

# TranSiGen: Deep representation learning of chemical-induced transcriptional profile

Xiaochu Tong<sup>1,2</sup>, Ning Qu<sup>1,2</sup>, Xiangtai Kong<sup>1,2</sup>, Shengkun Ni<sup>1,2</sup>, Kun Wang<sup>3,1</sup>, Lehan Zhang<sup>1,2</sup>, Yiming Wen<sup>1,2,4</sup>, Sulin Zhang<sup>1,2</sup>, Xutong Li<sup>1,2,\*</sup> & Mingyue Zheng<sup>1,2,4\*</sup>

\*Correspondence: [lixutong@simmm.ac.cn](mailto:lixutong@simmm.ac.cn); [myzheng@simmm.ac.cn](mailto:myzheng@simmm.ac.cn)

<sup>1</sup>Drug Discovery and Design Center, State Key Laboratory of Drug Research, Shanghai Institute of Materia Medica, Chinese Academy of Sciences, 555 Zuchongzhi Road, Shanghai 201203, China

<sup>2</sup>University of Chinese Academy of Sciences, No. 19A Yuquan Road, Beijing 100049, China

<sup>3</sup>School of Life Sciences, Division of Life Sciences and Medicine, University of Science and Technology of China, Hefei 230026, China

<sup>4</sup>School of Pharmaceutical Science and Technology, Hangzhou Institute for Advanced Study, University of Chinese Academy of Sciences, Hangzhou 310024, China

## Abstract

With the advancement of high-throughput RNA sequencing technologies, the use of chemical-induced transcriptional profiling has greatly increased in biomedical research. However, the usefulness of transcriptomics data is limited by inherent random noise and technical artefacts that may cause systematical biases. These limitations make it challenging to identify the true signal of perturbation and extract knowledge from the data. In this study, we propose a deep generative model called Transcriptional Signatures Generator (TranSiGen), which aims to denoise and reconstruct transcriptional profiles through self-supervised representation learning. TranSiGen uses cell basal gene expression and compound molecular structure representation to

26 infer the chemical-induced transcriptional profile. Results demonstrate the effectiveness of  
27 TranSiGen in learning and predicting differential expression genes. The representation derived  
28 from TranSiGen can also serve as an alternative phenotype information, with applications in  
29 ligand-based virtual screening, drug response prediction, and phenotype-based drug repurposing.  
30 We envisage that integrating TranSiGen into the drug discovery and mechanism research pipeline  
31 will promote the development of biomedicine.

32

## 33 **Introduction**

34 Transcriptional profiling is a powerful tool that is widely used to characterize the phenotype  
35 of various cells and organisms. It captures the landscape of thousands of genes in different  
36 biological circumstance at a holistic level, providing abundant information on cellular and  
37 organismal status. Transcriptomics data analysis helps us further understand the changes in  
38 cellular functionality during the pathogenesis of disease and identify the regulatory mechanisms  
39 of cells in response to different perturbations.

40 With the significant advancements in high-throughput RNA sequencing (RNA-seq)  
41 technologies, the availability of comprehensive and systematic perturbational gene expression  
42 profiles has greatly expanded. Connectivity Map (CMap)<sup>1</sup>, a public database of perturbational  
43 profiles, initially contained data of 164 different bioactive small molecules in four cell lines. To  
44 obtain large-scale CMap data in a cost-effective and high-throughput manner, a novel gene  
45 expression profiling platform called L1000 was introduced by the Library of Integrated Network-  
46 based Cell-Signature (LINCS) program. L1000 measures the expression of 978 landmark genes  
47 to capture most of the information in the full transcriptome<sup>2</sup>. In 2017, the CMap-L1000v1 dataset

48 (Phase I, GEO: GSE92742) was reported to contain L1000 profiles from 42,080 genetic and  
 49 small-molecule perturbations across approximately 80 cell lines, with ongoing updates. The  
 50 PANACEA data, being developed by the Columbia Cancer Target Discovery and Development  
 51 Center, is expected to contain dose-responses and RNA-seq profiles from perturbations with  
 52 approximately 400 clinical cancer drugs across 25 cell lines<sup>3</sup>. In addition, there are ongoing  
 53 efforts to collect and organize gene expression profiles from publicly available databases to  
 54 create standardized and unified transcriptional data. ARCHS4 is a web resource containing RNA-  
 55 seq data from human and mouse<sup>4</sup>, while ChemPert focuses on RNA-seq data for non-cancer cell  
 56 lines<sup>5</sup>.

57 These large-scale perturbational gene expression profiles play an indispensable role in drug  
 58 discovery and mechanism research. The CMap project proposed a pattern-matching strategy to  
 59 identify compounds that share a mechanism of action (MOA). This approach led to the discovery  
 60 of a novel mechanism for gedunin, which acts as an inhibitor of HSP90 function<sup>1</sup>. Associating  
 61 chemical-induced transcriptional profiles with diseases has several benefits, such as screening  
 62 candidate compounds for diseases and revealing mechanisms of drug resistance. Chen et al.  
 63 proposed a method to measure the potency of reversing disease gene expression and quantify the  
 64 reversal relationship between disease and drug gene expression profiles. Through their systematic  
 65 exploration, they validated that pyrvinium pamoate, a FDA-approved drug for the treatment of  
 66 pinworms, also has in vivo efficacy for liver cancer<sup>6</sup>. In addition, Wei et al. identified that  
 67 rapamycin can reverse glucocorticoid resistance by screening chemical-induced profiles. They  
 68 discovered that MCL1 is an important regulator of glucocorticoid-induced apoptosis<sup>7</sup>.

69 Machine learning algorithms have been used to automatically capture relationships between

different perturbational profiles. Pabon et al. employed a random forest (RF) model to learn the relationship between chemical-induced profiles and gene knockdown-induced profiles, thereby identifying potential targets of compounds<sup>8</sup>. We proposed a graph convolutional network model SSGCN, which integrates protein-protein interaction networks and detects intricate relationships behind perturbational gene expression profiles. This approach facilitates the inference of protein targets for compounds<sup>9</sup>.

While numerous drug-like molecules have undergone high-throughput transcriptional perturbation experiments, it is not feasible to explore all possible perturbation-cell combinations due to the vast size of the combinatorial space. Therefore, a promising approach is to develop deep learning models capable of learning from high-dimensional data extracted from publicly accessible transcriptional profiles. By utilizing these models to predict gene expression profiles, it becomes possible to access the unexplored space of perturbation transcriptomics. Recently, several studies have explored the use of deep learning models to predict transcriptional profiles for new chemicals. Zhu et al. proposed a deep neural network DLEPS, which is capable of fitting chemical-induced gene expression values without considering the types of cell lines. This model has been utilized to identify potential candidates for obesity, hyperuricemia and nonalcoholic steatohepatitis<sup>10</sup>. On the other hand, DeepCE<sup>11</sup> and CIGER<sup>12</sup> leveraged one-hot encoding to discriminate different cell types, and learned from perturbation profiles of various cells. These two models focused on predicting differential gene expression profiles disrupted by novel chemicals, and have been applied to the drug repurposing pipeline using drugs from the DrugBank database<sup>11,12</sup>.

However, existing models that directly fit gene expression values via supervised learning

may not effectively distinguish among the true perturbation signals, confounding factors, and high level of noises in expression profiles. Recent studies have explored the use of variational autoencoder (VAE) and its variants for denoising, dimensionality reduction, imputing missing values<sup>13</sup>, and extracting meaningful biological signals<sup>14</sup> from gene expression profiles. While these studies have demonstrated the impressive capabilities of VAE-based models in processing high-dimensional and noisy transcriptomics data, these analyses are limited to the reconstruction of experimental perturbation profiles and could not be extended to novel perturbations.

To address these limitations, we propose a VAE-based framework called Transcriptional Signatures Generator (TranSiGen) for learning and reconstructing the transcriptional profiles. TranSiGen is designed to alleviate inherent limitations of data through self-supervised representation learning and can be used to infer new perturbational profiles. We evaluated the performance of TranSiGen in fitting basal profiles  $X_1$ , perturbational profiles  $X_2$ , and the corresponding differential expression genes (DEGs). Furthermore, we analyzed TranSiGen-derived DEGs to assess its ability to effectively learn cellular and compound features from perturbation profiles. In addition, we compared the performance of TranSiGen with other baseline models in the task of inferring DEGs. Finally, we demonstrated the effectiveness of TranSiGen-derived representation as a type of phenotype information by applying it to various downstream tasks.

The highlights of TranSiGen are as follows: (1) TranSiGen uses self-supervised representation learning strategy on a VAE-based model to denoise and reconstruct transcriptomics data. This approach yields impressive performance in inferring chemical-induced transcriptional profiles. (2) The perturbational expression profiles obtained by TranSiGen

effectively learn cellular and compound features from data, and it can serve as a new, unified and standardized representation for characterizing chemical-induced phenotype information. (3) TranSiGen-derived representation has demonstrated its applications in a range of downstream tasks, including ligand-based virtual screening, drug response prediction, and phenotype-based drug repurposing for disease.

## Results

### The Overview of TranSiGen

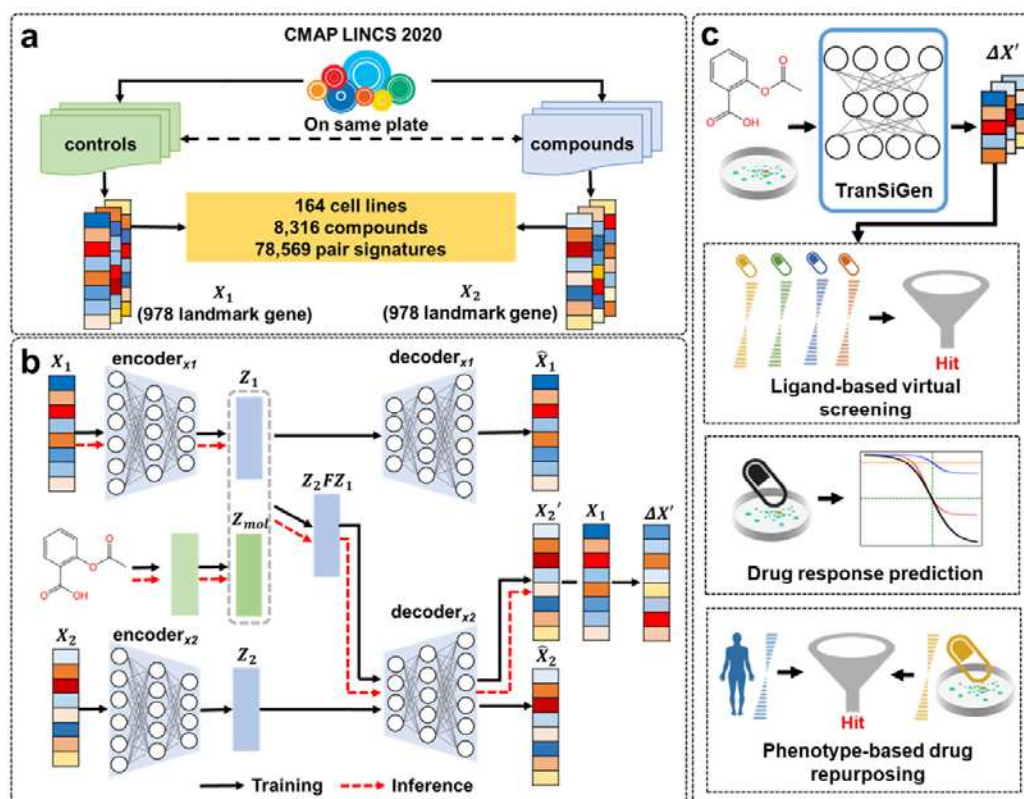
TranSiGen is a VAE-based model that simultaneously learns three distributions: basal profiles without perturbation, perturbational profiles, and the mapping relationship between them. It utilizes a self-supervised representation learning strategy to mitigate noise effects in the transcriptional profile and uncover the signal of perturbation.

The transcriptional profiles used in the model are obtained from level3 data of the newly released CMAP LINCS 2020 dataset<sup>2,15</sup>. These profiles consist of 978 measured landmark genes per profile. Specifically, basal profiles ( $X_1$ ) represent control profiles treated with DMSO, while perturbational profiles ( $X_2$ ) represent transcriptional profiles treated with compounds. For each plate, the DMSO-treated control profile from the same plate is selected as  $X_1$ , forming a paired  $X_1 \sim X_2$ . The dataset includes 219,650  $X_1 \sim X_2$  pairs for 8,316 compounds across 164 cell lines. Since L1000 assays are typically conducted with three or more biological replicates, there may be multiple  $X_1 \sim X_2$  pairs for a perturbation-cell combination in the dataset. To ensure only one  $X_1 \sim X_2$  pair per perturbation on each cell line, the repeated  $X_1$  and  $X_2$  pairs were further processed using the moderated-Z weighted averages algorithm (MODZ). The processed data consists of

transcriptional profiles for 8,316 compounds on 164 cell lines, including 78,569  $X_1 \sim X_2$  pairs (Fig. 1a).

TranSiGen consists of two VAE models: one encodes basal profiles  $X_1$ , and another encodes perturbational profiles  $X_2$  (Fig. 1b and Supplementary Fig. 1). It learns to map from  $X_1$  and the perturbation representation to  $X_2$ , which is denoted as  $X_2'$ . During inference, TranSiGen generates  $X_2'$  from the input  $X_1$  and the perturbation representation (Fig. 1b), and finally obtains the inferred DEGs  $\Delta X'(X_2' - X_1)$  of the compound. A complete list of the symbols and notations used here were summarized in Table 1.

In downstream applications, TranSiGen can generate perturbational profiles for numerous compounds, allowing exploration of a larger space that is not covered by training data. The perturbational representation derived from TranSiGen can be applied to ligand-based virtual screening, drug response prediction in cells, and phenotypic screening of candidate compounds for disease (Fig. 1c).



**Fig. 1** TranSiGen's architecture and application. **a** The data processing flow for TranSiGen. **b** The architecture and inference process of TranSiGen. **c** The applications of TranSiGen-derived representation.

### TranSiGen enables effective learning for transcriptional profiling

In this study, TranSiGen was used to simultaneously fit the basal profile and the perturbational profile. The model's performance in learning , and the corresponding DEGs was evaluated individually. As shown in Fig. 2a, TranSiGen exhibits excellent performance in reconstructing and , denoted as and , with the Pearson's correlation coefficients (PCC) close to 1 (between and , between and ). It also performs well in inferring , which is predicted by and compound representation. Compared to directly evaluating the performance of learning and , the corresponding performance in



161 fitting DEGs is slightly decreased, with the PCC 0.734 and 0.619 in reconstructing  $\Delta\hat{X}$  (between  
162  $\Delta\hat{X}$  and  $\Delta X$ ) and predicting  $\Delta X'$  (between  $\Delta X'$  and  $\Delta X$ ), respectively. Additionally, the  
163 relationship between TranSiGen's performance and  $X_1 \sim X_2$  correlation coefficient ( $R^2$ ) was  
164 analyzed. As shown in Fig. 2b, the sample size of the profiles increases with  $X_1 \sim X_2 R^2$ , as well  
165 as the prediction performance for DEGs. For  $X_1 \sim X_2 R^2 > 0.8$ , there is a slight decrease in  
166 performance, possibly due to the perturbation effects being too subtle for the model to fully  
167 capture. Overall, the model has learned the meaningful mapping from  $X_1$  and compound to  $X_2$ .

168 Furthermore, we evaluated the profiling capabilities of TranSiGen by analyzing its  
169 effectiveness in learning cellular and compound representations in  $\Delta X'$ . Fig. 2c presents a  
170 visualization of dimensionality reduction for both experimental  $\Delta X$  and TranSiGen-derived  $\Delta X'$ ,  
171 with each point color-coded by cell type. In the case of  $\Delta X$ , there was some clustering of the  
172 same cells, but also significant mixing between different cell types. In contrast, TranSiGen-  
173 derived  $\Delta X'$  exhibited clear clustering of the same cells and sharper distinctions between different  
174 cell types. This suggests that the representation derived from TranSiGen can more effectively  
175 differentiate between various cell types compared to experimental profiling, which is subject to  
176 high level of noise. Moreover, for compounds like decitabine, hydroxyurea and fludarabine, the  
177  $\Delta X'$  for different cells are closely grouped together, indicating their similar perturbation effects  
178 across different cells, as they all directly induce cell death due to cytotoxicity<sup>16</sup>. In addition to  
179 cytotoxic compounds, we expect other compounds sharing the same MOA to display similar  
180 effects on transcriptional profiling. The correlation between compounds with the same MOA was  
181 analyzed and is showed in Fig. 2d. TranSiGen-derived representations have higher PCC for  
182 compounds with the same MOA than  $\Delta X$ . Meanwhile, when compared to random MOA,

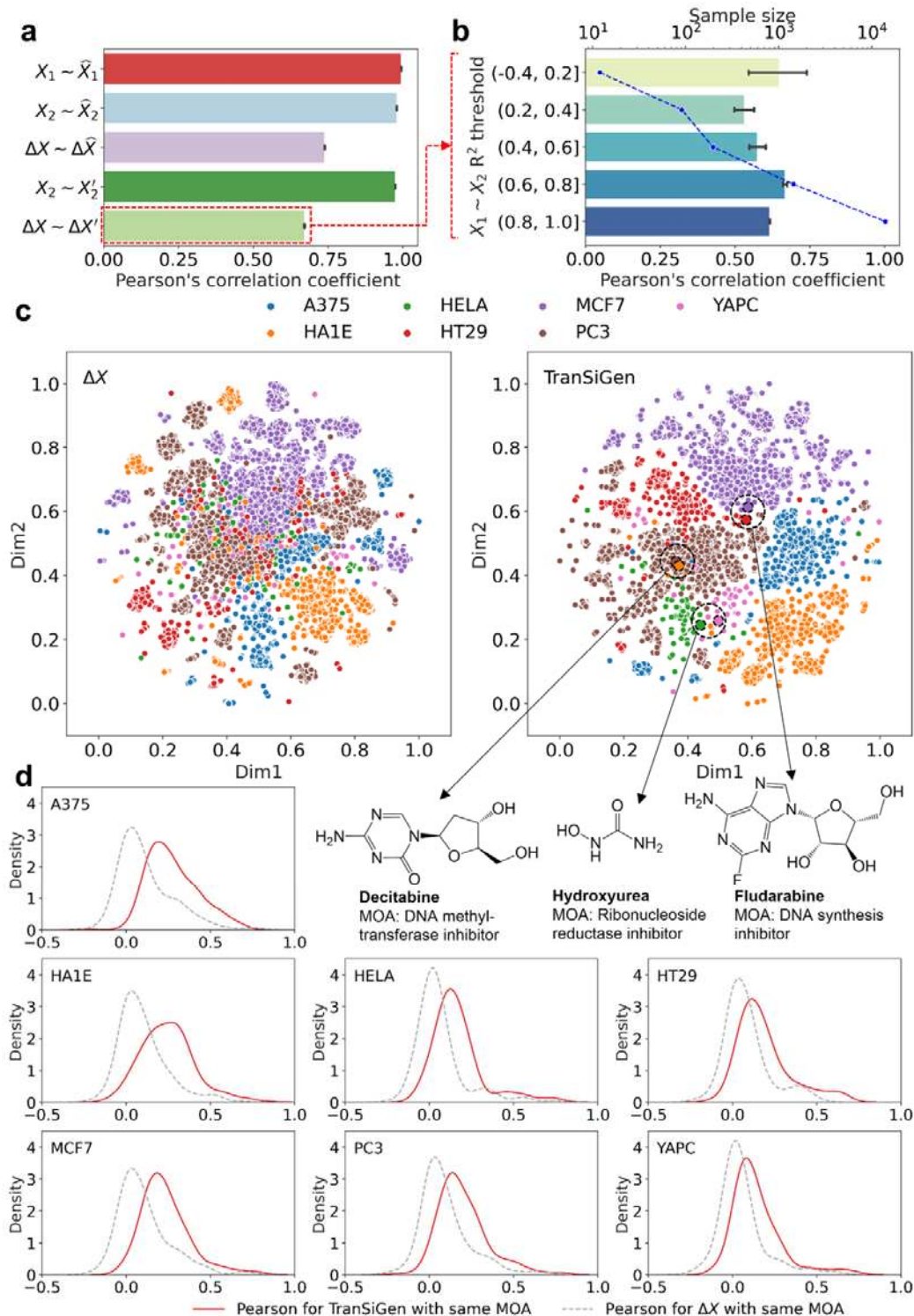
183 TranSiGen-derived representations of the same MOA also exhibit relatively high PCC

184 (Supplementary Fig. 2).

185 Overall, TranSiGen's self-supervised representation learning helps denoise and reconstruct

186 transcriptional profiles, effectively identifying and learning meaningful cellular and compound

187 representations from data.



188

189 **Fig. 2** Transcriptional profiling representation learning by TranSiGen. **a** Performance of

190 TranSiGen in transcriptional profiling reconstruction and prediction. **b** The change of

191 TranSiGen's performance for  $\Delta X'$  with the correlation between  $X_1$  and  $X_2$ . **c** Dimensionality  
 192 reduction visualization using  $\Delta X$  and TranSiGen-derived  $\Delta X'$  for different cell lines. **d**  
 193 Distribution of Pearson's correlation coefficients of profiles for the same MOA by  $\Delta X$  and  
 194 TranSiGen-derived  $\Delta X'$ .

## 195 **Comparison with existing models in inferring differential expression genes**

196 In this section, we will compare the performance of TranSiGen with other baseline models  
 197 in inferring differential expression genes. The models include DLEPS<sup>10</sup>, DeepCE<sup>11</sup> and CIGER<sup>12</sup>.  
 198 Among them, DLEPS<sup>10</sup> predicts one DEG profiling for each compound without considering cell  
 199 lines, while DeepCE<sup>11</sup>, CIGER<sup>12</sup>, and our model TranSiGen are all capable of inferring cell-  
 200 specific DEG profiling.

201 We evaluated the performance of the models using three data splitting methods: random,  
 202 chemical-blind, and cell-blind. Since not all pairwise combinations of cells and perturbations in  
 203 the L1000 dataset have experimental profiles, there are missing values in the perturbation-cell  
 204 combination space. In random splitting (scenario 1), the model predicts DEGs for new  
 205 combinations of compounds and cells, allowing for the imputation of missing values. In  
 206 chemical-blind splitting (scenario 2), the model's capacity to extrapolate profiles for novel  
 207 compounds (using test compounds not included in the training set) is evaluated. This scenario  
 208 includes two tests: In scenario 2-1, a dataset with 355 compounds across 7 cells was used to  
 209 ensure comparability among multiple models. In scenario 2-2, a complete dataset with 8,316  
 210 compounds across 164 cells was used to evaluate TranSiGen's performance. In cell-blind  
 211 splitting (scenario 3), cells not included in the training set were used to infer chemical-induced  
 212 profiles on new cells. TranSiGen characterizes cells using transcriptional profiles without

213 compound treatment, enabling exploration of performance in cell-blind splitting, a scenario that  
214 other models cannot address. The model was trained on 10, 50, and 150 cells, and evaluated on 7  
215 new cells. The diagram illustrating the data splitting for chemical-blind and cell-blind scenarios is  
216 shown in Fig. 3a.

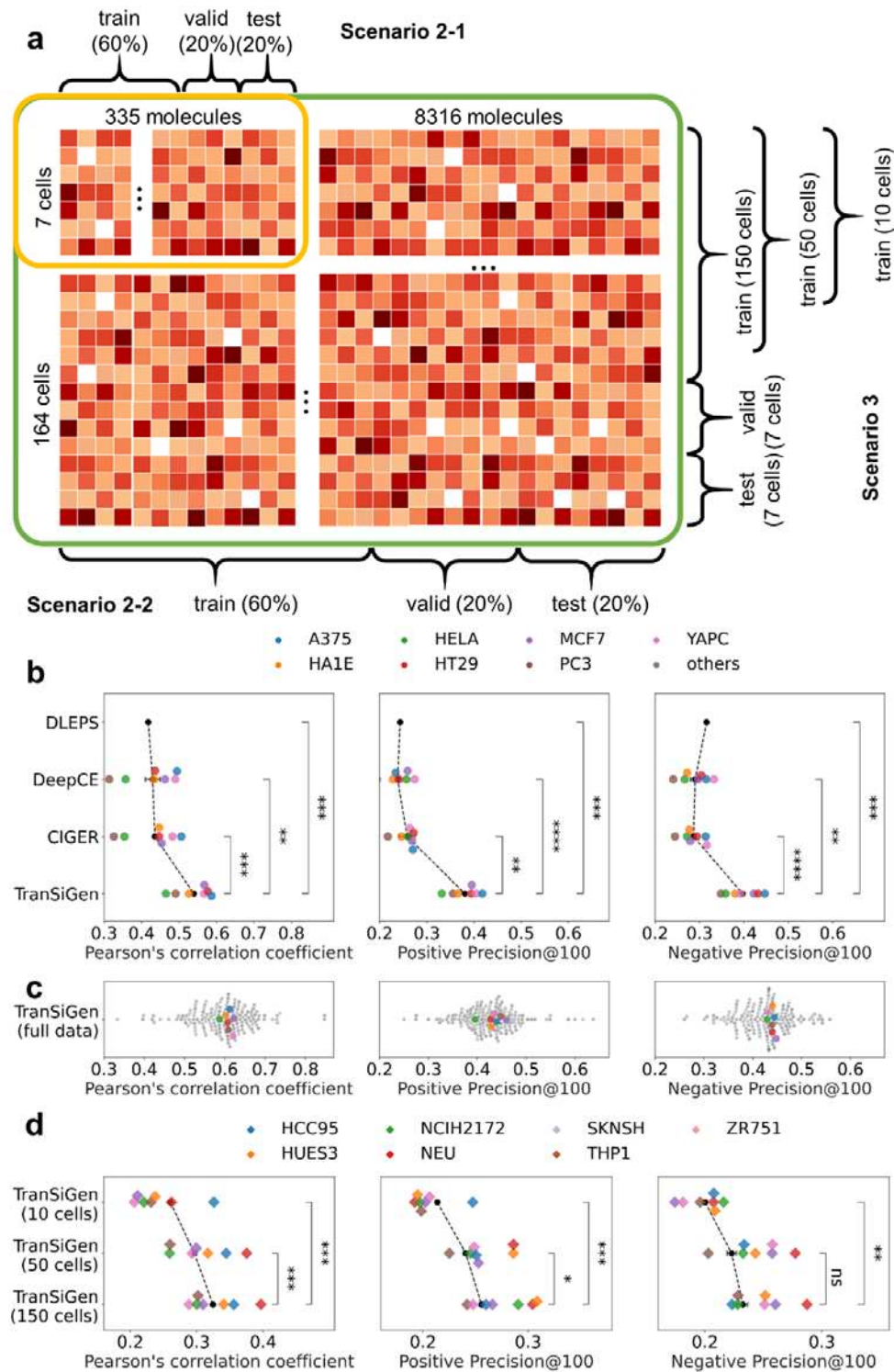
217 According to the results shown in Fig. 3b, TranSiGen performs better than the other three  
218 models in scenario 2-1, in terms of the average PCC (represented by the black dots) evaluated on  
219 seven cell lines. Additionally, TranSiGen, along with the other two models (DeepCE and CIGER)  
220 that consider cell lines contexts, shows similar trends in the prediction performance of DEGs on  
221 different cell lines. TranSiGen also outperforms the other three models in terms of Positive  
222 Precision@100 and Negative Precision@100 when evaluating metrics focusing on the most  
223 significantly up- and down-regulated expressed genes. Furthermore, TranSiGen achieves state-of-  
224 the-art results on the full dataset (scenario 2-2), which includes more cell lines and molecules  
225 (Fig. 3c and Supplementary Table 1). Compared to the significant performance differences across  
226 seven different cells in Fig. 2b, TranSiGen (full data) demonstrates considerable capabilities in  
227 inferring DEGs for these seven cell lines by benefiting from more training data. A similar trend is  
228 observed in random splitting, where TranSiGen performs better than other models, and it also  
229 demonstrates decent performance on the full dataset (Supplementary Table 2). Meanwhile, all  
230 models exhibit better performance in random splitting than in chemical-blind splitting due to the  
231 inclusion of compounds and cell lines seen during training (Supplementary Table 1 and 2).

232 In the context of chemical-blind splitting, different molecular representation and model  
233 initialization methods were tested. It was found that initializing the parameters of TranSiGen  
234 using perturbational profiles by gene knockdown leads to better performance compared to

235 random initialization. Additionally, the pre-training representation using Knowledge-guided Pre-  
 236 training of Graph Transformer (KPGT)<sup>17</sup> further enhances the performance of inferring DEGs,  
 237 surpassing the molecular fingerprint ECFP4 (Supplementary Table 1).

238 For cell-blind splitting, as the number of cells in the training set increases, TranSiGen  
 239 demonstrates improved performance in inferring DEGs on seven unseen cells during training  
 240 (Fig. 3d and Supplementary Table 3). However, its ability to generalize to new cells is not as  
 241 evident as its ability to generalize to novel compounds that are not in the training set  
 242 (Supplementary Table 1 and 3). This finding suggests that cross-cell prediction is challenging  
 243 because both the cell type and its own status have a significant influence on the transcriptional  
 244 profile. Effective prediction requires considering the basal profiles of the cells. Furthermore, this  
 245 finding highlights the limitations in disregarding the impact of cell type on transcriptional profile  
 246 prediction<sup>10,18</sup>.

247 In summary, TranSiGen stands out among other models in predicting DEGs in chemical-  
 248 blind splitting and random splitting tests. It can also be utilized for inferring DEGs in cell-blind  
 249 splitting. This further highlights the effectiveness of TranSiGen, which is based on self-  
 250 supervised representation learning for transcriptional profiling.



**Fig. 3** The diagram of data splitting and the performance of inferring DEGs in different scenarios.

**a** The diagram of chemical-blind splitting and cell-blind splitting. In scenario 2-1, a dataset of



254 355 compounds on 7 cell lines is split by compounds, ensuring that test compounds do not seen  
255 in the training set. In scenario 2-2, a complete dataset of 8,316 compounds on 164 cell lines is  
256 split by compounds. In scenario 3, the complete dataset of 8,316 compounds on 164 cell lines is  
257 split by cell lines. The model was trained using the profiling data of 10, 50, and 150 cell lines,  
258 and the prediction performance was evaluated on 7 new cell lines. **b** Model performance  
259 comparison in chemical-blind splitting. **c** The performance of TranSiGen in chemical-blind  
260 splitting by using the full dataset. **d** The performance of TranSiGen in cell-blind splitting by  
261 using different numbers of cell lines in the training set. Statistical t-test was applied between the  
262 models. (Note: \*\*\*\*,  $p < 0.0001$ ; \*\*\*,  $0.0001 < p \leq 0.001$ ; \*\*,  $0.001 < p \leq 0.01$ ; \*,  $0.01 < p \leq$   
263  $0.05$  and ns,  $0.05 < p \leq 1.0$ )

## 264 **Ligand-based virtual screening with TranSiGen-derived representation**

265 Compounds with shared mechanisms cause similar changes in cellular signaling, resulting in  
266 similar gene expression profiles<sup>1,7,19</sup>. As a simulation of chemical-induced transcriptional  
267 profiling, the TranSiGen-derived representation shows higher PCC among active compounds  
268 with the same target compared to the PCC between active and inactive compounds  
269 (Supplementary Fig. 4). In this section, the TranSiGen-derived representation was used to assess  
270 whether a compound is active against a specific target. Specifically, DEGs inferred by TranSiGen  
271 and other baseline models were used as the representations of compounds. Then, RF models  
272 based on these representations were constructed to screen active compounds for a target of  
273 interest.

274 Fig. 4a shows the performance of screening 5-hydroxytryptamine receptor 2A (HTR2A)  
275 active compounds on seven cell lines. The model based on TranSiGen-derived representation



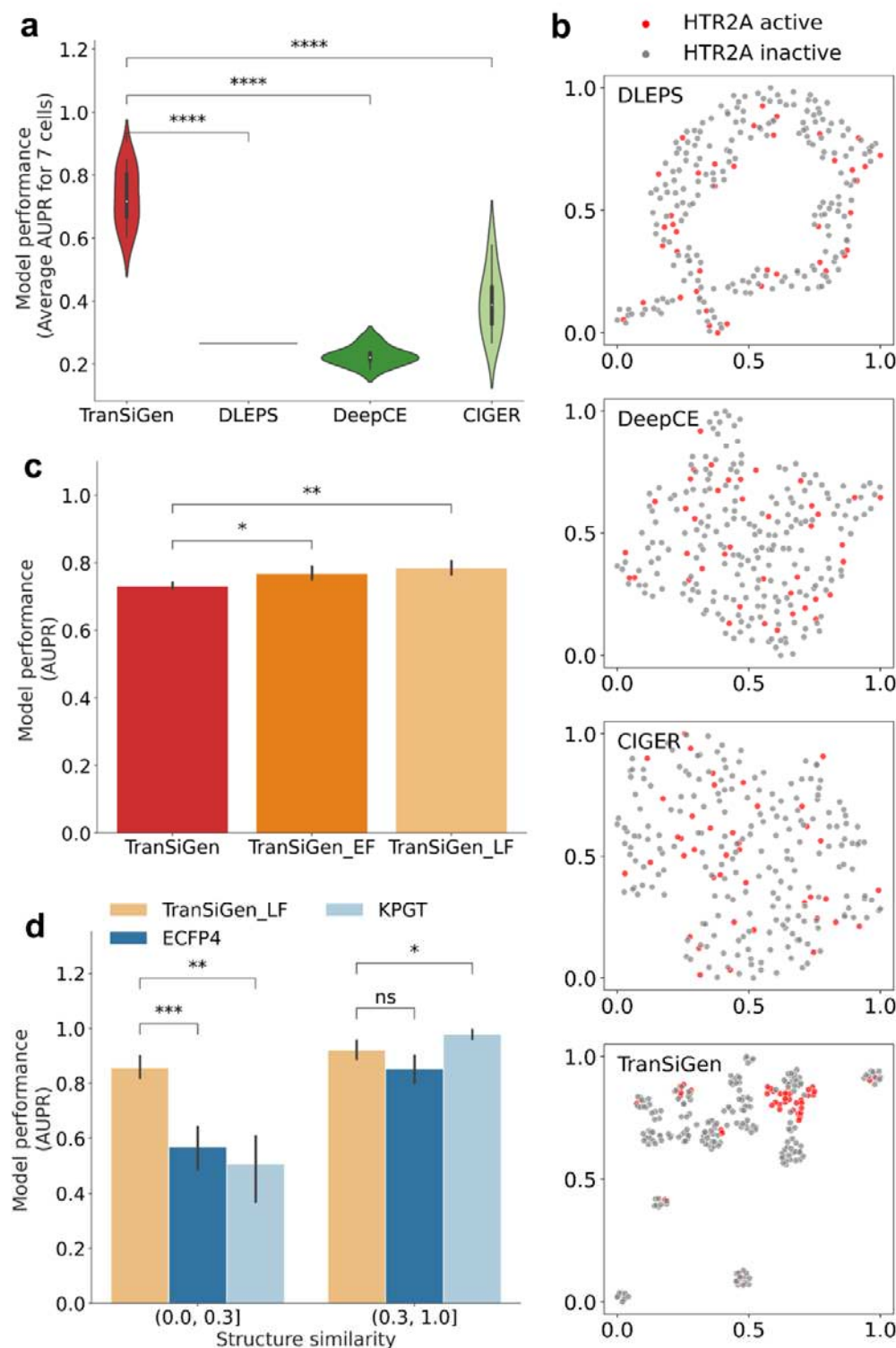
276 outperforms other perturbational representations by a significant margin. This result is further  
277 supported by the dimensionality reduction distribution of active/inactive compounds, where  
278 TranSiGen-derived representations clearly distinguish between the two, while other  
279 perturbational representations exhibit overlapped distributions (Fig. 4b and Supplementary Fig.  
280 5).

281 TranSiGen can be used to obtain the characteristics of each compound in different cellular  
282 backgrounds. In this study, TranSiGen-derived representations from seven cell lines were fused to  
283 evaluate their impact on compound screening performance. Early fusion involves concatenating  
284 TranSiGen-derived representations from seven cells into one single feature, while late fusion  
285 merges the prediction results from seven cells. These two models are denoted as TranSiGen\_EF  
286 and TranSiGen\_LF for early and late fusion, respectively. It was observed that fusing TranSiGen-  
287 derived representations from different cell lines further enhances the screening performance of  
288 active compounds compared to individual cells alone (Fig. 4c). However, the performance  
289 improvement of TranSiGen\_EF is not as significant as that of TranSiGen\_LF, possibly due to the  
290 “curse of dimensionality”<sup>20</sup>. High-dimensional input features in TranSiGen\_EF make it difficult  
291 to learn meaningful patterns. Similar phenomena are also observed in ligand-based virtual  
292 screening for other four targets evaluated (Supplementary Fig. 6).

293 As a molecular representation method, the TranSiGen-derived representation was compared  
294 to other molecular structural representations such as molecular fingerprint ECFP4 and the pre-  
295 trained representation KGPT. The maximum Tanimoto similarities of test molecules relative to  
296 training molecules were calculated using ECFP4. The performance of screening active  
297 compounds was evaluated at different maximal similarity thresholds. For compounds that are

298 dissimilar to the training set (chemical structure similarity  $\in (0.0, 0.3]$ ), the TranSiGen-based  
 299 model demonstrates better predictive ability than structure representation-based model (Fig. 4d  
 300 and Supplementary Fig. 6). This suggests that using transcriptional profiling, such as TranSiGen-  
 301 derived representation, may have advantages in screening for new scaffold compounds that differ  
 302 from known compound structures.

303 Therefore, TranSiGen-derived representation can be used as a new form of molecular  
 304 representation for describing the characteristics of compounds from various cell contexts. It can  
 305 also complement the structure-based representation and offer advantages in ligand-based virtual  
 306 screening.



307

308 **Fig. 4** Model performance of ligand-based virtual screening on target HTR2A. **a** Performance of

309 active compound prediction using different perturbational representations. **b** Dimensionality

310 reduction visualization of HTR2A active and inactive compounds based on various inferred  
 311 perturbational representations. **c** Performance of active compound prediction by applying early  
 312 fusion and late fusion for TranSiGen-derived representation from seven different cell lines. **d**  
 313 Performance of active compounds prediction within different thresholds of max similarity of test  
 314 molecules relative to train data. Statistical t-test was applied between the models. (Note: \*\*\*\*,  $p$   
 315  $< 0.0001$ ; \*\*\*,  $0.0001 < p \leq 0.001$ ; \*\*,  $0.001 < p \leq 0.01$ ; \*,  $0.01 < p \leq 0.05$  and ns,  $0.05 < p \leq 1.0$ )

### 316 **Drug response prediction with TranSiGen-derived representation**

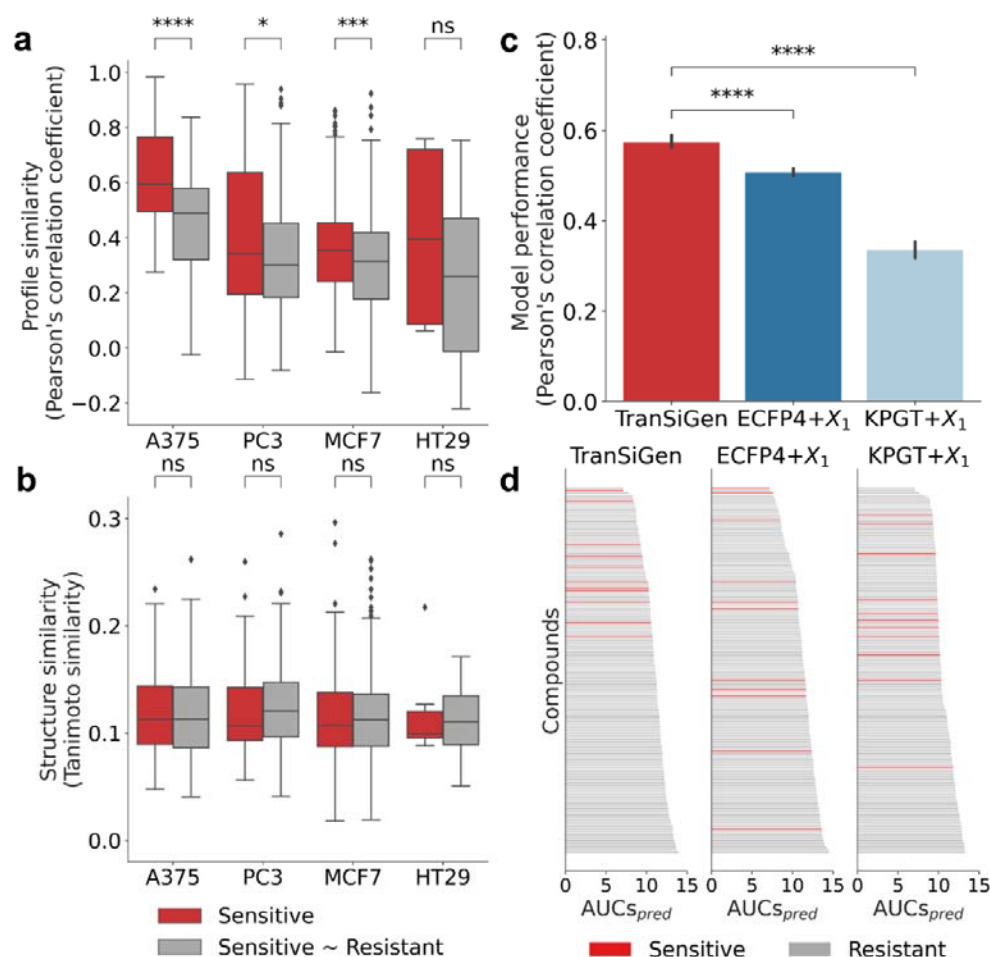
317 Chemical-induced transcriptional profiles directly associate molecular features with the  
 318 cellular effect of a particular drug. This association is beneficial for characterizing drug response  
 319 in different cells<sup>21–23</sup>. Here, we applied the TranSiGen-derived representation to predict the area  
 320 under the dose-response curve (AUC) of a compound on a specific cell line. The AUCs were  
 321 obtained from the cancer treatment response portal (CTRP)<sup>24,25</sup>. We defined compounds with  
 322 AUCs  $\geq 5.5$  as resistant to cell lines, while those with AUCs  $< 5.5$  were considered sensitive<sup>24</sup>.  
 323 More details about the dataset can be found in Supplementary Table 4.

324 To determine whether compounds can be classified as sensitive or resistant to a specific cell  
 325 line based on TranSiGen-derived representation, we first assessed the profiling similarity among  
 326 different compounds. In Fig. 5a, we calculated the PCC within a group of sensitive compounds  
 327 (denoted as Sensitive), as well as the PCC between sensitive and resistant compounds (denoted as  
 328 Sensitive~Resistant). Additionally, we compared the structural similarities of the two groups  
 329 using Tanimoto similarity based on the molecular fingerprint ECFP4 (Fig. 5b). The results  
 330 indicate that the Tanimoto similarities within Sensitive group and the Tanimoto similarities within  
 331 Sensitive~Resistant group are not significantly different on each cell line, suggesting that the

332 structural representation ECFP4 cannot distinguish sensitive and resistant compounds (Fig. 5b).  
 333 In contrast, we observed that the profiling similarities of Sensitive are significantly higher than  
 334 those of Sensitive~Resistant on most cell lines (Fig. 5a). This finding demonstrates the effective  
 335 discrimination between sensitive and resistant compounds achieved through TranSiGen-derived  
 336 representation.

337 Furthermore, the TranSiGen-derived representation was used for drug response prediction in  
 338 downstream task using a RF model. Its performance was compared with RF models based on  
 339 other alternative representations, including perturbational representations generated by baseline  
 340 models (DLEPS, DeepCE, and CIGER), as well as representations combining molecular  
 341 structures and cell information (ECFP4+  $X_1$  and KPGT+  $X_1$ ). As shown in Fig. 5c and  
 342 Supplementary Fig. 7a, the TranSiGen-based model demonstrates significantly better  
 343 performance than other models. Additionally, to evaluate the screening performance, compounds  
 344 were ranked by their predicted AUCs ( $AUC_{pred}$ ), and classified as sensitive or resistant  
 345 according to their true AUCs. The results showed that the TranSiGen-based model predicted  
 346 sensitive compounds with smaller  $AUC_{pred}$  and higher rankings, while other models ranked the  
 347 sensitive compounds randomly (Fig. 5d and Supplementary Fig. 7b). This indicates that the  
 348 TranSiGen-based model has superior screening ability for sensitive compounds.

349 In summary, the TranSiGen-derived representation, simulating DEGs of compounds on cell  
 350 lines, exhibits a distinguishable feature for sensitive and resistant compounds and demonstrates  
 351 outstanding performance on drug response prediction.



352

353 **Fig. 5** Model performance of drug response prediction. **a** The Pearson's correlation coefficients

354 within a group of sensitive compounds and the Pearson's correlation coefficients between

355 sensitive and resistant compounds based on TranSiGen-derived representation. **b** The Tanimoto

356 similarity within a group of sensitive compounds and the similarity between sensitive and

357 resistant compounds based on molecular fingerprint ECFP4, and the Mann-Whitney test was

358 used to analyze the data. **c** Performance of predicting drug response using various type of

359 representations, and statistical t-test was applied between the models. **d** Ranking results of

360 compounds by  $AUCs_{pred}$  of models based on various type of representations. (Note: \*\*\*\*,  $p <$

361 0.0001; \*\*\*,  $0.0001 < p \leq 0.001$ ; \*\*,  $0.001 < p \leq 0.01$ ; \*,  $0.01 < p \leq 0.05$  and ns,  $0.05 < p \leq 1.0$ )

## 362 **Phenotype-based drug repurposing for the treatment of pancreatic cancer**

363 Associating chemical-induced transcriptional profiles with diseases can help identify  
364 potential compounds for treating specific diseases<sup>2,7</sup>. TranSiGen-derived transcriptional profiles  
365 can be used alongside the profiles derived from chemical-treated and -untreated disease states to  
366 screen candidate compounds for disease treatment.

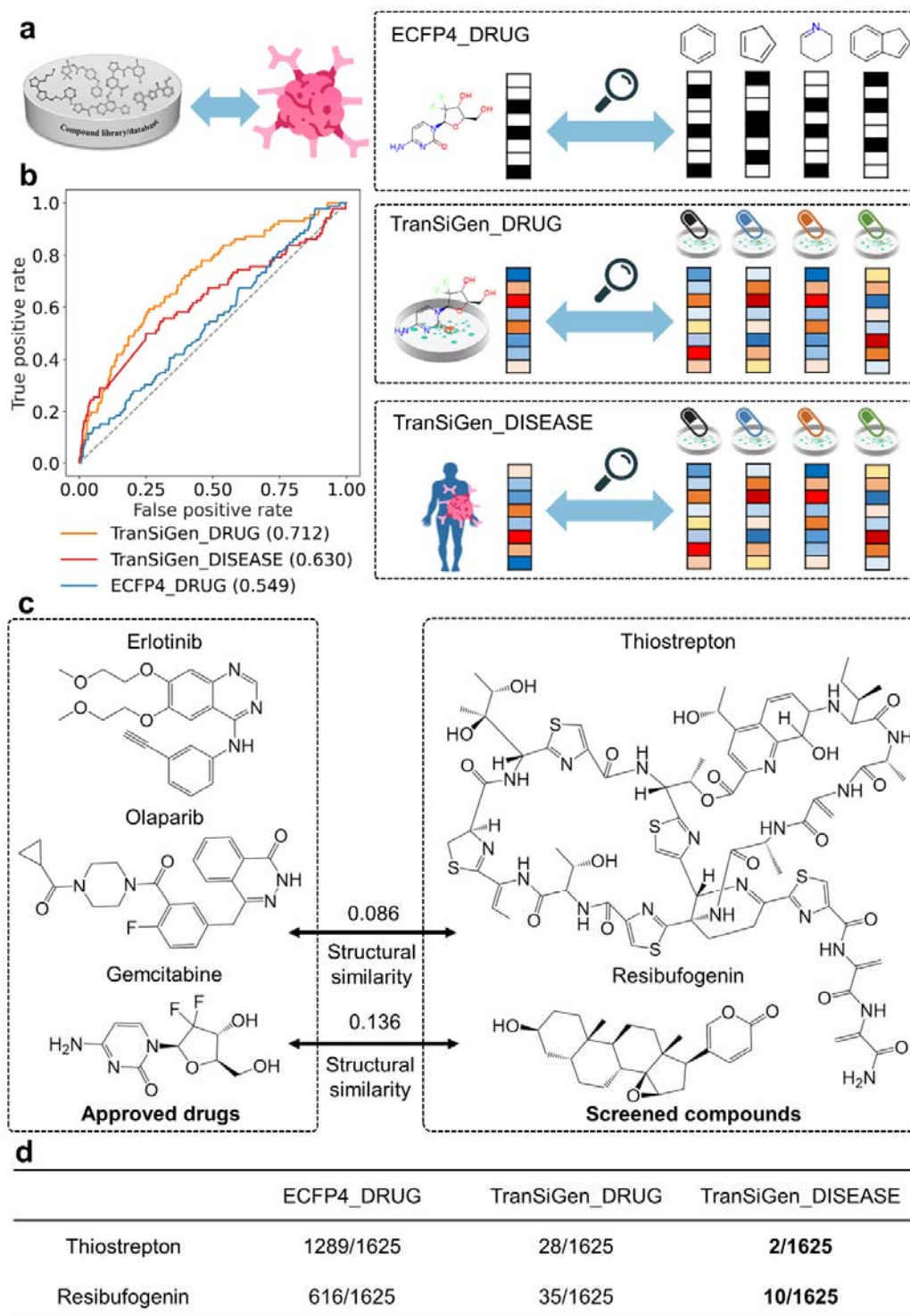
367 In this study, we integrated TranSiGen into a phenotype-based drug repurposing pipeline for  
368 pancreatic cancer<sup>26</sup> to assess its ability to prioritize sensitive compounds for the YAPC pancreatic  
369 cancer cell line from a pool of 1,625 compounds in the PRISM Repurposing dataset<sup>27</sup>. We used  
370 two phenotype-based strategies and compared them to a conventional structural similarity-based  
371 protocol (Fig. 6a). TranSiGen\_DRUG used the real DEGs of approved pancreatic cancer drugs to  
372 identify compounds with similar perturbation effects. Conversely, TranSiGen\_DISEASE looked  
373 for compounds that can reverse the DEGs of pancreatic cancer. Both strategies used connectivity  
374 scores<sup>28</sup> to measure the relationship between DEGs. For comparison, ECFP4\_DRUG was  
375 implemented to find compounds structurally similar to the approved drugs using ECFP4-based  
376 Tanimoto similarity. For more information, please refer to the section “Phenotype-based drug  
377 repurposing for pancreatic cancer” in the Methods section.

378 The screening performance of the three methods is shown in Fig. 6b. ECFP4\_DRUG  
379 yielded the worst predictive classification performance, followed by TranSiGen\_DISEASE, and  
380 the best was TranSiGen\_Drug. Notably, the TranSiGen\_DISEASE approach doesn’t require any  
381 chemical-treated profiles, simulating scenarios where diseases lack known therapeutic drugs.  
382 This is a challenge not addressed by the structural similarity-based strategy. Even without using  
383 perturbed profiles of known drugs, TranSiGen\_DISEASE effectively enriched hits among the

384 top-ranking compounds (Supplementary Table 5). Phenotype-based strategies can identify  
 385 compounds less similar to the approved drugs than those screened by ECFP4\_DRUG (Fig. 6c  
 386 and Supplementary Fig. 8). For instance, nature products thiostrepton and resibufogenin (Fig. 6c)  
 387 ranked in the top 10 without sensitive annotations in PRISM dataset (Supplementary Table 6).  
 388 Their abilities to inhibit pancreatic cancer cells have been confirmed by a literature survey<sup>29,30</sup>.  
 389 Thiostrepton, a natural cyclic oligopeptide, reduces the viability and clonogenicity of pancreatic  
 390 cancer cell lines and induces ferroptosis via STAT3/GPX4 signalling<sup>30</sup>. Resibufogenin, a steroid  
 391 lactone from the skin venom gland of toads, demonstrated potent anti-pancreatic cancer effects in  
 392 vivo and in vitro, and can induce caspase-dependent apoptosis<sup>29</sup>. Fig. 6d summarizes the rankings  
 393 of these two compounds in different screening strategies. Both phenotype-based strategies,  
 394 TranSiGen\_DISEASE and TranSiGen\_DRUG, consistently prioritized them. In contrast, the  
 395 structure-based strategy ECFP4\_DRUG failed to effectively prioritize these nature products,  
 396 ranking them at 1,289 and 616, respectively. This may be attributed to the large structural  
 397 differences between them and approved drugs, highlighting the inherent limitation of a structure-  
 398 based strategy.

399 These results highlight the effectiveness of the phenotype-based strategies that use  
 400 TranSiGen-derived representation in identifying potent candidate compounds, including those  
 401 with unique structures. Overall, TranSiGen expands the range of compounds that can be screened  
 402 with predicted transcriptional perturbational profiles. It can be easily integrated into a phenotype-  
 403 based drug repurposing pipeline, improving drug discovery efficiency and minimizing costs.





404

405 **Fig. 6** Phenotype-based drug repurposing for the treatment of pancreatic cancer. **a** The flow chart

406 of drug repurposing strategy. **b** The screening performance of phenotype-based strategy and

407 structural similarity-based strategy. **c** TranSiGen\_DISEASE screened compounds that are capable  
 408 of inhibiting pancreatic cancer cells, and their max structural similarities to approved drugs. **d**  
 409 The rankings of thiostrepton and resibufogenin from different screening strategies.

## 410 **Discussion**

411 Given the widespread application of perturbational gene expression profiling in biomedical  
 412 research and the limitations arising from random noise and systemic biases in profiling data, we  
 413 propose TranSiGen, a VAE-based model. TranSiGen employs a self-supervised representation  
 414 learning strategy to overcome these inherent limitations of data and infer novel chemical-induced  
 415 gene expression profiles. TranSiGen excels in denoising and reconstructing the transcriptional  
 416 profiling, as well as effectively learning cellular and compound features from data. It outperforms  
 417 existing models in both random splitting and chemical-blind splitting, and also shows potential  
 418 for inferring DEGs in new cell lines, a capability absent in existing models. Moreover,  
 419 TranSiGen-derived profiles have served as a unified and standardized representation of  
 420 phenotypic information, demonstrating its effectiveness in various downstream tasks, including  
 421 ligand-based virtual screening, drug response prediction, and phenotype-based drug repurposing.

422 Therefore, TranSiGen efficiently learns the representation of transcriptional data and can  
 423 predict novel chemical-induced transcriptional changes in a given cell line. We believe that  
 424 integrating it into the pipeline of drug discovery and mechanism of action research can enhance  
 425 the efficiency and reduce the costs, thereby promoting the development of biomedicine.

426

## 427 **Methods**

### 428 **Transcriptional data processing**

LINCS<sup>2</sup> has made publicly available resources on high-throughput gene expression profiles of different perturbagens, such as small molecule compounds and shRNAs. Using the L1000 assay, it is possible to measure the expression values of only 978 landmark genes, and while still recovering most of the full transcriptome. The latest CMAP LINCS 2020 dataset was used in this study<sup>15</sup>.

In LINCS, there are 5 levels of data. For this study, we utilized level 3 data, which included both raw perturbation profiles and control profiles (denoted as  $X_1$  and  $X_2$ , respectively). To filter the profiles, we used the most common condition with a duration of 24 h and a dosage concentration of 10 M for perturbed expression profiles by compounds. Additionally, we matched the expression profile with DMSO vehicle to the perturbed profiles on the same plate to create paired profiles  $X_1 \sim X_2$ , minimizing batch effects between cases and controls. The extracted dataset contained 219,650  $X_1 \sim X_2$  pairs for 8,316 compounds on 164 cell lines. Furthermore, MODZ was applied to ensure that only one  $X_1 \sim X_2$  pair per compound was included for each cell line with multiple  $X_1 \sim X_2$  pairs. The processed dataset contained the transcriptional profiles of 8,316 compounds on 164 cell lines, including 78,569  $X_1 \sim X_2$  pairs.

The gene expression profiles induced by shRNA were processed using the same method described above. Profiles from the 10 most common cell lines (A375, A549, ASC, HA1E, HCC515, HT29, MCF7, NPC, PC3, and VCAP) measured after 24 h were selected. The control profile with an empty vector in the same plate was then paired with the perturbed profiles. The final dataset contained 188,509  $X_1 \sim X_2$  pairs consisting of 4,112 shRNAs on 10 cell lines, which were used to initialize two VAEs in TranSiGen.

## Compound representations

451 Considering that the current number of compounds with experimentally measured gene  
452 expression profiles is still limited compared to the vast chemical space, TranSiGen utilized the  
453 pre-trained molecular representation KPGT<sup>17</sup> for compounds. KPGT is a novel self-supervised  
454 learning framework for molecular graph representation. It leveraged a knowledge-guided pre-  
455 training strategy to capture rich structural and semantic information from large-scale unlabeled  
456 molecular graphs. In this study, the 2034-dimensional representation obtained from the KPGT  
457 pre-trained model was used as the molecular input for TranSiGen.

458 Alternatively, chemical fingerprints, widely used as a form of molecular representation in  
459 machine learning, was also used as molecular input for TranSiGen. There are represented as  
460 binary vectors indicating the presence or absence of particular substructures in compounds.  
461 Specifically, the molecular fingerprint ECFP4<sup>31</sup> with a radius of 2 and a length of 2048 was used  
462 here.

### 463 **TranSiGen architecture**

464 The VAE<sup>32</sup> is a deep generative model consisting of an encoder and a decoder. VAE is  
465 capable of learning an efficient and meaningful latent space from high-dimensional data by  
466 compressing and reconstructing the original input. Unlike the standard autoencoder, which maps  
467 the input to a point in the latent space and trains by minimizing the reconstruction error, VAE  
468 encodes the input to a distribution. This requires the addition of a Kullback-Leibler (KL)  
469 divergence term to the reconstruction loss, which constrains the latent vectors to match a  
470 Gaussian distribution.

471 The architecture of TranSiGen consists of two VAEs: one for encoding the basal profiles  $X_1$   
472 and the other for encoding the perturbation profiles  $X_2$ . TranSiGen minimizes the loss of learning

the representations of  $X_1$  and  $X_2$ . Additionally, a linear function is used to map from the latent representation  $Z_1$  of  $X_1$  and the hidden representation  $Z_{mol}$  of the compound representation  $C_{mol}$  to the perturbed latent representation  $Z_2 F Z_1$  of  $X_2'$ , mimicking the chemical-induced transcription changes. The layer and dimension details of TranSiGen are shown in Supplementary Fig. 1. During the training process, TranSiGen also minimizes the loss of predicting the differential expression genes  $\Delta X'$ . This involves minimizing the reconstruction loss between  $X_2' - X_1$  and  $X_2 - X_1$ , as well as constraining the predicted perturbed latent representation  $Z_2 F Z_1$  match to the latent representation  $Z_2$ . The loss function of TranSiGen is defined as follow:

$$\text{Loss} = \text{MSE}(X_1, \hat{X}_1) + \text{MSE}(X_2, \hat{X}_2) + \text{MSE}(\Delta X, \Delta X') + \text{KL}(q(Z_1|X_1)||p(Z_1)) + \text{KL}(q(Z_2|X_2)||p(Z_2)) + \text{KL}(q(Z_2|X_1, C_{mol})||q(Z_2|X_2)) \quad (1)$$

$$\text{KL}(q(Z_1|X_1)||p(Z_1)) = -\frac{1}{2}(1 + \log \sigma_1^2 - \sigma_1^2 - \mu_1^2) \quad (2)$$

$$\text{KL}(q(Z_2|X_2)||p(Z_2)) = -\frac{1}{2}(1 + \log \sigma_2^2 - \sigma_2^2 - \mu_2^2) \quad (3)$$

$$\text{KL}(q(Z_2|X_1, C_{mol})||q(Z_2|X_2)) = -\frac{1}{2}(1 + \log \frac{\sigma_2'^2}{\sigma_2^2} - \frac{\sigma_2'^2 + (\mu_2' - \mu_2)^2}{\sigma_2^2}) \quad (4)$$

where  $\mu_1$  and  $\sigma_1^2$  represent the mean and variance for  $q(Z_1|X_1)$ ,  $\mu_2$  and  $\sigma_2^2$  represent the mean and variance for  $q(Z_2|X_2)$ ,  $\mu_2'$  and  $\sigma_2'^2$  represent the mean and variance for  $q(Z_2|X_1, C_{mol})$ .

## Performance evaluation metrics

As shown in Table 2, the model's prediction performance was evaluated using following metrics: Root mean squared error (RMSE), Pearson's correlation coefficient, and Precision@K. RMSE and Pearson coefficient were used to measure the prediction performance on the overall landmark genes. Precision@k, on the other hand, focused on the most significantly up- and down-regulated expressed genes. In this study, Positive Precision@100 was evaluated for the top 100 up-regulated genes, while Negative Precision@100 was evaluated for the top 100 down-

496 regulated genes.

# 497 **Ligand-based virtual screening**

498 For compounds with transcriptional profiles, the target annotations were collected from the  
 499 compound information file provided by LINCS2020<sup>15</sup> and PubChem<sup>33</sup>  
 500 (<https://pubchem.ncbi.nlm.nih.gov>). Specifically, the targets of the compound were determined  
 501 according to IC<sub>50</sub>, K<sub>i</sub>, K<sub>d</sub> less than 10<sup>-6</sup> M in PubChem, and the target information from these two  
 502 sources was deduplicated and integrated.

503 In the task of ligand-based virtual screening, compounds with target annotations were  
 504 selected from 7 cell lines (A375, HELA, PC3, MCF7, HT29, YAPC, and HA1E). Among these,  
 505 the target HTR2A had the most active compounds (Supplementary Fig. 3), and was selected for  
 506 the task. Inactive compounds were sampled from the remaining compounds at a ratio of 1:5 to  
 507 construct a dataset for the active compound screening. The compounds from each cell line were  
 508 randomly split into training and test set with a ratio of 4:1.

509 Given the limited availability of dataset for active compound screening, we used RF for  
 510 active compound prediction. To construct the RF classifiers, we used two types of features:  
 511 inferred perturbational representations (TranSiGen, DLEPS, DeepCE, and CIGER) and structural  
 512 representations (molecular fingerprint ECFP4, and pretrained representation KPGT). We  
 513 conducted a hyperparameter search for “n\_estimators”, “max\_depth”, “criterion” and “obb\_score”  
 514 to find the optimal model. To evaluate the model performance, we used the area under the  
 515 Precision–Recall curve (AUPR). The training-evaluation procedure was repeated five times with  
 516 different random seeds to determine the model performance. These processes were implemented  
 517 using scikit-learn<sup>34</sup>.

## 518 **Drug response prediction**

519       The CTRP<sup>24,25</sup> is a widely used cancer cell response dataset that associate genetic, lineage,  
520 and other cellular molecular characteristics of cancer cell lines with drug sensitivity. It  
521 quantitatively profiles the sensitivity of cancer cell lines to small molecules. The AUC label is a  
522 dose-independent measure of compound sensitivity. Smaller AUCs indicate greater sensitivity of  
523 cells to the drugs. A subset of the drug response dataset for 267 compounds on four cell lines  
524 (A375, PC3, MCF7, and HT29) was obtained from CTRP, and the details of the dataset are  
525 shown in Supplementary Table 4. The processed dataset was split into training and test set at 4:1  
526 ratio by compounds.

527       Similarly, RF regression models were used to predict drug response. Inferred perturbational  
528 representations (TranSiGen, DLEPS, DeepCE, and CIGER) and representations combining  
529 molecular structures and cell information (ECFP4+  $X_1$  and KPGT+  $X_1$ ) were used. Four  
530 hyperparameters, including “n\_estimators”, “max\_depth”, “criterion” and “obb\_score”, were  
531 considered to obtain the optimal model. The model’s performance was evaluated by the Pearson’s  
532 correlation coefficient. The model’s performance was assessed by repeating training-evaluation  
533 procedure five times with different random seeds.

## 534 **Phenotype-based drug repurposing for pancreatic cancer**

### 535 *Differential gene expression profiles of approved drugs*

536       The approved drugs for pancreatic cancer were downloaded from  
537 <https://www.cancer.gov/about-cancer/treatment/drugs/pancreatic>. Among them, the DEGs of  
538 erlotinib, olaparib and gemcitabine were obtained from LINCS 2020 dataset<sup>15</sup>. These profiles  
539 were used for subsequent phenotype-based drug repurposing for pancreatic cancer.

## 540 *Differential gene expression profile of disease*

541 The pancreatic adenocarcinoma cohort of the The Cancer Genome Atlas (TCGA)<sup>35</sup> was  
 542 downloaded from UCSC Xena (<https://xenabrowser.net/>). This cohort includes RNA-seq  
 543 expression data of tumor samples and normal samples. The DESeq2<sup>36</sup> method was used to  
 544 analyze the differential gene expression for pancreatic cancer. DEGs for pancreatic cancer were  
 545 selected based on the following criteria:  $|\log_2\text{Foldchange}| \geq 1.5$ ,  $p\text{-value} \leq 0.05$  and false  
 546 discovery rate  $< 0.25$ . A total of 293 up-regulated genes and 168 down-regulated genes were  
 547 identified.

## 548 *Inferring perturbation gene expression profiles of compounds*

549 This study utilized the PRISM Repurposing dataset<sup>27</sup>, which includes primary and  
 550 secondary screening datasets, for phenotype-based repurposing for pancreatic cancer. The  
 551 compounds from PRISM secondary screen were evaluated based on their AUC values, which  
 552 indicate compound sensitivities on cells and serve as labels for screening performance  
 553 assessment.

554 The dataset was downloaded from <https://depmap.org/repurposing/>. Compounds with  
 555 ground-truth expression profiles in the LINCS 2020 dataset were excluded, resulting a dataset  
 556 contains 1,625 compounds. TranSiGen inferred the DEGs  $\Delta X'$  of 978 landmark genes associated  
 557 with these compounds in YAPC pancreatic cancer cell. Additionally, and the expression values of  
 558 9,196 best inferred genes were inferred from the generated 978 landmark genes to obtain the  
 559 predicted expression values of 10,174 genes. The inference weight matrix was obtained from the  
 560 L1000 project<sup>2</sup>.

## 561 **Connectivity score**



562 The connectivity score, obtained from the gene set enrichment analysis<sup>28</sup>, is used to measure  
563 the relationship between transcriptional profiles. The connectivity score ranges from -1 to 1,  
564 where -1 indicates a complete reversal of the query profile to the reference profile, while 1  
565 indicates a complete similarity of the query and the reference profile.

566 Firstly, the enrichment score (ES) is used to evaluate the enrichment of a predefined gene set  
567 at the top or bottom of the reference differential gene list. The enrichment scores for up-regulated  
568 gene set and down-regulated gene set are denoted as  $a$  and  $b$ , respectively:

$$569 \quad a = \max_{j=1 \text{ to } t} \left[ \frac{j}{t} - \frac{V(j)}{n} \right] \quad (5)$$

$$570 \quad b = \max_{j=1 \text{ to } t} \left[ \frac{V(j)}{n} - \frac{(j-1)}{t} \right] \quad (6)$$

$$571 \quad ES = \begin{cases} a, & \text{if } a > b \\ -b, & \text{if } b > a \end{cases} \quad (7)$$

572 where  $n$  represents the number of genes in the expression profiles,  $t$  represents the number of  
573 genes in the predefined gene set, and  $V(j)$  represents the rank of a specific gene in the rank list.

574 Next, the above equations are used to calculate the enrichment scores of the predefined up-  
575 regulated and down-regulated genes by the query profile, resulting in the values  $ES_{up}$  and  
576  $ES_{down}$ . Finally, considering these two enrichment scores together, the connectivity score of the  
577 query profile relative to the reference profile is calculated as follows:

$$578 \quad \text{Connectivity score} = \begin{cases} ES_{up} - ES_{down}, & \text{if } \text{sign}(ES_{up}) \neq \text{sign}(ES_{down}) \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

579 Specifically, the reference profiles consist of the inferred DEGs of all compounds from  
580 TranSiGen. The DEGs of approved drugs for pancreatic cancer were used to identify compounds  
581 with positive connectivity scores, while the DEGs of pancreatic cancer were employed to identify  
582 compounds with negative connectivity scores. In the screening dataset, the top 20% compounds  
583 having the lowest AUCs on the YAPC pancreatic cancer cell line were identified as hits, and the

584 area under the ROC curve was used to evaluate the screening performance.

585

## 586 **Data availability**

587 The expanded CMap LINCS Resource 2020 is available at

588 <https://clue.io/data/CMap2020#LINCS2020>. The PRISM Repurposing dataset is available at

589 <https://depmap.org/repurposing/>, and the pancreatic adenocarcinoma cohort of the TCGA is

590 available <https://xenabrowser.net/>.

591

## 592 **Code availability**

593 The code for model training and analysis is available at: <https://github.com/myzheng->

594 [SIMM/TranSiGen](#), and have been deposited fully in the Zenodo under accession code

595 <https://zenodo.org/records/10056517>.

596

## 597 **References**

598 1. Lamb, J. *et al.* The connectivity map: using gene-expression signatures to connect small  
599 molecules, genes, and disease. *Science* **313**, 1929–1935 (2006).

600 2. Subramanian, A. *et al.* A next generation connectivity map: L1000 platform and the first  
601 1,000,000 profiles. *Cell* **171**, 1437–1452.e17 (2017).

602 3. Douglass, E. F. *et al.* A community challenge for a pancancer drug mechanism of action  
603 inference from perturbational profile data. *Cell Rep. Med.* **3**, 100492 (2022).

604 4. Lachmann, A. *et al.* Massive mining of publicly available RNA-seq data from human and  
605 mouse. *Nat. Commun.* **9**, 1366 (2018).

- 606 5. Zheng, M. *et al.* ChemPert: mapping between chemical perturbation and transcriptional  
607 response for non-cancer cells. *Nucleic Acids Res.* **51**, D877–D889 (2023).
- 608 6. Chen, B. *et al.* Reversal of cancer gene expression correlates with drug efficacy and reveals  
609 therapeutic targets. *Nat. Commun.* **8**, 16022 (2017).
- 610 7. Wei, G. *et al.* Gene expression-based chemical genomics identifies rapamycin as a modulator  
611 of MCL1 and glucocorticoid resistance. *Cancer Cell* **10**, 331–342 (2006).
- 612 8. Pabon, N. A. *et al.* Predicting protein targets for drug-like compounds using transcriptomics.  
613 *PLOS Comput. Biol.* **14**, e1006651 (2018).
- 614 9. Zhong, F. *et al.* Drug target inference by mining transcriptional data using a novel graph  
615 convolutional network framework. *Protein Cell* **13**, 281–301 (2022).
- 616 10. Zhu, J. *et al.* Prediction of drug efficacy from transcriptional profiles with deep learning. *Nat.*  
617 *Biotechnol.* **39**, 1444–1452 (2021).
- 618 11. Pham, T.-H., Qiu, Y., Zeng, J., Xie, L. & Zhang, P. A deep learning framework for high-  
619 throughput mechanism-driven phenotype compound screening and its application to COVID-19  
620 drug repurposing. *Nat. Mach. Intell.* **3**, 247–257 (2021).
- 621 12. Pham, T.-H. *et al.* Chemical-induced gene expression ranking and its application to  
622 pancreatic cancer drug repurposing. *Patterns* **3**, 100441 (2022).
- 623 13. Qiu, Y. L., Zheng, H. & Gevaert, O. Genomic data imputation with variational auto-encoders.  
624 *GigaScience* **9**, giaa082 (2020).
- 625 14. Way, G. P. & Greene, C. S. Extracting a biologically relevant latent space from cancer  
626 transcriptomes with variational autoencoders. *Pac. Symp. Biocomput.* **23**, 80–91 (2018).
- 627 15. Expanded CMap LINCS Resource 2020. <https://clue.io/data/CMap2020#LINCS2020>

628 (accessed February 20, 2022).

629 16. Szalai, B. *et al.* Signatures of cell death and proliferation in perturbation transcriptomics  
630 data—from confounding factor to effective prediction. *Nucleic Acids Res.* **47**, 10010–10026  
631 (2019).

632 17. Li, H., Zhao, D. & Zeng, J. KPGT: knowledge-guided pre-training of graph transformer for  
633 molecular property prediction. in *Proceedings of the 28th ACM SIGKDD Conference on*  
634 *Knowledge Discovery and Data Mining* 857–867 (2022).

635 18. Lagunin, A., Ivanov, S., Rudik, A., Filimonov, D. & Poroikov, V. DIGEP-Pred: web service  
636 for *in silico* prediction of drug-induced gene expression profiles based on structural formula.  
637 *Bioinformatics* **29**, 2062–2063 (2013).

638 19. Hieronymus, H. *et al.* Gene expression signature-based chemical genomic prediction  
639 identifies a novel class of HSP90 pathway modulators. *Cancer Cell* **10**, 321–330 (2006).

640 20. Bellman, R. Dynamic programming. *Science* **153**, 34–37 (1966).

641 21. Dong, H., Xie, J., Jing, Z. & Ren, D. Variational autoencoder for anti-cancer drug response  
642 prediction. Preprint at <https://doi.org/10.48550/arXiv.2008.09763> (2021).

643 22. Emdadi, A. & Eslahchi, C. DSPLMF: a method for cancer drug sensitivity prediction using a  
644 novel regularization approach in logistic matrix factorization. *Front. Genet.* **11**, 75 (2020).

645 23. Brubaker, D. *et al.* Drug intervention response predictions with PARADIGM (DIRPP)  
646 identifies drug resistant cancer cell lines and pathway mechanisms of resistance. *Pac. Symp.*  
647 *Biocomput. Pac. Symp. Biocomput.* 125–135 (2014).

648 24. Basu, A. *et al.* An interactive resource to identify cancer genetic and lineage dependencies  
649 targeted by small molecules. *Cell* **154**, 1151–1161 (2013).

- 650 25. Seashore-Ludlow, B. *et al.* Harnessing connectivity in a large-scale small-molecule  
651 sensitivity dataset. *Cancer Discov.* **5**, 1210–1223 (2015).
- 652 26. Singh, K., Shishodia, G. & Koul, H. K. Pancreatic cancer: genetics, disease progression,  
653 therapeutic resistance and treatment strategies. *J. Cancer Metastasis Treat.* **7**, 60 (2021).
- 654 27. Corsello, S. M. *et al.* Discovering the anticancer potential of non-oncology drugs by  
655 systematic viability profiling. *Nat. Cancer* **1**, 235–248 (2020).
- 656 28. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for  
657 interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci.* **102**, 15545–15550 (2005).
- 658 29. Liu, L. *et al.* Resibufogenin suppresses transforming growth factor- $\beta$ -activated kinase  
659 1-mediated nuclear factor- $\kappa$ B activity through protein kinase C-dependent inhibition of  
660 glycogen synthase kinase 3. *Cancer Sci.* **109**, 3611–3622 (2018).
- 661 30. Zhang, W. *et al.* Thiostrepton induces ferroptosis in pancreatic cancer cells through  
662 STAT3/GPX4 signalling. *Cell Death Dis.* **13**, 1–12 (2022).
- 663 31. Rogers, D. & Hahn, M. Extended-connectivity fingerprints. *J. Chem. Inf. Model.* **50**, 742–754  
664 (2010).
- 665 32. Kingma, D. P. & Welling, M. Auto-encoding variational bayes. Preprint at  
666 <https://doi.org/10.48550/arXiv.1312.6114> (2022).
- 667 33. Kim, S. *et al.* PubChem in 2021: new data content and improved web interfaces. *Nucleic  
668 Acids Res.* **49**, D1388–D1395 (2021).
- 669 34. Pedregosa, F. *et al.* Scikit-learn: machine learning in python. Preprint at  
670 <https://doi.org/10.48550/arXiv.1201.0490> (2018).
- 671 35. The Cancer Genome Atlas Research Network *et al.* The cancer genome atlas pan-cancer

672 analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).

673 36. Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome*

674 *Biol.* **11**, R106 (2010).

675

## 676 **Acknowledgements**

677 We gratefully acknowledge financial support from National Natural Science Foundation of China

678 (T2225002, 82273855 to M.Y.Z. and 82204278 to X.T.L), SIMM-SHUTCM Traditional Chinese

679 Medicine Innovation Joint Research Program (E2G805H), Shanghai Municipal Science and

680 Technology Major Project, National Key Research, and Development Program of China

681 (2022YFC3400504 to M.Y.Z.), the Youth Innovation Promotion Association CAS (2023296 to

682 S.L.Z).

683

## 684 **Author information**

### 685 **Authors and Affiliations**

686 **Drug Discovery and Design Center, State Key Laboratory of Drug Research, Shanghai**

687 **Institute of Materia Medica, Chinese Academy of Sciences, 555 Zuchongzhi Road, Shanghai,**

688 **201203, China**

689 Xiaochu Tong, Ning Qu, Xiangtai Kong, Shengkun Ni, Kun Wang, Lehan Zhang, Yiming Wen,

690 Sulin Zhang, Xutong Li & Mingyue Zheng

691 **University of Chinese Academy of Sciences, No. 19A Yuquan Road, Beijing, 100049, China**

692 Xiaochu Tong, Ning Qu, Xiangtai Kong, Shengkun Ni, Lehan Zhang, Yiming Wen, Sulin Zhang,

693 Xutong Li & Mingyue Zheng

694 **School of Life Sciences, Division of Life Sciences and Medicine, University of Science and**

695 **Technology of China, Hefei 230026, China**

696 Kun Wang

697 **School of Pharmaceutical Science and Technology, Hangzhou Institute for Advanced Study,**

698 **University of Chinese Academy of Sciences, Hangzhou 310024, China**

699 Yiming Wen

## 700 Contributions

701 M.Z. and X.L. conceived the project and were responsible for the decision to submit the  
702 manuscript; X.T. implemented the TranSiGen model, conducted computational analysis and wrote  
703 the paper; N.Q, X.K., S.N., K.W., L.Z, Y.W. and S.Z. discussed the results and commented on the  
704 manuscript.

## 705 Corresponding authors

706 Correspondence to Xutong Li or Mingyue Zheng.

707

## 708 Ethics declarations

## 709 Competing interests

710 The authors declare that they have no competing interests.

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739



740

741

742

743 **Table 1.** List of symbols and notations used in the paper.

Symbol	Description
$X_1$	The control profiles treated with DMSO
$X_2$	The transcriptional profiles treated with compounds
$\Delta X$	The differential expression genes ( $X_2 - X_1$ )
$\hat{X}_1$	The reconstructed control profiles
$\hat{X}_2$	The reconstructed transcriptional profiles
$\Delta \hat{X}$	The reconstructed differential expression genes ( $\hat{X}_2 - \hat{X}_1$ )
$X_2'$	The predicted transcriptional profiles from $X_1$ and perturbation representation
$\Delta X'$	The predicted differential expression genes ( $X_2' - X_1$ )
$Z_1$	The latent representation of $X_1$
$Z_2$	The latent representation of $X_2$
$Z_2 FZ_1$	The predicted transcriptional profiles from $X_1$ and perturbation representation
$C_{mol}$	The input representation of the compounds
$Z_{mol}$	The hidden representation of the compounds

744

745

746

747

748

749

750

751

752

753

754

755

756

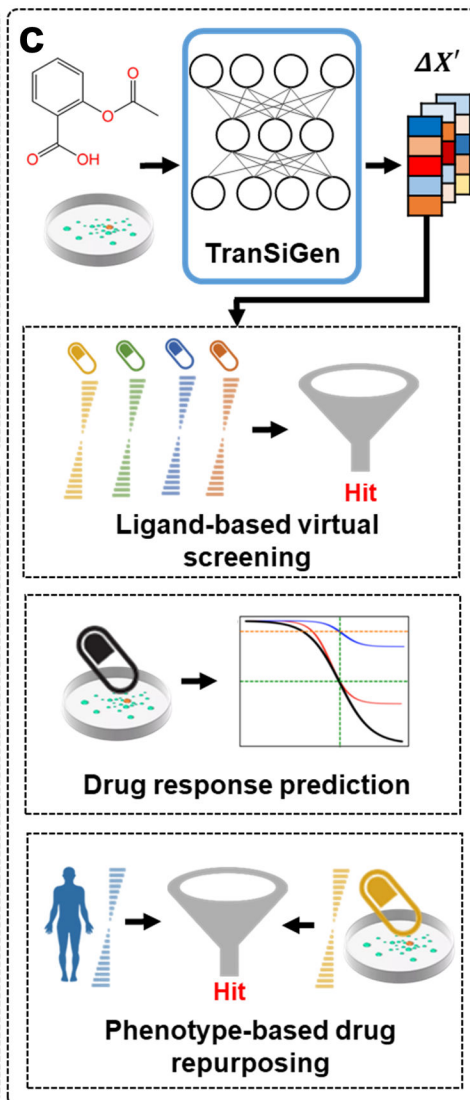
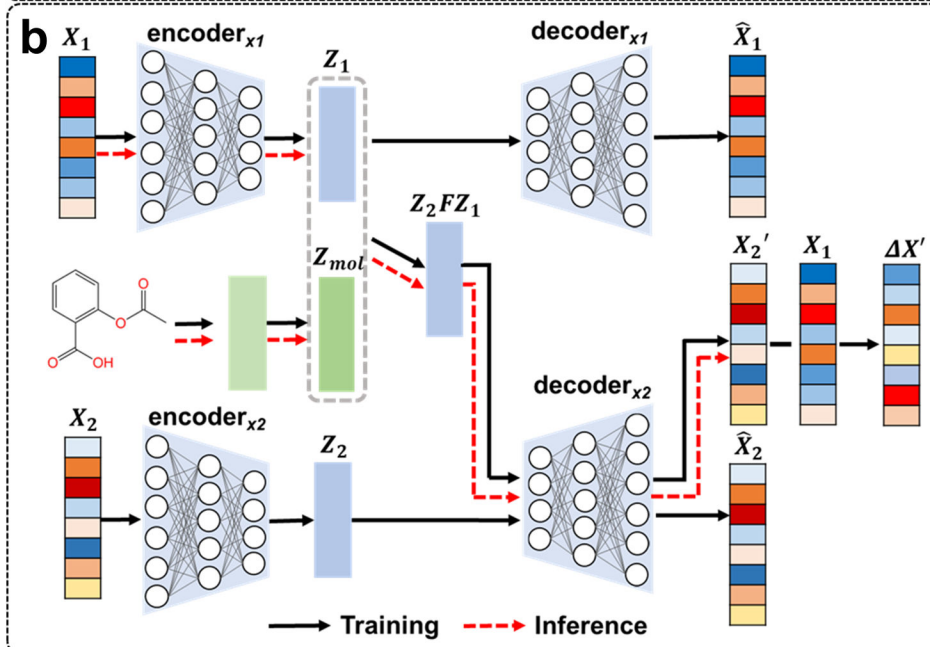
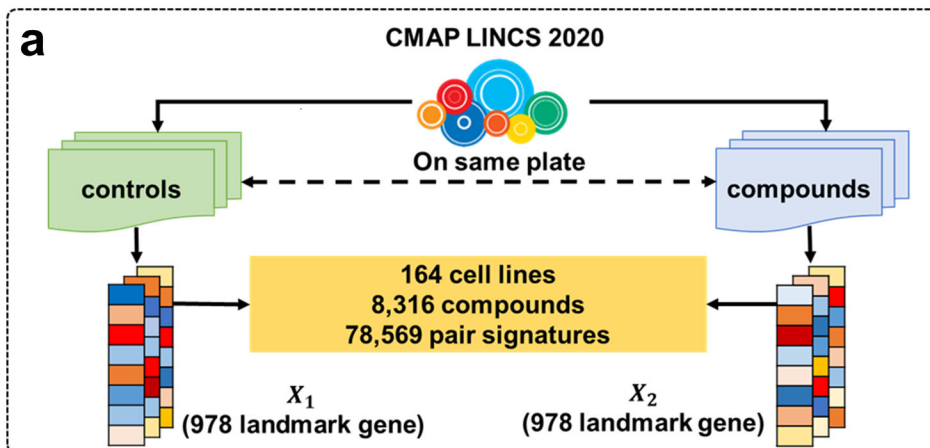
757  
758  
759  
760  
761  
762

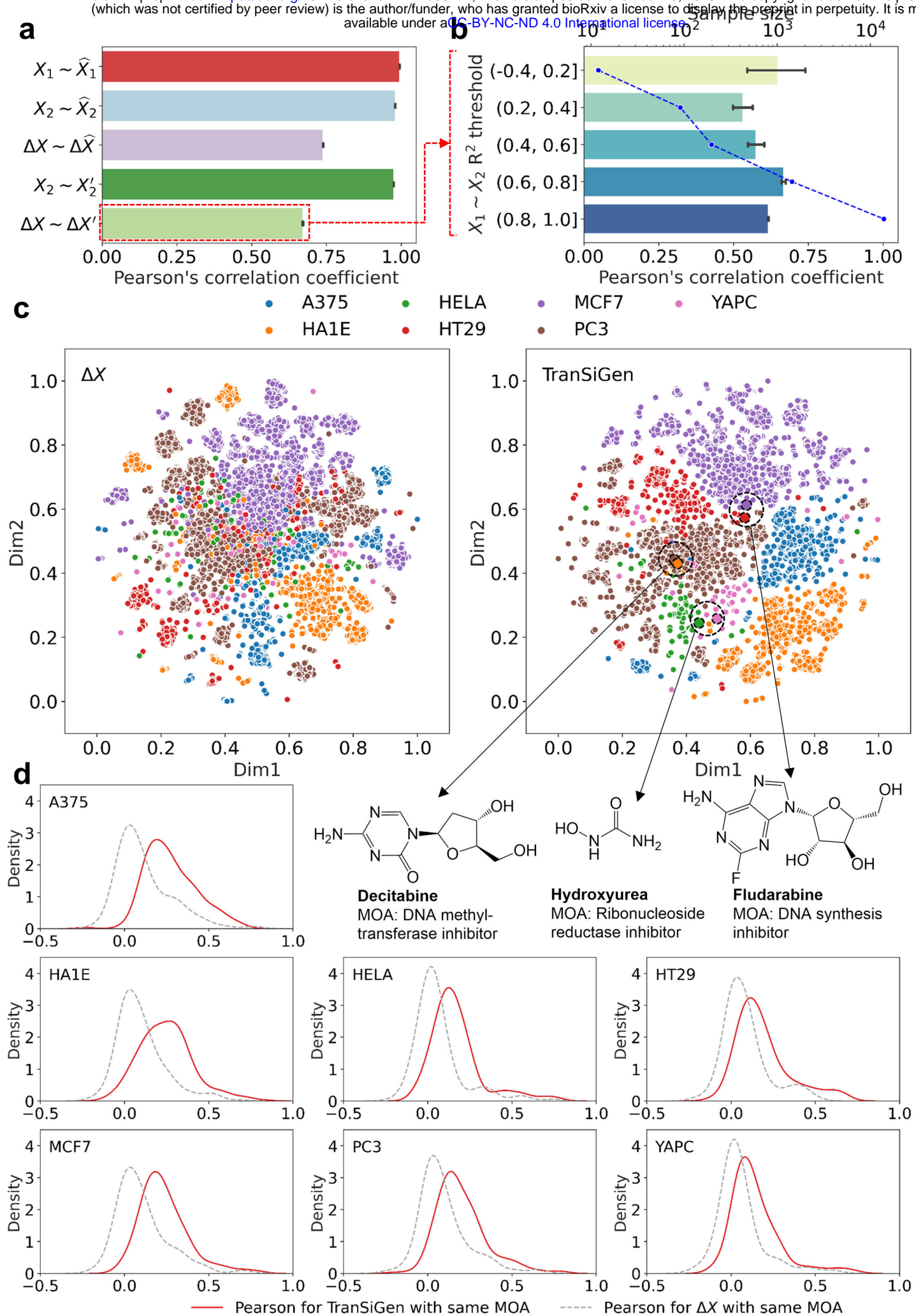
763 **Table 2.** Description of the evaluation metrics.

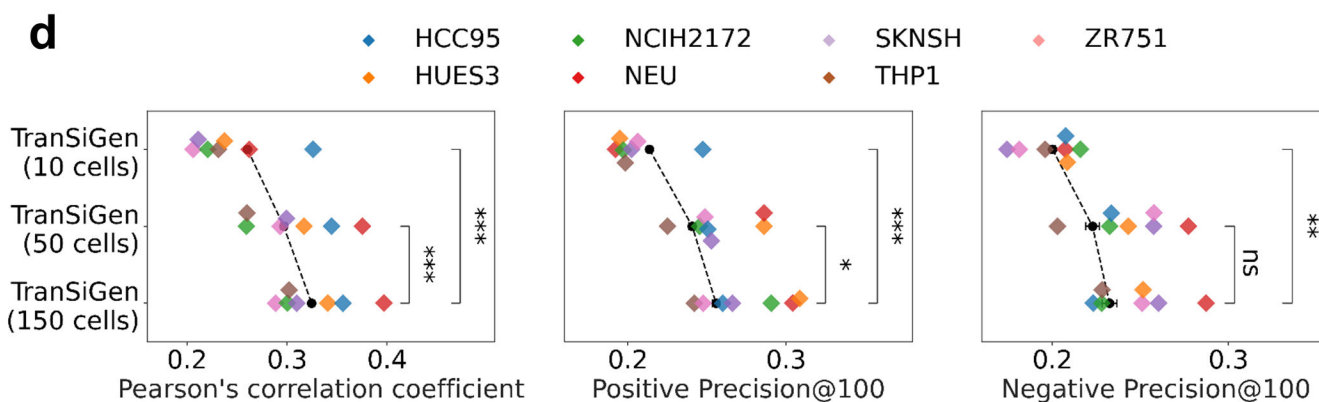
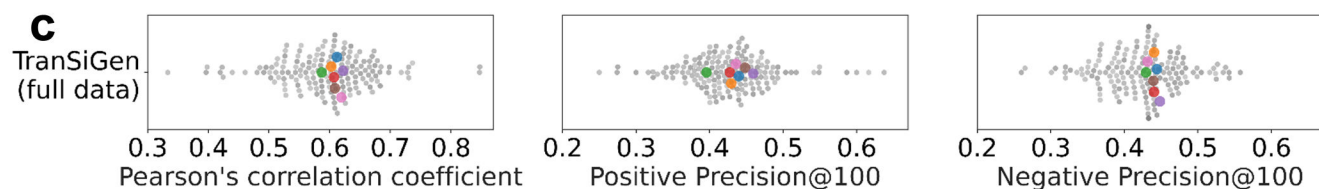
Evaluation metric	Equation*
RMSE	$\sqrt{\frac{1}{n} \sum_{i=1}^n (\Delta X_i - \Delta X'_i)^2}$
Pearson's correlation coefficient	$\frac{cov(\Delta X - \Delta X')}{\sigma_{\Delta X} \sigma_{\Delta X'}}$
Positive Precision @100	$\frac{G_{100-positive} \cap G'_{100-positive}}{G'_{100-positive}}$
Negative Precision@100	$\frac{G_{100-negative} \cap G'_{100-negative}}{G'_{100-negative}}$

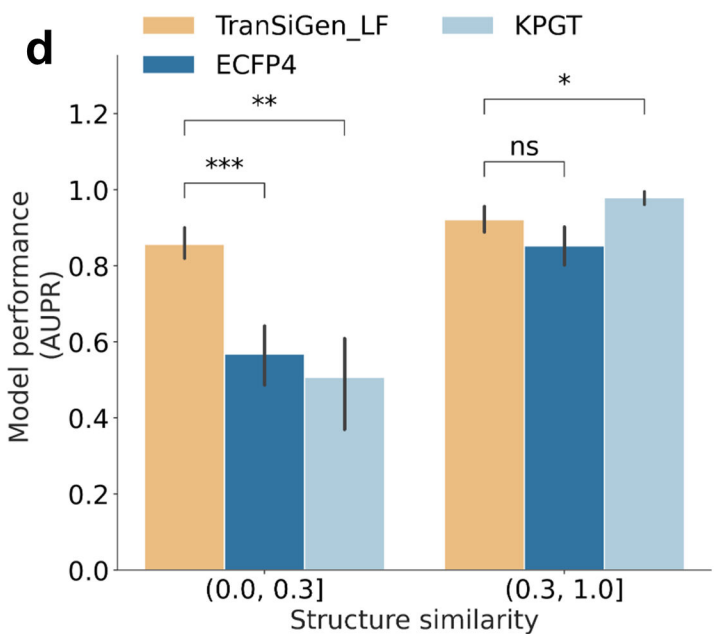
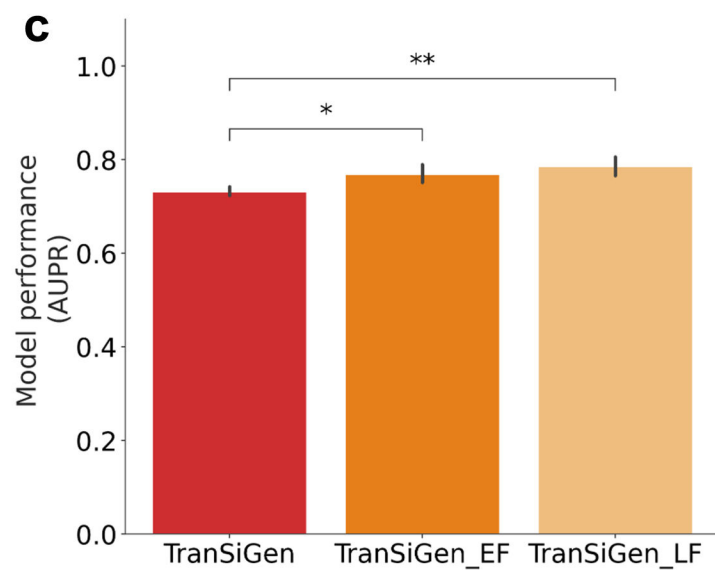
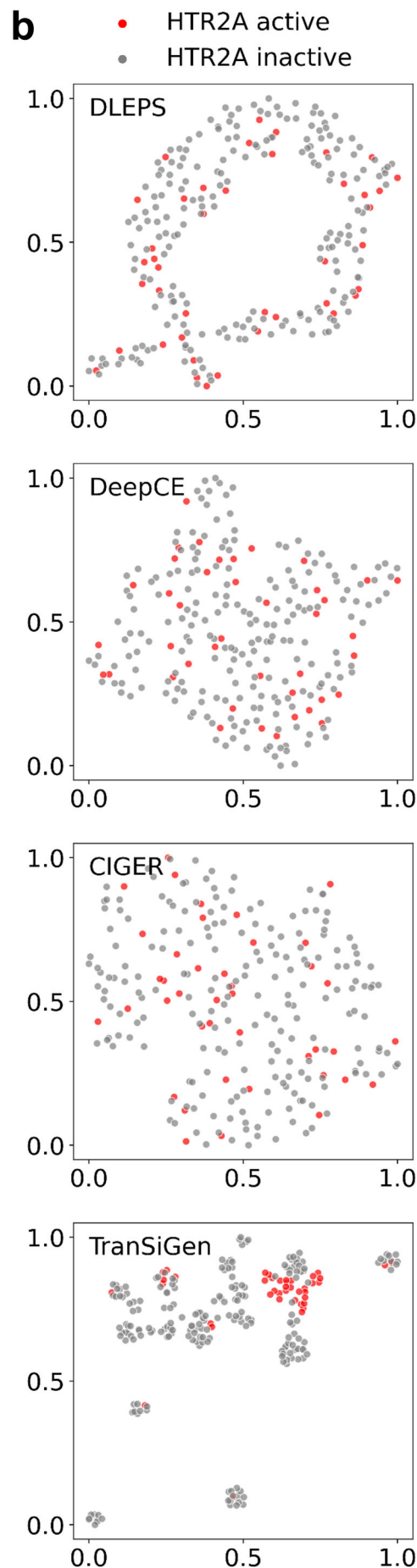
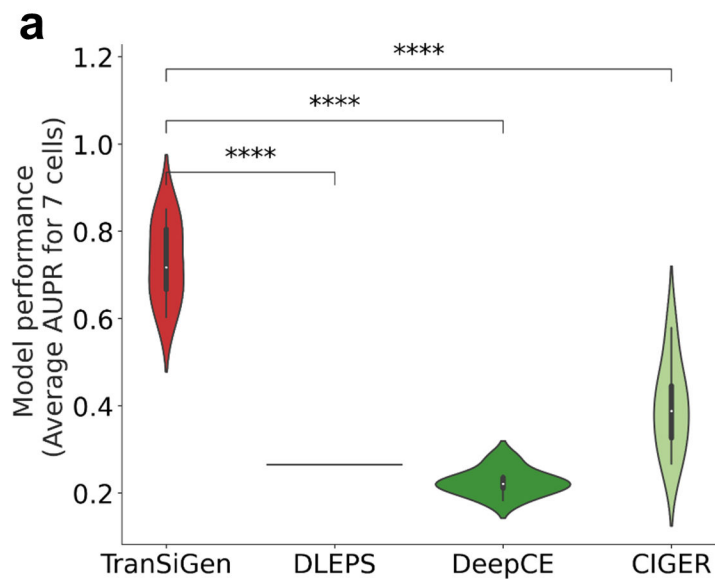
764 \*  $n$  represents the number of landmark genes in expression profiles,  $G$  represents the sets of top

765 100 positive/negative genes,  $G'$  represents the sets of top 100 predicted genes.

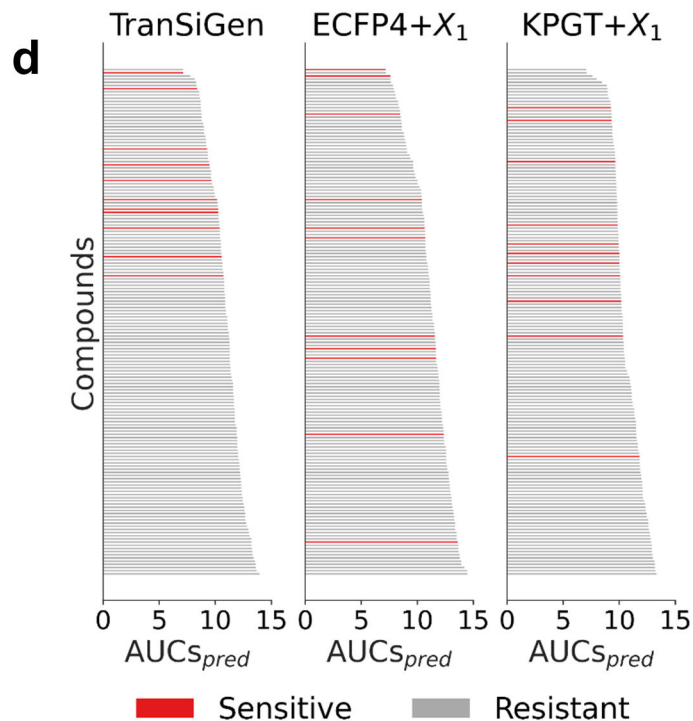
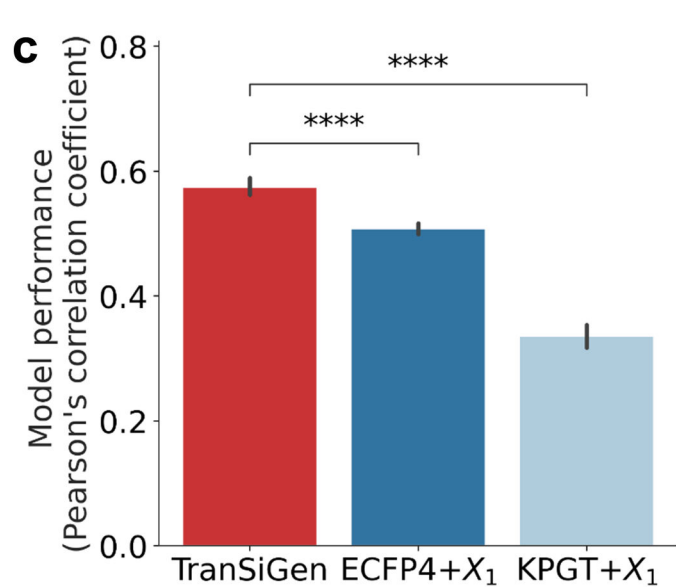
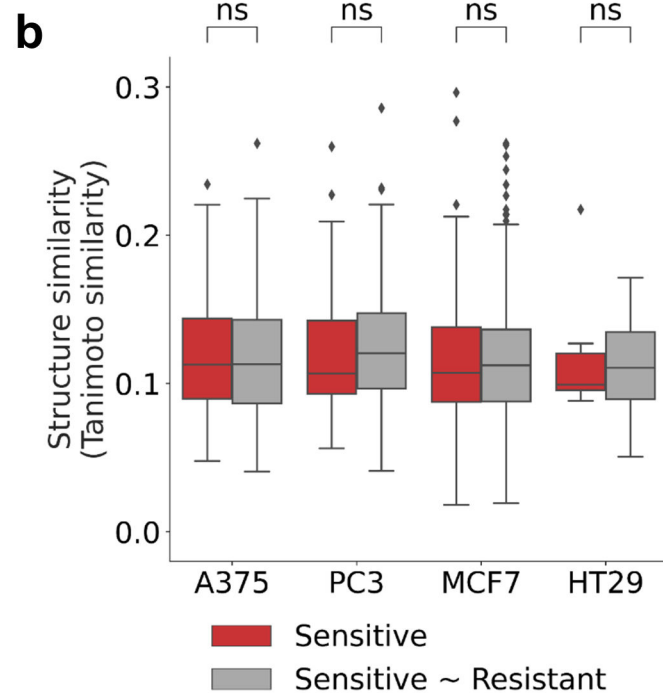
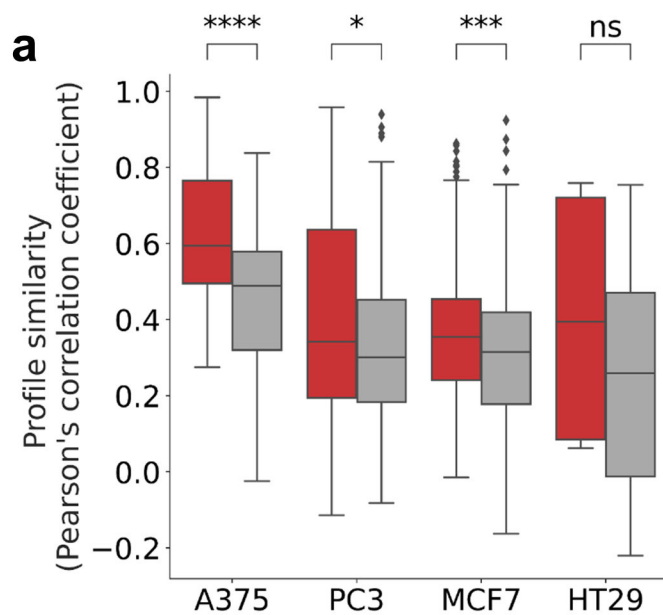


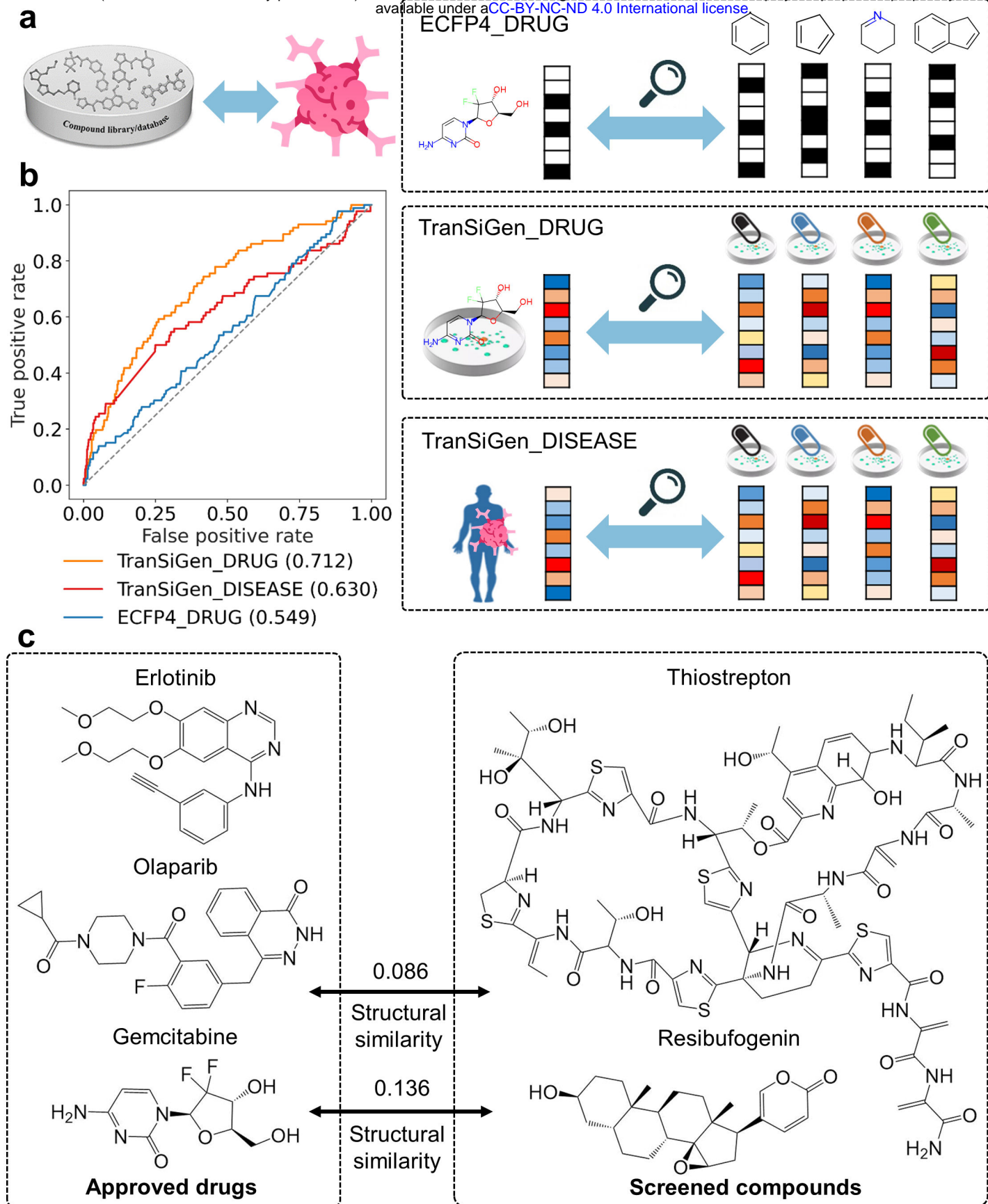












**d**

	ECFP4_DRUG	TranSiGen_DRUG	TranSiGen_DISEASE
Thioestrepton	1289/1625	28/1625	<b>2/1625</b>
Resibufogenin	616/1625	35/1625	<b>10/1625</b>