



Graph in Drug Discovery

Advanced Training
2020. 11. 26

Erkhembayar J. Ph.D

*Korea AI Center for
Drug Discovery and Development*



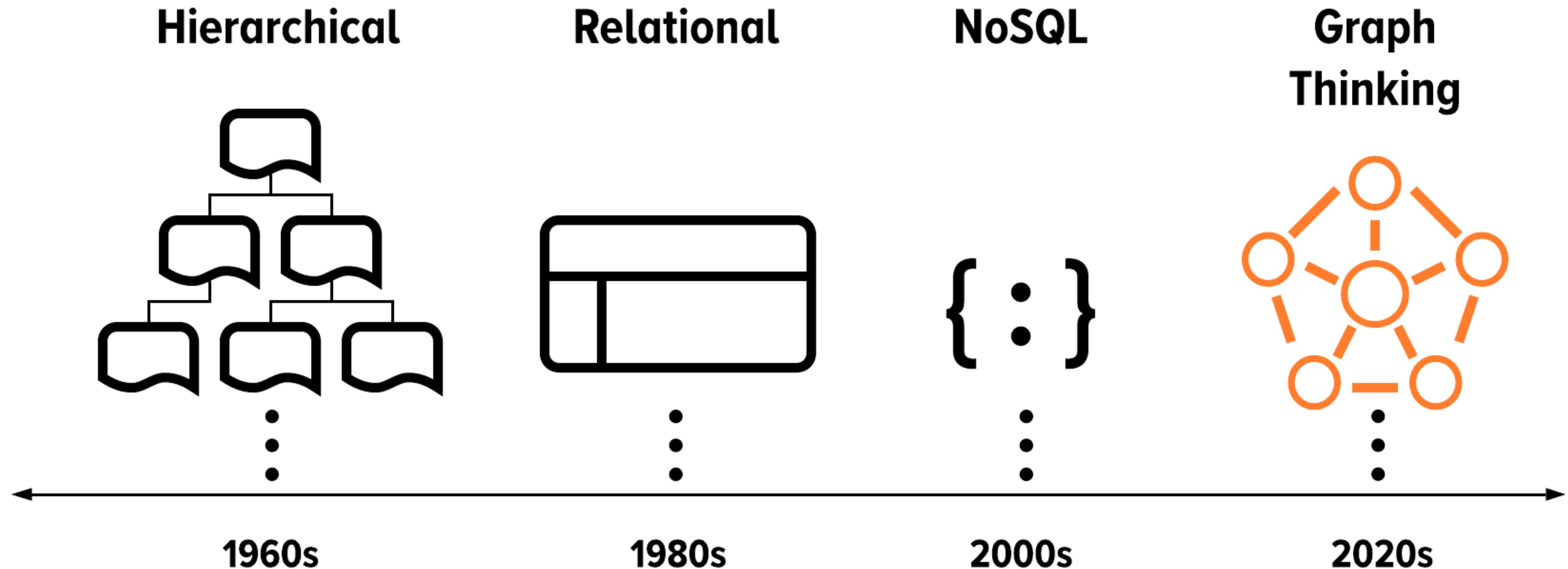
AGENDA

	HOURS	CONTENT
GRAPH IN DRUG DISCOVERY	GENERAL OVERVIEW 09.00 – 9.40	Why Graph is important ? How the graphs are used in Drug Discovery? COVID19: Drug Repositioning Knowledge Graph
CODING SESSION (colab)	CODING SESSION 10:00 - 11.40	How to Construct your own Graph? Building Compound to Graph from Scratch
CODING SESSION (colab)	CODING SESSION 11:40 - 12.00	Knowledge -based Drug Repositioning Source Code Understanding
	LUNCH	

GRAPH IN DRUG DISCOVERY

General Overview of Graph Application in the Field

Era of Data Representation



Denise et. al. 2020

Graphs

- Complex problems
 - are the individual problems that are observable and measurable within **complex systems**
- Complex systems
 - System that composed of many **individual components** that interact each other

GRAPHS are good at solving complex problem in complex systems

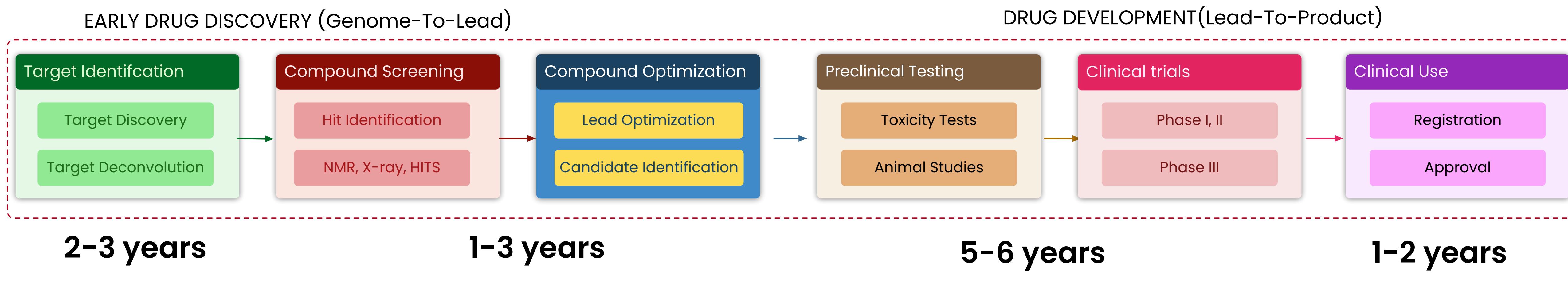
- Each component represented as graph
- Interaction between component could defined as edges
- ML/AI applied to individual components/or whole system to answer different questions

Discipline	System	Network Representation	Network Specifics
Systems Biology			Metabolic pathway map drawn from biochemical knowledge and biological intuition vs. metabolite-centric network representation of this pathway map. Node: gene or regulatory enzyme Links represent regulation, gene enzyme association or metabolic reaction.
Neuroscience & Computational Neuroscience			Node: individual neuron or cortical area Links between nodes. In spatially embedded networks distance between nodes is important.
Geomorphology			Node: spatial location within geomorphic network Link: unidirectional weighted links Sub-regions of the network represent spatially discrete sub-catchments
Ecology			Node: representing a patch, organisms or population Link: bi-directional links representing connected pathways between patches
Social Network Science			Node: representing a person (or interaction) Links: non-directional or bi-directional links Groups of nodes forming a community

Connectivity and complex systems: learning from a multi-disciplinary perspective 2018.

Drug Discovery is Complex System

- Drug Discovery is complex system with complex problem



Expensive

~\$2.8B

DiMasi et al.: Innovation in the Pharmaceutical Industry: New Estimates of R&D Costs. J. Health Econ. 2016, 47, 20-33

Time Consuming

~11-16 years

Matthews, et al.: "Omics"-Informed Drug and Biomarker Discovery: Opportunities, Challenges and Future Perspectives. Proteomes 2016, 4, 28.

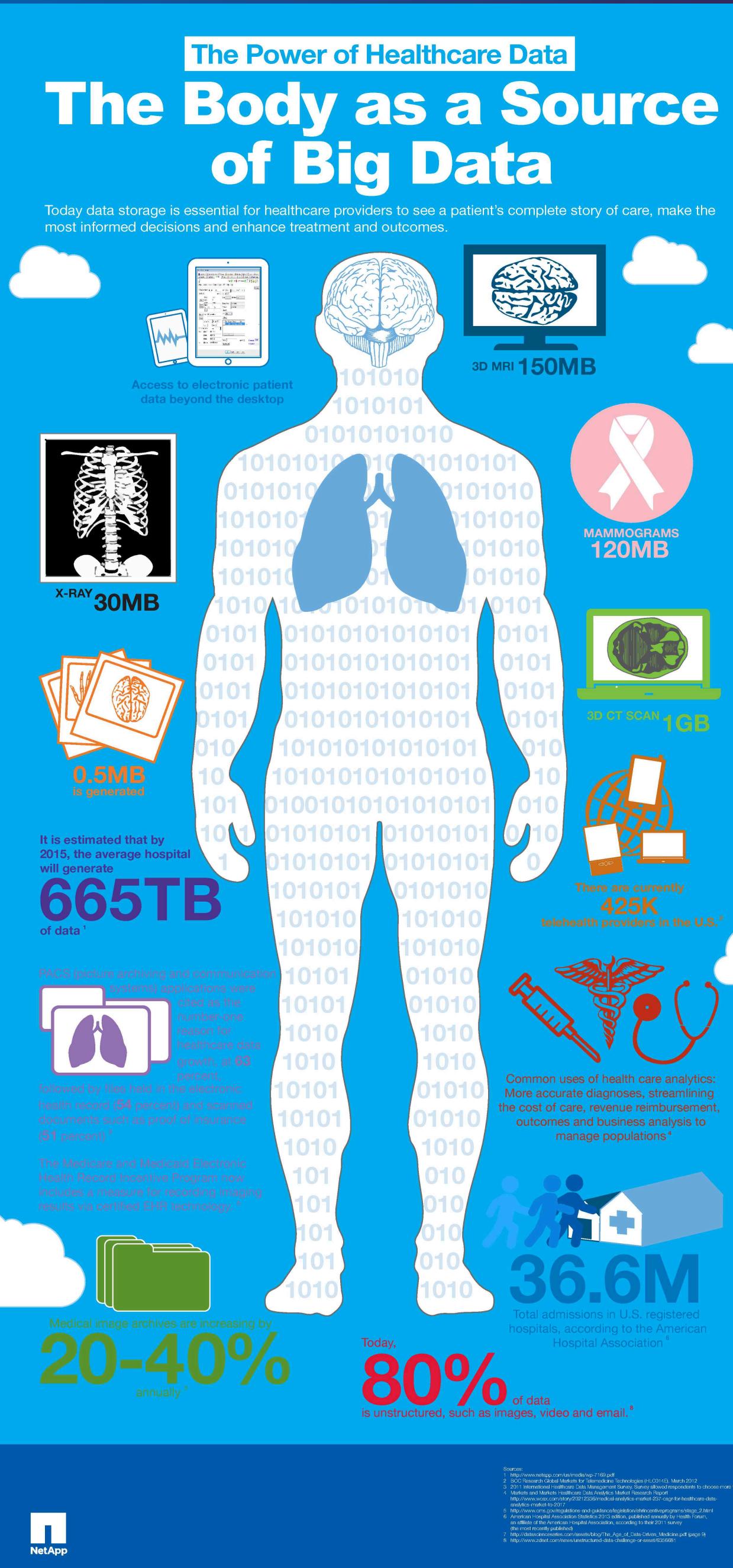
High Failure Rate

90%

Dowden, et al.: Trends in Clinical Success Rates and Therapeutic Focus. Nat. Rev. Drug Discovery 2019, 18, 495-496.

GRAPHS are good at solving complex problem in complex systems

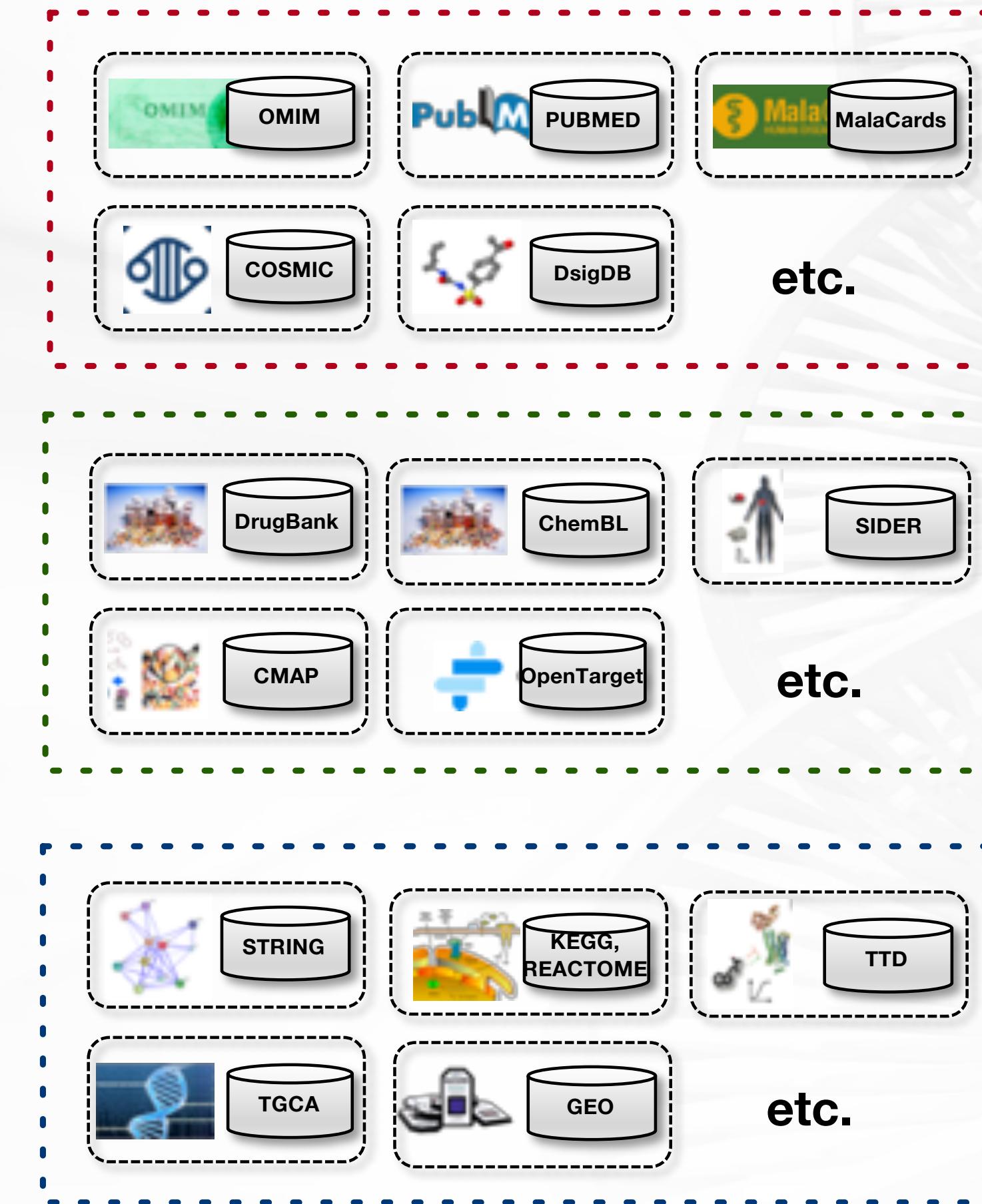
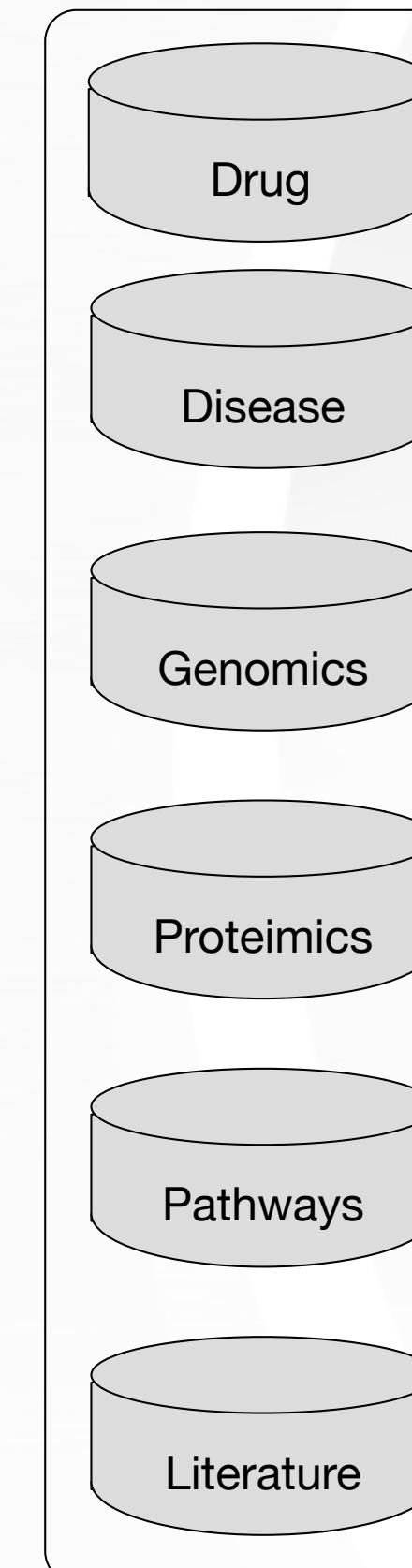
RECENTLY GRAPH + DL HAVE PROMISING RESULT in **EACH DRUG DISCOVERY COMPONENTS**



Drug Discovery Data Space

KAICD
한국인공지능신약개발지원센터
Korea AI Center for Drug Discovery and Development

high-content, high-dimensional biological data



Disease

Drug

Protein

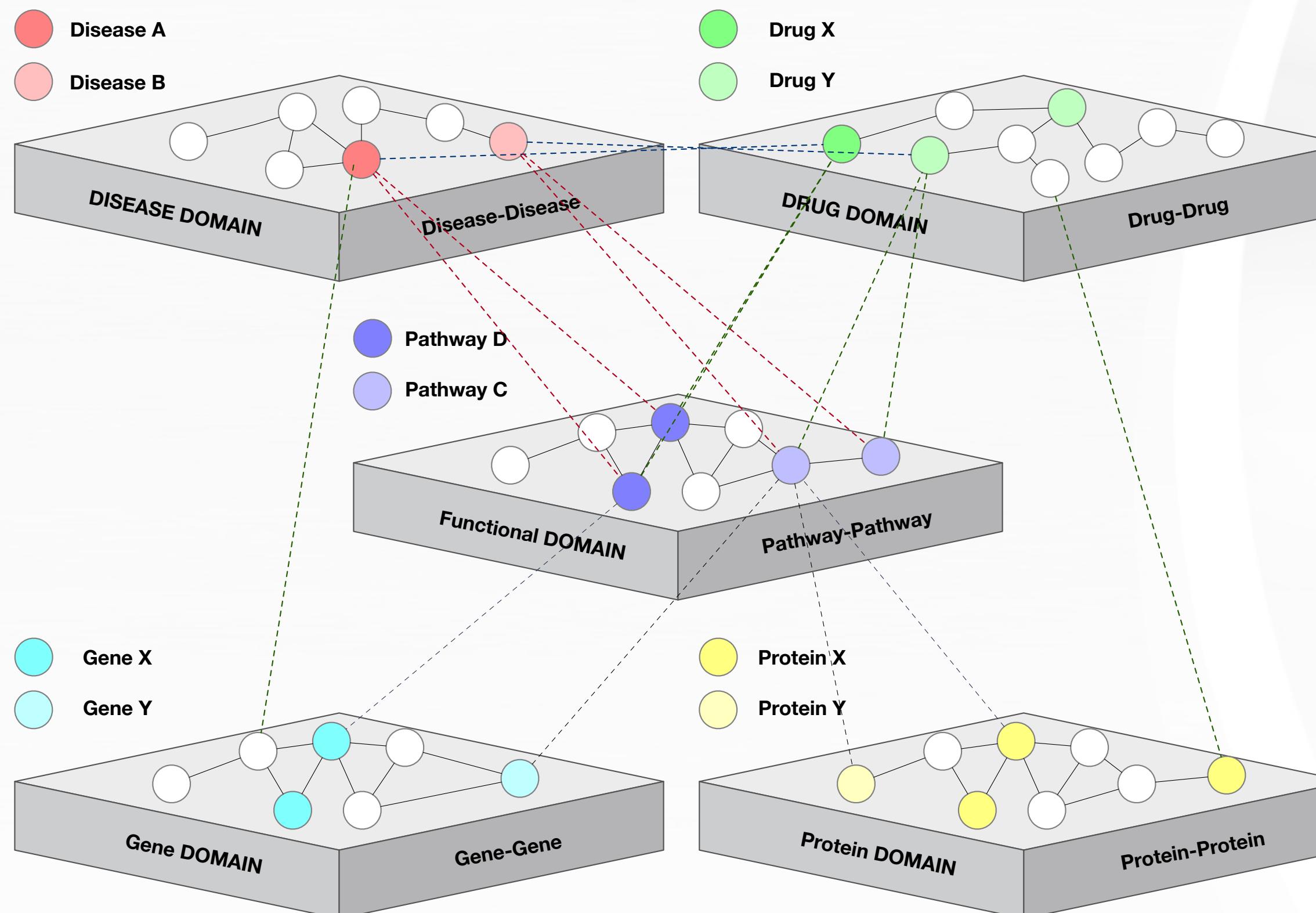
Genomics

Pathways

GRAPHS are representing complex data domain

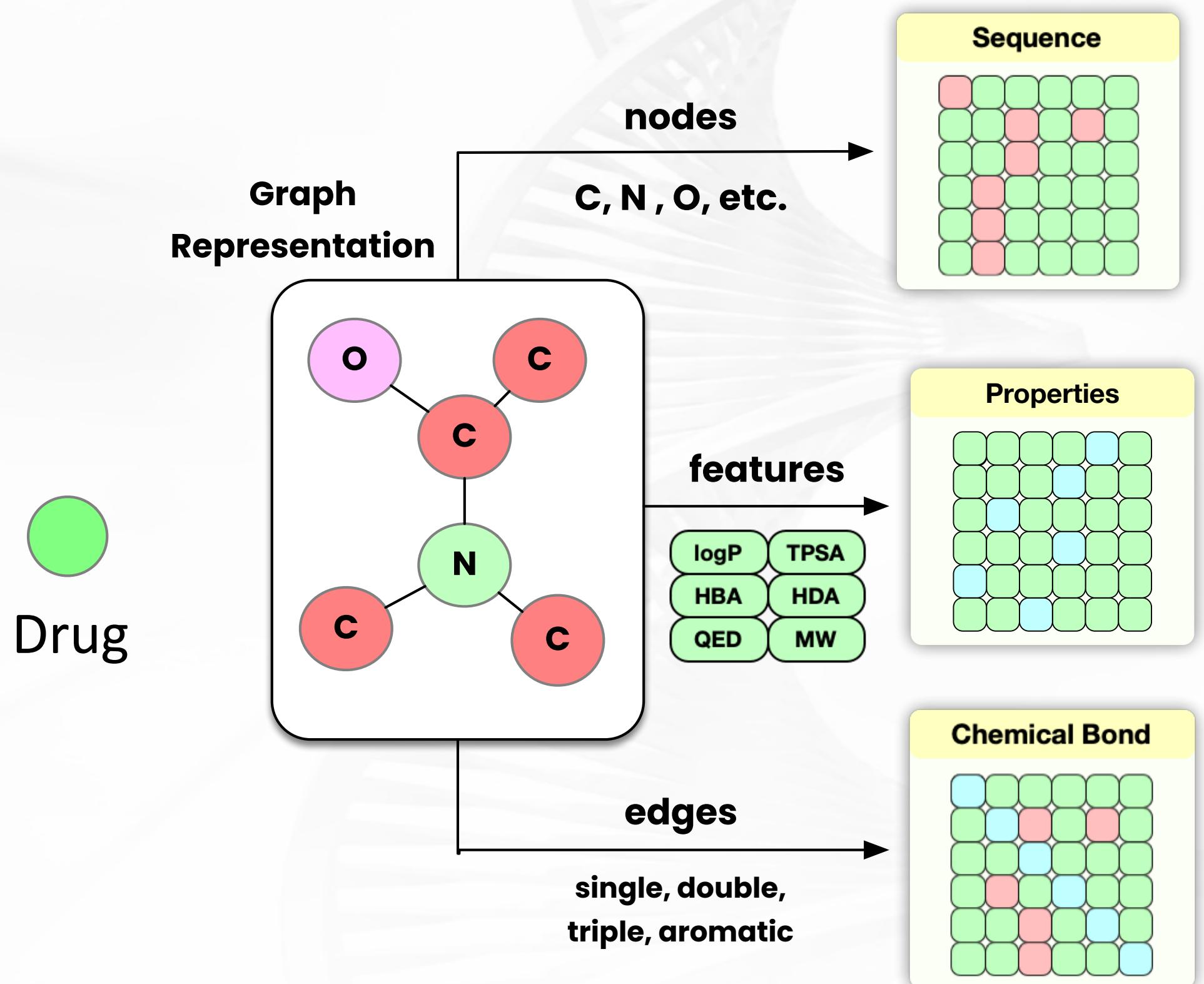
Graph Representation

■ Systematic Representation



Target Identification
Drug Repositioning

■ Individual Representation



Drug Property Prediction
DTI, Drug Design

Graphs in Drug Discovery

Data Scientist Perspective

Early Drug Discovery

TARGET IDENTIFICATION

- Drug-Target Identification
- Target Validation
- Protein and pathway screening in-vitro



HIT IDENTIFICATION

- DRUG-TARGET INTERACTION
- QSAR
- Compound Screening



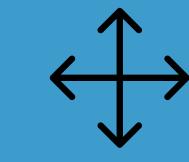
REPURPOSE EXISTING DRUGS

- Use Existing Drugs
- New Indication for Drugs



HIT TO LEAD

- Lead Identification
- Lead Optimization



Input

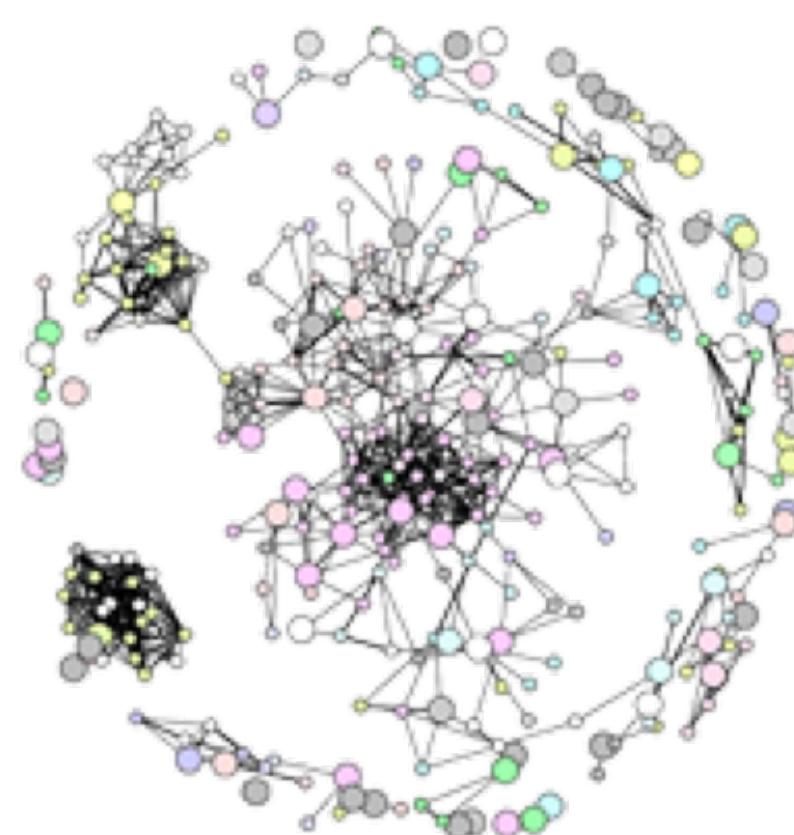
DISEASE, TARGET(Omics, Pathways, PPI)

COMPOUND, TARGET

DRUG, TARGET RELATIONSHIP

COMPOUND, TARGET

GRAPH Representation and Graph Learning



INDIVIDUAL REPRESENTATION
SYSTEMATIC REPRESENTATION

Output

NEW TARGET

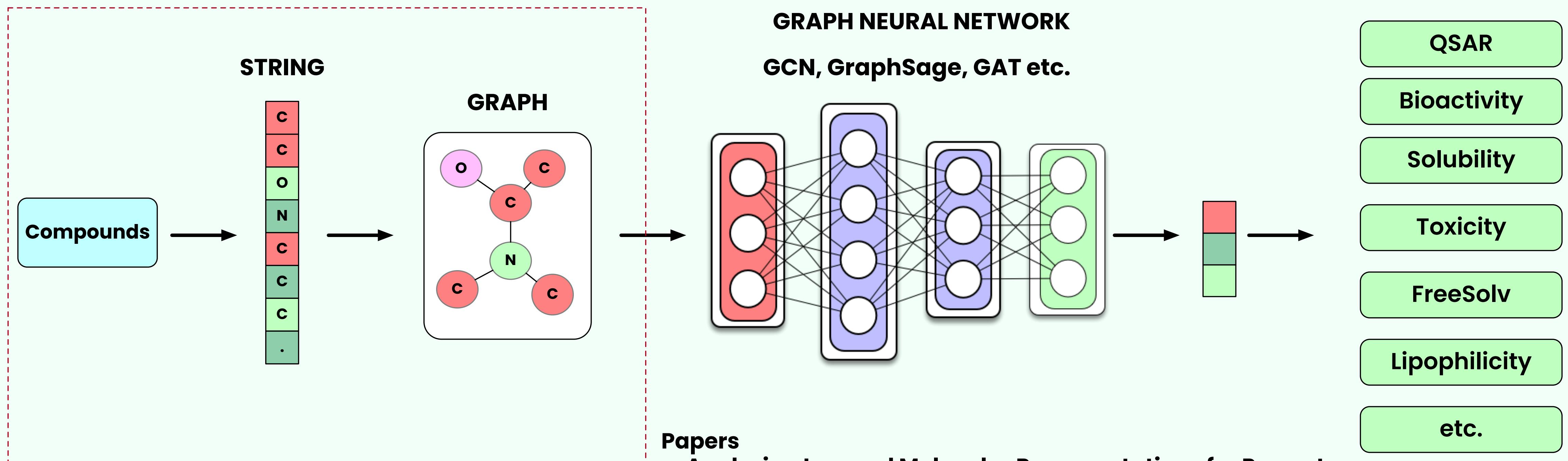
COMPOUND (HITs)

New Indications,
Drug Candidates

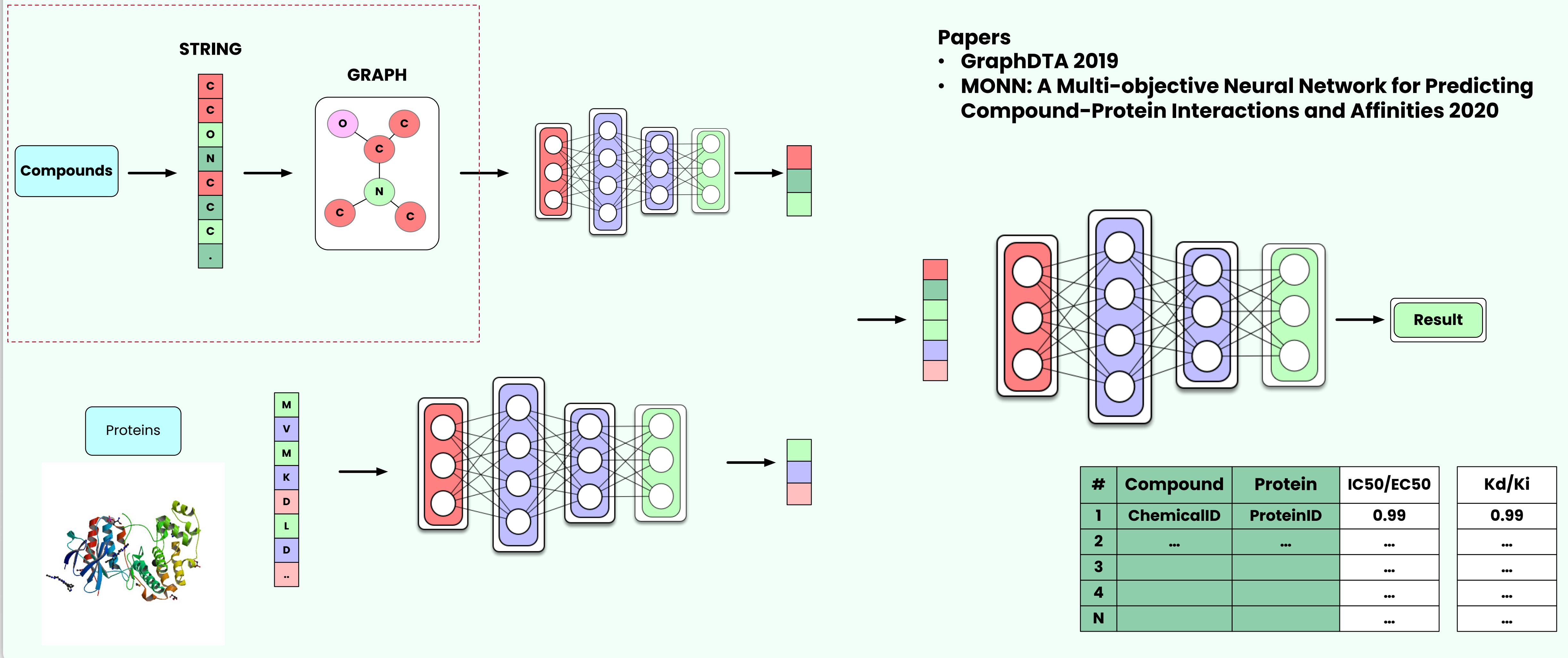
Drug Candidates (LEADS)

DRUG PROPERTY PREDICTION

Data Augmentation

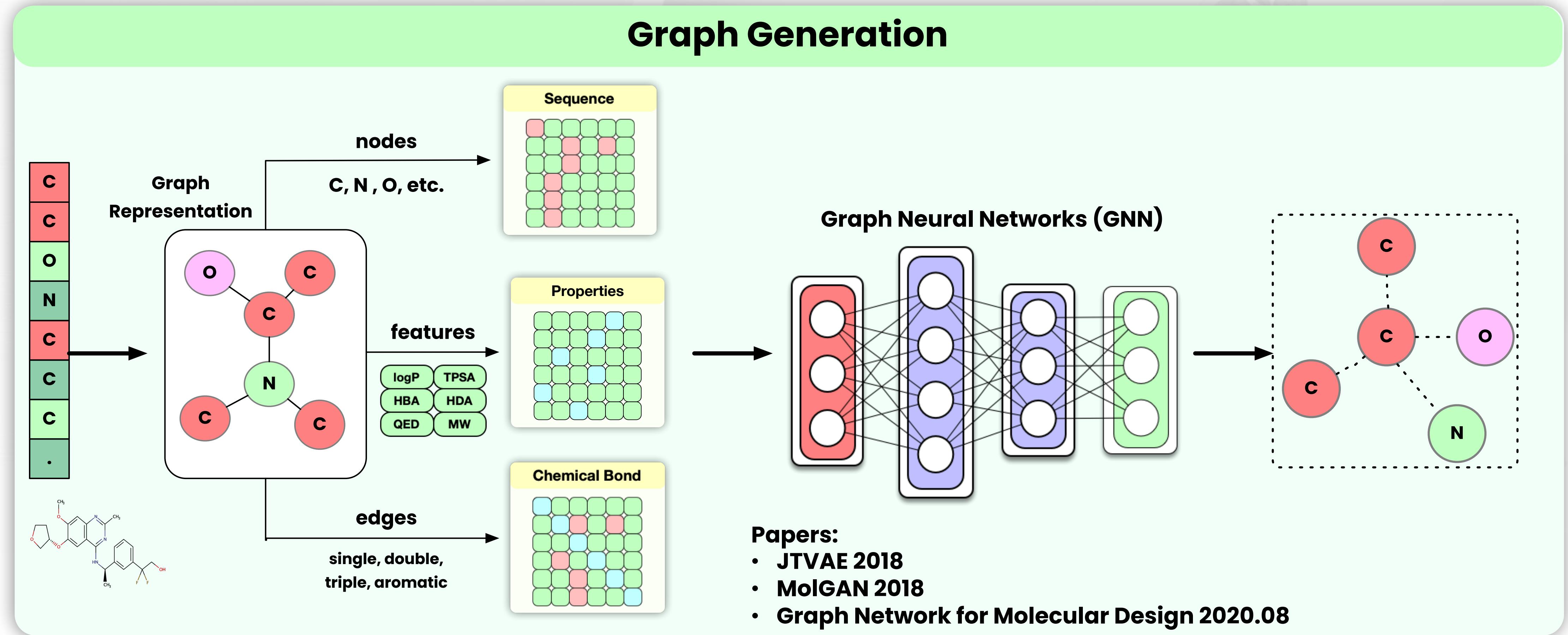


DRUG-TARGET INTERACTION PREDICTION

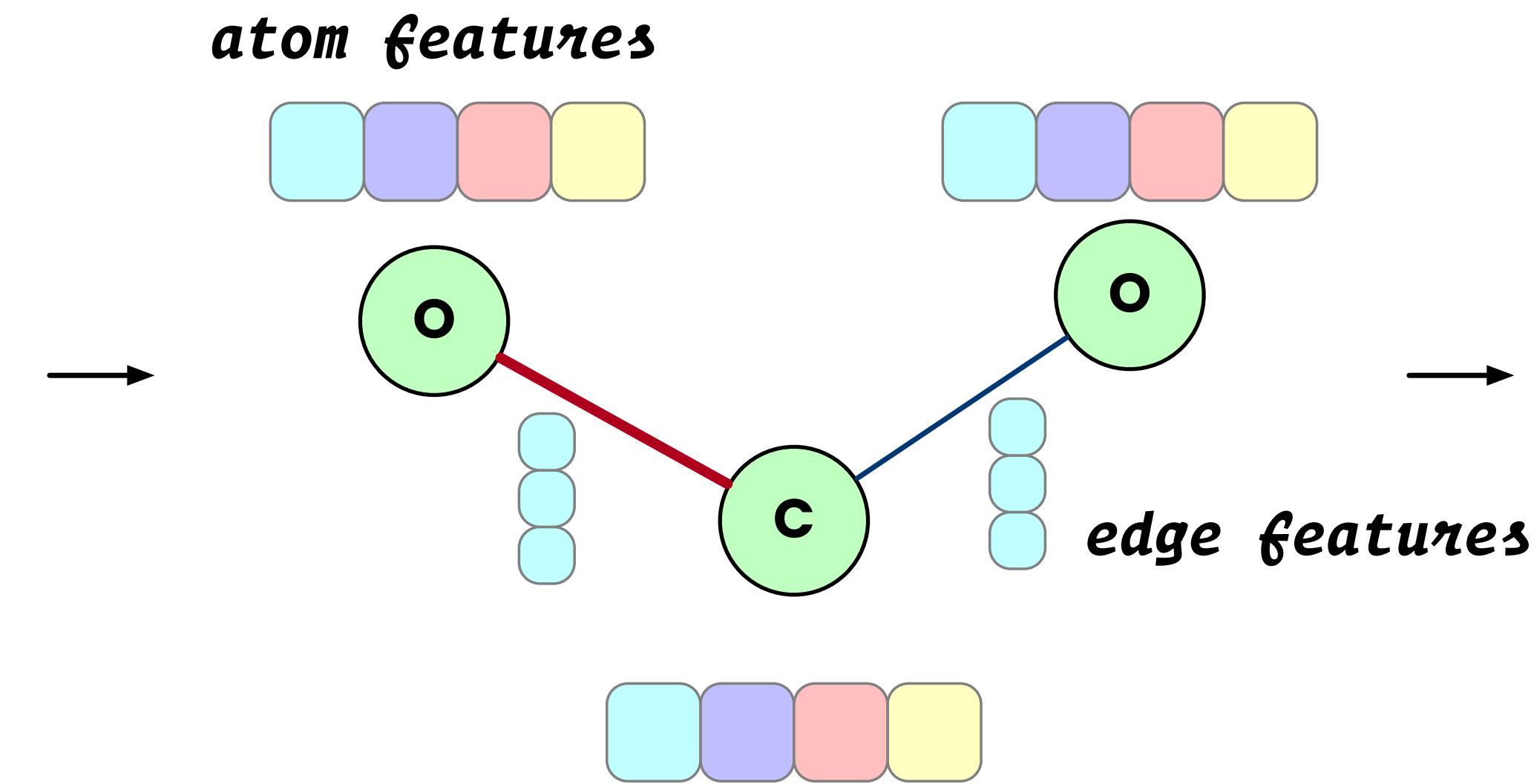
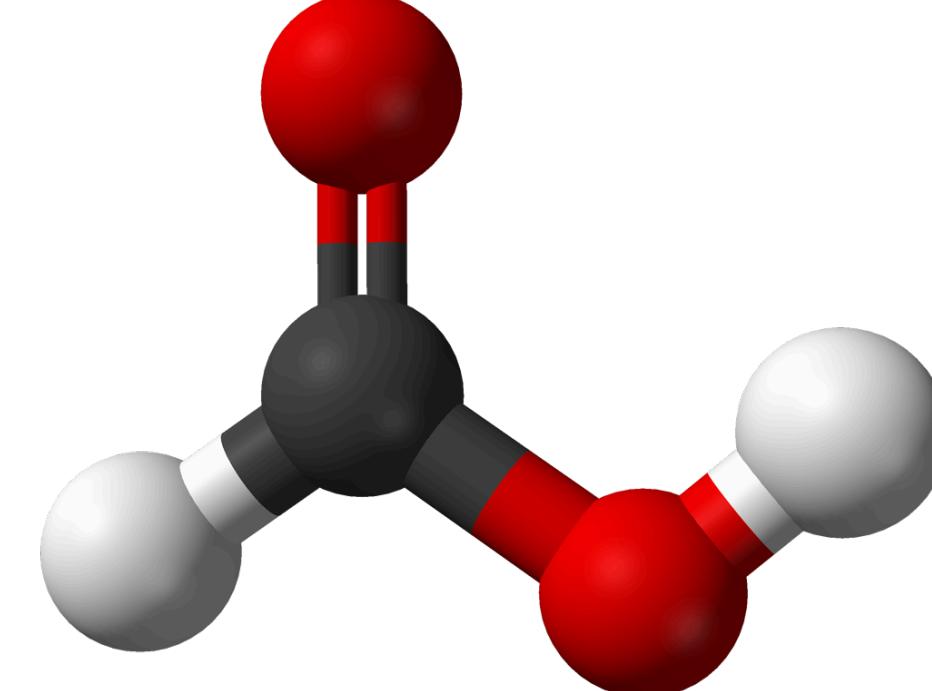


Graph Generation

- Graph Neural Networks
 - Input: Graph; Output: Graph



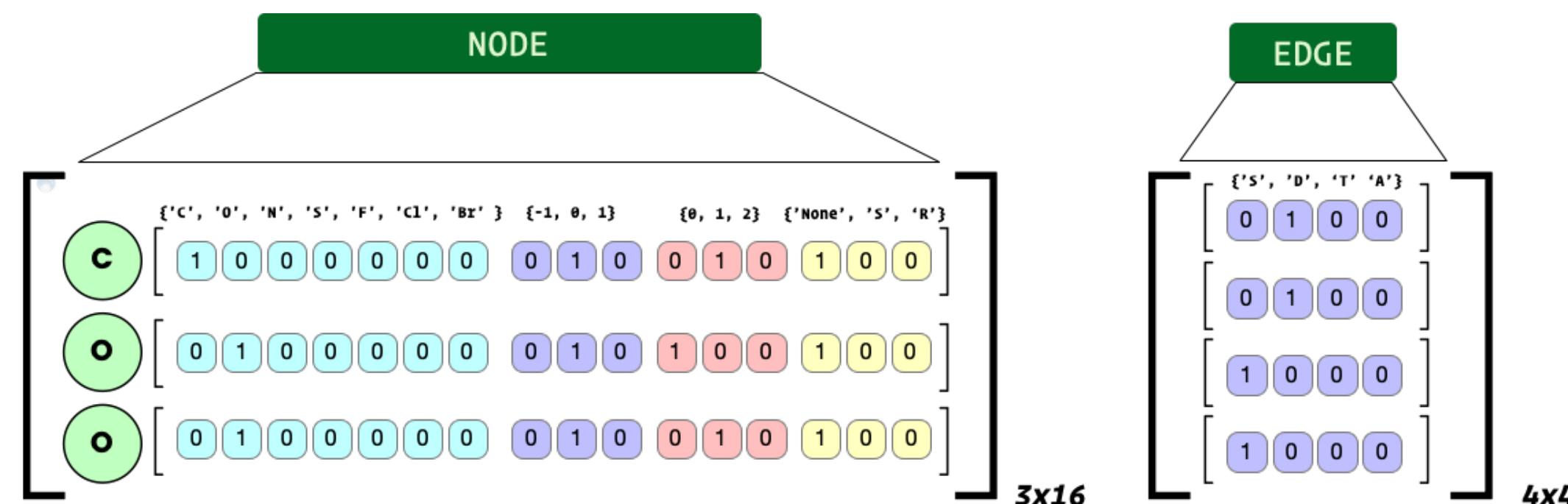
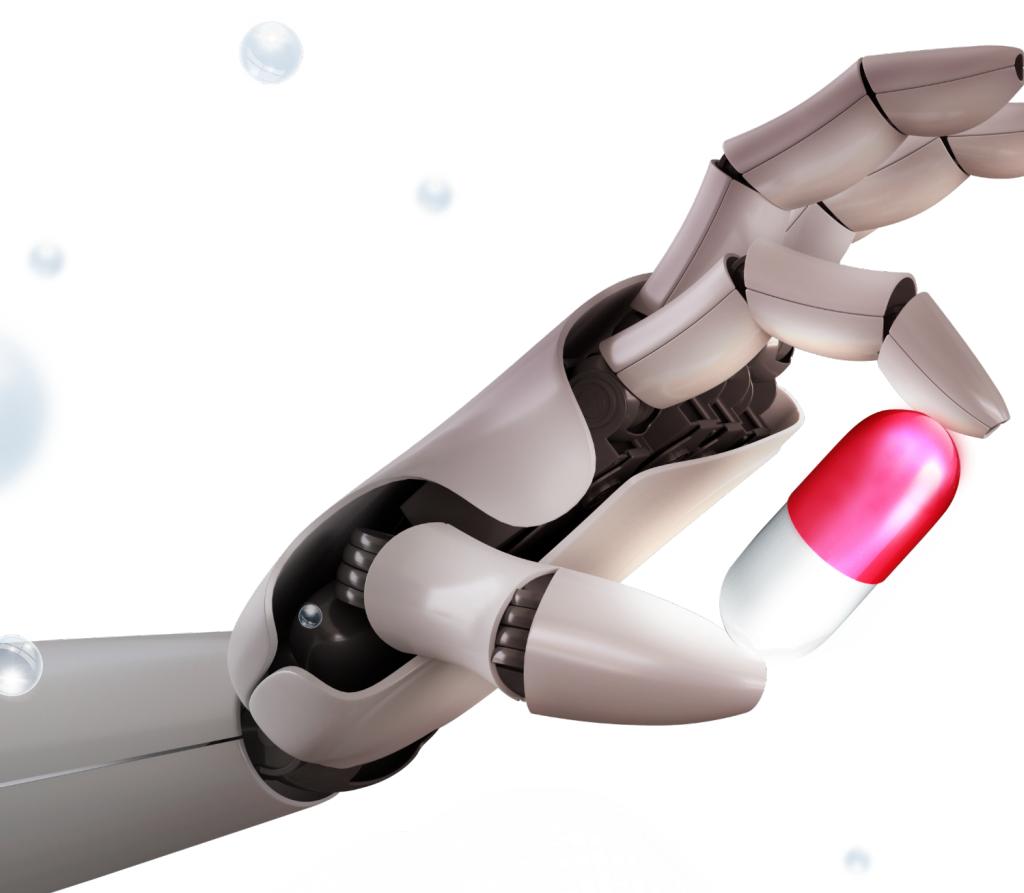
How Graph is Used In Chemical Space ?



PROPERTY PREDICTION

DRUG-TARGET
INTERACTION

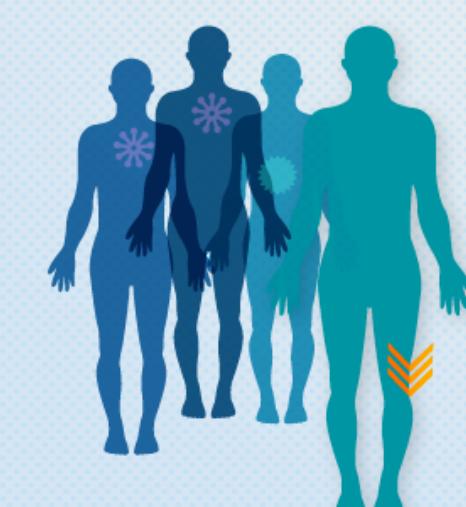
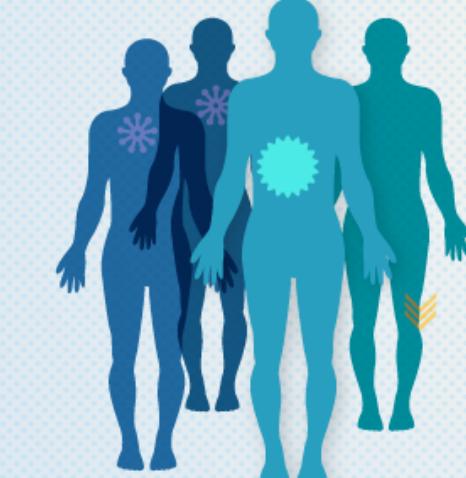
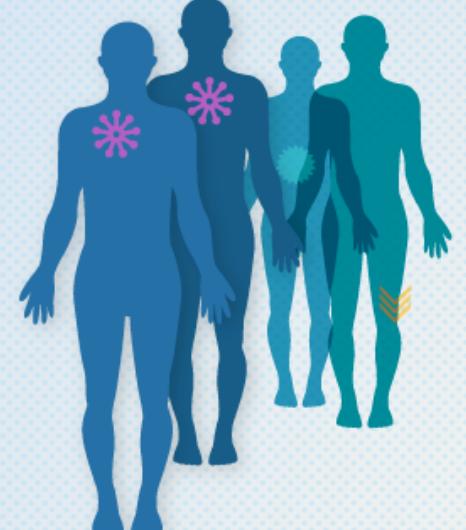
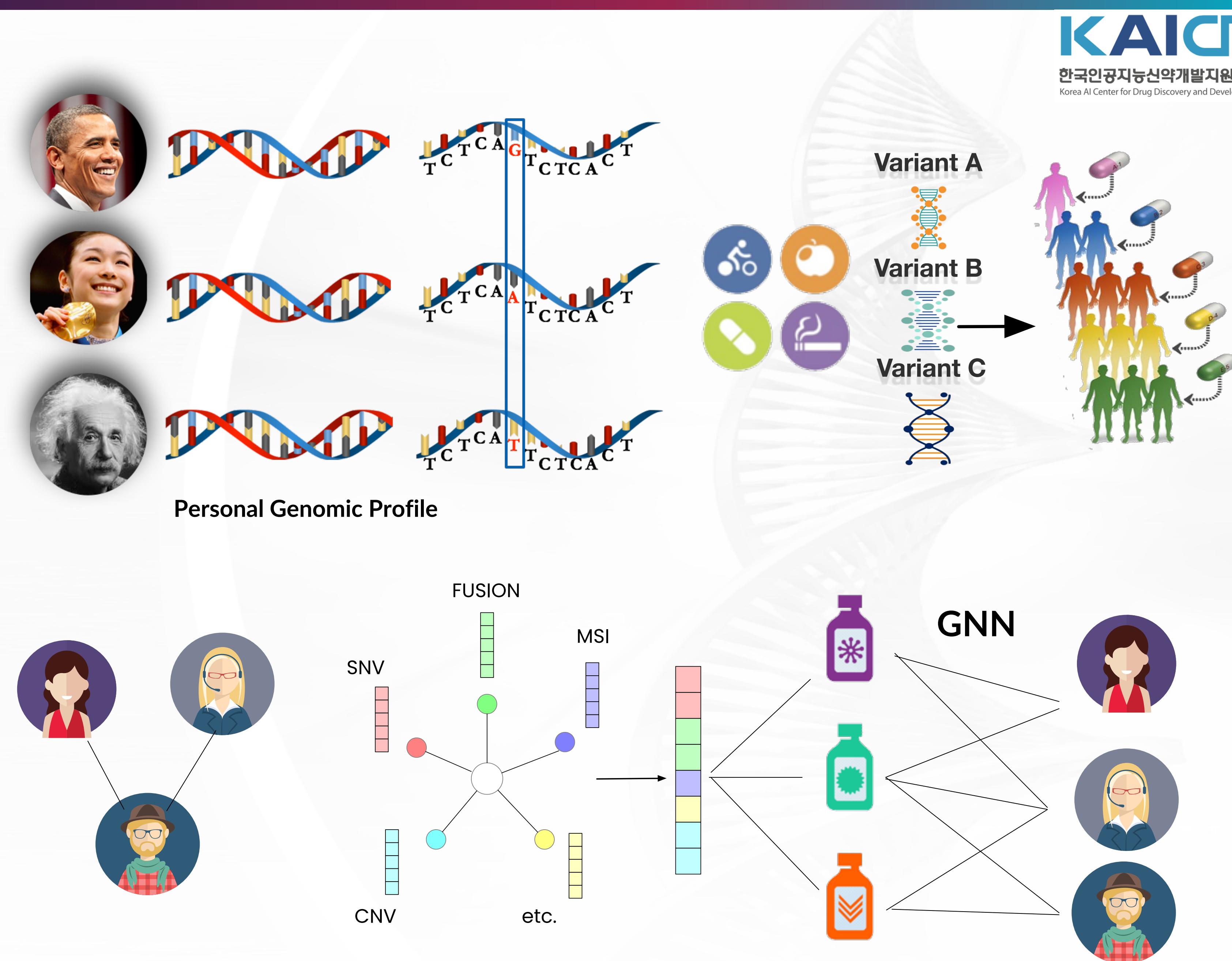
DE NOVO DRUG DESIGN



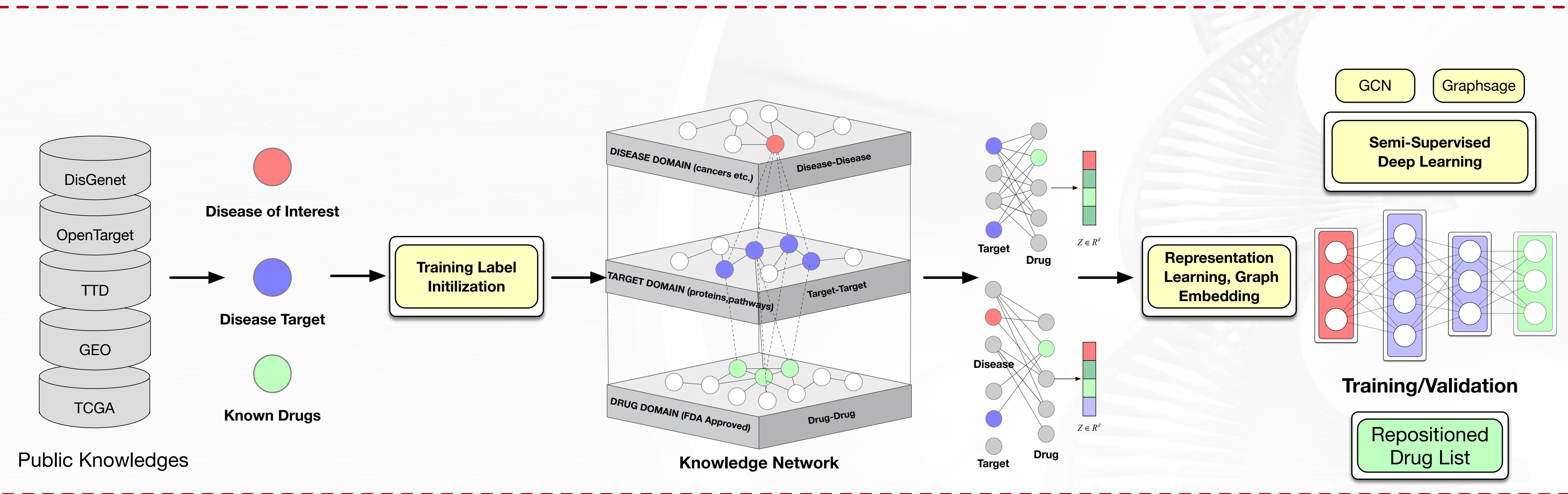
CODING SESSION: How to Make Graph Augmentation

NATIONAL CANCER INSTITUTE
PRECISION MEDICINE
IN CANCER TREATMENT

Discovering unique therapies that treat an individual's cancer based on the specific genetic abnormalities of that person's tumor.


www.cancer.gov


Drug Repositioning



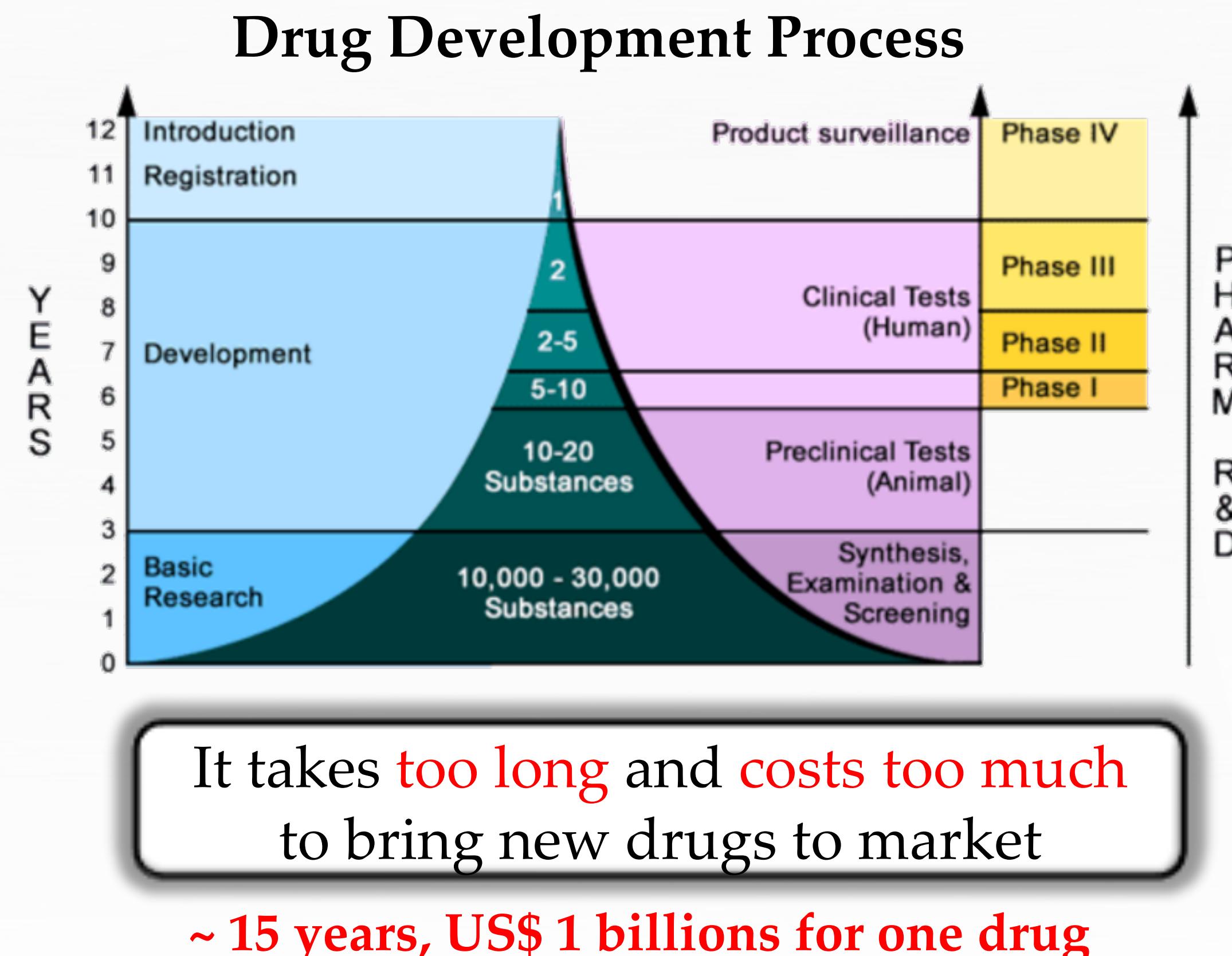
DRUG REPOSITIONING

“The most fruitful basis for the discovery of a new drug is to start with an old drug”

Nobel Prize-winning pharmacologist Sir James Black

Drug Repositioning

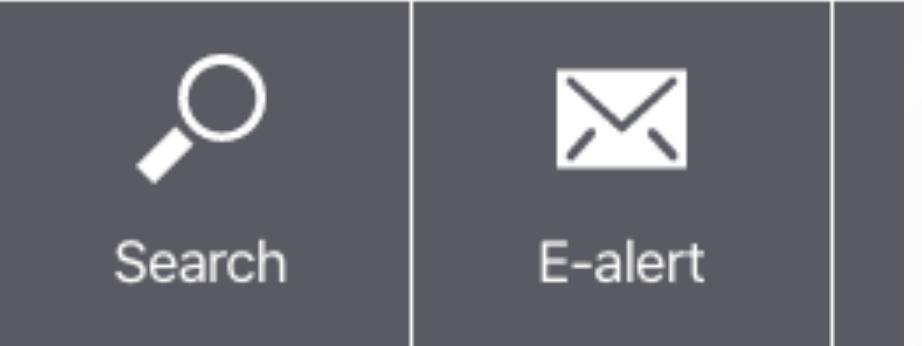
- A method for identifying and discovering new uses for existing drugs



New Use

Let's beef up efforts to screen existing drugs for new uses

nature
New uses for old drugs 2007
Curtis R. Chong¹ & David J. Sullivan, Jr²



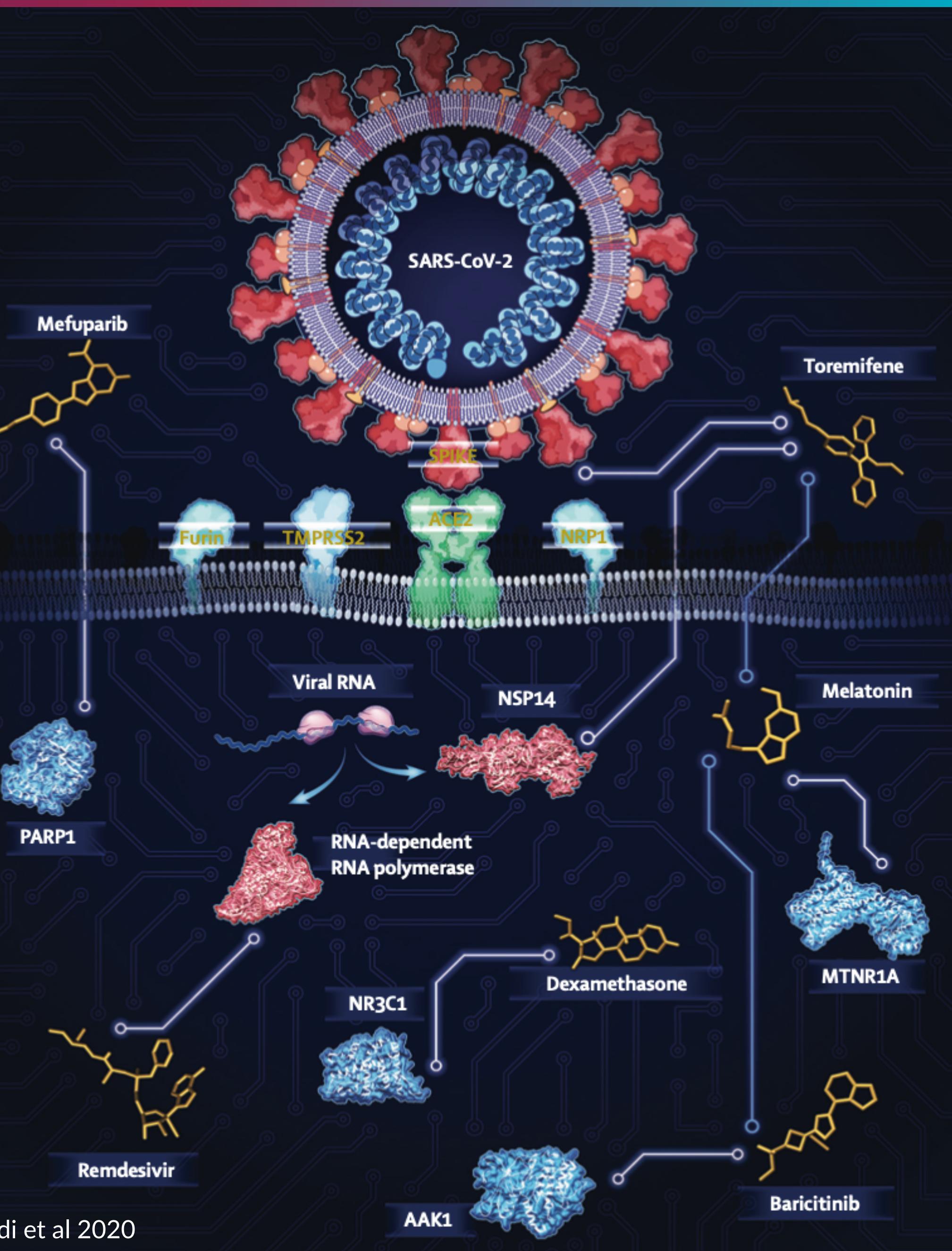
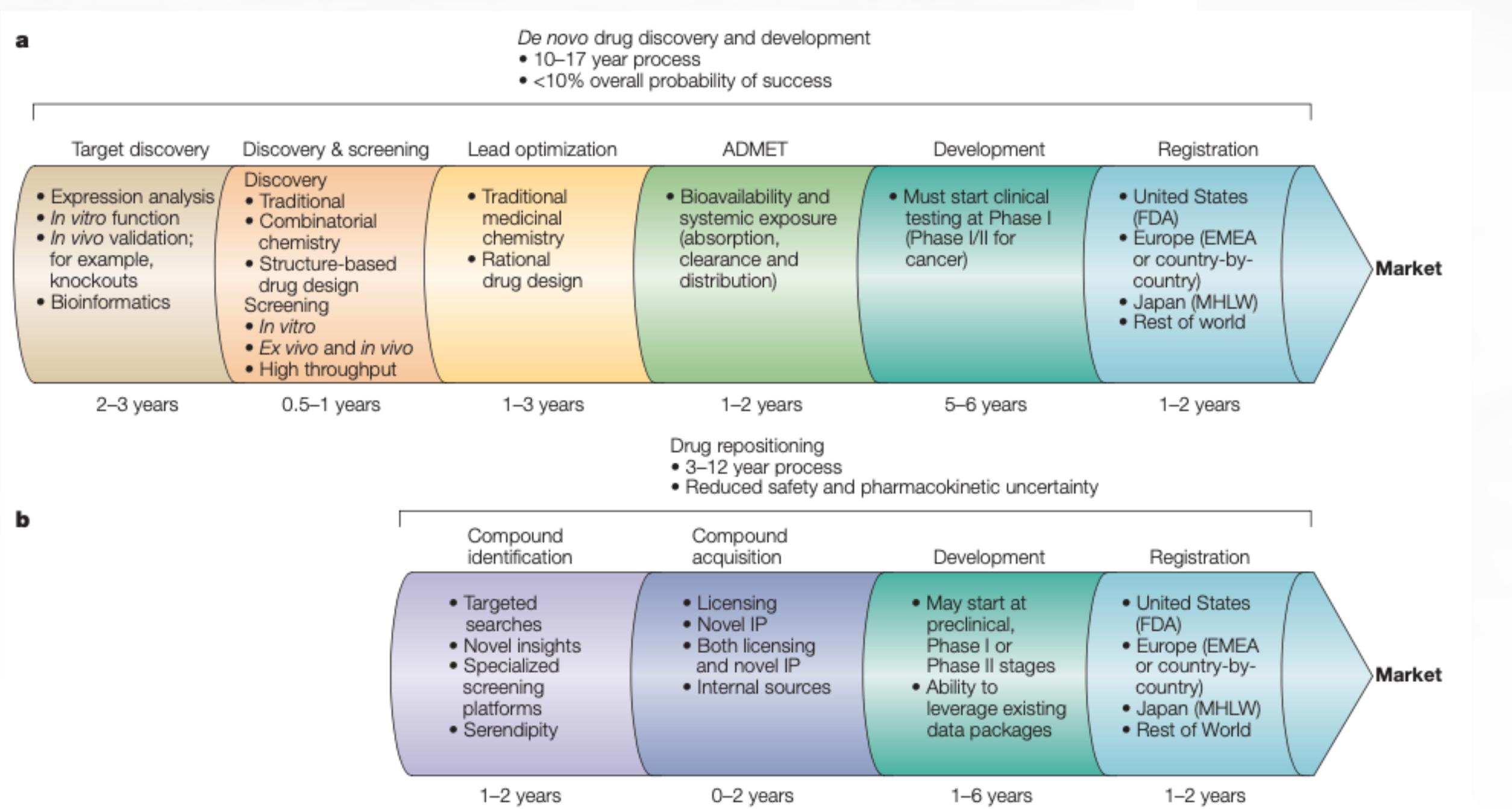
NEWS · 27 FEBRUARY 2020

Coronavirus puts drug repurposing on the fast track

Existing antivirals and knowledge gained from the SARS and MERS outbreaks gain traction as the fastest route to fight the current coronavirus epidemic.

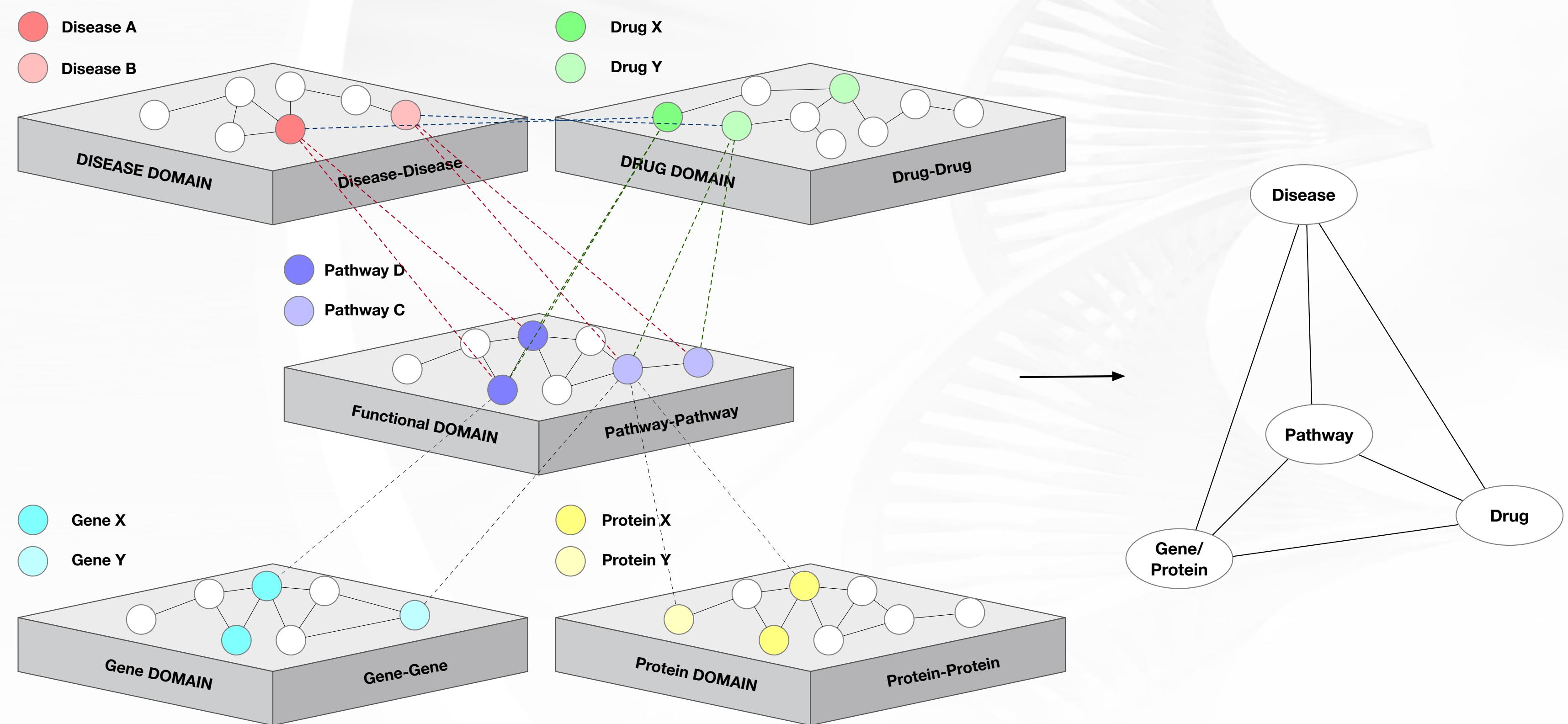
Advantages

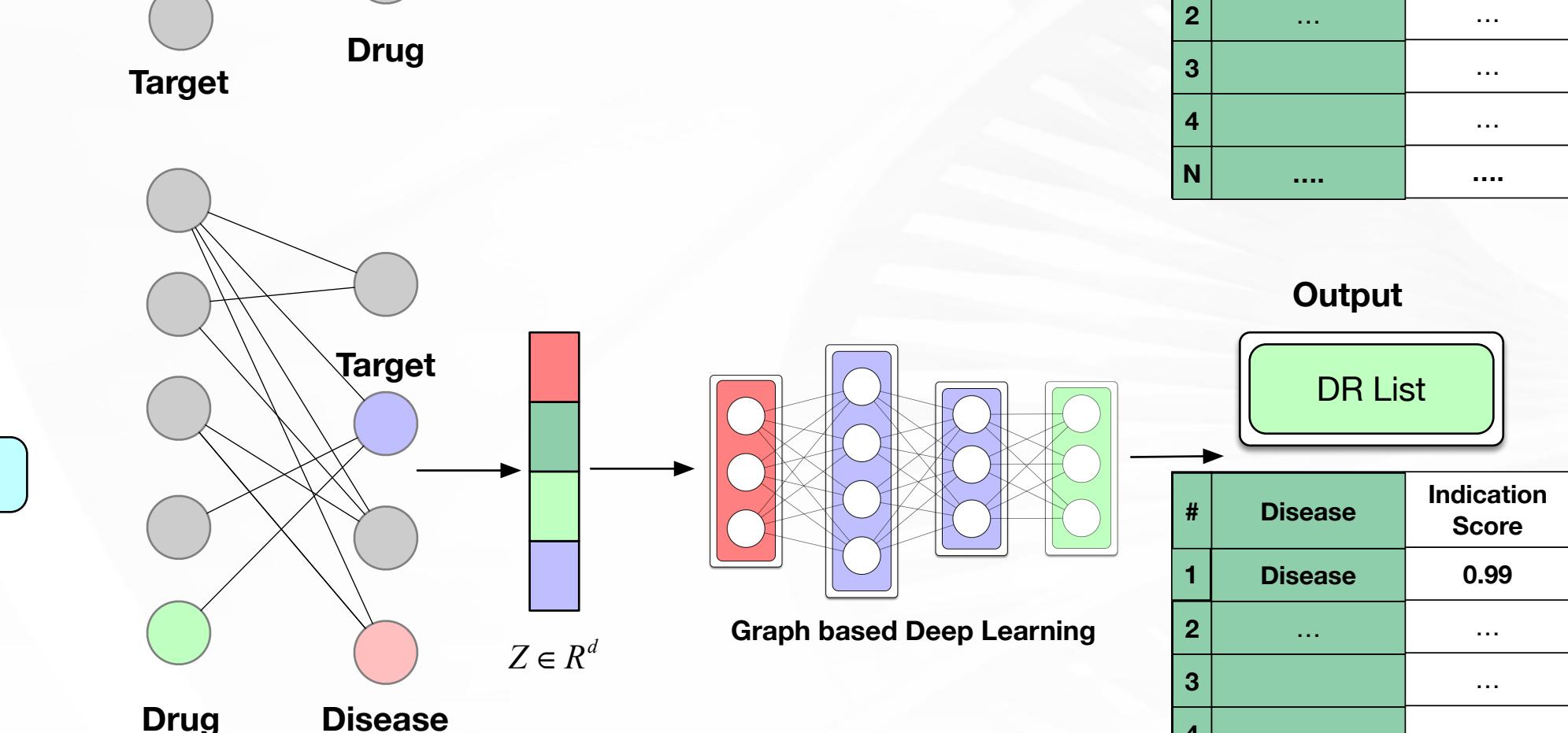
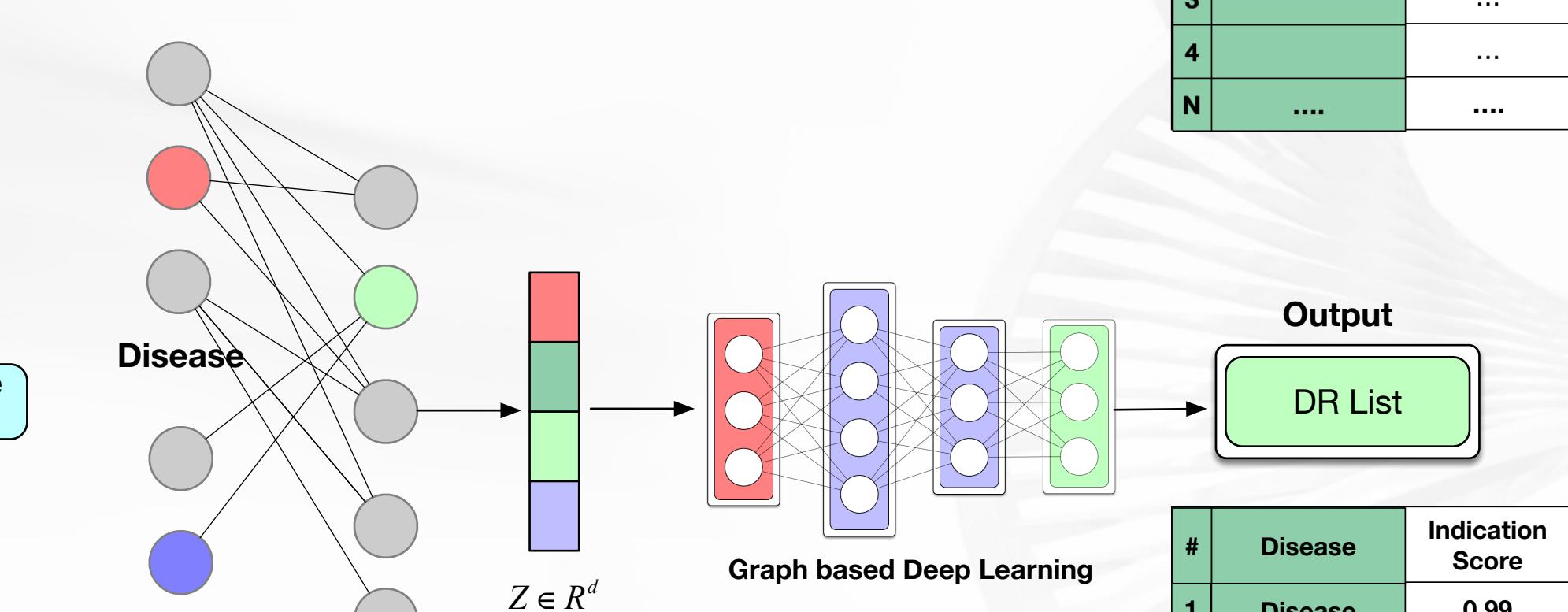
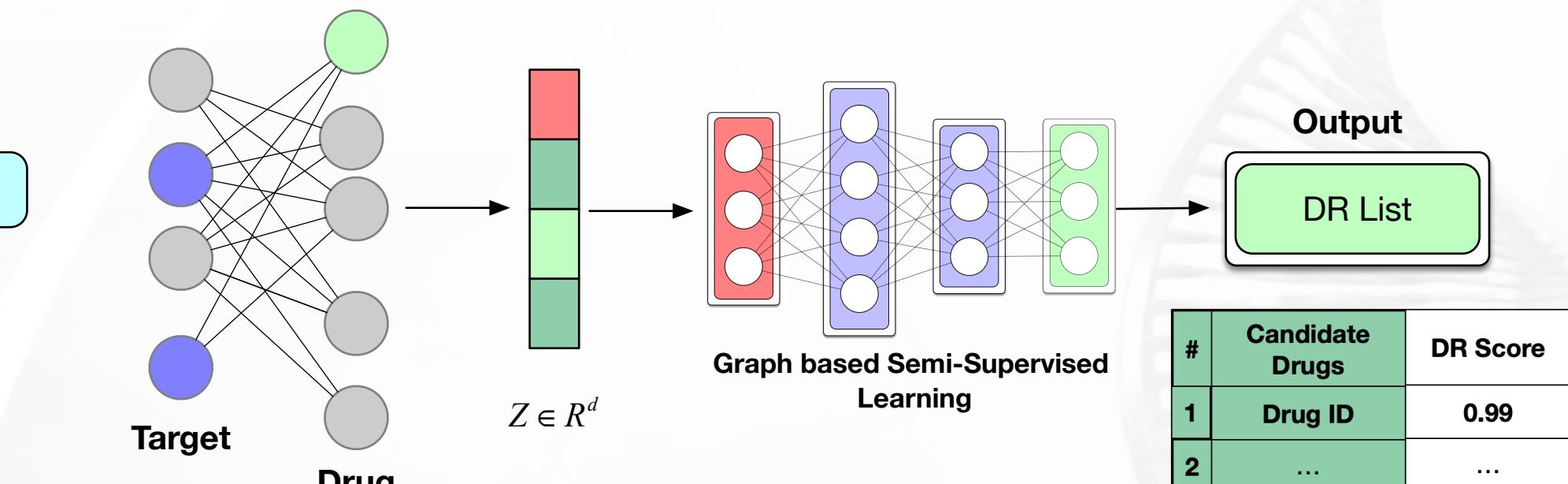
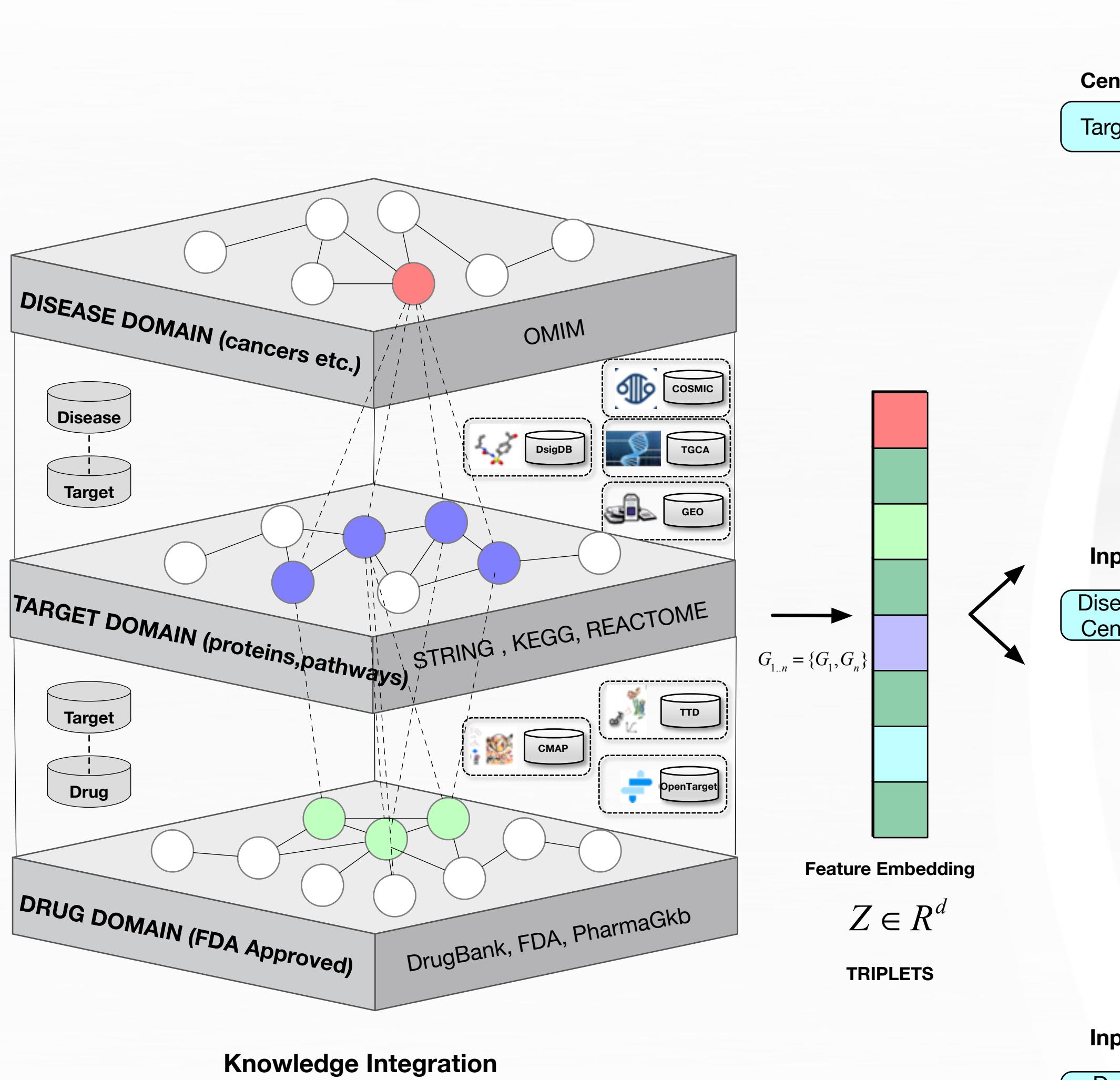
- Opportunity:
 - Reduced development timelines and overall cost
 - **Fast response to pandemic/epidemic (e.g: COVID19,)**



Traditional Method

- A classic way to repurpose drugs is through network medicine, which includes the construction of medical knowledge graphs
 - Target-centric
 - Disease-centric
 - Drug-centric





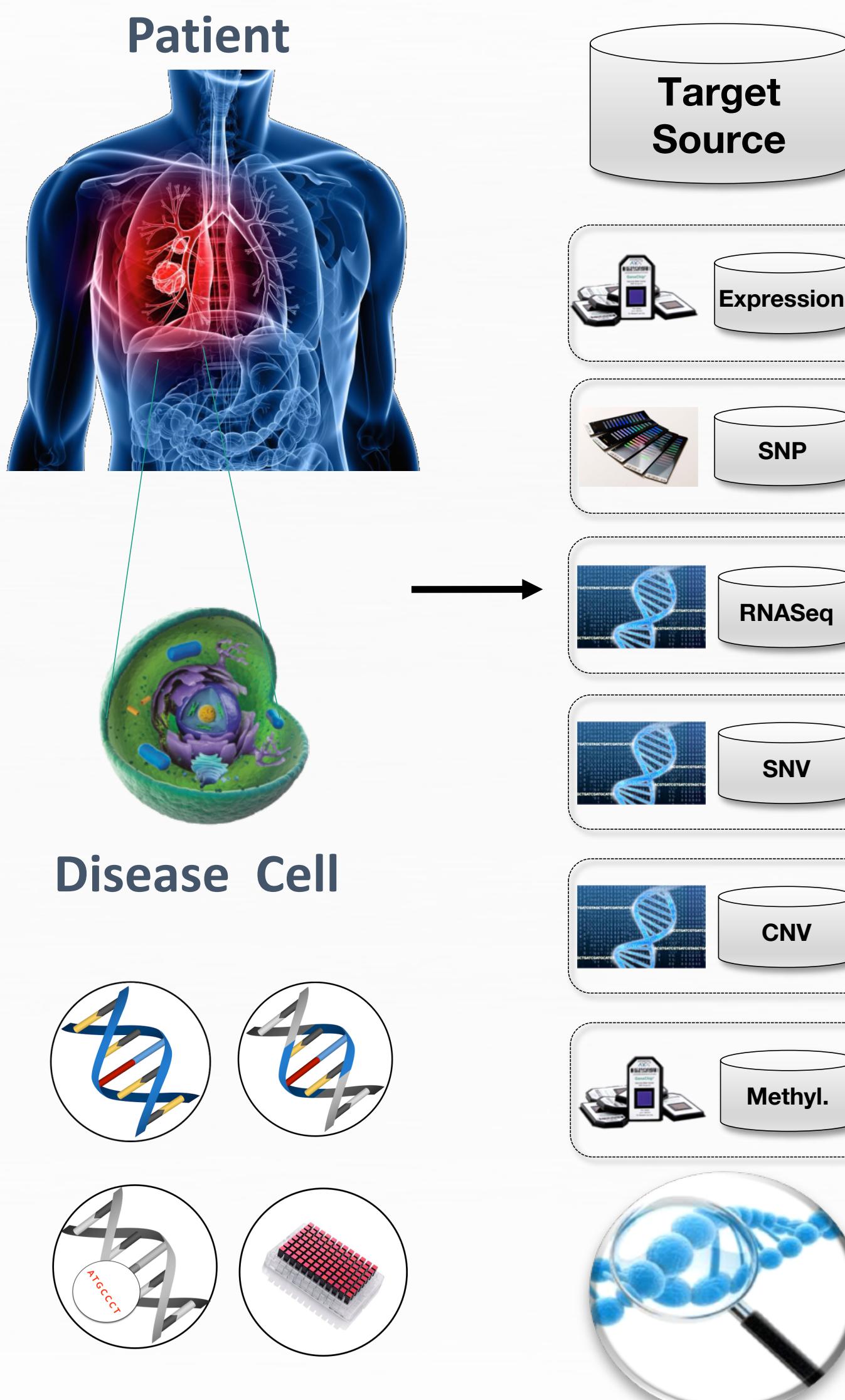
#	Candidate Drugs	DR Score
1	Drug ID	0.99
2
3
4
N

#	Disease	Indication Score
1	Disease	0.99
2
3
4
N

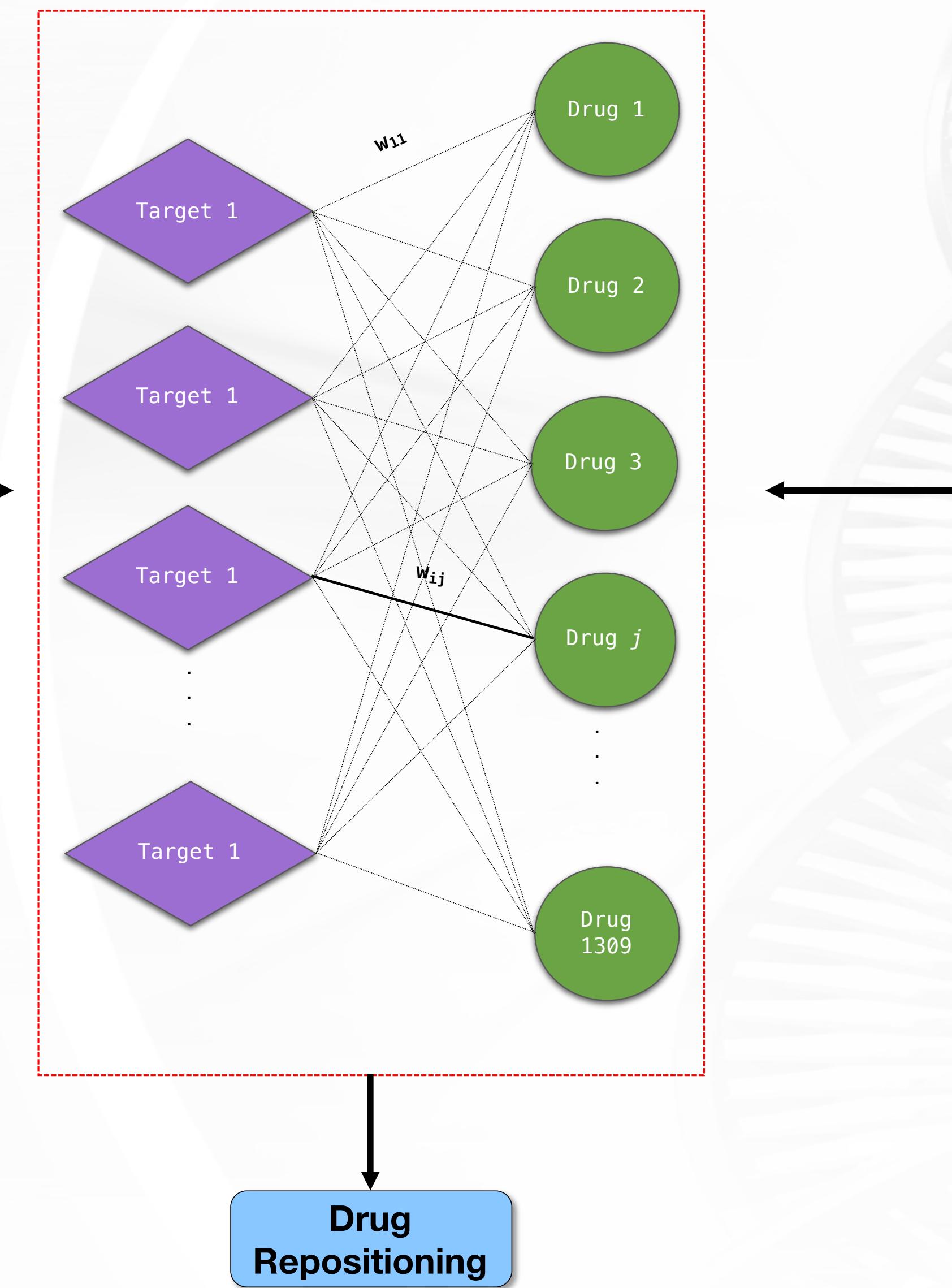
#	Disease	Indication Score
1	Disease	0.99
2
3
4
N

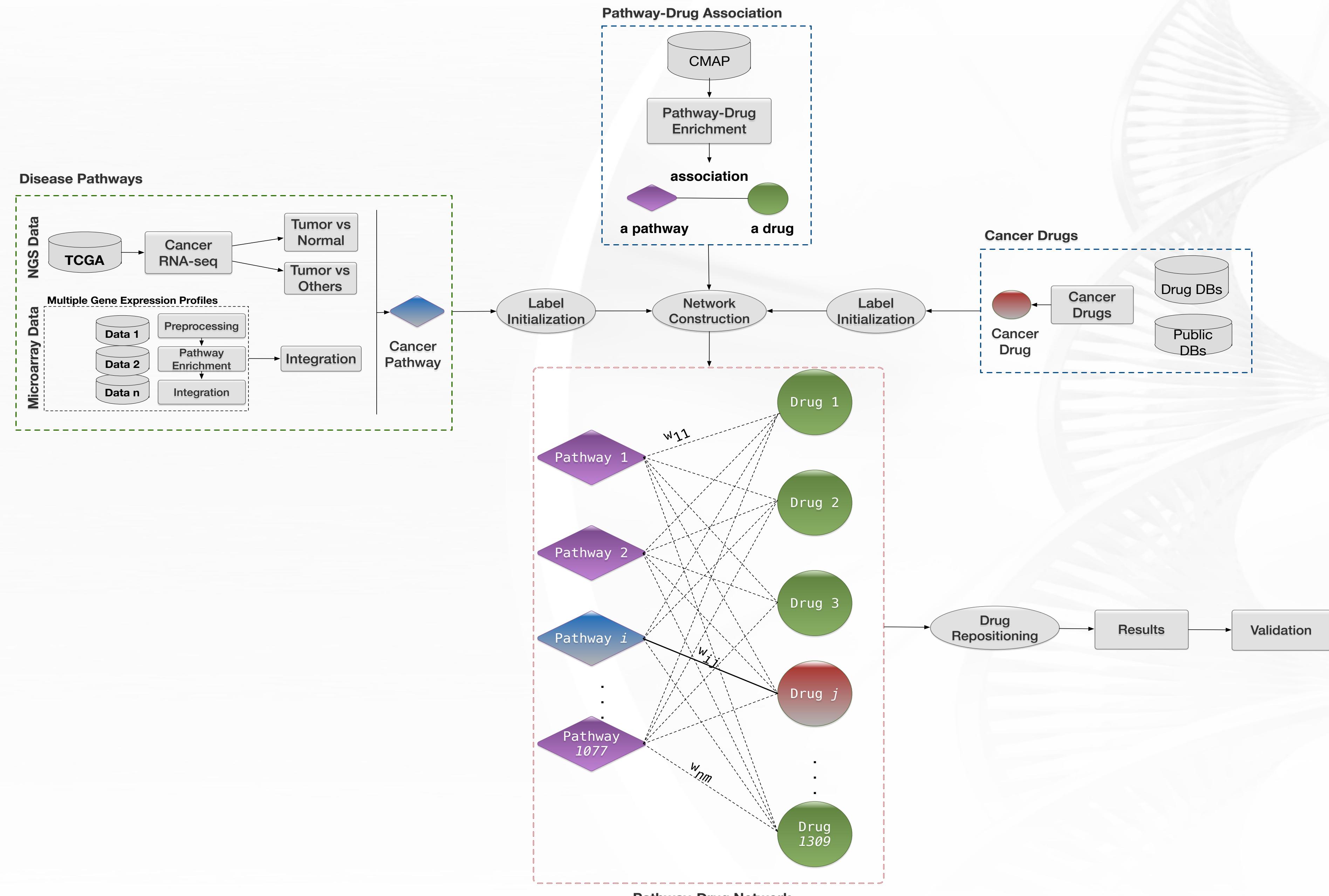
#	Drug	Indication Score
1	Drug ID	0.99
2
3
4
N

#	Drug	Indication Score
1	Drug ID	0.99
2
3
4
N



EMBEDDED TRIPLET

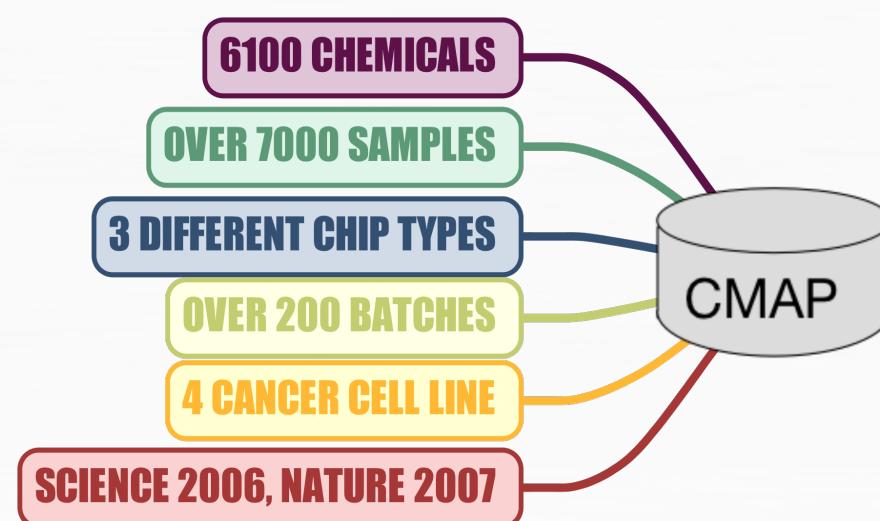




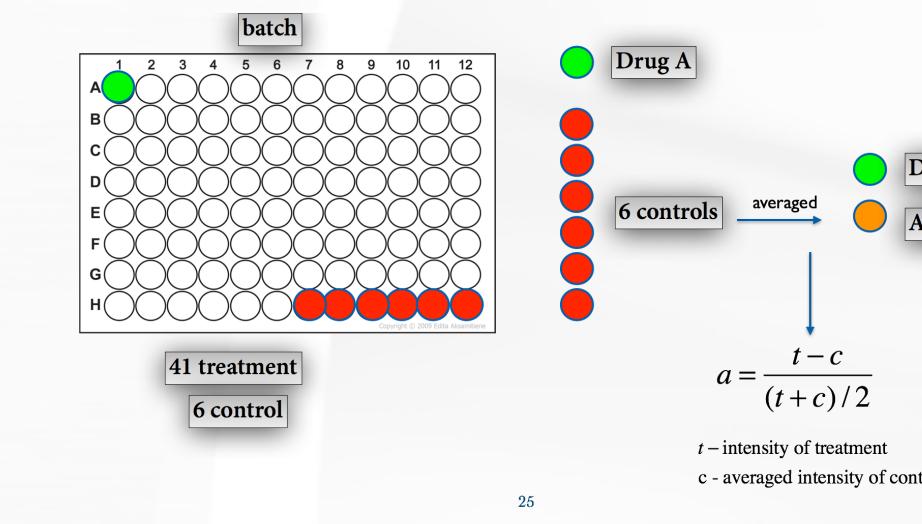
Drug Induced Omics

- CMAP was mostly used for DR, but now LINCS

	HT_HG-U133A	HG-U133A	HT_HG-U133A_EA	Total
Treated	5242	674	184	6100
Control	787	133	36	956
Total	6029	807	220	7056



- ✓ HT-HG-U133A platform
- ✓ 106 batches
- ✓ 3906 samples
- ✓ 1309 drugs



Gene Expression Matrix

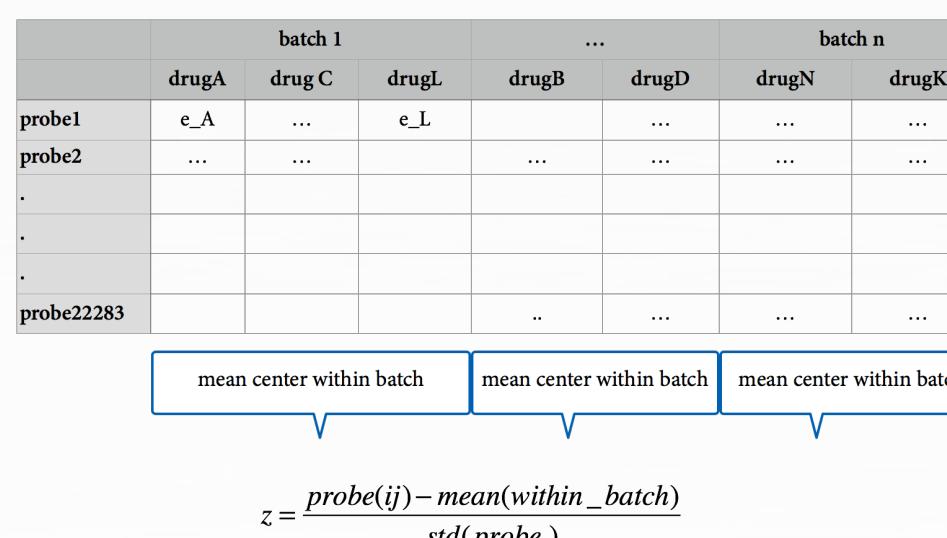
- ✓ Each cel file is are normalised by RMA
- ✓ All probe-sets less than 64 are set to 64
- ✓ Log2 transformed
- ✓ All controls for same batches averaged
- ✓ Batch effect removed by mean centering

- ✓ 222277 probes converted 19930 gene symbol

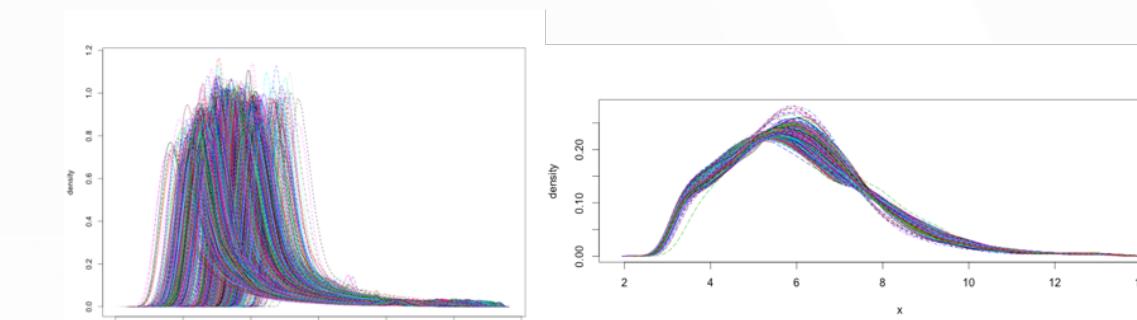
- ✓ D-score enrichment
- ✓ GSEA

Drug-Pathway Enrichment Matrix

- ✓ Each Pathway enrichment scores for each drug



Batch effect removal



RMA preprocessing

	Drug 1	Drug 2	Drug 3	Drug 1309
Pathway 1	Enrichment score				
Pathway 2					
Pathway 3					
...					
Pathway 1077					

Drug-Pathway Matrix

Zhou et al. *Cell Discovery* (2020)6:14
<https://doi.org/10.1038/s41421-020-0153-3>

Cell Discovery
www.nature.com/celldisc

ARTICLE

Open Access

Network-based drug repurposing for novel coronavirus 2019-nCoV/SARS-CoV-2

Yadi Zhou¹, Yuan Hou¹, Jiayu Shen¹, Yin Huang¹, William Martin¹ and Feixiong Cheng^{1,2,3}

Abstract

Human coronaviruses (HCoVs), including severe acute respiratory syndrome coronavirus (SARS-CoV) and 2019 novel coronavirus (2019-nCoV, also known as SARS-CoV-2), lead global epidemics with high morbidity and mortality. However, there are currently no effective drugs targeting 2019-nCoV/SARS-CoV-2. Drug repurposing, representing as an effective drug discovery strategy from existing drugs, could shorten the time and reduce the cost compared to de novo drug discovery. In this study, we present an integrative, antiviral drug repurposing methodology implementing a systems pharmacology-based network medicine platform, quantifying the interplay between the HCoV-host interactome and drug targets in the human protein–protein interaction network. Phylogenetic analyses of 15 HCoV whole genomes reveal that 2019-nCoV/SARS-CoV-2 shares the highest nucleotide sequence identity with SARS-CoV (79.7%). Specifically, the envelope and nucleocapsid proteins of 2019-nCoV/SARS-CoV-2 are two evolutionarily conserved regions, having the sequence identities of 96% and 89.6%, respectively, compared to SARS-CoV. Using network proximity analyses of drug targets and HCoV-host interactions in the human interactome, we prioritize 16 potential anti-HCoV repurposable drugs (e.g., melatonin, mercaptopurine, and sirolimus) that are further validated by enrichment analyses of drug-gene signatures and HCoV-induced transcriptomics data in human cell lines. We further identify three potential drug combinations (e.g., sirolimus plus dactinomycin, mercaptopurine plus melatonin, and toremifene plus emodin) captured by the “Complementary Exposure” pattern: the targets of the drugs both hit the HCoV-host subnetwork but target separate neighborhoods in the human interactome network. In summary, this study offers powerful network-based methodologies for rapid identification of candidate repurposable drugs and potential drug combinations targeting 2019-nCoV/SARS-CoV-2.

Introduction

Coronaviruses (CoVs) typically affect the respiratory tract of mammals, including humans, and lead to mild to severe respiratory tract infections¹. In the past two decades, two highly pathogenic human CoVs (HCoVs), including severe acute respiratory syndrome coronavirus (SARS-CoV) and Middle East respiratory syndrome coronavirus (MERS-CoV), emerging from animal reservoirs, have led to global epidemics with high morbidity and

mortality². For example, 8098 individuals were infected and 774 died in the SARS-CoV pandemic, which cost the global economy with an estimated \$30 to \$100 billion^{3,4}. According to the World Health Organization (WHO), as of November 2019, MERS-CoV has had a total of 2494 diagnosed cases causing 858 deaths, the majority in Saudi Arabia². In December 2019, the third pathogenic HCoV, named 2019 novel coronavirus (2019-nCoV/SARS-CoV-2), as the cause of coronavirus disease 2019 (abbreviated as COVID-19)⁵, was found in Wuhan, China. As of 24 February 2020, there have been over 79,000 cases with over 2600 deaths for the 2019-nCoV/SARS-CoV-2 outbreak worldwide; furthermore, human-to-human transmission has occurred among close contacts⁶. However, there are currently no effective medications against 2019-

Correspondence: Feixiong Cheng (chengf@ccf.org)

¹Genomic Medicine Institute, Lerner Research Institute, Cleveland Clinic, Cleveland, OH 44195, USA

²Department of Molecular Medicine, Cleveland Clinic Lerner College of Medicine, Case Western Reserve University, Cleveland, OH 44195, USA

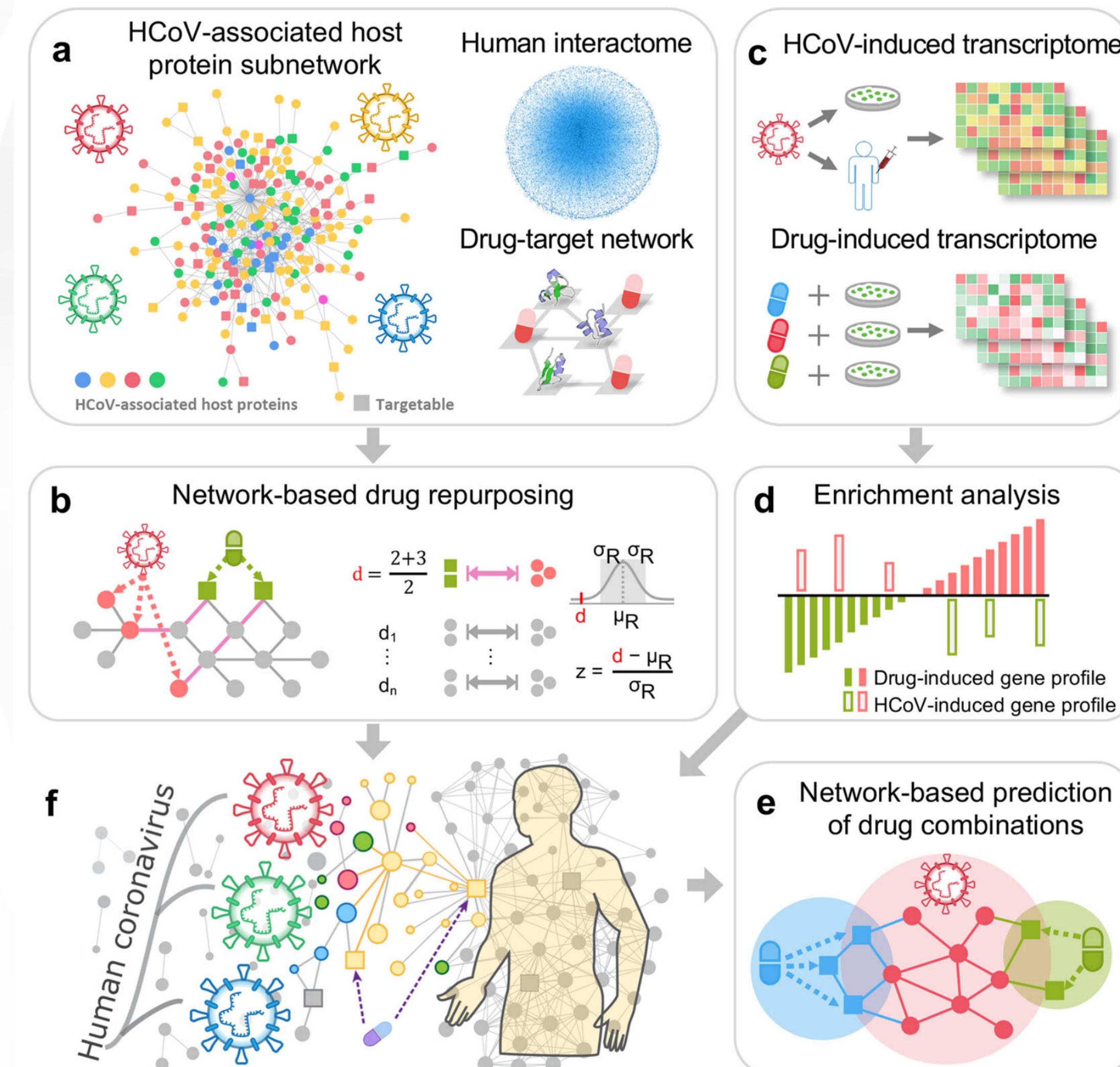
Full list of author information is available at the end of the article

These authors contributed equally: Yadi Zhou, Yuan Hou

© The Author(s) 2020

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

Nature / Cell Discovery, 16 March 2020



bioRxiv preprint doi: <https://doi.org/10.1101/2020.03.11.986836>. The copyright holder for this preprint (which was not peer-reviewed) is the author/funder. All rights reserved. No reuse allowed without permission.

A data-driven drug repositioning framework discovered a potential therapeutic agent targeting COVID-19

Yiyue Ge^{1,2,†}, Tingzhong Tian^{1,2,†}, Suling Huang^{3,†}, Fangping Wan^{1,†}, Jingxin Li^{2,†}, Shuya Li¹, Hui Yang¹¹, Lixiang Hong¹, Nian Wu¹, Enming Yuan¹, Lili Cheng⁴, Yipin Lei¹¹, Hantao Shu¹, Xiaolong Feng^{6,7}, Ziyuan Jiang⁵, Ying Chi², Xiling Guo², Lumbiao Cui², Liang Xiao¹⁰, Zeng Li¹⁰, Chunhao Yang³, Zehong Miao³, Haidong Tang⁴, Ligong Chen⁴, Haimian Zeng¹¹, Dan Zhao^{1,*}, Fengcai Zhu^{2,8,*}, Xiaokun Shen^{10,*}, Jianyang Zeng^{1,9,*}

¹Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing, 100084, China.

²NHC Key laboratory of Enteric Pathogenic Microbiology, Jiangsu Provincial Center for Diseases Control and Prevention, Nanjing, Jiangsu Province, 210009, China.

³Shanghai Institute of Materia Medica, Chinese Academy of Sciences, Shanghai, 201203, China.

⁴School of Pharmaceutical Sciences, Beijing Advanced Innovation Center for Structural Biology, Tsinghua University, Beijing, 100084, China.

⁵Department of Automation, Tsinghua University, Beijing, 100084, China.

⁶School of Electronic Information and Communications, Huazhong University of Science and Technology, Wuhan, Hubei Province, 430074, China.

⁷Institute of Pathology, Tongji Hospital, Tongji Medical College, Huazhong University of Science and Technology, Wuhan, Hubei Province, 430030, China.

⁸Center for Global Health, Nanjing Medical University, Nanjing, Jiangsu Province, 210009, China.

⁹MOE Key Laboratory of Bioinformatics, Tsinghua University, Beijing, 100084, China.

¹⁰Convalife (Shanghai) Co., Ltd., Shanghai, 201203, China.

¹¹Silexon AI Technology Co., Ltd., Nanjing, Jiangsu Province, 210033, China.

[†]These authors contributed equally to this work.

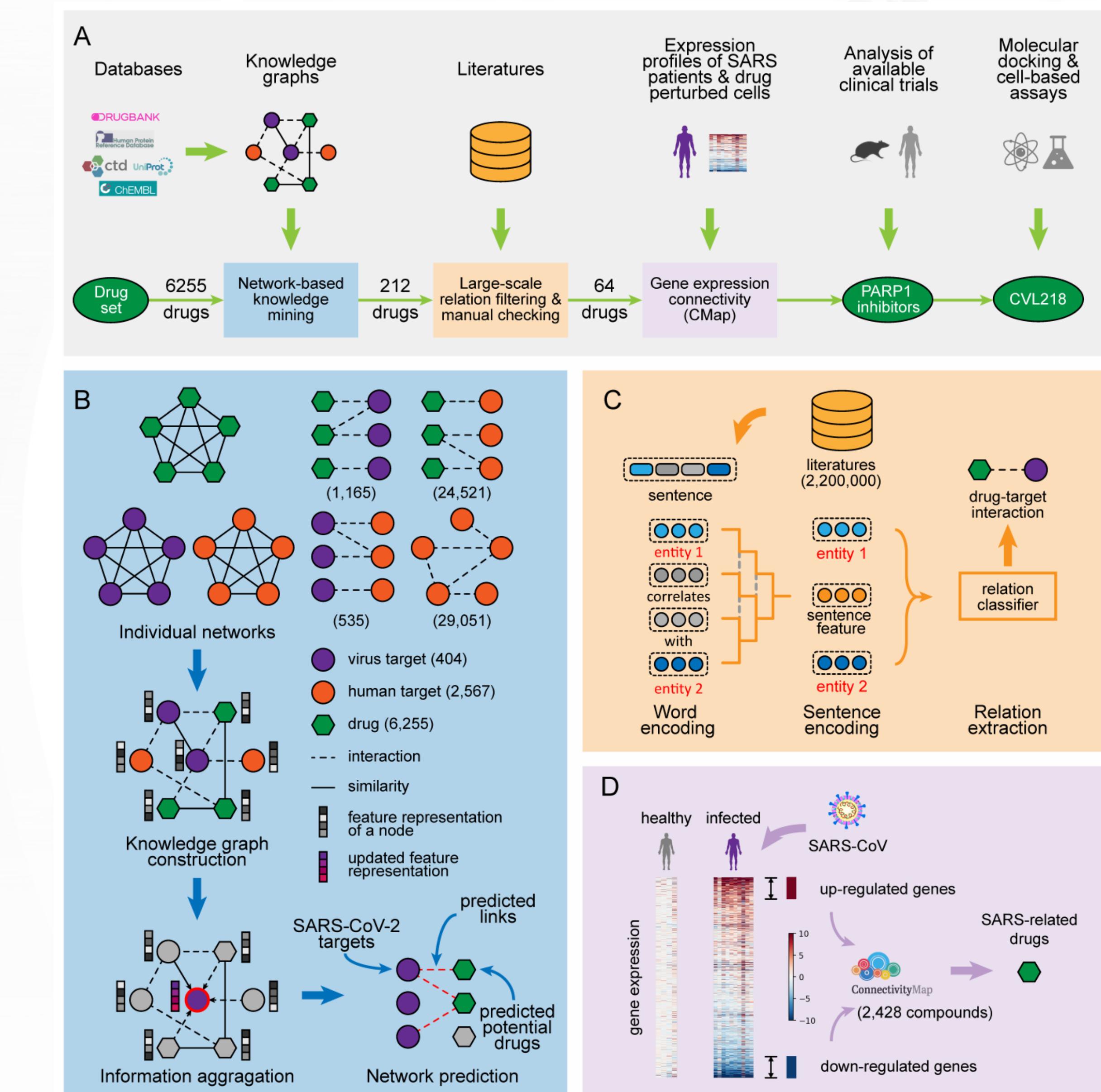
^{*}Corresponding authors.

Abstract

The global spread of SARS-CoV-2 requires an urgent need to find effective therapeutics for the treatment of COVID-19. We developed a data-driven drug repositioning framework, which applies both machine learning and statistical analysis approaches to systematically integrate and mine large-scale knowledge graph, literature and transcriptome data to discover the potential drug candidates against SARS-CoV-2. The retrospective study using the past SARS-CoV and MERS-CoV data demonstrated that our machine learning based method can successfully predict effective drug candidates

Email addresses: zhaodan2018@tsinghua.edu.cn (Dan Zhao), jszfc@vip.sina.com (Fengcai Zhu), steve.shen@convalife.com (Xiaokun Shen), zengjy321@tsinghua.edu.cn (Jianyang Zeng)

March 11, 2020



Knowledge-based DR

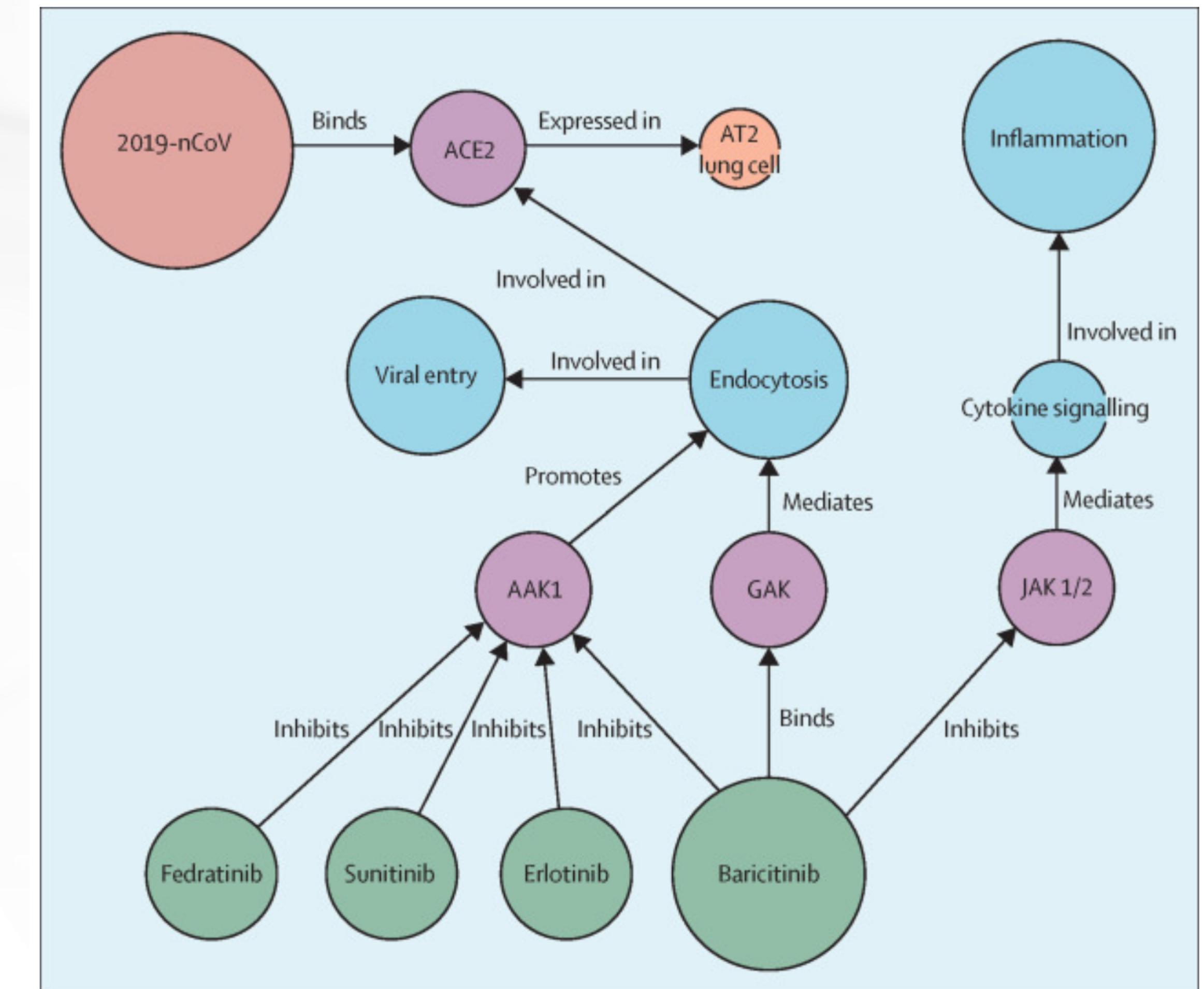
- BenevolentAI

CORRESPONDENCE | VOLUME 395, ISSUE 10223, PE30-E31, FEBRUARY 15, 2020

Baricitinib as potential treatment for 2019-nCoV acute respiratory disease

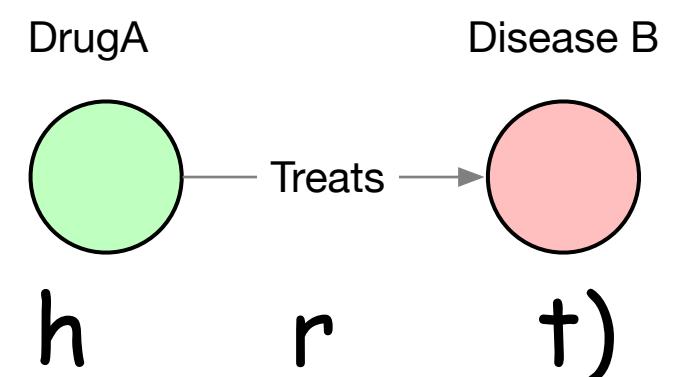
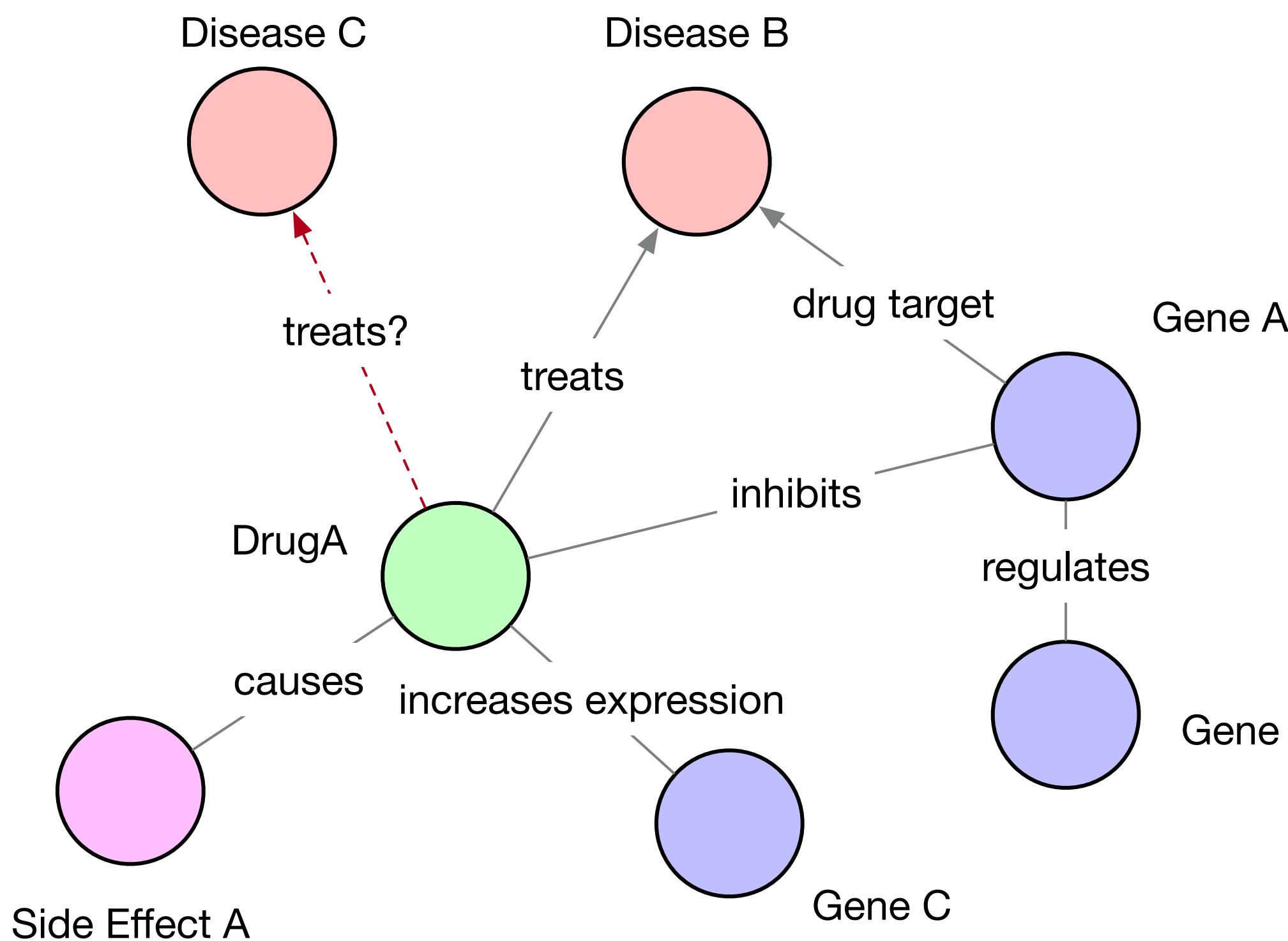
Peter Richardson • Ivan Griffin • Catherine Tucker • Dan Smith • Olly Oechsle • Anne Phelan • et al. [Show all authors](#)

Published: February 04, 2020 • DOI: [https://doi.org/10.1016/S0140-6736\(20\)30304-4](https://doi.org/10.1016/S0140-6736(20)30304-4)



What is Knowledge Graph?

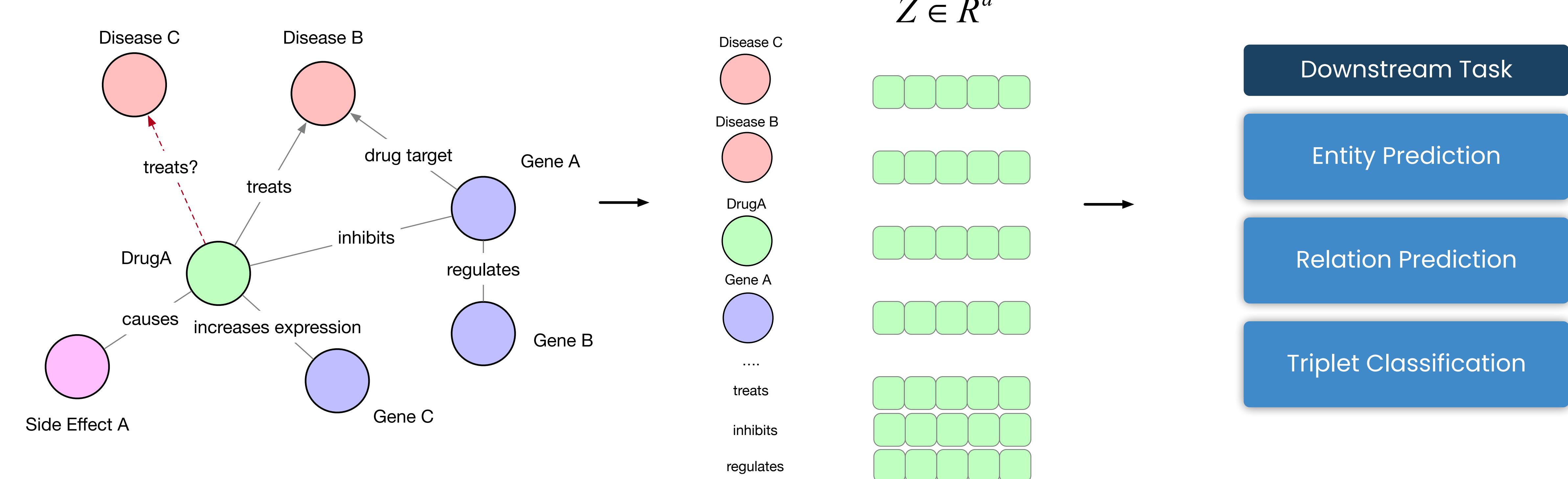
- Specific type **directed** graph whose nodes and edges represent **entities** and **relationship** between entities
 - The relationships organized in form of (head, relation, tail) -> (h, r, t) - TRIPLETS*



- Entity prediction:
 - (?, r, t)
- Relation prediction:
 - (h, ?, t) (e.g: DrugA-?-DiseaseB)
- Triplet classification
 - (h, r, t) is correct or not

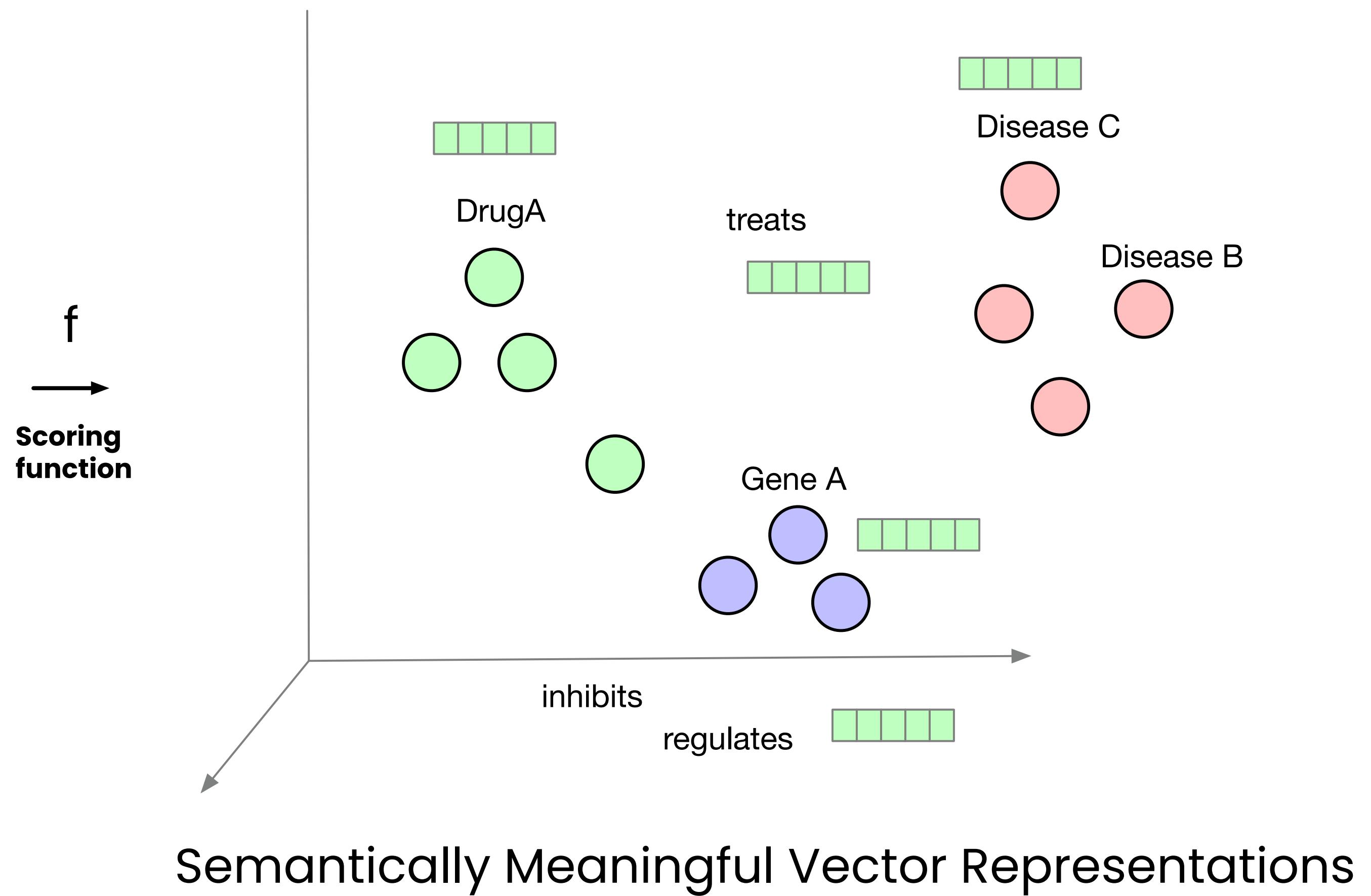
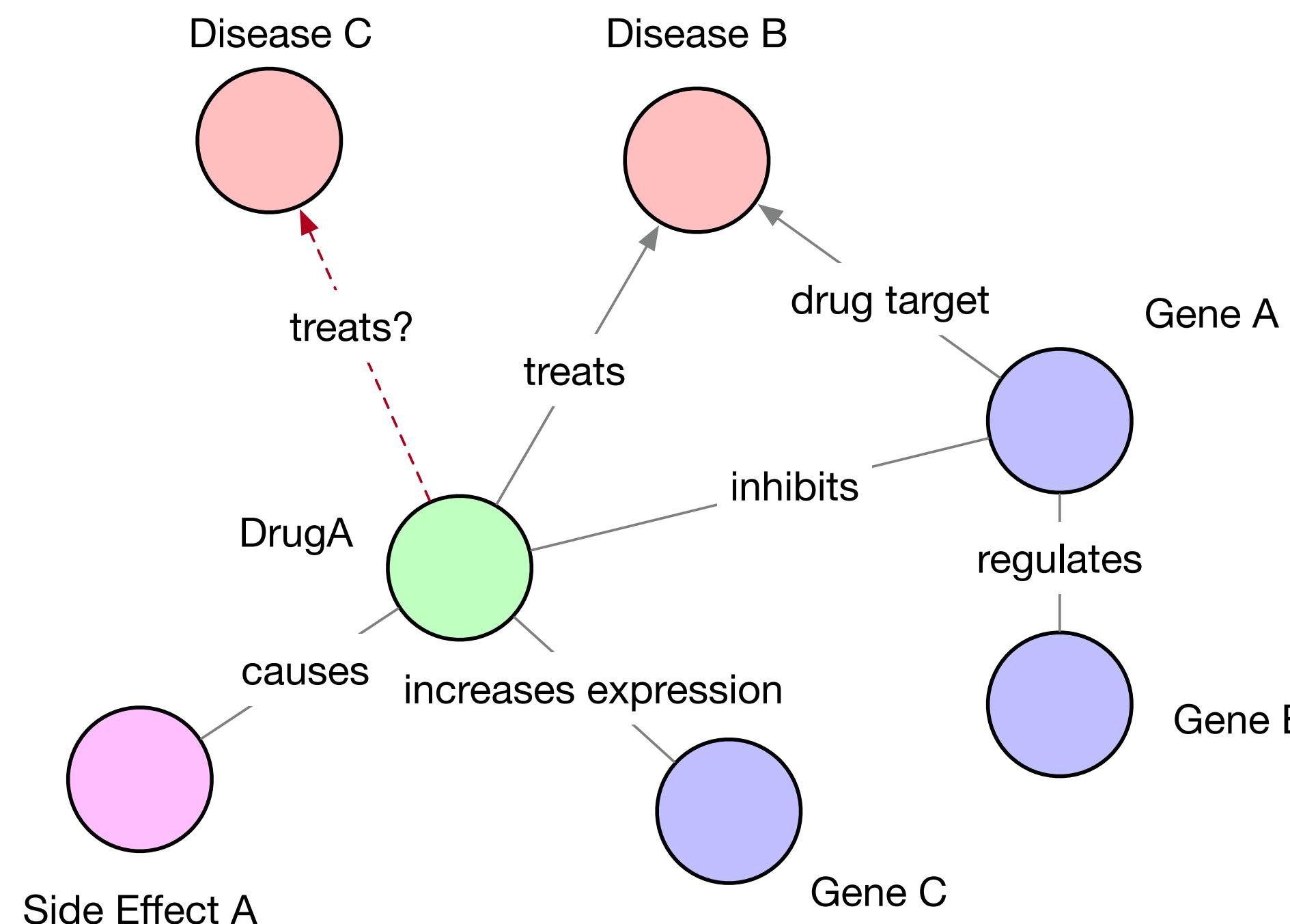
KG Embedding

projections of entities and relations into a continuous low-dimensional space.



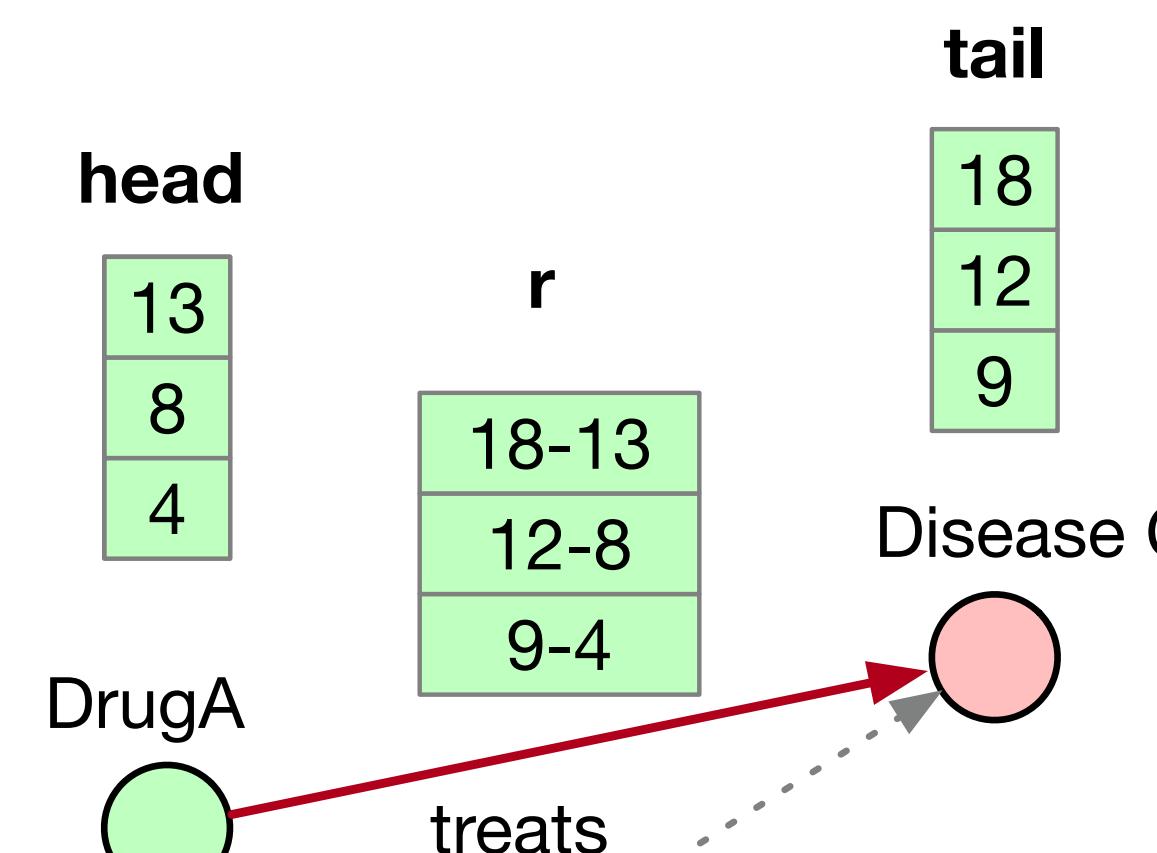
KG Embedding

projections of entities and relations into a continuous low-dimensional space.

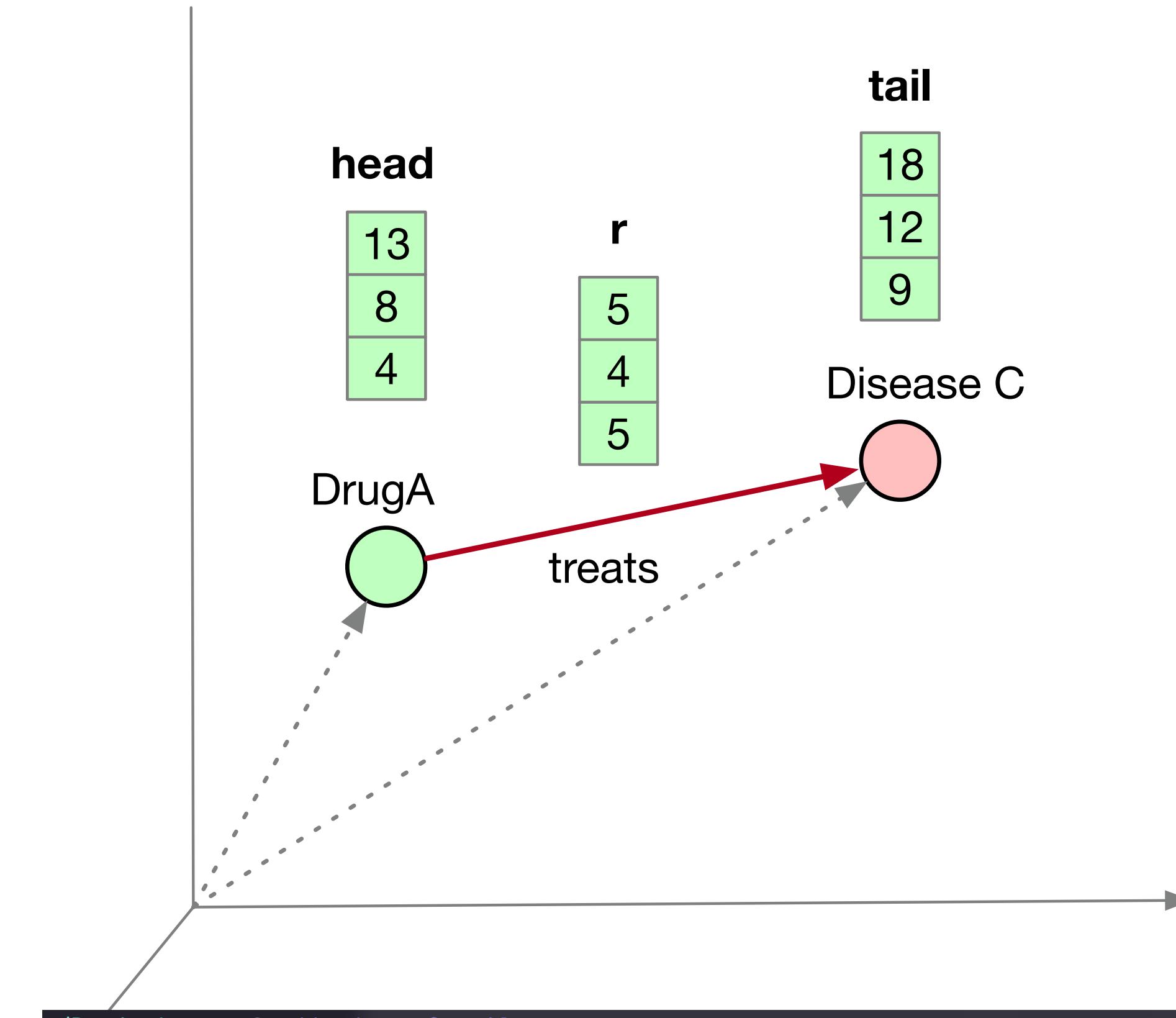


Scoring Rule

$$h+r = t$$



$$\mathcal{L} = \sum_{(h,r,t \in S)} \sum_{(h',r',t')} [\gamma + d(h+r, t) - d(h'+r', t')]$$



```
~/Projects on ⛅erkhembayar@gmail.com
> DGLBACKEND=pytorch dglke_train --dataset DRKG --data_path ./train \
--data_files drkg_train.tsv drkg_valid.tsv drkg_test.tsv \
--format 'raw_utt_hrt' --model_name TransE_12 --batch_size 2048 \
--neg_sample_size 256 --hidden_dim 400 --gamma 12.0 --lr 0.1 \
--max_step 100000 --log_interval 1000 --batch_size_eval 16 -adv \
--regularization_coef 1.00E-07 --test --num_thread 1 \
--gpu 0 1 2 3 4 5 6 7 --num_proc 8 \
--neg_sample_size_eval 10000 --async_update
```

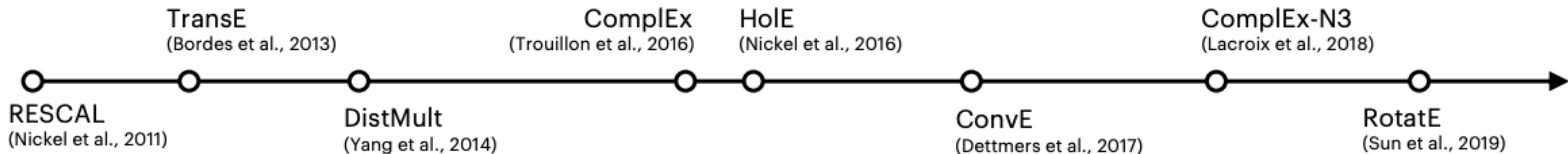
Embedding Methods

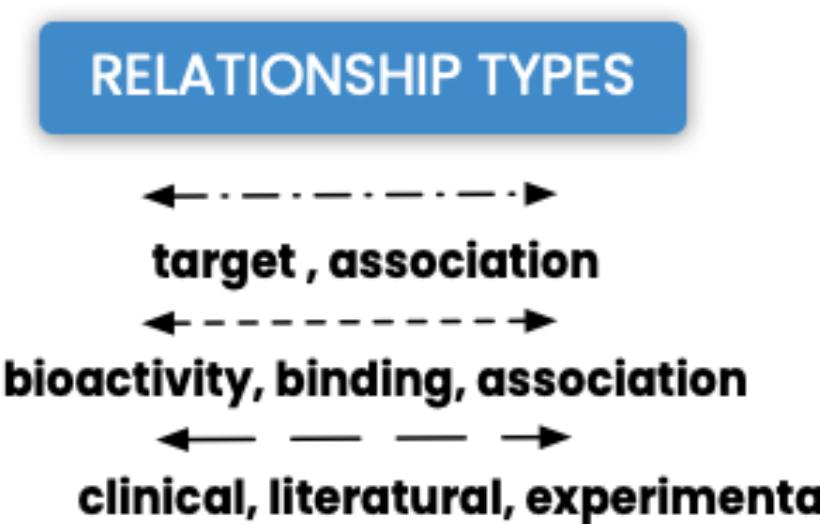
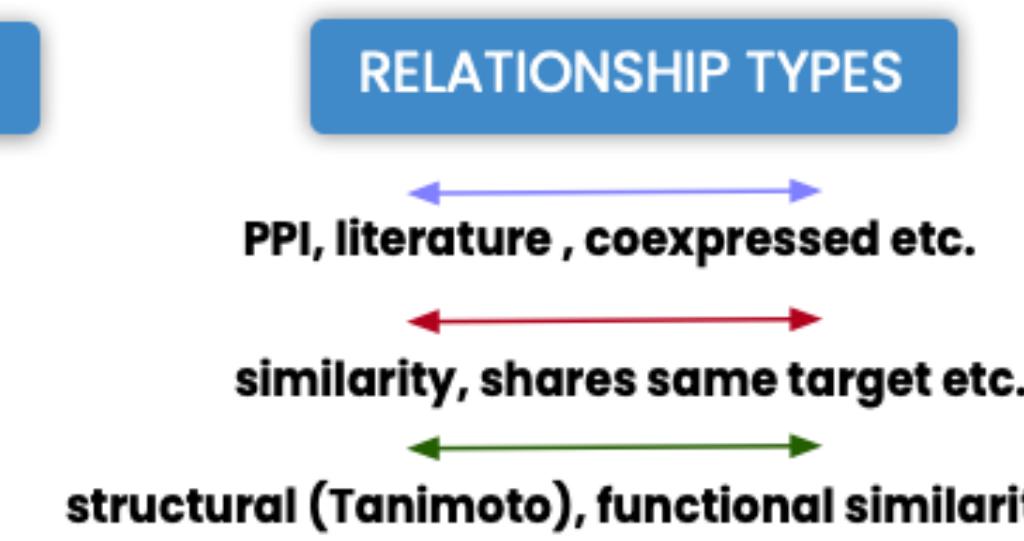
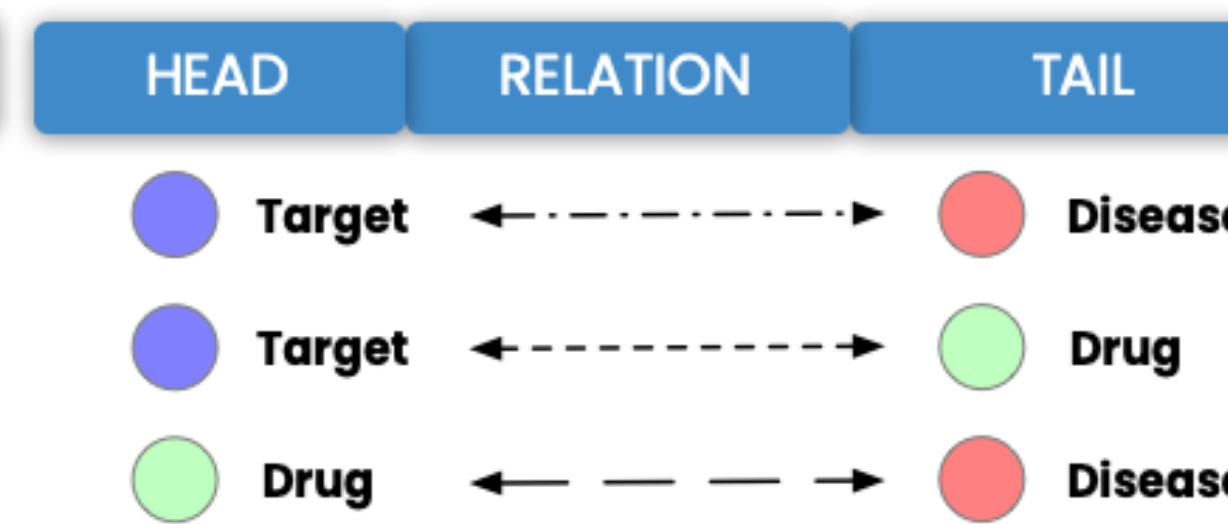
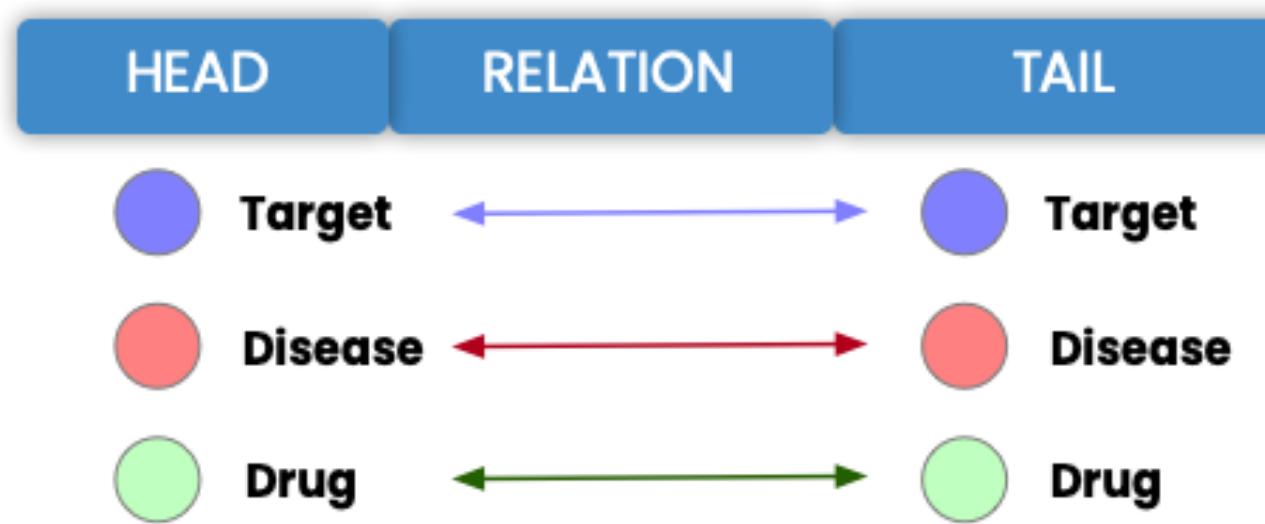
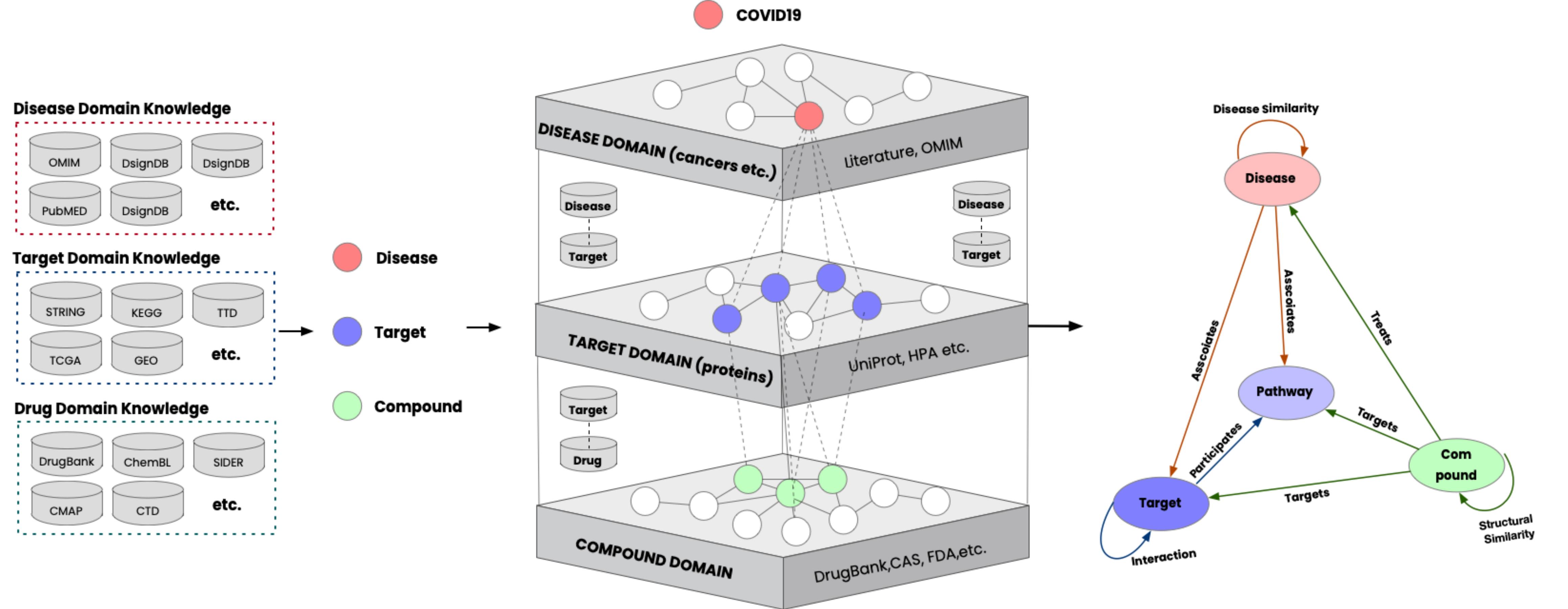
Relation Properties

Symmetry	<Alice marriedTo Bob>
Asymmetry	<Alice childOf Jack>
Inversion	<Alice childOf Jack> <Jack fatherOf Alice>
Composition	<Alice childOf Jack> <Jack siblingOf Mary> <Alice nieceOf Mary>

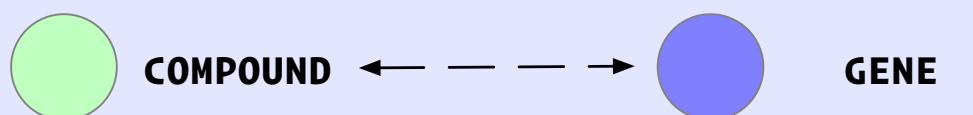
Model	Symmetry	Antisymmetry	Inversion	Composition
SE	X	X	X	X
TransE	X	✓	✓	✓
TransX	✓	✓	X	X
DistMult	✓	X	X	X
ComplEx	✓	✓	✓	X
RotatE	✓	✓	✓	✓

Models	score function $f(\mathbf{h}, \mathbf{r}, \mathbf{t})$
TransE [2]	$- \mathbf{h} + \mathbf{r} - \mathbf{t} _{1/2}$
TransR [10]	$- M_r \mathbf{h} + \mathbf{r} - M_r \mathbf{t} _2^2$
DistMult [20]	$\mathbf{h}^\top \text{diag}(\mathbf{r}) \mathbf{t}$
ComplEx [16]	$\text{Real}(\mathbf{h}^\top \text{diag}(\mathbf{r}) \bar{\mathbf{t}})$
RESCAL [12]	$\mathbf{h}^\top M_r \mathbf{t}$
RotatE [15]	$- \mathbf{h} \circ \mathbf{r} - \mathbf{t} ^2$





READING RELATIONSHIP TYPES (ONTOLOGY)



COMPOUND-DISEASE

Compound-To-Disease →

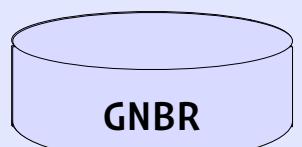
- (T) treatment/therapy (including investigatory)
- (C) inhibits cell growth (esp. cancers)
- (Sa) side effect/adverse event
- (Pr) prevents, suppresses
- (Pa) alleviates, reduces
- (J) role in disease pathogenesis

Disease-To-Compound ←

- (Mp) biomarkers (of disease progression)

33.3GB

LITERATURE (TEXT)


 Biomedical Relationships (PUBCHEM)
 ~ 24 million research articles

Compound:1234 - GNBR::(T) Compound:Disease - Disease:MESH:D01234



GENE-DISEASE

Gene-To-Disease →

- (U) causal mutations
- (Ud) mutations affecting disease course
- (D) drug targets
- (J) role in pathogenesis
- (Te) possible therapeutic effect
- (Y) polymorphisms alter risk
- (G) promotes progression

Disease-To-Gene ←

- (Md) biomarkers (diagnostic)
- (X) overexpression in disease
- (L) improper regulation linked to disease

Gene:1234 - GNBR::(D) Gene:Disease - Disease:MESH:D01234

COMPOUND-GENE

Compound-To-Gene →

- (A+) agonism, activation
- (A-) antagonism, blocking
- (B) binding, ligand (esp. receptors)
- (E+) increases expression/production
- (E-) decreases expression/production
- (E) affects expression/production (neutral)
- (N) inhibits

Gene-To-Compound ←

- (O) transport, channels
- (K) metabolism, pharmacokinetics
- (Z) enzyme activity

Compound:1234 - GNBR::(B) Compound:Gene - Gene:4321



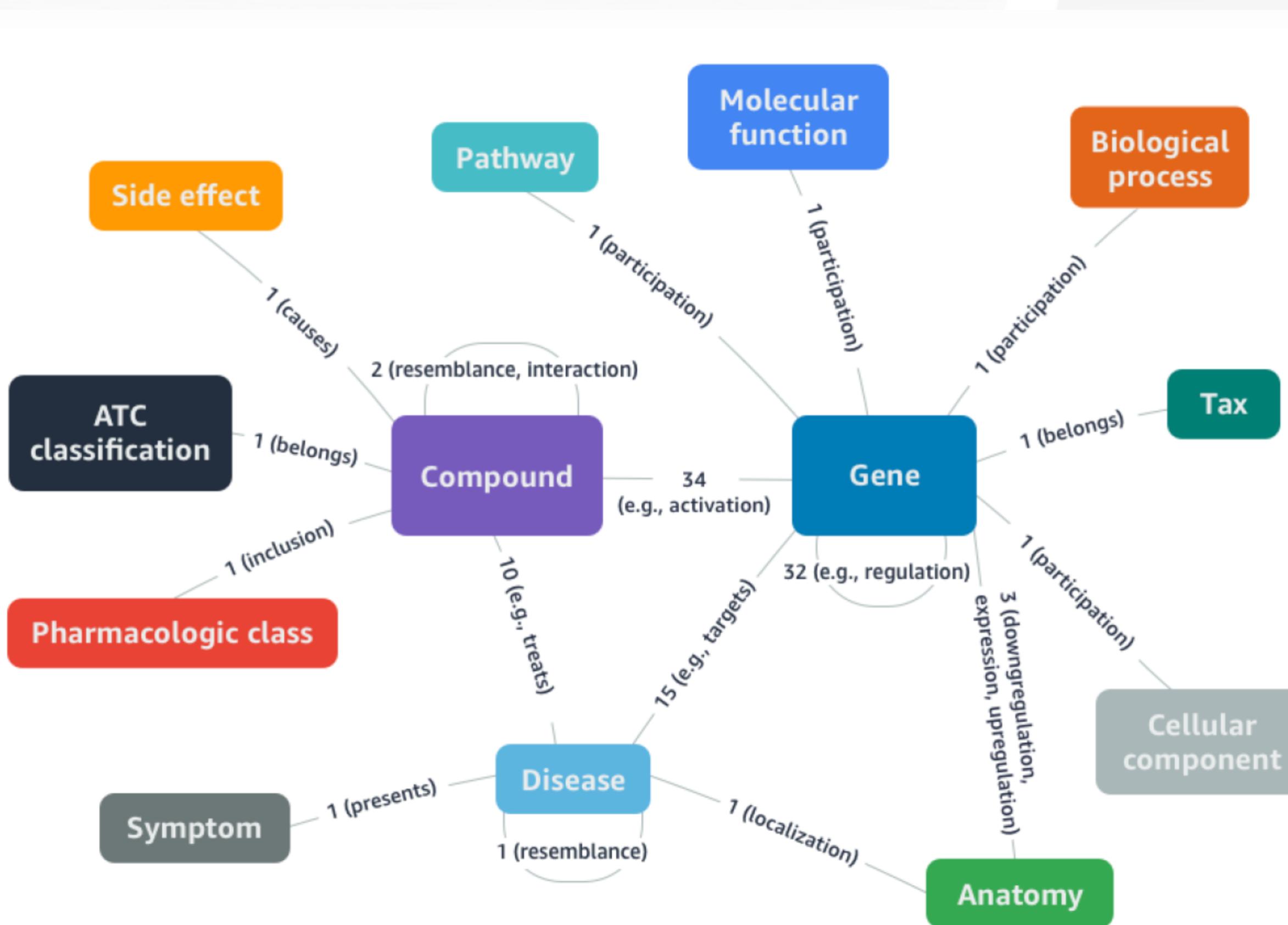
GENE-GENE

- (B) binding, ligand (esp. receptors)
- (W) enhances response
- (V+) activates, stimulates
- (E+) increases expression/production
- (E) affects expression/production (neutral)
- (I) signaling pathway
- (H) same protein or complex
- (Rg) regulation
- (Q) production by cell population

Gene:1234 - GNBR::(Rg) Gene:Gene - Gene:4321

Knowledge Graph Database

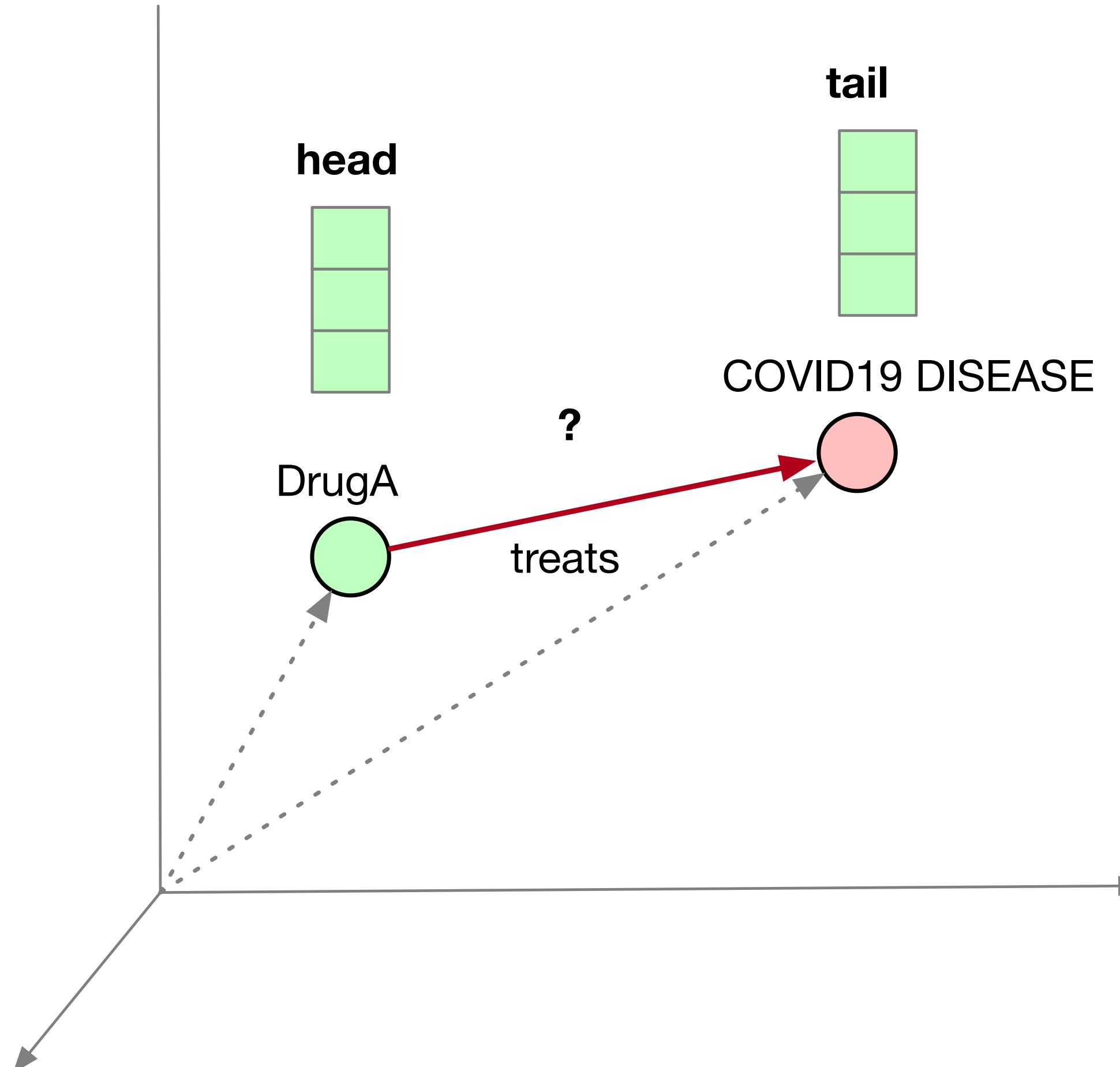
- Amazon Deep Engine Science Team
 - University of Minnesota, The Ohio State University and
 - **Drug Repurposing Knowledge Graph (DRKG)**



- 6 DBs
 - DrugBank, GNBR,Hetionet (Knowledge Database),STRING (PPI), IntAct (PPI), Bibliography
- 97,328 entities
 - 13 types of entities
- 5,874,461 triplets
 - 107 types of relations
- Embedded vectors

Knowledge-based DR for COVID19

Link Prediction:
[**h, ?, t**]



$$\mathbf{d} = \gamma - \|\mathbf{h} + \mathbf{r} - \mathbf{t}\|_2$$

$$\text{score} = \log\left(\frac{1}{1+\exp(-\mathbf{d})}\right)$$

**Drugs will be repositioned (ranked)
according to this score**

Summary

- **Graph natural** representative power which helps to find patterns in data that scientist might not see
- The solution to many applications in biomedical field can be formulated as **GRAPH LEARNING PROBLEMs**
- **Graphs** are used to
 - Detect new target by taking advantages of complementeriness in variety of biomedical resources
 - Represent chemical compounds and showed promising performance
- **KNOWING GRAPH REPRESENTATION is IMPORTANT in DRUG DISCOVERY**

LET`S DO SOME CODING



Thank You