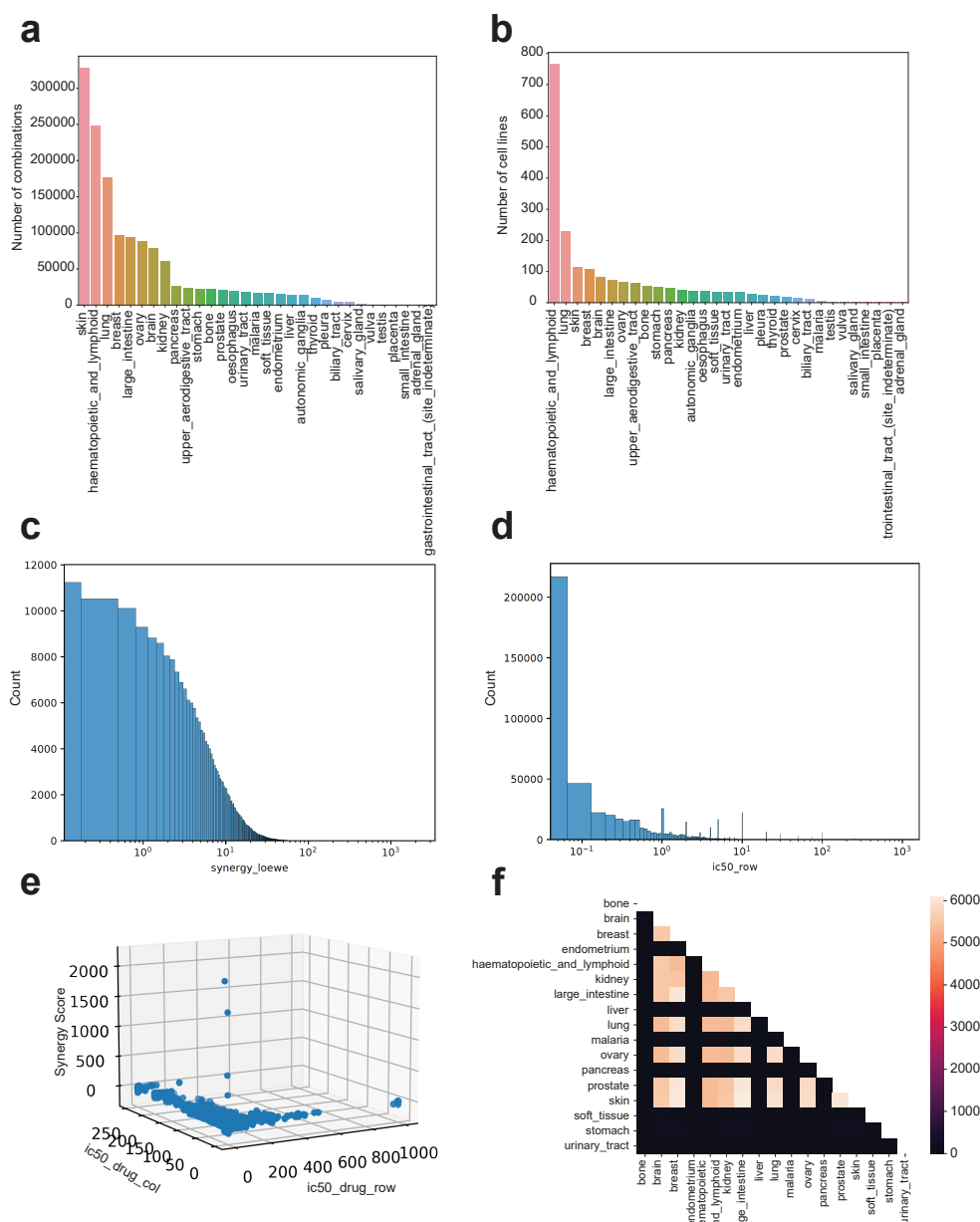


Fig. 4 Statistics of the pre-training dataset from DrugComb. (a) The barplot for the number of *drug-cell lines* combinations by different tissues. (b) The barplot for the number of *cell lines* by different tissues. (c) The histogram for the distribution of synergy score computed based on Loewe [64]. The x-axis is transferred into log scale. (d) The histogram for the distribution of single-drug IC₅₀ levels. The x-axis is transferred into log scale. (e) 3D plot for the relation between synergy score and synergy score. (f) The heatmap for the overlap of combinations across different tissues. Source data are provided as a Source Data file.

the Loewe score. For other synergy scores including Zip score [65], HSA score [66], and Bliss score [67], we pre-trained specific models and restored the pre-training weights. Instead of using the simple average of loss functions from different tasks during the training process, we introduced the Uncertainty Weighting (UW) method [68] advocated by the performance evaluations of different multi-task learning strategies from [69] and improved the numerical stability and the validation strategy of this method.

We first determined the tasks that can help each other in the multi-task learning framework by constructing the Help-Harm matrix. We sampled 1% of the pre-training dataset and trained task-specific models as well as multi-task models with paired tasks, and constructed the Help-Harm matrix shown in Figure 5 (a). According to this figure, joint training always boosts the classification task, while joint training with the classification task can help predict the synergy scores as well as inhibition levels for a single drug. Moreover, the relative inhibition (RI) information from one of the drugs in drug pairs did not show a significant contribution to other tasks, and incorporating this information reduced performance for the classification task. Since we had RI levels for both drug pairs, we removed the information of RI_col in the training process, and collected three tasks in the pre-training stage. After finishing pre-training based on the sampled and full datasets, we plot the metrics for comparing the performance between BAITSAO under the STL framework and our final MTL framework in Figure 5 (b). MTL can improve the performance of BAITSAO for solving all regression-based tasks. We show the outputs from the hidden layers of BAITSAO by ground truth synergistic labels and predicted synergistic labels in Supplementary Figures 9 (a) and (b). According to these two figures, the learned drug embeddings for drugs with no synergistic effect tended to be co-embedded. Therefore, BAITSAO with the MTL framework is reasonable and superior in drug synergy analysis. Finally, we consider the generalization ability of BAITSAO with pre-trained weights. We conducted experiments based on three datasets we used in the subsection *Selection of the model structure by evaluation without pre-training* and visualized the results in Figure 5 (c). We report the metrics based on five-fold cross-validation results. According to this figure, BAITSAO with the pre-training design after fine-tuning (BAITSAO-FT) is comparable or better for the regression and classification tasks, compared with BAITSAO without pre-training (BAITSAO-ZS). When evaluating the ZS mode, we ensured that the combinations used in the pre-training stage are not used for testing. Moreover, our fine-tuning stage used fewer epochs and we froze the shared layers during the fine-tuning process, thus our fine-tuning approach was more efficient. We note the potentials of BAITSAO under the zero-shot learning framework for solving this task. For example, BAITSAO-ZS showed a high ACC score in the evaluation based on D1. Moreover, for the metrics related to classification, BAITSAO-ZS had results higher than 0.5, and thus BAITSAO under the zero-shot learning framework was better than random guessing. We also performed Rank-sum tests [62] between the pre-training dataset and fine-tuning datasets and the results are shown in Supplementary Figure 10, which demonstrated that samples in the fine-tuning datasets satisfied the OOD cases. Finally, we compared BAITSAO with other LLM-based models, discussed in Appendix B, which shows that BAITSAO also has unique advantages. We also pre-trained other deep-learning-based synergy predictors, such as DeepSynergy,

DeepDDs, and MARSY based on their designed tasks and compared the fine-tuned version of these models with BAITSAO (ft). According to Supplementary Figure 11, BAITSAO still shows better performances than other baselines with either fine-tuning mode or from-scratch mode. Therefore, the multi-task pre-training strategy of BAITSAO is unique and contributive, which leads to consistent improvement across different datasets. In summary, the combination of MTL and the pre-training process can improve the performance of BAITSAO on tasks related to drug synergy analysis.

We then predicted the synergistic effect for the combination of three drugs (tri-drugs) and one cell line, with two examples shown in Figures 5 (d) and (e). The drug names and cell line names were extracted from DrugCombDB [71], which did not provide the observed synergistic information for the existing combinations. To enhance the reliability of our prediction results, we relied on Monte Carlo Dropout (MC Dropout) [72, 73] and ran inference 100 times to generate the prediction interval of different drug combinations. According to [74], MC Dropout was the only method considered in this benchmarking paper to estimate the mean and variance without extra hyper-parameters. Our full prediction results are summarized in Supplementary file 3. Here we compared the difference between the two combinations by changing the third drug. We found that the combination with I-BET151 was predicted to have a positive sign in the synergy score under Loewe, while the combination with I-BET was predicted to have a negative synergistic effect. As an explanation, although these two drugs can both combine with bromodomain and extra terminal domain (BET) with the same major targeted proteins [75, 76], I-BET151 was reported as an optimized version with excellent BET target potency and selectivity [76]. Therefore, we expected I-BET151 to have better efficacy and thus a higher synergy score. Another example from Figure 5 (e) presents the difference between PF562271 [77] and Saracatinib [78] as a third drug under the cell line MZ7-mel. The combination with PF562271 had a higher predicted synergy score compared with Saracatinib, which was supported by the experimental results from [79] as PF562271 generated higher growth inhibition. Therefore, the results from BAITSAO can help researchers to optimize drugs with higher synergistic effect and better clinical outcomes.

Sensitivity analysis. Here we investigated the sensitivity of model training based on the statistics we collected. Figure 6 (a) displays the ablation results by considering different types of embeddings as well as different types of combination rules for embeddings as model input. *BAITSAO* denotes our final choice for pre-training and fine-tuning. *BAITSAO-v3* denotes that we utilized the updated embeddings from OpenAI in 2024 [80]. *Mean* denotes that we took the mean of drug embeddings and cell embeddings as input for training. *Sum* denotes that we took the sum of drug embeddings and cell embeddings as input for training. *SentStack* denotes that we stacked the descriptions of different drugs and used the modified description to generate drug embeddings, and then stacked such drug embeddings with cell-line embeddings. *Stack* represents that we stacked the drug embeddings and cell embeddings by rows. *Rdkit* [41, 42] represents that we generated embeddings from Rdikt with SMILES and stacked the embeddings with cell embeddings from LLMs. This figure shows that averaging the drug embeddings and stacking them with cell embeddings by rows generated the best performance for all tasks. These results suggest the most effective way to incorporate

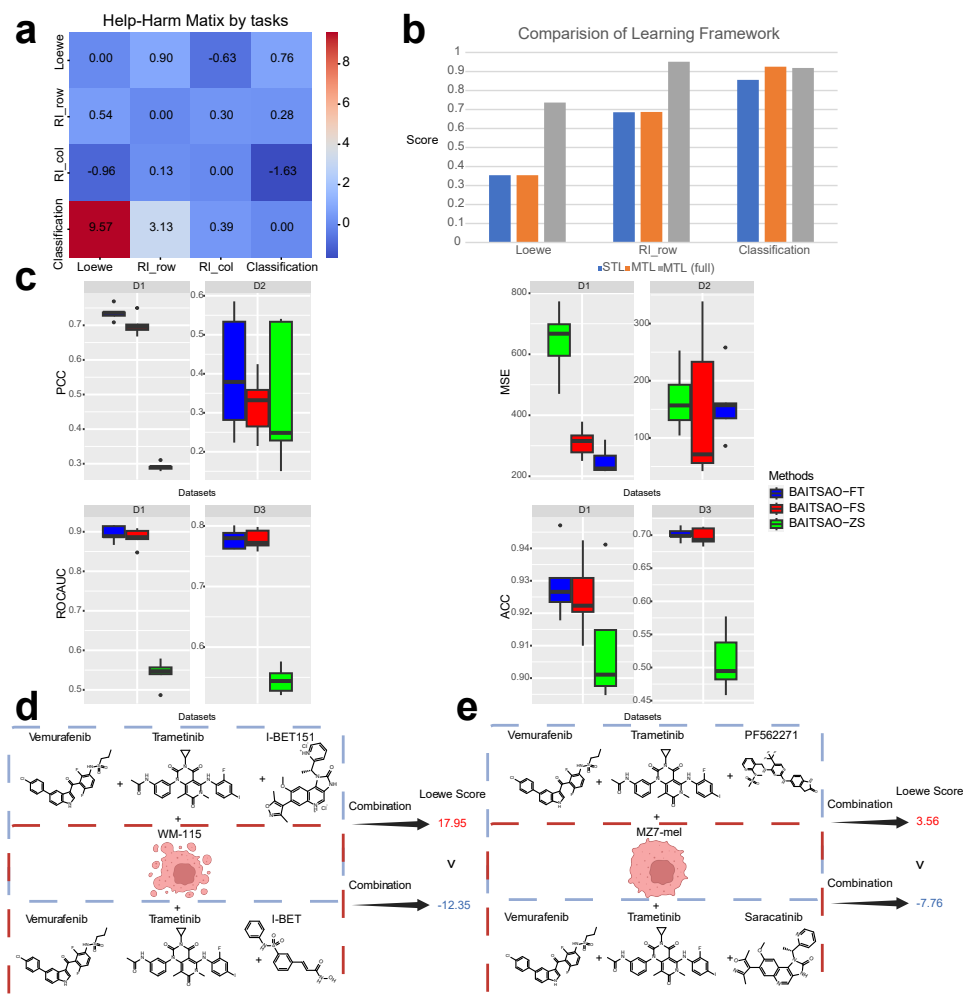


Fig. 5 Results under the multi-task learning framework. (a) The Help-Harm matrix for different combinations of tasks. The values indicate the percentage (unit: %) of improvement using multi-task learning compared to single-task learning (STL) defined by the tasks in rows. The columns represent the paired tasks. We boldfaced blocks with increments larger than 0.5%, which is a threshold reported in [70] as acceptable improvement and half of the natural threshold 1%. (b) Comparisons for the results under MTL and STL. The metric for regression tasks, including Loewe and RI_row, is PCC. The metric for the classification task, including Classification, is ROCAUC. (c) Comparisons for the results under different training settings. Data are presented in boxplots (n=5 per group; center line, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers). Here *BAITSAO-FT* represents that we fine-tuned the pre-trained model, *BAITSAO-ZS* represents that we applied the pre-trained model for these tasks under zero-shot learning framework, and *BAITSAO-FS* represents that we did not use the pre-trained weights for these tasks. Here FT means fine-tuning, ZS means zero-shot learning and FS means from scratch. We included four metrics across three datasets for comparisons. (d) The first example of tri-drug cases for drug synergy prediction with BAITSAO. (e) The second example of tri-drug cases for drug synergy prediction with BAITSAO. Source data are provided as a Source Data file.

embeddings from different sources to construct the datasets for training and testing. Moreover, our approach strikes a good balance between efficiency and performance. According to Figure 6 (b), the running time of BAITSAO without pre-training was significantly lower than the classical methods DeepSynergy and SVM for drug synergistic effect prediction. Moreover, the pre-trained BAITSAO with the fine-tuning framework converges at a much faster rate, thus pre-trained BAITSAO achieved an even faster running speed by comparing with MARSY and DeepDDs. Therefore, our training framework strikes a good balance between runtime and model performance. Both pre-training and fine-tuning stages can be finished with only one GPU, presenting no hardware barrier to deploy BAITSAO.

We performed ablation tests for the MTL strategy, shown in Supplementary Figure 12 for ablation of methods and Supplementary Figure 13 for ablation of task-specific layers. We compared the gradient matching-based approaches including PCGrad [81], GradVac [82], CAGrad [83], Nash-MTL [84] and the linear MTL framework LinearMTL [85, 86] with our revised UW approach and found that our choice generally had comparable or better results, especially for the classification task. Moreover, LinearMTL performed much worse than deep learning based methods on the regression type tasks. Therefore, we chose the revised UW as the method for the pre-training stage. Moreover, our final choice with one task-specific layer for each task had the best overall performance, and increasing the number of layers required more computing resources, thus we chose our design shown in Figure 1 (c).

We also analyzed the relation between the size of the training dataset and model performance. We adjust the proportion we used for model training and visualize the relation between proportion and metrics in Figures 6 (c) for regression and (d) for classification. From these figures, a larger proportion tended to increase the model performance, with its limit for proportion ≥ 0.9 for these two tasks. Moreover, using only 0.1% training dataset to train a model for a classification task can still generate relatively high ROCAUC, thus the classification task may not be difficult for BAITSAO.

In Figures 6 (e) and (f), we examined the scaling law [87, 88] of BAITSAO. We adjusted the layer width of our model and plotted the relation between the layer width in the hidden layer and model performance for the regression task and the classification task. These figures show that we can model the relation between model parameters and model performance to predict performance, where more parameters lead to better performance. Therefore, the performance improvement of our model with scaling can be explained by the scaling law, and our model has good scalability. Our findings can help us understand the model training process in a better approach and determine the optimized source allocation of a fixed compute budget. For example, for machines that cannot support the version of BAITSAO with layer width as 10240, the version of BAITSAO with layer width as 4096 can also have acceptable performances and can be considered to deploy.

3 Discussion

Predicting drug synergistic effect is important for drug development and patient treatment. In the past, limited by available experimental data, information on drugs/cell lines, and pipelines to predict drug synergistic effect, there are few approaches to predicting drug synergistic effect for general use. With the help of large-scale drug synergy information databases, LLMs, and an MTL framework, we introduced BAITSAO as a unified model with a general pipeline for drug synergistic effect prediction as well as single-drug inhibition prediction. BAITSAO optimized the network architecture through comprehensive benchmarking analysis and was pre-trained based on the latest large-scale databases. It achieved top-tier performance in both regression tasks and classification tasks for drug synergistic effect prediction.

There are two major contributions of our work. Firstly, we presented a unified pipeline to construct datasets for synergistic effect analysis for both drugs and cell lines based on the embeddings from LLMs, thus we mitigated the difference caused by aliases for drugs and cell lines of different datasets. We demonstrated that the embeddings contained functional information for drugs and cell lines. We proposed a new design to construct training datasets, thus we only need to utilize the overlapped information across datasets for drug synergy analysis. Secondly, we pre-trained a unified model with a MTL framework for drug synergy analysis and single-drug inhibition analysis supported by rigorous task-selection steps. We demonstrated that BAITSAO benefited from the pre-training process and had good generalization ability with fine-tuning in fewer steps compared with the training process from scratch. Moreover, pre-trained BAITSAO showed its potential as a good zero-shot reasoner for drug synergy prediction under the classification settings. Therefore, we overcame the generalization issue in previous work based on transfer learning [21] and proposed a new avenue for the construction of BAITSAO for drug synergy analysis.

We conducted a sensitivity analysis to offer guidance for future model deployment. We showed that our current hyper-parameter settings and data construction methods are the optimal choices by hyper-parameter tuning and ablation tests. We also analyzed the relation between the proportions of data we used for training and model performance. While increasing training data proportions tended to improve prediction, BAITSAO performed well for the classification task for small data scales. Finally, we investigated the scaling law of BAITSAO and showed that the model performance is predictable and we could increase the model performance by scaling up BAITSAO for drug synergy prediction.

In conclusion, we have developed BAITSAO, an explainable model for drug synergy prediction, and demonstrated the superiority of BAITSAO over other methods by comprehensive benchmarking analysis and rigorous sensitivity analysis. We hope that BAITSAO can help researchers to better understand the process of drug synergistic effect prediction and further help in optimizing drug structures for drug design and discovering novel drug combinations with synergistic effects for clinical usage.

Furthermore, we also found that BAITSAO might not work well for drugs without a clear functional or chemical description in the early stage of drug development, which is a potential limitation of our application scenarios of all functional-based synergy predictors. In the future, we plan to incorporate more updated drug synergy

databases to keep this model updated, and we also plan to combine this model with information from genomics including single-cell data [89] and genome-wide association studies (GWAS) [90], especially for early-stage and novel drugs.

4 Methods

Problem definition. In this manuscript, we intend to construct a dataset $\mathcal{D} = (X, Y)$ and pre-train a model known as \mathcal{M} for the prediction of values in $Y^{n \times t}$, where n represents the number of combinations between drug pairs and the cell line, and t represents the number of tasks. Here $X^{n \times p}$ represents the feature space with n samples and p features. We then split the dataset \mathcal{D} into $\mathcal{D}_{train} = (X, Y)_{i=1}^{n_0}$ for training and $\mathcal{D}_{val} = (X, Y)_{i=1}^{n_1}$ for validation. Our target is to train a model \mathcal{M}^* based on \mathcal{D}_{train} and then select the optimal model based on \mathcal{D}_{val} . That is,

$$\theta^* = \operatorname{argmin}_{\theta} \mathcal{L}_m(\mathcal{M}(X_{val}, \theta), Y_{val}), \quad (1)$$

where $\mathcal{M}(\theta)$ represents the pre-trained model with parameter θ , and θ^* represents the optimal model parameters. \mathcal{L}_m represents multi-task learning loss. After obtaining the optimal model, we apply the model $\mathcal{M}^*(\theta^*)$ for a new dataset containing out-of-distribution (OOD) data, known as \mathcal{D}_{test} .

Construction of pre-training datasets and testing datasets. One major contribution of our work is to unify the features we need to predict the drug-related information for both the synergistic effect and inhibition effect. We at least need the names of drugs and cell lines. Considering we have a drug pair (d_1, d_2) and a cell line (c_1) , our idea is to generate the description of both drugs as $W(d_1), W(d_2)$ and the cell line as $W(c_1)$ based on LLMs such as GPT 3.5, and then utilize the embeddings tool of GPT 3.5 to transfer the text description into embeddings with e dimensions. Therefore, our final sample $x \in X$ is defined as:

$$x = \operatorname{AVG}(\operatorname{emb}(W(d_1)), \operatorname{emb}(W(d_2))) || \operatorname{emb}(W(c_1)) || \#d, \quad (2)$$

where $\operatorname{AVG}()$ represents the functions to compute the mean of the given variables, and $\operatorname{emb}()$ is the function to obtain the embeddings of the input. $\#d$ represents the number of drugs we used, which can be encoded as embeddings [91]. We take the unbiased estimation of the drug combination in the feature levels by computing the average value of embeddings, and we show that this approach works better than other types of feature integration in the sensitivity analysis section of the manuscript. Notably, this approach also scales for more drug combinations. Considering the case of k drugs with the cell line c_i , we define one sample $x \in X$ as:

$$x = \operatorname{AVG}(\operatorname{emb}(W(d_1)), \dots, \operatorname{emb}(W(d_k))) || \operatorname{emb}(W(c_i)) || k. \quad (3)$$

Therefore, for an arbitrary input dataset with feature space X containing drug information and cell-line information, we can transfer the samples in the given dataset from the text space to the numerical space, thus we unify the input data format for this task. Furthermore, to predict the drug synergistic effect, we consider both regression

and classification. In the case of regression, we intend to predict the specific synergy score of samples. To compute the synergistic effect based on IC₅₀ information, under different rules, we can have different scores. Here we consider four methods to model the synergy scores, known as HSA, Bliss, Loewe and ZIP. If we consider N drugs with multi-drug combination effect as $E_{A,B,\dots,N}$ and we intend to compute the synergy scores S_{HSA} , S_{Bliss} , S_{Loewe} , and S_{ZIP} , according to [5], we have:

$$S_{HSA} = E_{A,B,\dots,N} - \max(E_A, E_B, \dots, E_N). \quad (4)$$

$$S_{Bliss} = E_{A,B,\dots,N} - (E_A + E_B + \dots + E_N - E_A E_B \dots E_N). \quad (5)$$

$$-E_A E_N - E_B E_N - \dots - E_A E_B \dots E_N). \quad (6)$$

$$S_{Loewe} = \frac{a}{E_A} + \frac{b}{E_B} + \dots + \frac{n}{E_N}. \quad (7)$$

$$S_{ZIP} = E_{A,B,\dots,N} - \left(\frac{\left(\frac{x_A}{m_A}\right)^{\lambda_A}}{1 + \left(\frac{x_A}{m_A}\right)^{\lambda_A}} + \frac{\left(\frac{x_B}{m_B}\right)^{\lambda_B}}{1 + \left(\frac{x_B}{m_B}\right)^{\lambda_B}} + \dots \right) \quad (8)$$

$$+ \frac{\left(\frac{x_N}{m_N}\right)^{\lambda_N}}{1 + \left(\frac{x_N}{m_N}\right)^{\lambda_N}} - \frac{\left(\frac{x_A}{m_A}\right)^{\lambda_A}}{1 + \left(\frac{x_A}{m_A}\right)^{\lambda_A}} \frac{\left(\frac{x_B}{m_B}\right)^{\lambda_B}}{1 + \left(\frac{x_B}{m_B}\right)^{\lambda_B}} \quad (9)$$

$$- \frac{\left(\frac{x_A}{m_A}\right)^{\lambda_A}}{1 + \left(\frac{x_A}{m_A}\right)^{\lambda_A}} \frac{\left(\frac{x_N}{m_N}\right)^{\lambda_N}}{1 + \left(\frac{x_N}{m_N}\right)^{\lambda_N}} \quad (10)$$

$$- \frac{\left(\frac{x_B}{m_B}\right)^{\lambda_B}}{1 + \left(\frac{x_B}{m_B}\right)^{\lambda_B}} \frac{\left(\frac{x_N}{m_N}\right)^{\lambda_N}}{1 + \left(\frac{x_N}{m_N}\right)^{\lambda_N}} - \dots \quad (11)$$

$$- \frac{\left(\frac{x_A}{m_A}\right)^{\lambda_A}}{1 + \left(\frac{x_A}{m_A}\right)^{\lambda_A}} \frac{\left(\frac{x_B}{m_B}\right)^{\lambda_B}}{1 + \left(\frac{x_B}{m_B}\right)^{\lambda_B}} \dots \frac{\left(\frac{x_N}{m_N}\right)^{\lambda_N}}{1 + \left(\frac{x_N}{m_N}\right)^{\lambda_N}} \right). \quad (12)$$

Here we have E_A, E_B, \dots, E_N as measured responses of different drugs, and a, b, \dots, n represent the doses of the single drugs we need to produce the combination effect. Moreover, to compute S_{ZIP} , we have x_N as the dose of drug N fitted with the four-parameter log-logistic model, and m_N represents the dose we need to produce the half-maximum effect (IC₅₀). We also have λ_N as the shape parameter to indicate the slope of the dose-response curve. In the MTL case, we consider S_{Loewe} for the targets

of regression. We also pre-train models to predict the other three scores. All of the synergy scores are extracted from the database of DrugComb.

In order to characterize the inhibitory effects of individual drugs, we introduce the RI score in the prediction task. RI score is the normalized area under the \log_{10} -transformed dose-response curves. RI scores of all drugs are also extracted from the database of DrugComb.

In the case of classification, we intend to predict whether the given drug pair has a synergistic effect under a specific cell line, which is a binary classification problem. To construct the dataset for this task, we set the threshold of S_{Loewe} to binarize the synergistic effect of different drug combinations. Since not all of the testing datasets in the real world contain data for both regression and classification, thus introducing a classification task is meaningful.

Investigation of embeddings. We set up different methods to ensure that embeddings from LLMs contain the necessary information to describe the properties of drugs and cell lines. We consider two prompt engineering approaches for description generation, including MetaPrompt [47] and Chain-of-Thought (COT) [47]. MetaPrompt introduces a system prompt for LLMs and generates the outputs conditioned on the context. COT allows LLMs to obtain complex reasoning capabilities by forcing models to address the problem with intermediate steps. We also generate text descriptions and embeddings for drugs and cell lines from the dataset used by Deep-synergy (D1). We check the correctness of the all descriptions and the similarity of 10 sampled embeddings across different drugs to evaluate the correctness. Moreover, we change the random seed to generate different descriptions as well as embeddings to check the variance of drug embeddings from the same drug. We also record the description of cell lines in Supplementary file 1. Furthermore, we modify CPA to predict gene expression under different perturbations enhanced by drug embeddings from LLMs. In this step, we replace the original drug embeddings used in the CPA with our new embeddings. This approach allows us to check the correctness of embeddings from the application perspective. We use R^2 scores to evaluate the performances. To compute the R^2 score, we have the ground truth synergy score y and predicted synergy score \hat{y} and follow its definition:

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}, \quad (13)$$

where \bar{y} means that we compute the average value of the input variable y . R^2 represents the explanation of independent variables for the dependent variable. Higher R^2 means better model performance.

Therefore, our assessment of the quality of embeddings takes into account meanings, variance, and applications.

Hyper-parameter searching. We summarize the search space for hyper-parameters of each method in Table 1. The best hyper-parameter setting is determined by the performance of models based on the validation dataset.

Here lr means learning rate, Dropout means dropout rate (the ratio of neurons we intend to close during the training process), max_depth means the maximal depth for

Table 1 Hyper-parameter search space for each method ranked in alphabetical order.

Methods	Searching space
BAITSAO	lr:[1e-5,1e-4]; Dropout:[0.1,0.3]
BERT	epochs:[5,10]; lr: [1e-5,1e-4]
DeepDDs	Optimal hyper-parameters from [19]
DeepSynergy	Optimal hyper-parameters from [14]
Lasso	alpha:[1,10]
MARSY	Optimal hyper-parameters from [18]
SVM	C:[1,5]
TabNet	n_steps:[1,5]; n_a:[8,64]; n_d:[8,64]; gamma:[1.1,1.5]
TreeComb	max_depth:[10, 100]; n_estimators:[10,50]; min_child_weight:[1, 3]

tree-based models, n_estimators means the number of estimators for tree-based models, min_child_weight means the minimal weights of child nodes in tree-based models, C means the regularization weight for SVM, n_steps means the number of decision steps in the model architecture, n_a means the width of the attention embedding for each masked choice, n_d means the width of the prediction layer, gamma means the coefficient for the feature re-usage in the masking process, epochs mean the number of epochs we used to train the model, alpha represents the regularization coefficient for Lasso. We present the results under different hyper-parameters for BAITSAO in Supplementary Figures 14 (a) and (b). We find that lr plays a more important role in the training process while adjusting the dropout rate does not affect the model performance much.

Selection of model architecture. After setting up the pre-training dataset, we seek a suitable model architecture. Since deep neural networks (DNNs) related methods have shown impressive performance as a base model for large-scale models [92, 93], we construct the pre-training architecture of BAITSAO based on DeepSynergy. To assess the strength of our model architecture, we remove the pre-training step and compare BAITSAO with other methods for both the regression task and classification task with three different datasets. We also determine the hyper-parameters of model training in this stage. The superiority of BAITSAO is shown in the Results section, and we expect to see its similar performance at both the pre-training and fine-tuning stages.

In the model architecture selection stage, we utilize Adam [94] as the optimizer and ReduceLROnPlateau [91] as the learning rate scheduler. The starting learning rate for D1 and D3 is 1e-5, while it is 1e-4 for D2. The dropout rate is 0.2, and the patience for the scheduler is 10. Our patience for early-stopping step is 100, and the maximum number of epochs is 1000.

Explainability. The design of BAITSAO allows us to characterize the relevance between the specific gene and drug combinations across different cell lines. To perform this analysis, the input format of one combination becomes:

$$x' = AVG(emb(W(d_1)), ..., emb(W(d_k))) || exp(c_i) || emb(W(c_i)) || k, \quad (14)$$

where $exp(c_i)$ represents the gene expression profile for the cell line i . After the training process, we can extract the importance of different genes based on SHAP. For gene j , its importance for the synergistic effect of drug combinations (d_1, \dots, d_k) for cell line c_i can be calculated as:

$$I_j = ShapValue(\mathcal{M}, x'), \quad (15)$$

where I_j represents the importance and x' was defined above. $ShapValue()$ is a function to compute the importance with model \mathcal{M} and input x' . Here larger I_j represents more importance in the prediction process.

We select 1000 highly-variable genes for the analysis of explainability. This number is determined by adjusting the number of genes to achieve the best model performance. Our tuning results are shown in Supplementary Figure 15.

For the bulk RNA-seq datasets of different cell lines, we use DESeq2 to identify DEGs by comparing groups with and without predicted drug synergistic effects.

For the two scRNA-seq datasets used for validating our selected genes, we follow the pre-processing pipeline of Scanpy [95] and run the Wilcoxon rank-sum test to access the list of DEGs.

Pre-training BAITSAO under the multi-task learning framework. Here we explain our settings for the multi-task learning framework. After filtering tasks based on the constructed help-harm matrix, we consider three tasks: 1. Prediction of S_{Loewe} as a regression task. 2. Prediction of RI for one drug as a regression task. 3. Prediction of drug synergistic effect as a binary classification task. Therefore, we have two regression tasks and one classification task, and their loss functions are represented as $\mathcal{L}_1, \mathcal{L}_2$, and \mathcal{L}_3 . Traditionally, we construct the final loss function as a linear combination:

$$\mathcal{L}_m = \sum_{i=1}^{n_t} w_i \mathcal{L}_i, \quad (16)$$

where w_i represents the pre-defined weights for the loss function \mathcal{L}_i and $n_t = 3$. However, determining the values of the weights is difficult. Moreover, it is a strong assumption that the weights do not change during training is also a very strong assumption. Therefore, we introduce the uncertainty of loss function in this process and make the weights learnable. Typically, by choosing mean squared error (MSE) as the loss function in the training process for the regression task, we have the equivalent maximum likelihood framework of a Gaussian distribution for prediction output y and model \mathcal{M} . Therefore, the log-likelihood of the regression task can be represented as:

$$\log(p(y|\mathcal{M}(x, \theta))) \propto -\frac{1}{2\sigma^2} \|y - \mathcal{M}(x, \theta)\|^2 - \log\sigma, \quad (17)$$

where the uncertainty is defined as σ , and x represents model input. σ is a learnable parameter. Similarly, for a classification problem, we can represent the log-likelihood based on a Softmax function, that is:

$$\log(p(y|\mathcal{M}(x, \theta))) = \log(\text{Softmax}(\frac{1}{\sigma^2} \mathcal{M}(x, \theta))), \quad (18)$$

where $\text{Softmax}(x)_i = \frac{exp(x_i)}{\sum_j exp(x_j)}$, and p represents the length of x .

Therefore, the final loss function can be represented as maximizing the joint distribution of three tasks. In the validation stage, we minimize the maximal term in the loss function group rather than the original weighted loss function design from UW. We also add a constant term ϵ to ensure the numerical stability, so our final loss function is:

$$\mathcal{L}_m = -\log(p(y_1, y_2, y_3 | M(x, \theta))) \quad (19)$$

$$\approx \frac{1}{2\sigma_1^2 + \epsilon} \mathcal{L}_1(y_1, M(x, \theta)) + \frac{1}{2\sigma_2^2 + \epsilon} \mathcal{L}_2(y_2, M(x, \theta)) \quad (20)$$

$$+ \frac{1}{\sigma_3^2 + \epsilon} \mathcal{L}_3(y_3, M(x, \theta)) + \log(\sigma_1 \sigma_2 \sigma_3). \quad (21)$$

In the pre-training stage, we utilize Adam [94] as optimizer and ReduceLROnPlateau [91] as the learning rate scheduler. The starting learning rate is 1e-4, the dropout rate is 0.2, and the patience for the scheduler is 100. Our patience of early-stopping step is 500, and the maximum number of epochs is 1000. The number of combinations we used for pre-training is 739,652, including 4268 unique drugs and 288 unique cell lines.

After finishing the pre-training step, we test the model performance on the testing datasets under both the zero-shot learning case and the fine-tuning case with a parameter-freezing design. We also extend the prediction of the synergistic effect to the case of n ($n \geq 3$) drug combinations. Finally, we include a tutorial in our code repository for both the fine-tuning approach and the zero-shot inference approach.

We also pre-train other baselines, including DeepSynergy [14], DeepDDs [19], and MARSY [18], based on the same dataset. Details of model comparison are discussed in the Results section.

Zero-shot query and multi-drug prediction. Our model is capable of zero-shot synergy effect prediction. By transferring the knowledge and information of drugs and cell lines into embeddings through GPT 3.5 and the embedding layer, users can generate embeddings of arbitrary combinations as input for querying the synergy effects with a pre-trained BAITSAO.

For the combinations with three or more drugs, we directly generate the synergy score under the pre-trained model with the zero-shot learning framework. The three-drug case we used in the main text is from a known database, while it is possible to explore combinations with a larger number of drugs as long as the combinations are practical and meaningful. To access the determined predicted value, we do not use the dropout layers in the testing process. To access the predicted value with uncertainty, we keep the dropout layers in the testing process and repeat the prediction process for 100 times to access the estimation of mean and standard deviation for each combination. Such approach is known as MC Dropout.

We summarize the details of zero-shot query as a tutorial in our code repository.

Model evaluation. We consider four different metrics to evaluate the performance of different models for the drug synergistic effect prediction task, with two metrics for regression and two metrics for classification.

For the regression task, we consider two metrics: Pearson correlation coefficient (PCC) and Mean Squared Error (MSE).

1. PCC: Since we know the ground truth synergy score y and predicted synergy score \hat{y} , we can directly compute the PCC as:

$$PCC(y, \hat{y}) = \frac{COV(y, \hat{y})}{\sigma(y)\sigma(\hat{y})}, \quad (22)$$

where $COV()$ is the function to compute the covariance of two variables, and $\sigma()$ is the function to compute the standard deviation of the input variable. Higher PCC means better model performance.

2. MSE: To compute the mean squared error, we have the ground truth synergy score y and predicted synergy score \hat{y} and follow its definition:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad (23)$$

where i represents the index of samples, and lower MSE means better model performance.

For the classification task, we consider two metrics: Area under the ROC Curve (ROCAUC) and Accuracy (ACC).

1. ROCAUC: To compute this metric, we construct the relation between the true-positive rate and the false-positive rate under different probability thresholds. Such relation can be reflected in the ROC curve. We then compute the area under the ROC curve, and this area represents ROCAUC. Higher ROCAUC means better model performance.

2. ACC: To compute this metric, we have the ground truth synergistic effect condition $y^{n \times 1}$ and predicted binary value \hat{y} , we then compute the ACC as:

$$ACC = \frac{\sum_{i=1}^n \mathbb{1}_{y_i=\hat{y}_i}}{n}, \quad (24)$$

where $\mathbb{1}_{y_i=\hat{y}_i}$ is a indicator function and only takes 1 when $y_i = \hat{y}_i$. Higher ACC means better model performance.

We report the mean and standard deviation of these metrics by using five-fold cross-validation for each dataset.

Overview of other methods. In this section, we summarize the benchmarking methods used in our work. These methods (ranked in alphabetical order) include:

- BERT: BERT is a pre-trained bi-directional transformer for language understanding. For this model, we construct the training datasets and testing datasets directly from drug descriptions and cell-line descriptions. The problem is then formalized as Question-answering case for both classification and regression tasks.
- DeepDDs [19]: DeepDDs is a Graph Neural Network (GNN)-based method for drug synergistic effect prediction. This method can only handle the classification task.

- 726 The training dataset of DeepDDs is constructed base on features of drugs as graphs
 727 from chemical information and gene expression levels from cell lines.
- 728 • DeepSynergy [14]: DeepSynergy is a DNN-based method for drug synergistic effect
 729 prediction. This method can handle both the regression task and the classification
 730 task, by changing the loss function and the activation function of the last network
 731 layer. The training dataset of DeepSynergy follows its default mode, including fea-
 732 tures of drugs from chemical information and cell-line features from gene expression
 733 levels.
 - 734 • Lasso [39, 52]: Lasso is a regularized regression method for drug synergistic effect
 735 prediction. This method can handle both the regression task and the classification
 736 task, by using the default mode or logistic regression mode with L1 penalty. The
 737 training dataset of Lasso is constructed based on the drug embeddings and cell-line
 738 embeddings from LLMs.
 - 739 • MARSY [18]: MARSY is a DNN-based method with multi-task learning framework
 740 for drug synergistic effect prediction. This method can only handle the regression
 741 task. The training dataset of MARSY is constructed based on features of drugs from
 742 chemical information, gene expression levels from cell lines and tissue information.
 - 743 • SVM [39, 49]: SVM is a machine learning method based on constructing decision-
 744 making boundaries for drug synergistic effect prediction. This method can handle
 745 both the regression task and the classification task, by using SVR or SVC. The
 746 training dataset of SVM is constructed based on the drug embeddings and cell-line
 747 embeddings from LLMs.
 - 748 • TabNet [50]: TabNet is a DNN-based method with transformer architecture for drug
 749 synergistic effect prediction. TabNet combines the ideas from both neural network
 750 design and tree-model design. This method can handle both the regression task and
 751 the classification task, by changing the loss function and the activation function of
 752 the last network layer. The training dataset of TabNet is constructed based on the
 753 drug embeddings and cell-line embeddings from LLMs.
 - 754 • TreeComb [15]: TreeComb is an explainable machine learning method based on
 755 XGBoost for drug synergistic effect prediction. This method can handle both the
 756 regression task and the classification task, by using XGBREGRESSOR or XGB-
 757 CLASSIFIER. The training dataset of TreeComb is constructed based on the drug
 758 embeddings and cell-line embeddings from LLMs.

759 **Datasets preparation.** We utilize public datasets from DrugComb v1.5 for pre-
 760 training. For the regression task, we have one dataset from DeepSynergy (as D1, which
 761 is processed in the original paper) using Loewe as the synergy score computation
 762 method. We also have one dataset from MARSY (as D2, which is processed in the
 763 original paper) using ZIP as the synergy score computation method. For the classifi-
 764 cation task, we have one dataset from DeepSynergy (as D1) using the Loewe synergy
 765 score with a threshold. We also have one dataset from DeepDDs with a known binary
 766 synergistic effect condition (as D3 [96]). For multi-drug synergistic effect inference, we
 767 utilize one dataset from DrugCombDB. Every dataset at least contains the names of
 768 drugs and cell lines.

769 5 Data availability

770 We do not generate new data in this research and all data used in this manuscript are
771 publicly available. The DrugComb data used in this study are available under acces-
772 sion code [DrugCombDownload](#). The training and testing data used in this study are
773 available under accession code [D1](#), [D2](#), and [D3](#). The scRNA-seq data used in this study
774 are available under accession code [GSE215121](#) and [SCP109](#). We collect the informa-
775 tion of downloading training datasets as well as their statistics in Supplementary file
776 4. Source data are provided with this paper.

777 6 Reproductivity and codes availability

778 We used the resources from the Yale High Performance Center (Yale HPC) and
779 UCLA Computing Servers to conduct all of the experiments. Our maximum run-
780 ning time for each dataset was 24 hours and maximum RAM was 100 GB. The
781 version of GPU we used is NVIDIA A5000 (24 GB) for fine-tuning and single task
782 learning, and NVIDIA A100 (40GB) for pre-training. The codes of BAITSAO can
783 be found in <https://github.com/HelloWorldLTY/BAITSAO> and <https://doi.org/10.5281/zenodo.15105815> with MIT license. The pre-trained weights can be found in
784 <https://huggingface.co/iLOVE2D/BAITSAO>. The version of softwares used for data
785 collection and model training is summarized in Supplementary file 4.
786

787 References

- 788 [1] Clercq, E.D.: The design of drugs for hiv and hcv. *Nature reviews Drug discovery*
789 **6**(12), 1001–1018 (2007)
- 790 [2] Mokhtari, R.B., Homayouni, T.S., Baluch, N., Morgatskaya, E., Kumar, S., Das,
791 B., Yeger, H.: Combination therapy in combating cancer. *Oncotarget* **8**(23),
792 38022 (2017)
- 793 [3] Al-Lazikani, B., Banerji, U., Workman, P.: Combinatorial drug therapy for
794 cancer in the post-genomic era. *Nature biotechnology* **30**(7), 679–692 (2012)
- 795 [4] Holbeck, S.L., Camalier, R., Crowell, J.A., Govindharajulu, J.P., Hollingshead,
796 M., Anderson, L.W., Polley, E., Rubinstein, L., Srivastava, A., Wilsker, D., *et*
797 *al.*: The national cancer institute almanac: a comprehensive screening resource
798 for the detection of anticancer drug pairs with enhanced therapeutic activity.
799 *Cancer research* **77**(13), 3564–3576 (2017)
- 800 [5] Ianevski, A., Giri, A.K., Aittokallio, T.: Synergyfinder 2.0: visual analytics
801 of multi-drug combination synergies. *Nucleic acids research* **48**(W1), 488–493
802 (2020)
- 803 [6] Law, M., Wald, N., Morris, J., Jordan, R.: Value of low dose combination treat-
804 ment with blood pressure lowering drugs: analysis of 354 randomised trials. *Bmj*
805 **326**(7404), 1427 (2003)

- 806 [7] Wood, K.B., Wood, K.C., Nishida, S., Chuzel, P.: Uncovering scaling laws to
807 infer multidrug response of resistant microbes and cancer cells. *Cell reports* **6**(6),
808 1073–1084 (2014)
- 809 [8] Hetzel, L., Boehm, S., Kilbertus, N., Günnemann, S., Theis, F., *et al.*: Pre-
810 dicting cellular responses to novel drug perturbations at a single-cell resolution.
811 *Advances in Neural Information Processing Systems* **35**, 26711–26722 (2022)
- 812 [9] Lotfollahi, M., Klimovskaia Susmelj, A., De Donno, C., Hetzel, L., Ji, Y., Ibarra,
813 I.L., Srivatsan, S.R., Naghipourfar, M., Daza, R.M., Martin, B., *et al.*: Pre-
814 dicting cellular responses to complex perturbations in high-throughput screens.
815 *Molecular Systems Biology*, 11517 (2023)
- 816 [10] Wilding, J.L., Bodmer, W.F.: Cancer cell lines for drug discovery and develop-
817 ment. *Cancer research* **74**(9), 2377–2384 (2014)
- 818 [11] Ianevski, A., Giri, A.K., Aittokallio, T.: Synergyfinder 3.0: an interactive analysis
819 and consensus interpretation of multi-drug synergies across multiple samples.
820 *Nucleic Acids Research* **50**(W1), 739–743 (2022)
- 821 [12] Roemer, T., Boone, C.: Systems-level antimicrobial drug and drug synergy
822 discovery. *Nature chemical biology* **9**(4), 222–231 (2013)
- 823 [13] Sun, W., Sanderson, P.E., Zheng, W.: Drug combination therapy increases
824 successful drug repositioning. *Drug discovery today* **21**(7), 1189–1195 (2016)
- 825 [14] Preuer, K., Lewis, R.P., Hochreiter, S., Bender, A., Bulusu, K.C., Klam-
826 bauer, G.: Deepsynergy: predicting anti-cancer drug synergy with deep learning.
827 *Bioinformatics* **34**(9), 1538–1546 (2018)
- 828 [15] Janizek, J.D., Celik, S., Lee, S.-I.: Explainable machine learning prediction of
829 synergistic drug combinations for precision cancer medicine. *BioRxiv*, 331769
830 (2018)
- 831 [16] Janizek, J.D., Dincer, A.B., Celik, S., Chen, H., Chen, W., Naxerova, K., Lee, S.-
832 I.: Uncovering expression signatures of synergistic drug responses via ensembles
833 of explainable machine-learning models. *Nature Biomedical Engineering* **7**(6),
834 811–829 (2023)
- 835 [17] Kuru, H.I., Tastan, O., Cicek, A.E.: Matchmaker: a deep learning framework
836 for drug synergy prediction. *IEEE/ACM transactions on computational biology*
837 *and bioinformatics* **19**(4), 2334–2344 (2021)
- 838 [18] El Khili, M.R., Memon, S.A., Emad, A.: Marsy: a multitask deep-learning frame-
839 work for prediction of drug combination synergy scores. *Bioinformatics* **39**(4),
840 177 (2023)

- 841 [19] Wang, J., Liu, X., Shen, S., Deng, L., Liu, H.: Deepdds: deep graph neural
842 network with attention mechanism to predict synergistic drug combinations.
843 Briefings in Bioinformatics **23**(1), 390 (2022)
- 844 [20] Zheng, S., Aldahdooh, J., Shadbahr, T., Wang, Y., Aldahdooh, D., Bao, J.,
845 Wang, W., Tang, J.: Drugcomb update: a more comprehensive drug sensitivity
846 data repository and analysis portal. Nucleic acids research **49**(W1), 174–184
847 (2021)
- 848 [21] Kim, Y., Zheng, S., Tang, J., Jim Zheng, W., Li, Z., Jiang, X.: Anticancer drug
849 synergy prediction in understudied tissues using transfer learning. Journal of the
850 American Medical Informatics Association **28**(1), 42–51 (2021)
- 851 [22] Bommasani, R., Hudson, D.A., Adeli, E., Altman, R., Arora, S., Arx, S., Bern-
852 stein, M.S., Bohg, J., Bosselut, A., Brunskill, E., et al.: On the opportunities
853 and risks of foundation models. arXiv preprint arXiv:2108.07258 (2021)
- 854 [23] Min, B., Ross, H., Sulem, E., Veyseh, A.P.B., Nguyen, T.H., Sainz, O., Agirre,
855 E., Heintz, I., Roth, D.: Recent advances in natural language processing via large
856 pre-trained language models: A survey. ACM Computing Surveys **56**(2), 1–40
857 (2023)
- 858 [24] Zhao, W.X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B.,
859 Zhang, J., Dong, Z., et al.: A survey of large language models. arXiv preprint
860 arXiv:2303.18223 (2023)
- 861 [25] Edwards, C.N., Naik, A., Khot, T., Burke, M.D., Ji, H., Hope, T.: Synergpt:
862 In-context learning for personalized drug synergy prediction and drug design.
863 bioRxiv, 2023–07 (2023)
- 864 [26] Li, T., Shetty, S., Kamath, A., Jaiswal, A., Jiang, X., Ding, Y., Kim, Y.: Can-
865 cergpt for few shot drug pair synergy prediction using large pretrained language
866 models. npj Digital Medicine **7**(1), 40 (2024)
- 867 [27] Luecken, M.D., Büttner, M., Chaichoompu, K., Danese, A., Interlandi, M.,
868 Müller, M.F., Strobl, D.C., Zappia, L., Dugas, M., Colomé-Tatché, M., et
869 al.: Benchmarking atlas-level data integration in single-cell genomics. Nature
870 methods **19**(1), 41–50 (2022)
- 871 [28] Zhang, Y., Yang, Q.: A survey on multi-task learning. IEEE Transactions on
872 Knowledge and Data Engineering **34**(12), 5586–5609 (2021)
- 873 [29] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P.,
874 Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models
875 are few-shot learners. Advances in neural information processing systems **33**,
876 1877–1901 (2020)

- 877 [30] Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L.,
878 Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical
879 report. arXiv preprint arXiv:2303.08774 (2023)
- 880 [31] Liu, T., Chen, T., Zheng, W., Luo, X., Zhao, H.: scelmo: Embeddings from
881 language models are good learners for single-cell data analysis. bioRxiv, 2023–12
882 (2023)
- 883 [32] Bai, G., Chai, Z., Ling, C., Wang, S., Lu, J., Zhang, N., Shi, T., Yu, Z., Zhu,
884 M., Zhang, Y., et al.: Beyond efficiency: A systematic survey of resource-efficient
885 large language models. CoRR (2024)
- 886 [33] Yang, Z., Jin, Y., Xu, X.: Hades: Hardware accelerated decoding for efficient
887 speculation in large language models. arXiv preprint arXiv:2412.19925 (2024)
- 888 [34] Anthropic: The Claude 3 Model Family: Opus, Sonnet, Haiku
- 889 [35] Team, G., Anil, R., Borgeaud, S., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk,
890 J., Dai, A.M., Hauth, A., Millican, K., et al.: Gemini: a family of highly capable
891 multimodal models. arXiv preprint arXiv:2312.11805 (2023)
- 892 [36] McInnes, L., Healy, J., Melville, J.: Umap: Uniform manifold approximation and
893 projection for dimension reduction. arXiv preprint arXiv:1802.03426 (2018)
- 894 [37] Knox, C., Wilson, M., Klinger, C., Franklin, M., Oler, E., Wilson, A., Pon, A.,
895 Cox, J., Chin, N.E., Strawbridge, S., Garcia-Patino, M., Kruger, R., Sivaku-
896 maran, A., Sanford, S., Doshi, R., Khetarpal, N., Fatokun, O., Doucet, D.,
897 Zubkowski, A., Rayat, D., Jackson, H., Harford, K., Anjum, A., Zakir, M., Wang,
898 F., Tian, S., Lee, B., Liigand, J., Peters, H., Wang, R.Q., Nguyen, T., So, D.,
899 Sharp, M., da Silva, R., Gabriel, C., Scantlebury, J., Jasinski, M., Ackerman, D.,
900 Jewison, T., Sajed, T., Gautam, V., Wishart, D.: Drugbank 6.0: the drugbank
901 knowledgebase for 2024. Nucleic Acids Research **52**(D1), 1265–1275 (2023)
- 902 [38] Schoch, C.L., Ciufu, S., Domrachev, M., Hutton, C.L., Kannan, S., Khovanskaya,
903 R., Leipe, D., Mcveigh, R., O'Neill, K., Robbertse, B., et al.: Ncbi taxonomy:
904 a comprehensive update on curation, resources and tools. Database **2020**, 062
905 (2020)
- 906 [39] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel,
907 O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al.: Scikit-learn:
908 Machine learning in python. the Journal of machine Learning research **12**, 2825–
909 2830 (2011)
- 910 [40] Luecken, M.D., Theis, F.J.: Current best practices in single-cell rna-seq analysis:
911 a tutorial. Molecular systems biology **15**(6), 8746 (2019)
- 912 [41] Landrum, G.: Rdkit documentation. Release **1**(1-79), 4 (2013)

- 913 [42] Landrum, G.: RDKit: Open-source Cheminformatics. <http://www.rdkit.org>
- 914 [43] Weininger, D.: Smiles, a chemical language and information system. 1. introduc-
915 tion to methodology and encoding rules. Journal of chemical information and
916 computer sciences **28**(1), 31–36 (1988)
- 917 [44] Burkhardt, D., Benz, A., Lieberman, R., Gigante, S., Chow, A., Holbrook, R.,
918 Cannoodt, R., Luecken, M.: Open problems – single-cell perturbations. Kaggle,
919 (2023)
- 920 [45] Malyutina, A., Majumder, M.M., Wang, W., Pessia, A., Heckman, C.A., Tang,
921 J.: Drug combination sensitivity scoring facilitates the discovery of synergistic
922 and efficacious drug combinations in cancer. PLoS computational biology **15**(5),
923 1006752 (2019)
- 924 [46] Baptista, D., Ferreira, P.G., Rocha, M.: A systematic evaluation of deep learning
925 methods for the prediction of drug synergy in cancer. PLoS Computational
926 Biology **19**(3), 1010200 (2023)
- 927 [47] Suzgun, M., Kalai, A.T.: Meta-prompting: Enhancing language models with
928 task-agnostic scaffolding. arXiv preprint arXiv:2401.12954 (2024)
- 929 [48] Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q.V., Zhou,
930 D., *et al.*: Chain-of-thought prompting elicits reasoning in large language models.
931 Advances in neural information processing systems **35**, 24824–24837 (2022)
- 932 [49] Cortes, C., Vapnik, V.: Support-vector networks. Machine learning **20**, 273–297
933 (1995)
- 934 [50] Arik, S.Ö., Pfister, T.: Tabnet: Attentive interpretable tabular learning. In: Pro-
935 ceedings of the AAAI Conference on Artificial Intelligence, vol. 35, pp. 6679–6687
936 (2021)
- 937 [51] Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: BERT: Pre-training of deep
938 bidirectional transformers for language understanding. In: Burstein, J., Doran,
939 C., Solorio, T. (eds.) Proceedings of the 2019 Conference of the North American
940 Chapter of the Association for Computational Linguistics: Human Language
941 Technologies, Volume 1 (Long and Short Papers), pp. 4171–4186. Association
942 for Computational Linguistics, Minneapolis, Minnesota (2019). [https://doi.org/](https://doi.org/10.18653/v1/N19-1423)
943 [10.18653/v1/N19-1423](https://doi.org/10.18653/v1/N19-1423) . <https://aclanthology.org/N19-1423>
- 944 [52] Tibshirani, R.: Regression shrinkage and selection via the lasso. Journal of
945 the Royal Statistical Society Series B: Statistical Methodology **58**(1), 267–288
946 (1996)
- 947 [53] Lundberg, S.M., Lee, S.-I.: A unified approach to interpreting model predictions.
948 Advances in neural information processing systems **30**, 4768–4777 (2017)

- [54] Barretina, J., Caponigro, G., Stransky, N., Venkatesan, K., Margolin, A.A., Kim, S., Wilson, C.J., Lehár, J., Kryukov, G.V., Sonkin, D., *et al.*: The cancer cell line encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature* **483**(7391), 603–607 (2012)
- [55] Blatkiewicz, M., Białas, P., Taryma-Leśniak, O., Hukowska-Szematowicz, B.: Pan-cancer analysis of vim expression in human cancer tissues. Preprint available at Research Square (2021)
- [56] Muzellec, B., Teleńczuk, M., Cabeli, V., Andreux, M.: Pydeseq2: a python package for bulk rna-seq differential expression analysis. *Bioinformatics* **39**(9), 547 (2023)
- [57] Love, M.I., Huber, W., Anders, S.: Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome biology* **15**(12), 1–21 (2014)
- [58] Lee, S.Y., Lee, J.H., Bae, Y.C., Suh, K.T., Jung, J.S., *et al.*: Bmp2 increases adipogenic differentiation in the presence of dexamethasone, which is inhibited by the treatment of tn α in human adipose tissue-derived stromal cells. *Cellular Physiology and Biochemistry* **34**(4), 1339–1350 (2014)
- [59] Parveen, S., Ashfaq, H., Shahid, M., Kanwal, A., Tayyeb, A.: Emerging therapeutic role of cdk inhibitors in targeting cancer stem cells. *Journal ISSN* **2766**, 2276 (2021)
- [60] Thomas, P.D., Ebert, D., Muruganujan, A., Mushayahama, T., Albou, L.-P., Mi, H.: Panther: Making genome-scale phylogenetics accessible to all. *Protein Science* **31**(1), 8–22 (2022)
- [61] Liberzon, A., Birger, C., Thorvaldsdóttir, H., Ghandi, M., Mesirov, J.P., Tamayo, P.: The molecular signatures database hallmark gene set collection. *Cell systems* **1**(6), 417–425 (2015)
- [62] Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., *et al.*: Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods* **17**(3), 261–272 (2020)
- [63] Traag, V.A., Waltman, L., Van Eck, N.J.: From louvain to leiden: guaranteeing well-connected communities. *Scientific reports* **9**(1), 5233 (2019)
- [64] Loewe, S.: The problem of synergism and antagonism of combined drugs. *Arzneimittel-forschung* **3**(6), 285–290 (1953)
- [65] Yadav, B., Wennerberg, K., Aittokallio, T., Tang, J.: Searching for drug synergy in complex dose-response landscapes using an interaction potency model. *Computational and structural biotechnology journal* **13**, 504–513 (2015)

- 985 [66] Me, B.: What is synergy. *Pharmacol Rev* **41**, 93–141 (1989)
- 986 [67] Bliss, C.I.: The toxicity of poisons applied jointly 1. *Annals of applied biology*
987 **26**(3), 585–615 (1939)
- 988 [68] Kendall, A., Gal, Y., Cipolla, R.: Multi-task learning using uncertainty to weigh
989 losses for scene geometry and semantics. In: *Proceedings of the IEEE Conference*
990 *on Computer Vision and Pattern Recognition*, pp. 7482–7491 (2018)
- 991 [69] Lin, B., Zhang, Y.: Libmtl: A python library for deep multi-task learning. *Journal*
992 *of Machine Learning Research* **24**(1-7), 18 (2023)
- 993 [70] Zhao, Y., He, L.: Deep learning in the eeg diagnosis of alzheimer’s disease. In:
994 *Computer Vision-ACCV 2014 Workshops: Singapore, Singapore, November 1-2,*
995 *2014, Revised Selected Papers, Part I 12*, pp. 340–353 (2015). Springer
- 996 [71] Liu, H., Zhang, W., Zou, B., Wang, J., Deng, Y., Deng, L.: Drugcombdb: a com-
997 prehensive database of drug combinations toward the discovery of combinatorial
998 therapy. *Nucleic acids research* **48**(D1), 871–881 (2020)
- 999 [72] Gal, Y., Ghahramani, Z.: Dropout as a bayesian approximation: Representing
1000 model uncertainty in deep learning. In: Balcan, M.F., Weinberger, K.Q. (eds.)
1001 *Proceedings of The 33rd International Conference on Machine Learning. Pro-*
1002 *ceedings of Machine Learning Research*, vol. 48, pp. 1050–1059. PMLR, New
1003 York, New York, USA (2016). <https://proceedings.mlr.press/v48/gal16.html>
- 1004 [73] Lemay, A., Hoebel, K., Bridge, C.P., Befano, B., De Sanjosé, S., Egemen, D.,
1005 Rodriguez, A.C., Schiffman, M., Campbell, J.P., Kalpathy-Cramer, J.: Improv-
1006 ing the repeatability of deep learning models with monte carlo dropout. *npj*
1007 *Digital Medicine* **5**(1), 174 (2022)
- 1008 [74] Kirchoff, M., Mucsányi, B., Oh, S.J., Kasneci, D.E.: Url: A representation
1009 learning benchmark for transferable uncertainty estimates. *Advances in Neural*
1010 *Information Processing Systems* **36** (2024)
- 1011 [75] Sigma-Aldrich’s I-BET. [https://www.emdmillipore.com/US/en/product/](https://www.emdmillipore.com/US/en/product/I-BET-CAS-1260907-17-2-Calbiochem,EMD_BIO-401010)
1012 [I-BET-CAS-1260907-17-2-Calbiochem,EMD_BIO-401010](https://www.emdmillipore.com/US/en/product/I-BET-CAS-1260907-17-2-Calbiochem,EMD_BIO-401010). Accessed:
1013 2024-01-17
- 1014 [76] Selleck’s I-BET151 (GSK1210151A). [https://www.selleckchem.com/products/](https://www.selleckchem.com/products/i-bet151-gsk1210151a.html)
1015 [i-bet151-gsk1210151a.html](https://www.selleckchem.com/products/i-bet151-gsk1210151a.html). Accessed: 2024-01-17
- 1016 [77] Selleck’s PF-562271. <https://www.selleckchem.com/products/pf-562271.html>.
1017 Accessed: 2024-01-29
- 1018 [78] Selleck’s Saracatinib. <https://www.selleckchem.com/products/AZD0530.html>.
1019 Accessed: 2024-01-29

- 1020 [79] Saatci, O., Kaymak, A., Raza, U., Ersan, P.G., Akbulut, O., Banister, C.E.,
1021 Sikirzhyski, V., Tokat, U.M., Aykut, G., Ansari, S.A., *et al.*: Targeting lysyl
1022 oxidase (lox) overcomes chemotherapy resistance in triple negative breast cancer.
1023 Nature communications **11**(1), 2416 (2020)
- 1024 [80] New embedding models and API updates. [https://openai.com/blog/
1025 new-embedding-models-and-api-updates](https://openai.com/blog/new-embedding-models-and-api-updates). Accessed: 2024-01-27
- 1026 [81] Yu, T., Kumar, S., Gupta, A., Levine, S., Hausman, K., Finn, C.: Gradient
1027 surgery for multi-task learning. In: Larochelle, H., Ranzato, M., Hadsell, R.,
1028 Balcan, M.F., Lin, H. (eds.) Advances in Neural Information Processing Systems,
1029 vol. 33, pp. 5824–5836 (2020). [https://proceedings.neurips.cc/paper_files/paper/
1030 2020/file/3fe78a8acf5fda99de95303940a2420c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/3fe78a8acf5fda99de95303940a2420c-Paper.pdf)
- 1031 [82] Wang, Z., Tsvetkov, Y., Firat, O., Cao, Y.: Gradient vaccine: Investigating and
1032 improving multi-task optimization in massively multilingual models. In: Interna-
1033 tional Conference on Learning Representations (2021). [https://openreview.net/
1034 forum?id=F1vEjWK-lH_](https://openreview.net/forum?id=F1vEjWK-lH_)
- 1035 [83] Liu, B., Liu, X., Jin, X., Stone, P., Liu, Q.: Conflict-averse gradient descent
1036 for multi-task learning. Advances in Neural Information Processing Systems **34**,
1037 18878–18890 (2021)
- 1038 [84] Navon, A., Shamsian, A., Achituve, I., Maron, H., Kawaguchi, K., Chechik,
1039 G., Fetaya, E.: Multi-task learning as a bargaining game. In: International
1040 Conference on Machine Learning, pp. 16428–16446 (2022). PMLR
- 1041 [85] LinearMTL. GitHub (2018)
- 1042 [86] Kim, S., Xing, E.P.: Tree-guided group lasso for multi-response regression with
1043 structured sparsity, with an application to eqtl mapping. The Annals of Applied
1044 Statistics **6**(3), 1095 (2012)
- 1045 [87] Frantar, E., Ruiz, C.R., Houlsby, N., Alistarh, D., Evci, U.: Scaling laws for
1046 sparsely-connected foundation models. In: The Twelfth International Confer-
1047 ence on Learning Representations (2024). [https://openreview.net/forum?id=
1048 i9K2ZWkYIP](https://openreview.net/forum?id=i9K2ZWkYIP)
- 1049 [88] Kaplan, J., McCandlish, S., Henighan, T., Brown, T.B., Chess, B., Child, R.,
1050 Gray, S., Radford, A., Wu, J., Amodei, D.: Scaling laws for neural language
1051 models. arXiv preprint arXiv:2001.08361 (2020)
- 1052 [89] Heumos, L., Schaar, A.C., Lance, C., Litinetskaya, A., Drost, F., Zappia, L.,
1053 Lücken, M.D., Strobl, D.C., Henao, J., Curion, F., et al.: Best practices for
1054 single-cell analysis across modalities. Nature Reviews Genetics, 1–23 (2023)
- 1055 [90] Sun, N., Zhao, H.: Statistical methods in genome-wide association studies.

- 1056 Annual Review of Biomedical Data Science **3**, 265–288 (2020)
- 1057 [91] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen,
1058 T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-
1059 performance deep learning library. Advances in neural information processing
1060 systems **32** (2019)
- 1061 [92] Awais, M., Naseer, M., Khan, S., Anwer, R.M., Cholakkal, H., Shah, M., Yang,
1062 M.-H., Khan, F.S.: Foundational models defining a new era in vision: A survey
1063 and outlook. arXiv preprint arXiv:2307.13721 (2023)
- 1064 [93] Campos Zabala, F.J.: Neural networks, deep learning, foundational models.
1065 In: Grow Your Business with AI: A First Principles Approach for Scal-
1066 ing Artificial Intelligence in the Enterprise, pp. 245–275. Apress, Berkeley,
1067 CA (2023). https://doi.org/10.1007/978-1-4842-9669-1_10 . https://doi.org/10.1007/978-1-4842-9669-1_10
- 1069 [94] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: Bengio,
1070 Y., LeCun, Y. (eds.) 3rd International Conference on Learning Representations,
1071 ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings
1072 (2015). <http://arxiv.org/abs/1412.6980>
- 1073 [95] Wolf, F.A., Angerer, P., Theis, F.J.: Scanpy: large-scale single-cell gene expres-
1074 sion data analysis. Genome biology **19**, 1–5 (2018)
- 1075 [96] O’Neil, J., Benita, Y., Feldman, I., Chenard, M., Roberts, B., Liu, Y., Li, J.,
1076 Kral, A., Lejnine, S., Loboda, A., *et al.*: An unbiased oncology compound screen
1077 to identify novel combination strategies. Molecular cancer therapeutics **15**(6),
1078 1155–1162 (2016)
- 1079 [97] Sounni, N.E., Baramova, E.N., Munaut, C., Maquoi, E., Frankenhe, F., Foidart,
1080 J.-M., Noël, A.: Expression of membrane type 1 matrix metalloproteinase
1081 (mt1-mmp) in a2058 melanoma cells is associated with mmp-2 activation and
1082 increased tumor growth and vascularization. International journal of cancer
1083 **98**(1), 23–28 (2002)
- 1084 [98] Vizuet, A.F.K., Hansen, F., Negri, E., Leite, M.C., Oliveira, D.L., Gonçalves,
1085 C.-A.: Effects of dexamethasone on the li-pilocarpine model of epilepsy: pro-
1086 tection against hippocampal inflammation and astrogliosis. Journal of Neuroin-
1087 flammation **15**, 1–14 (2018)
- 1088 [99] Kohlmeyer, J.L., Kaemmer, C.A., Pulliam, C., Maharjan, C.K., Samayoa, A.M.,
1089 Major, H.J., Cornick, K.E., Knepper-Adrian, V., Khanna, R., Sieren, J.C., *et*
1090 *al.*: Rabl6a is an essential driver of mpnsts that negatively regulates the rb1
1091 pathway and sensitizes tumor cells to cdk4/6 inhibitors. Clinical cancer research
1092 **26**(12), 2997–3011 (2020)

- 1093 [100] Moon, H., Donahue, L.R., Choi, E., Scumpia, P.O., Lowry, W.E., Grenier, J.K.,
1094 Zhu, J., White, A.C.: Melanocyte stem cell activation and translocation initiate
1095 cutaneous melanoma in response to uv exposure. *Cell stem cell* **21**(5), 665–678
1096 (2017)
- 1097 [101] Bono, A., La Monica, G., Alamia, F., Mingoia, F., Gentile, C., Peri, D., Lauria,
1098 A., Martorana, A.: In silico mixed ligand/structure-based design of new cdk-
1099 1/parp-1 dual inhibitors as anti-breast cancer agents. *International Journal of*
1100 *Molecular Sciences* **24**(18), 13769 (2023)
- 1101 [102] Ramon, J., Engelen, Y., De Keersmaecker, H., Goemaere, I., Punj, D., Morales,
1102 J.M., Bonte, C., Berx, G., Hoste, E., Stremersch, S., *et al.*: Laser-induced
1103 vapor nanobubbles for b16-f10 melanoma cell killing and intracellular delivery
1104 of chemotherapeutics. *Journal of Controlled Release* **365**, 1019–1036 (2024)
- 1105 [103] Tirosh, I., Izar, B., Prakadan, S.M., Wadsworth, M.H., Treacy, D., Trombetta,
1106 J.J., Rotem, A., Rodman, C., Lian, C., Murphy, G., *et al.*: Dissecting the
1107 multicellular ecosystem of metastatic melanoma by single-cell rna-seq. *Science*
1108 **352**(6282), 189–196 (2016)
- 1109 [104] Jerby-Arnon, L., Shah, P., Cuoco, M.S., Rodman, C., Su, M.-J., Melms, J.C.,
1110 Leeson, R., Kanodia, A., Mei, S., Lin, J.-R., *et al.*: A cancer cell program
1111 promotes t cell exclusion and resistance to checkpoint blockade. *Cell* **175**(4),
1112 984–997 (2018)
- 1113 [105] Zhang, C., Shen, H., Yang, T., Li, T., Liu, X., Wang, J., Liao, Z., Wei, J., Lu,
1114 J., Liu, H., *et al.*: A single-cell analysis reveals tumor heterogeneity and immune
1115 environment of acral melanoma. *Nature Communications* **13**(1), 7250 (2022)

1116 7 Acknowledgements

1117 We thank Yijia Xiao for the fruitful discussion about training the model. We thank
1118 Wangjie Zheng, and Tianqi Chen for helpful discussion about designing drug embed-
1119 dings. We thank Dr. Ning Sun, Xinyi Chen, Dr. Gefei Wang, Dr. Yingxin Lin and
1120 Chenyu Wang for the suggestions about approaches to improve the experimental
1121 design. We also thank Jiawei He for helping us prepare the molecular file. This project
1122 was partly supported by NIH grants U24HG012108 and P50 CA196530 awarded to Dr.
1123 Hongyu Zhao, and OpenAI Researcher Access Program and Google Cloud Research
1124 Program awarded to Tianyu Liu.

1125 8 Author Contributions Statement

1126 T.L. proposed this study. T.L., T.C., and X.L. designed the model. T.L. ran all the
1127 experiments. T.L., X.L. and H.Z. wrote the manuscript. H.Z. supervised this study.

1128 9 Ethics declarations

1129 9.1 Competing Interests Statement

1130 The authors declare no competing interests.

1131 9.2 Ethics and Inclusion

1132 Although BAITSAO is not biased on gender, races, and other factors, the users are
1133 solely responsible for the content they generate with models in BAITSAO, and there
1134 are no mechanisms in place for addressing harmful, unfaithful, biased, and toxic con-
1135 tent disclosure. Any modifications of the models should be released under different
1136 version numbers to keep track of the original models related to this manuscript. The
1137 users must comply with the laws of the country in which they are located.

1138 The target of current BAITSAO only serves for academic research. The users
1139 cannot use it for other purposes. Finally, we are not responsible for any effects produced
1140 by other users.

1141 10 Figure Legends

1142 Figure 1: An overview of BAITSAO as a FM under the pre-training and fine-
1143 tuning/zero-shot learning pipeline. (a) The pre-processing steps we used to transfer
1144 the meta information into embeddings to construct training datasets. (b) The model
1145 architecture of BAITSAO under a multi-task learning framework. (c) Comparisons of
1146 different methods for drug synergy analysis. (d) Different functions of BAITSAO.

1147 Figure 2: Investigation of drug embeddings. (a) The heatmap for the similarity of
1148 embeddings across all the drugs. (b) Exploration of drug similarity related to MK-
1149 4541. The drugs above the red line represent the two most similar drugs, while the
1150 drugs below the red line represent the most different drug. We list five types of clinical
1151 trial information ranked by the phases. Source data are provided as a Source Data file.

Figure 3: Results of evaluations for model structure, reliability and explainability. (a) Evaluations for BAITSAO and other methods. Each panel represents one metric with two datasets. The ranks of methods are averaged by datasets. Data are presented as boxplots (n=5 per group; center line, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers). The explanations of datasets D1-D3 are summarized in the Methods section. (b) The explainability of BAITSAO for the combination DEXAMETHASONE (drug)-DINACICLIB (drug) for different cell lines. We also list the ranks based on variance (HVG rank) and the adjusted p-value based on DESeq2 analysis results for each gene. The genes with adjusted p-value for multiple comparisons smaller or close to 0.05 are boldfaced. (c) The violin plot (n=6299 for non-synergy group; n=6116 for synergy group; center point, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers) for the outputs of BAITSAO (Probability) across the synergistic labels. We also present the two-side p-value in this figure. This panel supports the reliability of selected features from SHAP. (d) The rank-based plot between the number of drugs-cell line combinations with synergy (Count) and the index of drug combination (Index of drug combination). The index is ranked by the value of Count. We present the Pearson correlation (corr) and two-side p-value in this figure. Source data are provided as a Source Data file.

Figure 4: Statistics of the pre-training dataset from DrugComb. (a) The barplot for the number of *drug-cell lines* combinations by different tissues. (b) The barplot for the number of *cell lines* by different tissues. (c) The histogram for the distribution of synergy score computed based on Loewe [64]. The x-axis is transferred into log scale. (d) The histogram for the distribution of single-drug IC₅₀ levels. The x-axis is transferred into log scale. (e) 3D plot for the relation between single-drug IC₅₀ levels and synergy score. (f) The heatmap for the overlap of combinations across different tissues. Source data are provided as a Source Data file.

Figure 5: Results under the multi-task learning framework. (a) The Help-Harm matrix for different combinations of tasks. The values indicate the percentage (unit: %) of improvement using multi-task learning compared to single-task learning (STL) defined by the tasks in rows. The columns represent the paired tasks. We boldfaced blocks with increments larger than 0.5%, which is a threshold reported in [70] as acceptable improvement and half of the natural threshold 1%. (b) Comparisons for the results under MTL and STL. The metric for regression tasks, including Loewe and RLrow, is PCC. The metric for the classification task, including Classification, is ROCAUC. (c) Comparisons for the results under different training settings. Data are presented in boxplots (n=5 per group; center line, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers). Here *BAITSAO-FT* represents that we fine-tuned the pre-trained model, *BAITSAO-ZS* represents that we applied the pre-trained model for these tasks under zero-shot learning framework, and *BAITSAO-FS* represents that we did not use the pre-trained weights for these tasks. Here FT means fine-tuning, ZS means zero-shot learning and FS means from scratch. We included four metrics across three datasets for comparisons. (d) The first example of tri-drug cases for drug synergy prediction with BAITSAO. (e) The second

1196 example of tri-drug cases for drug synergy prediction with BAITSAO. Source data are
1197 provided as a Source Data file.

1198 Figure 6: Statistics of model training. (a) Ablation test results for BAITSAO with
1199 different input formats. (n=5 per group; center line, median; box limits, upper and
1200 lower quartiles; whiskers, up to $1.5 \times$ interquartile range; points, outliers). (b) The
1201 comparison of running time for different methods. We highlight the running time of
1202 BAITSAO and use the regression task as an example. (c) Plot for the proportion
1203 of training dataset and PCC for BAITSAO under the regression task. We reported
1204 $(\mu - \sigma, \mu + \sigma)$ for each proportion, where μ represents the mean and σ represents the
1205 standard deviation. (d) Plot for the proportion of training dataset and ROCAUC for
1206 BAITSAO under the classification task. We reported $(\mu - \sigma, \mu + \sigma)$ for each proportion.
1207 (e) Plot for layer width and PCC for BAITSAO under the regression task. We reported
1208 $(\mu - \sigma, \mu + \sigma)$ for each proportion. (f) Plot for layer width and ROCAUC for BAITSAO
1209 under the classification task. We reported $(\mu - \sigma, \mu + \sigma)$ for each proportion. Source
1210 data are provided as a Source Data file.

A Analysis of important genes based on synergy prediction and single-cell datasets for melanoma cells

In Figure 17, we show the explainability of different genes based on the drug combination: DEXAMETHASONE-DINACICLIB on the cell line A2058, which is a cell line for melanoma cells [97]. We repeated the training process three times with different random seeds and computed the average importance of genes for the analysis of a single cell line, thus we reduced the negative effect of randomness. We also found support for these two drugs as CDK inhibitors for affecting the expression levels of the 1st gene S100B [98, 99] and 2nd gene TYR [100, 101] ranked by the importance. Therefore, by analyzing the case of a single cell line, we linked the targets of individual drugs to genes, thus offering more information on treatments for the specific cancer type. Furthermore, we list the rank of selected genes based on their variance in Figure 17, and the top 5 genes have lower rank compared with important genes from all cell lines. Moreover, we considered running enrichment analysis based on GO and MsigDB for this set of genes, and the results are shown in Supplementary Figures 6 (b) and (d). From these two figures, we found that the collected genes were significantly enriched in pathways related to melanosome, which also has high relevance with melanoma [102]. Moreover, compared to the analysis of multiple cell lines, we obtained fewer cancer-specific signals for the analysis based on the single cell line.

We also considered two approaches to validate the selected genes from cell line A2058 with the help of scRNA-seq datasets. We utilized one scRNA-seq dataset combined by [103, 104], known as the disease-control scRNA-seq dataset, to verify whether the selected genes are differentially expressed in malignant cells from patients. We analyzed its differentially expressed genes (DEGs) by running the Wilcoxon rank-sum test between the diseased cells and the control cells. We found that 16 out of the 20 genes were significant in the malignant cells. For four genes that were not significant, their ranks of expression levels across all genes are 13000+ by descending, so their contributions might be hidden by their measured expression levels in this dataset. Therefore, BAITSAO successfully captured the difference between tumor microenvironments and healthy cells.

Moreover, we utilized the other scRNA-seq dataset [105] labeled as acral melanoma (AM) and cutaneous melanoma (CM), known as AM-CM scRNA-seq dataset, to analyze its differentially expressed genes (DEGs) by running the test of the different diseased cases from patients. We show the UMAP plots of this scRNA-seq dataset in Supplementary Figures 18 (a)-(b). We found that 19 out of the 20 genes had low adjusted p-values, which meant that most of them were DEGs by comparing AM cells with CM cells, thus supporting the potential contribution of our model to uncover tumor heterogeneity.

We summarized our statistical analysis for these two datasets in Supplementary file 2. Our analyses therefore suggest that BAITSAO can provide information about potential targets of combined drugs for melanoma or other cancer types.

1253 **B What are the advantages of BAITSAO**
1254 **comparing with other LLM-based methods?**

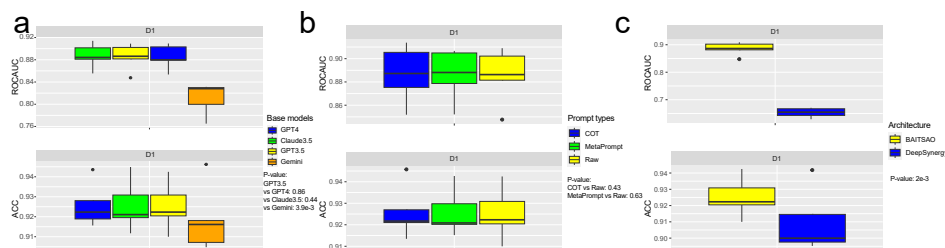
1255 In the area of drug synergy predictions based on Large language Models (LLMs),
1256 there are two existing methods, SynerGPT [25] and CancerGPT [26]. We discussed
1257 their limitations and differences from BAITSAO in the Introduction section of our
1258 submitted manuscript. Here we provide more details on the advantages of our method
1259 over these two other approaches by first summarizing the functions and challenges of
1260 these two methods in the Extended Data Table 1.

Methods	Functions	Drawbacks
SynerGPT	SynerGPT is a tool based on in-context learning. It models different examples containing drug pairs and one cell line based on their overlap to generate a synergy graph then uses the synergy graph as input to infer the synergetic effect of new examples.	1. SynerGPT can only handle simple binary classification. 2. It does not consider prior biological knowledge. 3. It cannot handle multi-drug cases (with # of drugs \geq 3). 4. Its performance is close to DeepDDS according to their benchmarking results in Table 2. 5. It cannot be used to analyze gene-drug interactions and gene-gene interactions. 6. It is closed-source.
CancerGPT	CancerGPT is a tool based on pre-training and fine-tuning. It first pre-trains a model based on GPT-2 with biomedical context, and then incorporates the drug pairs and one cell line into the input sentence to form a question. The model's output, as the answer, will be treated as the prediction result.	1. CancerGPT can only handle simple binary classifications. 2. It cannot handle multi-drug cases (with # of drugs \geq 3). 3. It cannot be used to analyze gene-drug interactions and gene-gene interactions. 4. It is closed-source.

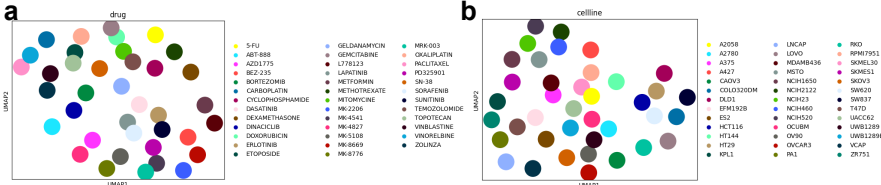
Supplementary Tab. 1 Comparison of SynerGPT and CancerGPT

1261 According to this table, both SynerGPT and CancerGPT have limitations that are
1262 addressed by BAITSAO, which is capable of using prior biological knowledge to pre-
1263 dict synergetic effects for multi-drug combinations under both the classification and
1264 regression frameworks. According to our benchmarking results, BAITSAO outper-
1265 forms baselines including DeepDDS, TreeComb, and MARSY significantly. Moreover,
1266 because SynerGPT and CancerGPT are not open-source (for CancerGPT we even can-
1267 not access their datasets for evaluation), we include BERT as a baseline to evaluate the
1268 performance of applying LLMs directly to address this problem, and BERT performs
1269 worse than BAITSAO. Moreover, we downloaded the datasets used by SynerGPT and
1270 performed prediction using BAITSAO, achieving an AUCROC score of 0.94, which is
1271 significantly higher than their reported score of 0.78. The results of our benchmarking
1272 analysis are shown in Figure 3. We are also willing to compare BAITSAO with Syn-
1273 erGPT and CancerGPT if they release their models in the future. Our method also
1274 allows researchers to explore gene-gene interactions as well as gene-drug interactions,
1275 which are not discussed in these two methods. Our method also introduces the idea
1276 of using embeddings from LLMs to analyze the similarity of drug functions. Finally,
1277 BAITSAO is not on a large scale (50M trainable parameters during the fine-tuning
1278 step), and it is fully open-source, which is friendly for the development of science and
1279 the extension from possible follow-up work.

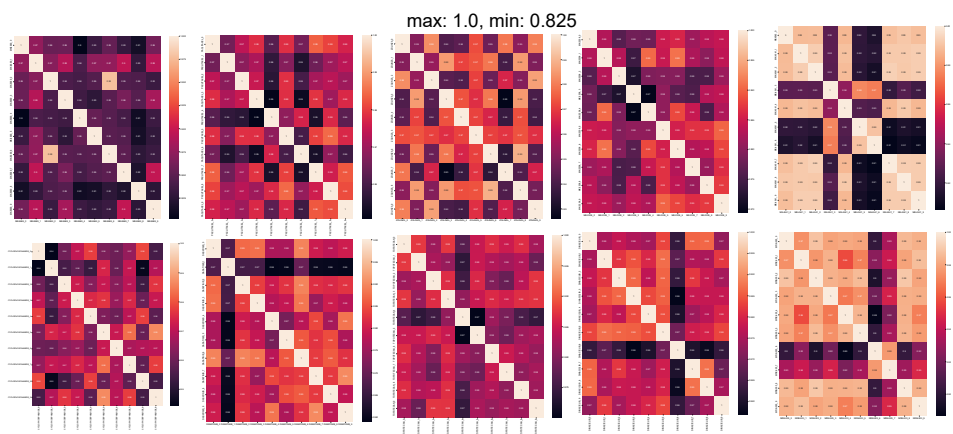
1280 C Supplementary figures



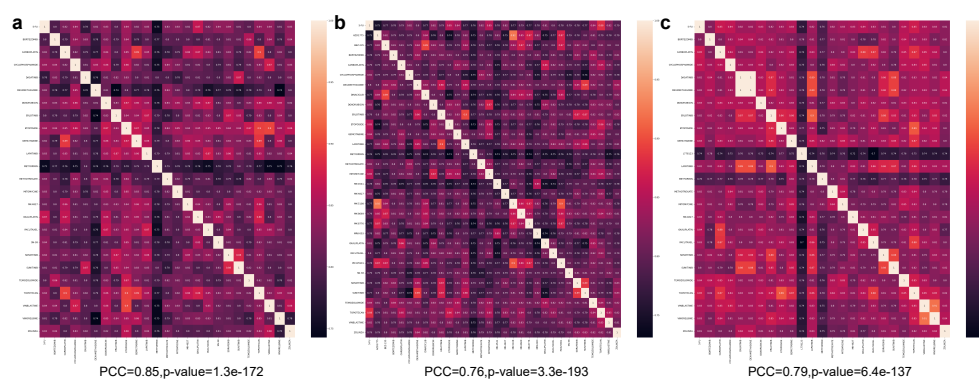
Supplementary Fig. 1 Evaluating the contributions of inputs and model architecture for drug synergetic effects. Data are presented in boxplots (n=5 per group; center line, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers). All of the experiments are performed based on dataset D1 and synergetic effect prediction is a classification task. The test is performed jointly for two metrics, based on Wilcoxon rank-sum test. (a) The performance comparison based on LLM embeddings from different base LLMs. (b) The performance comparison based on LLM embeddings from the drug descriptions generated by different prompts. (c) The performance comparison based on different model architectures.



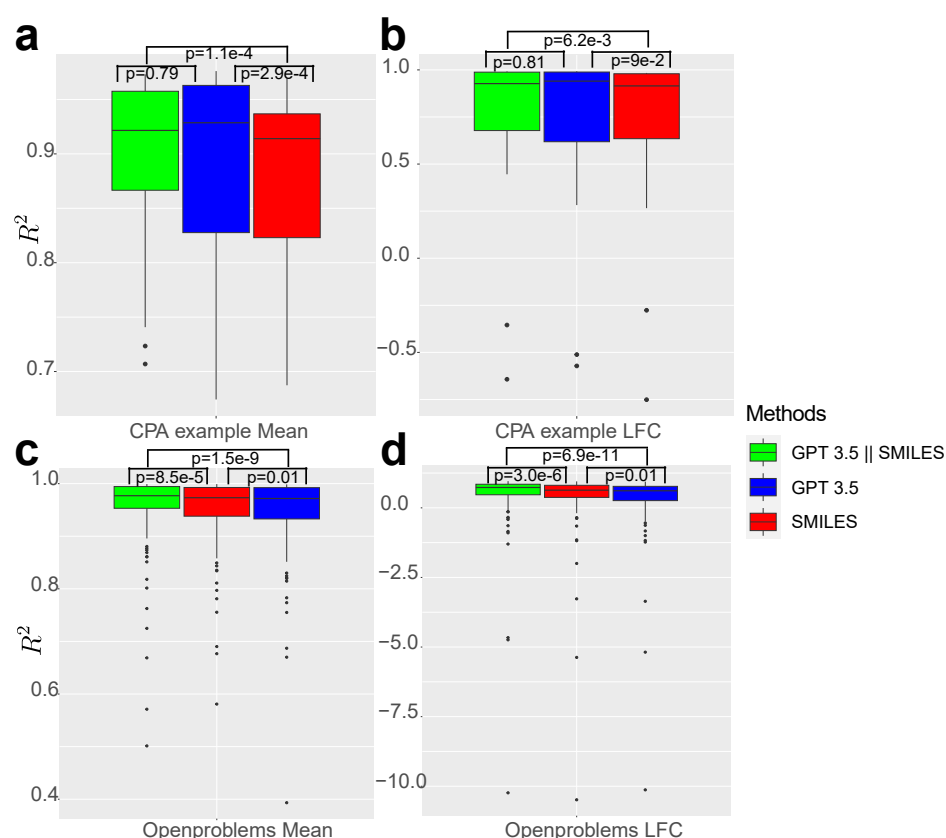
Supplementary Fig. 2 Visualization for the similarity of embeddings. (a) The UMAP plot for the drug embeddings from D1 colored by drug names. (b) The UMAP plot for the cell line embeddings from D1 colored by cell line names.



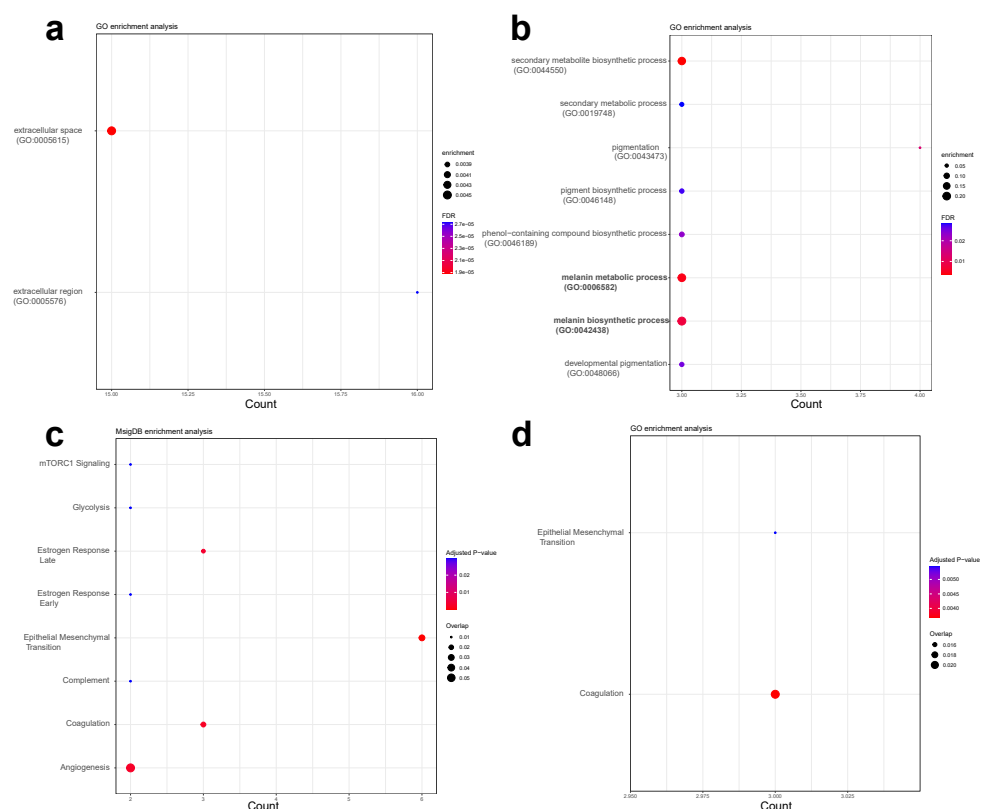
Supplementary Fig. 3 The group of heatmap for the similarity of embeddings under 10 random seeds of the 10 sampled drugs.



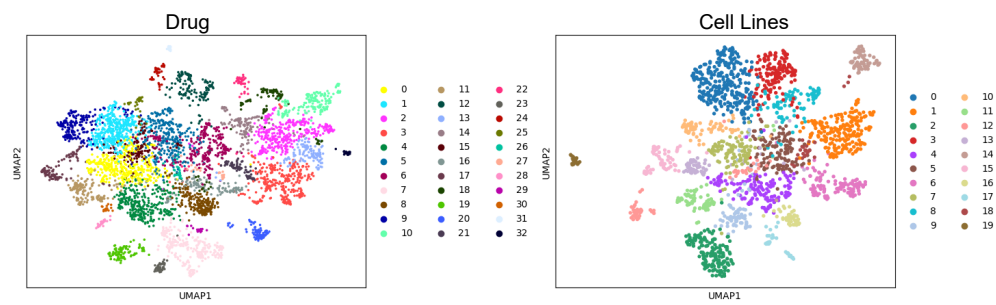
Supplementary Fig. 4 The group of heatmap for the similarity of embeddings from DrugBank under different drug properties. (a) The heatmap based on indication information. (b) The heatmap based on summary information. (c) The heatmap based on background information.



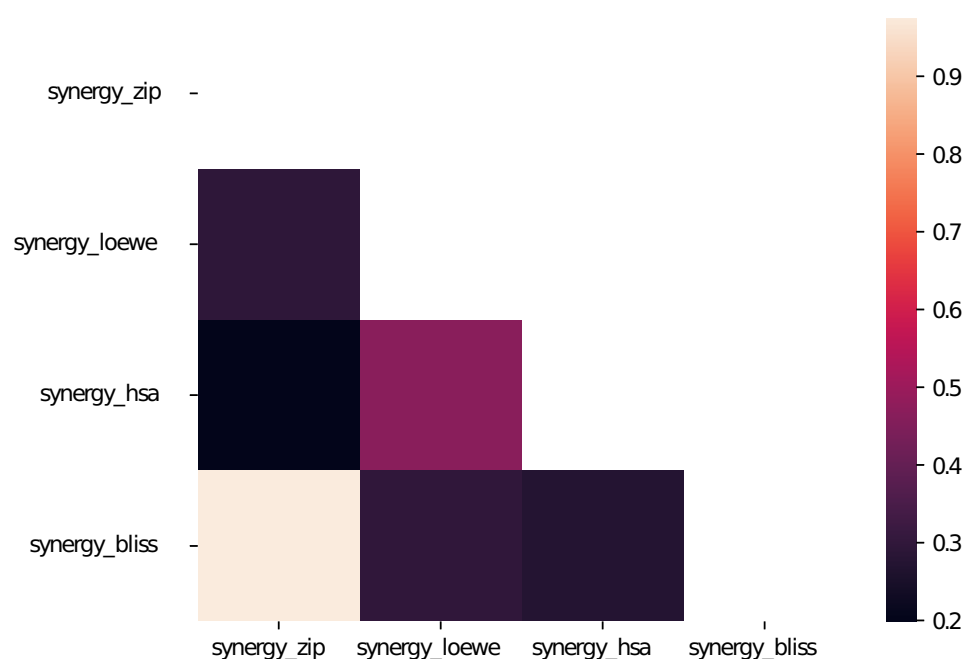
Supplementary Fig. 5 Exploration of drug-embeddings-augmented CPA results. Data are presented in boxplots (center line, median; box limits, upper and lower quartiles; whiskers, up to $1.5 \times$ interquartile range; points, outliers). Higher R^2 score means better performance. In panels (a)-(d), we also present the two-side p-values (p) based on Wilcoxon rank-sum test for each group. (a) R^2 of the expression levels of differentially expressed genes (DEGs) for the CPA example dataset colored by different settings (n=31). (b) R^2 of the log-fold change (LFC) of DEGs for CPA example dataset colored by different settings (n=31). (c) R^2 of the expression levels of DEGs for Openproblems dataset colored by different settings (n=137). (d) R^2 of the LFC of DEGs for Openproblems dataset colored by different settings (n=137).



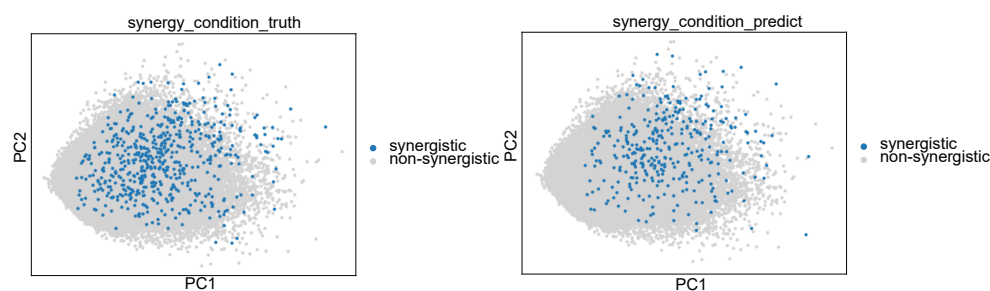
Supplementary Fig. 6 Results of gene enrichment analysis based on different databases. (a) Results of GO enrichment analysis using important genes from all cell lines. (b) Results of MsigDB enrichment analysis using important genes from all cell lines. (c) Results of GO enrichment analysis using important genes from the cell line A2058. We boldfaced the pathway related to melanosome. (d) Results of MsigDB enrichment analysis using important genes from the cell line A2058.



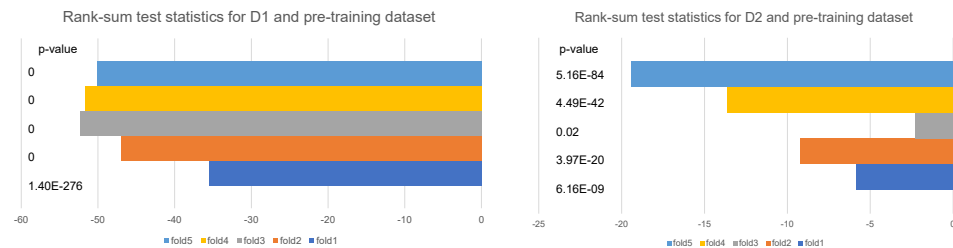
Supplementary Fig. 7 Visualization for the pre-training datasets. (a) The UMAP plot for the drug embeddings we used in the pre-training step, colored by Leiden clusters. (b) The UMAP plot for the cell-line embeddings we used in the pre-training step, colored by Leiden clusters.



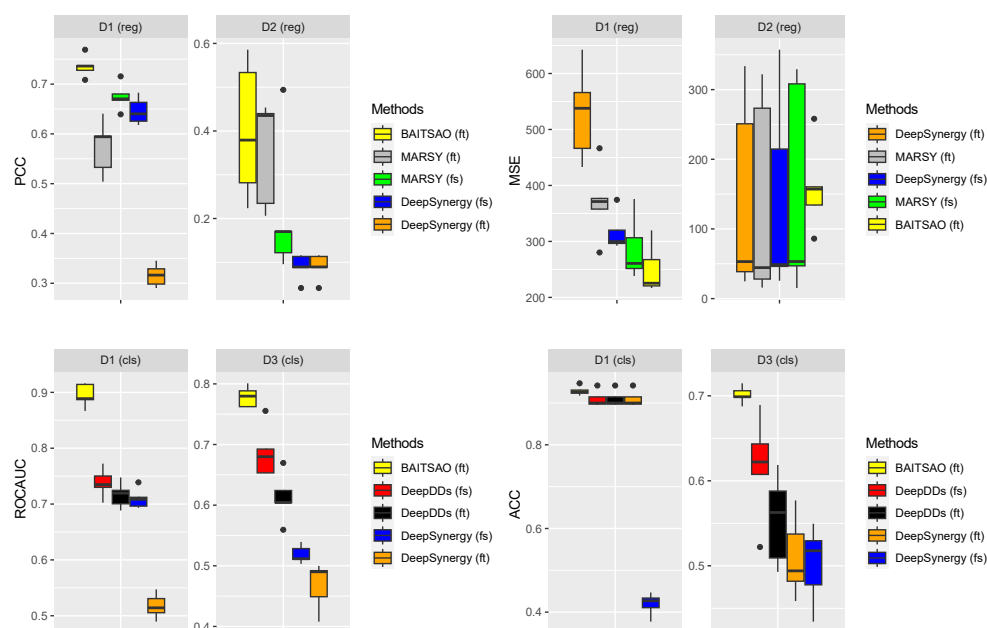
Supplementary Fig. 8 The heatmap for the PCC of different synergy scores from different computation methods.



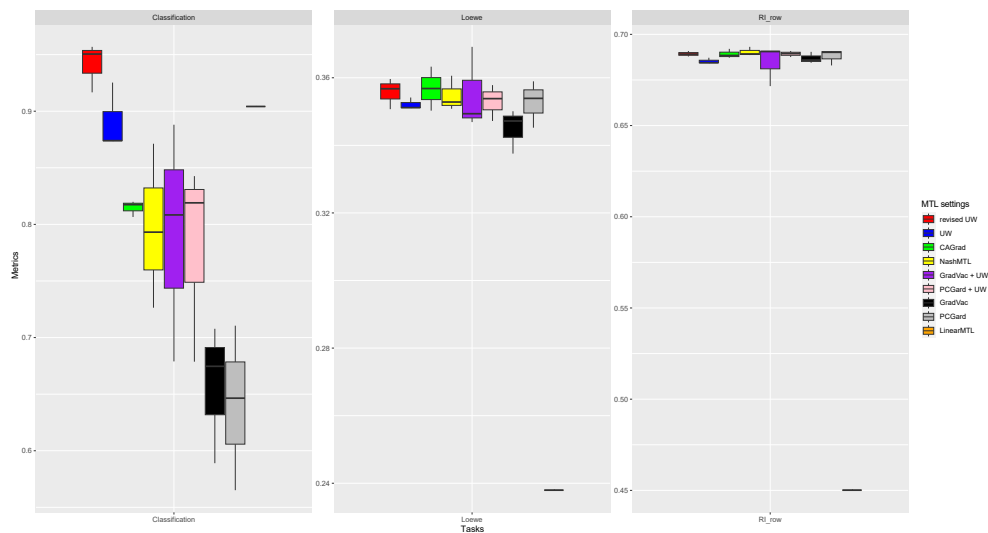
Supplementary Fig. 9 Visualization for the trained feature space of pre-training datasets (sub-sampled to 10 %). (a) The PCA plot from the outputs based on hidden layers of BAITSAO, colored by the ground truth synergistic information. (b) The PCA plot from the outputs based on hidden layers of BAITSAO, colored by the predicted synergistic information.



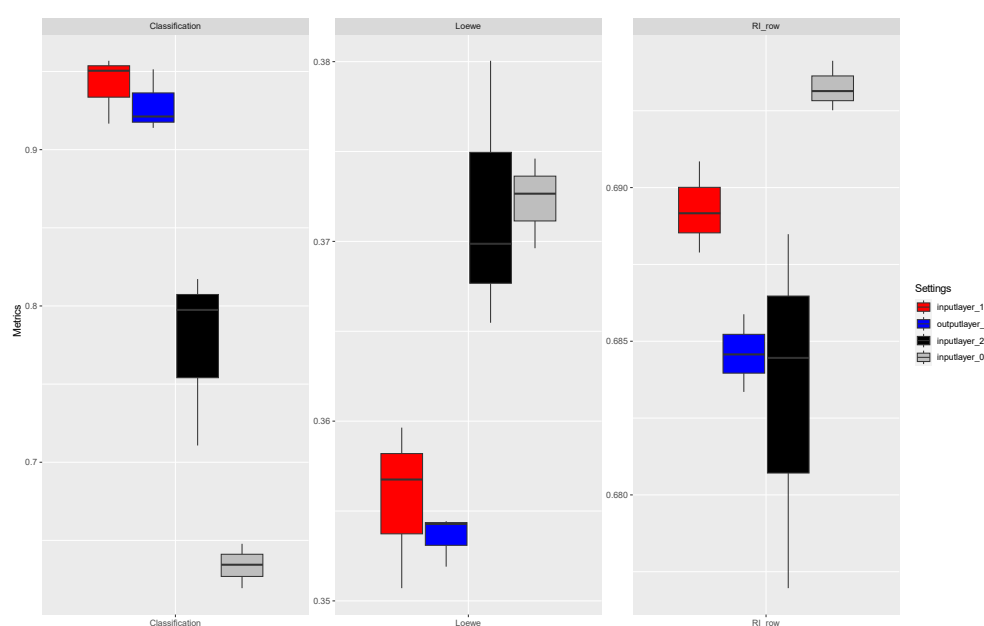
Supplementary Fig. 10 Statistics of the Rank-sum tests between the pre-training dataset and fine-tuning datasets. We included the two-side p-value for each comparisons. (a) The statistics of different folds by checking whether samples from D1 and the pre-training dataset come from the same distribution or not. (b) The statistics of different folds by checking whether samples from D2 and the pre-training dataset come from the same distribution or not.



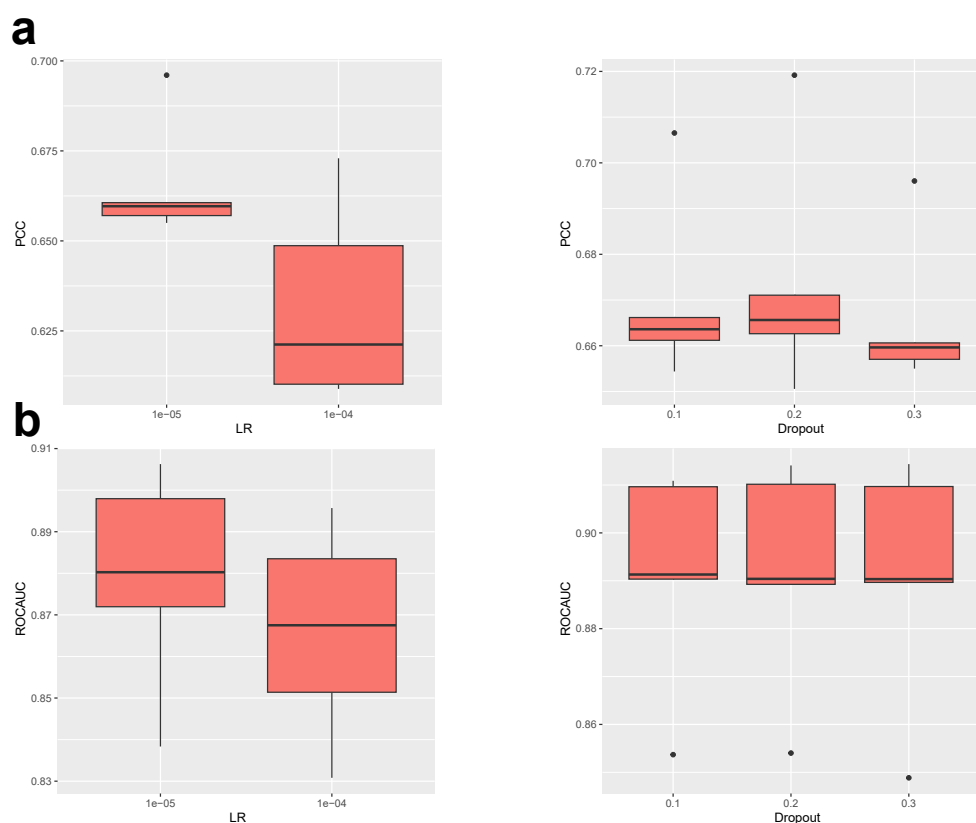
Supplementary Fig. 11 Comparisons of deep-learning-based models with pre-training (ending with ft) and from scratch (ending with fs) for two tasks across three different datasets. Data are presented in boxplots (n=5 per group; center line, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers).



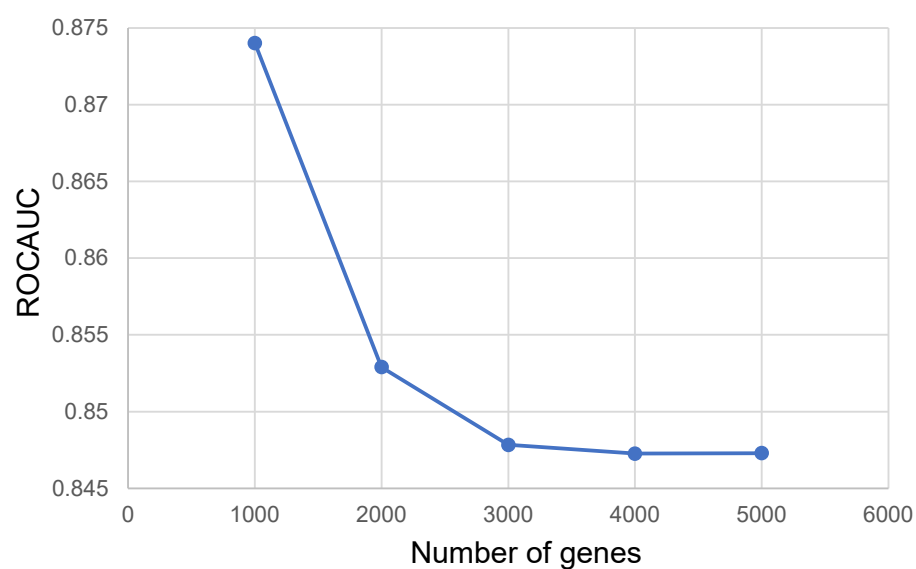
Supplementary Fig. 12 Ablation tests for MTL strategies. Here *PCGrad+UW* means we combine PCGrad with UW, and *GradVac+UW* means we combine GradVac with UW. *revised UW* represents the modified UW in this manuscript. For Loewe and RI_{row}, we present scores of PCC. For classification, we present scores of ROCAUC. Data are presented in boxplots (n=3 per group; center line, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers).



Supplementary Fig. 13 Ablation tests for the number of task-specific layers. Here *inputlayer_1* represents using one layer for processing input data and one layer for model's output, and it is our final choice. *outputlayer_2* represents using one layer for processing input data and two layers for model's output. *inputlayer_2* represents using two layers for processing input data and one layer for model's output. *inputlayer_0* represents we did not set task-specific layers for input data. Data are presented in boxplots (n=3 per group; center line, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers).



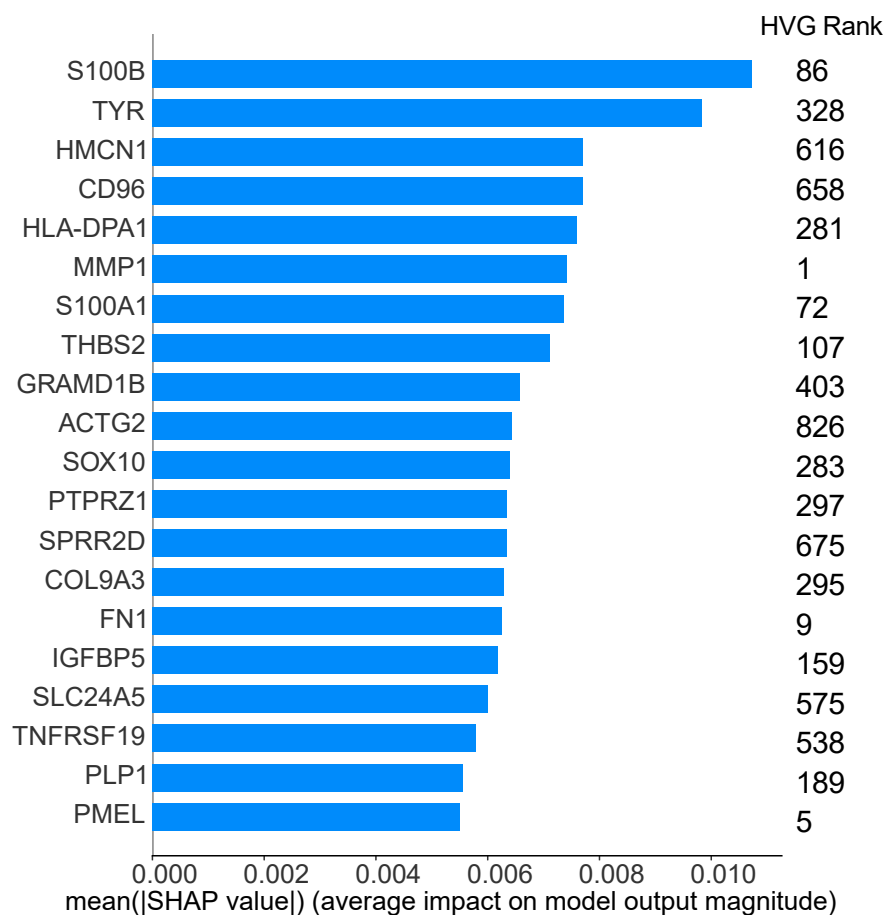
Supplementary Fig. 14 Results of hyper-parameter tuning of BAITSAO for D1. (a) The tuning results for the regression task. (b) The tuning results for the classification task. Data are presented in boxplots (n=5 per group; center line, median; box limits, upper and lower quartiles; whiskers, up to 1.5×interquartile range; points, outliers).



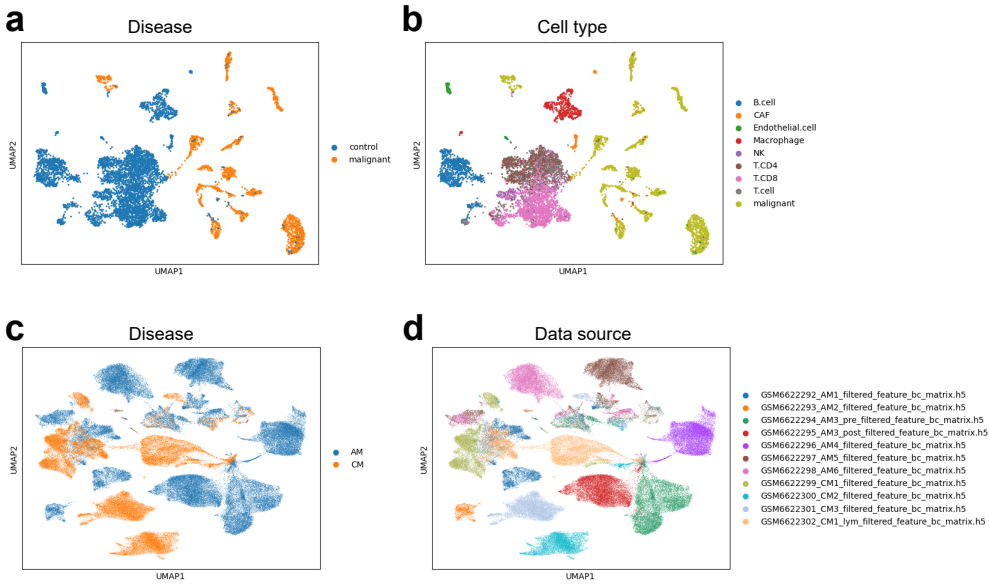
Supplementary Fig. 15 Results of tuning number of highly-variable genes for the analysis of explainability.



Supplementary Fig. 16 Visualization of important genes across different cell lines. We also annotate the rank based on variance for each gene.



Supplementary Fig. 17 The explainability of BAITSAO for the combination: DEXAMETHASONE (drug)-DINACICLIB (drug)-A2058 (cell line). We also annotate the rank based on variance for each gene.



Supplementary Fig. 18 Visualization for two scRNA-seq melanoma datasets. (a) The UMAP plot for the scRNA-seq dataset colored by the conditions of cells from the diseased-control scRNA-seq dataset. (b) The UMAP plot for the scRNA-seq dataset colored by cell types from the diseased-control scRNA-seq dataset. (c) The UMAP plot for the scRNA-seq dataset colored by the conditions of cells from the AM-CM scRNA-seq dataset. (d) The UMAP plot for the scRNA-seq dataset colored by the sources of cells from the AM-CM scRNA-seq dataset.