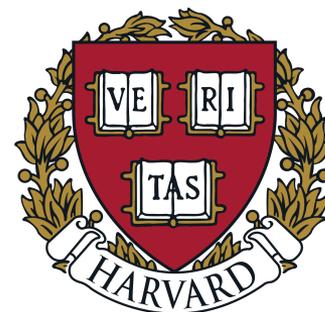# Machine Learning for Drug Development

## Marinka Zitnik

Department of Biomedical Informatics
Broad Institute of Harvard and MIT
Harvard Data Science Initiative

marinka@hms.harvard.edu
https://zitniklab.hms.harvard.edu

# Outline

✓ Overview and introduction

Part 1: Virtual drug screening and drug repurposing 👉

Part 2: Adverse drug effects, drug-drug interactions

Part 3: Clinical trial site identification, patient recruitment

Part 4: Molecule optimization, molecular graph generation, multimodal graph-to-graph translation

Part 5: Molecular property prediction and transformers

Demos, resources, wrap-up & future directions

# Method:
## Subgraph Neural Networks

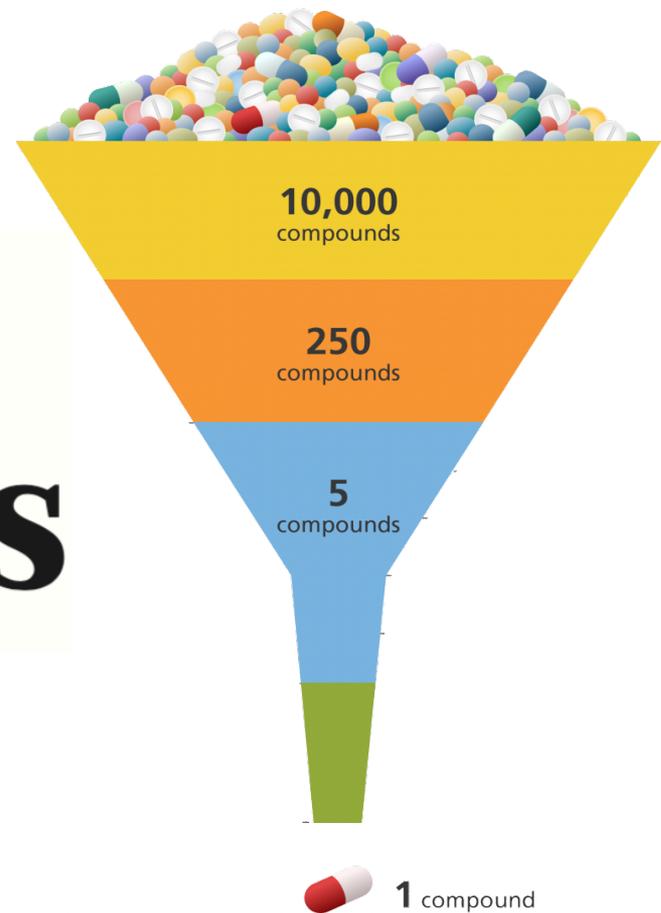Alsentzer, Finlayson, Li, and Zitnik, Subgraph Neural Networks, *NeurIPS* 2020

# Application:
## Finding Effective Drug Treatments

In submission

# New tricks for old drugs

*Faced with skyrocketing costs for developing new drugs, researchers are looking at ways to repurpose older ones — and even some that failed in initial trials.*

10,000 compounds

250 compounds

5 compounds

12–16 years, ~$1 billion to $2 billion

**1** compound

| Drug discovery | Preclinical testing | Phase I Phase II | Phase III | FDA approval |
|---|---|---|---|---|
| 3–6 years | 3 years | 3 years | 2 years | 1–2 years |

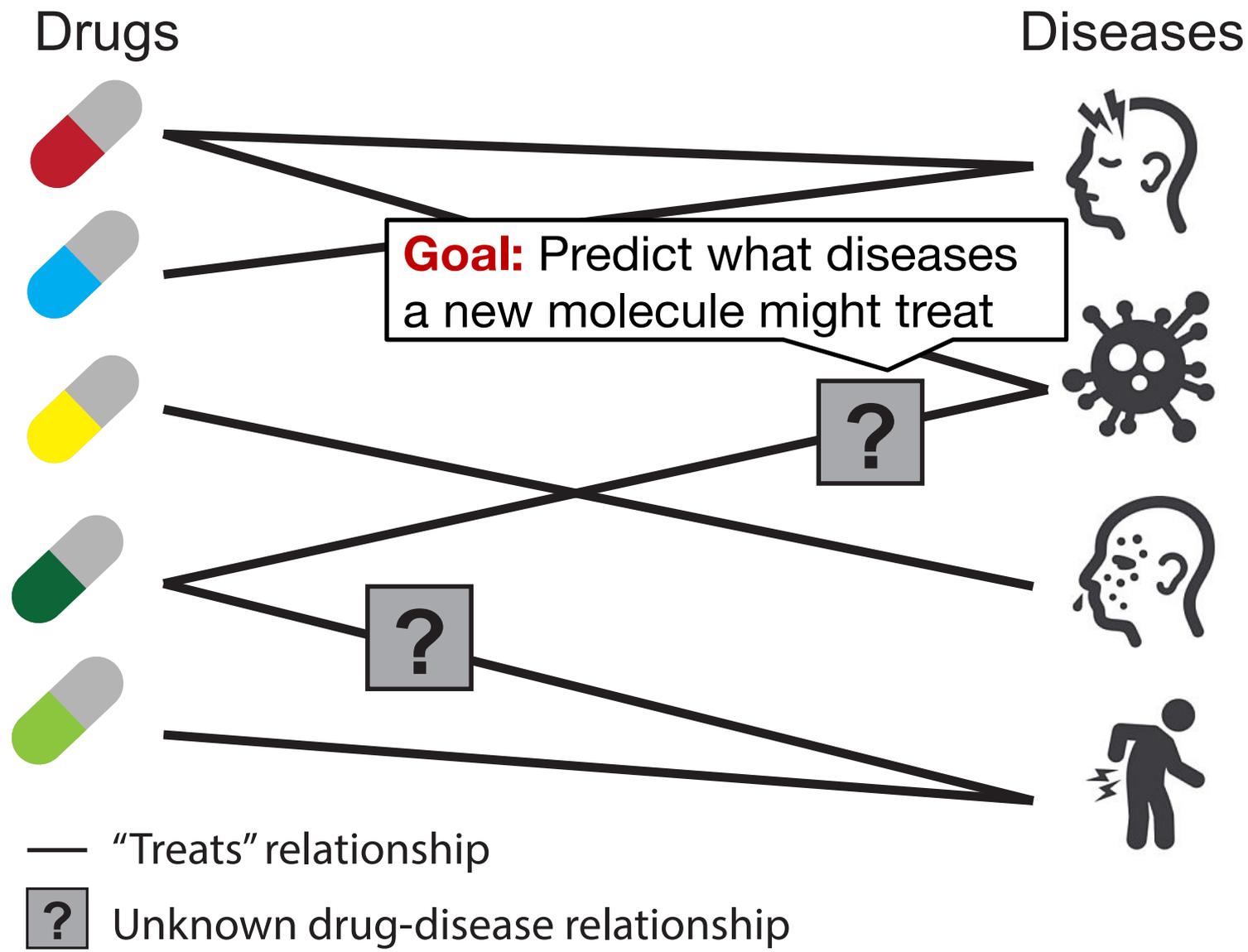12–16 years, **~$1 billion to $2 billion**

## A SHORTER TIMESCALE

Because most repositioned drugs have already passed the early phases of development and clinical testing, they can potentially win approval in less than half the time and at one-quarter of the cost.

**Drug repositioning**
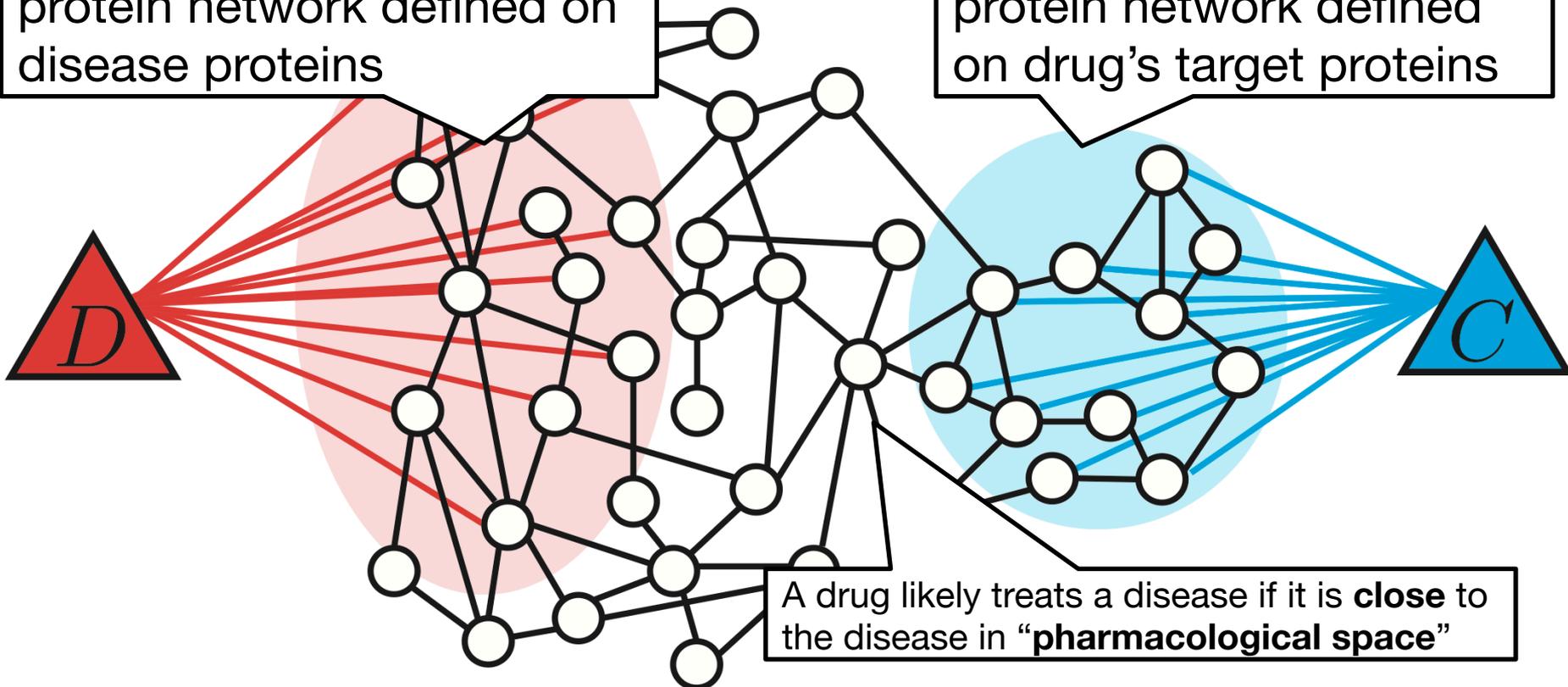
~6 years, **~$300 million**

# What drug treats what disease?

**Drugs**                                                    **Diseases**



**Goal:** Predict what diseases a new molecule might treat

**?**

**?**

—— "Treats" relationship

**?** Unknown drug-disease relationship

# Key Insight: Subgraphs



**Disease:** Subgraph of rich protein network defined on disease proteins

**Drug:** Subgraph of rich protein network defined on drug's target proteins

A drug likely treats a disease if it is **close** to the disease in "**pharmacological space**"

**Idea:** Use the paradigm of embeddings to operationalize the concept of closeness in pharmacological space

# Why Subgraphs? – Part #1

- Analysis of 238 drugs used in 78 diseases
- **Key result:** Therapeutic effect of drugs is localized in a small network neighborhood of disease genes



Guney, E., Menche, J., Vidal, M. and Barábasi, A.L., Network-based in silico drug efficacy screening. Nature Communications, 2016
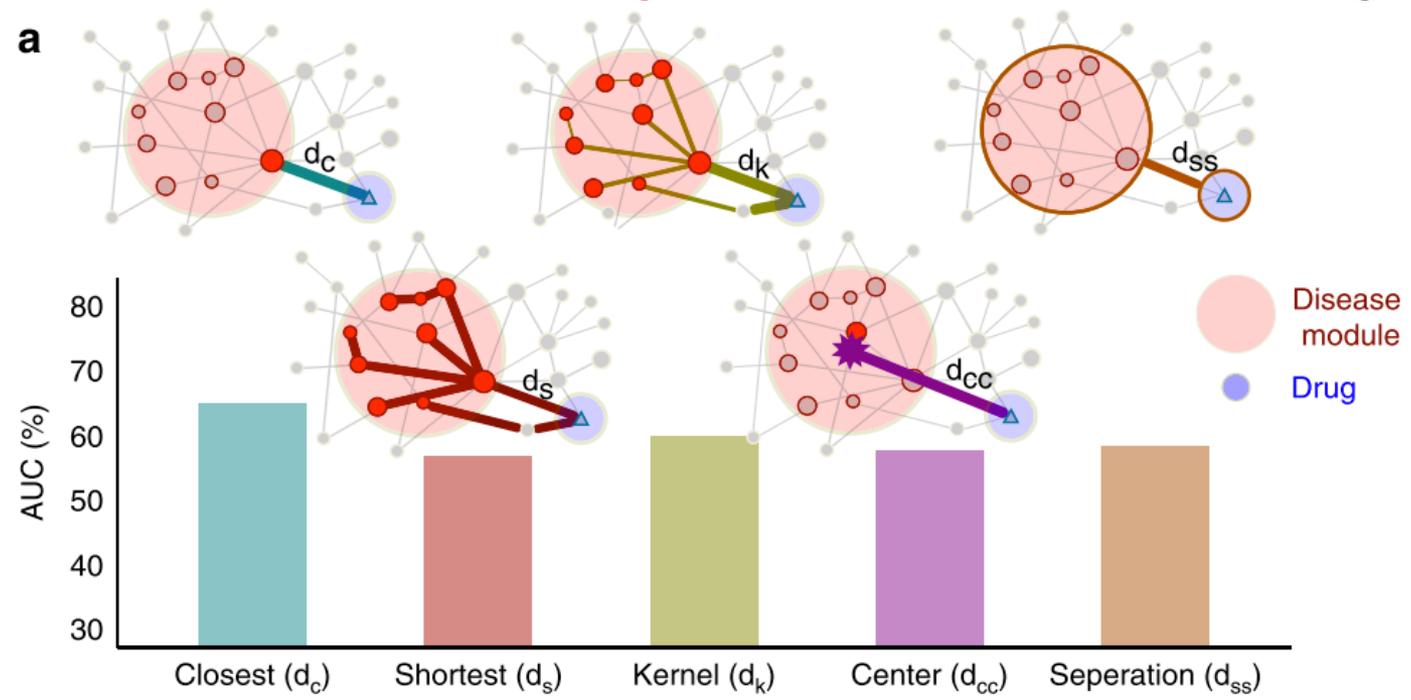
# Why Subgraphs? – Part #2

- Analysis of 238 drugs used in 78 diseases
- **Key result:** Therapeutic effect of drugs is localized in a small network neighborhood of disease genes



Guney, E., Menche, J., Vidal, M. and Barábasi, A.L., Network-based in silico drug efficacy screening. Nature Communications, 2016

# Why Subgraphs? – Part #3

- Analysis of 238 drugs used in 78 diseases
- **Key result:** Therapeutic effect in a small network neighborhood

Negative *z*-values: **Drug targets are close (i.e., proximal) to disease genes** in the PPI network → Successful repurposing
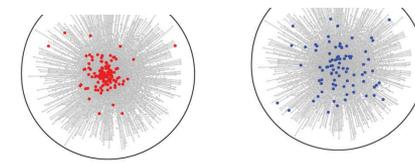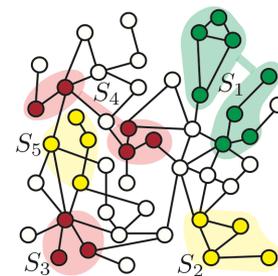
Positive *z*-values: **Drug targets are far away (i.e., not proximal) from disease genes** in the PPI network → Drug failure due to lack of efficacy

**Table 1 | Proximity values for several repurposed and failed drugs.**

| | | Phenotype | Proximity (z) |
|---|---|---|---|
| | | Non-Hodgkin's lymphoma | −2.4 |
| | | Restless legs syndrome | −1.1 |
| | Erectile dysfunction | −1.0 |
| | | Endometrial cancer | −1.1 |
| orgestrel | Confer protection against endometrial cancer | Endometrial cancer | −1.6 |
| **Failures due to lack of efficacy** | | | |
| Tabalumab | Showed lack of efficacy for systemic lupus erythematosus | Systemic lupus erythematosus | 1.8 |
| Preladenant | Discontinued trials for Parkinson due to lack of improvement compared with placebo | Parkinson's disease | 0.2 |
| Iniparib | Failed to achieve improvement while being tested for squamous non-small-cell lung cancer | Squamous cell cancer | 0.0 |
| **Failures due to adverse effetcs** | | | |
| Semagacestat | Failed trials due to worsening AD | AD | −5.6 |
| Terfenadine | Withdrawn due to inducing cardiac arrhythmia | Cardiac arrhythmia | −2.2 |
| | arrhythmia | Arrhythmia (side effect) | −2.6 |

Guney, E., Menche, J., Vidal, M. and Barábasi, A.L., Network-based in silico drug efficacy screening. Nature Communications, 2016

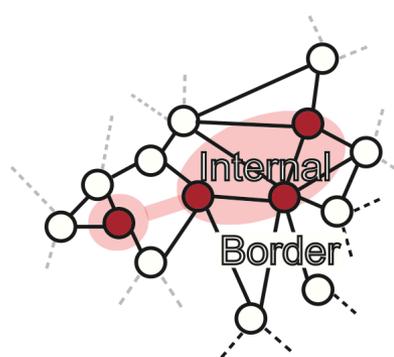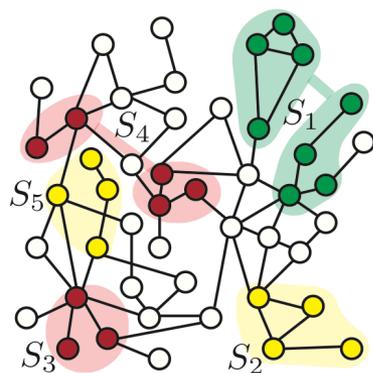# Why are subgraphs challenging?

- Need to predict over structures of varying size:

  - How to represent subgraphs that <u>are not</u> $k$-hop neighborhoods?

- Rich connectivity patterns, both internally a externally through interactions with the rest of $G$:

  - How to inject this information into a GNN?

- Subgraphs can be:

  - Localized and reside in our region of the graph

  - Distributed across multiple local neighborhoods

# Problem Formulation

- **Goal:** Learn subgraph embeddings such that the likelihood of preserving subgraph topology is maximized in the embedding space
  - $S_i$ and $S_j$ with similar subgraph topology should be embedded close together in the embedding space
- SubGNN: Representation learning framework for all key properties of subgraph topology

# SubGNN: Overview

- **SubGNN:** Representation learning framework for all key properties of subgraph topology

- Two key parts:

  - **Part 1:** Hierarchical propagation of information in $G$:

    - Propagate messages from anchor patches to subgraphs

    - Aggregate messages into a final subgraph embedding

  - **Part 2:** Routing of messages through 3 channels, each capturing a distinct property of subgraph topology: position, neighborhood, and structure channels

Emily Alsentzer    Sam Finlayson    Michelle Li

# Part 1: Neural Message Passing

- Property $x$-specific messages $m_x$ are propagated from anchor patch $A_x^q$ to subgraph component $S_i^c$

- Anchor patches are helper subgraphs randomly sampled from $G$; patches $A_P$, $A_N$, and $A_S$ for position, neighborhood and structure
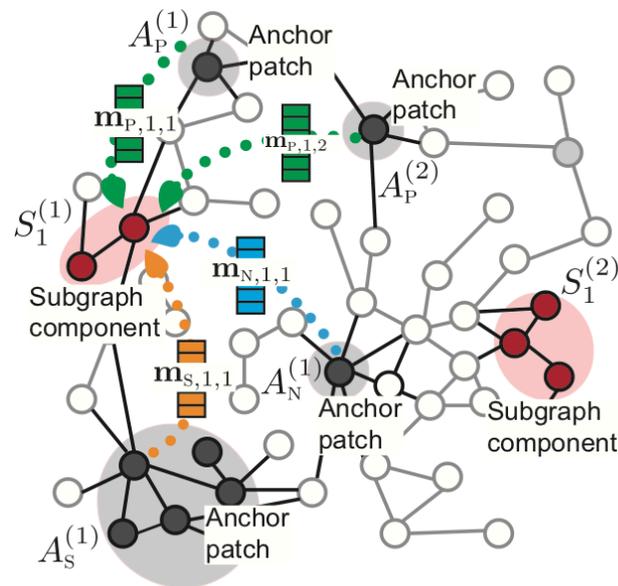
similarity function between a subgraph component and an anchor patch

$$\text{MSG}_\text{X} = \boxed{\gamma_\text{X}} \left( S^{(C)}, A_\text{X} \right) \cdot p_\text{X}$$

$$\mathbf{a}_{\text{X},c} = \text{AGG}_M \left( \left\{ \text{MSG}_\text{X}(S^{(C)}, A_\text{X}, p_\text{X}), \forall A_\text{X} \in \mathcal{A}_\text{X} \right\} \right),$$
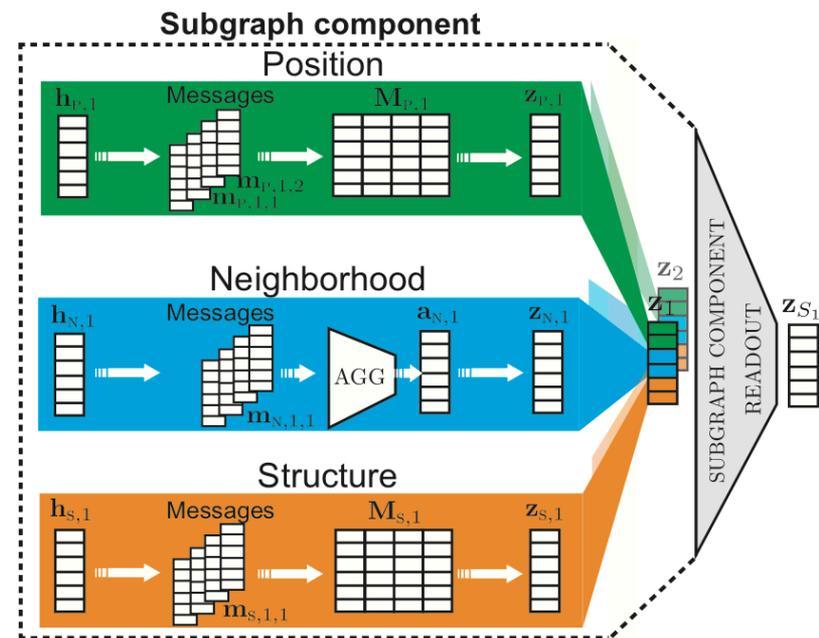
$$\boxed{\mathbf{h}_{\text{X},c}^{(l)}} = \sigma \left( \mathbf{W}_h \cdot [\mathbf{a}_{\text{X},c}; \mathbf{h}_{\text{X},c}^{(l-1)}] \right),$$

property-specific representation of a subgraph component; passed to the next layer

# Part 2: Property-aware Routing

- **SubGNN specifies three channels, each designed to capture a distinct subgraph property**
  - Position, neighborhood, and structure

- **Channel $x$ has three key elements:**
  - Similarity function $\gamma_x$ to weight messages sent between anchor patches and subgraph components
  - Sampling function $\varphi_x$ to generate anchor patches
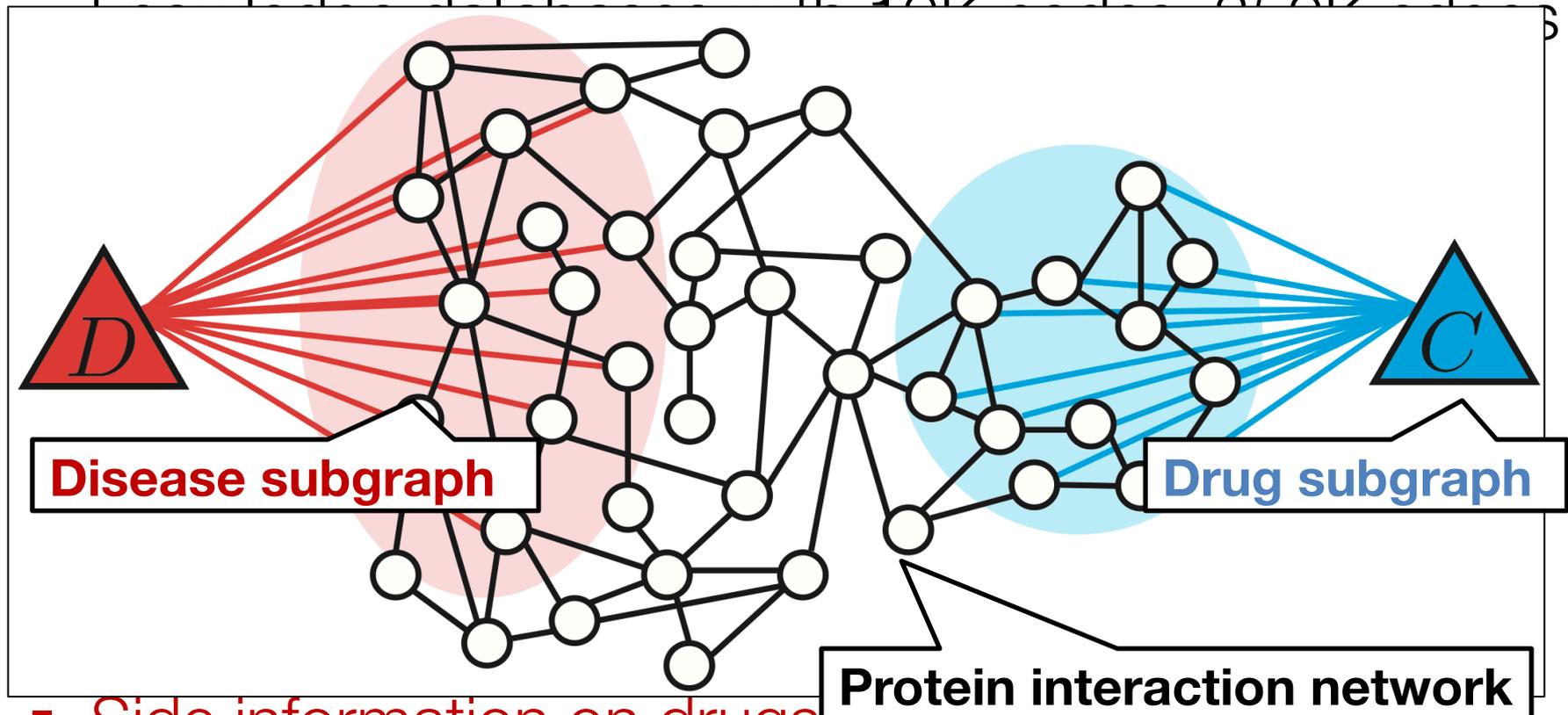  - Anchor patch encoder $\psi_x$



Channel outputs $z_x$ are concatenated to produce **a final subgraph representation $z_S$**

# SubGNN: Overview

# Setup: Drug Repurposing dataset

- **Protein-protein interaction network** culled from 15 knowledge databases with 19K nodes, 350K edges

- **Drug-protein** and **disease-protein** links:
  - DrugBank, OMIM, DisGeNET, STITCH DB and others
  - 20K drug-protein links, 560K disease-protein links

- **Medical indications and contra-indications:**
  - DrugBank, MEDI-HPS, DailyMed, Drug Central, RepoDB
  - 6K drug-disease indications

- **Side information on drugs, diseases, proteins, etc.:**
  - Molecular pathways, disease symptoms, side effects

# Setup: Drug repurposing dataset

- Protein-protein interaction network culled from 15
  knowledge databases with 18K nodes, 350K edges

**Disease subgraph**

**Drug subgraph**

**Protein interaction network**

- Side information on drugs, diseases, proteins, etc..
  - Molecular pathways, disease symptoms, side effects

# Predict links between drug and disease subgraphs



**Task:**
1) Learn embeddings for $C$'s and $D$'s subgraphs
2) Predict whether $C$ should be **indicated** or **contra-indicated** for $D$

# Results: Drug Repurposing

**Drug** | **Disease**

| Drug | Disease | | Rank |
|------|---------|---|------|
| N-acetyl-cysteine | cystic fibrosis | | |
| Xamoterol | neurodegenerat | | |
| Plerixafor | cancer | | |
| Sodium selenite | cancer | Rank: | 36/5000 |
| Ebselen | C difficile | Rank: | 10/5000 |
| Itraconazole | cancer | Rank: | 26/5000 |
| Bestatin | lymphedema | Rank: | 11/5000 |
| Bestatin | pulmonary arterial hypertension | Rank: | 16/5000 |
| Ketaprofen | lymphedema | Rank: | 28/5000 |
| Sildenafil | lymphatic malformation | Rank: | 26/5000 |
| Tacrolimus | pulmonary arterial hypertension | Rank: | 46/5000 |
| Benzamil | psoriasis | Rank: | 114/5000 |
| Carvedilol | Chagas' disease | Rank: | 9/5000 |
| Benserazide | BRCA1 cancer | Rank: | 41/5000 |
| Pioglitazone | interstitial cystitis | Rank: | 13/5000 |
| Sirolimus | dystrophic epidermolysis bullosa | Rank: | 46/5000 |

Stanford MEDICINE | SPARK Translational Research Program *From Bench to Bedside*

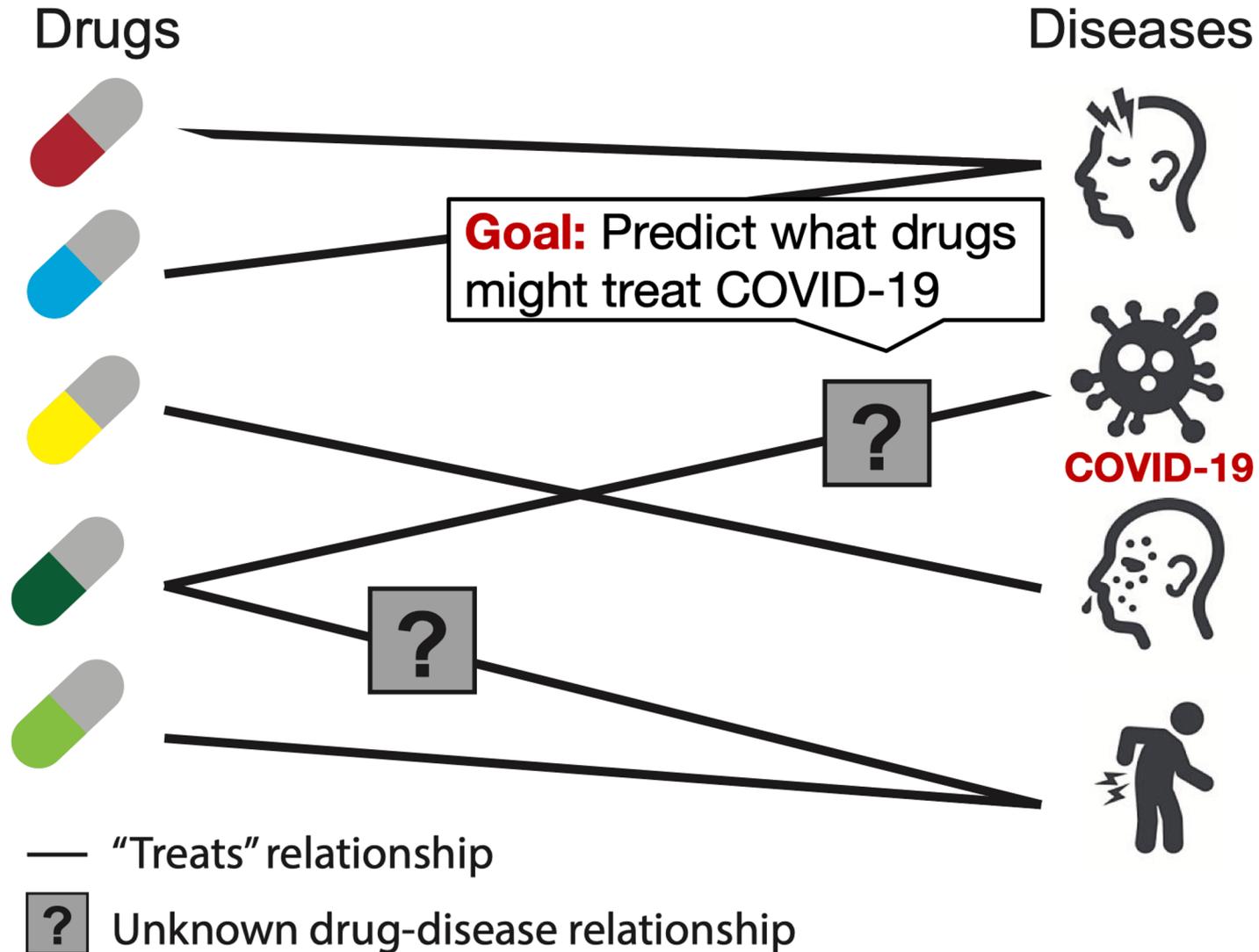**Task:** Predict if an existing drug can be repurposed for a new disease

# Drug Repurposing for Emerging Pathogens

**<span style="color:red">Paper:</span>**

Deisy Morselli Gysi, Ítalo Do Valle, Marinka Zitnik, Asher Ameli, Xiao Gan, et al. Network Medicine Framework for Identifying Drug Repurposing Opportunities for COVID-19, *arXiv:2004.07229*
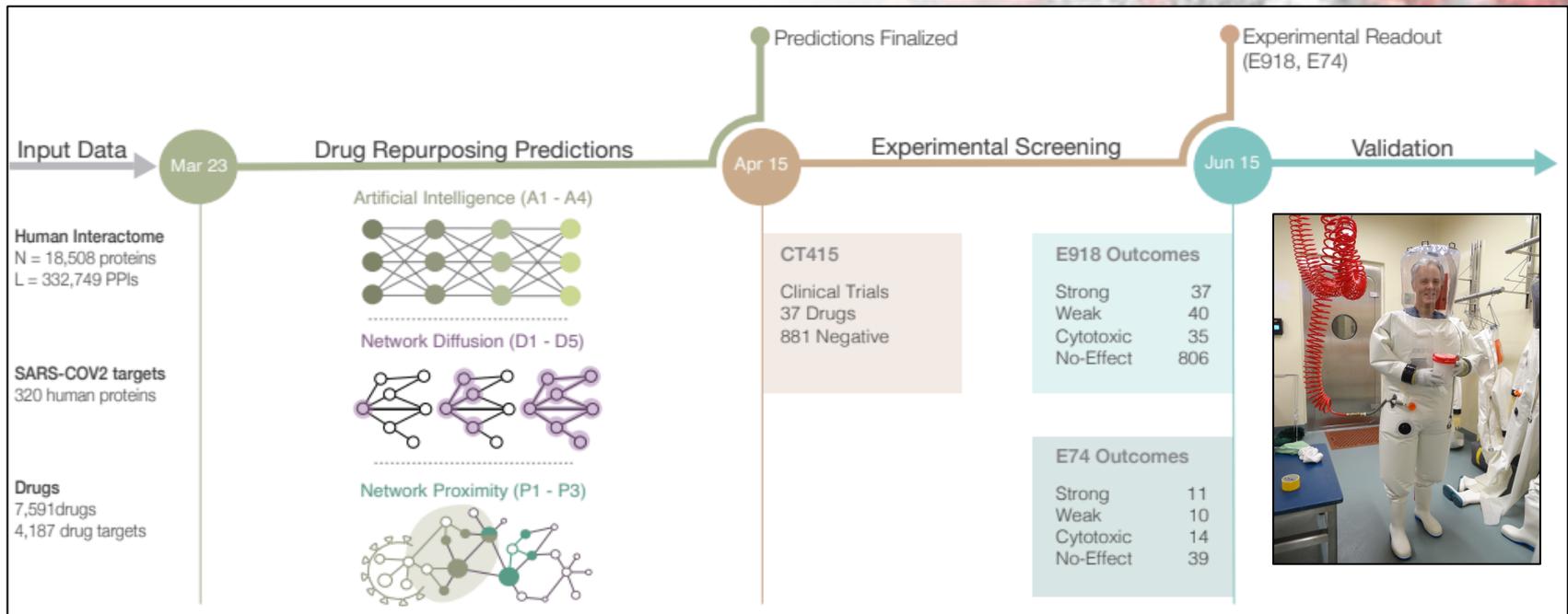
# Emerging Pathogens



Drugs — Diseases

**Goal:** Predict what drugs might treat COVID-19

COVID-19

— "Treats" relationship

? Unknown drug-disease relationship

# Never-Before-Seen Disease

The traditional approach of iterative development, experimental testing, clinical validation, and approval of new drugs are not feasible
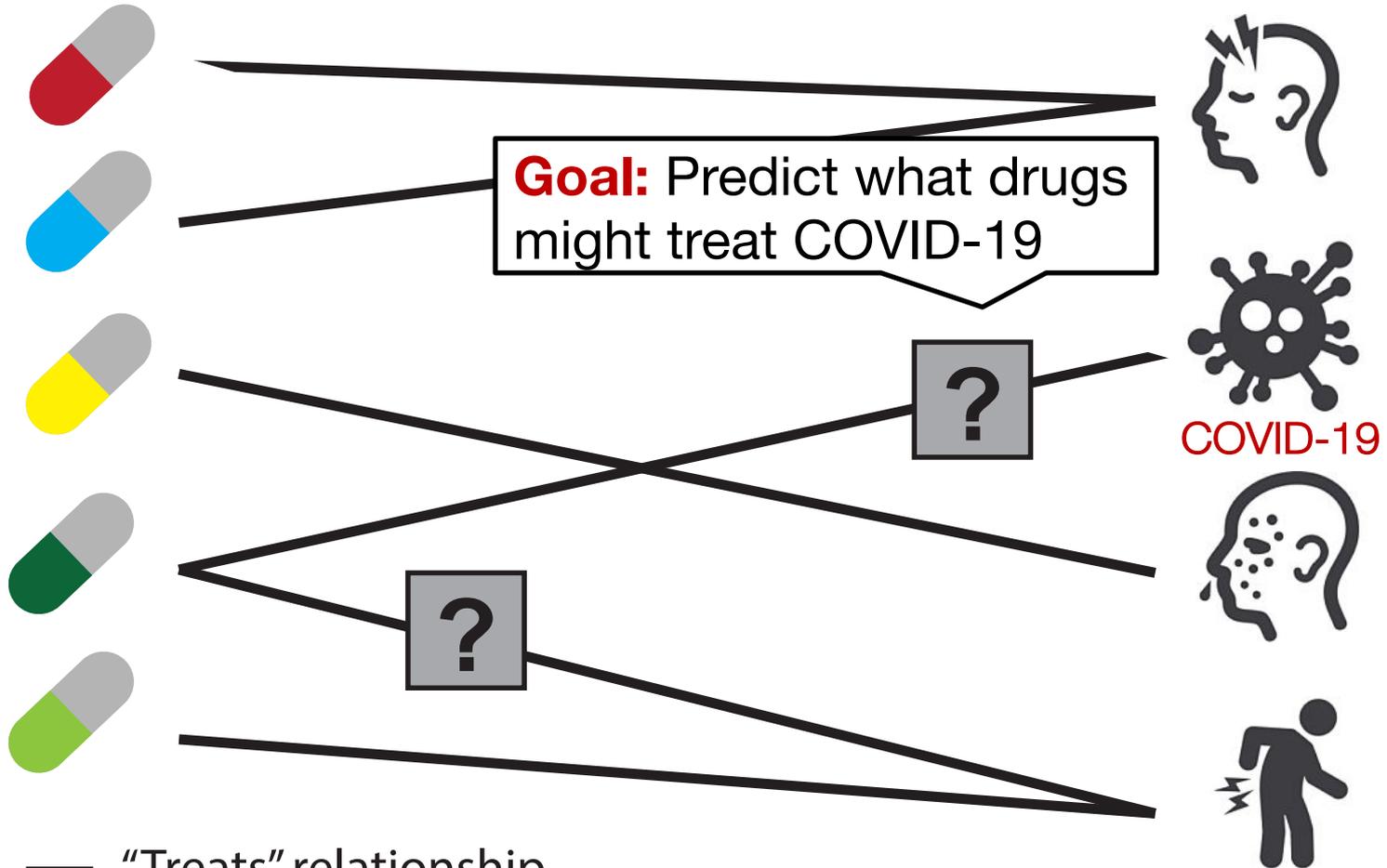
A more realistic strategy relies on drug repurposing, requiring us to identify clinically approved drugs that have a therapeutic effect in COVID-19 patients



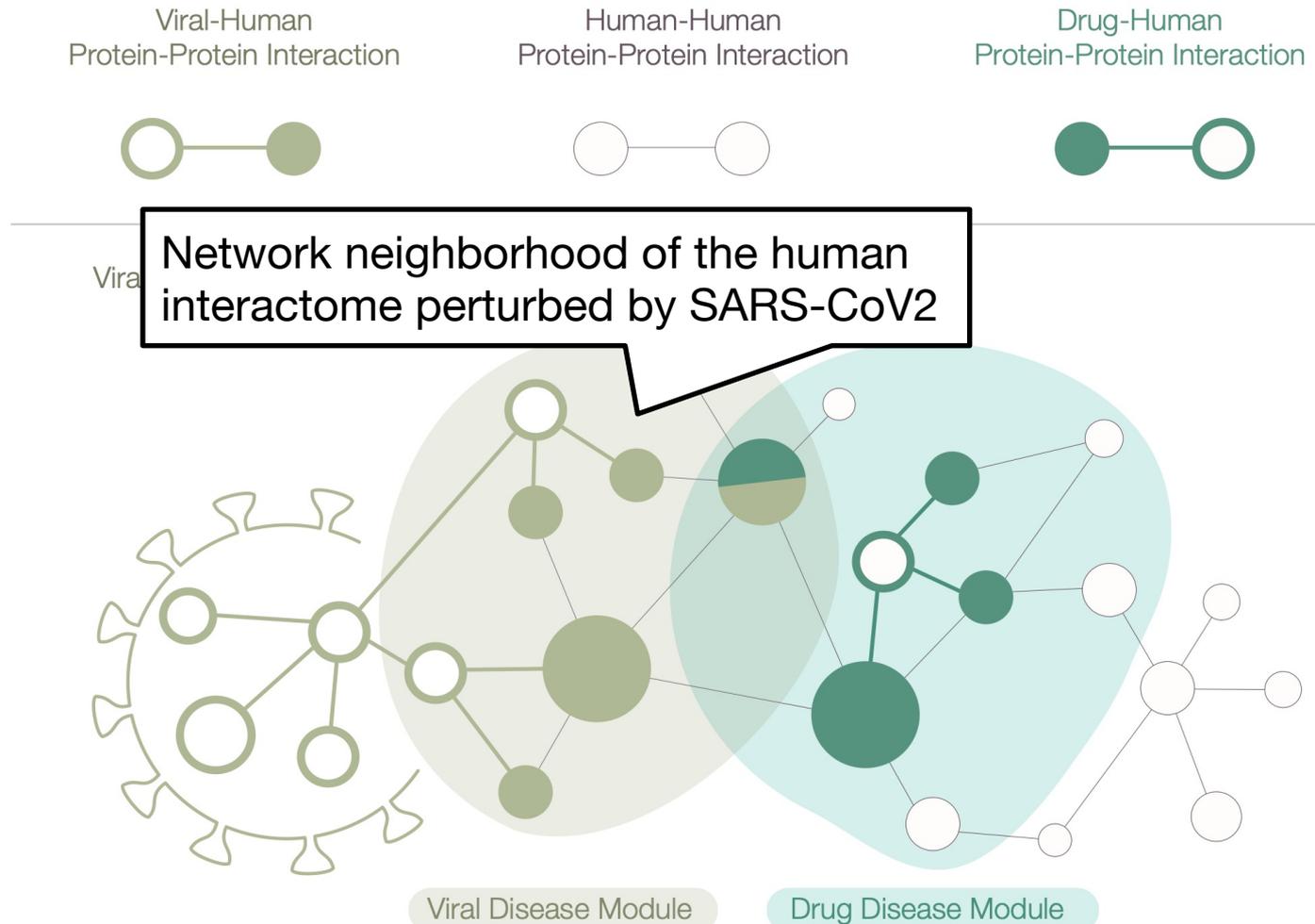Network Medicine Framework for Identifying Drug Repurposing Opportunities for Covid-19, *arXiv:2004.07229*

# Never-before-seen disease

Drugs

Diseases

**Goal:** Predict what drugs might treat COVID-19

COVID-19

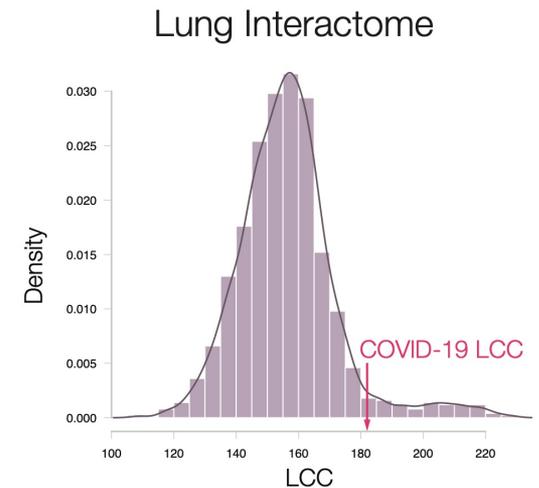—— "Treats" relationship

? Unknown drug-disease relationship
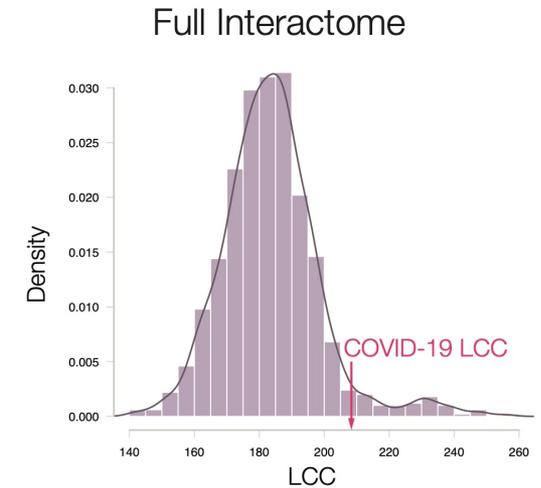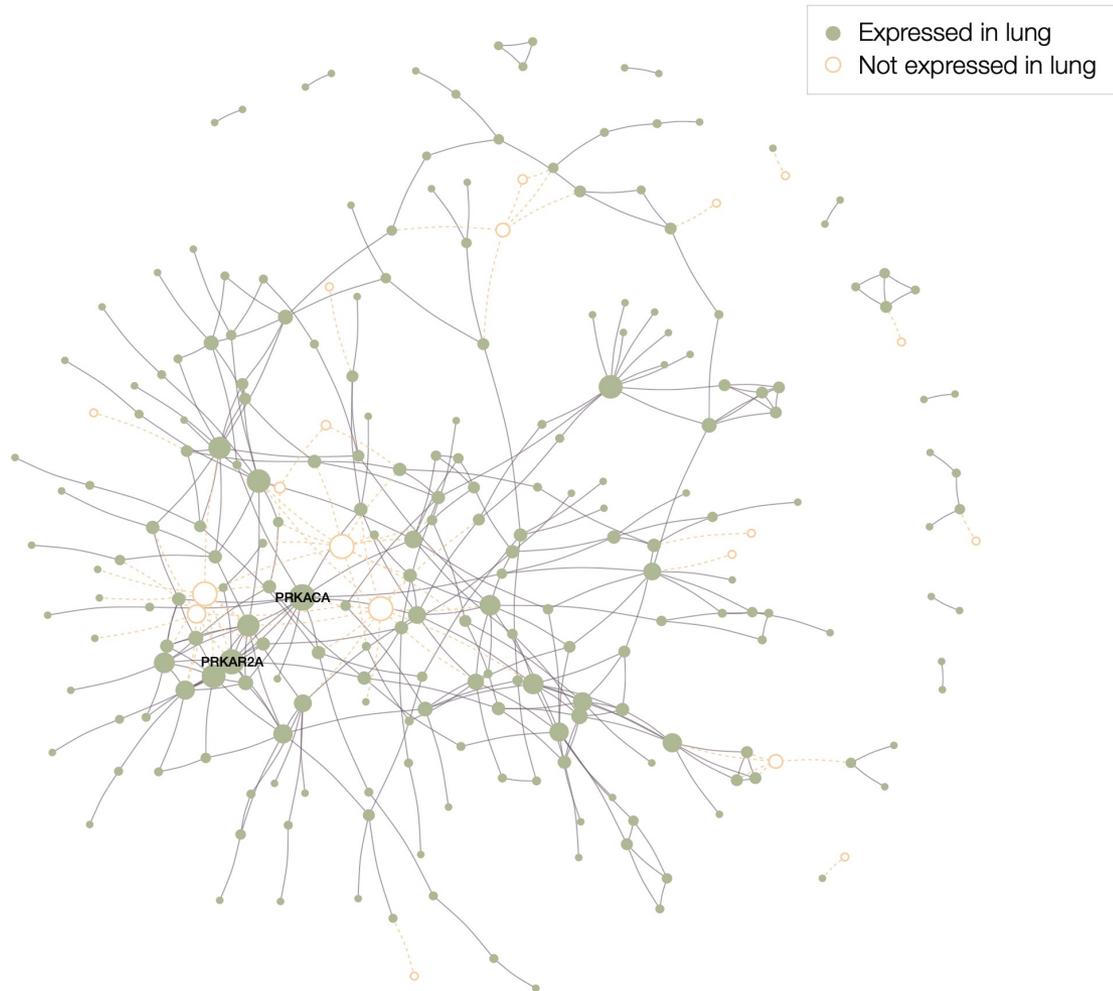
# How to represent COVID-19? Map SARS-CoV2 targets to the human interactome

# COVID-19 Subgraph



Gordon et al., Nature 2020 expressed 26 of the 29 SARS-CoV2 proteins and used AP-MS to identify 332 human proteins to which viral proteins bind

# Key Insight: Subgraphs



**Disease:** Subgraph of rich protein network defined on disease proteins

**Drug:** Subgraph of rich protein network defined on drug's target proteins

A drug likely treats a disease if it is **close** to the disease in **pharmacological space** [Paolini et al., Nature Biotech.'06; Menche et al., Science'15]

**Idea:** Use the paradigm of embeddings to operationalize the concept of closeness in pharmacological space

# Computational Setup

- Proxy for ground-truth information:
  - Monitor drugs under clinical trials
  - Capture the medical community's assessment of drugs

# Embedding Space



Closest drugs in the embedding space

| | |
|---|---|
| Atovaquone | Teriflunomide |
| Rifapentine | Ixekizumab |
| Chloroquine | Praziquantel |
| Mifepristone | Ritonavir |
| Lindane | Troleandomycin |
| Secukinumab | Budesonide |
| Elbasvir | Loxoprofen |
| Cobicistat | Fludrocortisone |
| Idelalisib | Crizotinib |
| Daclatasvir | Elvitegravir |

# Results: COVID-19 Repurposing

## Individual ROC



| | |
|---|---|
| **GNN** | |
| A1: | 0.86 |
| A2: | 0.86 |
| A3: | 0.87 |
| A4: | 0.86 |
| D1: | 0.56 |
| D2: | 0.56 |
| D3: | 0.55 |
| D4: | 0.56 |
| D5: | 0.55 |
| P1: | 0.68 |
| P2: | 0.58 |
| P3: | 0.70 |
| Random: | 0.50 |

We test each pipeline's ability to recover drugs currently in clinical trials for COVID-19

The best individual ROC curves are obtained by the GNN methods

The second-best performance is provided by the proximity P3. Close behind is P1 with AUC = 0.68 and AUC = 0.58

Diffusion methods offer ROC between 0.55-0.56

# Final Prediction Model – Part #1

| Input Data | Methods | Outcomes |
|---|---|---|
| Human Interactome N = 18,508 proteins L = 332,749 PPIs | Network Proximity 3 pipelines | Infected Tissues/Organs |
| SARS-COV2 targets 320 human proteins Gordon et al, 2020 | Network Diffusion 5 pipelines | Comorbidity |
| Drug Targets 7,591drugs 4,187 drug targets DrugBank | AI Prioritization 4 pipelines | Drug Repurposing & Validation |

# Final Prediction Model – Part #2

## Methods

- A COVID-19 treatment can not be derived from the arsenal of therapies approved for specific diseases

- Repurposing strategies focus on drugs previously approved for other pathogens, or on drugs that target the human proteins to which viral proteins bind.

- Most approved drugs do not target directly disease proteins but bind to proteins in their network vicinity
- [Yildirim, Nature Biotech. 2007]

- Identify drug candidates that have the potential to perturb the network vicinity of the COVID-19 disease module.

- Implement 3 Network Repurposing Methods.

**Network Proximity**
3 pipelines

**Network Diffusion**
5 pipelines

**AI Prioritization**
4 pipelines

# Final Prediction Model – Part #3

**Rank Aggregation Algorithm:** Maximize the number of pairwise agreements between the final ranking and each input ranking.

The combined performance of the AI methods is 0.87, the same as A3.

Improvement for proximity pipelines: 0.70 → 0.72.

Combined diffusion pipelines have lower performance (0.54 vs 0.56, for D1, D2, and D4).
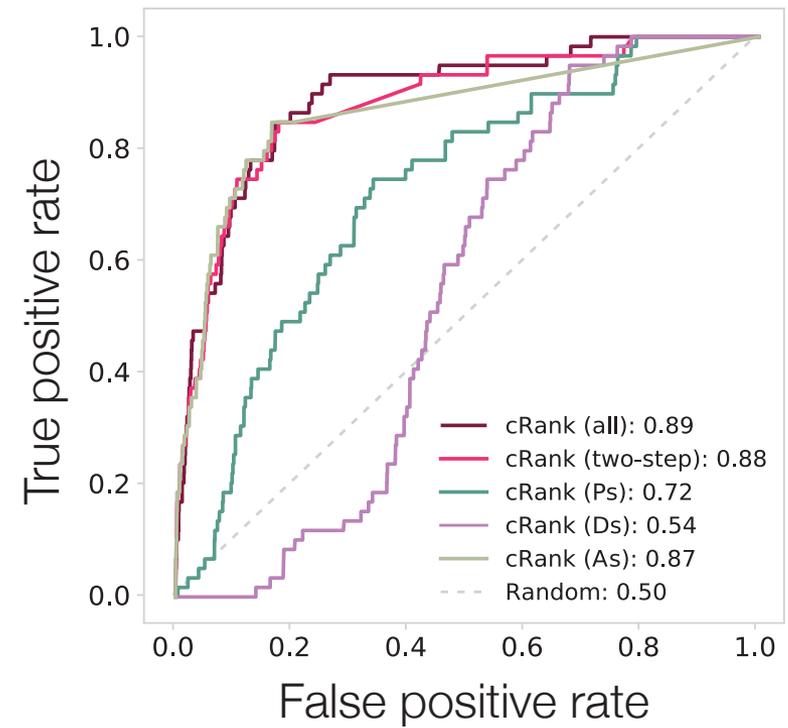
**Combining all 12 pipelines, gives AUROC=0.89, the highest of any individual or combination-based pipelines,**

Individual pipelines offer complementary information harnessed by the combined ranking.



## Combined ROC

# Predicted Drug Candidates

Joseph Loscalzo

◯ # of Clinical trials from ClinicalTrials.gov

86 drugs selected from the top 10% of the rank list.

Respiratory drugs (e.g., theophylline, montelukast).

Cardiovascular systems (e.g., verapamil, atorvastatin).

Antibiotics used to treat viral (e.g., ribavirin, lopinavir), parasitic (e.g., hydroxychloroquine, ivermectin, praziquantel), bacterial (e.g., rifaximin, sulfanilamide), mycotic (e.g., fluconazole), and mycobacterial (e.g., isoniazid) infections.

Immunomodulating/anti-inflammatory drugs (e.g., interferon-β, auranofin, montelukast, colchicine)

Anti-proteasomal drugs (e.g., bortezomib, carfilzomib)

Less obvious choices: aminoglutethimide, melatonin, levothyroxine, calcitriol, selegiline, deferoxamine, mitoxantrone, metformin, nintedanib, cinacalcet, and sildenafil.

| | Drug | C-rank | | Drug | C-rank | | Drug | C-rank |
|---|---|---|---|---|---|---|---|---|
| (20) | Ritonavir | 1 | | Mesalazine | 69 | | Sulfanilamide | 265 |
| | Isoniazid | 2 | | Pentamidine | 92 | | Hydralazine | 269 |
| | Troleandomycin | 3 | | Verapamil | 98 | | Gemfibrozil | 281 |
| | Cilostazol | 4 | | Melatonin | 109 | (4) | Ruxolitinib | 284 |
| (76) | Chloroquine | 5 | | Griseofulvin | 112 | | Propranolol | 297 |
| | Rifabutin | 6 | | Auranofin | 118 | | Carbamazepine | 301 |
| | Flutamide | 7 | (1) | Atovaquone | 124 | | Doxorubicin | 309 |
| (2) | Dexamethasone | 8 | | Montelukast | 131 | | Levothyroxine | 329 |
| | Rifaximin | 9 | | Romidepsin | 138 | | Dactinomycin | 335 |
| | Azelastine | 10 | (1) | Cobicistat | 141 | | Tenofivir | 338 |
| | Folic Acid | 16 | (17) | Lopinavir | 146 | | Tadalafil | 339 |
| | Rabeprazole | 27 | | Pomalidomide | 155 | | Doxazosin | 367 |
| | Methotrexate | 32 | | Sulfinpyrazone | 157 | | Rosiglitazone | 397 |
| | Digoxin | 33 | (1) | Levamisole | 161 | | Aminolevulinic acid | 398 |
| | Theophylline | 34 | | Calcitriol | 164 | | Nitroglycerin | 418 |
| | Fluconazole | 41 | (1) | Interferon-β-1a | 173 | | Metformin | 457 |
| | Aminoglutethimide | 42 | | Praziquantel | 176 | (1) | Nintedanib | 466 |
| (67) | Hydroxychloroquine | 44 | (1) | Ascorbic acid | 195 | | Allopurinol | 471 |
| | Methimazole | 47 | | Fluvastatin | 199 | | Ponatinib | 491 |
| (1) | Ribavirin | 49 | (1) | Interferon-β-1b | 203 | (1) | Sildenafil | 493 |
| (1) | Omeprazole | 50 | | Selegiline | 206 | | Dapagliflozin | 504 |
| | Bortezomib | 53 | (1) | Deferoxamine | 227 | | Nitroprusside | 515 |
| | Leflunomide | 54 | | Ivermectin | 235 | | Cinacalcet | 553 |
| | Dimethylfumarate | 55 | (1) | Atorvastatin | 243 | | Mexiletine | 559 |
| (4) | Colchicine | 57 | | Mitoxantrone | 250 | | Sitagliptin | 706 |
| | Quercetin | 63 | | Glyburide | 259 | | Carfilzomib | 765 |
| | Mebendazole | 67 | (2) | Thalidomide | 262 | (1) | Azithromycin | 786 |

# Experimental Validation of Predictions



National Emerging Infectious Diseases Laboratories (NEIDL)

| CRank | Drug Name |
|-------|-----------|
| 1 | Ritonavir |
| 2 | Isoniazid |
| 3 | Troleandomycin |
| 4 | Cilostazol |
| 5 | Chloroquine |
| 6 | Rifabutin |
| 7 | Flutamide |
| 8 | Dexamethasone |
| 9 | Rifaximin |
| 10 | Azelastine |
| 11 | Crizotinib |

| | |
|-----|-----------|
| 17 | Celecoxib |
| 18 | Betamethasone |
| 19 | Prednisolone |
| 20 | Mifepristone |
| 21 | Budesonide |
| 22 | Prednisone |
| 23 | Oxiconazole |
| 24 | Megestrol acetate |
| 25 | Idelalisib |
| 26 | Econazole |
| 27 | Rabeprazole |

Ranked lists of drugs

New algorithms:
Prioritizing Network Communities, *Nature Communications* 2018
Subgraph Neural Networks, *NeurIPS* 2020
Graph Meta Learning via Local Subgraphs, *NeurIPS* 2020

**Results:** 918 compounds screened for their efficacy against SARS-CoV-2 in VeroE6 cells:
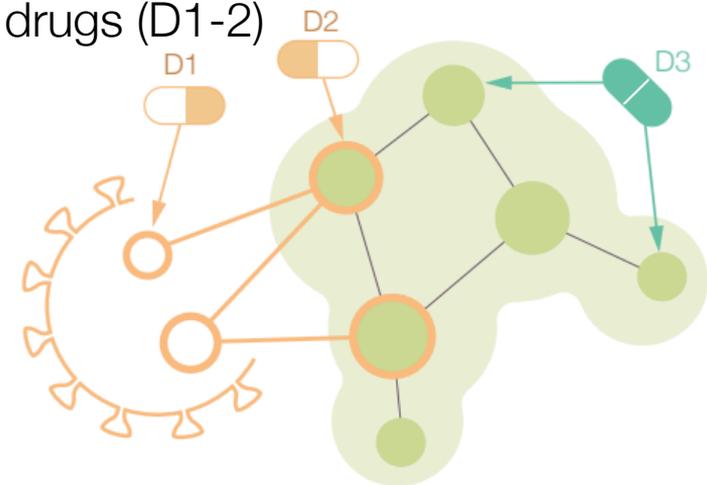
- **37 had a strong effect** being active over a broad range of concentrations

- **40 had a weak effect** on the virus

- An <u>order of magnitude higher hit rate</u> among top 100 drugs than <u>prior work</u>

# Results: Network Drugs

- **76/77 drugs that successfully reduced viral infection <u>do not bind</u> proteins targeted by SARS-CoV-2:**
  - These drugs rely on **network-based actions** that cannot be identified by docking-based strategies

Strong
Weak

| CRank | Drug Name |
|---|---|
| 5 | Chloroquine |
| 6 | Rifabutin |
| 9 | Rifaximin |
| 10 | Azelastine |
| 16 | Folic acid |
| 32 | Methotrexate |
| 33 | Digoxin |
| 44 | Hydroxychloroquine |
| 50 | Omeprazole |
| 113 | Clobetasol propionate |
| 118 | Auranofin |
| 120 | Vinblastine |
| 199 | Fluvastatin |
| 210 | Clomifene |
| 233 | Ibuprofen |
| 235 | Ivermectin |
| 243 | Atorvastatin |
| 253 | Pralatrexate |
| 263 | Cobimetinib |
| 269 | Hydralazine |
| 297 | Propranolol |
| 317 | Osimertinib |
| 348 | Vincristine |
| 367 | Doxazosin |
| 397 | Rosiglitazone |
| 398 | Aminolevulinic acid |

| CRank | Drug Name |
|---|---|
| 423 | Pitavastatin |
| 431 | Tenoxicam |
| 438 | Quinidine |
| 456 | Sertraline |
| 460 | Ingenol mebutate |
| 463 | Norelgestromin |
| 493 | Sildenafil |
| 499 | Eliglustat |
| 518 | Ulipristal |
| 553 | Cinacalcet |
| 556 | Perphenazine |
| 558 | Idarubicin |
| 564 | Perhexiline |
| 569 | Amiodarone |
| 577 | Duloxetine |
| 585 | Toremifene |
| 586 | Afatinib |
| 601 | Amitriptyline |
| 626 | Meclizine |
| 635 | Valsartan |
| 651 | Eletriptan |
| 673 | Sotalol |
| 678 | Thioridazine |
| 695 | Chlorcyclizine |
| 707 | Omacetaxine mepesuccinate |
| 721 | Candesartan |

| CRank | Drug Name |
|---|---|
| 742 | Mianserin |
| 755 | Clofazimine |
| 767 | Chlorpromazine |
| 772 | Imipramine |
| 830 | Promazine |
| 900 | L-Alanine |
| 917 | Moxifloxacin |
| 933 | Tasimelteon |
| 995 | Vandetanib |
| 1000 | Azilsartan medoxomil |
| 1020 | Frovatriptan |
| 1034 | Zolmitriptan |
| 1035 | Procarbazine |
| 1093 | Asenapine |
| 1107 | Dyclonine |
| 1140.5 | Clemastine |
| 1194 | Prochlorperazine |
| 1222 | Miglustat |
| 1224 | Prenylamine |
| 1276 | Dalfampridine |
| 1314 | Cinchocaine |
| 1355 | Methotrimeprazine |
| 1396 | Methylthioninium |
| 1403 | Metixene |
| 1443 | Trifluoperazine |

Direct target drugs (D1-2)

D1   D2   D3

SARS-CoV-2 Viral Interactome

Human Interactome

Network drugs (D3)

58/77 drugs with positive experimental outcome are among top 750 ranked drugs

# Outline

✓ Overview and introduction

✓ Part 1: Virtual drug screening and drug repurposing

Part 2: Adverse drug effects, drug-drug interactions

Part 3: Clinical trial site identification, patient recruitment

Part 4: Molecule optimization, molecular graph generation, multimodal graph-to-graph translation

Part 5: Molecular property prediction and transformers

Demos, resources, wrap-up & future directions