

Prescient Design

A Genentech Accelerator



**Combinatorial prediction of therapeutic
targets using a causally inspired GNN**

**Application of Deep Learning on Graphs
ETH**

Guadalupe Gonzalez
Senior ML Scientist, Frontier Research

Collaborators



Marinka Zitnik



Michael Bronstein

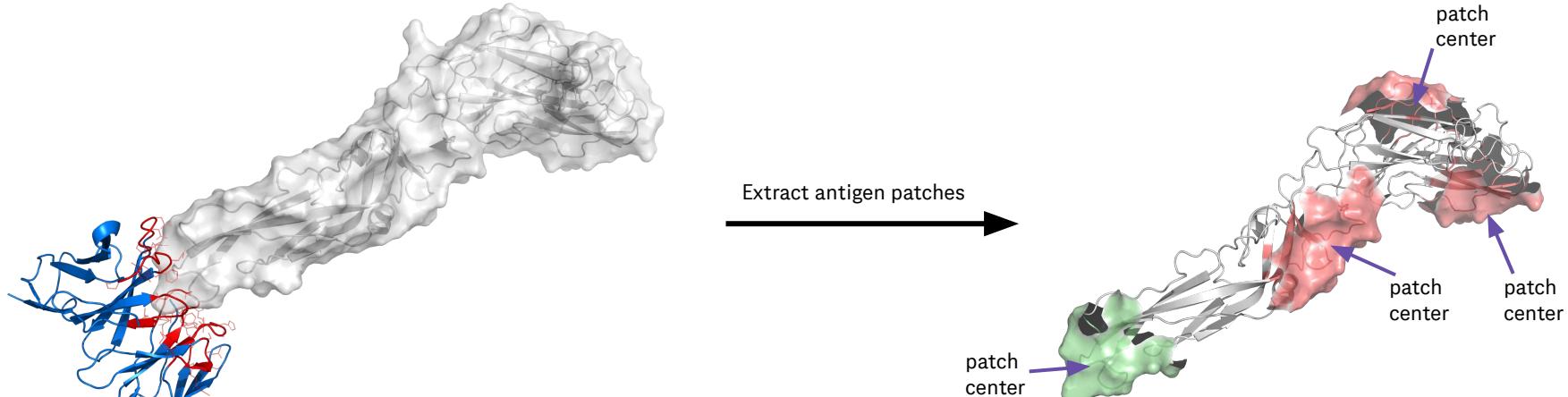


Kirill Veselkov

Isuru Herath
Domen Mohorcic

Introduction

A familiar use case | Epiphany Affinity



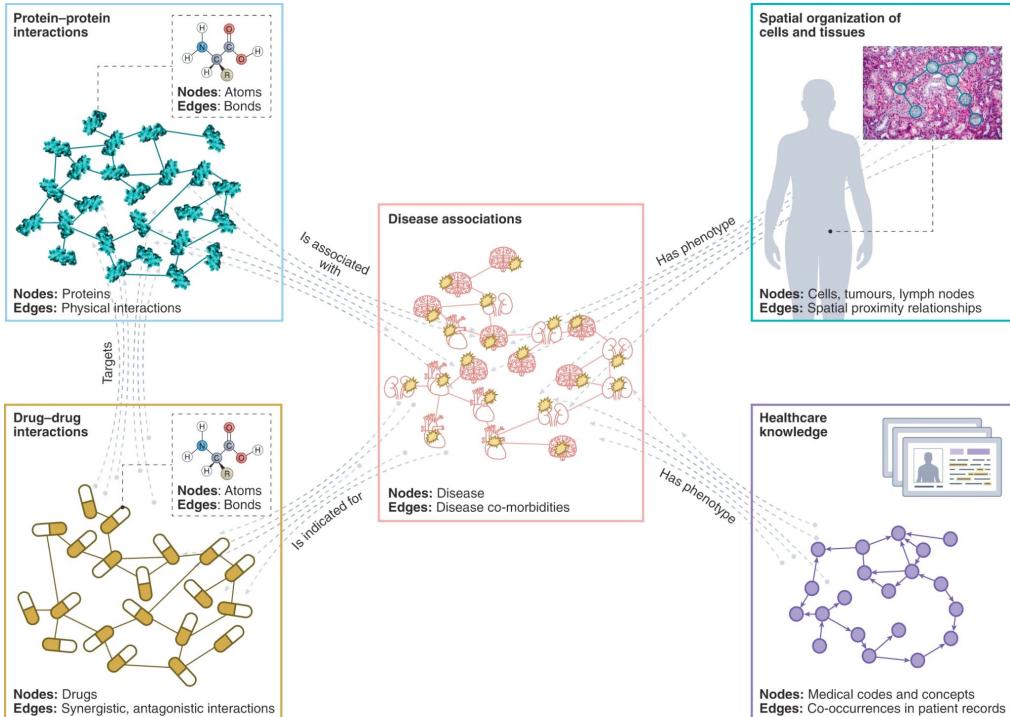
(1) do they bind?

(2) is this the epitope?

(3) is this the paratope?

Images courtesy of Jan Ludwiczak

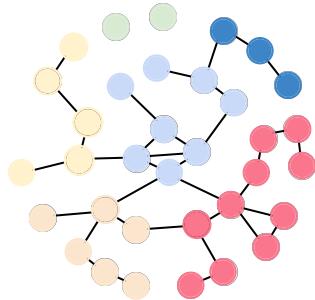
Zooming out | Graph representations are everywhere in biology



Current GNN paradigm | Message passing

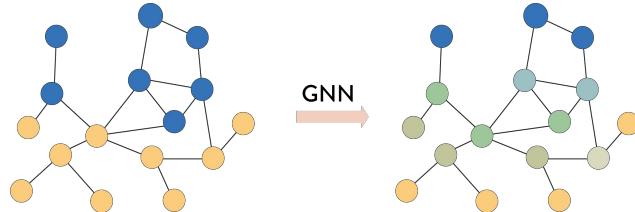
Biological principles that deep learning on graphs naturally includes:

Local hypothesis: interacting entities are more similar than non-interacting entities



PPI

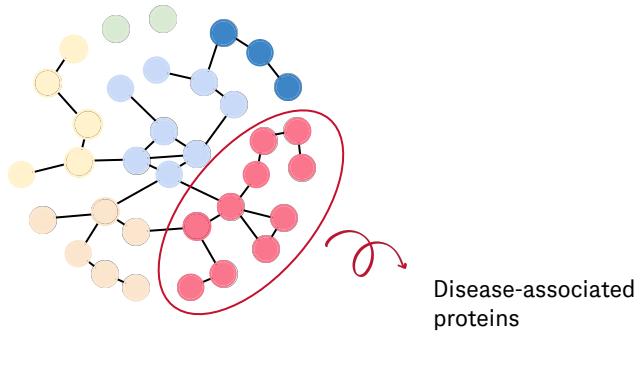
Homophily principle: aggregates messages from neighbors to compute enhanced representations under the assumption that neighbors are similar



Current GNN paradigm | Message passing

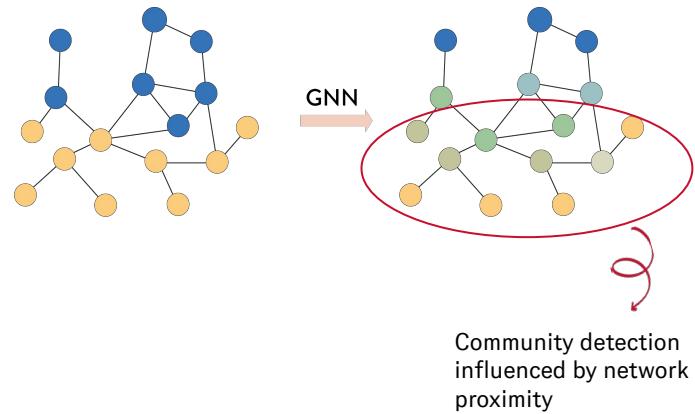
Biological principles that deep learning on graphs naturally includes:

Disease module hypothesis: cellular components tend to cluster in the same network neighborhood



PPI

Community detection: predict communities based on node embedding similarities which are higher for connected nodes
 → it will typically predict nodes closer in the network as belonging to the same community

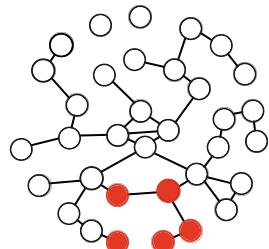


Li et al. 2022

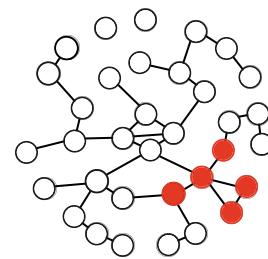
Current GNN paradigm | Message passing

Biological principles that deep learning on graphs naturally includes:

Shared-component hypothesis: diseases driven by perturbations of the same components (or close in the network) are phenotypically similar



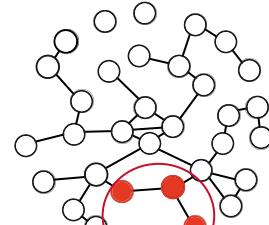
Disease A



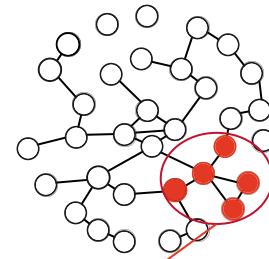
Disease B

PPI

Subgraph representations: similar representations for overlapping or similar subgraphs



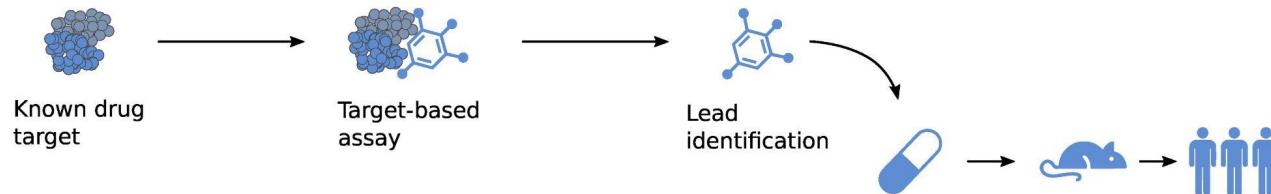
Similar representations



A brief history of therapeutics discovery: Phenotype-driven vs target-driven approaches

Phenotype-driven vs target-driven drug discovery

(A) Target-based drug discovery



(B) Phenotypic drug discovery

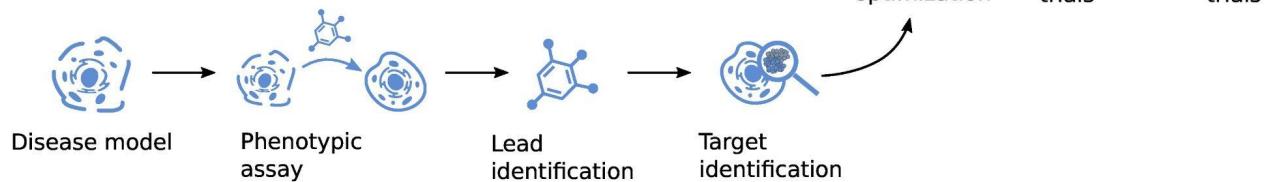
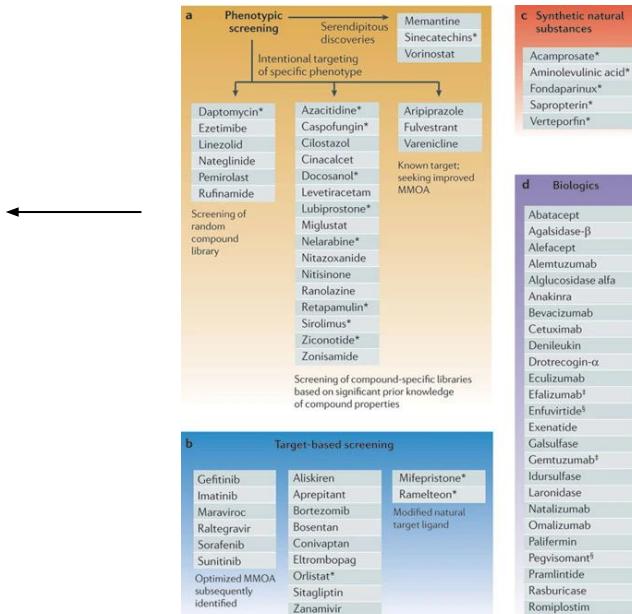


Image source = Krentzel, D., Shorte, S.L. and Zimmer, C. (2023) 'Deep learning in image-based phenotypic drug discovery', *Trends in Cell Biology*, 33(7), pp. 538–554. doi:10.1016/j.tcb.2022.11.011.

Phenotype-driven vs target-driven drug discovery

Majority of the first-in-class drugs approved by the US FDA between 1999-2008 were discovered empirically without a drug target hypothesis



Nature Reviews | Drug Discovery

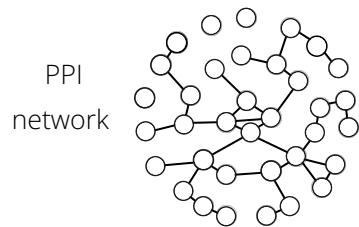
Image source = Swinney, D., Anthony, J. How were new medicines discovered?. *Nat Rev Drug Discov* 10, 507–519 (2011). <https://doi.org/10.1038/nrd3480>

Problem formulation

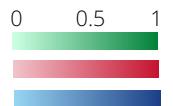
Problem formulation

Problem formulation

healthy

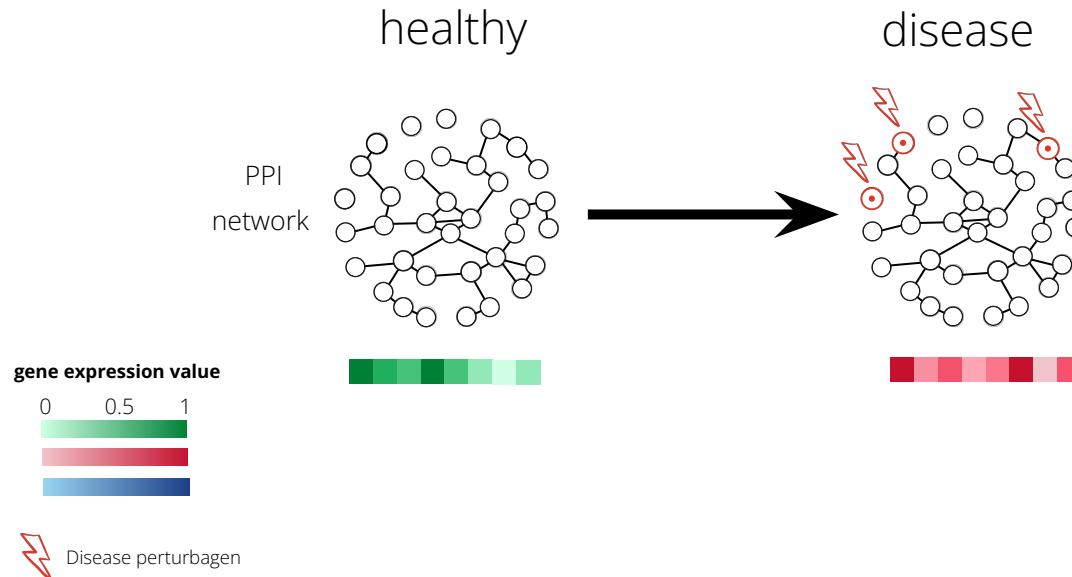


gene expression value



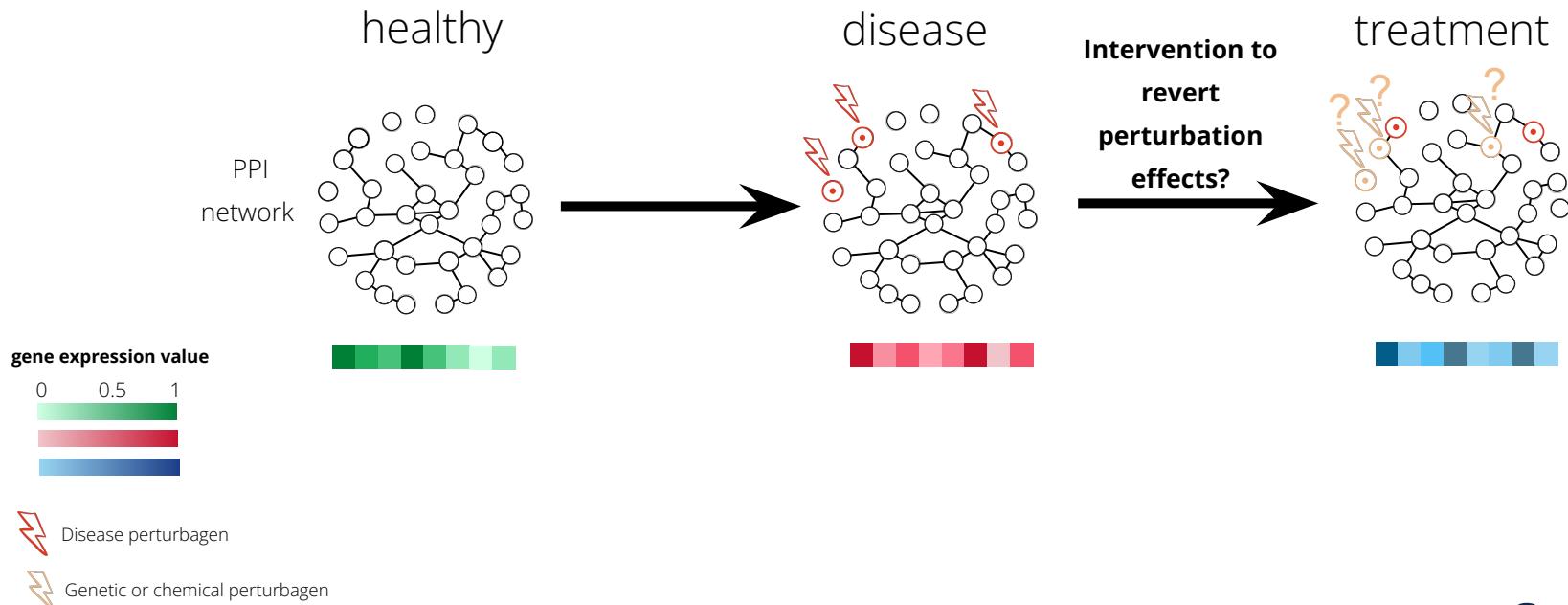
Problem formulation

Problem formulation



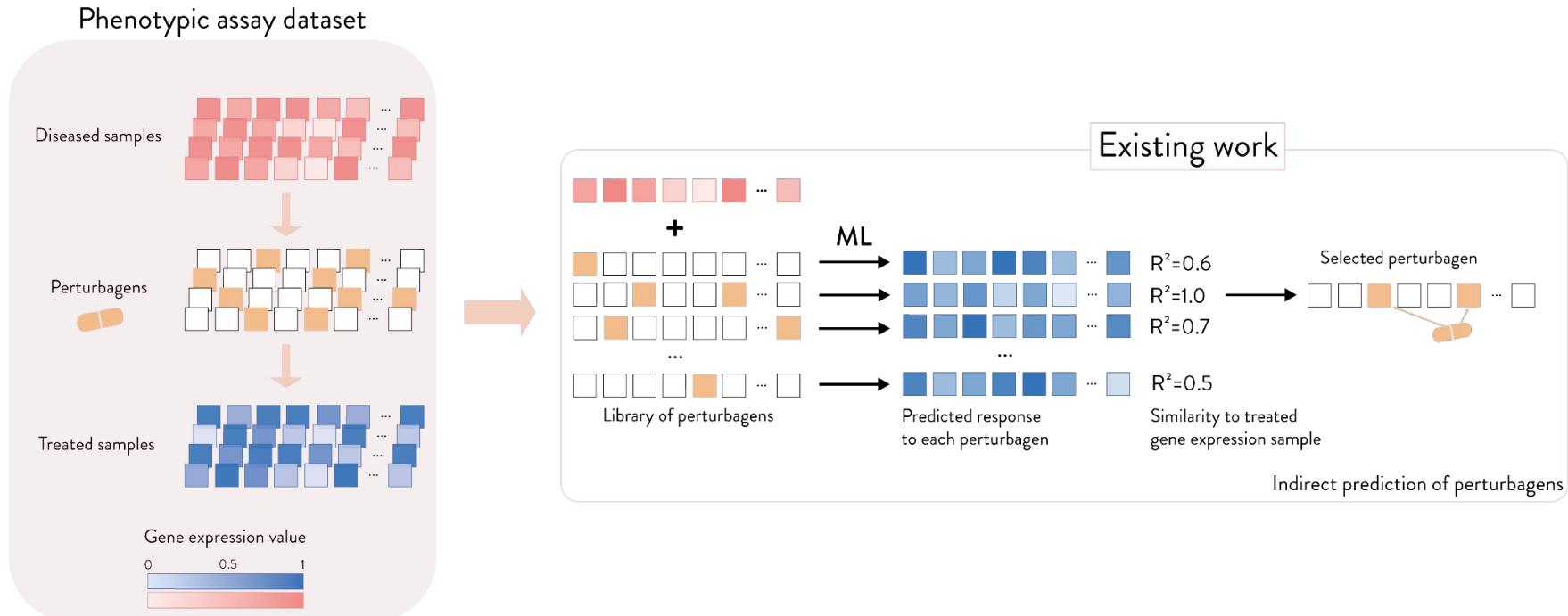
Problem formulation

Problem formulation

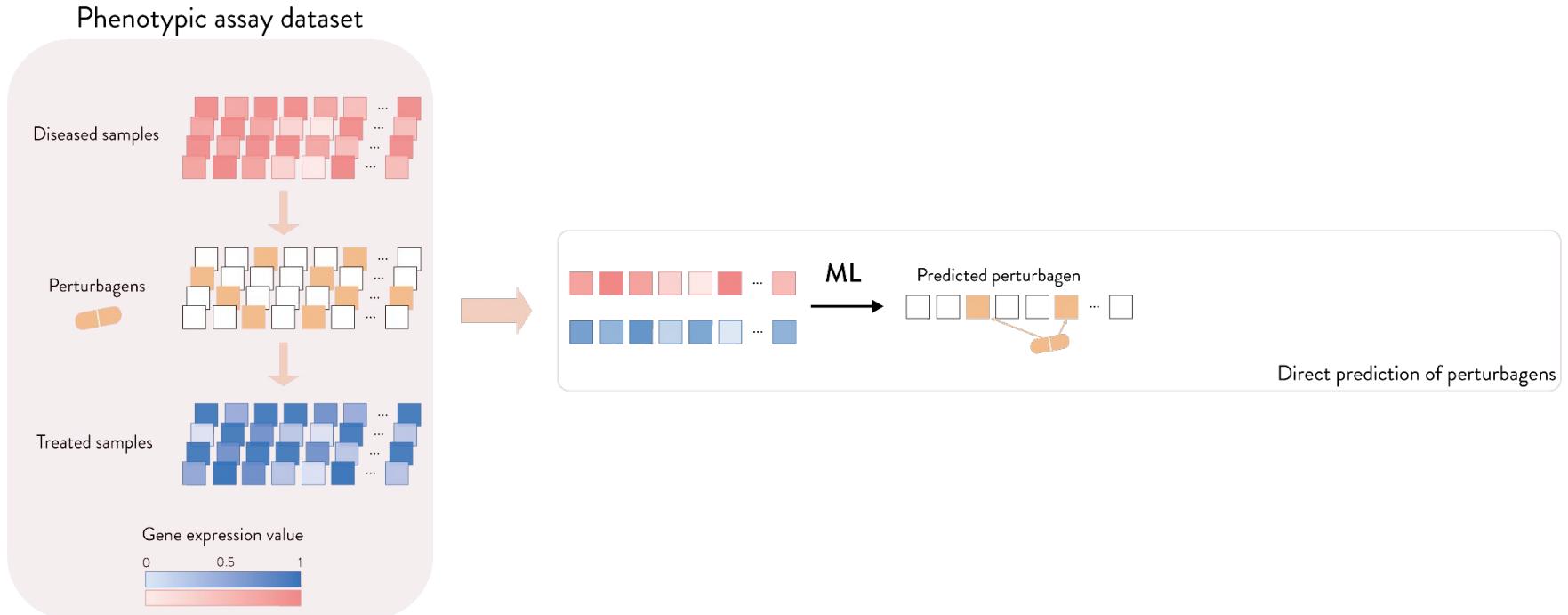


Existing work

Existing work



Existing work



Combinatorial prediction of therapeutic targets using a causally inspired neural network

Causal inference primer

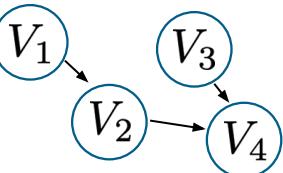
Structural Causal Models

Consists of a **causal graph** and **structural equations**

exogenous variables

U

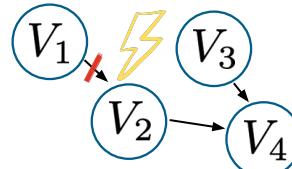
quantify causal
effects shown in
causal graph



$$\left\{ \begin{array}{l} V_1 = f(U_{V_1}) \\ V_2 = f(U_{V_2}, V_1) \\ V_3 = f(U_{V_3}) \\ V_4 = f(U_{V_4}, V_2, V_3) \end{array} \right.$$

intervention
→
 $do(V_2) = v_2$

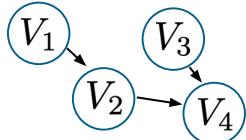
"do"
operator



**causal Markov
condition**

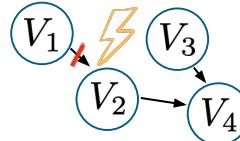
$$\begin{aligned} V_1 &= f(U_{V_1}) \\ V_2 &= v_2 \\ V_3 &= f(U_{V_3}) \\ V_4 &= f(U_{V_4}, V_2, V_3) \end{aligned}$$

Causal inference primer



$$U_{V_2} = U_{V_4} = 1$$

$$\begin{aligned} V_1 &= U_{V_1} \\ V_2 &= \text{AND}(U_{V_2}, V_1) \\ V_3 &= U_{V_3} \\ V_4 &= \text{AND}(U_{V_4}, V_2, V_3) \end{aligned}$$



$$U_{V_2} = U_{V_4} = 1$$

$$\begin{aligned} V_1 &= U_{V_1} \\ V_2 &= v_2 \\ V_3 &= U_{V_3} \\ V_4 &= \text{AND}(U_{V_4}, V_2, V_3) \end{aligned}$$

V_1	V_2	V_3	V_4
1	1	0	0
1	1	1	1
0	0	1	0
1	1	1	1
0	0	0	0
1	1	0	1

=

?

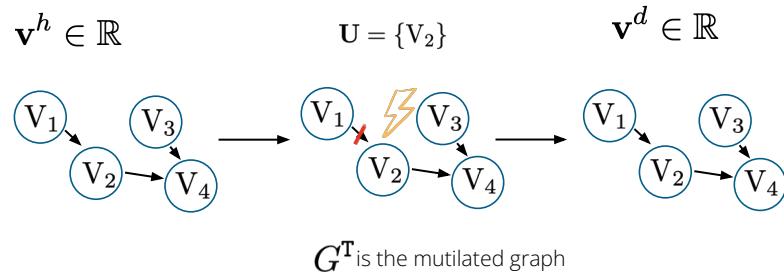
observational distribution
 $P^o(V_1, V_2, V_3, V_4)$

V_1	V_2	V_3	V_4
0	1	0	0
1	0	1	1
1	0	1	0
0	1	1	1
1	0	0	0
0	1	0	0

interventional distribution
 $P^i(V_1, V_2, V_3, V_4)$

Problem formulation

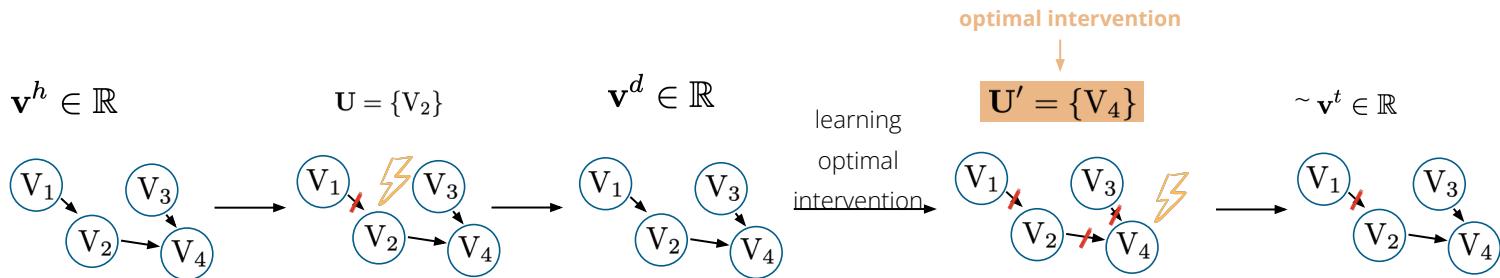
Let $\mathbf{M} = \langle \mathbf{E}, \mathbf{V}, \mathcal{F}, P(\mathbf{E}) \rangle$ be an SCM associated with causal graph G . Dataset: $\mathcal{T} = \{T_1, \dots, T_m\}$



$\mathbf{T} = \langle \mathbf{v}^h, \mathbf{U}, \mathbf{v}^d \rangle$ **disease intervention data**

Problem formulation

Let $\mathbf{M} = \langle \mathbf{E}, \mathbf{V}, \mathcal{F}, P(\mathbf{E}) \rangle$ be an SCM associated with causal graph G . Dataset: $\mathcal{T} = \{T_1, \dots, T_m\}$



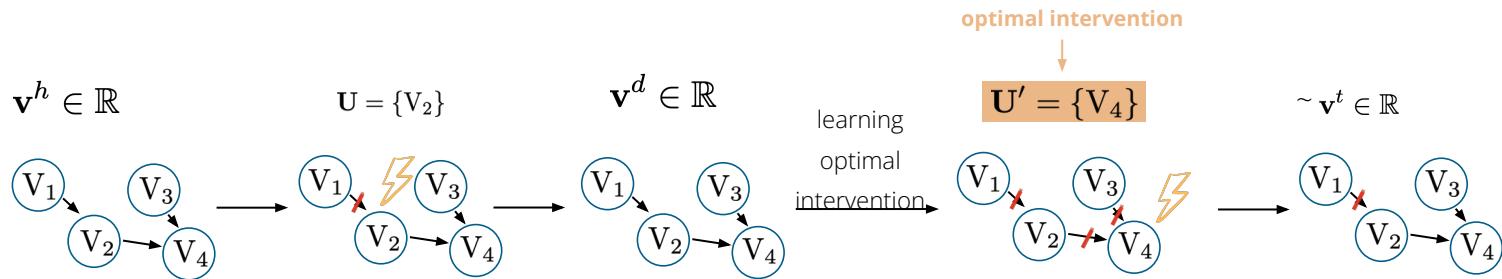
$\mathcal{T} = \langle \mathbf{v}^h, \mathbf{U}, \mathbf{v}^d \rangle$

disease intervention data

combinatorial prediction

Problem formulation

Let $\mathbf{M} = \langle \mathbf{E}, \mathbf{V}, \mathcal{F}, P(\mathbf{E}) \rangle$ be an SCM associated with causal graph G . Dataset: $\mathcal{T} = \{T_1, \dots, T_m\}$



$\mathcal{T} = \langle \mathbf{v}^h, \mathbf{U}, \mathbf{v}^d \rangle$ **disease intervention data**

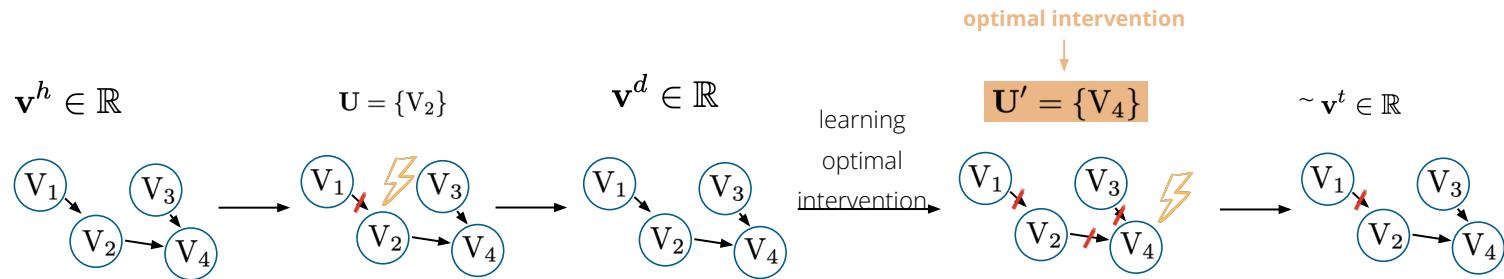
combinatorial prediction

$$\text{argmax}_{\mathbf{U}'} P^{G^{\mathcal{T}}}(\mathbf{V} = \mathbf{v}^t | \text{do}(\mathbf{U}'))$$

"do" operator

Problem formulation

Let $\mathbf{M} = \langle \mathbf{E}, \mathbf{V}, \mathcal{F}, P(\mathbf{E}) \rangle$ be an SCM associated with causal graph G . Dataset: $\mathcal{T} = \{T_1, \dots, T_m\}$



$T = \langle \mathbf{v}^h, U, \mathbf{v}^d \rangle$ disease intervention data

which in the absence of unobserved confounders is:

$$\begin{aligned} & argmax_{\mathbf{U}'} P^{G^T}(\mathbf{V} = \mathbf{v}^t | do(\mathbf{U}')) \\ & \quad \text{"do" operator} \\ & argmax_{\mathbf{U}'} P^{G^{T'}}(\mathbf{V} = \mathbf{v}^t | \mathbf{U}') \end{aligned}$$

Problem formulation

causal inference

The goal is to find the perturbation set \mathbf{U}'
 with the highest likelihood of achieving
 treated state $\mathbf{v}^t \in \mathbb{R}$

$$\operatorname{argmax}_{\mathbf{U}'} P^{G^T} (\mathbf{V} = \mathbf{v}^t | \operatorname{do}(\mathbf{U}'))$$

↑
 "do" operator

which under markovianity is:

our goal $\operatorname{argmax}_{\mathbf{U}'} P^{G^T} (\mathbf{V} = \mathbf{v}^t | \mathbf{U}')$

↑
 sets of random
 variables

Problem formulation

causal inference

The goal is to find the perturbation set \mathbf{U}' with the highest likelihood of achieving treated state $\mathbf{v}^t \in \mathbb{R}$

$$\operatorname{argmax}_{\mathbf{U}'} P^{G^T} (\mathbf{V} = \mathbf{v}^t | \operatorname{do}(\mathbf{U}'))$$

↑
"do" operator

which under markovianity is:

our goal $\operatorname{argmax}_{\mathbf{U}'} P^{G^T} (\mathbf{V} = \mathbf{v}^t | \mathbf{U}')$

↑
sets of random variables

graph representation learning

$$f : G^{T'}, \mathbf{x}^d, \mathbf{x}^t \rightarrow \operatorname{argmax}_{\mathbf{U}'} P^{G^{T'}} (\mathbf{x} = \mathbf{x}^t | \mathbf{x}^d, \mathbf{U}')$$

↑
vectors, sets

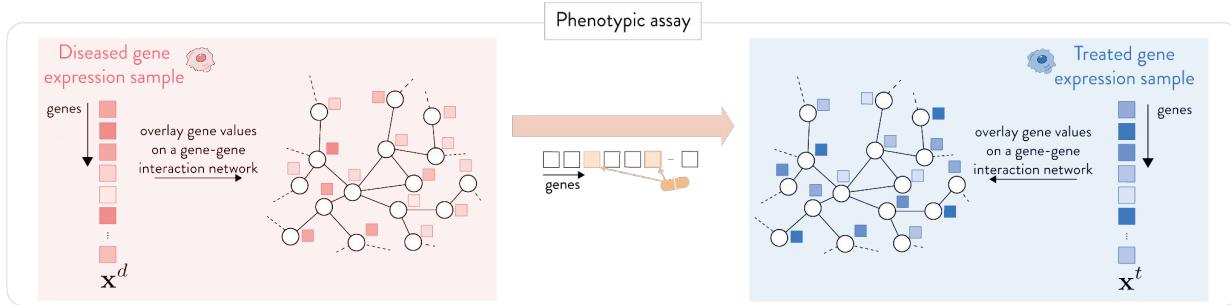
Intuitively, we need to search the space of all interventions and score how good they are

- intractable

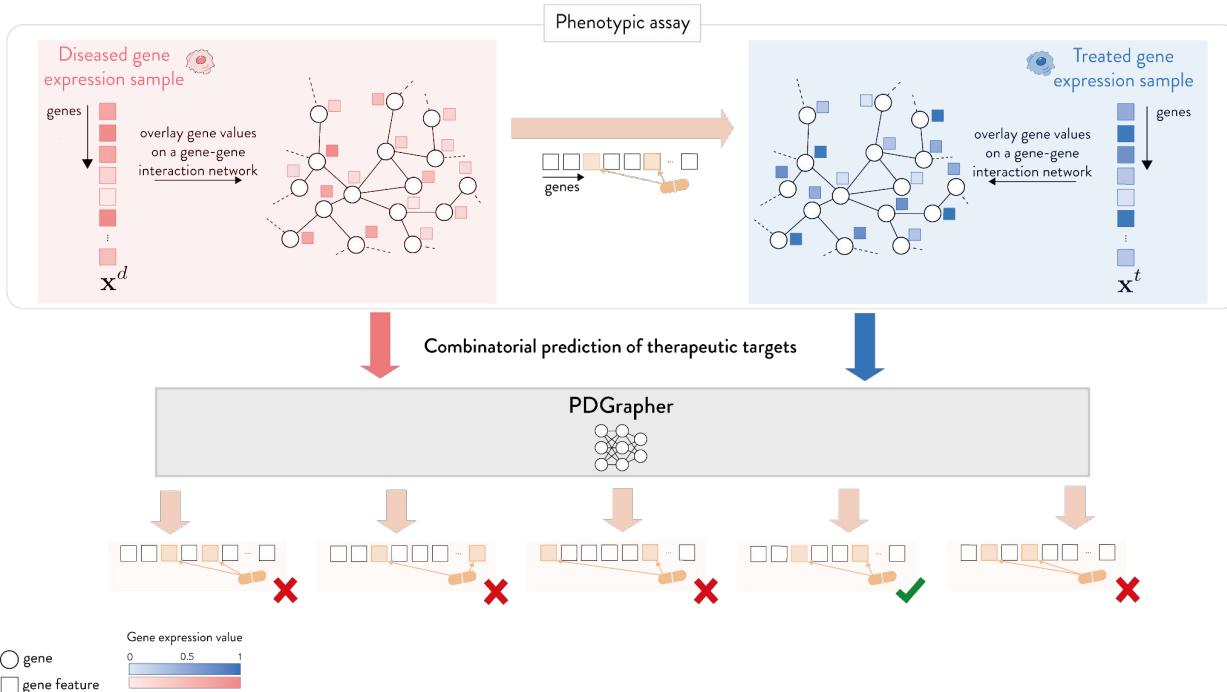
What we propose:

- a smart way to search the space of interventions (**intervention discovery module**)
- a way to score each intervention (**response prediction module**)

Our solution: PDGrapher

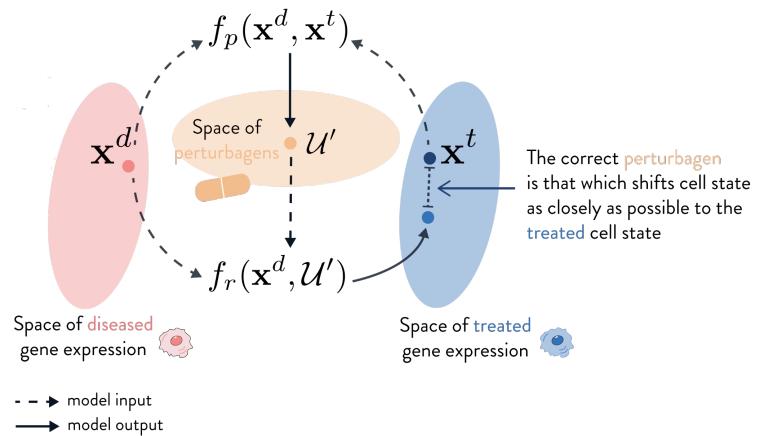


Our solution: PDGrapher

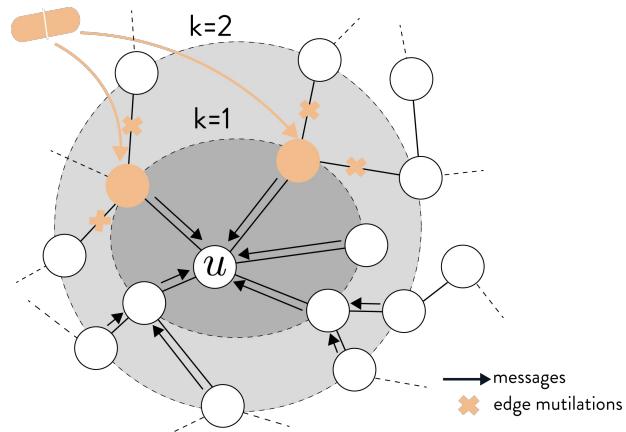


Our solution: PDGrapher

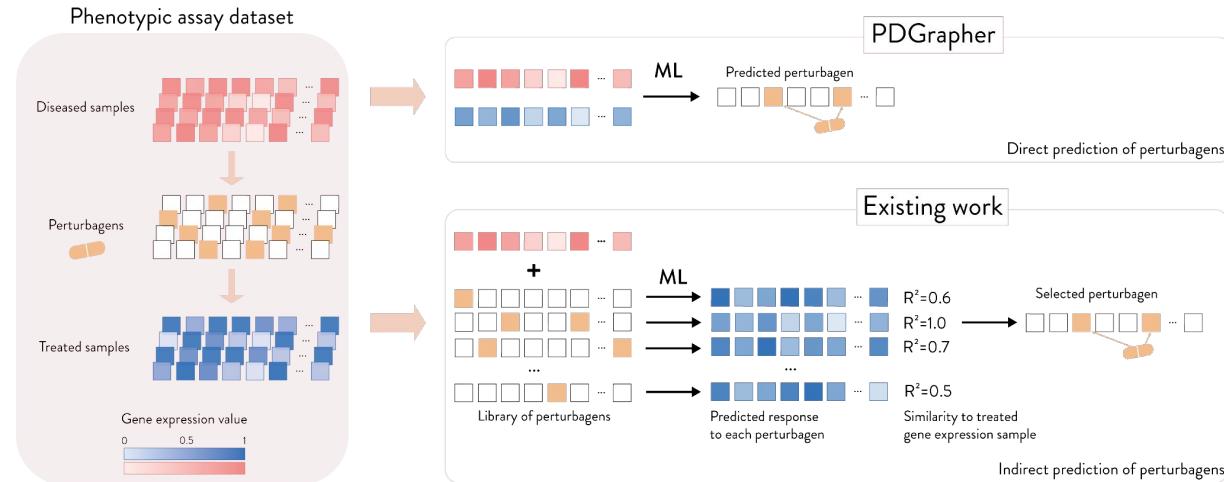
Optimization



Mutilations in message passing



Our solution: PDGrapher



PDGrapher

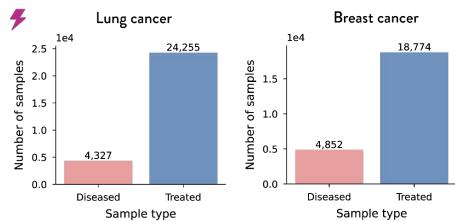
1. Models perturbations explicitly
2. Performs direct prediction of perturbagens

PDGrapher

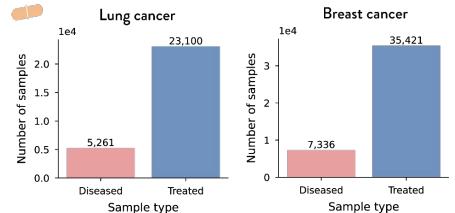
Datasets



single-gene knockouts



compounds



The Human Reference Protein Interactome Mapping Project

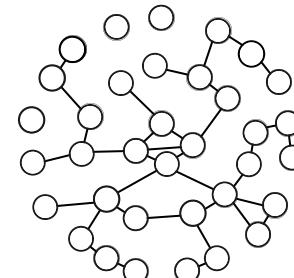
DISEASE NETWORKS

BioGRID 4.4

Uncovering disease-disease relationships through the incomplete interactome

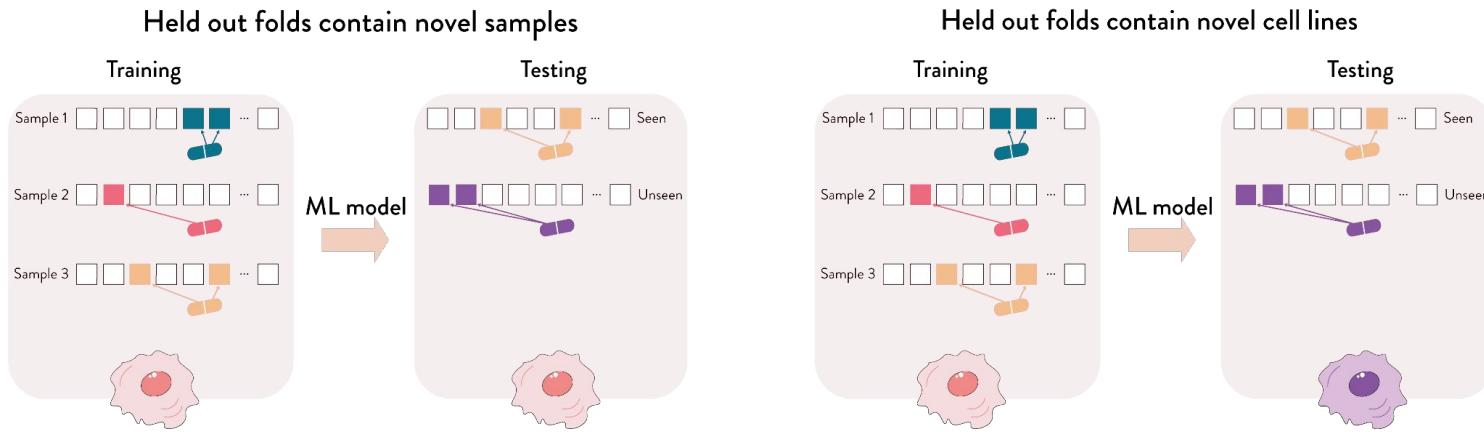
Jörg Menche, Amitabh Sharma, Maksim Kitsak, Susan Dina Ghiassian, Marc Vidal, Joseph Loscalzo, Albert-László Barabási*

Protein-Protein Interaction network



10,716 nodes
149,270 edges

Training settings



Baselines

- **Random:** Randomly ranks list of genes
- **Cancer targets:** First M genes are targets of existing cancer drugs, N-M genes are sorted randomly
- **Cancer genes:** First M genes are cancer-associated genes, N-M genes are sorted randomly
- **scGen:** Indirect prediction of perturbagen.

Results - Chemical-PPI, Random splitting

Prediction for a single perturbagen:

Gene 1
Gene 2
Gene 3
Gene 4
Gene 5
...
Gene 10,000

We want
ground truth
targets to be at
the top



Results - Chemical-PPI, Random splitting

Prediction for a single perturbagen:

Gene 1 ↗ If ground truth target
is predicted at the top
→ ranking = 1

Gene 2

Gene 3

Gene 4

Gene 5

...

Gene 10,000

Results - Chemical-PPI, Random splitting

Prediction for a single perturbagen:

Gene 1

Gene 2

Gene 3

Gene 4

Gene 5

...

Gene 10,000

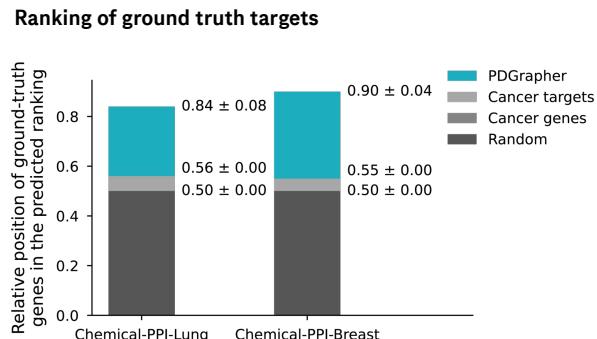


If ground truth target
is predicted at the
bottom → ranking = 0

Results - Chemical-PPI, Random splitting

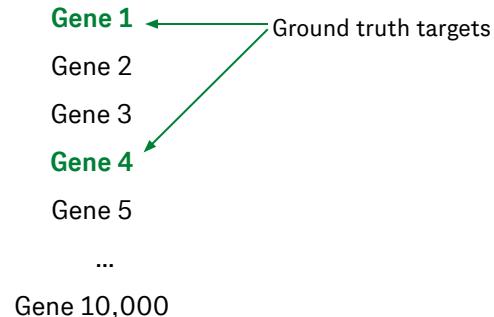
Prediction for a single perturbagen:

- Gene 1 ↗ If ground truth target is predicted at the top → ranking = 1
- Gene 2
- Gene 3
- Gene 4
- Gene 5
- ...
- Gene 10,000 ↗ If ground truth target is predicted at the bottom → ranking = 0



Results - Chemical-PPI, Random splitting

Prediction for a single perturbagen with 2 targets:



Results - Chemical-PPI, Random splitting

Prediction for a single perturbagen with 2 targets:

We look at the top 2 predicted targets

{ Gene 1 ←
Gene 2
Gene 3
Gene 4 → Ground truth targets
Gene 5

...

Gene 10,000

Results - Chemical-PPI, Random splitting

Prediction for a single perturbagen **with 2 targets**:

We look at the top 2 predicted targets { Gene 1
Gene 2
Gene 3



If at least one of the top K predicted genes is a target → partially accurate prediction

Gene 4

Gene 5

...

Gene 10,000

Results - Chemical-PPI, Random splitting

Prediction for a single perturbagen with 2 targets:

We look at the top 2 predicted targets { Gene 1
Gene 2
Gene 3

If at least one of the top K predicted genes is a target → partially accurate prediction

Gene 4

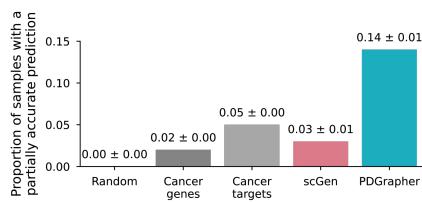
Gene 5

...

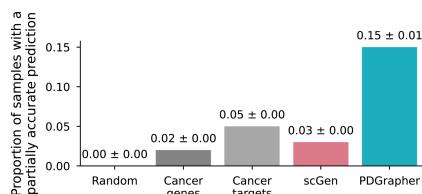
Gene 10,000

Proportion of samples with a partially accurate prediction

Chemical-PPI-Lung



Chemical-PPI-Breast



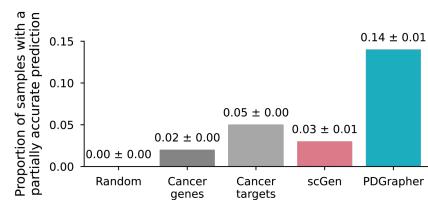
Results - Chemical-PPI, Random splitting

Prediction for a single perturbagen with 2 targets:

We look at the top 2 predicted targets { Gene 1 If at least one of the top K predicted genes is a target → partially accurate prediction
 Gene 2
 Gene 3
Gene 4
 Gene 5
 ...
 Gene 10,000

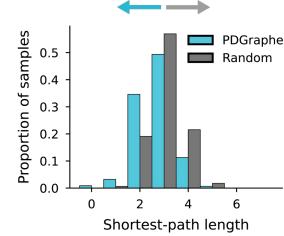
Proportion of samples with a partially accurate prediction

Chemical-PPI-Lung



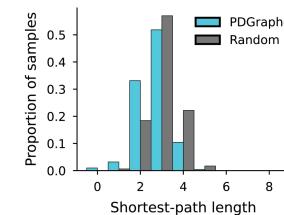
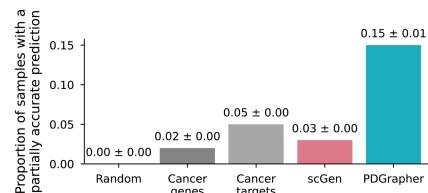
Shortest-path length between ground-truth targets and predicted targets

Chemical-PPI-Lung



Chemical-PPI-Breast

Chemical-PPI-Breast

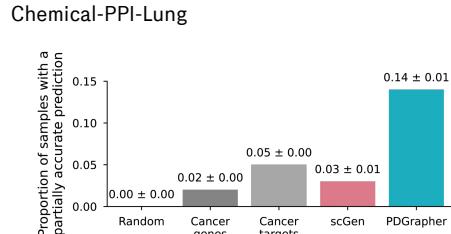


Results - Chemical-PPI, Random splitting

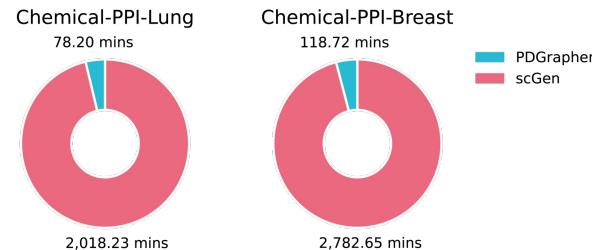
Prediction for a single perturbagen with 2 targets:

We look at the top 2 predicted targets { Gene 1 ↗ If at least one of the top K predicted genes is a target → partially accurate prediction
 Gene 2
 Gene 3
Gene 4
 Gene 5
 ...
 Gene 10,000

Proportion of samples with a partially accurate prediction



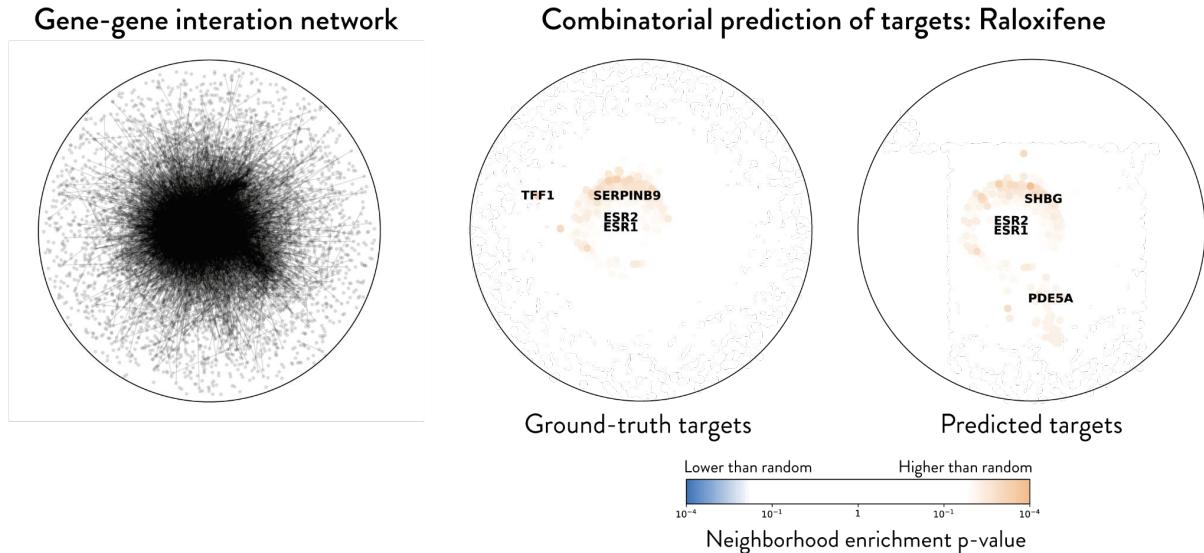
Training time of a single model



Our model needs 15x less time to train!

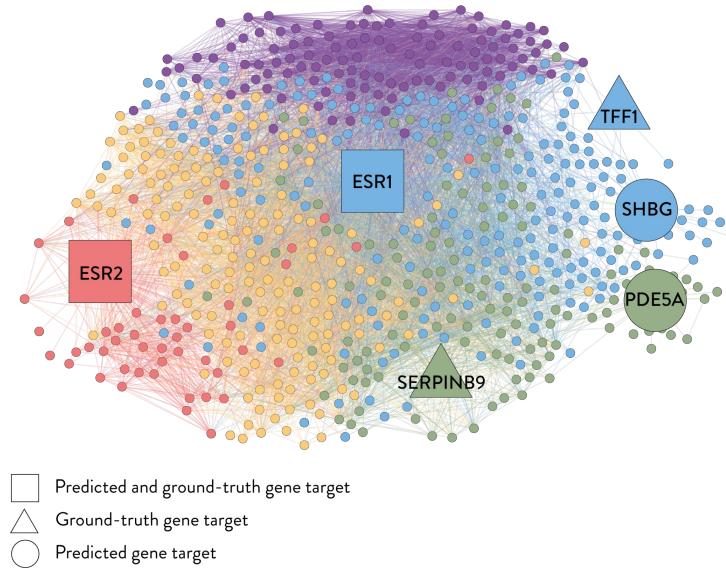
Results - Chemical-PPI, Random splitting

Raloxifene



Results - Chemical-PPI, Random splitting

Raloxifene

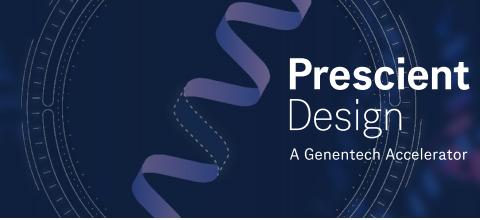


Raloxifene is a second-generation selective estrogen receptor modulator. Targets: ESR1, ESR2, SERPINB9, TFF1.

Additional predictions:

- Sex Hormone Binding Globulin (SHBG)
- Phosphodiesterase 5A (PDE5A)

Takeaways



- Reformulate phenotype-driven lead discovery
- Propose PDGrapher for combinatorial prediction of therapeutic targets, which:
 - Outperforms mechanistic and traditional baselines
 - Performs direct prediction of perturbagens
 - Is computationally efficient
 - Illuminates mode of action of predicted perturbagens

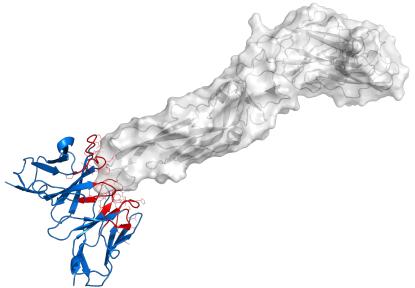


Misc

My research at Prescient

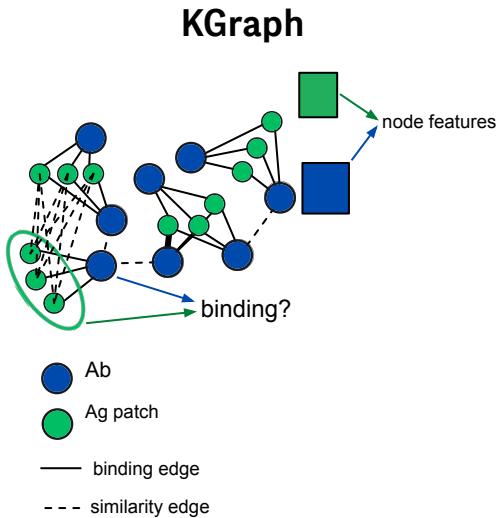
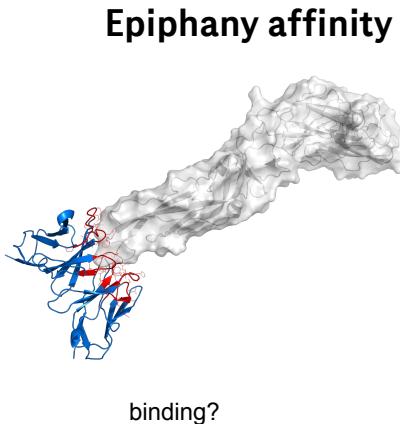


Epiphany affinity



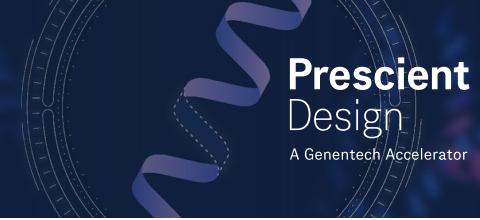
binding?

My research at Prescient

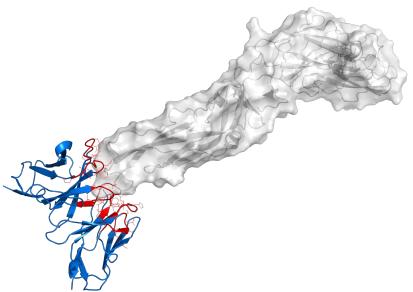


- Given an Ab node and an Ag patch node, predict link existence (binding)
- Given an Ab node and a set of Ag patch nodes, predict link existence (binding)

My research at Prescient

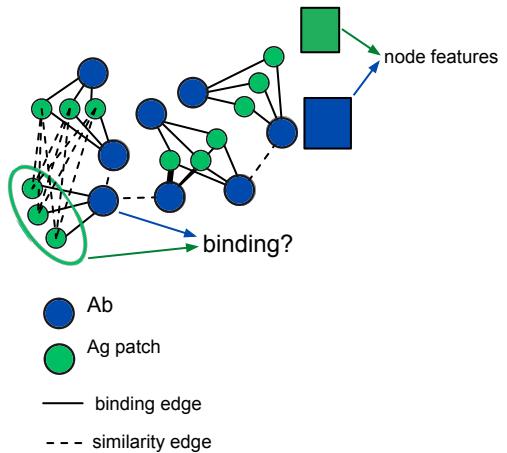


Epiphany affinity



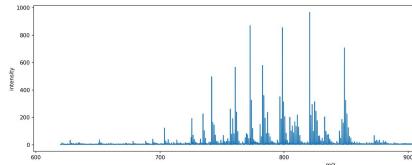
binding?

KGraph

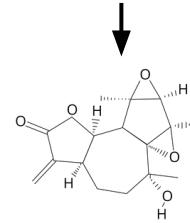


Structure elucidation

Mass spectrometry



Molecule



THE END