

Combinatorial prediction of therapeutic perturbations using causally inspired neural networks

Received: 1 January 2024

Accepted: 10 July 2025

Published online: 9 September 2025

 Check for updates

Guadalupe Gonzalez^{1,2,3,14}, Xiang Lin^{4,14}, Isuru Herath^{5,6,7}, Kirill Veselkov^{1,8}, Michael Bronstein^{9,10} & Marinka Zitnik^{4,11,12,13} ✉

Phenotype-driven approaches identify disease-counteracting compounds by analysing the phenotypic signatures that distinguish diseased from healthy states. Here we introduce PDGrapher, a causally inspired graph neural network model that predicts combinatorial perturbagens (sets of therapeutic targets) capable of reversing disease phenotypes. Unlike methods that learn how perturbations alter phenotypes, PDGrapher solves the inverse problem and predicts the perturbagens needed to achieve a desired response by embedding disease cell states into networks, learning a latent representation of these states, and identifying optimal combinatorial perturbations. In experiments in nine cell lines with chemical perturbations, PDGrapher identifies effective perturbagens in more testing samples than competing methods. It also shows competitive performance on ten genetic perturbation datasets. An advantage of PDGrapher is its direct prediction, in contrast to the indirect and computationally intensive approach common in phenotype-driven models. It trains up to 25× faster than existing methods, providing a fast approach for identifying therapeutic perturbations and advancing phenotype-driven drug discovery.

Target-driven drug discovery, which has been the dominant approach since the 1990s, focuses on designing highly specific compounds to act against targets, such as proteins or enzymes, that are implicated in disease, often through genetic evidence^{1–3}. An example of target-driven drug discovery is the development of small molecule kinase inhibitors like imatinib. Imatinib stops the progression of chronic myeloid leukaemia by inhibiting BCR-ABL tyrosine kinase, a mutated protein involved in uncontrolled proliferation of leukocytes in patients with chronic myeloid leukemia⁴. Other notable examples include monoclonal antibodies such as trastuzumab, which specifically targets the HER2 receptor, a protein overexpressed in certain types of breast cancer. Trastuzumab

inhibits cell proliferation while engaging the body's immune system to initiate an anti-cancer response⁵. These examples illustrate the success of target-driven drug discovery, yet the past decade has seen a revival of phenotype-driven approaches. This shift has been fuelled by the observation that many first-in-class drugs approved by the US Food and Drug Administration (FDA) between 1999 and 2008 were discovered without a drug target hypothesis⁶. Instead of the 'one drug, one gene, one disease' model of target-driven approaches, phenotype-driven drug discovery focuses on identifying compounds or, more broadly, perturbagens—combinations of therapeutic targets—that reverse disease phenotypes as measured by assays without predefined targets^{1,7}.

¹Imperial College London, London, UK. ²F. Hoffmann-La Roche Ltd, Basel, Switzerland. ³Prescient Design, Genentech, South San Francisco, CA, USA.

⁴Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA. ⁵Merck & Co., South San Francisco, CA, USA. ⁶Cornell University, Ithaca, NY, USA. ⁷University of Pittsburgh School of Medicine–Carnegie Mellon University, Pittsburgh, PA, USA. ⁸Yale School of Public Health, New Haven, CT, USA. ⁹University of Oxford, Oxford, UK. ¹⁰AITHYRA, Vienna, Austria. ¹¹Kempner Institute for the Study of Natural and Artificial Intelligence, Harvard University, Cambridge, MA, USA. ¹²Broad Institute of MIT and Harvard, Cambridge, MA, USA. ¹³Harvard Data Science Initiative, Cambridge, MA, USA.

¹⁴These authors contributed equally: Guadalupe Gonzalez, Xiang Lin. ✉ e-mail: marinka@hms.harvard.edu

Ivacaftor illustrates how these approaches can intersect. Although ivacaftor was developed through a target-driven strategy to modulate the cystic fibrosis transmembrane conductance regulator protein in individuals with specific mutations, its development relied on phenotypic assays to confirm functional improvements, such as increased chloride transport^{8–10}.

Phenotype-driven drug discovery has been bolstered by the advent of chemical and genetic libraries such as the Connectivity Map (CMap)¹¹ and the Library of Integrated Network-based Cellular Signatures (LINCS)¹². CMap and LINCS contain gene expression profiles of dozens of cell lines treated with thousands of genetic and chemical perturbagens. CMap introduced connectivity scores to quantify similarities between compound responses and disease gene expression signatures. Identifying compounds with gene expression signatures either similar to those of known disease-treating drugs or that counter disease signatures can help in selecting therapeutic leads^{13–16}. These strategies have successfully identified drugs with high in vitro efficacy¹⁶ across a range of diseases^{17–19}.

Deep learning methods have been used for lead discovery by predicting gene expression responses to perturbagens, including perturbagens that were not yet experimentally tested^{20–23}. However, these approaches rely on chemical and genetic libraries, meaning that they select perturbagens from predefined libraries and cannot identify perturbagens as new combinations of drug targets. Further, they are perturbation response methods that predict changes in phenotypes upon perturbations. Thus, they can identify perturbagens by exhaustively predicting responses to all perturbations in the library and then searching for perturbagens with the desired response. Unlike existing methods that learn responses to perturbations, phenotype-based approaches need to solve the inverse problem, which is to infer perturbagens necessary to achieve a specific response—that is, directly predicting perturbagens by learning which perturbations elicit a desired response.

In causal discovery, the problem of identifying which elements of a system should be perturbed to achieve a desired state is referred to as optimal intervention design^{24–26}. Using insights from causal discovery and geometric deep learning, here we introduce PDGrapher, an approach for the combinatorial prediction of therapeutic targets that can shift gene expression from an initial diseased state to a desired treated state. PDGrapher is formulated using a causal model in which genes represent the nodes in a causal graph and structural causal equations define their causal relationships. Given a genetic or chemical intervention dataset, PDGrapher pinpoints a set of genes that a pertubagen should target to facilitate the transition of node states from diseased to treated. PDGrapher uses protein–protein interaction (PPI) networks or gene regulatory networks (GRNs) as approximations of the causal graph, operating under the assumption of no unobserved confounders. PDGrapher tackles the optimal intervention design using representation learning, using a graph neural network (GNN) to represent structural equations.

PDGrapher is trained on a dataset of disease–treated sample pairs to predict therapeutic gene targets that can shift the gene expression phenotype from a diseased to a healthy or treated state. Once trained, PDGrapher processes a new diseased sample and outputs a pertubagen—a set of therapeutic targets—predicted to counteract the disease effects in that specific sample. We evaluate PDGrapher across 19 datasets, comprising genetic and chemical interventions across 11 cancer types and two proxy causal graphs. We also consider different evaluation set-ups, including settings where held-out folds contain new samples in the same cell line with the training samples, and settings where held-out folds contain new samples from a cancer type that PDGrapher has never encountered during training. In held-out folds that contain new samples, PDGrapher detects up to 13.37% and 1.09% more ground-truth therapeutic targets in chemical and genetic intervention datasets, respectively, than existing methods. We also

find that in chemical intervention datasets, candidate therapeutic targets predicted by PDGrapher are on average up to 11.58% closer to ground-truth therapeutic targets in the gene–gene interaction network than what would be expected by chance. Even in held-out folds containing new samples from a previously unseen cancer type, PDGrapher maintains robust performance. Unlike methods that indirectly identify perturbagens by predicting cell responses, PDGrapher directly predicts perturbagens that can shift gene expression from diseased to treated states. This feature of PDGrapher enables model training up to 25× faster than indirect prediction methods, such as scGen²² and CellOT²⁷. As these approaches build a separate model for each perturbation, they become increasingly ineffective when applied to datasets with a large number of perturbagens. For example, with its default settings, CellOT needs 10 h to train for a single pertubagen in a cell line from the LINCS dataset.

PDGrapher can aid in elucidating the mechanism of action of chemical perturbagens (Supplementary Fig. 3), which we show in the case of vorinostat, a histone deacetylase inhibitor used to treat cutaneous T cell lymphoma, and sorafenib, a multikinase inhibitor used in the treatment of several types of cancer (Supplementary Note 1). PDGrapher can also suggest potential anti-cancer therapeutic targets: it highlighted kinase insert domain receptor (*KDR*) as a top predicted target for non-small cell lung cancer (NSCLC; see ‘PDGrapher identifies therapeutic targets validated through clinical and biological evidence’ in Results). It identified associated drugs—vandetanib, sorafenib, catequantinib and rivocecranib—that inhibit the kinase activity of the protein encoded by *KDR*. These drugs block VEGF signalling, suppressing endothelial cell proliferation, migration and blood vessel formation, which tumours rely on for growth and metastasis^{28,29}. By predicting combinatorial therapeutic targets based on phenotypic transitions, PDGrapher provides a scalable approach to phenotype-driven perturbation modelling.

Results

Overview of intervention datasets and causal graphs

We evaluate our method across a total of 38 preprocessed datasets that span 2 types of intervention (genetic and chemical), 11 cancer types (lung, breast, prostate, colon, skin, cervical, head and neck, pancreatic, stomach, brain and ovarian), and 2 types of proxy causal graph: PPI networks and GRNs. Each dataset is uniquely defined by a combination of intervention type, causal graph type, cancer type, and cell line, and is denoted in the format: treatment type-graph type-cancer type-cell line. The chemical-PPI datasets include cell lines A549 (lung); MCF7, MDAMB231 and BT20 (breast); PC3 and VCAP (prostate); HT29 (colon); A375 (skin); and HELA (cervix). The genetic-PPI datasets include A549 (lung), MCF7 (breast), PC3 (prostate), HT29 (colon), A375 (skin), ES2 (ovary), BICR6 (head and neck), YAPC (pancreas), AGS (stomach) and U251MG (brain). Similarly, the chemical-GRN and genetic-GRN datasets span the same combinations of cancer types and cell lines as their PPI counterparts. This comprehensive collection of datasets enables systematic benchmarking of our model across diverse perturbation modalities, biological contexts and graph structures. Genetic interventions are single-gene knockout experiments by CRISPR–Cas9-mediated gene knockouts, while chemical interventions are multiple-gene treatments induced using chemical compounds. We use a PPI network from BIOGRID that has 10,716 nodes and 151,839 undirected edges. We additionally construct GRNs for each disease-treatment type pair using GENIE3 (ref. 30) (Supplementary Note 3), with GRNs on average having 10,000 nodes and 500,000 directed edges. The training data for PDGrapher consist of two components: disease intervention data and treatment intervention data. Disease intervention data include paired healthy and diseased gene expression profiles, along with associated disease genes; however, these data are only available for cell lines corresponding to lung, breast and prostate cancers. In contrast, the treatment intervention data comprise paired diseased and treated gene expression

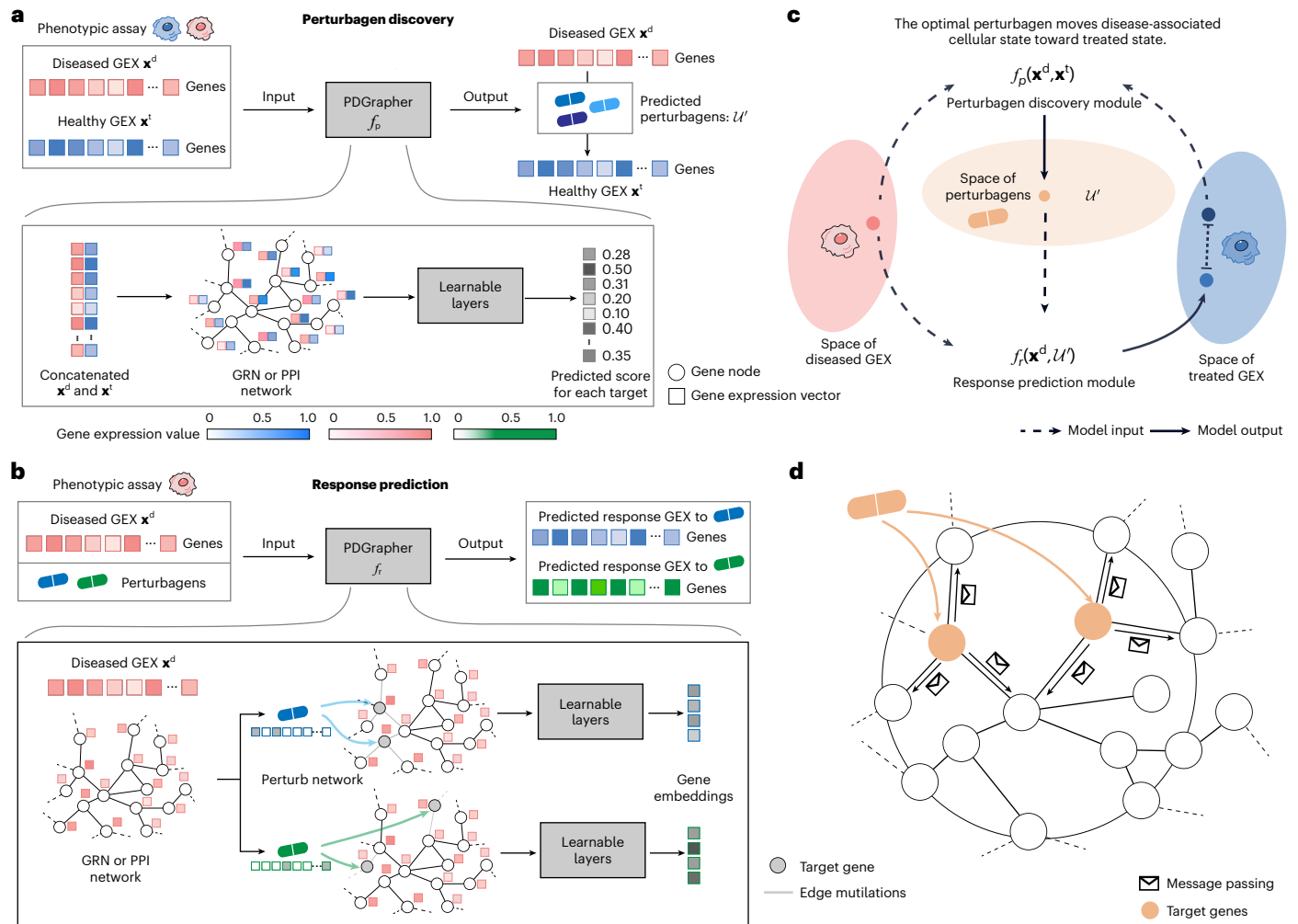


Fig. 1 | Overview of PDGrapher. a, Given a paired diseased and treated gene expression sample and a proxy causal graph, PDGrapher’s perturbagen discovery module, f_p , predicts a candidate set of therapeutic targets to shift cell gene expression from a diseased to a treated state. **b**, Given a disease sample’s gene expression, a proxy causal graph and a set of perturbagens, PDGrapher’s response prediction module, f_r , predicts the gene expression response of the sample to each perturbagen. f_r represents perturbagen’s effects in the graph as edge mutilations. **c**, f_p is optimized using two losses: a cross-entropy cycle loss to

predict a perturbagen U' that aims to shift the diseased cell state closely approximating the treated state, $CE(x^t, f_r(x^d, f_p(x^d, x^t)))$ (with f_r frozen), and a cross-entropy supervision loss that directly supervises the prediction of U' , $CE(U', f_p(x^d, x^t))$ (see Methods for more details). **d**, Both f_r and f_p follow the standard message-passing framework, where node representations are updated by aggregating the information from neighbours in the graph. GEX, gene expression.

samples, along with genetic or chemical perturbagens, and are available across all datasets. Supplementary Tables 11 and 12 summarize the number of samples for each cell line and intervention datasets.

Overview of PDGrapher model

Given a diseased cell state (gene expression profile), the goal of PDGrapher is to predict the genes that, if targeted by a perturbagen, would shift the cell to a treated state (Fig. 1a). Unlike methods for learning the response of cells to a given perturbation^{22,27,31,32}, PDGrapher focuses on the inverse problem by learning which perturbation elicits a desired response. PDGrapher predicts perturbagens that shift cellular states under the assumption that an optimal perturbagen is one that alters the gene expression profile of a cell to closely match a desired target state. Our approach comprises two modules (Fig. 1a–c). First, a perturbagen discovery module f_p takes the initial and desired cell states and outputs a candidate perturbagen as a set of therapeutic targets U' (Fig. 1a). Then, a response prediction module f_r takes the initial state and the predicted perturbagen U' and predicts the cell response upon perturbing genes in U' (Fig. 1b). Our response prediction and

perturbagen discovery modules are GNN models that operate on a proxy causal graph, where edge mutilations, or edge removals, represent the effects of interventions on the graph (Fig. 1c).

PDGrapher is trained using an objective function with two components, one for each module, f_r and f_p . The response prediction module f_r is trained using disease and treatment intervention data on cell state transitions so that the predicted cell states are close to the known perturbed cell states upon interventions. The perturbagen discovery module f_p is trained only using the treatment intervention data; given a diseased cell state, f_p predicts the set of therapeutic targets U' that caused the corresponding treated cell state. The objective function for the perturbagen discovery module consists of two elements: (1) a cycle loss that optimizes the parameters of f_p such that the response upon intervening on the predicted genes in U' , as measured by f_r , closely approximates the actual treated cellular state; and (2) a supervision loss on the therapeutic target set U' that directly pushes PDGrapher to predict the correct perturbagen. Both models are trained simultaneously using early stopping independently so that each model finishes training upon convergence.

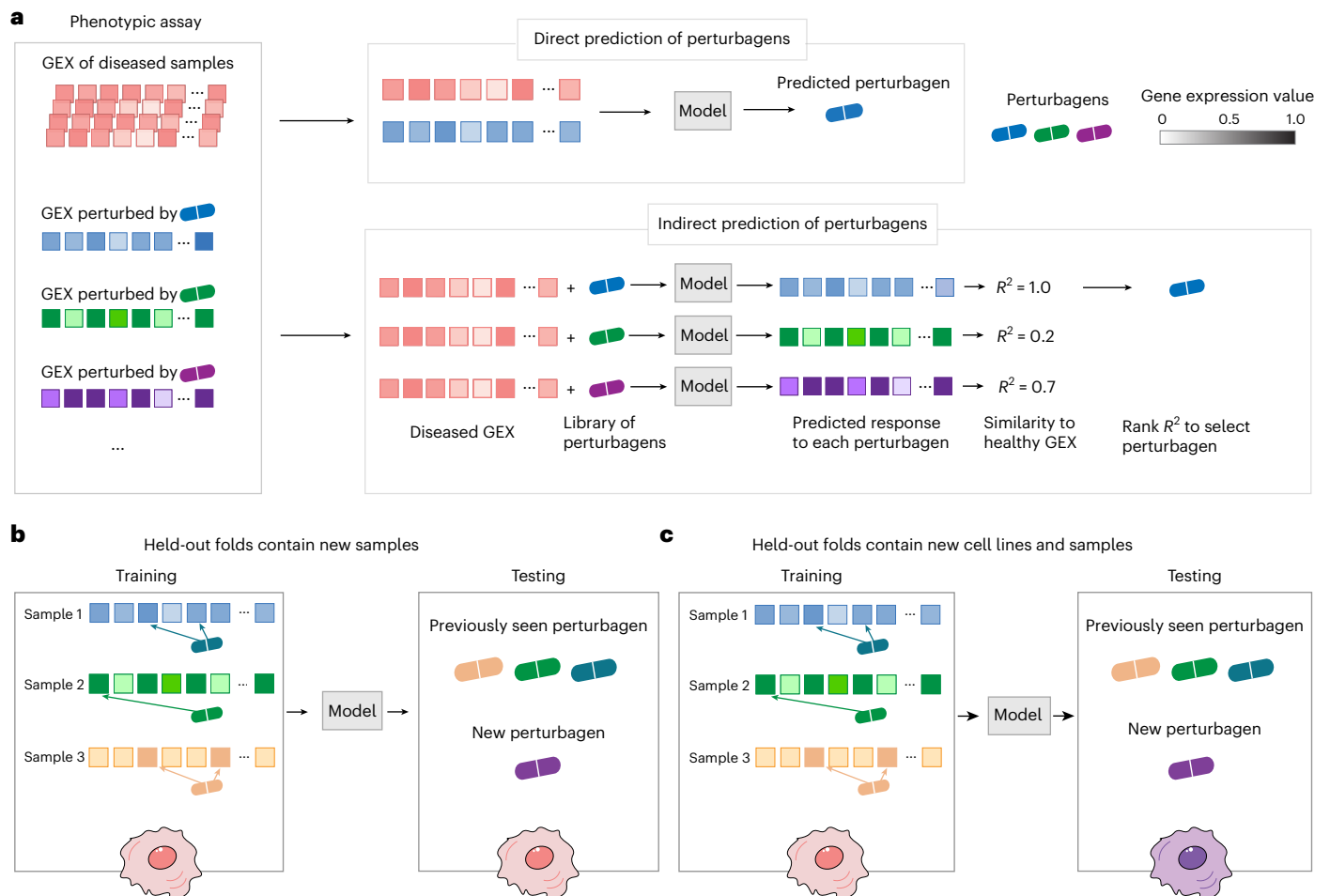


Fig. 2 | Overview of evaluation settings and data splits. a, Given a dataset with paired diseased and treated samples and a set of perturbagens, PDGrapher makes a direct prediction of candidate perturbagens that shift gene expression from a diseased to a treated state for each disease–treated sample pair. The direct prediction means that PDGrapher directly infers the perturbation necessary to achieve a specific response. In contrast to direct prediction of perturbagens, existing methods predict perturbagens only indirectly through a two-stage approach. For a given diseased sample, the model learns the response to each candidate perturbation from an existing library and identifies the perturbation whose induced response most closely approximates the desired treated state.

Existing methods learn the response of cells to a given perturbation^{22,27,31,32}, whereas PDGrapher focuses on the inverse problem by learning which perturbation elicits a given response, even in the most challenging cases when the combinatorial composition of perturbagens was never seen before. **b,c**, We evaluate PDGrapher’s performance across two settings: given a cell line, randomly splitting samples between training and testing set (**b**), and by splitting samples based on cell lines, where the model is trained on one cell line and evaluated on a different, previously unseen cell line to assess PDGrapher’s generalization performance (**c**).

When trained, PDGrapher predicts perturbagens—as sets of candidate target genes—to shift cells from diseased to treated. Given a pair of diseased and treated samples, PDGrapher directly predicts perturbagens by learning which perturbations elicit target responses. In contrast, existing approaches are perturbation response methods that predict changes in phenotype that occur upon perturbation; thus, they can only indirectly predict perturbagens (Fig. 2a). Given a disease–treated sample pair, a response prediction module (such as scGen²², ChemCPA²³, Biolord³³, GEARS³⁴ or CellOT²⁷) is used to predict the response of the diseased sample to a library of perturbagens. The predicted perturbation is the one that produces a response that is the most similar to the treated sample. We evaluate PDGrapher’s performance in two separate settings (Fig. 2b,c): (1) a random splitting setting, where the samples are split randomly between training and test sets within a cell line (denoted as random for convenience); and (2) a leave-cell-out setting, where PDGrapher is trained in one cell line, and its performance is evaluated in a cell line the model never encountered during training to test how well the model generalizes to a new disease. Supplementary Tables 1 and 2

show the numbers of unseen perturbagens in chemical perturbation datasets in the random and leave-cell-out splits, respectively; Supplementary Tables 3 and 4 show the number of unseen perturbagens in genetic perturbation datasets in the random and leave-cell-out splits, respectively.

PDGrapher predicts perturbagens to reverse disease states

In the random splitting setting, we assess the ability of PDGrapher for combinatorial prediction of therapeutic targets across chemical PPI datasets (chemical-PPI-lung-A549, chemical-PPI-breast-MCF7, chemical-PPI-breast-MDAMB231, chemical-PPI-breast-BT20, chemical-PPI-prostate-PC3, chemical-PPI-prostate-VCAP, chemical-PPI-colon-HT29, chemical-PPI-skin-A375 and chemical-PPI-cervix-HELA). Specifically, we measure whether, given paired diseased–treated gene expression samples, PDGrapher can predict the set of therapeutic genes targeted by the chemical compound in the diseased sample to generate the treated sample. Given paired diseased–treated gene expression samples, PDGrapher ranks genes in the dataset according to their likelihood of being therapeutic targets.

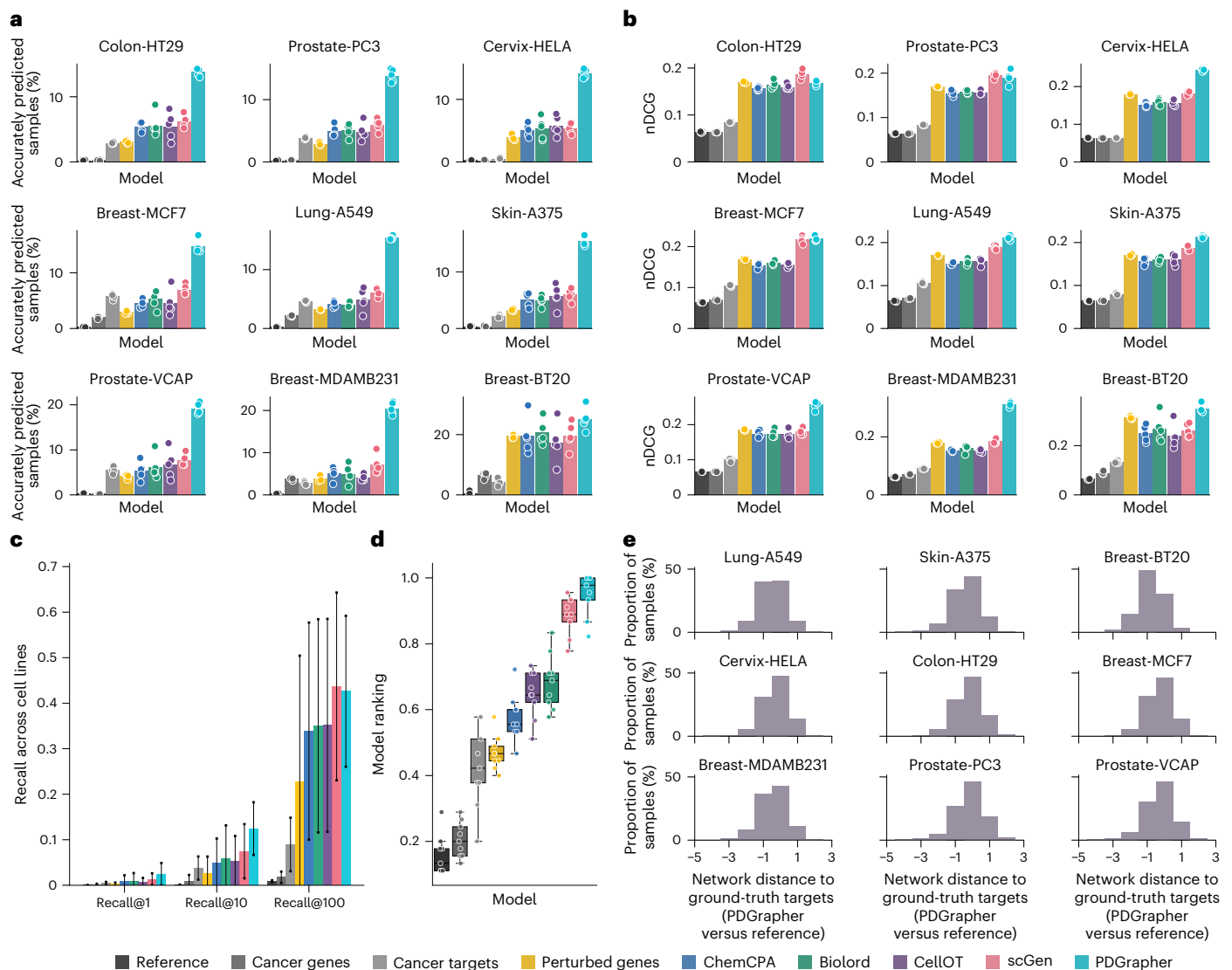


Fig. 3 | PDGrapher efficiently predicts chemical perturbagens to shift cells from diseased to treated states in held-out folds containing new samples.

a, b, PDGrapher shows improved performance across nine chemical perturbation datasets with various diseases, yielding up to 13.37% more accurately predicted samples in the testing sets compared with the second-best model (for example, for chemical-PPI-breast-MDAMB231, 20.43% versus 7.05% (**a**)) and up to 0.13 higher nDCG than the second-best model (for chemical-PPI-breast-MDAMB231, 0.31 versus 0.18 (**b**)). In **a** and **b**, the bars show the average performance across five cross-validation test splits for each of the nine chemical datasets. The overlaid points represent performance values from individual data splits ($n = 5$ per cell line). Each data split contains 20% of samples in the dataset, with each sample corresponding to a perturbation-response instance. Where replicates exist for a given drug, they are treated as independent inputs during training and evaluation. **c**, PDGrapher recovers ground-truth therapeutic targets at higher rates (evaluated by recall 1–100) compared with competing methods

for chemical-PPI datasets. **d**, Box plots show the distribution of average model rankings across 9 cell lines ($n = 9$); each dot corresponds to the aggregated ranking value across cross-validation splits and across all metrics for a distinct cell line. A higher value indicates better performance. The central line inside the box represents the median, while the top and bottom edges correspond to the first and third quartiles. The whiskers extend to the smallest and largest values within 1.5 \times the interquartile range from the quartiles. Each dot represents a data point for a specific cell line. P values from the statistical tests are provided in the Source data. **e**, Shown is the difference of shortest-path distances between ground-truth therapeutic genes and predicted genes by PDGrapher and a random reference across nine cell lines. Predominantly negative values indicate that PDGrapher predicts sets of therapeutic genes that are closer in the network to ground-truth therapeutic genes compared with what would be expected by chance (average shortest-path distances across cell lines for PDGrapher versus reference = 2.77 versus 3.11).

We quantify the ranking quality using normalized discounted cumulative gain (nDCG), where the gain reflects the ranking accuracy of the model. An nDCG value close to one indicates highly accurate predictions, with the top ranked gene targets closely matching the ground-truth targets, whereas lower nDCG values indicate poorer ranking performance. This metric provides a normalized and scalable measure of ranking quality, enabling consistent comparison across different datasets and models. PDGrapher outperforms competing methods in all cell lines, achieving nDCG values

that are higher than the second-best competing method by 0.02 (chemical-PPI-lung-A549), 0.13 (chemical-PPI-breast-MDAMB231), 0.03 (chemical-PPI-breast-BT20), 0.004 (chemical-PPI-breast-MCF7), 0.07 (chemical-PPI-prostate-VCAP), 0.005 (chemical-PPI-prostate-PC3), 0.03 (chemical-PPI-skin-A375), 0.06 (chemical-PPI-cervix-HELA) and 0.001 (chemical-PPI-colon-HT29) (Fig. 3b). In addition to evaluating the entire predicted target rank, it is even more practically crucial to assess the accuracy of the top ranked predicted targets. As perturbagens target multiple genes, we measure the fraction of samples in the test

set for which we obtain a partially accurate prediction, where at least one of the top N predicted gene targets corresponds to an actual gene target (denoted as the percentage of accurately predicted samples). Here, N represents the number of known target genes of a perturbation. PDGrapher consistently provides accurate predictions for more samples in the test set than competing methods. Specifically, it outperforms the second-best competing method by predicting ground-truth targets in an additional 7.73% (chemical-PPI-breast-MCF7), 9.32% (chemical-PPI-lung-A549), 13.37% (chemical-PPI-breast-MDAMB231), 4.50% (chemical-PPI-breast-BT20), 7.88% (chemical-PPI-prostate-PC3), 11.53% (chemical-PPI-prostate-VCAP), 7.56% (chemical-PPI-colon-HT29), 9.55% (chemical-PPI-skin-A375) and 8.41% (chemical-PPI-cervix-HELA) of samples (Fig. 3a). We also evaluate the performance of PDGrapher using recall@1, recall@10 and recall@100, which calculate the ratio of true target genes included in the top 1, top 10 and top 100 predicted gene targets, respectively. Although the absolute recall values are modest due to the inherent difficulty of the task, PDGrapher consistently outperforms all competing methods, showing its relative strength and robustness. Specifically, PDGrapher outperforms the second-best method in all the recall metrics with the averaged margin being 3.31% (chemical-PPI-lung-MCF7), 3.28% (chemical-PPI-lung-A549), 11.65% (chemical-PPI-breast-MDAMB231), 7.27% (chemical-PPI-breast-BT20), 2.50% (chemical-PPI-prostate-PC3), 9.53% (chemical-PPI-prostate-VCAP), 3.08% (chemical-PPI-colon-HT29), 2.87% (chemical-PPI-skin-A375) and 5.13% (chemical-PPI-cervix-HELA) (Fig. 3c). We then consolidated the results using the rankings from experiments across different cell lines and metrics for each method. PDGrapher achieved the best overall rankings, with a median significantly higher than all competing methods (Fig. 3d). P values of the chemical perturbation discovery tests are provided in the Source data.

PDGrapher not only provides accurate predictions for a larger proportion of samples and consistently predicts ground-truth therapeutic targets close to the top of the ranked list but it also predicts gene targets that are closer in the network (measured by the shortest-path distance) to ground-truth targets compared with what would be expected by chance (Fig. 3e). In all cell lines, the ground-truth therapeutic targets predicted by PDGrapher are significantly closer to the ground-truth targets compared with what would be expected by chance (Supplementary Table 6). For example, for chemical-PPI-lung-A549, the median distance between the predicted and ground-truth therapeutic targets is 3.0 for both PDGrapher and the random reference. However, the distributions show a statistically significant difference, with a 1-sided Mann–Whitney U -test that yields $P < 0.001$, an effect size (rank-biserial correlation) of 0.3531 (95% confidence interval (CI), [0.3515, 0.3549]) and a U -statistic of 1.29×10^{11} . Similarly, for chemical-PPI-breast-MCF7, the median distance is 3.0 for both groups, yet the distributions are significantly different ($P < 0.001$, effect size = 0.2160 (95% CI, [0.2146, 0.2174]), U -statistic = 3.91×10^{11}) (Supplementary Table 6). This finding suggests that PDGrapher predicts targets in a manner that reflects PPI network structure³². According to the local network hypothesis, which posits that genes in closer network proximity tend to be more functionally similar, PDGrapher's predictions are more functionally related to ground-truth targets than would be expected by chance^{35–37}.

PDGrapher also shows strong performance across genetic datasets, specifically genetic-PPI-lung-A549, genetic-PPI-breast-MCF7, genetic-PPI-prostate-PC3, genetic-PPI-skin-A375, genetic-PPI-colon-HT29, genetic-PPI-ovary-ES2, genetic-PPI-head-BICR6, genetic-PPI-pancreas-YAPC, genetic-PPI-stomach-AGS and genetic-PPI-brain-U251MG (Extended Data Fig. 2). Briefly, PDGrapher successfully detected ground-truth targets in 0.87% (genetic-PPI-lung-A549), 0.50% (genetic-PPI-breast-MCF7), 0.24% (genetic-PPI-prostate-PC3), 0.38% (genetic-PPI-skin-A375), 0.36% (genetic-PPI-colon-HT29), 1.09% (genetic-PPI-ovary-ES2), 0.54% (genetic-PPI-head-BICR6), 0.11% (genetic-PPI-pancreas-YAPC) and 0.92% (genetic-PPI-brain-U251MG) more samples compared with the

second-best competing method (Extended Data Fig. 2a). Its ability to effectively predict targets at the top of the ranks is further supported by the metrics recall@1 and recall@10 (Extended Data Fig. 2c). PDGrapher achieves the second-best overall rankings (Extended Data Fig. 2d), closely following scGEN, which obtained the highest nDCG values (Extended Data Fig. 2b) but showed weaker performance when evaluating only the top-ranked predicted targets (partially accurate prediction and recall@1). P values of the genetic perturbation discovery tests are provided in the Source data. PDGrapher and competing methods perform worse on genetic data than on chemical data. This may be due to knockout interventions generating weaker phenotypic signals than small molecule interventions. While gene knockouts are essential for understanding gene function, single-gene knockout studies can offer limited insights into complex cellular processes due to compensatory mechanisms^{38–40}. Despite the modest performance in genetic intervention datasets, PDGrapher outperforms competing methods in the combinatorial prediction of therapeutic targets.

PDGrapher achieves the best performance in response prediction for both chemical (Extended Data Fig. 1) and genetic perturbation (Extended Data Fig. 2e–g). P values of the response prediction tests are provided in the Source data. When using GRNs as proxy causal graphs, PDGrapher has comparable performance with GRNs and PPI networks across both genetic and chemical intervention datasets (Supplementary Figs. 4 and 5). One difference is that GRNs were constructed individually for each cell line, which makes leave-cell-out splitting setting prediction particularly challenging. Therefore, we only conducted random splitting setting experiments for GRN datasets. We also used PDGrapher to clarify the mode of action of the chemical perturbations vorinostat and sorafenib in chemical-PPI-lung-A549 (Supplementary Note 1).

PDGrapher generalizes to cell lines unseen during training

PDGrapher shows consistently strong performance on chemical intervention datasets in the leave-cell-out setting (Fig. 4). In this setting, we use the trained models in the random splitting setting for each cell line to predict therapeutic targets in the remaining cell lines. PDGrapher successfully predicts perturbagens that describe the cellular dynamics and shift gene expression phenotypes from a diseased to a treated state in 7.16% (chemical-PPI-breast-MCF7), 6.50% (chemical-PPI-lung-A549), 5.00% (chemical-PPI-breast-MDAMB231), 8.67% (chemical-PPI-prostate-PC3), 7.72% (chemical-PPI-prostate-VCAP), 7.31% (chemical-PPI-skin-A375), 7.08% (chemical-PPI-colon-HT29) and 7.13% (chemical-PPI-cervix-HELA) additional testing samples compared with the second-best competing method (Fig. 4a). PDGrapher also outperforms the competing methods in 8 of 9 cell lines by predicting nDCG values that are 0.02 (chemical-PPI-breast-MCF7), 0.01 (chemical-PPI-lung-A549), 0.01 (chemical-PPI-breast-MDAMB231), 0.03 (chemical-PPI-prostate-PC3), 0.02 (chemical-PPI-prostate-VCAP), 0.01 (chemical-PPI-skin-A375), 0.03 (chemical-PPI-colon-HT29) and 0.02 (chemical-PPI-cervix-HELA) higher than those of the second-best competing method (Fig. 4b). Its strong performance is further supported by the recall metrics, particularly recall@10 (Fig. 4c). Considering the overall performance across different cell lines and metrics, PDGrapher achieves the highest rank, with a median surpassing competing methods (Fig. 4d). Combinations of therapeutic targets predicted by PDGrapher in chemical datasets are closer to ground-truth targets than expected by chance (Fig. 4e and Supplementary Table 7). For example, for chemical-PPI-lung-A549, the median distance between predicted and ground-truth therapeutic targets is 3.0 for both PDGrapher and the random reference. However, the distributions show a statistically significant difference, with a 1-sided Mann–Whitney U -test yielding $P < 0.001$, an effect size (rank-biserial correlation) of 0.2191 (95% CI, [0.2182, 0.2200]) and a U -statistic of 2.46×10^{12} . Similarly, for chemical-PPI-breast-MCF7, the median distance is 3.0 for both groups, yet the distributions are significantly different ($P < 0.001$, effect size = 0.2457 (95% CI, [0.2451,

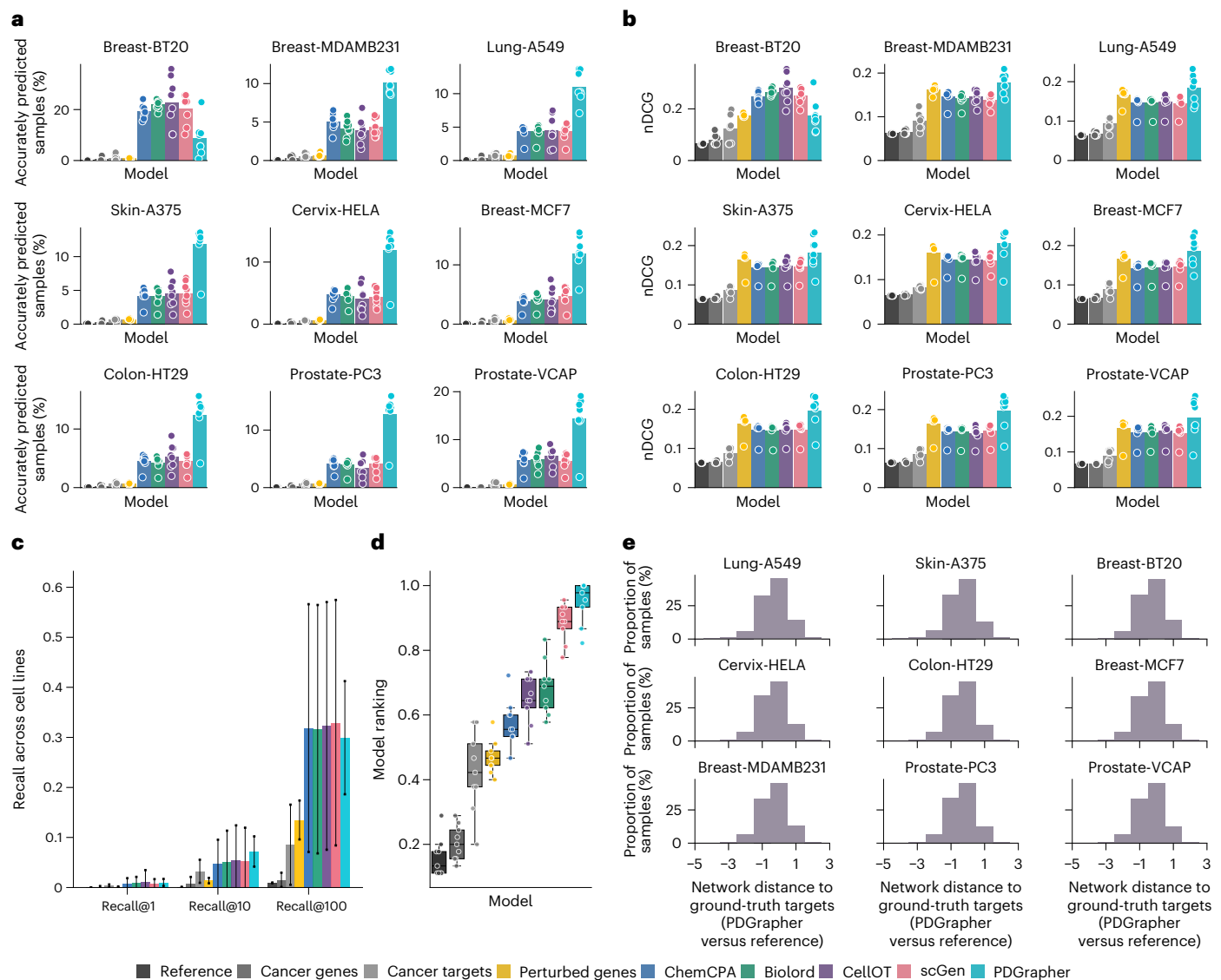


Fig. 4 | PDGrapher generalizes to new (previously unseen) cell lines and learns optimal chemical perturbagens in held-out folds that contain both new cell lines and new samples. a, b, PDGrapher shows improved performance when trained on nine chemical perturbation datasets spanning various diseases and evaluated on the remaining eight cell lines. It achieves up to 8.67% more accurately predicted samples in the testing sets compared with the second-best baseline (for example, when trained on chemical-PPI-prostate-PC3, 12.81% versus 4.13% (a)) and an nDCG value of up to 0.03 higher than the second-best baseline (for example, when trained on chemical-PPI-colon-HT29, 0.19 versus 0.16 (b)). In a and b, the bars show the average performance across five cross-validation test splits for each of the nine chemical datasets. The overlaid points represent performance values from individual data splits ($n = 5$ per cell line). Each data split contains 20% samples in the dataset, with each sample corresponding to a perturbation-response instance. Where replicates exist for a given drug, they are treated as independent inputs during training and evaluation. c, PDGrapher recovers ground-truth therapeutic targets at higher rates (evaluated by recall

1–100) compared with competing methods for chemical-PPI datasets. d, Box plots show the distribution of average model rankings across 9 cell lines ($n = 9$); each dot corresponds to the aggregated ranking value across cross-validation splits, train cell lines and across all metrics for a distinct cell line. A higher value indicates better performance. The central line inside the box represents the median, while the top and bottom edges correspond to the first and third quartiles. The whiskers extend to the smallest and largest values within $1.5 \times$ the interquartile range from the quartiles. Each dot represents a data point for a specific cell line and metrics. *P* values from the statistical tests are provided in the Source data. e, Shown is the difference of shortest-path distances between ground-truth therapeutic genes and predicted genes by PDGrapher and a random reference across nine cell lines. Predominantly negative values indicate that PDGrapher predicts sets of therapeutic genes that are closer in the network to ground-truth therapeutic genes compared with what would be expected by chance (average shortest-path distances across cell lines for PDGrapher versus random reference = 2.75 versus 3.11).

0.2464], U -statistic = 6.07×10^{12}) (Supplementary Table 7). PDGrapher also outperforms existing methods in genetic perturbagen prediction across cell lines, as measured by the top targets on the predicted gene ranks (Supplementary Fig. 2a–d). PDGrapher also shows superior performance in response prediction for both chemical (Supplementary Fig. 1) and genetic (Supplementary Fig. 2e–g) datasets. The *P* values of the leave-cell-out perturbagen discovery tests and response prediction tests are provided in the Source data.

Approaches that train individual models for each perturbagen (such as scGen and CellOT) generally achieve a better perturbagen prediction performance than those that use a single model for all perturbagens (Biolord, GEARS and ChemCPA). However, training individual models becomes infeasible for large-scale datasets with many perturbagens. For example, without parallelization, scGen would require about 8 years to complete the leave-cell-out experiments on the chemical and genetic perturbation data used in this study. PDGrapher

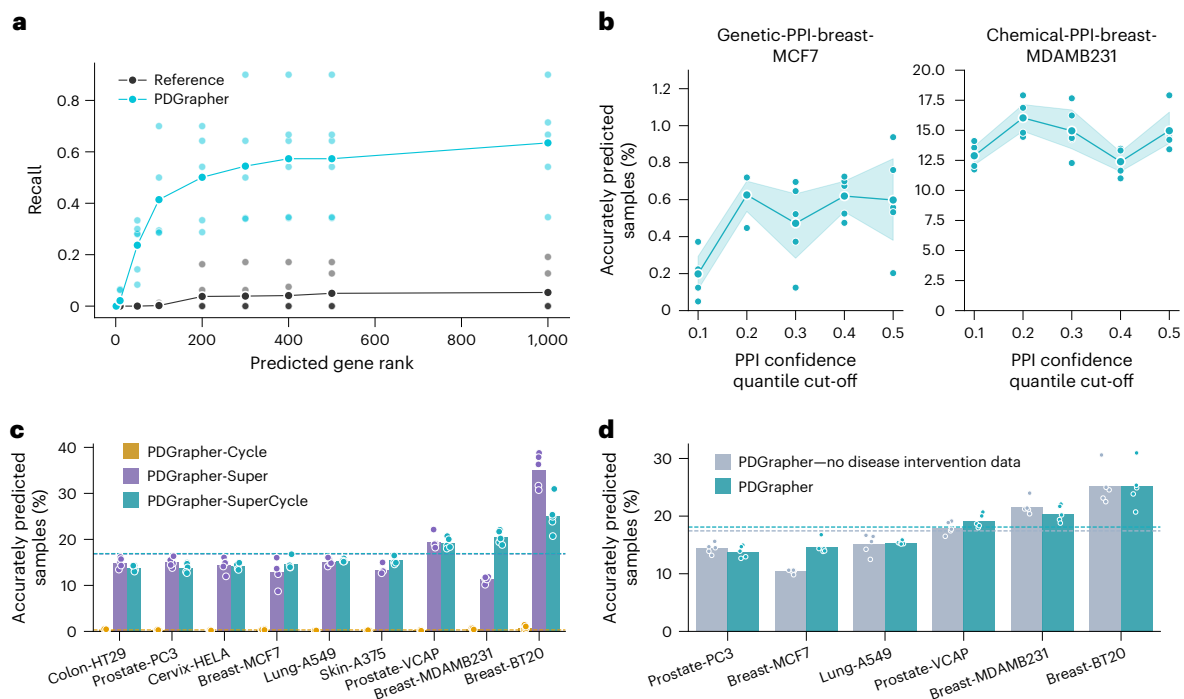


Fig. 5 | PDGrapher shows robust performance across training strategies, PPI networks and data availability settings. a, Performance of PDGrapher in the prediction of unseen approved drug targets to reverse disease effects across all cell lines with healthy counterparts in chemical perturbation datasets. Individual data points represent individual cell lines ($n = 6$). **b**, Performance of sensitivity analyses evaluated by the percentage of accurately predicted samples for cell lines MDAMB231 and MCF7 under chemical and genetic perturbations, respectively. The PPI network used here is from STRING (string-db.org) with a confidence score for each edge. The edges are filtered by the 0.1, 0.2, 0.3, 0.4 and 0.5 quantiles of the confidence scores as cut-offs, resulting in 5 PPI networks with 625,818, 582,305, 516,683, 443,051 and 296,451 edges, respectively. Data are presented as mean values across five cross-validation data splits per PPI confidence quantile. Shaded bands represent ± 1 s.d. from the mean ($n = 5$ computational replicates per quantile). Each point corresponds to performance on a specific data split. **c**, Performance metrics of the ablation study on PDGrapher's objective function components: PDGrapher-Cycle trained using only the cycle loss, PDGrapher-SuperCycle trained using the supervision

and cycle loss, and PDGrapher-Super trained using only the supervision loss, evaluated by percentage of accurately predicted samples. PDGrapher-Cycle shows inferior performance, resulting in limited visibility in the bar plot. **d**, Performance metrics of the second ablation study on PDGrapher's input data: PDGrapher—no disease intervention data using only treatment intervention data, and PDGrapher using both disease and treatment intervention data. The disease and treatment intervention data are organized as 'healthy, mutation, disease' and 'diseased, drug, treated', respectively. In **c** and **d**, bars show the average performance across five cross-validation test splits for each of the nine chemical datasets. The overlaid points represent performance values from individual data splits ($n = 5$ per cell line). The dashed horizontal lines represent the average performance across all cell lines. Each data split contains 20% samples in the dataset, with each sample corresponding to a perturbation-response instance. Where replicates exist for a given drug, they are treated as independent inputs during training and evaluation. P values from the statistical tests are provided in the Source data.

addresses this scalability challenge. Its training is up to $25\times$ faster than scGen and more than $100\times$ faster than CellOT when using the default setting of 100,000 epochs, substantially reducing computational costs. This efficiency highlights a key advantage of PDGrapher. This improved efficiency is due to PDGrapher's approach. Existing methods predict phenotypic responses to perturbations and identify perturbagens indirectly by searching through predicted responses for all candidates. In contrast, PDGrapher directly infers the perturbagen needed to achieve a specific response, learning which perturbations elicit a desired effect.

PDGrapher predicts therapeutic targets supported by clinical and biological evidence

We examined PDGrapher's ability to predict targets of anti-cancer drugs that were not encountered by the model during training time using chemical cell lines with matched healthy data: chemical-PPI-lung-A549, chemical-PPI-breast-MCF7, chemical-PPI-breast-MDAMB231, chemical-PPI-breast-BT20, chemical-PPI-prostate-PC3 and chemical-PPI-prostate-VCAP. PDGrapher was used to predict gene targets to shift these diseased cell lines into their healthy states. Figure 5a shows the recovery of targets of FDA-approved drugs for varying values of K (where K represents the number of predicted target genes considered in the predicted ranked list), indicating that PDGrapher can identify

targets of approved anti-cancer drugs not seen during training among the top predictions.

We analysed lung cancer by comparing the targets predicted by PDGrapher for lung cancer cell lines with the targets of candidate drugs in clinical development, curated from the Open Targets Platform⁴¹. This evaluation tested PDGrapher's ability to predict combinatorial chemical perturbagens. We compared the top ten targets predicted by PDGrapher for the A549 lung cancer cell line to ten randomly selected genes. The predicted targets had significantly higher Open Targets scores and more supporting resources than the random genes (Extended Data Fig. 7). Using an Open Targets evidence cut-off score of 0.5, 8 of 10 predicted targets had evidence supporting their association with lung cancer, compared with only 2 of 10 in the random gene set. Four drugs, tacedinaline (DrugBank:DB12291; clinical trial identifier, NCT00005093), selpercatinib (DrugBank:DB15685), pralsetinib (DrugBank:DB15822) and dexmedetomidine (DrugBank:DB00633; ref. 42), targeting these predicted genes were not included in the training set but have been identified as potential treatments for NSCLC.

We then evaluated PDGrapher's predictions by examining FDA-approved drugs that were not present in the training set of PDGrapher. Specifically, we assessed PDGrapher's performance using the chemical-PPI-lung-A549 dataset, focusing initially on pralsetinib, a

targeted cancer therapy primarily used to treat NSCLC⁴³. Pralsetinib is a selective Ret proto-oncogene (RET) kinase inhibitor designed to block the activity of RET proteins that have become aberrantly active due to mutations or fusions. Pralsetinib is known to target 11 key proteins: RET, DDR1, NTRK3, fms-related receptor tyrosine kinase 3 (FLT3), JAK1, JAK2, NTRK1, KDR, platelet derived growth factor receptor- β (PDGFRB), fibroblast growth factor receptor 1 (FGFR1) and FGFR2 (ref. 44). *RET*, the gene encoding pralsetinib's primary target, was ranked 11th of 10,716 genes in the predicted list. Half of the genes encoding pralsetinib's targets (5 of 11) were ranked within the top 100 predicted targets by PDGrapher, including *KDR* (ranked at 3), *FLT3* (ranked at 10), *RET* (ranked at 11), *PDGFRB* (ranked at 14) and *FGFR2* (ranked at 81). This substantial overlap highlights the potential of the candidate targets identified by PDGrapher for pralsetinib-based lung cancer treatment, given that pralsetinib was not included in the training set of PDGrapher.

Next, we examined *KDR* as a therapeutic target for lung cancer. *KDR*, also known as *VEGFR2*, has been identified as a critical therapeutic target in A549 lung adenocarcinoma cells. These cells express *KDR* at both mRNA and protein levels, facilitating autocrine signalling that promotes tumour cell survival and proliferation²⁸. Activation of *KDR* enhances tumour angiogenesis and growth by upregulating oncogenic factors such as enhancer of zeste homologue 2 (*EZH2*), which is associated with increased cell proliferation and migration. Inhibiting *KDR* has showed promising therapeutic effects, including reduced cell proliferation and induced apoptosis. For instance, *KDR* inhibitors have been shown to decrease the malignant potential of lung adenocarcinoma cells by downregulating *EZH2* expression and increasing sensitivity to chemotherapy²⁹. These findings underscore the importance of *KDR* as a therapeutic target in A549 lung adenocarcinoma cells, highlighting its role in tumour progression and the potential benefits of its inhibition in cancer treatment strategies. Importantly, PDGrapher has successfully identified *KDR* among the top 20 predicted targets in chemical-PPI-lung-A549, validating its precision in detecting key therapeutic targets for lung cancer.

Given that Open Targets offers more comprehensive evidence for targets currently under development, we conducted a second series of case studies using Open Targets data to evaluate PDGrapher's capability to identify candidate therapeutic targets and drugs. This analysis aims to identify targets for lung cancer. Figure 6a presents a bubble graph that illustrates the union of the top 10 predicted targets to transition cell states from diseased to healthy in the six cell lines of three types of cancer that have available healthy controls. In the plot, the colour intensity and size of the bubbles represent the number of evidence sources and the association scores for each type of evidence. Most predicted targets are supported by drugs, pathology and systemic biology, and somatic mutation databases, which were considered strong evidence sources. Two unique targets, DNA topoisomerase II- α (*TOP2A*) and cyclin-dependent kinase 2 (*CDK2*), are predicted exclusively for the lung cancer cell line (Extended Data Fig. 8). *TOP2A* is ranked as the top predicted target by PDGrapher. This gene encodes a crucial decatenating enzyme that alters DNA topology by binding to two double-stranded DNA molecules, introducing a double-strand break, passing the intact strand through the break, and repairing the broken strand. This mechanism is vital for DNA replication and repair processes. *TOP2A* could be a potential therapeutic target for anti-metastatic therapy of NSCLC because it promotes metastasis of NSCLC by stimulating the canonical Wnt signalling pathway and inducing epithelial-mesenchymal transition⁴⁵. Using the predicted target of *TOP2A*, PDGrapher then identified three drugs, aldoxorubicin, vosaroxin and doxorubicin hydrochloride, as candidate drugs. These drugs were not part of the training dataset of PDGrapher and are in the early stages of clinical development: aldoxorubicin and vosaroxin are in phase II trials (ClinicalTrials.gov), and doxorubicin hydrochloride is in phase I trials but has been shown to improve survival in patients with metastatic or surgically unresectable uterine or soft tissue leiomyosarcoma⁴⁶.

Given that PDGrapher can rank all genes based on PPI network or GRN data, we assessed two questions: whether top-ranked genes have stronger evidence from Open Targets compared with lower-ranked genes, and what rank threshold should be used to identify reliably predicted genes. Figure 6b shows the number of sources of evidence and the global scores for the predicted target genes within the rank ranges of 1–10, 11–20, 51–60, 101–110 and 1,001–1,010 for lung cancer (chemical-PPI-lung-A549). The analysis revealed a clear trend: both the number of supporting evidence sources and global scores decrease with increasing rank, validating the predictive accuracy of PDGrapher. Most targets ranked within the top 100 have strong evidence from Open Targets, indicating that a rank threshold of 100 could serve as a cut-off for selecting candidate targets.

Training PDGrapher models

We conducted an ablation study to evaluate the components of PDGrapher's objective function using chemical datasets. We trained PDGrapher under three configurations: with only the cycle loss (PDGrapher-Cycle), only the supervision loss (PDGrapher-Super) and with both losses combined (PDGrapher-SuperCycle). The experiments were performed in the random splitting setting across all nine PPI chemical datasets. We assessed performance using several metrics, including the percentage of accurately predicted samples (Fig. 5c), nDCG (Supplementary Fig. 8a), recall values (Supplementary Fig. 8b) and strength of evidence (Extended Data Fig. 9). The results showed that PDGrapher-Super achieves the highest performance in predicting correct perturbagens but performs the worst in reconstructing treated samples. In contrast, PDGrapher-Cycle performs poorly in identifying correct perturbagens but shows improved performance in predicting (reconstructing) held-out treated samples. PDGrapher-SuperCycle (the configuration used throughout this study) strikes a balance between these two objectives, achieving competitive performance in predicting therapeutic genes while showing the best performance in reconstructing treated samples from diseased samples after intervening on the predicted genes. This makes PDGrapher-SuperCycle the most effective choice for balancing accuracy in perturbagen prediction with reconstruction fidelity. The findings show that supervision loss is essential for PDGrapher's overall performance. The PDGrapher-Cycle model consistently underperforms in all cell lines and metrics. Although PDGrapher-Super often excels in ranking performance, including cycle loss (in PDGrapher-SuperCycle) proves its value by moderately improving top prediction metrics such as recall@1 and recall@10. In addition, when healthy cell line data are available, the top predictions of PDGrapher-SuperCycle show stronger evidence compared with those of PDGrapher-Super in more than half (four of six) of the cell lines (Extended Data Fig. 9). We chose PDGrapher-SuperCycle for this work because it provides accurate target gene predictions from the top-ranked genes in the predicted list and bases its predictions on the changes they would induce in diseased samples.

Recognizing the role of biological pathways in disease phenotypes, PDGrapher-SuperCycle can identify alternative gene targets with close network proximity that may produce similar phenotypic outcomes. The organization of genes with similar functions, where each gene contributes to specific biochemical processes or signalling cascades, allows perturbations in different genes to yield analogous effects⁴⁷. This function-based interconnectivity implies that targeting different genes with similar functions can achieve therapeutic outcomes, as these genes collectively influence cellular phenotypic states⁴⁸. Although PDGrapher-SuperCycle shows slightly lower performance than PDGrapher-Super in pinpointing targets (Fig. 5c), it excels in identifying sets of gene targets capable of transitioning cell states from diseased to treated conditions (Supplementary Fig. 8b and Extended Data Fig. 9).

We conducted four analyses to test the sensitivity of PDGrapher to the causal graph. The first analysis uses five PPI networks constructed

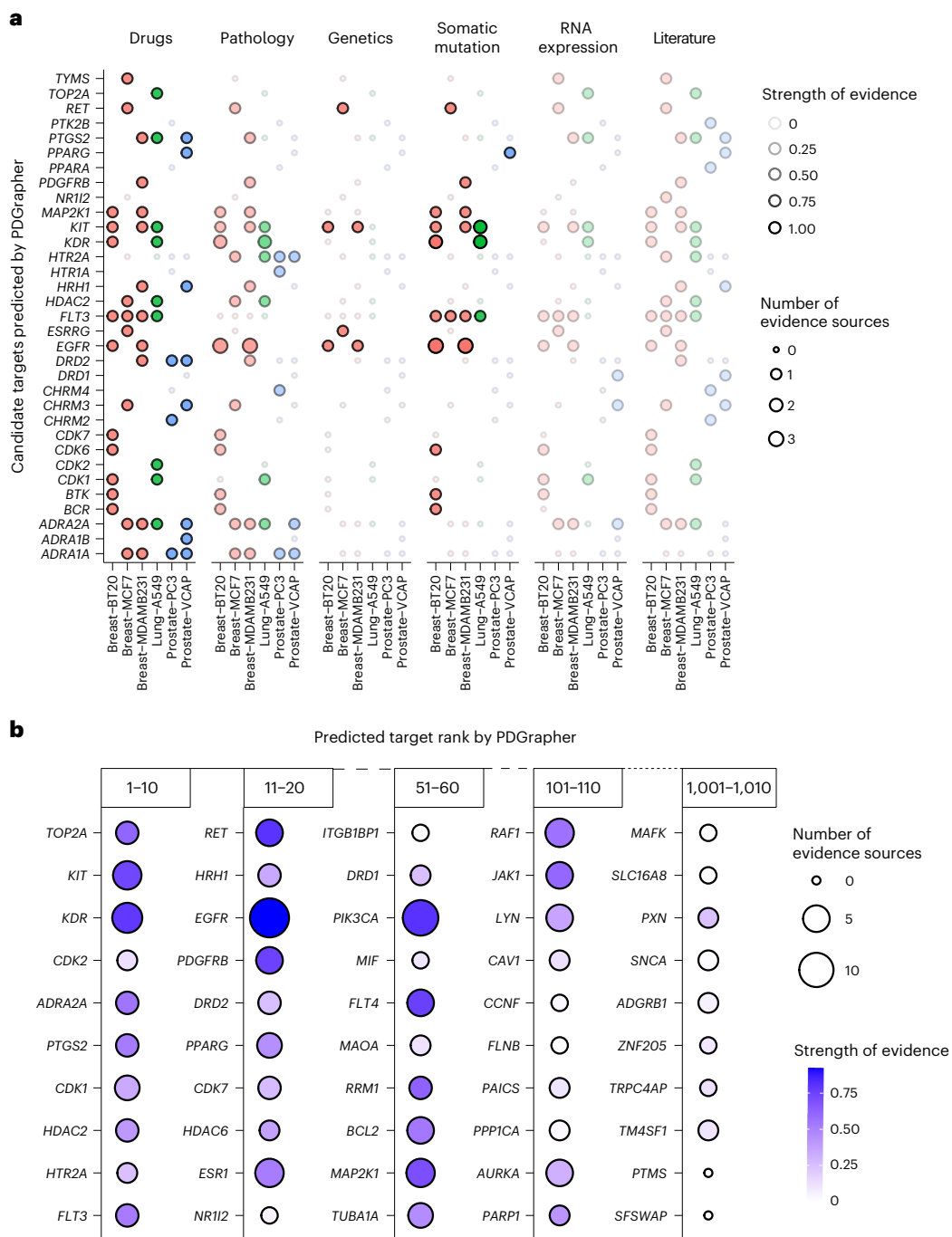


Fig. 6 | PDGrapher's prioritization of lung cancer targets is supported by Open Targets evidence. **a**, Union of the top 10 targets predicted by PDGrapher in lung, breast and prostate cancer. The colour intensity and size of the bubbles represent the number of evidence sources and the association scores for each type of evidence, respectively. Red, blue and green dots represent breast, lung and

prostate cancer, respectively. The details of the scoring system are provided in Supplementary Note 2. **b**, Predicted target rank from PDGrapher in 5 ranges, 1–10, 11–20, 51–60, 101–110, and 1,001–1,010 for lung cancer (chemical-PPI-lung-A549). The colour intensity and size of the bubbles represent the number of evidence sources and the global scores of targets from Open Targets, respectively.

with varying edge confidence cut-offs. The PPI network was obtained from STRING (<https://string-db.org/>)⁴⁹, which assigns a confidence score to each edge. To create networks with different levels of confidence, we filtered edges based on the quantiles 0.1, 0.2, 0.3, 0.4 and 0.5 of the confidence scores, resulting in 5 networks with decreasing numbers of edges. For this analysis, we selected two cell lines: chemical-PPI-breast-MDAMB231 and genetic-PPI-breast-MCF7. The LINCS perturbation data for each cell line were processed using the five PPI networks (Supplementary Note 3). We trained PDGrapher

with one, two and three GNN layers, selecting the best configuration based on the performance of the validation set. As shown in Fig. 5b and Extended Data Fig. 3, PDGrapher performs robustly at all levels of confidence in PPI networks. It maintains stable performance on both chemical and genetic intervention datasets, even as an increasing number of edges is removed from the PPI networks. The second to fourth analyses are based on the synthetic graphs. We created two sparse PPI networks using different edge removal strategies and one synthetic gene expression dataset with increasing levels of latent confounders.

We applied two edge removal strategies to the PPI network: removing increasing numbers of either bridge edges or random edges. Details of the data generation process are provided in the Methods. Results from the edge removal experiments indicate that although bridge edges are structurally critical, their limited number in the PPI graph reduces their overall impact on model predictions (Extended Data Fig. 5). In contrast, the removal of random edges, which include both high-confidence and redundant connections, has a more pronounced effect on performance, highlighting the model's sensitivity to network perturbations (Extended Data Fig. 6). The fourth dataset introduces latent confounders in the gene expression data. PDGrapher showed stable performance in perturbagen prediction, with only a slight decrease in performance as stronger confounders were introduced (Extended Data Fig. 4).

We then evaluated whether PDGrapher can maintain robust performance in the absence of disease intervention data. In our training datasets, some cell lines lacked associated healthy control samples from disease-relevant tissues and cell types. These cell lines contained only treatment intervention data (diseased cell state, perturbagen and treated cell state) without disease intervention data (healthy cell state, disease mutations or diseased cell state) for model training and inference. For cell lines with healthy controls, we trained the response prediction module using both intervention datasets. For cell lines without healthy controls, we trained PDGrapher using only treatment intervention data. To evaluate PDGrapher's dependency on healthy control data, we trained the model on cell lines with available disease intervention data under two conditions: one using the disease intervention data for training and one excluding it. This evaluation was conducted on six chemical perturbation datasets (chemical-PPI-lung-A549, chemical-PPI-breast-MCF7, chemical-PPI-breast-MDAMB231, chemical-PPI-breast-BT20, chemical-PPI-prostate-PC3 and chemical-PPI-prostate-VCAP) and three genetic perturbation datasets (genetic-PPI-lung-A549, genetic-PPI-breast-MCF7 and genetic-PPI-prostate-PC3). The results indicated that the two versions of PDGrapher perform consistently across cell types and data types (chemical and genetic; Fig. 5d and Supplementary Figs. 6 and 7). In half of the cell lines (four of nine), the model trained without disease intervention data outperformed the model trained with it. This shows that PDGrapher has a weak dependency on healthy control data and can perform well even when such data are unavailable.

Discussion

We formulate phenotype-driven lead discovery as a combinatorial prediction problem for therapeutic targets. Given a diseased sample, the goal is to identify genes that a genetic or chemical perturbagen should target to reverse disease effects and shift the sample towards a treated state that matches the distribution of a healthy state. This requires predicting a combination of gene targets, framing the task as combinatorial prediction. To address this, we introduce PDGrapher. Using a diseased cell state represented by a gene expression signature and a proxy causal graph of gene–gene interactions, PDGrapher predicts candidate target genes to transition cells to the desired treated state. PDGrapher includes two modules: a perturbagen discovery module that proposes a set of therapeutic targets based on the diseased and treated states, and a response prediction module that evaluates the effect of applying the predicted perturbagen to the diseased state. Both modules are GNN models that operate on gene–gene networks, which serve as approximations of noisy causal graphs. We use PPI networks and GRNs as two representations of these noisy causal graphs. PDGrapher predicts perturbagens that shift gene expression from diseased to treated states across 2 evaluation settings (random and leave-cell-out) and 19 datasets involving genetic and chemical interventions. Unlike alternative response prediction methods, which rely on indirect prediction to identify perturbagens, PDGrapher selects candidate gene targets to achieve the desired transformation^{36,37,50–52}.

PDGrapher has the potential to improve therapeutic lead design and expand the search space for perturbagens. It leverages large datasets of genetic and chemical interventions to identify sets of candidate targets that can shift cell line gene expression from diseased to treated states. By selecting sets of therapeutic targets for intervention instead of a single perturbagen, PDGrapher enhances phenotype-driven lead discovery. PDGrapher's approach to identifying therapeutic targets can enable personalized therapies by tailoring treatments to individual gene expression profiles. Its ability to output multiple genes is particularly relevant for diseases where dependencies among several genes affect treatment efficacy and safety.

PDGrapher operates under the assumption that there are no unobserved confounders, a stringent condition that is challenging to validate empirically. Future work could focus on re-evaluating and relaxing this assumption. Another limitation lies in the reliance on PPI networks and GRNs as proxies for causal gene networks, as these networks are inherently noisy and incomplete^{53–55}. PDGrapher posits that representation learning can overcome incomplete causal graph approximations. A valuable research direction is to theoretically examine the impact of such approximations, focusing on how they influence the accuracy and reliability of predicted likelihoods. Such analyses could uncover high-level causal variables with therapeutic effects from low-level observations and contribute to reconciling structural causality and representation learning approaches, which generally lack any causal understanding⁵⁶. We performed two experiments to evaluate the robustness of PDGrapher. First, we tested PDGrapher on a PPI network with weighted edges, progressively removing low-confidence edges to assess its performance under increasing network sparsity. Second, we applied PDGrapher to synthetic datasets with varying levels of missing graph components and confounding factors in gene expression data. In both experiments, PDGrapher maintained stable performance. PDGrapher also showed robust performance across PPI networks with different numbers of edges.

Phenotype-driven drug discovery using PDGrapher faces certain limitations, one of which is its reliance on transcriptomic data. Although transcriptomics is broadly applicable, including other data modalities, such as cell morphology screens, could produce more comprehensive models. Cell morphology screens, including cell painting, capture cellular responses by staining organelles and cytoskeletal components, generating image profiles that capture the effects of genetic or chemical perturbations^{57,58}. These screens allow identification of phenotypic signatures that correlate with compound activity, mechanisms of action and potential off-target effects. The recent release of the JUMP Cell Painting dataset⁵⁹ exemplifies how high-content morphological profiling can complement databases such as CMap and LINCS, creating integrated datasets for phenotype-driven discovery. By integrating multimodal data, including phenotypic layers from transcriptomic and image data, it becomes possible to uncover more comprehensive patterns of compound effects⁶⁰. Such integration would broaden the scope of PDGrapher, allowing it to capture wider mechanistic insights and support more effective therapeutic discovery^{61,62}.

A limitation of our study is the use of NL20 as a control cell line for A549 (refs. 63–65). Although NL20 is a normal human bronchial epithelial cell line and A549 is a human lung carcinoma cell line derived from the alveolar region, the two cell lines differ in anatomical origin and molecular characteristics. This mismatch could introduce biases in comparative analyses due to variations in baseline gene expression profiles and cellular behaviours. To mitigate this concern, we evaluate PDGrapher's performance across datasets with and without healthy control data. PDGrapher performs consistently regardless of the inclusion of healthy controls, indicating that its predictions are robust to the absence of matched control cells. Ablation analyses showed that incorporating cycle loss improved PDGrapher's performance in top target predictions for five of nine cell lines. On the basis of this improvement, we included the cycle loss in all experiments. Cycle loss helps

maintain the robustness and biological relevance of model predictions. PDGrapher learns to predict drug targets that shift cells from a diseased state to a healthy or treated state. It then uses the diseased gene expression profile and the predicted targets to estimate the gene expression after treatment. This bidirectional approach enforces the fidelity of predicted targets as they must contain sufficient information to reconstruct state B from state A. Cycle loss also serves as a regularizer that penalizes discrepancies between the original input and its reconstruction⁶⁶.

PDGrapher is a GNN approach for combinatorial prediction of perturbations that transition diseased cells to treated states. By leveraging causal reasoning and representation learning on gene networks, PDGrapher identifies perturbation necessary to achieve specific phenotypic changes. This approach enables the direct prediction of therapeutic targets that can reverse disease phenotypes, bypassing the need for exhaustive response simulations across large perturbation libraries. Its design and evaluation lay the groundwork for future advances in phenotype-based modeling of therapeutic perturbations by improving the precision and scalability of perturbation prediction methods.

Methods

Preliminaries

A calligraphic letter x indicates a set, an italic uppercase letter X denotes a graph, uppercase \mathbf{X} denotes a matrix, lowercase \mathbf{x} denotes a vector, and a monospaced letter \mathbf{x} indicates a tuple. Uppercase letter X indicates a random variable, and lowercase letter x indicates its corresponding value; bold uppercase \mathbf{X} denotes a set of random variables, and lowercase letter \mathbf{x} indicates its corresponding values. We denote $P(\mathbf{X})$ as a probability distribution over a set of random variables \mathbf{X} and $P(\mathbf{X} = \mathbf{x})$ as the probability of \mathbf{X} that is equal to the value of \mathbf{x} under the distribution $P(\mathbf{X})$. For simplicity, $P(\mathbf{X} = \mathbf{x})$ is abbreviated as $P(\mathbf{x})$. This section uses terminology and concepts from the framework of causal inference⁶⁷.

Problem formulation for combinatorial prediction of targets

Intuitively, given a diseased cell line sample, we would like to predict the set of therapeutic genes that need to be targeted to reverse the effects of disease, that is, the genes that need to be perturbed to shift the cell gene expression state as close as possible to the healthy state. Next, we formalize our problem formulation. Let $\mathbb{M} = \langle \mathbf{E}, \mathbf{V}, \mathcal{F}, P(\mathbf{E}) \rangle$ be a structural causal model (SCM; see the description of related works in Supplementary Note 4) associated with causal graph G , where \mathbf{E} is a set of exogenous variables affecting the system, \mathbf{V} are the system variables, \mathcal{F} are structural equations encoding causal relations between variables and $P(\mathbf{E})$ is a probability distribution over exogenous variables. Let $\mathcal{T} = \{\tau_1, \dots, \tau_m\}$ be a dataset of paired healthy and diseased samples (namely, disease intervention data), where each element is a triplet $\tau = \langle \mathbf{v}^h, \mathbf{U}, \mathbf{v}^d \rangle$ with $\mathbf{v}^h \in [0, 1]^N$ being normalized gene expression values of a healthy cell line (variable states before perturbation), \mathbf{V}_U being the disease-causing perturbed variable (gene) set in \mathbf{V} , and $\mathbf{v}^d \in [0, 1]^N$ being gene expression values of a diseased cell line (variable states after perturbation). Our goal is to find, for each sample $\tau = \langle \mathbf{v}^h, \mathbf{U}, \mathbf{v}^d \rangle$, the variable set \mathbf{U}' with the highest likelihood of shifting variable states from diseased \mathbf{v}^d to healthy \mathbf{v}^h state. To increase generality, we refer to the desired variable states as treated (\mathbf{v}^t). Our goal can then be expressed as:

$$\operatorname{argmax}_{\mathbf{U}'} P^{G^{\mathbf{U}'}}(\mathbf{V} = \mathbf{v}^t \mid \operatorname{do}(\mathbf{U}')), \quad (1)$$

where $P^{G^{\mathbf{U}'}}$ represents the probability on the graph G mutilated by perturbations in variables in \mathbf{U} . Under the assumption of no unobserved confounders, the above interventional probability can be expressed as a conditional probability on the mutilated graph $G^{\mathbf{U}'}$:

$$\operatorname{argmax}_{\mathbf{U}'} P^{G^{\mathbf{U}'}}(\mathbf{V} = \mathbf{v}^t \mid \mathbf{U}'), \quad (2)$$

which under the causal Markov condition is:

$$\operatorname{argmax}_{\mathbf{U}'} \prod_i P(\mathbf{v}_i = \mathbf{v}_i^t \mid \mathbf{Pa}_{\mathbf{v}_i}), \quad (3)$$

where $\mathbf{Pa}_{\mathbf{v}_i}$ represents parents of variable \mathbf{v}_i according to graph $G^{\mathbf{U}'}$ (that is, the mutilated graph upon intervening on variables in \mathbf{U}'). Here state of a variable $\mathbf{v}_j \in \mathbf{Pa}_{\mathbf{v}_i}$ will be equal to an arbitrary value \mathbf{v}_j' if $\mathbf{v}_j \in \mathbf{U}'$. Therefore, intervening on the variable set \mathbf{U}' modifies the graph used to obtain conditional probabilities and determine the state of variables in \mathbf{U}' .

Problem formulation from a representation learning perspective

In the previous section, we drew on the SCM framework to introduce a generic formulation for the task of combinatorial prediction of therapeutic targets. Instead of approaching the problem from a purely causal inference perspective, we draw upon representation learning to approximate the queries of interest to address the limiting real-world setting of a noisy and incomplete causal graph. Formulating our problem using the SCM framework allows for explicit modelling of interventions and formulation of interventional queries (see the description of related works in Supplementary Note 4). Inspired by this principled problem formulation, we next introduce the problem formulation using a representation learning paradigm.

We let $G = (\mathcal{V}, \mathcal{E})$ denote a graph with $|\mathcal{V}| = n$ nodes and $|\mathcal{E}|$ edges, which contains partial information on causal relationships between nodes in \mathcal{V} and some noisy relationships. We refer to this graph as a proxy causal graph. Let $\mathcal{T} = \{\tau_1, \dots, \tau_m\}$ be a dataset with an individual sample being a triplet $\tau = \langle \mathbf{x}^h, \mathbf{U}, \mathbf{x}^d \rangle$ with $\mathbf{x}^h \in [0, 1]^n$ being the node states (attributes) of a healthy cell sample (before perturbation), \mathbf{U} being the set of disease-causing perturbed nodes in \mathcal{V} , and $\mathbf{x}^d \in [0, 1]^n$ being the node states (attributes) of a diseased cell sample (after perturbation). We denote by $G^{\mathbf{U}} = (\mathcal{V}, \mathcal{E}^{\mathbf{U}})$ the graph resulting from the mutilation of edges in G as a result of perturbing nodes in \mathbf{U} (one graph per perturbation; we avoid using superindices for simplicity). Here again we refer to the desired variable states as treated (\mathbf{x}^t). Our goal is then to learn a function:

$$f : G^{\mathbf{U}'}, \mathbf{x}^d, \mathbf{x}^t \rightarrow \operatorname{argmax}_{\mathbf{U}'} P^{G^{\mathbf{U}'}}(\mathbf{x} = \mathbf{x}^t \mid \mathbf{x}^d, \mathbf{U}'). \quad (4)$$

That, given the graph $G^{\mathbf{U}'}$, the diseased node states \mathbf{x}^d and treated node states \mathbf{x}^t , predicts the combinatorial set of nodes \mathbf{U}' that if perturbed have the highest chance of shifting the node states to the treated state \mathbf{x}^t . We note here that $P^{G^{\mathbf{U}'}}$ represents probabilities over graph $G^{\mathbf{U}'}$ mutilated upon perturbations in nodes in \mathbf{U}' . Under causal Markov condition, we can factorize $P^{G^{\mathbf{U}'}}$ over graph $G^{\mathbf{U}'}$:

$$f : G^{\mathbf{U}'}, \mathbf{x}^d, \mathbf{x}^t \rightarrow \operatorname{argmax}_{\mathbf{U}'} \prod_i P(\mathbf{x}_i = \mathbf{x}_i^t \mid \mathbf{x}_{\mathcal{P}(\mathbf{v}_i)}), \quad (5)$$

that is, the probability of each node i depending only on its parents $\mathcal{P}(\mathbf{v}_i)$ in graph $G^{\mathbf{U}'}$.

We assume (1) real-valued node states, (2) G is fixed and given, and (3) atomic and non-atomic perturbagens (intervening on individual nodes or sets of nodes). Given that the value of each node should depend only on its parents in the graph $G^{\mathbf{U}'}$, a message-passing framework appears especially suited to compute the factorized probabilities P .

In the SCM framework, the conditional probabilities in equation (3) are computed recursively on the graph, each being an expectation over exogenous variables \mathbf{E} . Therefore, node states of the previous time point are not necessary. To translate this query into the representation learning realm, we discard the existence of noise variables and directly try to learn a function encoding the transition from an initial state to a desired state. An exhaustive approach to solving equation (5)

would be to search the space of all potential sets of therapeutic targets u' and score how effective each is in achieving the desired treated state. This is, how many cell response prediction approaches can be used for perturbagen discovery^{22,23,68}. However, with moderately sized graphs, this is highly computationally expensive, if not intractable. Instead, we propose to search for potential perturbagens efficiently with a perturbagen discovery module (f_p) and a way to score each potential perturbagen with a response prediction module (f_r).

Relationship to conventional graph prediction tasks

Given that the prediction for each variable is dependent only on its parents in a graph, GNNs appear especially suited for this problem. We can formulate the query of interest under a graph representation learning paradigm as follows: given a graph $G = (\mathcal{V}, \mathcal{E})$, paired sets of node attributes $x = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m\}$ and node labels $y = \{\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_m\}$, where each $\mathbf{Y} = \{y_1, \dots, y_n\}$, with $y_i \in [0, 1]$, we aim at training a neural message-passing architecture that given node attributes \mathbf{X}_i predicts the corresponding node labels \mathbf{Y}_i . There are, however, differences between our problem formulation and the conventional graph prediction tasks, namely, graph and node classification (summarized in Supplementary Table 13).

In node classification, a single graph G is paired with node attributes \mathbf{X} , and the task is to predict the node labels \mathbf{Y} . Our formulation differs in that we have m paired sets of node attributes x and labels y instead of a single set, yet they are similar in that there is a single graph in which GNNs operate. In graph classification, a set of graphs $\mathcal{G} = \{G_1, \dots, G_m\}$ is paired with a set of node attributes $x = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m\}$ and the task is to predict a label for each graph $\mathbf{Y} = \{y_1, \dots, y_m\}$. Here graphs have a varying structure, and both the topological information and node attributes predict graph labels. In our formulation, a single graph is combined with each node attribute \mathbf{X}_i , and the goal is to predict a label for each node, not for the whole graph.

PDGrapher model

PDGrapher is an approach for combinatorial prediction of therapeutic targets composed of two modules. First, a perturbagen discovery module f_p searches the space of potential gene sets to predict a suitable candidate u' . Next, a response prediction module f_r checks the goodness of the predicted set u' , that is, how effective intervening on variables in u' is to shift node states to the desired treated state \mathbf{x}^t .

$$(1) \mathbf{x}^d, \mathbf{x}^t \xrightarrow{f_p} u'$$

$$(2) \mathbf{x}^d, u' \xrightarrow{f_r} \hat{\mathbf{x}}^t$$

Model optimization

We optimize our response prediction module f_r using cross-entropy (CE) loss on known triplets of disease intervention $\langle \mathbf{x}^d, u, \mathbf{x}^d \rangle$ and treatment intervention $\langle \mathbf{x}^d, u', \mathbf{x}^t \rangle$:

$$\mathcal{L}_{f_r} = \text{CE}(\mathbf{x}^d, f_r(\mathbf{x}^d, u)) + \text{CE}(\mathbf{x}^t, f_r(\mathbf{x}^d, u')). \tag{6}$$

We optimize our intervention discovery module f_p using a cycle loss, ensuring that the response to the predicted intervention set u' closely matches the desired treated state (the first part of equation (7)). In addition, we provide a supervisory signal for predicting u' in the form of cross-entropy loss (the second part of equation (7)). So, the total loss is defined as:

$$\mathcal{L}_{f_p} = \text{CE}(\mathbf{x}^t, f_r(\mathbf{x}^d, f_p(\mathbf{x}^d, \mathbf{x}^t))) + \text{CE}(u', f_p(\mathbf{x}^d, \mathbf{x}^t)) \text{ (with } f_r \text{ frozen)}. \tag{7}$$

We train f_p and f_r in parallel and implement early stopping separately (see ‘Experimental set-up’ for more details). Trained module f_p is then

used to predict, for each diseased cell sample, which nodes should be perturbed (u') to achieve a desired treated state (Fig. 1a).

Response prediction module

Our response prediction module f_r should learn to map pre-perturbagen node values to post-perturbagen node values through learning relationships between connected nodes (equivalent to learning structural equations in SCMs) and propagating the effects of perturbations downstream in the graph (analogous to the recursive nature of query computations in SCMs).

Given a disease intervention triplet $\langle \mathbf{x}^d, u, \mathbf{x}^d \rangle$, we propose a neural model operating on a mutilated graph, G^u , where the node attributes are the concatenation of \mathbf{x}^d and \mathbf{x}'_u , predicting diseased node values \mathbf{x}^d . The first element is its gene expression value \mathbf{x}^d_i and the second is a perturbation flag, a binary label indicating whether a perturbation occurs at node i . So, each node i has a two-dimensional attribute vector $\mathbf{d}_i = [\mathbf{x}^d_i \parallel \mathbf{x}'_u]$. In practice, we embed each node feature into a high-dimensional continuous space by assigning learnable embeddings to each node based on the value of each input feature dimension. Specifically, for each node, we use the binary perturbation flag to assign a d -dimensional learnable embedding, which is different between nodes but shared across samples for each node. To embed the gene expression value $\mathbf{x}^d_i \in [0, 1]$, we first calculate thresholds using quantiles to assign the gene expression value into one of the B bins. We use the bin index to assign a d -dimensional learnable embedding, which is different between nodes but shared across samples for each node. To increase our model’s representation power, we concatenate a d -dimensional positional embedding (a d -dimensional vector initialized randomly following a normal distribution). Concatenating these three embeddings results in an input node representation of dimensionality $3d$. For each node $i \in \mathcal{V}$, an embedding \mathbf{z}_i is computed using a GNN operating on the node’s neighbours’ attributes. The most general formulation of a GNN layer is:

$$\mathbf{h}'_i = \phi \left(\mathbf{h}_i, \bigoplus_{j \in \mathcal{N}^i} \psi(\mathbf{h}_j, \mathbf{h}_j) \right), \tag{8}$$

where \mathbf{h}'_i represents the updated information of node i , and \mathbf{h}_i represents the information of node i in the previous layer, with embedded \mathbf{d}_i being the input to the first layer. ψ is a message function, \bigoplus a permutation-invariant aggregate function, and ϕ is an update function. We obtain an embedding \mathbf{z}_i for node i by stacking K GNN layers. The node embedding $\mathbf{z}_i \in \mathbb{R}$ is then passed to a multilayer feedforward neural network to obtain an estimate of the values of the post-perturbation nodes \mathbf{x}^d .

Perturbation discovery module

Our perturbagen prediction module f_p should learn the nodes in the graph that should be perturbed to shift the node states (attributes) from diseased \mathbf{x}^d to the desired treated state \mathbf{x}^t . Given a triplet $\langle \mathbf{x}^d, u', \mathbf{x}^t \rangle$, we propose a neural model operating on graph G^u with node features \mathbf{x}^d and \mathbf{x}^t that predicts a ranking for each node, where the top P ranked nodes should be predicted as the nodes in u' . Each node i has a two-dimensional attribute vector: $\mathbf{d}_i = [\mathbf{x}^d_i \parallel \mathbf{x}^t_i]$. In practice, we represent these binary features in a continuous space using the same approach as described for our response prediction module f_r .

For each node $i \in \mathcal{V}$, an embedding \mathbf{z}_i is computed using a GNN operating on the node’s neighbours’ attributes. We obtain an embedding \mathbf{z}_i for node i by stacking K GNN layers. The node embedding $\mathbf{z}_i \in \mathbb{R}$ is then passed to a multilayer feedforward neural network to predict a real-valued number for node i .

Model implementation and training

We implement PDGrapher using PyTorch 1.10.1 (ref. 69) and the Torch Geometric 2.0.4 Library⁷⁰. The implemented architecture yields a

neural network with the following hyperparameters: number of GNN layers and number of prediction layers. We set the number of prediction layers to two and performed a grid search over the number of GNN layers (one to three layers). We train our model using a 5-fold cross-validation strategy and report PDGrapher's performance resulting from the best-performing hyperparameter setting.

Further details on statistical analysis

We next outline the evaluation set-up, baseline models and statistical tests used to evaluate PDGrapher. We evaluate the performance of PDGrapher against the following existing methods:

- **Random reference:** Given a sample $\tau = \langle \mathbf{x}^d, u', \mathbf{x}^t \rangle$, the random reference baseline returns N random genes as the prediction of target genes in u' , where N is the number of genes in u' .
- **Cancer genes:** Given a sample $\tau = \langle \mathbf{x}^d, u', \mathbf{x}^t \rangle$, the cancer genes baseline returns the top N genes from an ordered list where the first M genes are disease associated (cancer-driver genes). The remaining genes are ranked randomly, and N is the number of genes in u' . The processing of cancer genes is described in 'Disease-genes information' in Supplementary Note 3.
- **Cancer drug targets:** Given a sample $\tau = \langle \mathbf{x}^d, u', \mathbf{x}^t \rangle$, the cancer targets baseline returns the top N genes from an ordered list where the first M genes are cancer drug targets and the remaining genes are ranked randomly, and N is the number of genes in u' . The processing of drug target information is described in 'Drug-targets information' and 'Cancer drug and target information' in Supplementary Note 3.
- **Perturbed genes:** Given a sample $\tau = \langle \mathbf{x}^d, u', \mathbf{x}^t \rangle$, the perturbed genes baseline returns the top N genes from an ordered list where the first M genes are all perturbed genes in the training set and the remaining genes are ranked randomly, and N is the number of genes in u' .
- **scGen²²:** scGen is a widely used gold-standard latent variable model for response prediction⁷¹⁻⁷⁴. Given a set of observed cell types in control and perturbed states, scGen predicts the response of a new cell type to the pertubagen seen in training. To use scGen as a baseline, we first fit it to our LINCS gene expression data for each dataset type to predict response to perturbagens, training one model per pertubagen (chemical or genetic). Then, given a sample of paired diseased-treated cell line states, $\tau = \langle \mathbf{x}^d, u', \mathbf{x}^t \rangle$, we compute the response of the cell line with gene expression \mathbf{x}^d to all perturbagens. The predicted pertubagen is that whose predicted response is closest to \mathbf{x}^t in R^2 score, which quantifies the proportion of variance in treated state explained by the prediction. As scGen trains one model per pertubagen, it needs an exhaustively long training time for datasets with a large number of perturbagens, especially in the leave-cell-out setting. Therefore, we set the maximum training epochs to 100 and only conducted leave-cell-out tests for one split of data for scGen.
- **Biolord³³:** Biolord can predict pertubagen response for both chemical and genetic datasets. We followed the official tutorial from the Biolord GitHub repository (<https://github.com/nitzanlab/biolord>), using the recommended parameters. To prevent memory and quota errors, we implemented two filtering steps: (1) instead of storing the entire response gene expression (rGEX) matrix of all input (control) cells for each pertubagen, we only store a vector of the averaged rGEX of the input cells per pertubagen, which is necessary for calculating R^2 for evaluation; and (2) during prediction, if the number of control cells exceeds 10,000, we randomly downsample the control cells to 10,000. Similar to scGen, we predict the responses gene expression \mathbf{x}^d for all perturbagens and use them to calculate R^2 to get the rank of predicted perturbagens.

- **ChemCPA²³:** ChemCPA is specifically designed for chemical perturbation. We followed the official tutorials on GitHub for running this model (<https://github.com/theislab/chemCPA>), with all parameters set following the authors' recommendations. Data processing was also conducted using the provided scripts. We constructed drug embedding using RDKit with canonical SMILES sequences, as this is the default setting in the model and the tutorial. As the original ChemCPA model lacks functionality to obtain the predicted rGEX for each drug (averaging over the dosages), we developed a custom script to perform this task. These predictions, \mathbf{x}^d , were subsequently used for calculating R^2 to get the rank of predicted perturbagens.
- **GEARS³⁴:** GEARS is capable of predicting pertubagen responses for genetic perturbation datasets, specifically for predicting the rGEX to unseen perturbagens. However, it is limited to predicting only those genes that are present in the gene network used as prior knowledge for model training. In addition, GEARS cannot process perturbagens with only one sample, so we filtered the data accordingly. We followed the official tutorial from the GEARS GitHub repository (<https://github.com/snap-stanford/GEARS>), using the recommended parameters. After confirming with the authors, we established that GEARS is suitable only for within-cell-line prediction. Consequently, our experiments with GEARS were conducted exclusively within this scenario.
- **CellOT²⁷:** CellOT is capable of working with both chemical and genetic datasets. We ran this model by following the official tutorial from GitHub (<https://github.com/bunnech/cellot>), ensuring that all parameters were set according to the provided guidelines. Due to CellOT's limitation in processing perturbagens with small sample sizes, we filtered the data to retain only those perturbagens with more than five samples or cells. We then used the predicted rGEX \mathbf{x}^d to calculate R^2 and the predicted pertubagen ranks. Similar to scGen, CellOT trains one model per pertubagen, which results in an exhaustively long training time for datasets with a large number of perturbagens. This issue becomes even more pronounced when doing leave-cell-out evaluations. Therefore, for this method, we set the maximum training epochs to 100 and only conduct one split in leave-cell-out tests.

Dataset splits and evaluation settings

We evaluate PDGrapher and competing methods on two different settings.

Systematic random dataset splits. For each cell line, the dataset is split randomly into train and test sets to measure our model performance in an independent and identically distributed setting.

Leave-cell-out dataset splits. To measure model performance on unseen cell lines, we train our model with random splits on one cell line and test on a new cell line. Specifically, for chemical perturbation data, we train a model for each random split per cell line and test it on the entire dataset of the remaining eight cell lines. For genetic data, we train a model for each random split per cell line and test it on the entire dataset of the remaining nine cell lines. For example, with nine cell lines with chemical perturbation (A549, MDAMB231, BT20, VCAP, MCF7, PC3, A375, HT29 and HELA), we conducted an experiment where each split of cell line A549 was used as the training set, and the trained model was tested on the remaining eight cell lines (MDAMB231, BT20, VCAP, MCF7, PC3, A375, HT29 and HELA). Similarly, for cell line MDAMB231, we trained the model on each split of it and tested the model on the other eight cell lines (A549, BT20, VCAP, MCF7, PC3, A375, HT29 and HELA). This process was repeated for all cell lines, providing a comprehensive evaluation of PDGrapher and all competing methods.

Evaluation set-up

For all dataset split settings, our model is trained using 5-fold cross-validation, and metrics are reported as the average on the test set. Within each fold, we further split the training set into training and validation sets (8:2) to perform early stopping. We train the model on the training set until the validation loss has not decreased at least 10^{-5} for 15 continuous epochs.

Evaluation metrics

We report average sample-wise R^2 score and average perturbagen-wise R^2 score to measure performance in the prediction of \mathbf{x}^t . The sample-wise R^2 score is computed as the square of the Pearson correlation between the predicted sample $\mathbf{x}^t \in \mathbb{R}^N$ and real sample $\mathbf{x}^t \in \mathbb{R}^N$. The perturbagen-wise R^2 score is adopted from scGen. It is computed as the square of the Pearson correlation of a linear least-squares regression between a set of predicted treated samples $\mathbf{X}^t \in \mathbb{R}^{N \times S}$ and a set of real treated samples $\mathbf{X}^t \in \mathbb{R}^{N \times S}$ for the same perturbagen. Here, S indicates the size of the sets. Higher values indicate better performance in predicting the treated sample \mathbf{x}^t given the diseased sample \mathbf{x}^d and predicted perturbagen. This is used for evaluating the performance of response prediction. For evaluating perturbagen discovery, when the competing methods cannot predict perturbagen ranks for chemical perturbation data, we first calculate the rank of drugs based on the R^2 score. We then build a target gene rank from the drug rank by substituting the drugs with their target genes acquired from DrugBank⁷⁵ (accessed in November 2022; see details in Supplementary Note 3). A single drug can have multiple target genes, which we place in the rank in random order. As some methods cannot predict unseen drugs, their predicted target gene lists are often short, introducing bias in evaluation. To address this, we shuffle the missing target genes and attach them to the predicted ranks to create a complete rank. For genetic perturbation data, we directly obtain the target gene rank from the results and then attach the shuffled missing genes to the rank.

To evaluate the performance of our model in ranking predicted therapeutic targets, we use the nDCG, a widely used metric in information retrieval adapted for our setting. The raw DCG score is computed by summing the relevance of each correct target based on its rank in the predicted list, with relevance weighted by a logarithmic discount factor to prioritize higher-ranked interventions. The gain function is defined as $1 - \text{ranking}/N$, ensuring that the score reflects the quality of the ranking relative to the total number of nodes in the system. To ensure comparability across datasets or experiments with different numbers of correct interventions, DCG is normalized by the ideal DCG, which represents the maximum possible score for a perfect ranking. This results in nDCG values in the range $[0, 1]$, where higher values indicate better ranking performance and alignment with the ground truth. This metric is particularly suited for our task as it emphasizes the accuracy of top-ranked interventions while accounting for the diminishing importance of lower-ranked predictions.

In addition, we report the proportion of test samples for which the predicted therapeutic targets set has at least one overlapping gene with the ground-truth therapeutic targets set (denoted as the percentage of accurately predicted samples). We also calculated the ratio of correct therapeutic targets that appeared in the top 1, top 10 and top 100 predicted therapeutic targets in the predicted rank, denoted as $\text{recall}@1$, $\text{recall}@10$ and $\text{recall}@100$, respectively.

To assess the overall performance across all experiments and metrics, we calculated an aggregated metric, averaging all metric values for each method.

Statistical tests

In the benchmarking experiments, we performed a one-tailed pairwise t -test to evaluate whether PDGrapher significantly outperforms the competing methods. For other experiments, such as ablation studies, we used a two-tailed t -test to determine whether there is a significant

difference in performance between the two models. A significance threshold of 0.05 was used for all tests. P values of perturbagen discovery and response prediction tests are presented in the Source data.

Ablation studies

In the ablation study, we evaluated PDGrapher by optimizing it with only the supervision loss (PDGrapher-Super) and with only the cycle loss (PDGrapher-Cycle) across all chemical datasets. We then compared the perturbagen prediction performance of these submodels with that of PDGrapher (PDGrapher-SuperCycle). To train PDGrapher-Super and PDGrapher-Cycle, for each cell line, we set the number of layers to that which was found optimal for the validation set in the random splitting setting for PDGrapher-SuperCycle.

Sensitivity studies

To test the sensitivity of PDGrapher on PPI networks, we used data from STRING (string-db.org), which provides a confidence score for each edge. The method for acquiring and preprocessing the PPI networks from STRING is detailed in Supplementary Note 3. For the sensitivity tests, we selected two cell lines: the chemical dataset MDAMB231 and the genetic dataset MCF7. For each cell line, we processed the data using the five PPI networks described in Supplementary Note 3. We optimized PDGrapher using 5-fold cross-validation as described in 'Evaluation set-up' and optimized the number of GNN layers using the validation set in each split.

Synthetic datasets

We generated three synthetic datasets:

1. Dataset with missing components removing bridge edges: this dataset is generated by progressively removing bridge edges from the existing PPI network. Bridge edges are those whose removal disconnects parts of the network. We vary the fraction of bridge edges removed in increments (from zero to one) and, for each fraction, we create a new edge list representing the modified network (Supplementary Table 5). This process ensures that different levels of network sparsity are introduced, affecting the overall structure and connectivity. We pair these networks with gene expression data from chemical-PPI-breast-MDAMB231.
2. Dataset with missing components removing random edges: this dataset is generated by progressively removing random edges from the existing PPI network. We vary the fraction of bridge edges removed in increments $[0, 0.1, \dots, 0.6]$ and, for each fraction, we create a new edge list representing the modified network. The number of remaining directed edges in the network upon random edge removal are 273,319, 242,912, 212,525, 182,177, 151,811 and 121,472.
3. Dataset with latent confounder noise: our starting point is the chemical-PPI-breast-MDAMB231 dataset. The synthetic datasets were constructed with varying levels of confounding bias introduced into the gene expression data. To simulate latent confounder effects, Gaussian noise with distinct means and variances was progressively added to random subsets of genes. Genes were grouped into 50 predefined subsets, each representing a latent confounder group. For each group, a Gaussian distribution was defined, with the mean drawn randomly from a uniform distribution in the range $[0.5, 0.5]$ and the standard deviation $[0.1, 0.5]$. A fraction $[0.2, 0.4, 0.6, 0.8, 1]$ of these subsets was randomly selected for perturbation and, for each gene in these subsets, its expression value was incremented by a value sampled from the respective Gaussian distribution. The perturbed gene expression values were then clamped between zero and one to ensure validity. This strategy ensures that different latent biases are introduced globally to gene

expression patterns while maintaining controlled variability. We pair the noisy version of the gene expression data with the global unperturbed PPI network.

Network proximity between predicted and true perturbagens

Let \mathcal{P} be the set of predicted therapeutic targets, \mathcal{R} be the set of ground-truth therapeutic targets, and $\text{spd}(p, r)$ be the shortest-path distance between nodes in \mathcal{P} and \mathcal{R} . We measure the closest distance between \mathcal{P} and \mathcal{R} as:

$$d(\mathcal{P}, \mathcal{R}) = \frac{1}{|\mathcal{R}||\mathcal{P}|} \sum_{r \in \mathcal{R}} \sum_{p \in \mathcal{P}} \text{spd}(p, r). \quad (9)$$

As part of our performance analyses, we measure the network proximity of PDGrapher and competing methods. We compared the distributions of network proximity values using a Mann–Whitney U -test, along with a rank-biserial correlation to measure effect size. To assess the uncertainty of effect sizes, we performed bootstrapping with 1,000 resamples to estimate 95% CIs.

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

Processed datasets, including cell line gene expression datasets, PPI networks, drug targets and disease-associated genes, are available via the project website at <https://zitniklab.hms.harvard.edu/projects/PDGrapher> or directly at <https://zenodo.org/records/15375990> (ref. 76) and <https://zenodo.org/records/15390483> (ref. 77). The PPI data were obtained from <https://downloads.thebiogrid.org/File/BioGRID/Release-Archive/BIOGRID-3.5.186/BIOGRID-MV-Physical-3.5.186.tab3.zip>, https://www.science.org/doi/suppl/10.1126/science.1257601/suppl_file/datasets_s1-s4.zip and <http://www.interactome-atlas.org/data/HuRI.tsv>. Raw gene expression datasets were obtained from <https://clue.io/releases/data-dashboard>. Disease-associated genes were obtained from COSMIC at https://cancer.sanger.ac.uk/cell_lines/archive-download#:~:text=Complete. Source data are provided with this paper.

Code availability

Python implementation of PDGrapher is available at <https://github.com/mims-harvard/PDGrapher> (ref. 78).

References

- Vincent, F. et al. Phenotypic drug discovery: recent successes, lessons learned and new directions. *Nat. Rev. Drug Discov.* **21**, 899–914 (2022).
- Moffat, J. G., Vincent, F., Lee, J. A., Eder, J. & Prunotto, M. Opportunities and challenges in phenotypic drug discovery: an industry perspective. *Nat. Rev. Drug Discov.* **16**, 531–543 (2017).
- Minikel, E. V., Painter, J. L., Dong, C. C. & Nelson, M. R. Refining the impact of genetic evidence on clinical success. *Nature* **629**, 624–629 (2024).
- Druker, B. J. et al. Effects of a selective inhibitor of the Abl tyrosine kinase on the growth of Bcr–Abl positive cells. *Nat. Med.* **2**, 561–566 (1996).
- Bange, J., Zwick, E. & Ullrich, A. Molecular targets for breast cancer therapy and prevention. *Nat. Med.* **7**, 548–552 (2001).
- Swinney, D. C. & Anthony, J. How were new medicines discovered? *Nat. Rev. Drug Discov.* **10**, 507–519 (2011).
- Musa, A. et al. A review of Connectivity Map and computational approaches in pharmacogenomics. *Brief. Bioinform.* **19**, 506–523 (2018).
- Davies, J. C., Alton, E. W. & Bush, A. Cystic fibrosis. *BMJ* **335**, 1255–1259 (2007).
- Van Goor, F. et al. Rescue of CF airway epithelial cell function in vitro by a CFTR potentiator, VX-770. *Proc. Natl Acad. Sci. USA* **106**, 18825–18830 (2009).
- Van Goor, F. et al. Correction of the F508del-CFTR protein processing defect in vitro by the investigational drug VX-809. *Proc. Natl Acad. Sci. USA* **108**, 18843–18848 (2011).
- Keenan, A. B. et al. Connectivity mapping: methods and applications. *Annu. Rev. Biomed. Data Sci.* **2**, 69–92 (2019).
- Keenan, A. B. et al. The Library of Integrated Network-based Cellular Signatures (LINCS) NIH program: system-level cataloging of human cells response to perturbations. *Cell Syst.* **6**, 13–24 (2018).
- Lamb, J. et al. The Connectivity Map: using gene-expression signatures to connect small molecules, genes, and disease. *Science* **313**, 1929–1935 (2006).
- Samart, K., Tuyishime, P., Krishnan, A. & Ravi, J. Reconciling multiple connectivity scores for drug repurposing. *Brief. Bioinform.* **22**, bbab161 (2021).
- Guney, E. Reproducible drug repurposing: when similarity does not suffice. *Pac. Symp. Biocomput.* **22**, 132–143 (2017).
- Chen, B. Reversal of cancer gene expression correlates with drug efficacy and reveals therapeutic targets. *Nat. Commun.* **8**, 16022 (2017).
- Chen, B. et al. Computational discovery of niclosamide ethanolamine, a repurposed drug candidate that reduces growth of hepatocellular carcinoma cells in vitro and in mice by inhibiting cell division cycle 37 signaling. *Gastroenterology* **152**, 2022–2036 (2017).
- Pesetto, Z. Y. et al. In silico and in vitro drug screening identifies new therapeutic approaches for Ewing sarcoma. *Oncotarget* **8**, 4079–4095 (2016).
- Morselli Gysi, D. et al. Network medicine framework for identifying drug-repurposing opportunities for COVID-19. *Proc. Natl Acad. Sci. USA* **118**, e2025581118 (2021).
- Zhu, J. et al. Prediction of drug efficacy from transcriptional profiles with deep learning. *Nat. Biotechnol.* **39**, 1444–1452 (2021).
- Pham, T. H., Qiu, Y., Zeng, J., Xie, L. & Zhang, P. A deep learning framework for high-throughput mechanism-driven phenotype compound screening and its application to COVID-19 drug repurposing. *Nat. Mach. Intell.* **3**, 247–257 (2021).
- Lotfollahi, M., Wolf, F. A. & Theis, F. J. scGen predicts single-cell perturbation responses. *Nat. Methods* **16**, 715–721 (2019).
- Hetzl, L. et al. Predicting cellular responses to novel drug perturbations at a single-cell resolution. *Adv. Neural Inf. Process. Syst.* **35**, 26711–26722 (2022).
- Hauser, A. & Bühlmann, P. Two optimal strategies for active learning of causal models from interventional data. *Int. J. Approx. Reason.* **55**, 926–939 (2014).
- Ghassami, A. E., Salehkaleybar, S., Kiyavash, N. & Bareinboim, E. Budgeted experiment design for causal structure learning. In *Proc. 35th International Conference on Machine Learning* (eds Dy, J. & Krause, A.) 2788–2801 (Proceedings of Machine Learning Research, 2018).
- Agrawal, R., Squires, C., Yang, K., Shanmugam, K. & Uhler, C. ABCD-strategy: budgeted experimental design for targeted causal structure discovery. In *Proc. 22nd International Conference on Artificial Intelligence and Statistics* (eds Chaudhuri, K. & Sugiyama, M.) 3400–3409 (Proceedings of Machine Learning Research, 2019).
- Bunne, C. et al. Learning single-cell perturbation responses using neural optimal transport. *Nat. Methods* **20**, 1759–1768 (2023).
- Barr, M. P. et al. Vascular endothelial growth factor is an autocrine growth factor, signaling through neuropilin-1 in non-small cell lung cancer. *Mol. Cancer* **14**, 45 (2015).

29. Riquelme, E. et al. VEGF/VEGFR-2 upregulates EZH2 expression in lung adenocarcinoma cells and EZH2 depletion enhances the response to platinum-based and VEGFR-2-targeted therapy. *Clin. Cancer Res.* **20**, 3849–3861 (2014).
30. Huynh-Thu, V. A., Irrthum, A., Wehenkel, L. & Geurts, P. Inferring regulatory networks from expression data using tree-based methods. *PLoS ONE* **5**, e12776 (2010).
31. Yuan, B. et al. CellBox: interpretable machine learning for perturbation biology with application to the design of cancer combination therapy. *Cell Syst.* **12**, 128–140 (2021).
32. Kamimoto, K. et al. Dissecting cell identity via network inference and in silico gene perturbation. *Nature* **614**, 742–751 (2023).
33. Piran, Z., Cohen, N., Hoshen, Y. & Nitzan, M. Disentanglement of single-cell data with biolord. *Nat. Biotechnol.* **42**, 1678–1683 (2024).
34. Roohani, Y., Huang, K. & Leskovec, J. Predicting transcriptional outcomes of novel multigene perturbations with gears. *Nat. Biotechnol.* **42**, 927–935 (2024).
35. Barabási, A. L., Gulbahce, N. & Loscalzo, J. Network medicine: a network-based approach to human disease. *Nat. Rev. Genet.* **12**, 56–68 (2011).
36. Ruiz, C., Zitnik, M. & Leskovec, J. Identification of disease treatment mechanisms through the multiscale interactome. *Nat. Commun.* **12**, 1796 (2021).
37. Eyuboglu, S., Zitnik, M. & Leskovec, J. Mutual interactors as a principle for phenotype discovery in molecular interaction networks. *Pac. Symp. Biocomput.* **28**, 61–72 (2023).
38. Gu, Z. et al. Role of duplicate genes in genetic robustness against null mutations. *Nature* **421**, 63–66 (2003).
39. Deutscher, D., Meilijson, I., Kupiec, M. & Ruppín, E. Multiple knockout analysis of genetic robustness in the yeast metabolic network. *Nat. Genet.* **38**, 993–998 (2006).
40. Deutscher, D., Meilijson, I., Schuster, S. & Ruppín, E. Can single knockouts accurately single out gene functions? *BMC Syst. Biol.* **2**, 50 (2008).
41. Ochoa, D. et al. The next-generation Open Targets Platform: reimaged, redesigned, rebuilt. *Nucleic Acids Res.* **51**, D1353–D1359 (2023).
42. Cai, Q. et al. The role of dexmedetomidine in tumor-progressive factors in the perioperative period and cancer recurrence: a narrative review. *Drug Des. Devel. Ther.* **16**, 2161–2175 (2022).
43. Gainor, J. F. et al. Pralsetinib for ret fusion-positive non-small-cell lung cancer (arrow): a multi-cohort, open-label, phase 1/2 study. *Lancet Oncol.* **22**, 959–969 (2021).
44. Faraat Ali, G. C. & Kumari, N. Pralsetinib: chemical and therapeutic development with FDA authorization for the management of ret fusion-positive non-small-cell lung cancers. *Arch. Pharm. Res.* **45**, 309–327 (2022).
45. Wu, J. et al. Expression and potential molecular mechanism of TOP2A in metastasis of non-small cell lung cancer. *Sci. Rep.* **14**, 12228 (2024).
46. Pautier, P. et al. Doxorubicin-trabectedin with trabectedin maintenance in leiomyosarcoma. *N. Engl. J. Med.* **391**, 789–799 (2024).
47. Hartwell, L. H., Hopfield, J. J., Leibler, S. & Murray, A. W. From molecular to modular cell biology. *Nature* **402**, C47–C52 (1999).
48. Menche, J. et al. Uncovering disease–disease relationships through the incomplete interactome. *Science* **347**, 1257601 (2015).
49. Szklarczyk, D. et al. The string database in 2023: protein–protein association networks and functional enrichment analyses for any sequenced genome of interest. *Nucleic Acids Res.* **51**, D638–D646 (2023).
50. Bagherian, M. et al. Machine learning approaches and databases for prediction of drug-target interaction: a survey paper. *Brief. Bioinform.* **22**, 247–269 (2021).
51. Pan, J. et al. Sparse dictionary learning recovers pleiotropy from human cell fitness screens. *Cell Syst.* **13**, 286–303 (2022).
52. Huang, K. et al. Deepurpose: a deep learning library for drug–target interaction prediction. *Bioinformatics* **36**, 5545–5547 (2020).
53. Hart, G. T., Ramani, A. K. & Marcotte, E. M. How complete are current yeast and human protein–interaction networks? *Genome Biol.* **7**, 120 (2006).
54. Stumpf, M. P. et al. Estimating the size of the human interactome. *Proc. Natl Acad. Sci. USA* **105**, 6959–6964 (2008).
55. Zitnik, M., Sosič, R., Feldman, M. W. & Leskovec, J. Evolution of resilience in protein interactomes across the tree of life. *Proc. Natl Acad. Sci. USA* **116**, 4426–4433 (2019).
56. Schölkopf, B. et al. Toward causal representation learning. *Proc. IEEE* **109**, 612–634 (2021).
57. Gustafsdottir, S. M. et al. Multiplex cytological profiling assay to measure diverse. *PLoS ONE* **8**, e80999 (2013).
58. Bray, M.-A. et al. Cell Painting, a high-content image-based assay for morphological profiling using multiplexed fluorescent dyes. *Physiol. Behav.* **176**, 139–148 (2017).
59. Chandrasekaran, S. N. et al. JUMP Cell Painting dataset: morphological impact of 136,000 chemical and genetic perturbations. Preprint at *bioRxiv* <https://doi.org/10.1101/2023.03.23.534023> (2023).
60. Akbarzadeh, M. et al. Morphological profiling by means of the Cell Painting assay enables identification of tubulin-targeting compounds. *Cell Chem. Biol.* **29**, 1053–1064 (2022).
61. Pruteanu, L. L. & Bender, A. Using transcriptomics and cell morphology data in drug discovery: the long road to practice. *ACS Med. Chem. Lett.* **14**, 386–395 (2023).
62. Huang, K. et al. Artificial intelligence foundation for therapeutic science. *Nat. Chem. Biol.* **18**, 1033–1036 (2022).
63. Thomas, K. J. & Jacobson, M. R. Defects in mitochondrial fission protein dynamin-related protein 1 are linked to apoptotic resistance and autophagy in a lung cancer model. *PLOS ONE* **7**, e45319 (2012).
64. Le, C. F., Yusof, M. Y. Y., Hassan, H. & Sekaran, S. D. In vitro properties of designed antimicrobial peptides that exhibit potent antipneumococcal activity and produce synergism in combination with penicillin. *Sci. Rep.* **5**, 9761 (2015).
65. Zhao, S., Song, P., Zhou, G., Zhang, D. & Hu, Y. METTL3 promotes the malignancy of non-small cell lung cancer by N6-methyladenosine modifying SFRP2. *Cancer Gene Ther.* **30**, 1094–1104 (2023).
66. Zhu, J.-Y., Park, T., Isola, P. & Efros, A. A. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proc. 2017 IEEE International Conference on Computer Vision (ICCV)* 2242–2251 (IEEE, 2017).
67. Neal, B. *Introduction to Causal Inference*. Course Lecture Notes (draft). (Department of Statistics, Harvard University, 2020); https://www.bradyneal.com/Introduction_to_Causal_Inference-Dec17_2020-Neal.pdf
68. Lotfollahi, M. et al. Predicting cellular responses to complex perturbations in high-throughput screens. *Mol. Syst. Biol.* **19**, e11517 (2023).
69. Paszke, A. et al. Automatic differentiation in PyTorch. In *Proc. NIPS Workshop on Autodiff* (2017).
70. Fey, M. & Lenssen, J. E. Fast graph representation learning with PyTorch Geometric. Preprint at <https://arxiv.org/abs/1903.02428> (2019).
71. Heumos, L. et al. Best practices for single-cell analysis across modalities. *Nat. Rev. Genet.* **24**, 550–572 (2023).
72. Gayoso, A. et al. A Python library for probabilistic analysis of single-cell omics data. *Nat. Biotechnol.* **40**, 163–166 (2022).
73. Jovic, D. et al. Single-cell RNA sequencing technologies and applications: a brief overview. *Clin. Transl. Med.* **12**, e694 (2022).

74. Luecken, M. D. et al. Benchmarking atlas-level data integration in single-cell genomics. *Nat. Methods* **19**, 41–50 (2021).
75. Wishart, D. S. et al. DrugBank 5.0: a major update to the DrugBank database for 2018. *Nucleic Acids Res.* **46**, D1074–D1082 (2018).
76. Gonzalez, G. et al. Combinatorial prediction of therapeutic perturbations using causally-inspired neural networks (genetic data). Zenodo <https://doi.org/10.5281/zenodo.15375990> (2025).
77. Gonzalez, G. et al. Combinatorial prediction of therapeutic perturbations using causally-inspired neural networks (chemical data). Zenodo <https://doi.org/10.5281/zenodo.15390483> (2025).
78. Gonzalez, G. et al. PDGrapher. GitHub <https://github.com/mims-harvard/PDGrapher/tree/main> (2025).

Acknowledgements

We thank D. Mohorcic for his help with the PDGrapher codebase. We acknowledge the support of NIH R01-HD108794, NSF CAREER 2339524, US DoD FA8702-15-D-0001, ARPA-H BDF programme, awards from Chan Zuckerberg Initiative, Bill & Melinda Gates Foundation INV-079038, Amazon Faculty Research, Google Research Scholar Program, AstraZeneca Research, Roche Alliance with Distinguished Scientists, Sanofi iDEA-iTECH, Pfizer Research, John and Virginia Kaneb Fellowship at Harvard Medical School, Biswas Computational Biology Initiative in partnership with the Milken Institute, Harvard Medical School Dean's Innovation Fund for the Use of Artificial Intelligence, Harvard Data Science Initiative, and Kempner Institute for the Study of Natural and Artificial Intelligence at Harvard University. I.H. was supported, in part, by the Summer Institute in Biomedical Informatics at Harvard Medical School. M.B. and G.G. were supported by the ERC Consolidator Grant number 724228 (LEMAN). Any opinions, findings, conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the funders.

Author contributions

G.G. retrieved, processed and analysed gene expression data, PPI data and disease gene datasets. I.H. retrieved, processed and analysed drug target data and processed GRN. G.G. and X.L. developed, implemented and benchmarked PDGrapher and performed detailed analyses of PDGrapher. G.G., X.L., K.V., M.B. and M.Z. designed the study. G.G., X.L., I.H. and M.Z. wrote the paper.

Competing interests

G.G. is currently employed by Genentech and I.H. was employed by Merck & Co. during the study. The other authors declare no competing interests.

Additional information

Extended data is available for this paper at <https://doi.org/10.1038/s41551-025-01481-x>.

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41551-025-01481-x>.

Correspondence and requests for materials should be addressed to Marinka Zitnik.

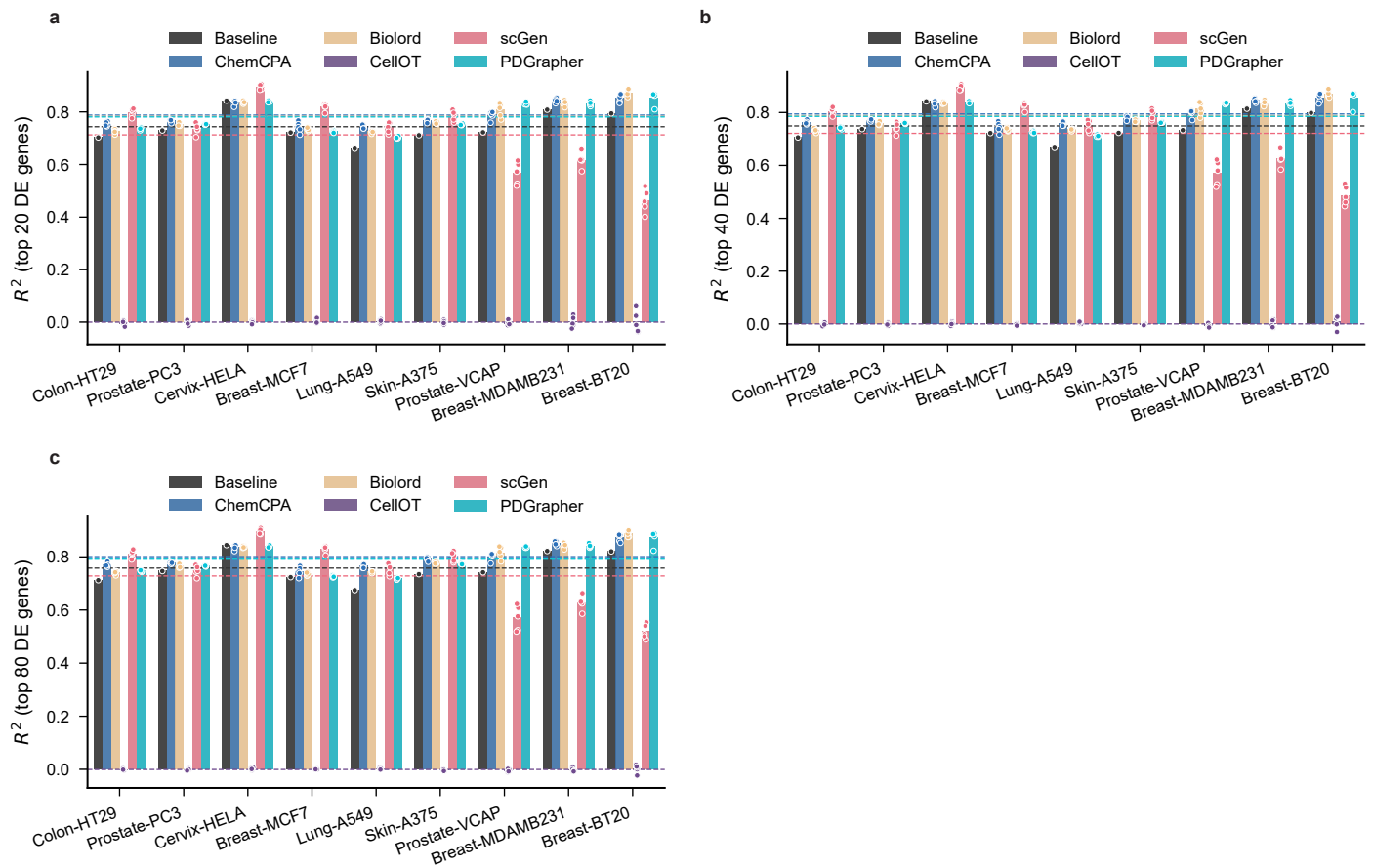
Peer review information *Nature Biomedical Engineering* thanks Ping Zhang and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

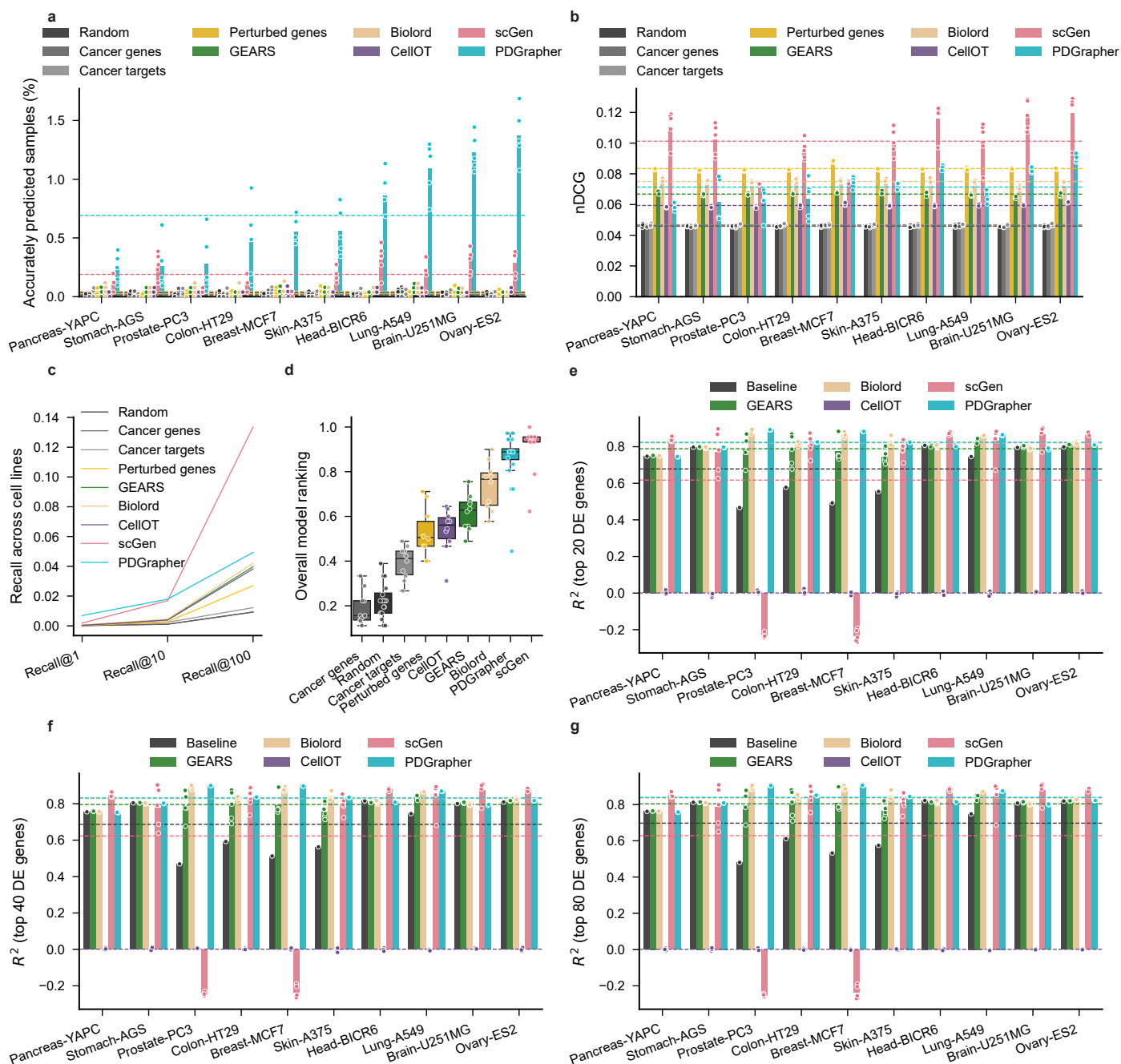
Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025



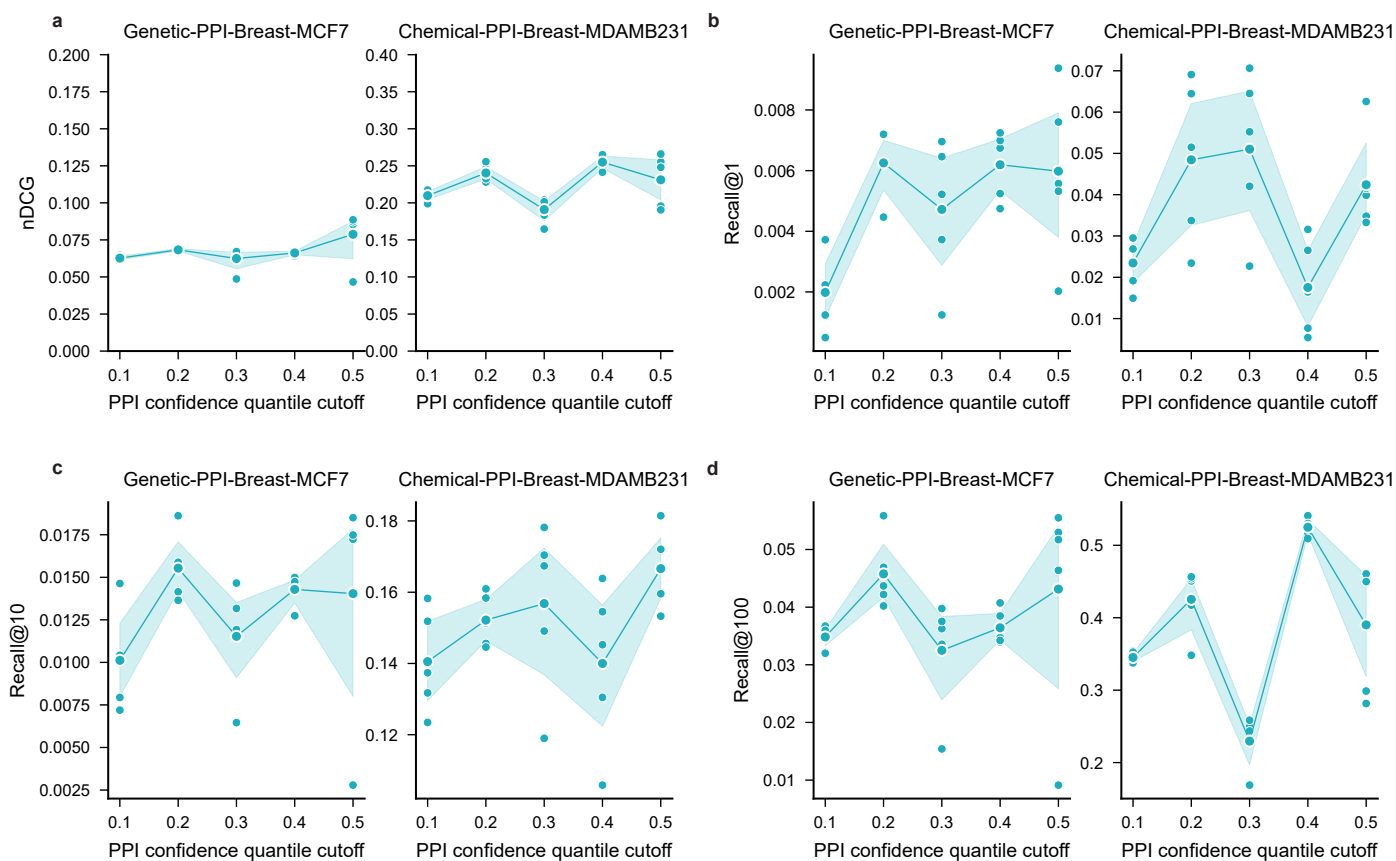
Extended Data Fig. 1 | The performance of response prediction within nine cell lines under chemical perturbation. The R^2 values are calculated between the predicted and actual gene expression for the top 20 (a), 40 (b), and 80 (c)

differentially expressed genes per cell line. Dotted lines represent the average performance across cell lines, dots indicate individual data points, and bars represent the average R^2 across five data splits.



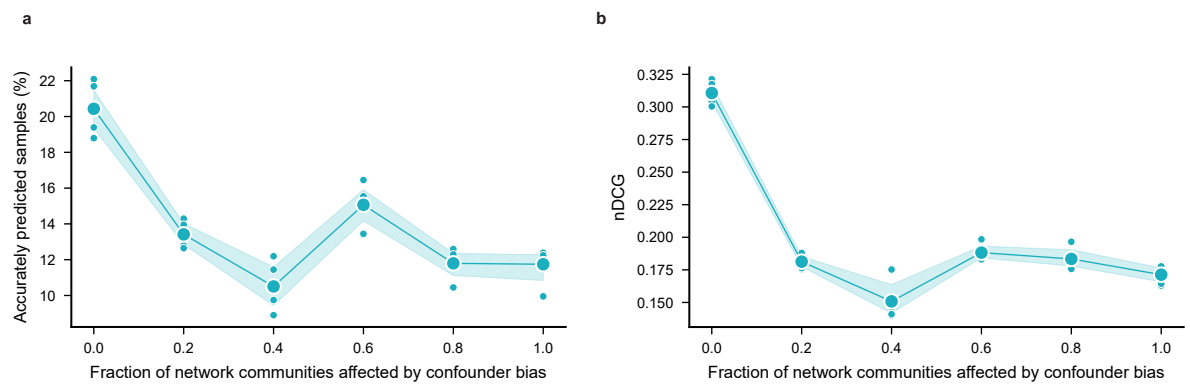
Extended Data Fig. 2 | PDGrapher efficiently predicts genetic perturbagens to shift cells from diseased to treated states in a random splitting setting within ten cell lines. (a) PDGrapher provides accurate predictions for up to 1.09% more samples in the test set compared to the second-best baseline across Genetic-PPI datasets (genetic-PPI-ovary-ES2: 1.37% vs 0.28%). (b) scGen takes the leading position in nDCG across genetic-PPI datasets. (c) PDGrapher recovers ground-truth therapeutic targets at comparable rates as competing methods for genetic-PPI datasets. (d) PDGrapher has the best overall performance in

perturbagen prediction within each cell line, evaluated by the averaged rank over multiple cell lines and metrics. The central line inside the box represents the median, while the top and bottom edges correspond to the first (Q1) and third (Q3) quartiles. The whiskers extend to the smallest and largest values within 1.5 times the interquartile range (IQR) from the quartiles. (e-g) Shown is the R^2 of the response prediction module of PDGrapher compared to competing baselines for the top 20 (e), 40 (f), and 80 (g) differentially expressed (DE) genes.



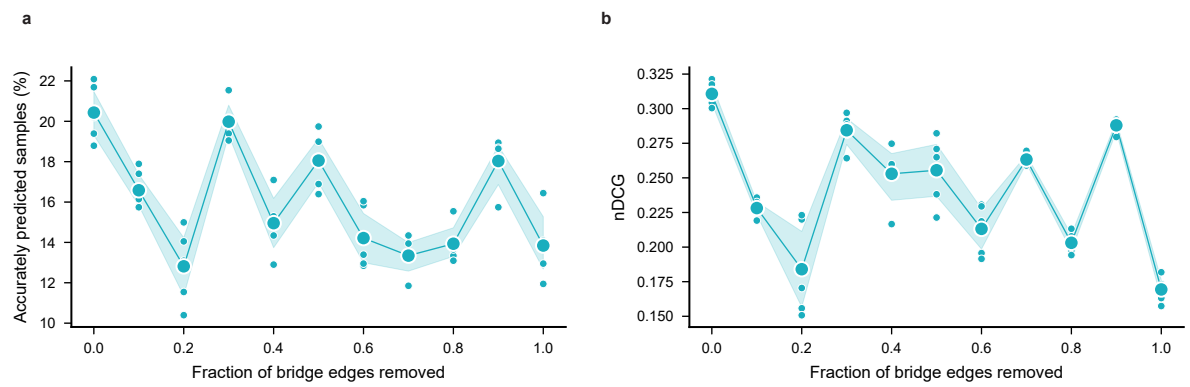
Extended Data Fig. 3 | Sensitivity analysis of the PPI used for training PDGrapher. Performance of sensitivity analyses evaluated by nDCG (a) and recalls (b–d) for datasets genetic-PPI-breast-MCF7 (left) and chemical-PPI-breast-MDAMB231 (right). The PPI used here is from STRING (string-db.org), which includes a confidence score for each edge. The edges are

filtered by the 0.1, 0.2, 0.3, 0.4, and 0.5 quantiles of the confidence scores as cutoffs, resulting in five PPI networks with 625,818, 582,305, 516,683, 443,051, and 296,451 edges, respectively. The results of the percentage of accurately predicted samples are shown in Figure 5b.



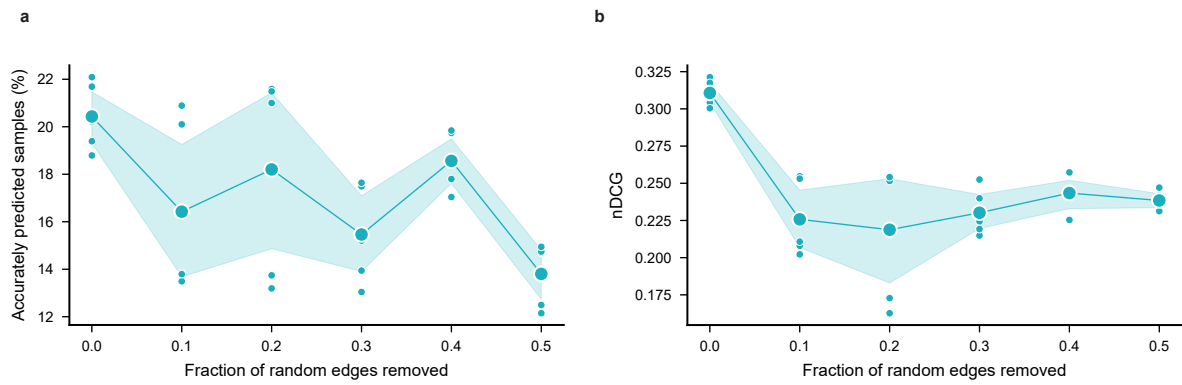
Extended Data Fig. 4 | PDGrapher has stable performance on the synthetic datasets with various intensities of confounders added on the gene expression. Performance of simulation analyses evaluated by percentage of accurately predicted samples (**a**) and nDCG (**b**) for the synthetic datasets with varying levels (0 to 1) of confounding bias introduced into the gene expression data. Gaussian noise, with distinct means and variances, was added progressively to random subsets of genes, simulating latent confounder effects in the treated

gene expression data. The intensity of the confounding bias increases as more gene groups (representing network communities) are affected. This approach creates global, controlled variability in the gene expression data, paired with an unperturbed PPI network, allowing for the evaluation of algorithmic performance across different degrees of confounder noise. See Online Methods for more details on data generation.



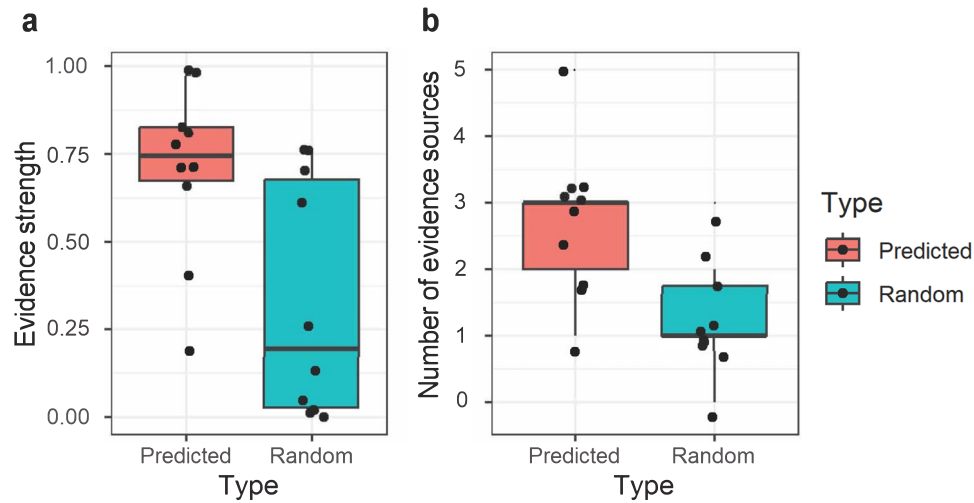
Extended Data Fig. 5 | PDGrapher has stable performance on the synthetic datasets with various fractions of bridge edges removed. Performance of simulation analyses evaluated by percentage of accurately predicted samples (a) and nDCG (b) for the synthetic datasets with a [0, 0.1, ..., 1] fraction of bridge edges removed in the simulated PPI. Bridge edges are those with

high connectivity in the network, which, if removed, increase the number of disconnected communities. The number of connected components in the network upon bridge edge removal is [90, 179, 268, 358, 447, 536, 626, 715, 804, 894]. See more information in Supplementary Table S5.



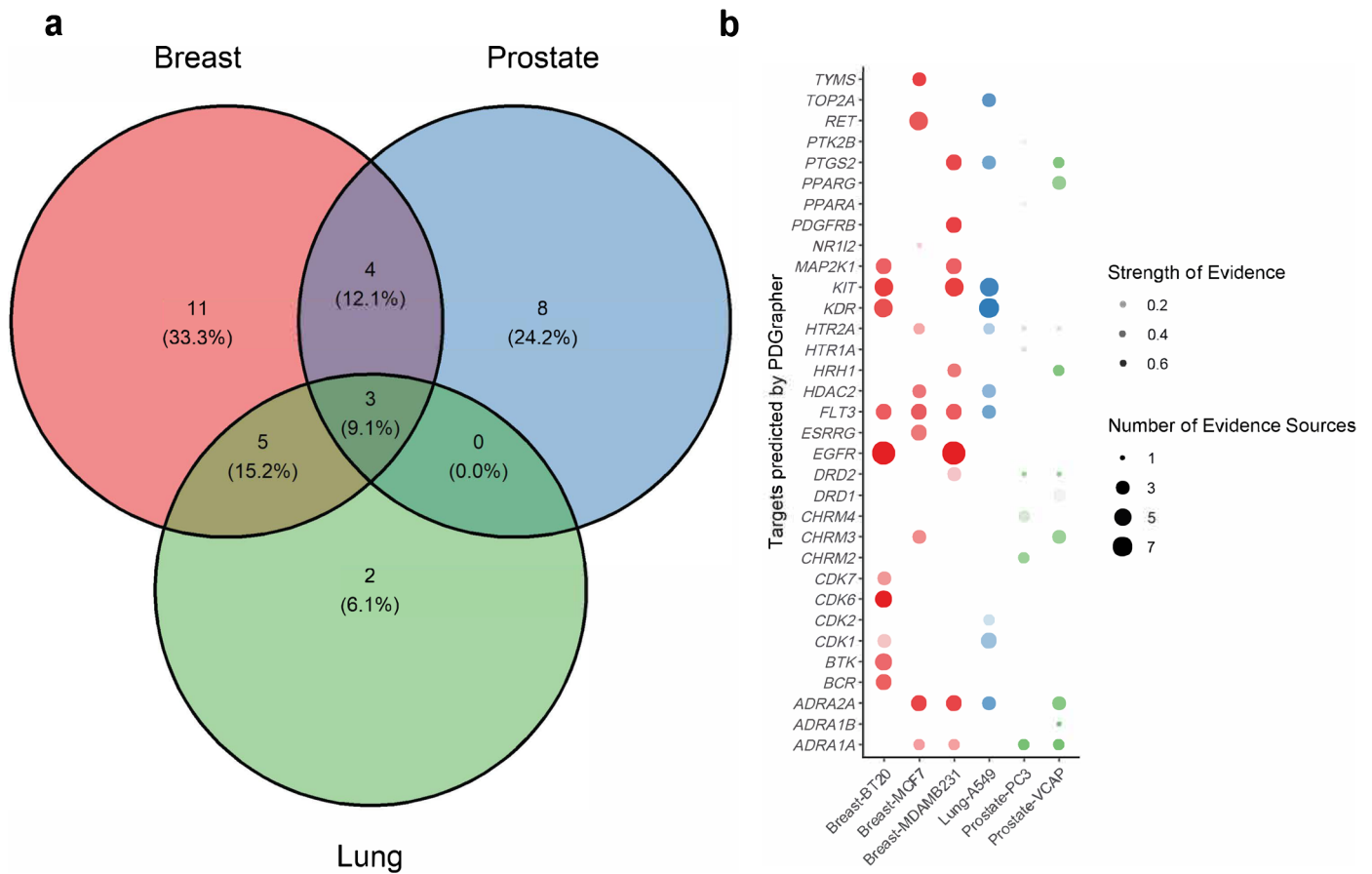
Extended Data Fig. 6 | PDGrapher's performance is influenced by network incompleteness. Performance in ablation studies evaluated by percentage of accurately predicted samples (a) and nDCG (b) for the synthetic datasets with

a [0, 0.1, ... 0.6] fraction of random edges removed in the PPI. The number of remaining edges in the network upon random edge removal are [273,319; 242,912; 212,525; 182,177; 151,811; 121,472]. See the Method section for more details.



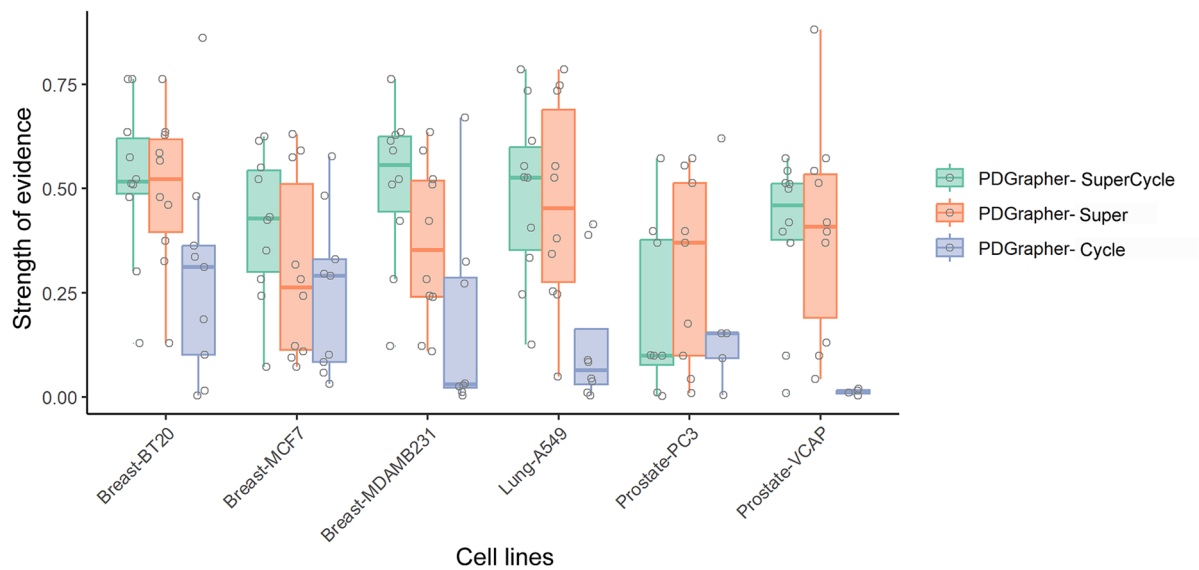
Extended Data Fig. 7 | Comparison of predicted targets from PDGrapher and a random model for lung cancer. Boxplots display the evidence strength (a) and the number of evidence sources (b) for the top 10 predicted targets from PDGrapher versus 10 randomly selected genes. The central line inside the box

represents the median, while the top and bottom edges correspond to the first (Q1) and third (Q3) quartiles. The whiskers extend to the smallest and largest values within 1.5 times the interquartile range (IQR) from the quartiles.



Extended Data Fig. 8 | Unique and common targets predicted by PDGrapher among three cancer types. The Venn diagram (a) shows the number and ratio of unique and common predicted targets, while the bubble plot (b) shows the strength

of evidence and the number of evidence sources for each cell line. The evidence strengths are the global association scores and overall evidence sources provided by Open Targets. Details of the scoring system are in Supplementary Note 2.



Extended Data Fig. 9 | Ablation studies for loss functions of PDGrapher evaluated by the evidence from Open Targets. The strength of evidence in the cell lines with healthy control data is shown in the box plots. The strength of the evidence is the global association scores, which are based on all evidence sources provided by Open Targets. See details of the scoring system in Supplementary

Note 2. The central line inside the box represents the median, while the top and bottom edges correspond to the first (Q1) and third (Q3) quartiles. The whiskers extend to the smallest and largest values within 1.5 times the interquartile range (IQR) from the quartile.

Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

- | n/a | Confirmed |
|-------------------------------------|--|
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The statistical test(s) used AND whether they are one- or two-sided
<i>Only common tests should be described solely by name; describe more complex techniques in the Methods section.</i> |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A description of all covariates tested |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
<i>Give P values as exact values whenever suitable.</i> |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated |

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

Data analysis

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

Processed data used in this paper, including the cell line gene expression dataset, protein-protein interaction network, drug targets, and disease-associated genes, are available via the project website at <https://zitniklab.hms.harvard.edu/projects/PDGrapher> or directly at <https://figshare.com/articles/dataset/>

Combinatorial_prediction_of_therapeutic_targets_using_a_causally-inspired_neural_network/24798855. The raw protein-protein interaction network data was obtained from <https://downloads.thebiogrid.org/File/BioGRID/Release-Archive/BIOGRID-3.5.186/BIOGRID-MV-Physical-3.5.186.tab3.zip>, https://www.science.org/doi/suppl/10.1126/science.1257601/suppl_file/datasets_sl-s4.zip and <http://www.interactome-atlas.org/data/HuRI.tsv> Raw gene expression datasets were obtained from <https://clue.io/releases/data-dashboard>. Disease-associated genes were obtained from COSMIC at https://cancer.sanger.ac.uk/cell_lines/archivedownload#:~:text=Complete%20mutation%20data and <https://cancer.sanger.ac.uk/cosmic/curation>. Drug targets were extracted from DrugBank at <https://go.drugbank.com/releases/5-1-9> and a list of cancer drugs was obtained from NCI at <https://www.cancer.gov/about-cancer/treatment/types/targeted-therapies/approved-drug-list>. STRING (v12.0) was downloaded from <https://stringdb-downloads.org/download/protein.physical.links.detailed.v12.0.txt.gz>

Research involving human participants, their data, or biological material

Policy information about studies with [human participants or human data](#). See also policy information about [sex, gender \(identity/presentation\), and sexual orientation](#) and [race, ethnicity and racism](#).

Reporting on sex and gender	N/A.
Reporting on race, ethnicity, or other socially relevant groupings	N/A.
Population characteristics	N/A.
Recruitment	N/A.
Ethics oversight	N/A.

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	No formal statistical method was used to predetermine sample size. Instead, we included all samples from the LINCS dataset that passed a set of stringent filtering criteria designed to ensure biological relevance and data quality (see the information below). As a result, the sample size reflects the maximum number of high-quality, biologically relevant perturbations available for each model. Tables S9–S10 detail the cell line selection process, gene filtering criteria, and rationale for inclusion/exclusion. Sample sizes can be found in Tables S11-12
Data exclusions	<ol style="list-style-type: none"> 1. Cell line filtering: Cell lines were filtered to keep those with sufficient perturbational coverage and the inclusion of healthy cell line counterparts. Figures S9 and S10 contain a description of the cell line selection criteria together with a list of cell lines with the largest number of perturbed samples and a reason for inclusion/exclusion. 2. Healthy counterpart selection is described in Tables S11 and S12 in the supplementary material. 3. Disease-associated genes: We extracted disease-associated genes from COSMIC (Accessed in September 2022) in addition to expert-curated genes available at https://cancer.sanger.ac.uk/cosmic/curation. Genes were represented using the HUGO Gene Nomenclature Committee ID. For each cell line in our dataset that has disease intervention data, we extracted cancer-causing mutations as the list of genes with "Verified" Mutation verification status in COSMIC and present in the list of genes curated by experts. Mapping the resulting genes to our list of genes in the PPI resulted in disease-associated genes. We excluded cell lines for which there were no disease-associated genes in COSMIC. 4. Gene matching: Treated samples were excluded if the targeted genes were not included in the protein-protein interaction (PPI) network. Genes in the gene expression dataset (LINCS) were matched to proteins in the PPI using the HUGO Gene Nomenclature Committee ID, which identified 10,716 overlapping genes. 5. Data level selection: Only level 3 gene expression data, which is quantile-normalized and can be compared across plates, was used. 6. Treatment and measurement specifics: Chemical interventions were included at all dose levels and time points.
Replication	We used 5-fold cross-validation to assess the stability of the model across different subsets of the data
Randomization	The datasets were split randomly into training and test sets to measure model performance. Additionally, a leave-cell-line-out setting was used for some tests to assess performance on unseen cell lines
Blinding	Blinding was not applicable because the study exclusively involved computational analysis of publicly available, pre-existing datasets (e.g., LINCS, COSMIC) where no new data collection or subjective labeling was performed. All labels (e.g., cell line identity, perturbation type, gene

expression readouts) were defined prior to analysis and were not influenced by the investigators, eliminating potential for observer or experimenter bias.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

- | n/a | Involvement in the study |
|-------------------------------------|--|
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Antibodies |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Eukaryotic cell lines |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Palaeontology and archaeology |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Animals and other organisms |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Clinical data |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Dual use research of concern |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Plants |

Methods

- | n/a | Involvement in the study |
|-------------------------------------|---|
| <input checked="" type="checkbox"/> | <input type="checkbox"/> ChIP-seq |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Flow cytometry |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> MRI-based neuroimaging |

Plants

Seed stocks

N/A.

Novel plant genotypes

N/A.

Authentication

N/A.