# AIM 2: Artificial Intelligence in Medicine II

## Harvard - BMIF 203 and BMI 702, Spring 2025

Lecture 9: Knowledge graph learning, Building multimodal knowledge graphs, Structure-inducing pre-training, Knowledge-based foundation models

**HARVARD**
MEDICAL SCHOOL

**Kempner** INSTITUTE — For the Study of Natural & Artificial Intelligence at Harvard University
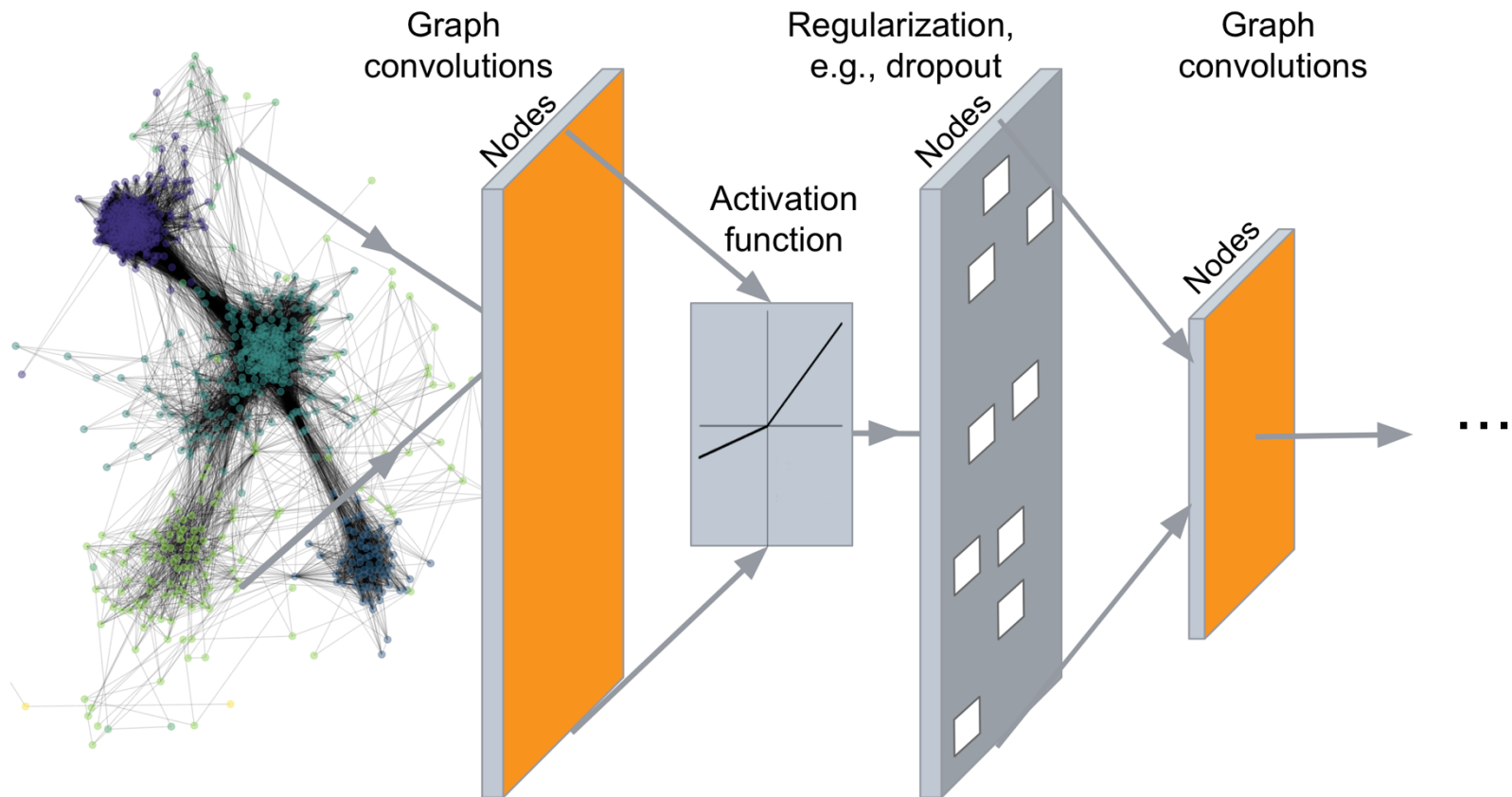
**BROAD** INSTITUTE

Marinka Zitnik
marinka@hms.harvard.edu

# Deep graph representation learning

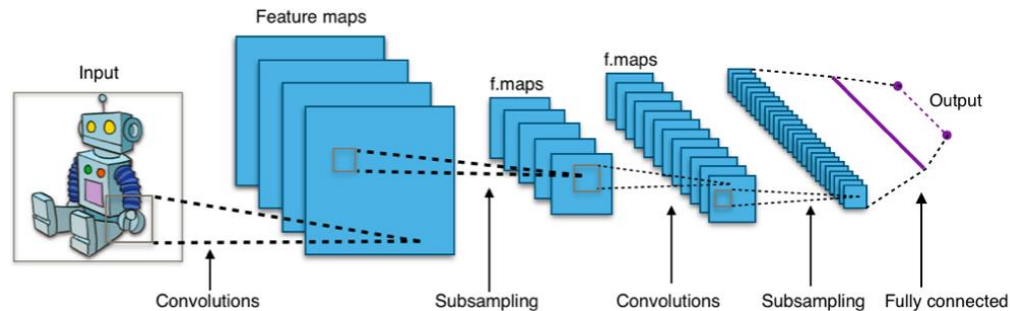## Recap of message passing neural network (MPNN) strategies

# Graph neural networks

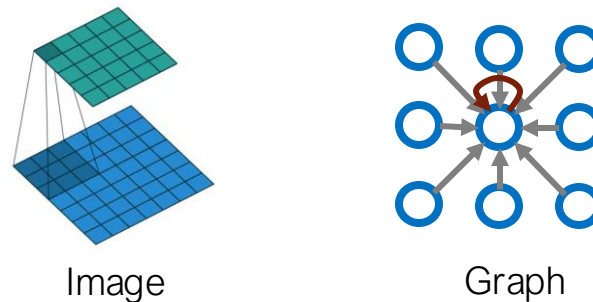- Encoder: Multiple layers of nonlinear transformation of graph structure
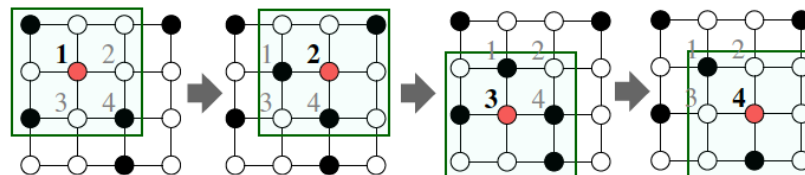
# Convolutional networks

- Let's start with convolutional networks on an image:
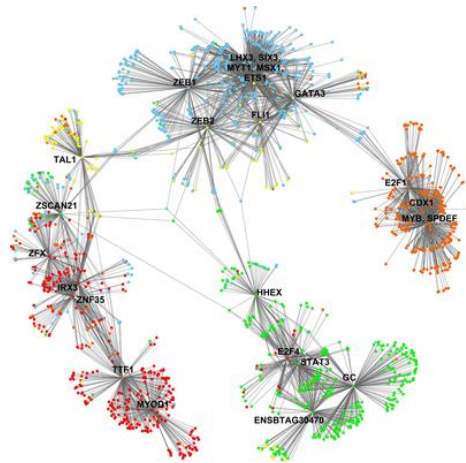


- Single convolutional network with a 3x3 filter:



Image                    Graph

- Transform information (or messages) from the neighbors and combine them: $\sum_i W_i \, h_i$
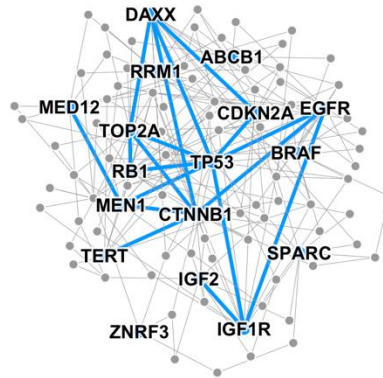
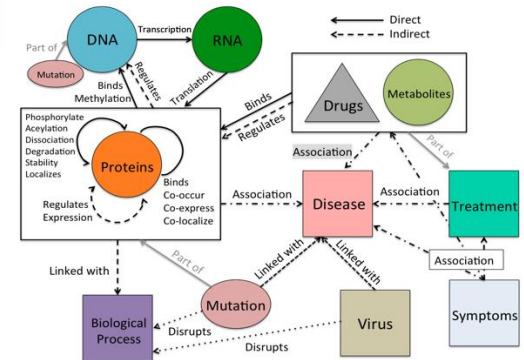# Real world graphs

- But what if your graphs look like this?



Gene interaction network          Disease pathways          Biomedical knowledge graphs
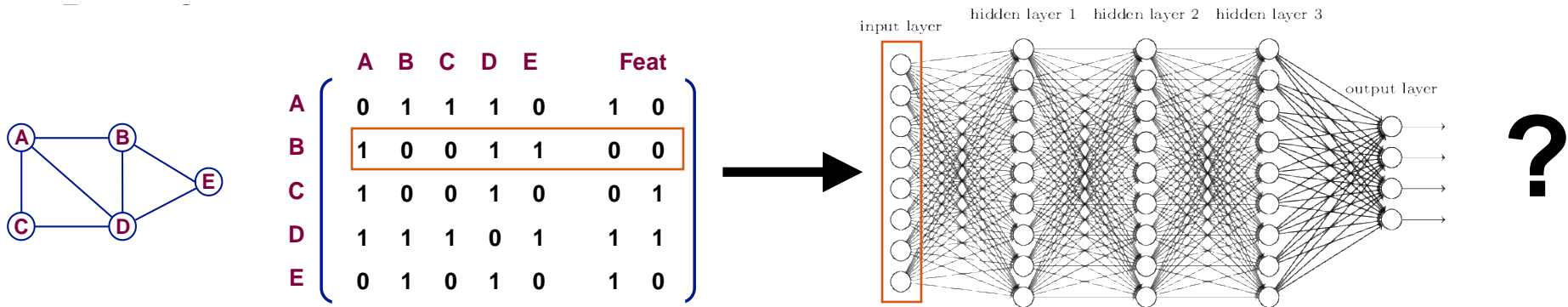
- Examples:
  - Biological or medical networks
  - Social networks
  - Information networks
  - Knowledge graphs
  - Communication networks
  - Web graphs
  - …

# Naïve approach

- Join adjacency matrix and features
- Feed them into a deep neural network:



- **Issues with this idea:**
  - $O(N)$ parameters
  - Not applicable to graphs of different sizes
  - Not invariant to node ordering

# Graph neural networks

- Intuition:
  - Each node's neighborhood defines a computational graph
  - Generate node embeddings based on local network neighborhoods
- Neighborhood aggregation:



- Model can be of arbitrary depth
  - Nodes have embeddings at each layer
  - Layer 0 embedding of node $u$ is its input features $X_u$
- Basic neighborhood aggregation: Average information from neighbors and apply a neural network

# Basic approach



| Graph convolutions | Regularization, e.g., dropout | Graph convolutions |

Nodes · Activation function · Nodes · Nodes · ...

Initial 0-th layer embeddings are equal to node features

Previous layer embedding of $v$

$$\mathbf{h}_v^0 = \mathbf{x}_v$$

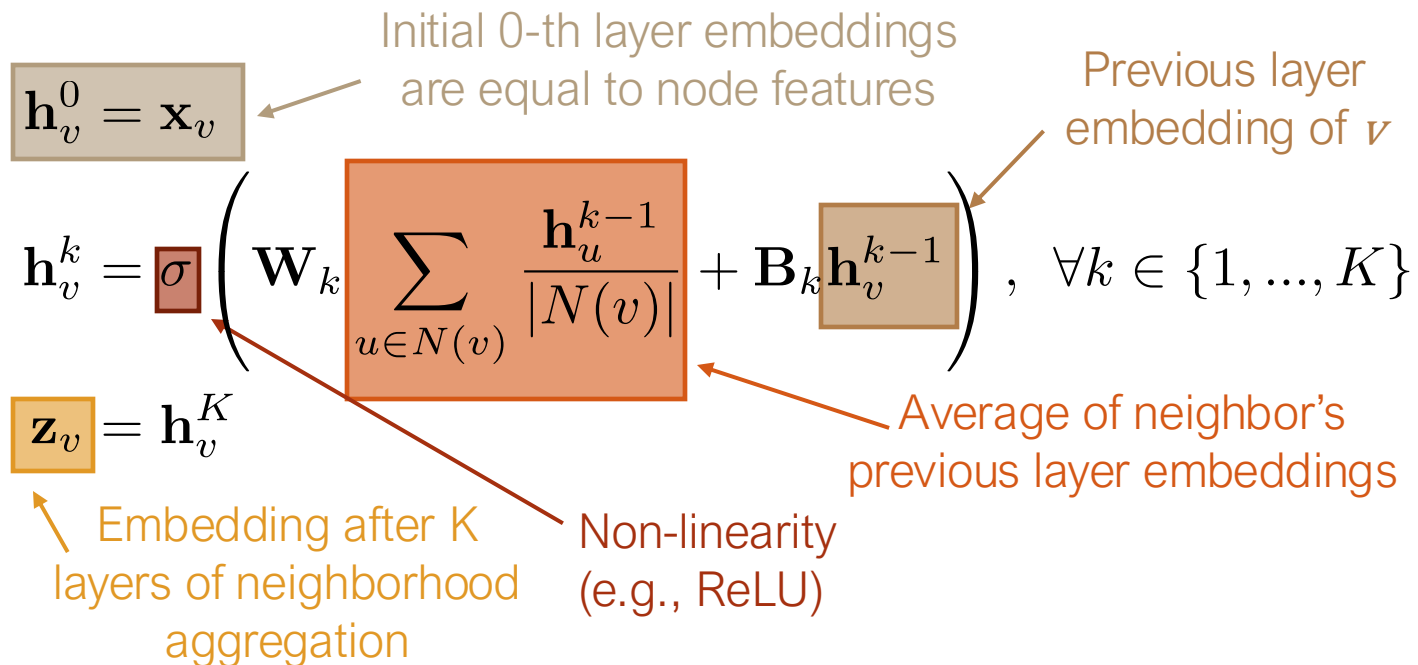$$\mathbf{h}_v^k = \sigma\left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1}\right), \ \ \forall k \in \{1, ..., K\}$$

$$\mathbf{z}_v = \mathbf{h}_v^K$$

Average of neighbor's previous layer embeddings

Embedding after K layers of neighborhood aggregation

Non-linearity (e.g., ReLU)

8

# Basic approach



Graph convolutions · Regularization, e.g., dropout · Graph convolutions · Nodes · Activation function

$$\mathbf{h}_v^0 = \mathbf{x}_v$$

trainable weight matrices
(i.e., what we learn)

$$\mathbf{h}_v^k = \sigma \left( \mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right), \quad \forall k \in \{1, ..., K\}$$

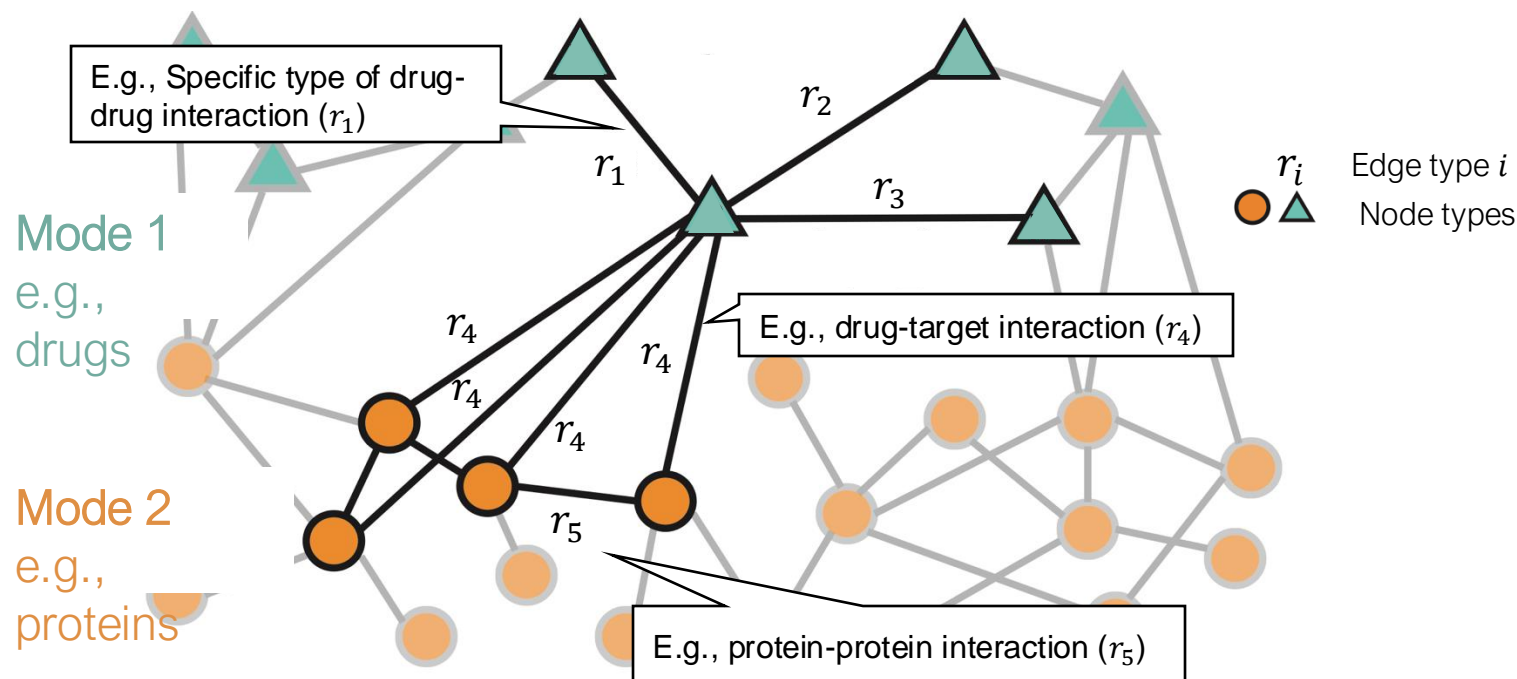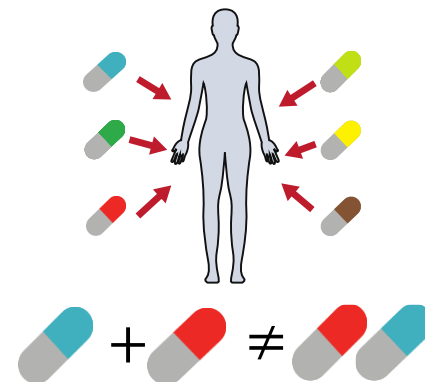$$\mathbf{z}_v = \mathbf{h}_v^K$$

We can feed these into any loss function and run stochastic gradient descent to train the weight parameters

# Polypharmacy modeling and antibiotic discovery

# Application: Drug combinations

- Combinatorial explosion
  - >13 million possible combinations of 2 drugs
  - >20 billion possible combinations of 3 drugs
- Non-linear & non-additive interactions
  - Different effect than the additive effect of individual drugs
- Small subsets of patients
  - Side effects are interdependent
  - No info on drug combinations not yet used in patients



E.g., Specific type of drug-drug interaction ($r_1$)

$r_2$

$r_1$

$r_3$

$r_i$ — Edge type $i$

Node types

Mode 1
e.g.,
drugs

$r_4$

$r_4$

$r_4$

E.g., drug-target interaction ($r_4$)

$r_4$

Mode 2
e.g.,
proteins

$r_5$

E.g., protein-protein interaction ($r_5$)

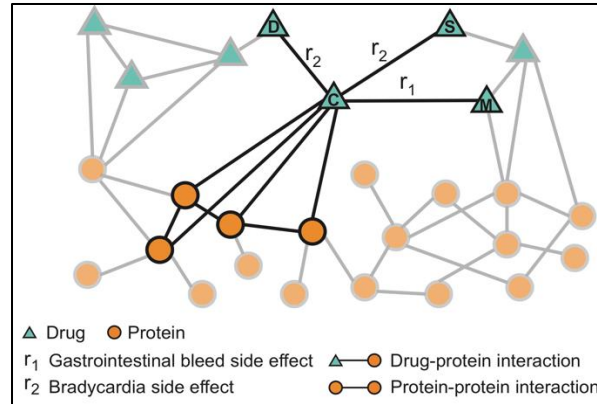# Polypharmacy dataset



- Molecular, drug, and patient data for all US-approved drugs
  - 4,651,131 drug-drug edges: Patient data from adverse event system, tested for confounders [FDA]
  - 18,596 drug-protein edges
  - 719,402 protein-protein edges: Physical, metabolic enzyme-coupled, and signaling interactions
  - Drug and protein features: drugs' chemical structure, proteins' membership in pathways
- This is a multimodal network with over 5 million edges separated into 1,000 different edge types

# Experimental setup



- Drug   ○ Protein
- $r_1$ Gastrointestinal bleed side effect   ▲——○ Drug-protein interaction
- $r_2$ Bradycardia side effect   ○——○ Protein-protein interaction
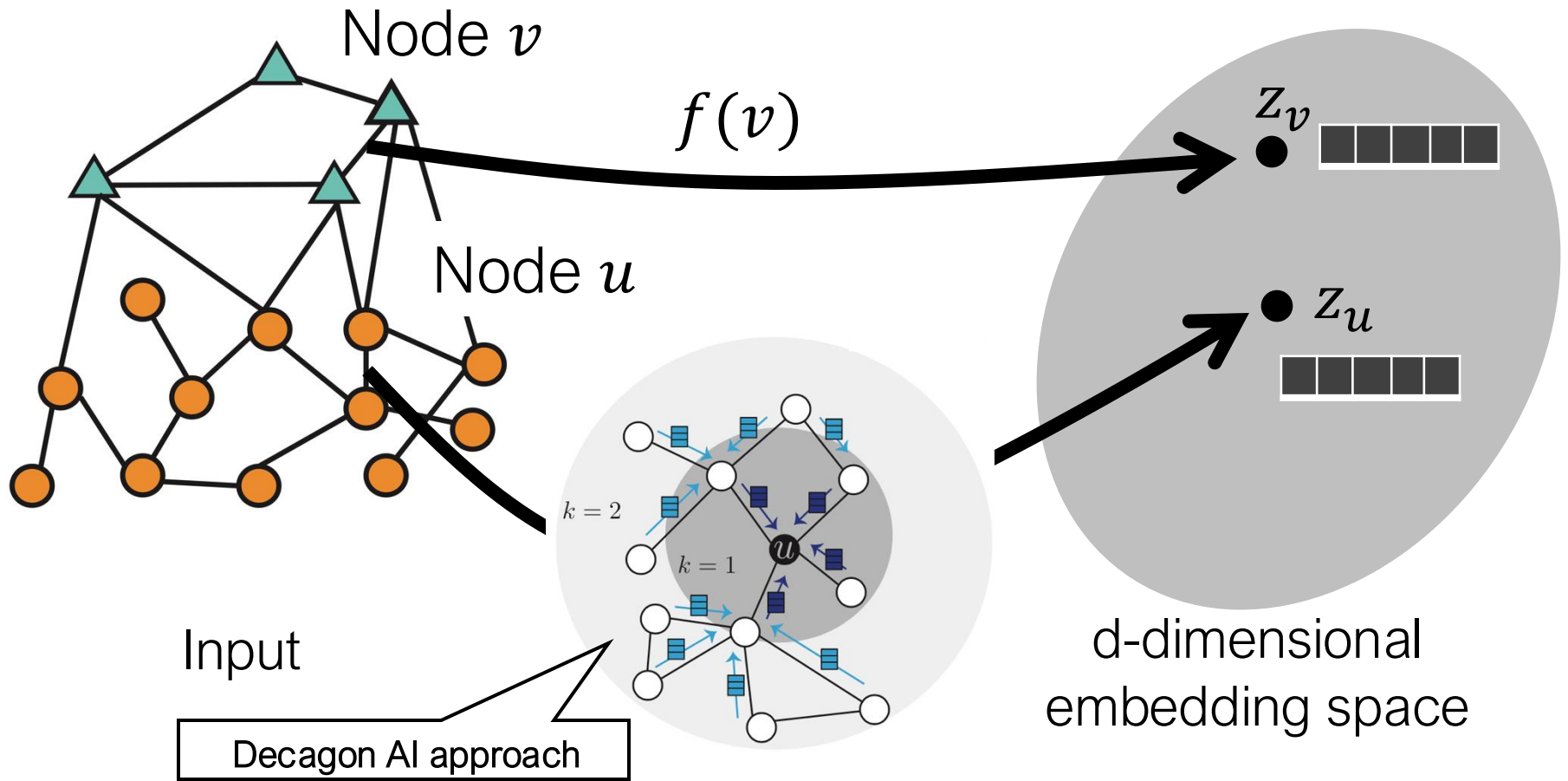
■ Two main stages:

1. Learn an embedding for every node in polypharmacy network
2. Predict a score for every drug-drug, drug-protein, protein-protein pair in the test set based on the embeddings



$r_2$ (breakdown of muscle tissue)

Simvastatin

Ciprofloxacin
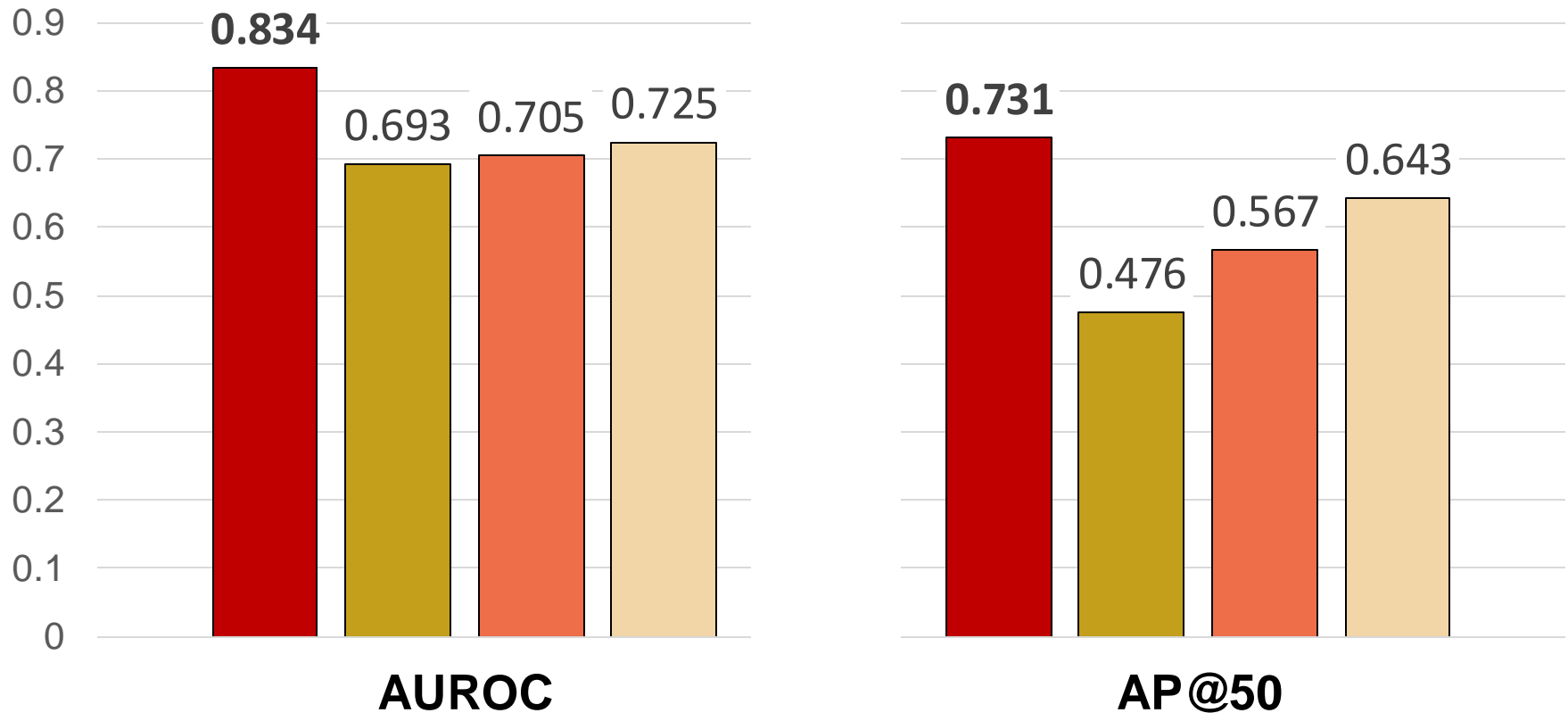
Example: How likely will Simvastatin and Ciprofloxacin, when taken together, break down muscle tissue?

# Approach: Graph Neural Network



Node $v$

$f(v)$

$z_v$

Node $u$

$z_u$

$k = 2$

$k = 1$

Input

Decagon AI approach

d-dimensional embedding space

Map nodes to d-dimensional embeddings such that nodes with similar network neighborhoods are embedded close together

# Results: Polypharmacy side effects



0.9 0.8 0.7 0.6 0.5 0.4 0.3 0.2 0.1 0

**0.834** 0.693 0.705 0.725

**AUROC**

**0.731** 0.476 0.567 0.643

**AP@50**

■ Decagon
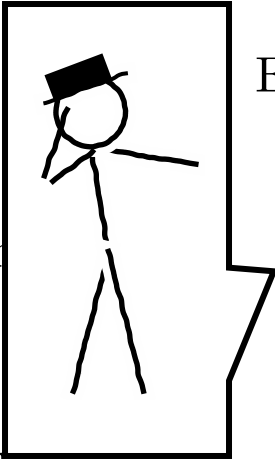■ RESCAL Tensor Factorization [Nickel et al., ICML'11]
■ Multi-relational Factorization [Perros, Papalexakis et al., KDD'17]
□ Shallow Network Embedding [Zong et al., Bioinformatics'17]
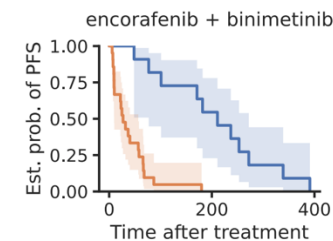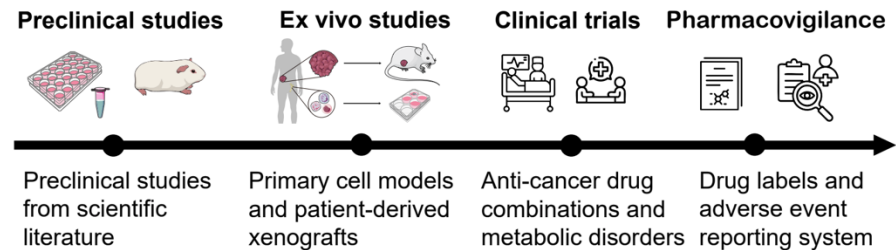
# Results: Polypharmacy side effects

**Approach:**

1) Train deep model on data generated **prior to 2012**
2) How many **predictions** have been **confirmed after 2012**?

| Rank | Drug | Drug | Side effect | Evidence found |
|------|------|------|-------------|----------------|
| 1 | Pyrimethamine | Aliskiren | Sarcoma | |
| 2 | Tigecycline | Bimatoprost | Autonomic | |
| 3 | Telangiectases | Omeprazole | Dacarbazine | |
| 4 | Tolcapone | Pyrimethamine | Blood brain | |
| | | | eadache | |
| | | | ular acidosis | |
| 7 | Anag | Azelaic acid | Cerebral thrombosis | |
| 8 | Atorvastatin | Amlodipine | Muscle inflammation | |
| 9 | Aliskiren | Tioconazole | Breast inflammation | |
| 10 | Estradiol | Nadolol | Endometriosis | |

*Case Report*
**Severe Rhabdomyolysis due to Presumed Drug Interactions between Atorvastatin with Amlodipine and Ticagrelor**
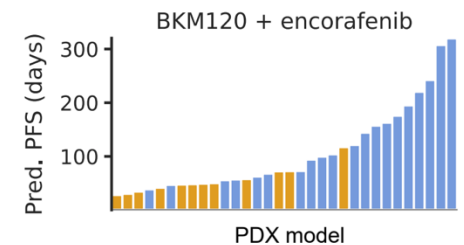
# Multimodal AI predicts clinical outcomes of drug combinations from preclinical data

- **Personalized oncology therapy:** Predicts leukemia drug combination responses using patient genomics and xenograft models

- **Drug safety & transporter interactions:** Identifies organ-specific toxicities and transporter-based risks for early drug development

- **Oncology drug combinations & polypharmacy:** Assesses PARP inhibitor safety, differentiating approved vs. investigational regimens

- **Metabolic disease insights:** Ranked Resmetirom among the safest candidates for MASH, supporting FDA approval
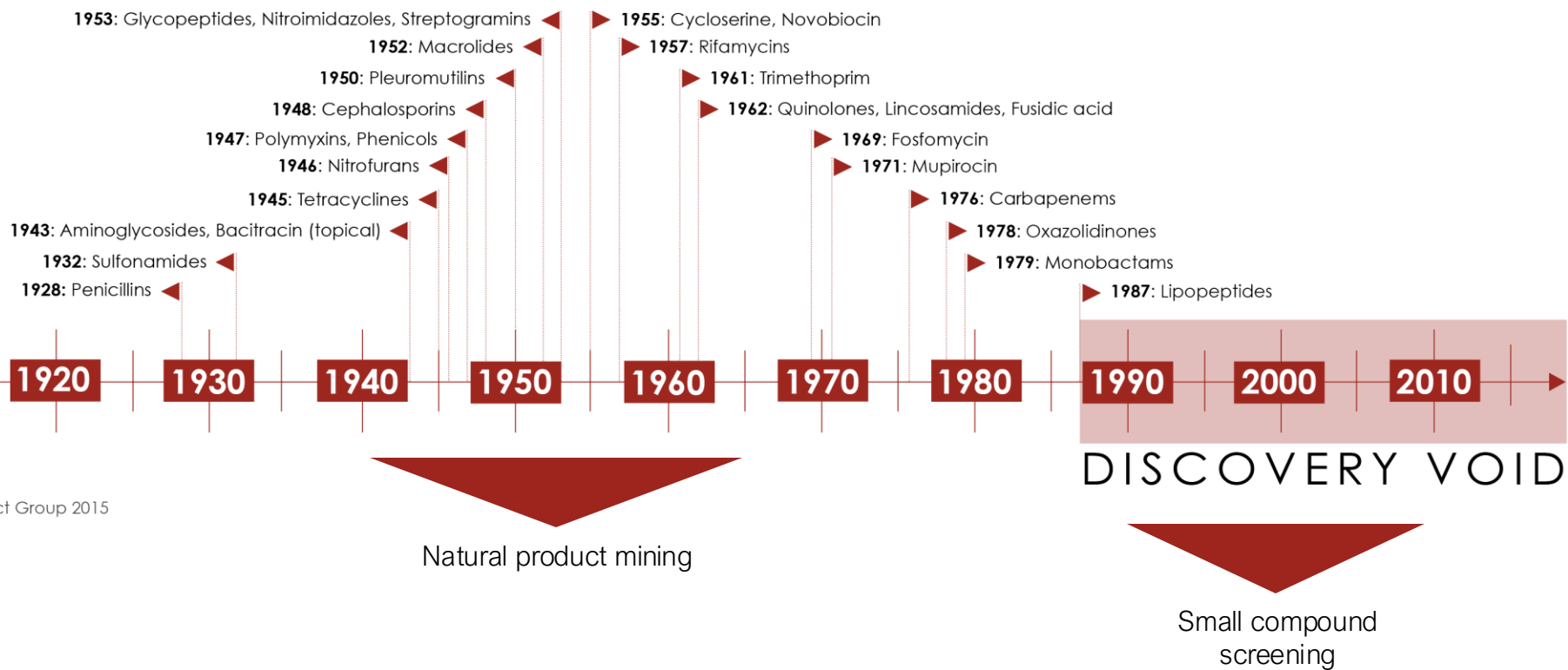
Multimodal AI predicts clinical outcomes of drug combinations from preclinical data, arXiv 2025

17

# Application: Antibiotic discovery



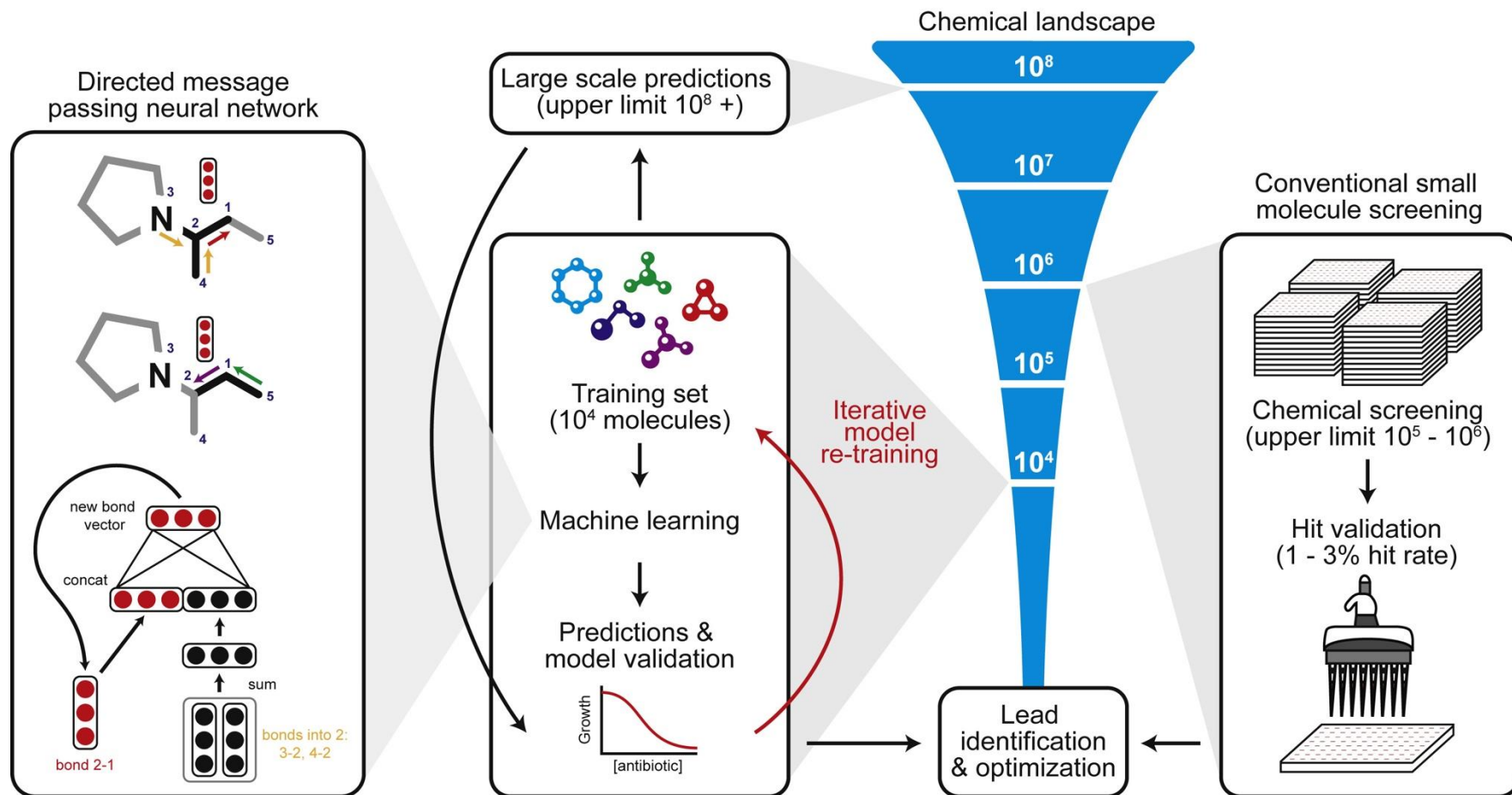1953: Glycopeptides, Nitroimidazoles, Streptogramins
1952: Macrolides
1950: Pleuromutilins
1948: Cephalosporins
1947: Polymyxins, Phenicols
1946: Nitrofurans
1945: Tetracyclines
1943: Aminoglycosides, Bacitracin (topical)
1932: Sulfonamides
1928: Penicillins

1955: Cycloserine, Novobiocin
1957: Rifamycins
1961: Trimethoprim
1962: Quinolones, Lincosamides, Fusidic acid
1969: Fosfomycin
1971: Mupirocin
1976: Carbapenems
1978: Oxazolidinones
1979: Monobactams
1987: Lipopeptides

1920  1930  1940  1950  1960  1970  1980  1990  2000  2010

DISCOVERY VOID

© ReAct Group 2015

Natural product mining
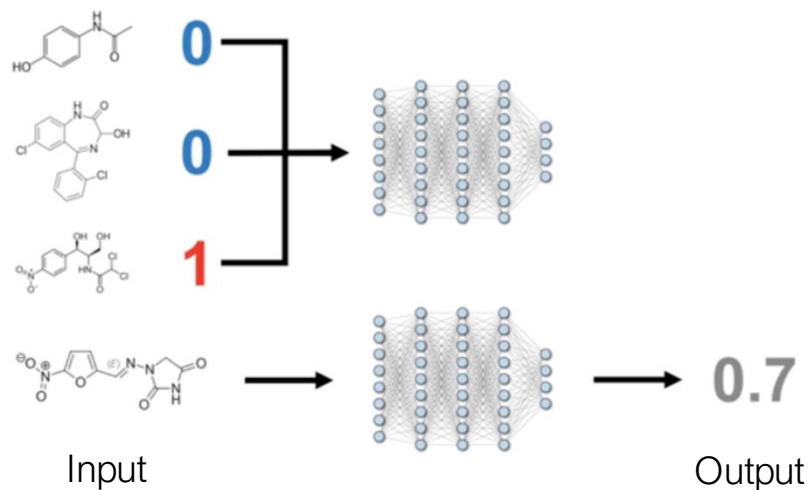
Small compound screening

18

# GNNs to learn molecular structure



Directed message passing neural network model iteratively (1) learns representations of molecules and (2) optimizes the representations for predicting growth inhibition
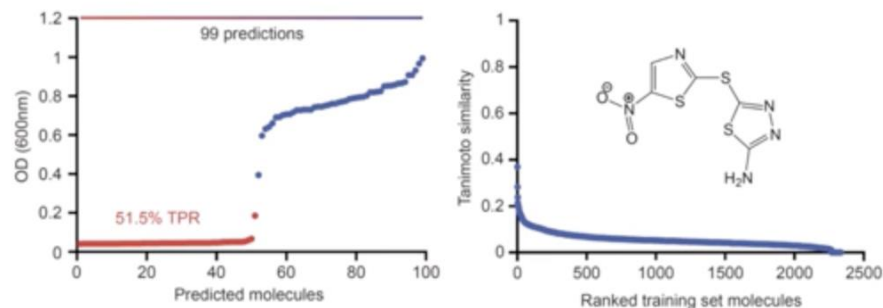
# Experimental setup

## Training Dataset
(Human Medicines and Natural Products)



Input — Output

Data: 2,335 molecules (human medicines and natural products) screened for growth inhibition
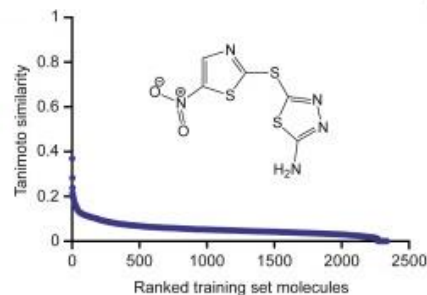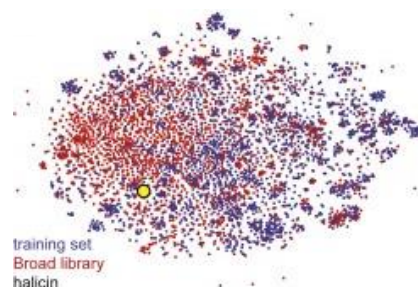
## Empirical Validation
(Broad Repurposing Hub)



Data: 6,111 molecules (at various stages of investigation for human diseases) in Broad Repurposing Hub

Task: Test top 99 predictions & prioritize based on similarity to known antibiotics or predicted toxicity

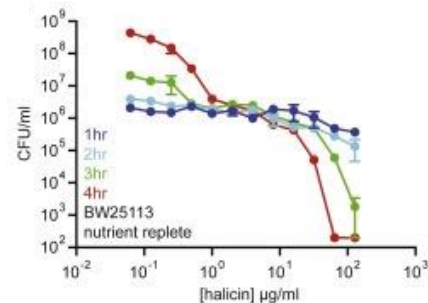# Results
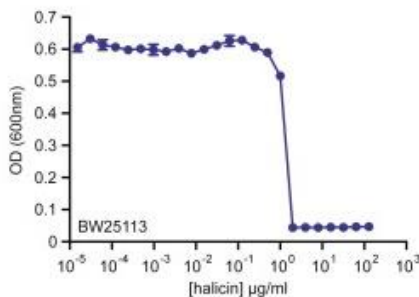
Halicin was developed to be an anti-diabetic drug, but the development was discontinued due to poor results in testing.

Halicin predicted to be antibacterial

Halicin against *E. coli*

Halicin against *M. tuberculosis*

# Results

Halicin's efficacy in murine models of infection



Validated against ~6K molecules to identify halicin, a novel candidate antibiotic

# Rare disease diagnosis

# Rare disease diagnosis

- Rare diseases affect between 300-400 million or 1 in 20 people worldwide, yet each disease affects no more than 50 per 100,000 individuals

- Diagnosis is challenging due to the heterogeneity of clinical presentations and small patient populations

# Rare disease diagnosis

- Many patients suffering from rare diseases are undiagnosed. It currently takes 4-5 years on average for patients to receive a diagnosis.



**4-5 years on average to diagnosis**

## Can AI help shorten diagnostic odysseys for rare disease patients?

Haendel et al. How many rare diseases are there? *Nature Review Drug Discovery* (2020).
Wakap et al. Estimating cumulative point prevalence of rare diseases: analysis of the Orphanet database. *EJHG* (2020).

# Diagnostic odysseys

- Over 7,000 rare diseases, each affects < 200,000 patients in the US
  - Most diseases are phenotypically heterogeneous
  - Front-line clinicians might lack disease experience, resulting in expensive clinical workups for patients across multiple years
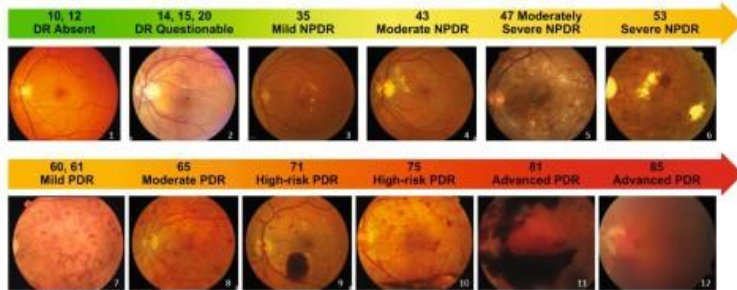  - Diagnosis often requires a specialist, sub-specialist, or multi-disciplinary referrals

- On average, the long search for a **rare disease diagnosis takes 5 to 7 years, 4 up to 8 physicians, and 2 to 3 misdiagnoses**

- Diagnostic delay is so pervasive that it leads to problems for patients:
  - Undergoing **redundant testing and procedures**
  - Substantial delay in obtaining disease-appropriate management and **inappropriate therapies**
  - **Irreversible disease progression**—time window for intervention can be missed leading to disease progression
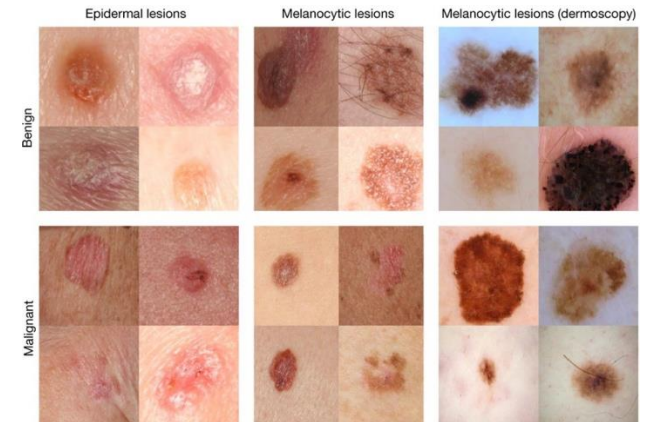
## Can AI help shorten diagnostic odysseys for rare disease patients?

# AI models for disease diagnosis

Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopahty in Retinal Fundus Photographs (*JAMA*)



Dermatologist-level Classification of Skin Cancer (*Nature*)



Evaluation and Accurate Diagnoses of Pediatric Diseases Using AI (*Nature Medicine*)

# AI models for disease diagnosis
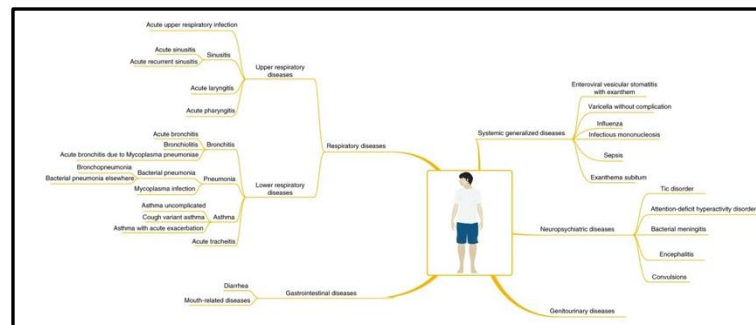
Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopahty in Retinal Fundus Photographs (*JAMA*)

128,175 retinal images
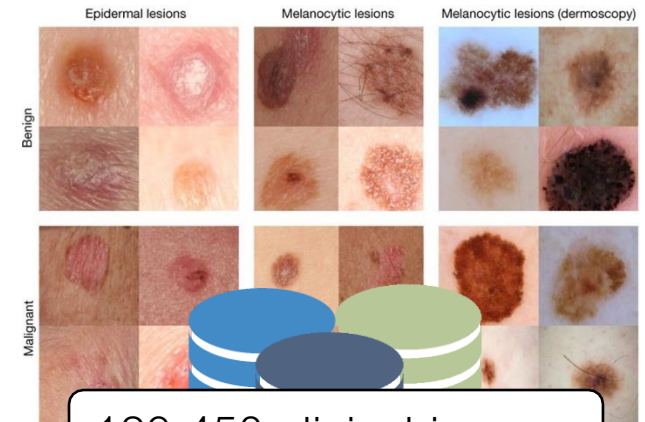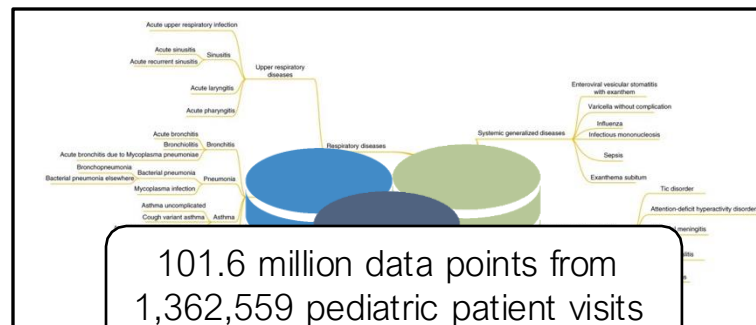
Dermatologist-level Classification of Skin Cancer (*Nature*)

129,450 clinical images

Evaluation and Accurate Diagnoses of Pediatric Diseases Using AI (*Nature Medicine*)

101.6 million data points from 1,362,559 pediatric patient visits

# Rare disease diagnosis is hard!

- Deep learning models trained (via supervised learning) on large labeled datasets can achieve **near-expert clinical accuracy for common diseases**

- Existing models require **labeled datasets with thousands of diagnosed patients per disease**:
    - Diabetic retinopathy: deep neural net on 128 K retinal images
    - Skin lesions: deep neural net on 129 K clinical images of skin cancers
    - Childhood diseases: deep neural net on 1 M pediatric patient visits

The challenge with rare diseases is fundamental — **datasets are three orders of magnitude smaller than in other uses of AI for medical diagnosis**
Needed is an entirely new approach to making AI-based rare disease diagnosis possible. This is for two primary reasons:
- Rare disease diagnosis cannot simply be solved by recruiting/labeling more patients because of high disease heterogeneity and low disease prevalence
- Rare disease diagnosis cannot be solved by supervised deep learning because the models cannot extrapolate to novel genetic diseases and atypical disease presentations

# Rare disease diagnosis is hard!

1. Need to extrapolate beyond training distribution to never-before-seen genetic conditions
2. Approaches must be able to learn from limited data given the lack of large annotated datasets of patients with rare genetic diseases & low prevalence of each disease

Low overlap of phenotypes, causal genes, and diseases across patients

UDN Undiagnosed Diseases Network

Of 465 diagnosed patients in the UDN, there are 378 unique causal genes and 299 unique diseases.
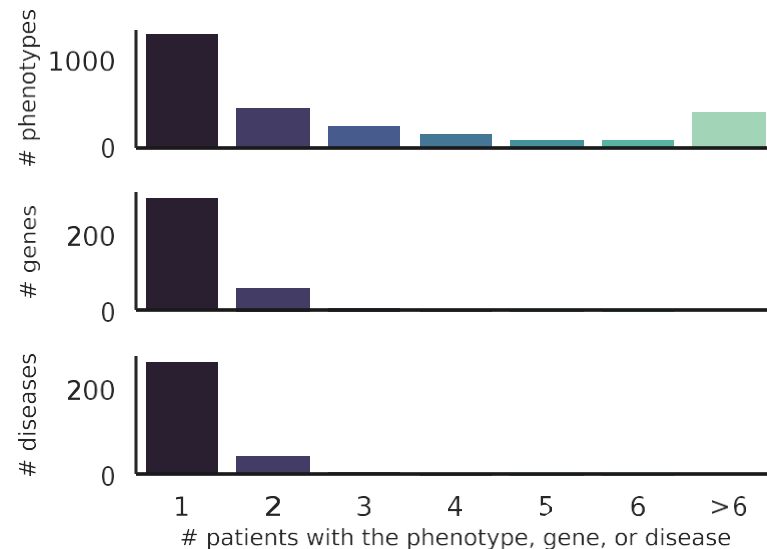


# patients with the phenotype, gene, or disease

# Rare disease diagnosis is hard!

1. Need to extrapolate beyond training distribution to never-before-seen genetic conditions
2. Approaches must be able to learn from limited data given the lack of large annotated datasets of patients with rare genetic diseases & low prevalence of each disease

UDN Undiagnosed Diseases Network

Of 465 diagnosed patients in the UDN, there are 378 unique causal genes and 299 unique diseases.

**Phenotypic heterogeneity**

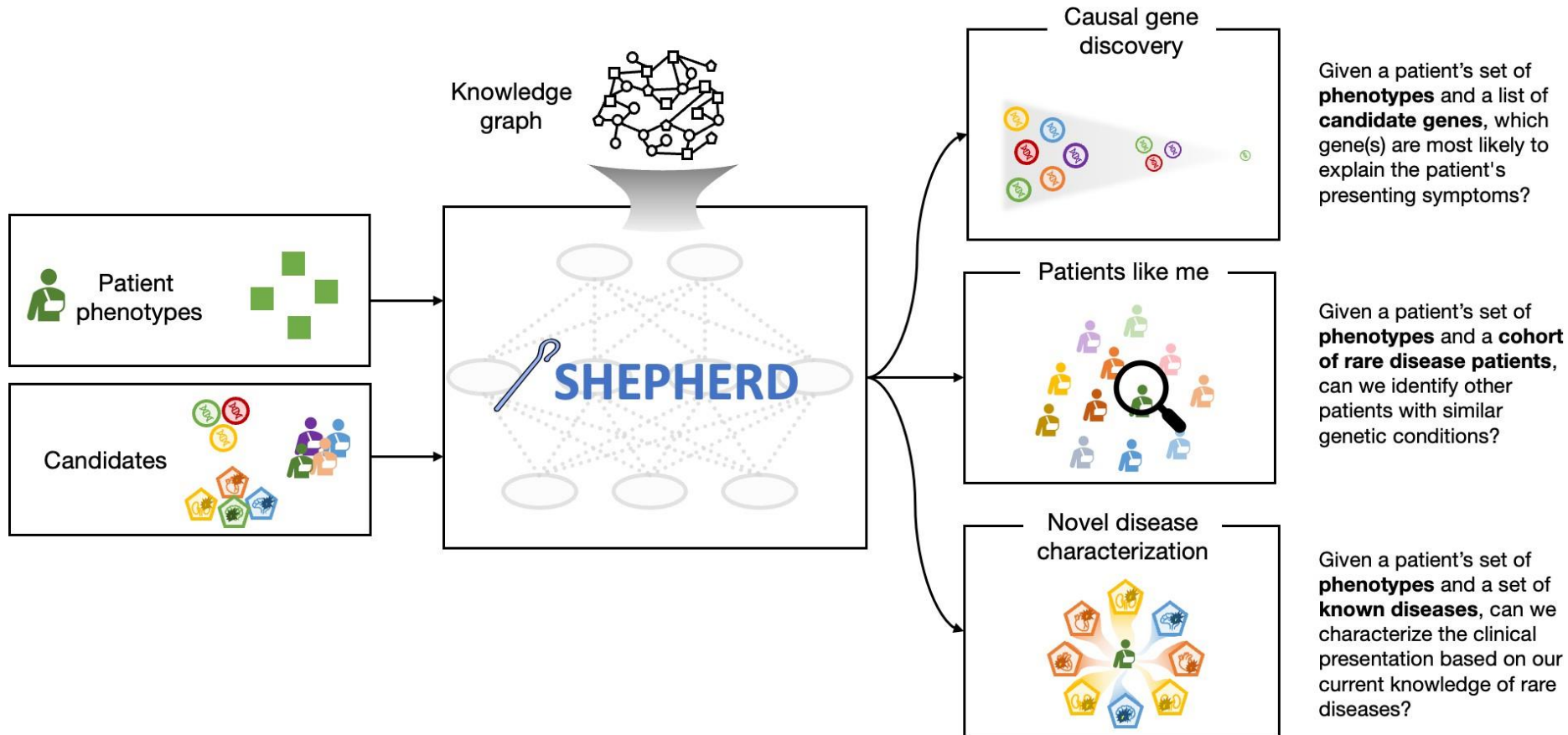% phenotypic overlap in patients with the same diseases

67%  +/-  43%

**Novel / atypical conditions**

% patient phenotypes with known association to causal gene

28% +/-  21%

# SHEPHERD: KG-based AI for rare disease diagnosis



**Causal gene discovery**

Given a patient's set of **phenotypes** and a list of **candidate genes**, which gene(s) are most likely to explain the patient's presenting symptoms?

**Patients like me**

Given a patient's set of **phenotypes** and a **cohort of rare disease patients**, can we identify other patients with similar genetic conditions?

**Novel disease characterization**

Given a patient's set of **phenotypes** and a set of **known diseases**, can we characterize the clinical presentation based on our current knowledge of rare diseases?
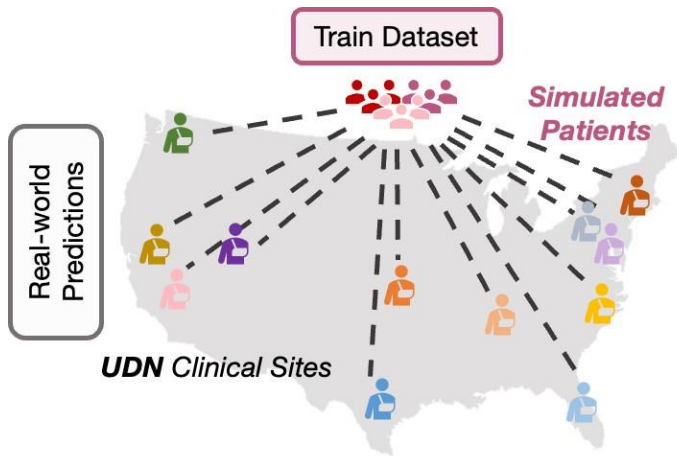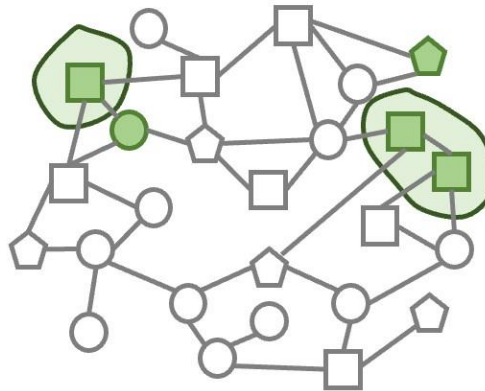
# AI for hard-to-diagnose diseases

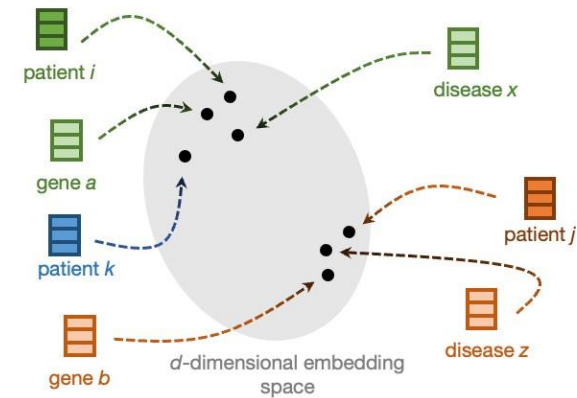# Key features of SHEPHERD



**Train on simulated patients, evaluate on UDN patients**

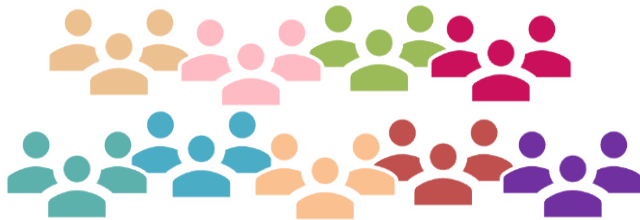**Model patients as subgraphs in knowledge graph**

**Perform label-efficient model training**

# Training data: Simulated patients

42,680 simulated patients across 2,134 diseases in Orphanet
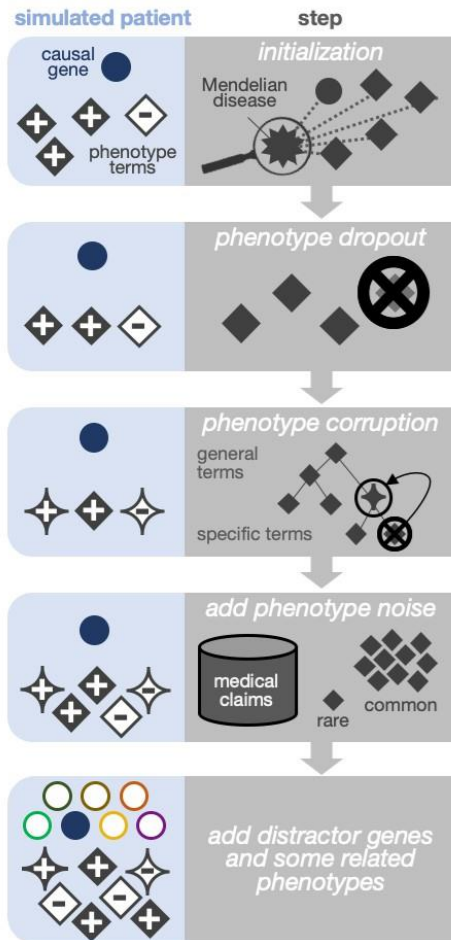
Train set
(N = 36,224)



Disease-split training and validation to select for generalizable models
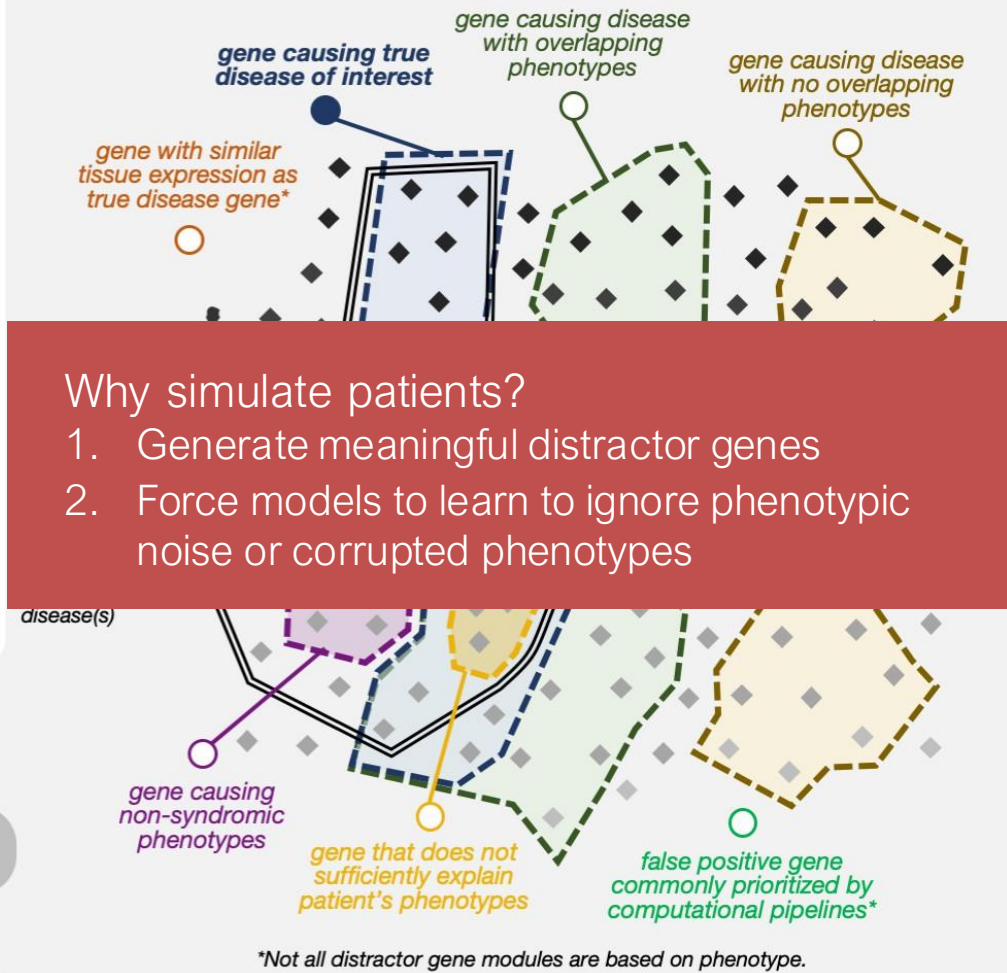
Validation set
(N = 6,400)

# Simulation process



**a Patient Simulation Process**

simulated patient | step

initialization — Mendelian disease

phenotype dropout

phenotype corruption — general terms, specific terms

add phenotype noise — medical claims, rare, common

add distractor genes and some related phenotypes

**b Distractor Gene Modules**

gene causing true disease of interest

gene causing disease with overlapping phenotypes

gene causing disease with no overlapping phenotypes

gene with similar tissue expression as true disease gene*

disease(s)

gene causing non-syndromic phenotypes

gene that does not sufficiently explain patient's phenotypes

false positive gene commonly prioritized by computational pipelines*

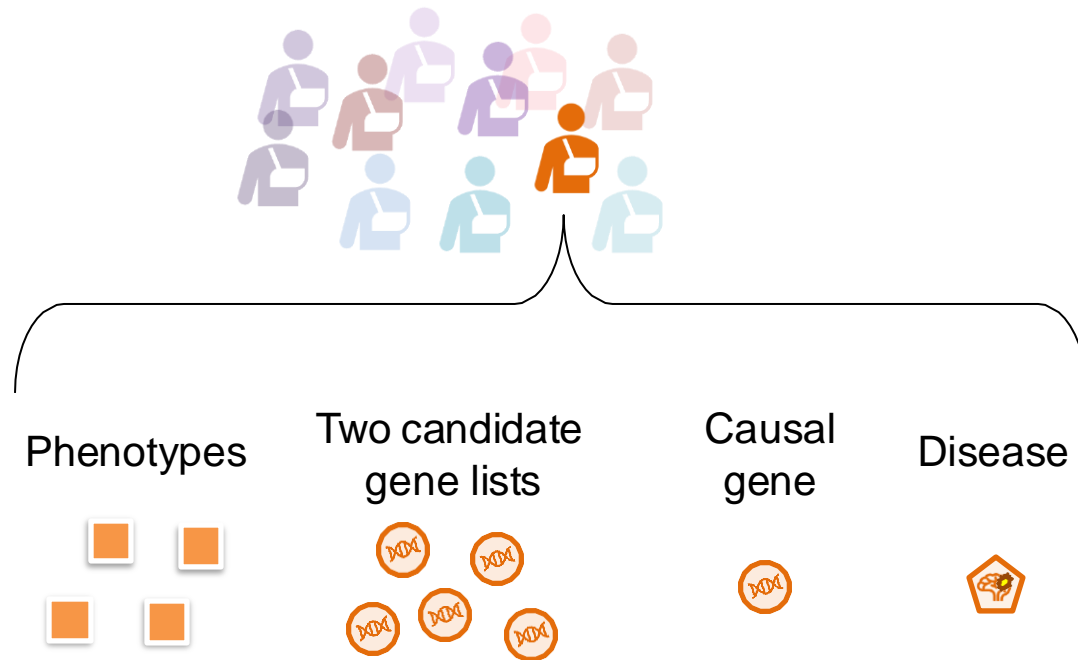*Not all distractor gene modules are based on phenotype.

Why simulate patients?
1. Generate meaningful distractor genes
2. Force models to learn to ignore phenotypic noise or corrupted phenotypes

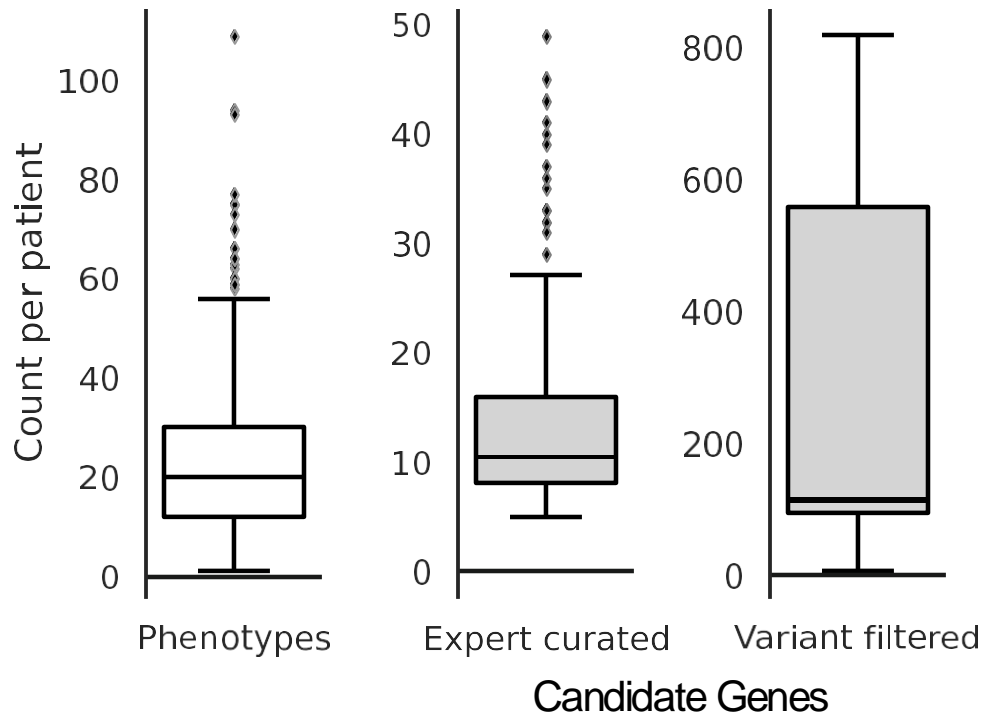# Undiagnosed Disease Network (UDN) cohort

465 patients who have received a molecular diagnosis



Phenotypes

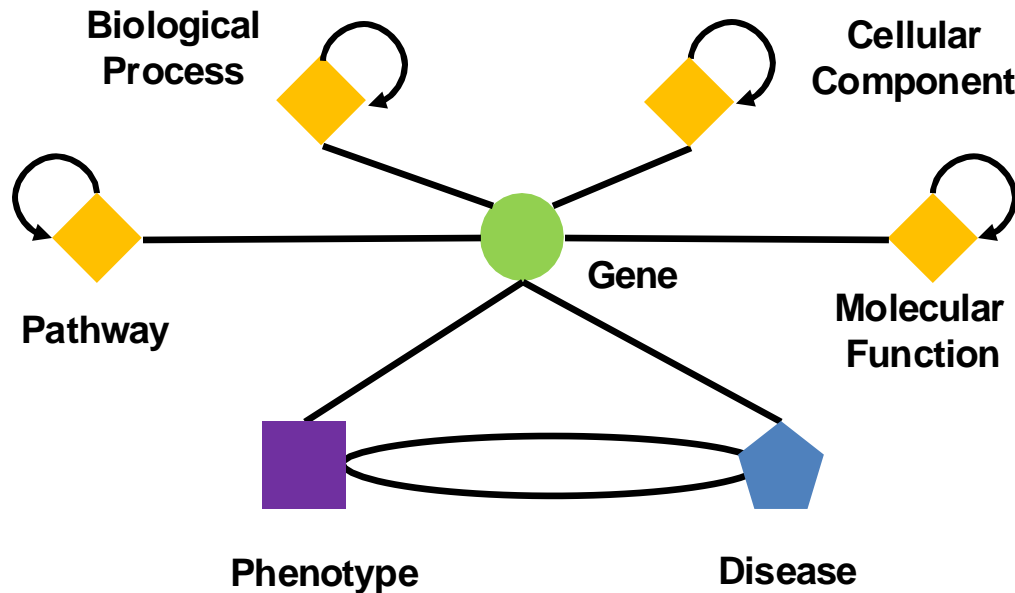Two candidate gene lists

Causal gene

Disease

# Undiagnosed Disease Network (UDN) cohort

465 patients who have received a molecular diagnosis

Number of phenotypes and
candidate genes per patient

# Rare disease knowledge graph (KG)



**Biological Process**

**Cellular Component**

**Pathway**

**Gene**

**Molecular Function**

**Phenotype**

**Disease**

Protein-Protein Interaction — STRING, TRANSFAC 2.0

Pathway Membership — reactome

Functional Similarity — GENE ONTOLOGY Unifying Biology

Phenotypic Similarity — human phenotype ontology, orphanet, OMIM Human Genetics Knowledge for the World, DisGeNET

| KG | # Types | Count |
|---|---|---|
| Nodes | 7 | 100,272 |
| Edges | 15 | 2,092,690 |

KG Modified from zitniklab.hms.harvard.edu/projects/PrimeKG/

# Knowledge graph learning



- **Step 1:** Incorporate knowledge of known phenotype, gene, and disease relationships via GNN
  - Knowledge-guided learning is achieved by self-supervised pre-training on our precision-medicine knowledge graph

- **Step 2:** Pre-trained GNN from Step 1 is fine-tuned using synthetic patients
  - Training exclusively on synthetic rare disease patients without the use of any real-world labeled cases
  - Synthetic patients used for training are created using an adaptive simulation approach
  - Realistic rare disease patients with varying numbers of phenotypes and candidate genes

# SHEPHERD's model



**Embed Biomedical Knowledge**

Sample nodes in external knowledge graph

Embed biomedical knowledge

Multi-layer Graph Attention Network

Edge exists
Edge does not exist
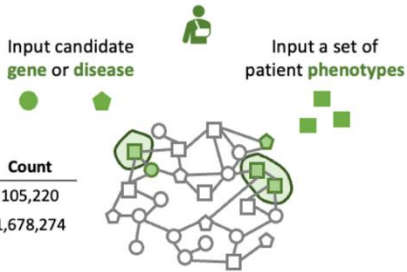
Self-supervised learning via link prediction on the knowledge graph

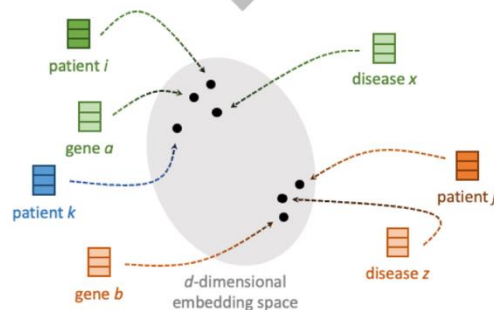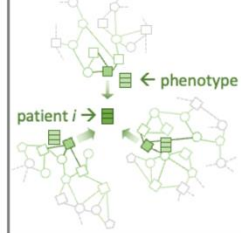**Embed Rare Disease Patient Information**

Input a set of patient **phenotypes**

Input candidate **genes** or **diseases** or **patients**

Embed candidate gene or disease

Embed patient phenotypes

Multi-layer Graph Attention Network

Align embedding space

patient i

gene a

patient k

gene b

disease x

patient j

disease z

$d$-dimensional embedding space

Embed **patient** <u>closer</u> to the **correct gene**, **disease**, or **patients with the same gene/disease**, and <u>farther</u> from the **incorrect gene**, **disease**, or **patients with a different gene/disease**.
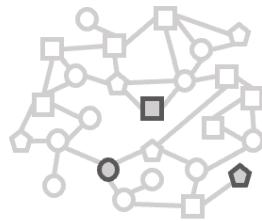
# Experimental setup

SHEPHERD's model training:

- 42K synthetic patients

SHEPHERD's model evaluation

- **UDN patient cohort:** 465 rare disease patients with labeled diagnoses, spanning 299 diseases
  - 79% of genes and 83% of diseases are represented in only a single patient
- **MyGene2 patient cohort:** 146 rare disease patients, spanning 55 diseases



| Patient dataset | Train cohort | Validation cohort | Test cohort |
|---|---|---|---|
| Simulated | N = 36,224 | N = 6,400 | --- |
| UDN | --- | --- | N = 465 |
| MyGene2 | --- | --- | N = 146 |

# Diagnostic tasks

- Three diagnostic tasks:

    - **Causal gene discovery:** Given a patient's set of phenotypes and a list of genes in which the patient has mutations, **prioritize genes** harboring mutations that cause the disease (phenotypes)

    - **Patients-like-me:** Given a patient, **find other patients** with similar genetic and phenotypic features suitable for clinical follow-up

    - **Characterization of novel diseases:** Given a patient's phenotypes, **provide an interpretable NLP name** for the patient's disease based on its similarity to each disease in the KG

# Diagnostic tasks



Causal Gene Discovery

$$similarity(\ \blacksquare \blacksquare \ ,\ \oplus \ ) \propto d(\ z_P\ ,\ z_g\ )$$

Patients-like-me

$$similarity(\ \blacksquare \blacksquare \ ,\ \blacksquare \blacksquare \ ) \propto d(\ z_{P_i}\ ,\ z_{P_j}\ )$$

Novel disease characterization

$$similarity(\ \blacksquare \blacksquare \ ,\ \oplus \ ) \propto d(\ z_P\ ,\ z_d\ )$$

Legend

- Patient phenotypes $P$
- Gene $g$
- Disease $d$

# Causal gene discovery: Results



Given a patient's set of **phenotypes** and a list of **candidate genes**, which gene(s) are most likely to explain the patient's presenting symptoms?

# Causal gene discovery: Results

Causal gene discovery

SHEPHERD

| | 0.40 | 0.69 | 0.85 |

Error bars denote standard deviation over models trained with 5 random seeds

Top k
- k = 1
- k = 3
- k = 5

# of causal genes retrieved in top k ranked genes on average

Average Recall at k

0.0    0.2    0.4    0.6    0.8    1.0

* LR = logistic regression

† Jagadeesh et al. Phrank measures phenotype sets similarity to greatly improve Mendelian diagnostic disease prioritization. Genetics in Medicine.

‡ Peng et al. CADA: phenotype-driven gene prioritization based on a case-enriched knowledge graph. NAR Genom Bioinform.

# Causal gene discovery: Results



Causal gene discovery

** p-value < 0.005
**** p-value < 0.00005

**SHEPHERD**: 0.40 | 0.69 | 0.85
†Information Theoretic: 0.35 | 0.65 | 0.81
Network Science: 0.34 | 0.63 | 0.77
‡Shallow Embedding: 0.16 | 0.24 | 0.51
*LR (PCA): 0.12 | 0.35 | 0.54
*LR (Embed): 0.11 | 0.35 | 0.55
Random: 0.09 | 0.27 | 0.50

Top k
k = 1
k = 3
k = 5

0.0    0.2    0.4    0.6    0.8    1.0
Average Recall at k

* LR = logistic regression
† Jagadeesh et al. Phrank measures phenotype sets similarity to greatly improve Mendelian diagnostic disease prioritization. Genetics in Medicine.
‡ Peng et al. CADA: phenotype-driven gene prioritization based on a case-enriched knowledge graph. NAR Genom Bioinform.

# Causal gene discovery: Results

SHEPHERD generalizes across…



Performance by Clinical Site

Performance by Evaluation Year

Performance by Primary Symptoms

# Atypical disease presentation

**Patient:** UDN-1
**Admitted**: 2016    **Diagnosed**: 2019
**Causal gene:** *POLR3A*
**Disease**: POLR3-Related Leukodystrophy
**Atypical Phenotypes**: – lack of tear production, premature adrenarche, laryngeal cleft, hearing loss, and high blood pressure

Only **28.3%** of the patient's 46 phenotypes are directly connected to *POLR3A*

**94%** of the 205 phenotypes directly connected to *POLR3A* are <u>not</u> associated with the patient

Subset of Rare Disease Knowledge Graph

# Atypical disease presentation



Expert Curated (N = 17)

Genes:
KAT6A, POLR3A, ORC4, WDFY4, ZFYVE26, GMPPA, NCOR2, APC, NDUFAF5, ANO3, INSL3, DST, TYMP, TOPORS, SLK, DYNAP, PIWIL3

Variant Filtered (N = 86)

UBE3A, POLR3A, KMT2E, TNIK, ORC4, CTU2, TGIF1, TBP, MED16, DVL3

FHDC1, MUC3A, OVGP1, DPY19L2, KCNJ18, KLF18, UMODL1, NPY4R, GAGE12J, ANKRD36C

Score

| Attention | Phenotype (N = 46) |
|---|---|
| 0.037 | Short stature |
| 0.034 | Failure to thrive |
| 0.033 | Central hypotonia |
| 0.032 | Microcephaly |
| 0.032 | Prominent eyelashes |
| 0.032 | Respiratory insufficiency |
| 0.031 | Gastrostomy tube feeding in infancy |
| 0.031 | Chronic lung disease |
| 0.028 | Ventriculomegaly |
| 0.027 | Growth delay |
| … | … |
| ⭐ 0.014 | Alacrima |
| 0.014 | Premature loss of primary teeth |
| ⭐ 0.013 | Moderate sensorineural hearing impairment |
| 0.013 | Pancreatitis |
| 0.012 | Abnormal sternum morphology |
| 0.011 | T2 hypointense basal ganglia |
| 0.011 | Febrile seizure (within the age range of 3 months to 6 years) |
| 0.009 | Chronic pancreatitis |
| ⭐ 0.006 | Laryngeal cleft |
| 0.0003 | T2 hypointense brainstem |

Top 10 phenotypes

Bottom 10 phenotypes

# Results: Patients-like-me

UMAP plot of SHEPHERD's embedding space of all simulated (circle), UDN (up-facing triangle), and MyGene2 (down-facing triangle) patients colored by their Orphanet disease category



**a**

**Patient: UDN-P3**            *Patient Card*
**Causal gene:** *RPS6KA3*
**Disease:** Coffin-Lowry syndrome

| Patient Rank | Gene | Disease |
|---|---|---|
| 1 | GRIA3 | X-linked intellectual disability due to GRIA3 anomalies |
| 2 | **RPS6KA3** | **Coffin-Lowry syndrome** |
| 3 | THOC2 | X-linked intellectual disability-short stature-overweight syndrome |
| 4 | AP1S2 | Fried syndrome |
| 5 | SMS | Syndromic X-linked intellectual disability Snyder type |

**Patient: UDN-P5**            *Patient Card*
**Causal gene:** *NLRP12, RAPGEFL1*
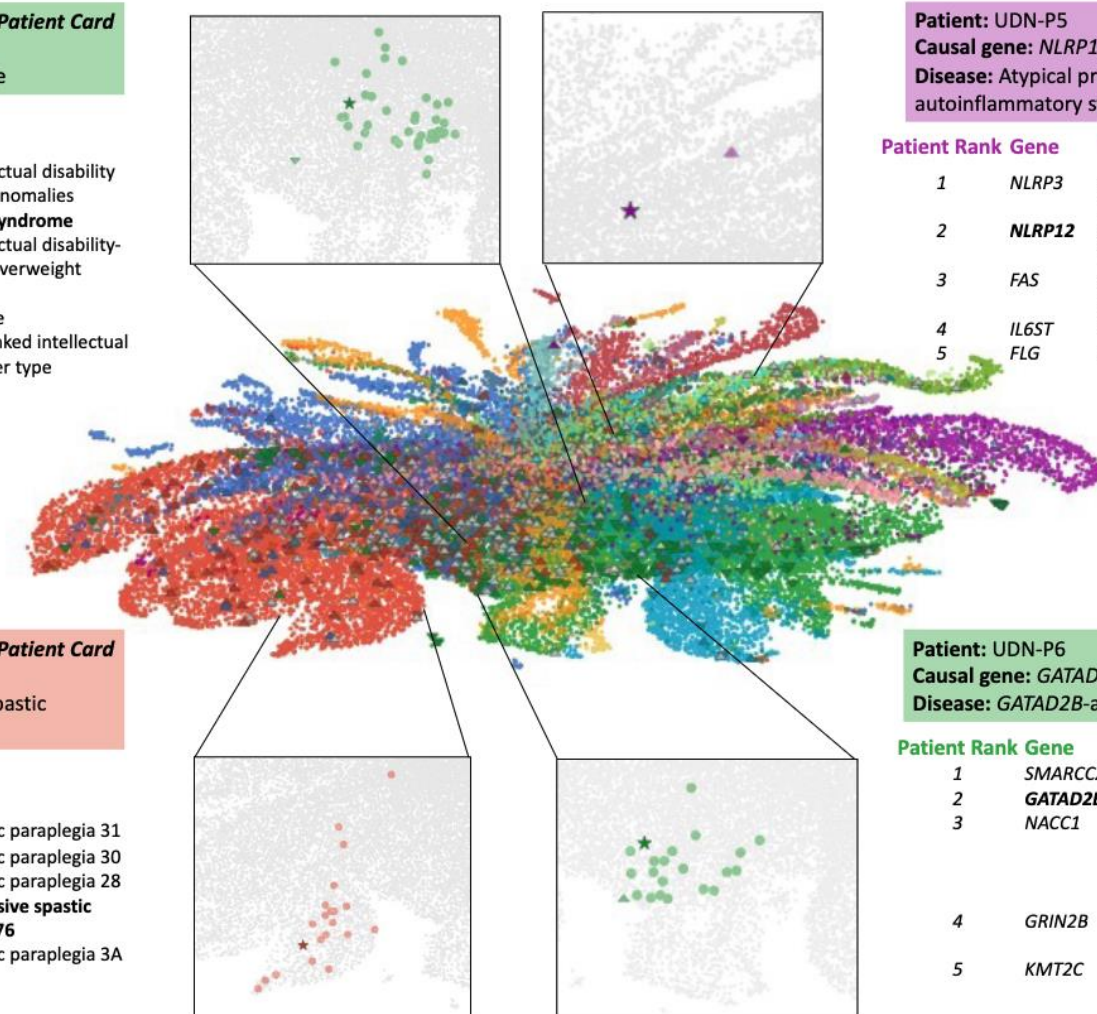**Disease:** Atypical presentation of familial cold autoinflammatory syndrome

| Patient Rank | Gene | Disease |
|---|---|---|
| 1 | NLRP3 | Familial cold-induced autoinflammatory syndrome 1 |
| 2 | **NLRP12** | **Familial cold-induced autoinflammatory syndrome 2** |
| 3 | FAS | autoimmune lymphoproliferative syndrome type 1 |
| 4 | IL6ST | GP130-deficient hyper-IgE syndrome |
| 5 | FLG | atopic dermatitis 2 |

**Patient: UDN-P4**            *Patient Card*
**Causal gene:** *CAPN1*
**Disease:** autosomal recessive spastic paraplegia type 76

| Patient Rank | Gene | Disease |
|---|---|---|
| 1 | REEP1 | hereditary spastic paraplegia 31 |
| 2 | KIF1A | hereditary spastic paraplegia 30 |
| 3 | DDHD1 | hereditary spastic paraplegia 28 |
| 4 | **CAPN1** | **autosomal recessive spastic paraplegia type 76** |
| 5 | MTPAP | hereditary spastic paraplegia 3A |

**Patient: UDN-P6**            *Patient Card*
**Causal gene:** *GATAD2B*
**Disease:** *GATAD2B*-associated syndrome

| Patient Rank | Gene | Disease |
|---|---|---|
| 1 | SMARCC2 | Coffin-Siris syndrome 8 |
| 2 | **GATAD2B** | **GATAD2B-associated syndrome** |
| 3 | NACC1 | neurodevelopmental disorder with epilepsy, cataracts, feeding difficulties, and delayed brain myelination syndrome |
| 4 | GRIN2B | intellectual disability, autosomal dominant 6 |
| 5 | KMT2C | Kleefstra syndrome |

# Results: New disease naming

**a** Rank Disease

1. AR limb-girdle muscular dystrophy type 2B
2. GNE myopathy
3. MYH7-related late-onset scapuloperoneal muscular dystrophy
4. Emery-Dreifuss muscular dystrophy 2, AD
5. AR limb-girdle muscular dystrophy type 2G
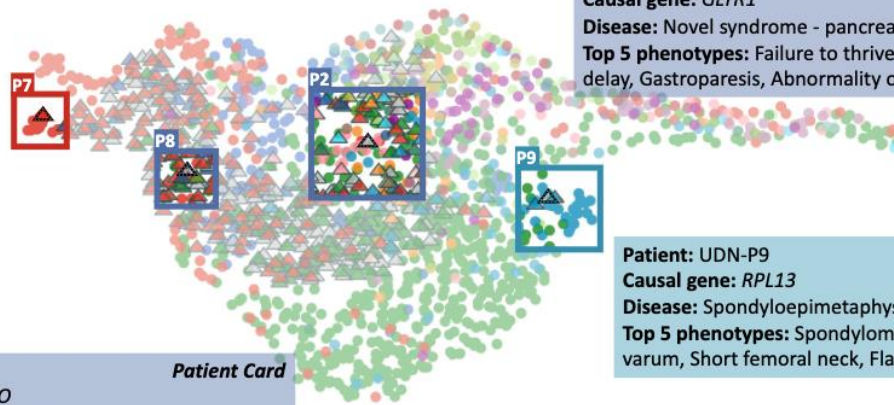
Rank Disease

1. Methylmalonic aciduria & homocystinuria type cblF
2. Neonatal hemochromatosis
3. Homozygous 11P15-p14 deletion syndrome
4. ALG8-CDG
5. Congenital anemia

**Patient: UDN-P7** *Patient Card*
**Causal gene:** *SGCA*
**Disease:** AR limb-girdle muscular atrophy type 2D
**Top 5 phenotypes:** Toe walking, Calf muscle pseudohypertrophy, Elevated serum creatine kinase, Proximal muscle weakness, Generalized muscle weakness

**Patient: UDN-P2** *Patient Card*
**Causal gene:** *GLYR1*
**Disease:** Novel syndrome - pancreatic insufficiency & malabsorption
**Top 5 phenotypes:** Failure to thrive in infancy, Global developmental delay, Gastroparesis, Abnormality of vision, Duodenal atresia
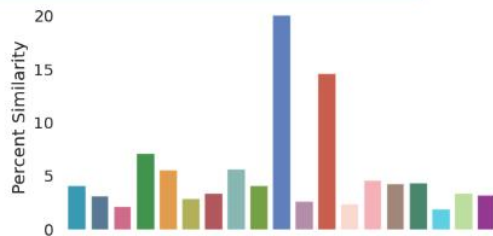
**Patient: UDN-P9** *Patient Card*
**Causal gene:** *RPL13*
**Disease:** Spondyloepimetaphyseal dysplasia, Isidor-Toutain type
**Top 5 phenotypes:** Spondylometaphyseal dysplasia, Genu varum, Short femoral neck, Flat glenoid fossa, Platyspondyly

Rank Disease

1. Combined oxidative phosphorylation deficiency 39
2. Hypomyelinating leukodystropy-20
3. Pyruvate dehydrogenase E3-binding protein deficiency
4. Intellectual disability-epilepsy-extrapyramidal syndrome
5. Combined oxidative phosphorylation defect type 27

**Patient: UDN-P8** *Patient Card*
**Causal gene:** *ATP5PO*
**Disease:** *ATP5PO*-related Leigh syndrome
**Top 5 phenotypes:** Profound global developmental delay, cerebral hypomyelination, limb hypertonia, hypoplasia of the corpus callosum, infantile spasms

Rank Disease

1. Multiple epiphyseal dysplasia type 1
2. Progressive pseudorheumatoid arthropathy of childhood
3. Multiple epiphyseal dysplasia type 5
4. Metaphyseal chondrodysplasia, Spahr type
5. Multiple epiphyseal dysplasia

# SHEPHERD: KG-based AI for rare disease diagnosis



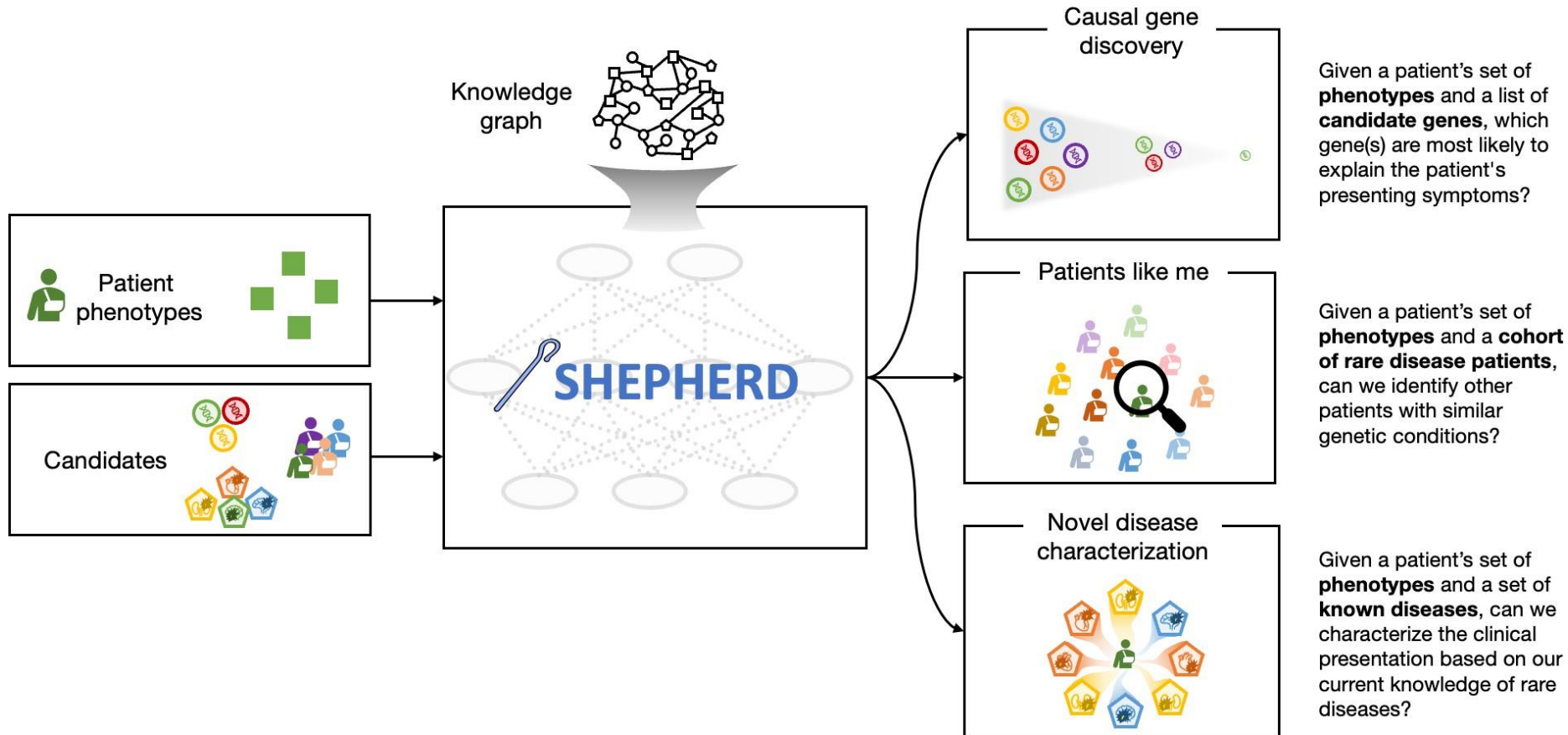**Causal gene discovery**

Given a patient's set of **phenotypes** and a list of **candidate genes**, which gene(s) are most likely to explain the patient's presenting symptoms?

**Patients like me**

Given a patient's set of **phenotypes** and a **cohort of rare disease patients**, can we identify other patients with similar genetic conditions?

**Novel disease characterization**

Given a patient's set of **phenotypes** and a set of **known diseases**, can we characterize the clinical presentation based on our current knowledge of rare diseases?

# Take-away messages

- SHEPHERD overcomes limitations of standard machine learning:
  - Model inputs as **KG subgraphs** (i.e., clinic-genetic subgraphs of patients)
  - Use **self-supervised pre-training on biomedical knowledge**
  - Train the model on a large cohort of **synthetic patients**

- SHEPHERD generalizes to novel phenotypes, genes, and diseases:
  - Performs well on patients whose **subgraphs are of varying size**
  - Performs well on **diagnosing patients with novel diseases**

- Implications:
  - Implications for **generalist models applicable across diagnostic process**
  - New opportunities to **shorten the diagnostic odyssey for rare disease**
  - Implications for using **deep learning on medical datasets with very few labels**

**First deep learning approach for individualized diagnosis
of rare genetic diseases**

**Graph learning approach is not only helpful but necessary**

# Quick check

https://forms.gle/AfRT7pdXGa7MoJxJA



## AIM 2: Artificial Intelligence in Medicine II

*Artificial Intelligence in Medicine II, Spring 2025*

Lecture 9: Knowledge graph learning, Building multimodal knowledge graphs, Structure-inducing pre-training, Knowledge-based foundation models

Course website and slides: **https://zitniklab.hms.harvard.edu/AIM2**

* Indicates required question

First and last name *

Your answer

Harvard email address *

Your answer

SHEPHERD model was evaluated on three diagnostic tasks: causal gene discovery, patient-like-me retrieval, and characterization of new diseases. Suggest another use case (application) for SHEPHERD for rare diseases. *

Your answer

List two reasons why the SHEPHERD model was trained on a dataset of simulated patients. *

Your answer

# Towards foundation models for knowledge graphs
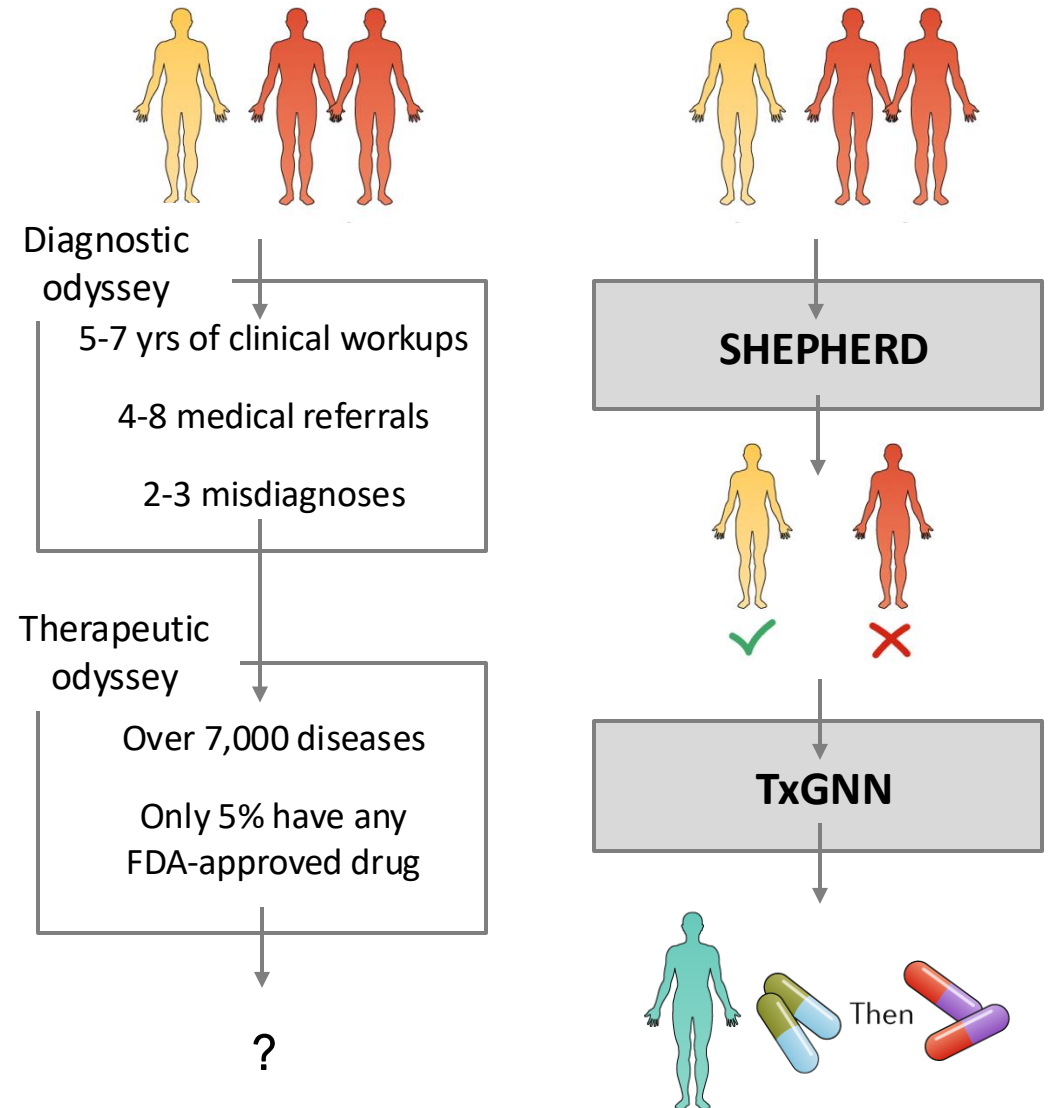
# Future with AI: From mysteries to therapies



Knowledge graph models for diagnosing rare disease patients

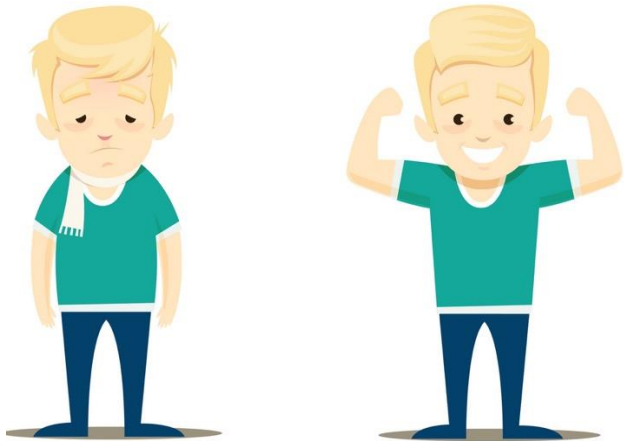SHEPHERD: Deep learning for diagnosing patients with rare genetic diseases, medRxiv 2025

Knowledge graph models for universal drug repurposing

TxGNN: A foundation model for clinician-centered drug repurposing, *Nature Medicine* 2024

Diagnostic odyssey

5-7 yrs of clinical workups

4-8 medical referrals

2-3 misdiagnoses

Therapeutic odyssey

Over 7,000 diseases

Only 5% have any FDA-approved drug

?

SHEPHERD

TxGNN

Then

# Precision medicine (treatments)

**Measure phenotype and mechanisms**

**Design therapeutic agents or select optimal perturbations**

**+**

**Provide each patient with the right drug, at the right dose, at the right time**

| *Clinical phenotypes and diseases* | |
|---|---|
| 17,000 | Diseases |
| 7,000 | Rare diseases |
| 5-7% | Rare diseases with treatments |
| No | Treatment options for many disease subtypes |

| *Medicines and drugs* | |
|---|---|
| 40-50 | New molecules per year |
| 30% | Drugs are issued at least one post-approval new indication |
| Many | Drugs have accrued over 10 drug indications over the years |

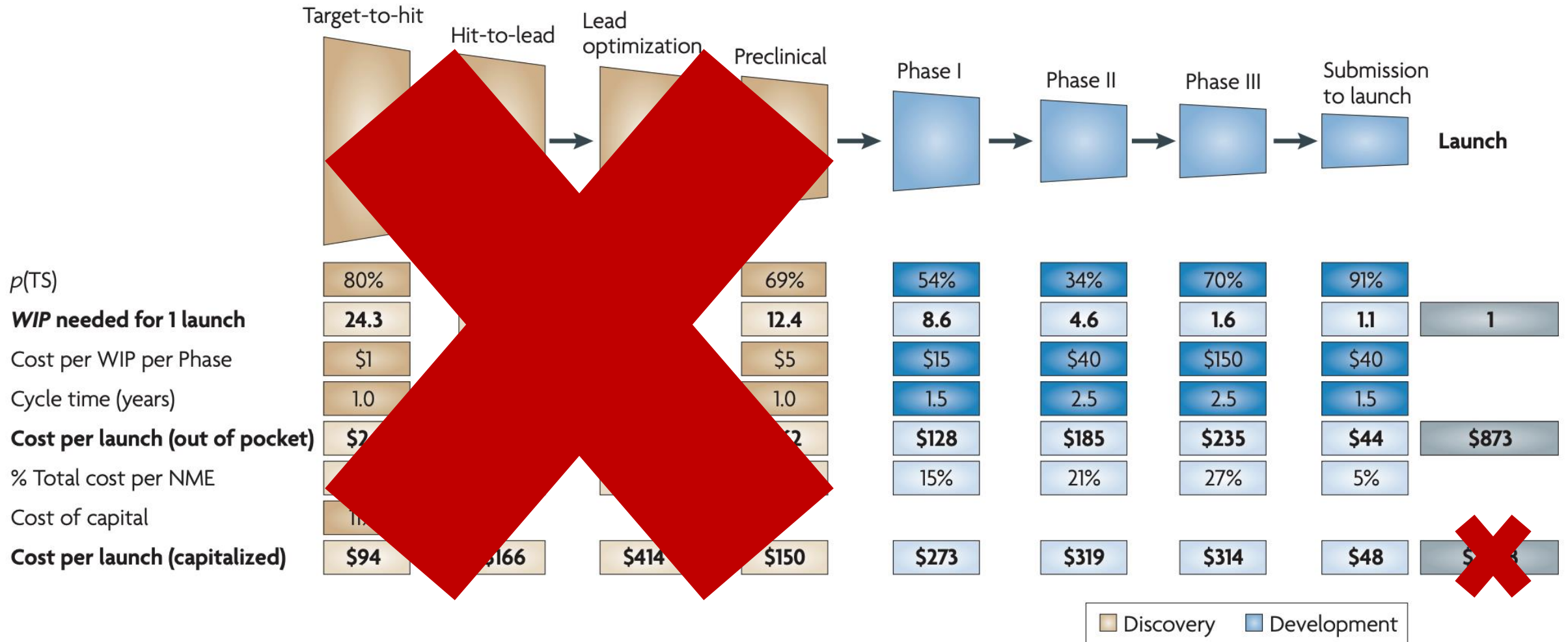**Matching drugs to clinical outcomes across thousands of diseases**

# Drug repurposing as an effective drug development strategy for many diseases

**1** No effective treatments for rare and even many complex diseases:

- Over 7,000 rare diseases affect 300-400 million people worldwide. Only 5% of rare diseases have FDA-approved drugs

- Even for diseases with approved treatments, new drugs can offer alternative options that cause fewer side effects and replace drugs that are ineffective for patient subpopulations

**2** Faster translation to the clinic and lower development costs

- 30% of drugs approved were issued at least one post-approval new indication. Many drugs have accrued over 10 indications over years

- Most repurposed drugs are the results of serendipity (luck is not a strategy!)

# Phases of drug discovery from initial stage (target-to-hit) to final stage (launch)



p(TS) – probability of successful transition from one stage to the next; NME – new molecular entity; WIP – work in process

# All-disease model for drug repurposing

Biomedical data span multiple scales and multiple data modalities

Once trained, models are adapted to an array of tasks, with no or minimal training



Transcriptomics

Physical contacts

Molecular pathways and patient subtypes

Treatment information

**TxGNN:
All-disease drug
repurposing model**

What patient populations will respond to treatment?

What candidate therapeutics will have an acceptable safety profile for patients with metastatic melanoma?

What small-molecule compounds will inhibit a kinase?

# All-disease model for drug repurposing

Multimodal knowledge graph
of 17,080 disease phenotypes

Process therapeutic tasks and predict candidate
indications and contraindications



txgnn.org

TxGNN
AI Model

TxGNN
Explainer

Mechanistic path from drug to disease

"indication"

"contraindication"

"contraindication"

Semi-automatic KG rebuild when new datasets
become available
Building a knowledge graph to enable precision
medicine, Scientific Data  2023

Structure-inducing pre-training, *Nature Machine Intelligence* 2023; Multimodal learning with graphs, *Nature Machine Intelligence* 2023;
Graph Representation Learning in Biomedicine and Healthcare, *Nature Biomedical Engineering* 2022; Multimodal Learning with Graphs, *Nature Machine Intelligence* 2023; A foundation
model for clinician-centered drug repurposing, *Nature Medicine* 2024

# All-disease model for drug repurposing



TxGNN Predictor

Likelihood of Indication

Likelihood of Contraindication

TxGNN Explainer

Multi-hop interpretable path from a drug to a disease

# Building knowledge graphs: Medical data are multimodal and scattered across databases



**VAST
UNORGANIZED
KNOWLEDGE**

**CURATED
KNOWLEDGE
GRAPH**

Ayush Noori

9

# Building knowledge graphs: Medical data are multimodal and scattered across databases

**ChEMBL**

includes 1.6M assays covering 2.4M compounds

**Bgee**

includes 31,467 bulk and single-cell RNA-seq libraries

**STRING**

includes 20B interactions between 59.3M proteins

**GENEONTOLOGY** Unifying Biology

includes 6M gene annotations derived from 150K publications

**reactome**

includes 2,711 pathways manually curated by PhDs

**DRUGBANK**

includes 17K FDA-approved and experimental drugs

**NIH National Library of Medicine** National Center for Biotechnology Information

includes annotations for 192K human genetic elements

**SIDER**

includes 139K adverse reactions for marketed drugs

**human phenotype ontology**

includes 13K phenotypes and 156K disease annotations

Ayush Noori

# Building knowledge graphs: Medical data are multimodal and scattered across databases

**ChEMBL**

includes 1.6M assays covering 2.4M compounds

## ChEMBL evidence integration pipeline

*repeat for all 36 databases*

| | |
|---|---|
| **1** | Include only drugs approved for marketing by the FDA or clinical candidates. |
| **2** | Machine learning analysis to evaluate clinical trials that ended earlier than scheduled. |
| **3** | Score all drug-disease edges by clinical precedence (*i.e.*, Phase 0, I, II, III, IV). |
| **4** | Down-weight scores for trials that stop early due to negative outcomes or safety concerns. |
| **5** | Construct typed edges in knowledge graph based on strength of ChEMBL evidence. |

Ayush Noori

# Knowledge graph based TxGNN model enables transfer learning across 17,080 disease phenotypes



Disease pooling is a module that identifies diseases similar to a query disease and transfers information from related diseases to the query disease

# TxGNN identifies candidate drugs for diseases with no treatment options

Once trained, TxGNN can perform zero-shot prediction on new diseases without additional parameters or fine-tuning on labeled data



Scenario A: Current state-of-the-art
- Disease with existing treatments
- Easier to predict

Scenario B: Zero-shot prediction
- Diseases with no existing treatments
- Much harder to predict

7,000+ rare diseases affect 300-400M globally; only 5% have FDA-approved drugs. New drugs can offer better, side-effect reduced options for specific patients

# Benchmarking TxGNN on challenging dataset splits across disease areas



Held-out folds contain diseases …

… with zero approved drugs

… from distinct disease areas

… with limited molecular data

Held-out folds contain **cancer diseases**

Diseases in this area include:
- Leydig cell tumor
- Neurofibroma
- Acute myeloid leukemia

Held-out folds contain **anemia-related diseases**
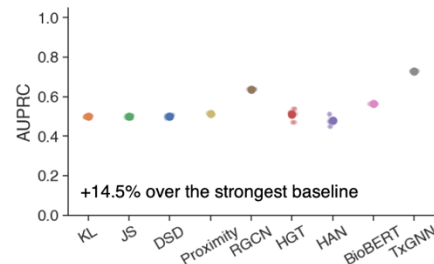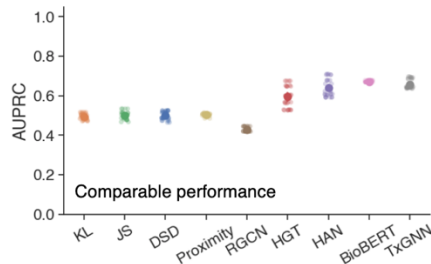
Diseases in this area include:
- Thalassemia
- Aplastic anemia
- Hemoglobin C disease

Held-out folds contain **cardiovascular diseases**
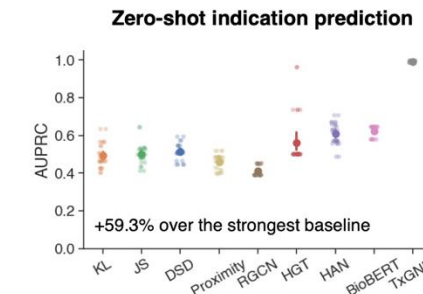
Diseases in this area include:
- Mitral valve stenosis
- Congestive heart failure
- Long QT syndrome
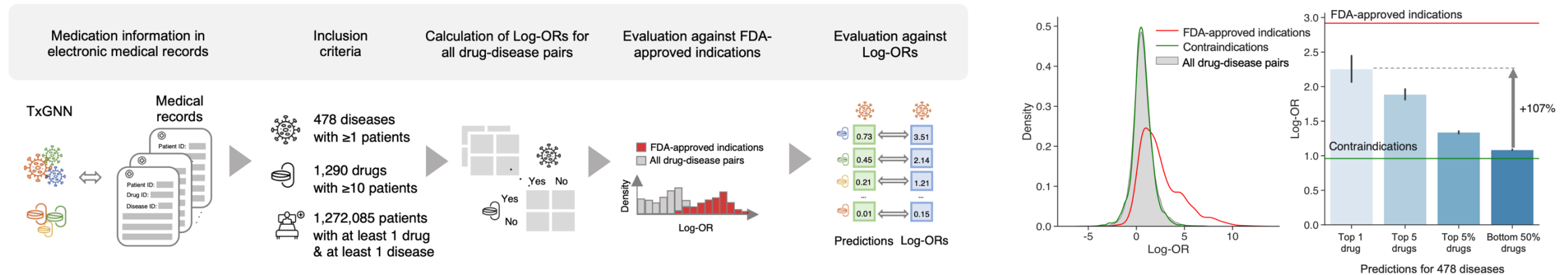
Held-out folds contain **adrenal gland diseases**

Diseases in this area include:
- Hyperaldosteronism
- Addison disease
- Ectopic cushing syndrome

+10.2% over the strongest baseline

+11.8% over the strongest baseline

+42.3% over the strongest baseline

+13.7% over the strongest baseline

Comparable performance

+14.5% over the strongest baseline

**Zero-shot indication prediction**

+59.3% over the strongest baseline

**Zero-shot contraindication prediction**

+17.8% over the strongest baseline

# Evaluating new drug repurposing predictions

- TxGNN's novel predictions are consistent with off-label prescription decisions made by clinicians in a large healthcare system
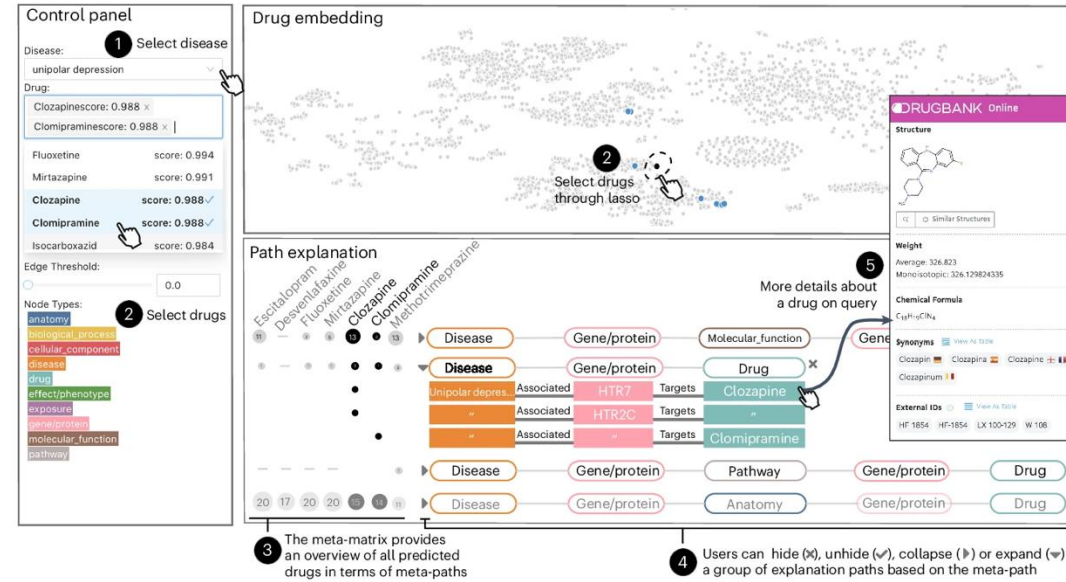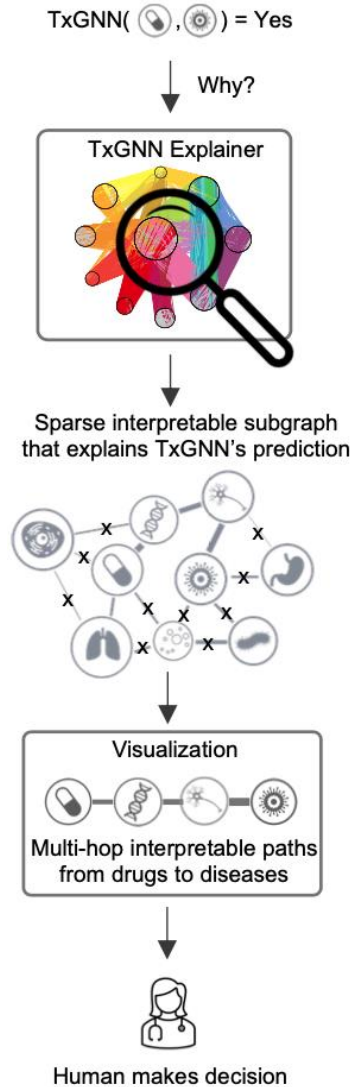


- TxGNN predicts therapeutic use for recent FDA approvals and informs laboratory testing

PNAS'21

| Drug name | Ingredient | Disease | Approval date | Company | FDA Number | Orphan | Prediction | Percentile |
|---|---|---|---|---|---|---|---|---|
| Welireg | Belzutifan | von Hippel-Lindau disease | 08/13/2021 | Merck | NDA215383 | Yes | 0.720 | 4.11% |
| Livtencity | Maribavir | Cytomegalovirus infection | 11/23/2021 | Takeda | NDA215596 | Yes | 0.033 | 66.37% |
| Tezspire | Tezepelumab-Ekko | Asthma | 12/17/2021 | Astrazeneca | BLA761224 | No | 0.233 | 32.41% |
| Leqvio | Inclisiran Sodium | Familial hypercholesterolemia | 12/22/2021 | Novartis | NDA214012 | No | 0.301 | 19.32% |
| Adbry | Tralokinumab | Atopic dermatitis | 12/27/2021 | Leo Pharma | BLA761180 | No | 0.040 | 50.37% |
| Vabysmo | Faricimab-Svoa | Macular degeneration | 01/28/2022 | Genentech | BLA761235 | No | 0.938 | 2.25% |
| Vonjo | Pacritinib Citrate | Myelofibrosis | 02/28/2022 | Cti Biopharma | NDA208712 | Yes | 0.011 | 63.14% |
| Ztalmy | Ganaxolone | CDKL5 disorder | 03/18/2022 | Marinus | NDA215904 | Yes | 0.335 | 18.73% |
| Mounjaro | Tirzepatide | Type 2 diabetes mellitus | 05/13/2022 | Eli Lilly | NDA215866 | No | 0.286 | 12.50% |
| Vtama | Tapinarof | Psoriasis | 05/23/2022 | Dermavant | NDA215272 | No | 0.261 | 32.70% |

A foundation model for clinician-centered drug repurposing, *Nature Medicine* 2024
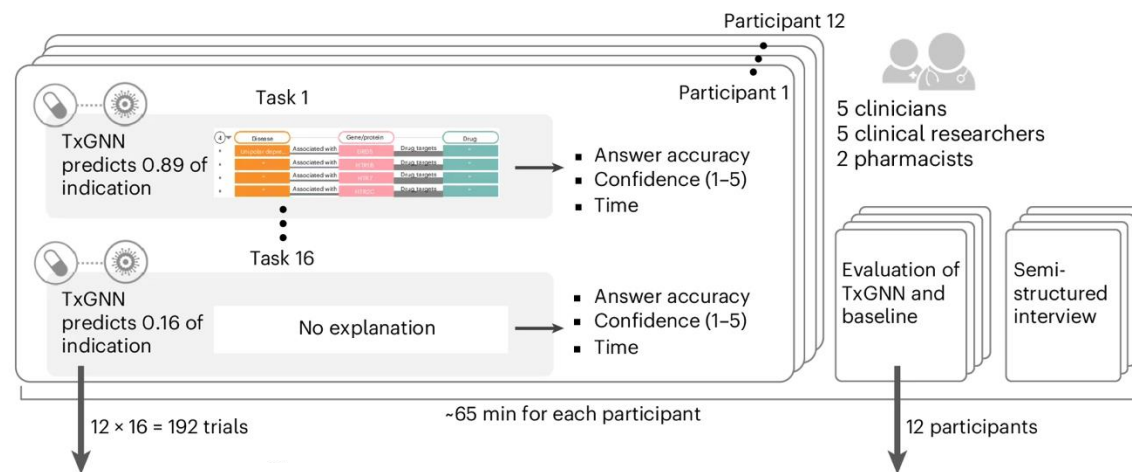
# Clinician-centered design: txgnn.org



Panels of clinicians, clinical researchers and pharmacists test usability of TxGNN:
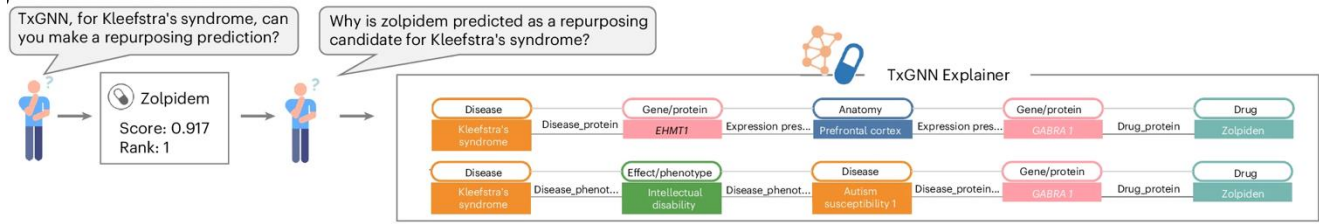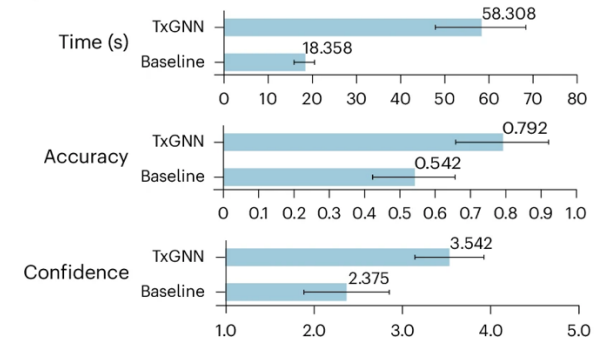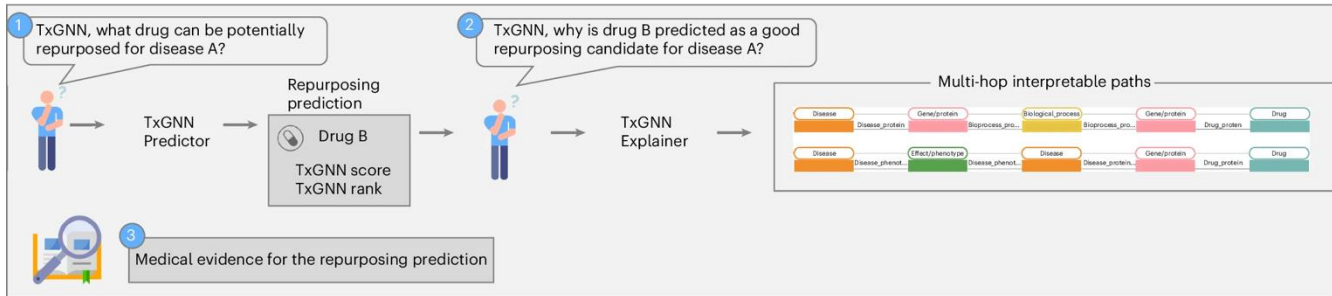- Scientific and medical consensus
- User confidence and trust
- User agreement
- Time used for exploring predictions

**Path-based explanations** perform significantly better than **node-based explanations** and **subgraph-based explanations** across three usability metrics: accuracy, confidence, time
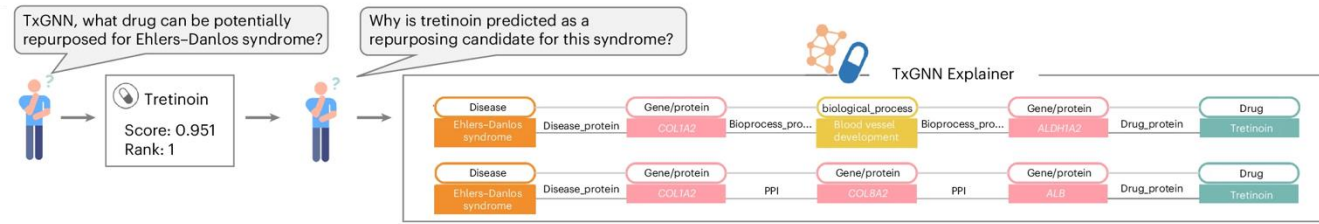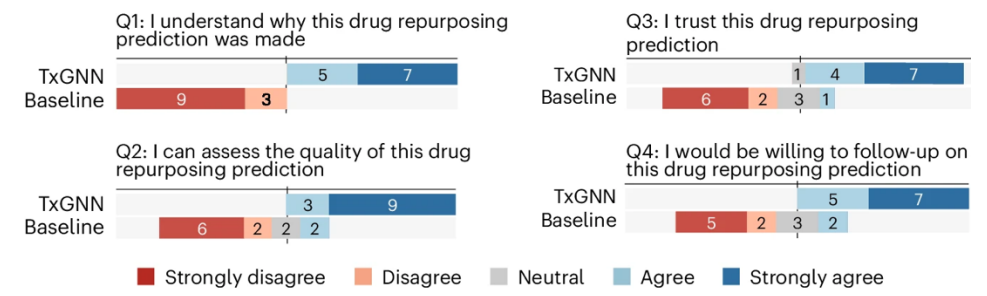
# Clinician-centered design: txgnn.org



- Better accuracy (+46%) and confidence (+49%) when explanations provided
- Support scientists in interacting with TxGNN and interpreting TxGNN predictions

# Open models, open datasets, and evaluations



The Harvard Gazette

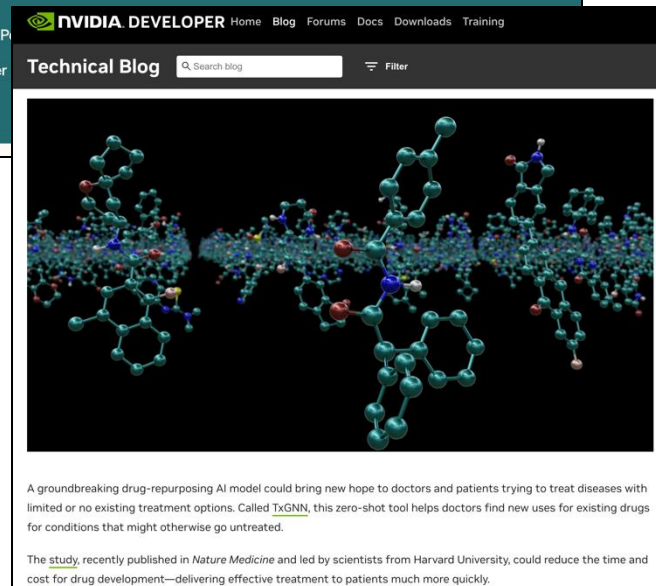Findings | Campus & Community | Health | Science & Tech | Nation & World | Arts & Culture | ≡ Menu | 🔍

**HEALTH**

## Using AI to repurpose existing drugs for treatment of rare diseases

Identifies possible therapies for thousands of diseases, including ones with no current treatments
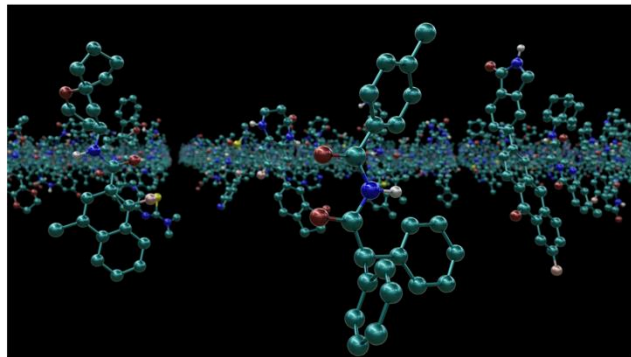
Ekaterina P
September

---

Forbes

Subscribe: Less than $1.50/wk    Sign In

## AI Tool Speeds Drug Repurposing: And It's Free

By Greg Licholai MD, Contributor. Greg Licholai writes and teaches about...    Follow Author

Sep 26, 2024, 06:00am EDT

Save Article    Comment 0

---

NVIDIA. DEVELOPER    Home    Blog    Forums    Docs    Downloads    Training

**Technical Blog**    Search blog    Filter

A groundbreaking drug-repurposing AI model could bring new hope to doctors and patients trying to treat diseases with limited or no existing treatment options. Called TxGNN, this zero-shot tool helps doctors find new uses for existing drugs for conditions that might otherwise go untreated.

The study, recently published in *Nature Medicine* and led by scientists from Harvard University, could reduce the time and cost for drug development—delivering effective treatment to patients much more quickly.

---

Kempner INSTITUTE    About Us    People    Research    Compute    Education    Careers & Oppor

**PRESS RELEASE**

## With TxGNN, Kempner Researchers Introduce an AI "Dr. House" to Find Treatments for Rare Diseases

By Yohan J. John, Ph.D. | September 30, 2024

SHARE ON 𝕏 in f ✉

Kempner scientists are using powerful AI technology to identify potential drug–disease pairings that could help advance treatment for rare diseases.

---

AI for drug repur

Innovative a
the myriad
expensive, w

- Real-world implementation
- Clinical collaborations for 20+ diseases, including neurology, cancer, and rare diseases

18