

# Knowledge-driven AI enables mechanistic predictive toxicology

*Presented by Joseph D. Romano, PhD*

*November 2, 2023*

DEPARTMENT of  
**BIOSTATISTICS**  
**EPIDEMIOLOGY &**  
**INFORMATICS**

Center of Excellence in  
Environmental Toxicology

Penn Institute for  
Biomedical Informatics



# The Romano Lab @ Penn



**Joe Romano, PhD**  
Assistant Professor of Informatics and Pharmacology



**Chloé Paris**  
PhD Student (GCB)



**Yunchae Kim, MS**  
Software Developer



**Estefania Morales**  
Summer Undergraduate Internship Program (SUIP) Fellow

**Not Pictured:**

**Kevin Shen**  
Undergraduate Researcher

**Tom Pan**  
Undergraduate Researcher

# We study...

- Computational toxicology
  - Designing interpretable AI models to predict toxicity
  - Discovering *mechanisms* linking specific chemicals to human diseases
- Autoimmunity
  - Biomarker discovery for immune-related adverse events (irAEs)
- Clinical informatics
  - How to computationally model and extract knowledge from electronic health record data
- Artificial Intelligence and Machine Learning
  - Multimodal graph machine learning
  - Automated machine learning
  - Large language models (LLMs) and their applications in biomedicine

# **Computational Toxicology and Artificial Intelligence**

# The role of AI in toxicology

- AI should help experimental toxicologists by:
  - **Predicting** new associations between chemicals and endpoints of toxicity
  - **Explaining** mechanisms that may underlie those predictions
- It does not replace experimental validation; rather, it helps us to prioritize our time and effort

## Drugs (incl. withdrawn, investigational...)

DRUGBANK Online

Browse ▾

Search ▾

Interaction Checker

 Small Molecule Drugs

 Biotech Drugs

FILTER BY GROUP

Approved Nutraceutical Illicit

Investigational Withdrawn Experimental

FILTER BY MARKET AVAILABILITY

U.S. Canada E.U.

Displaying drugs 1 - 25 of 12220 in total

1 2 3 4 5 ... > »

## Environmental chemicals of toxicological concern

CompTox Chemicals Dashboard

Search 1,200,059 Chemicals

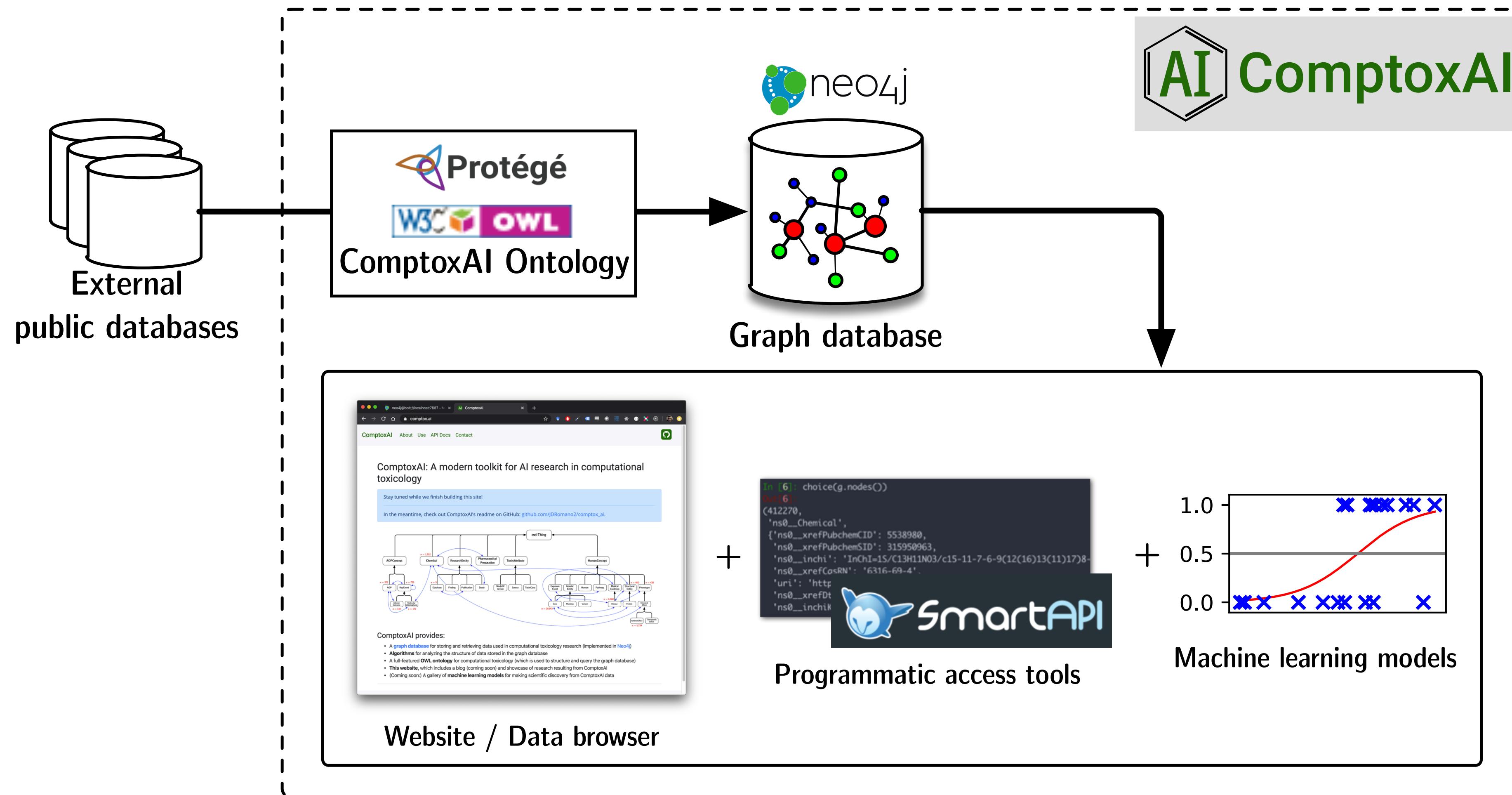
Chemicals Products/Use Categories Assay/Gene

Search for chemical by systematic name, synonym, CAS number, DTXSID or InChIKey

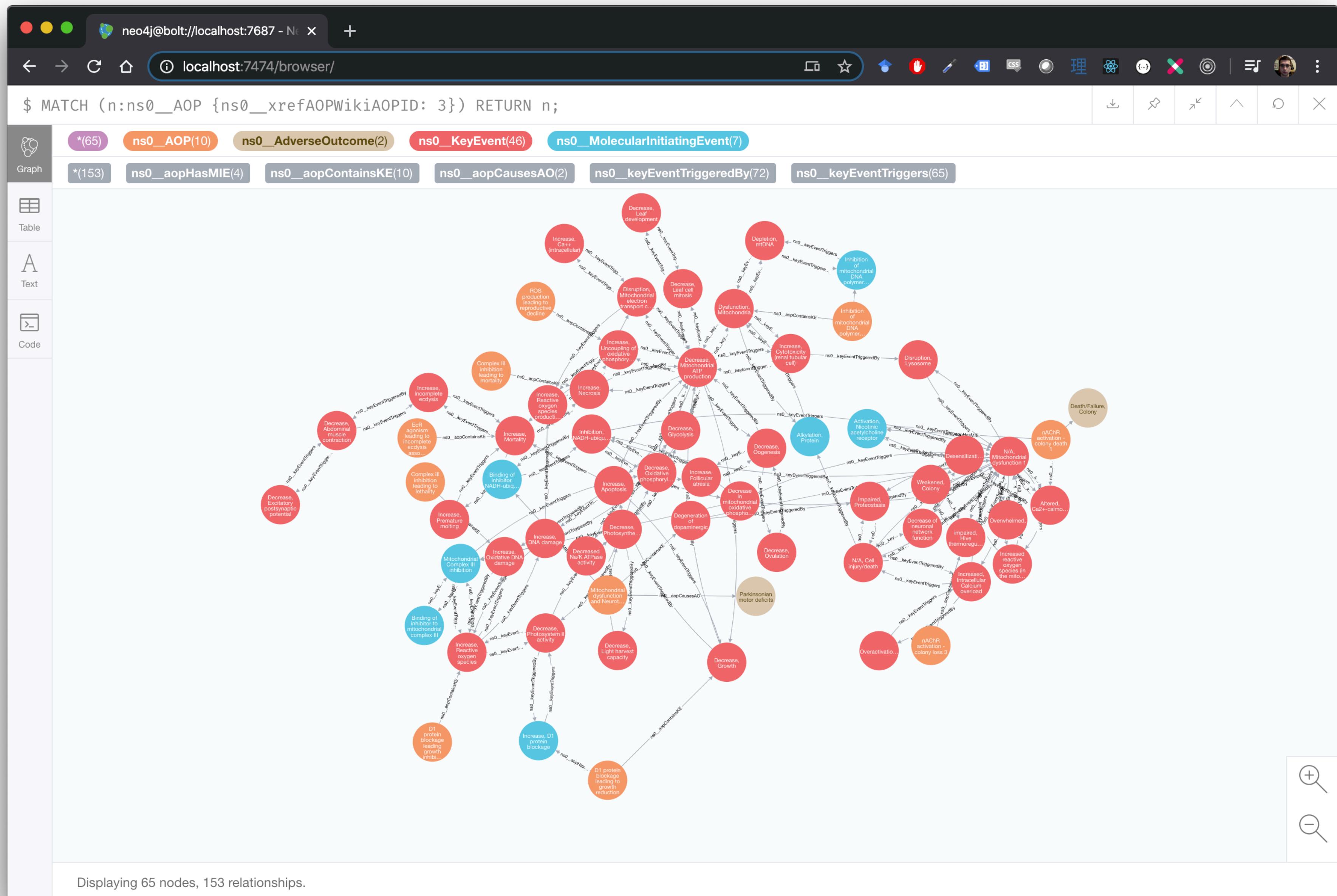
Start typing to search.

Identifier substring search

# ComptoxAI



<https://comptox.ai>



# ComptoxAI: Data Interfaces

- Data browser / dataset generator tools on website
- Direct access to graph database (local or remote)
- Web API (Programmatic access to data)
- Python package (Access data and construct machine learning models from the Python programming language)

# Simplified data browser

## ComptoxAI interactive data portal

From this page, you can search for individual entities (nodes) in ComptoxAI's graph database. When you select a query result, adjacent nodes (related data elements) are loaded and displayed below.

For detailed usage instructions, please see [this page](#).

### Nodes

Search [LOAD EXAMPLE QUERY](#)

Node Type  
Gene

Node Field  
Gene Symbol

Value  
CYP2E1

[SEARCH](#) [CLEAR FORM](#)

#### Search Results

[CLEAR NODE SEARCH RESULTS](#)

Node details:

 Gene cytochrome P450 family 2 subfamily E member 1

[COPY JSON](#)

External Identifiers:

Database	Identifier	Feature name	Value
OMIM ID	124040	typeOfGene	protein-coding
HGNC ID	2631	geneSymbol	CYP2E1

Other node features:

Feature name	Value
typeOfGene	protein-coding
geneSymbol	CYP2E1

# Graph database: Neo4j browser

The screenshot shows the Neo4j browser interface with a dark theme. On the left, a sidebar contains various navigation icons and sections like 'Database Information' and 'Relationship Types'. The main area displays a graph visualization of 'AdverseOutcome' nodes. A central node is labeled 'Lung fibrosis; Lung fibro...'. Other visible nodes include 'Hepatot...', 'Necrosis', 'Increase...', 'Decrease...', 'Reprodu...', 'Apoptosi...', 'Chemical induced', and 'Facial cartilage struct...'. A tooltip for 'KEYEVENT...' points to a relationship between 'Apoptosi...' and 'Reprodu...'. The top navigation bar shows the URL 'neo4j://neo4j.comptox.ai:7687/' and the bottom navigation bar includes links for 'Getting started', 'Try Neo4j with live', and 'Cypher basics'.

neo4j\$

```
neo4j$ MATCH (n:AdverseOutcome) RETURN n LIMIT 25
```

neo4j\$ :play start

neo4j

Getting started Try Neo4j with live Cypher basics

# Web API

The screenshot shows a web browser displaying the ComptoxAI REST API documentation via the Swagger UI. The title "ComptoxAI REST API" is at the top, followed by a version "1.0.0" and an "OAS3" badge. A sub-header states: "A REST Web API providing programmatic access to ComptoxAI's graph database." Below this, under the "Servers" section, the URL "https://comptox.ai/api - ComptoxAI's public REST API" is listed. The main content area is divided into sections: "nodes" and "paths". The "nodes" section contains five GET methods:

- GET /nodes/listNodeTypes** Get a list of all node types in ComptoxAI
- GET /nodes/listNodeTypeProperties/{type}** Get a list of properties defined for a particular node type
- GET /nodes/{type}/search** Search for a node using string matching on a specific field
- GET /nodes/{type}/searchContains** Search for a node where a certain field contains a query string
- GET /nodes/fetchById/{id}** Fetch a single node using its Neo4j ID

The "paths" section contains one GET method:

- GET /paths/findByIds** Use a start node ID and end node ID to retrieve a shortest path connecting those nodes

# Python package

```
Python 3.7.4 (default, Aug 13 2019, 15:17:50)
Type 'copyright', 'credits' or 'license' for more information
IPython 7.13.0 -- An enhanced Interactive Python. Type '?' for help.
```

```
In [1]: from comptox_ai.db import GraphDB
```

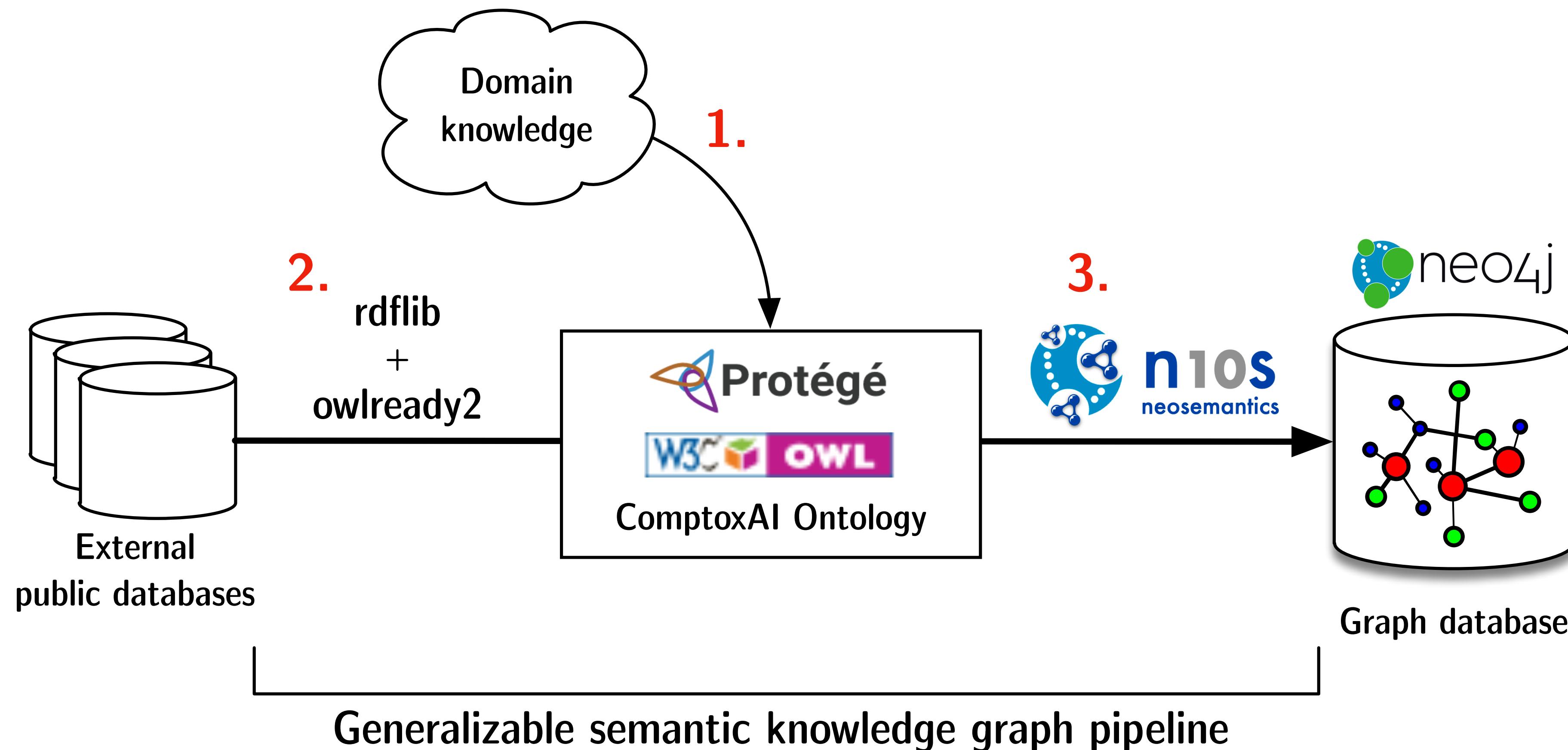
RDFLib Version: 5.0.0

```
In [2]: db = GraphDB(hostname="neo4j.comptox.ai")
```

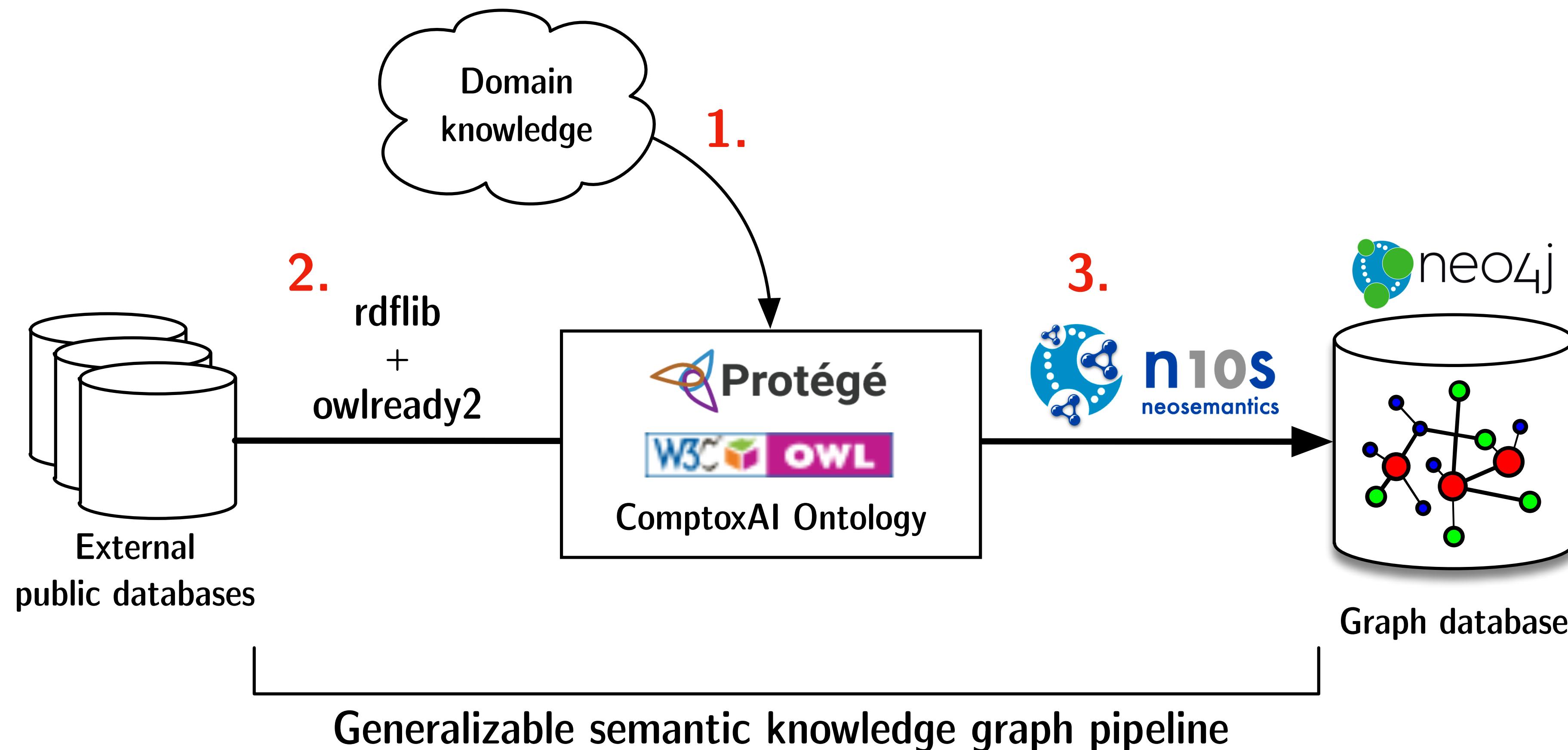
```
In [3]: result = db.run_cypher("MATCH (c:Chemical {commonName: 'PFOA'}) RETURN c;")
```

```
In [4]: print(result)
```

# Fragmented data → Knowledge graph



# Fragmented data → Knowledge graph



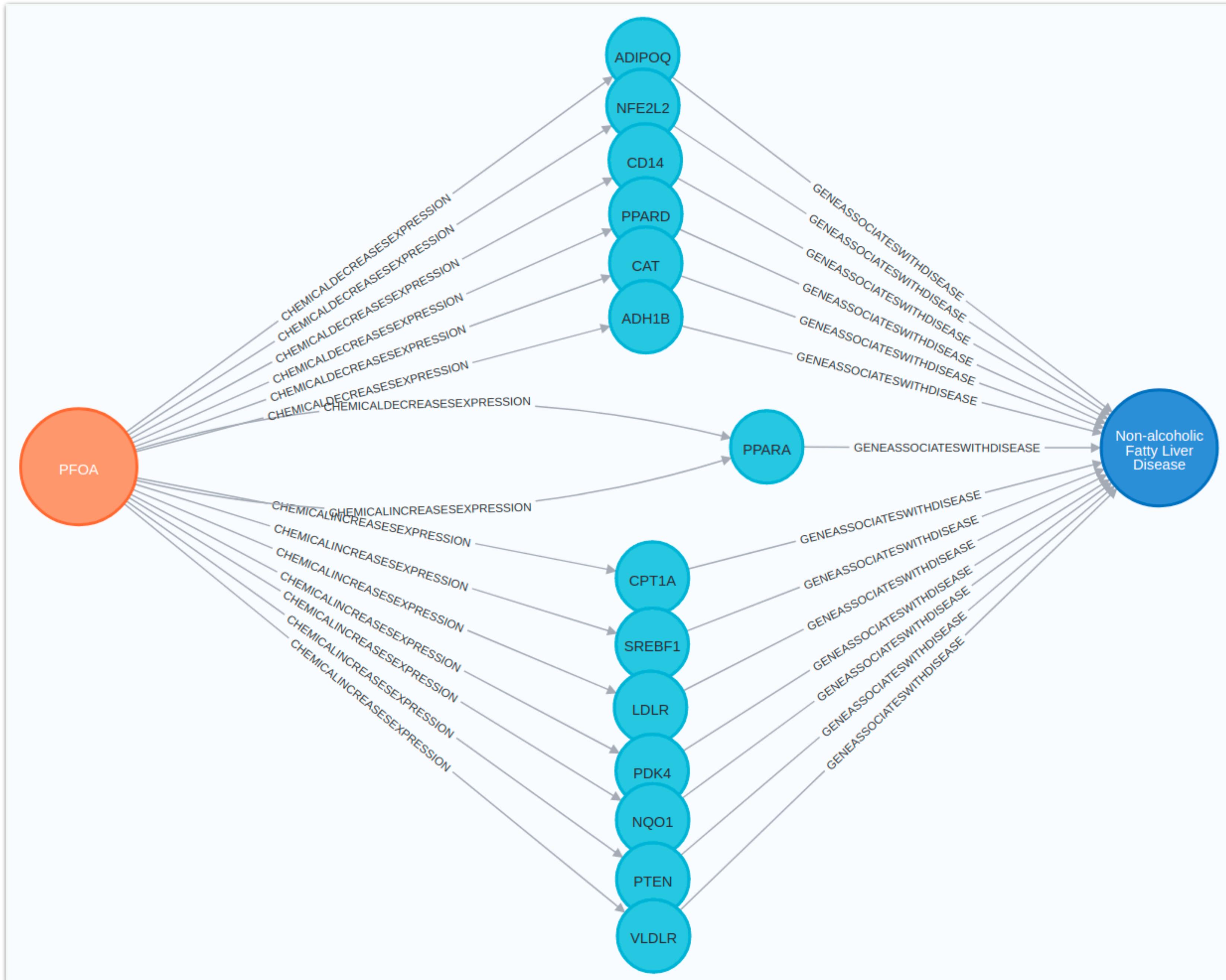
Now used in ComptoxAI, AlzKB, and other KGs under development!

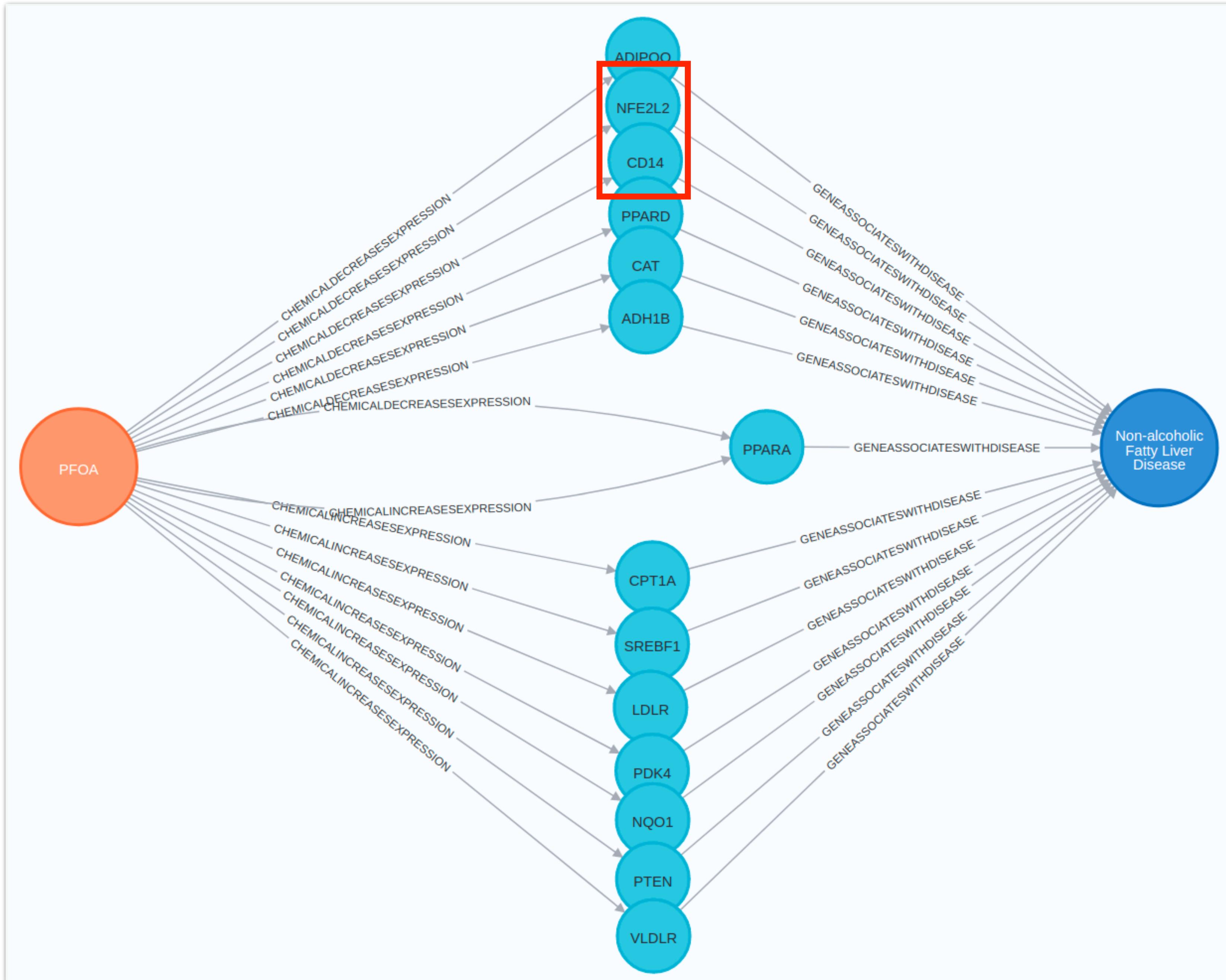
# 2 paradigms for discovery:

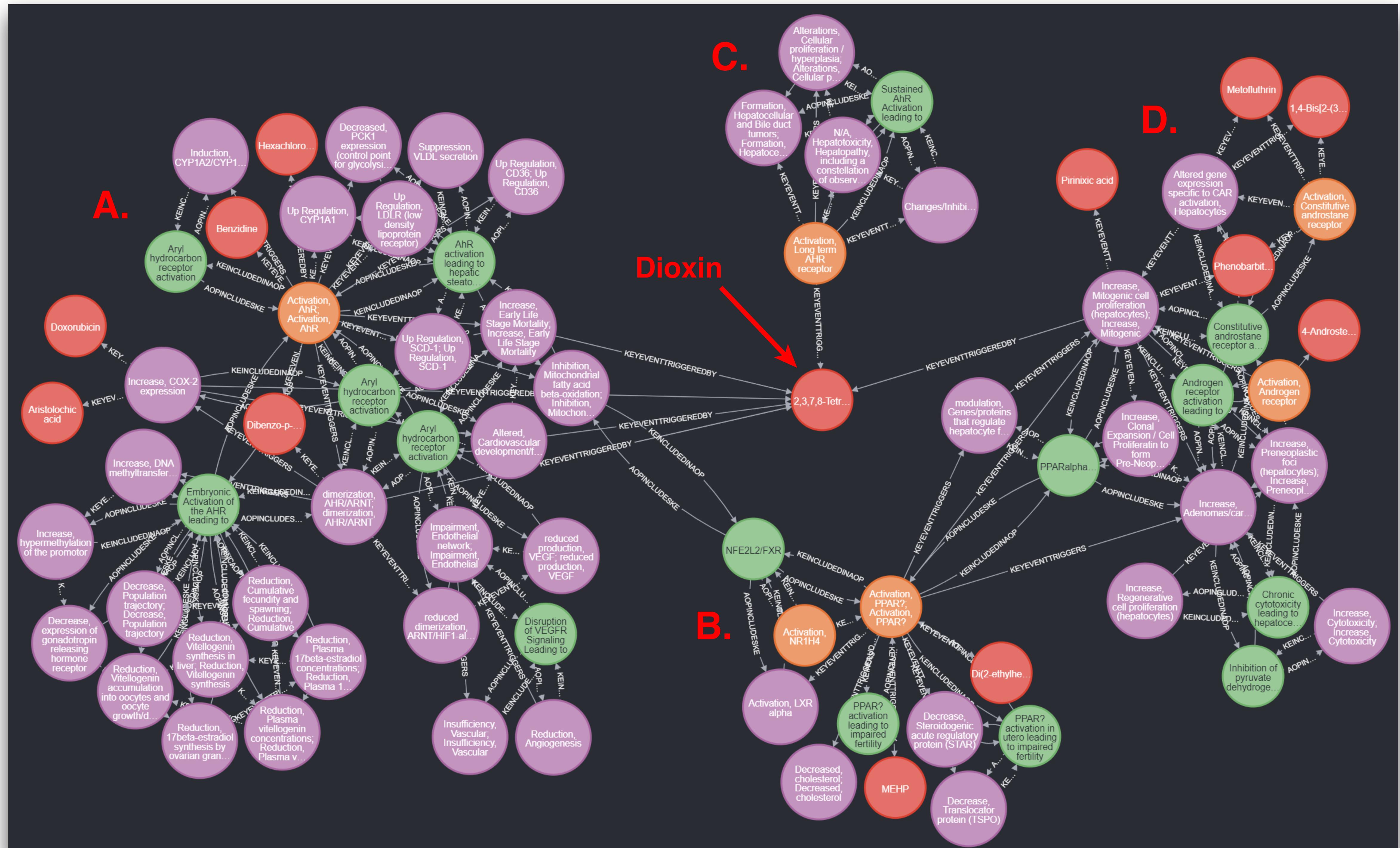
- **Directly inferring relationships from network structure**
  - Stitching together chains of relationships originally fragmented across multiple sources
  - Inspecting higher-order network architecture
- **Predicting new, unobserved relationships**
  - Machine learning (especially graph machine learning)

# 2 paradigms for discovery:

- Directly inferring relationships from network structure
  - Stitching together chains of relationships originally fragmented across multiple sources
  - Inspecting higher-order network architecture
- Predicting new, unobserved relationships
  - Machine learning (especially graph machine learning)







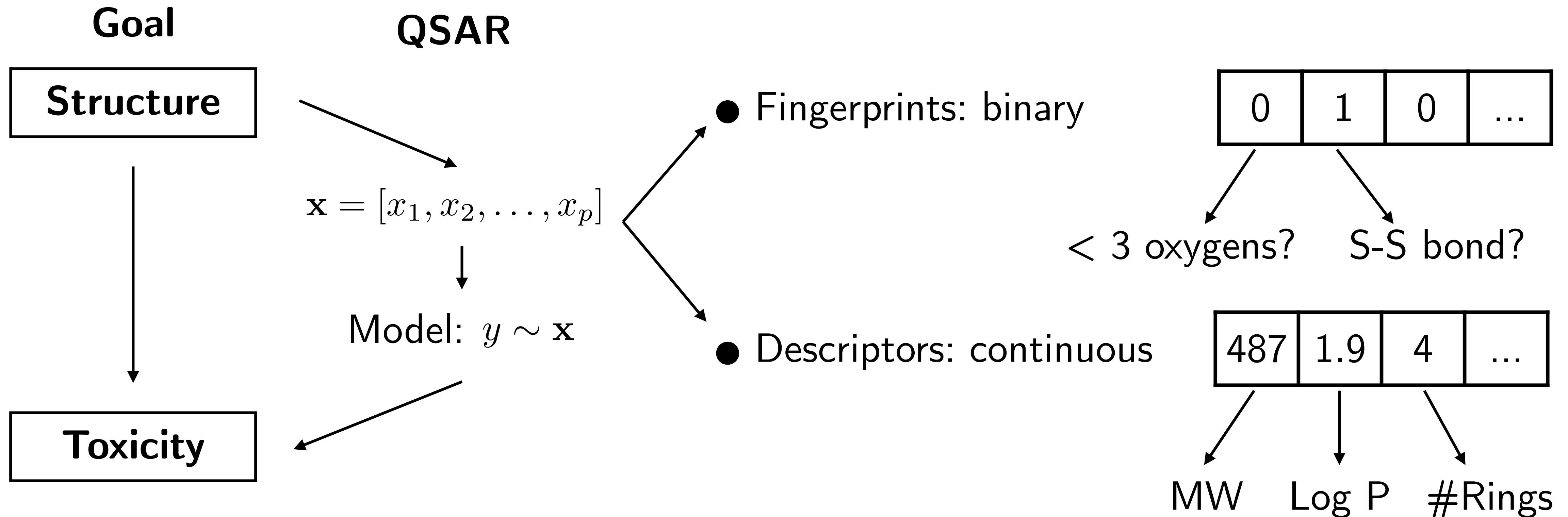
# 2 paradigms for discovery:

- **Directly inferring relationships from network structure**
  - Stitching together chains of relationships originally fragmented across multiple sources
  - Inspecting higher-order network architecture
- **Predicting new, unobserved relationships**
  - Machine learning (especially graph machine learning)

# 2 paradigms for discovery:

- Directly inferring relationships from network structure
  - Stitching together chains of relationships originally fragmented across multiple sources
  - Inspecting higher-order network architecture
- Predicting new, unobserved relationships
  - Machine learning (especially graph machine learning)

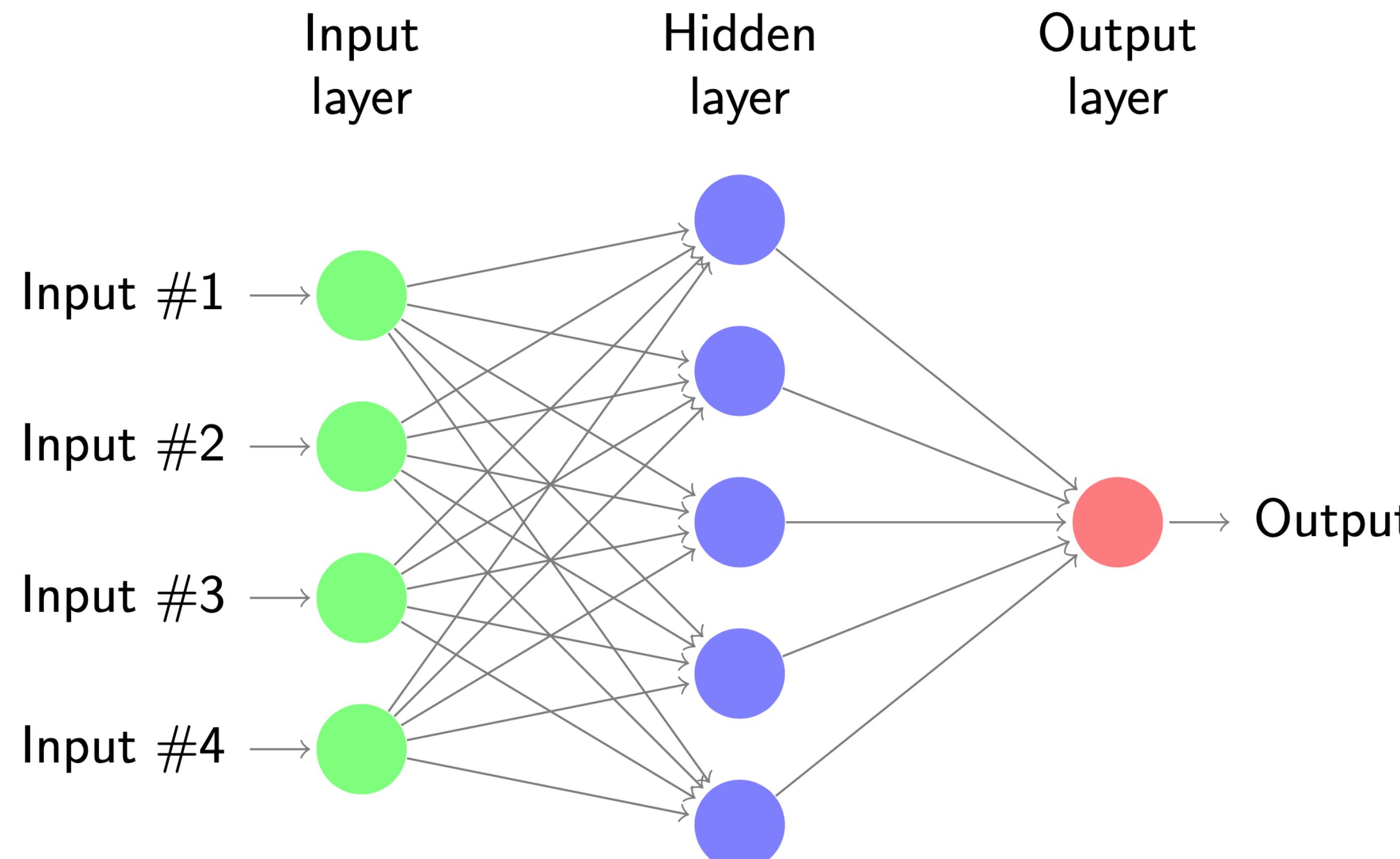
# QSAR



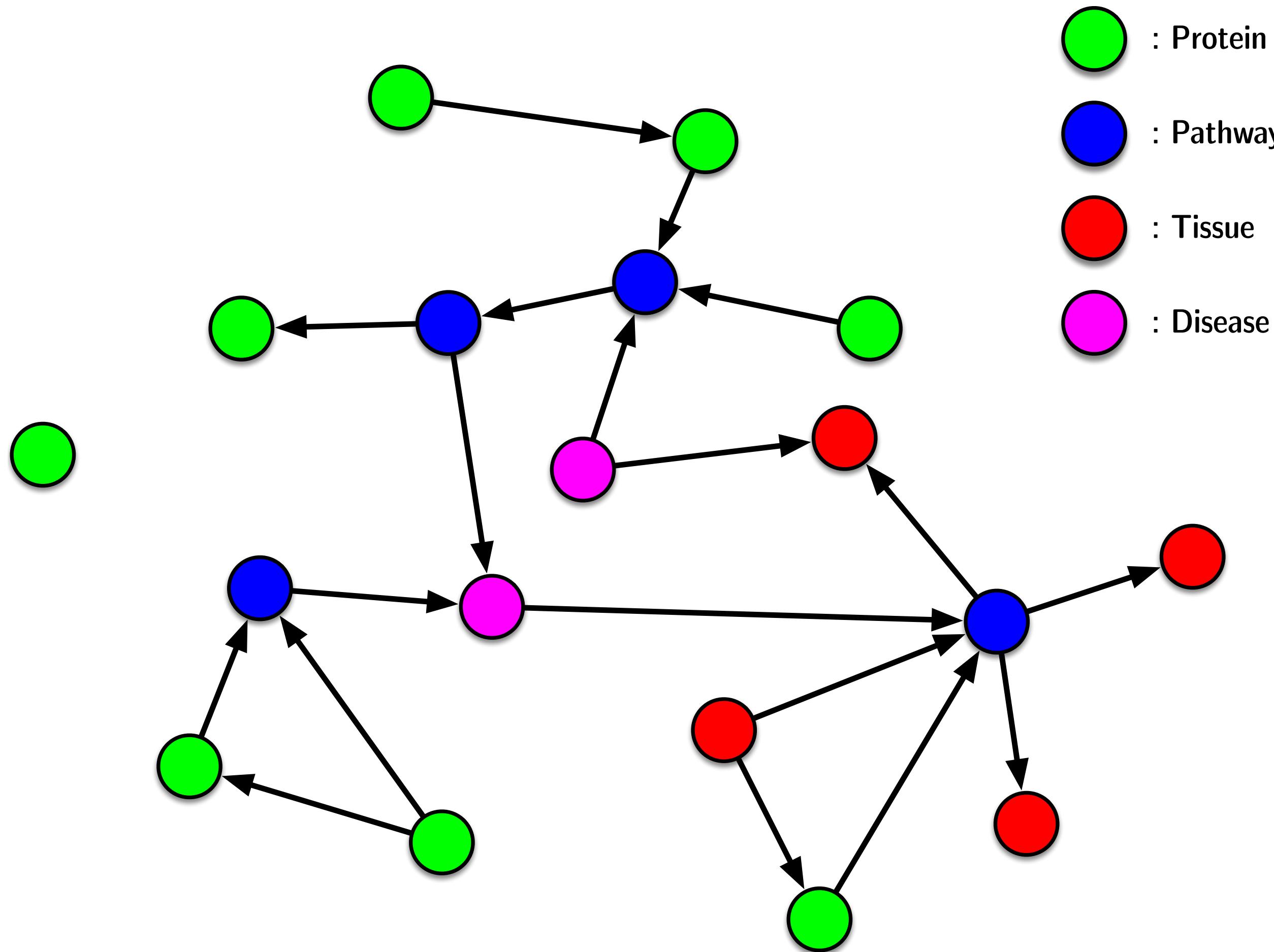
- QSAR: Quantitative Structure-Activity Relationship

# (Artificial) Neural Networks

- Consist of **nodes** organized into **layers**, which are usually stacked
- Deep learning —> NN with tens or hundreds of layers

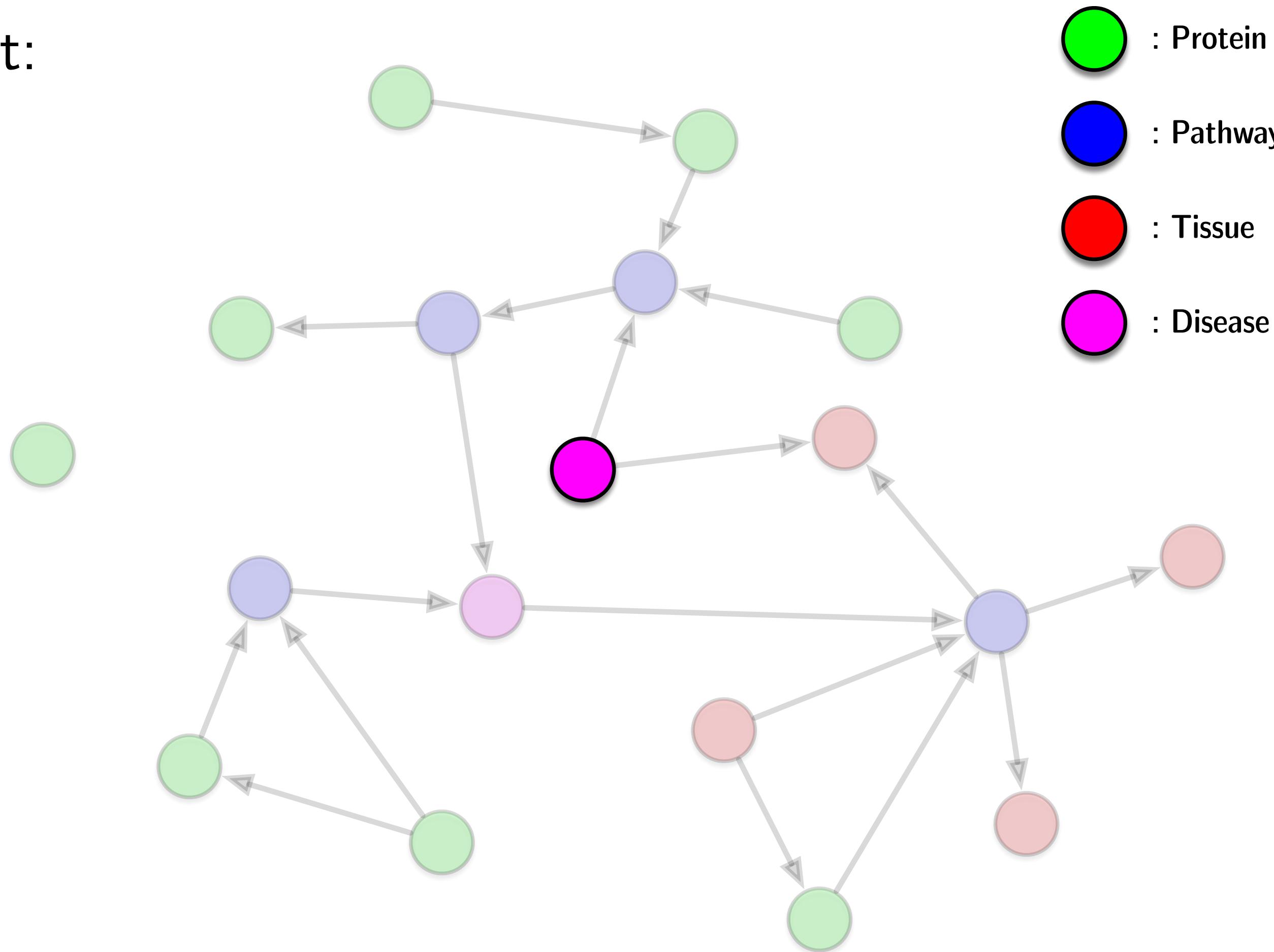


# Graph Neural Networks



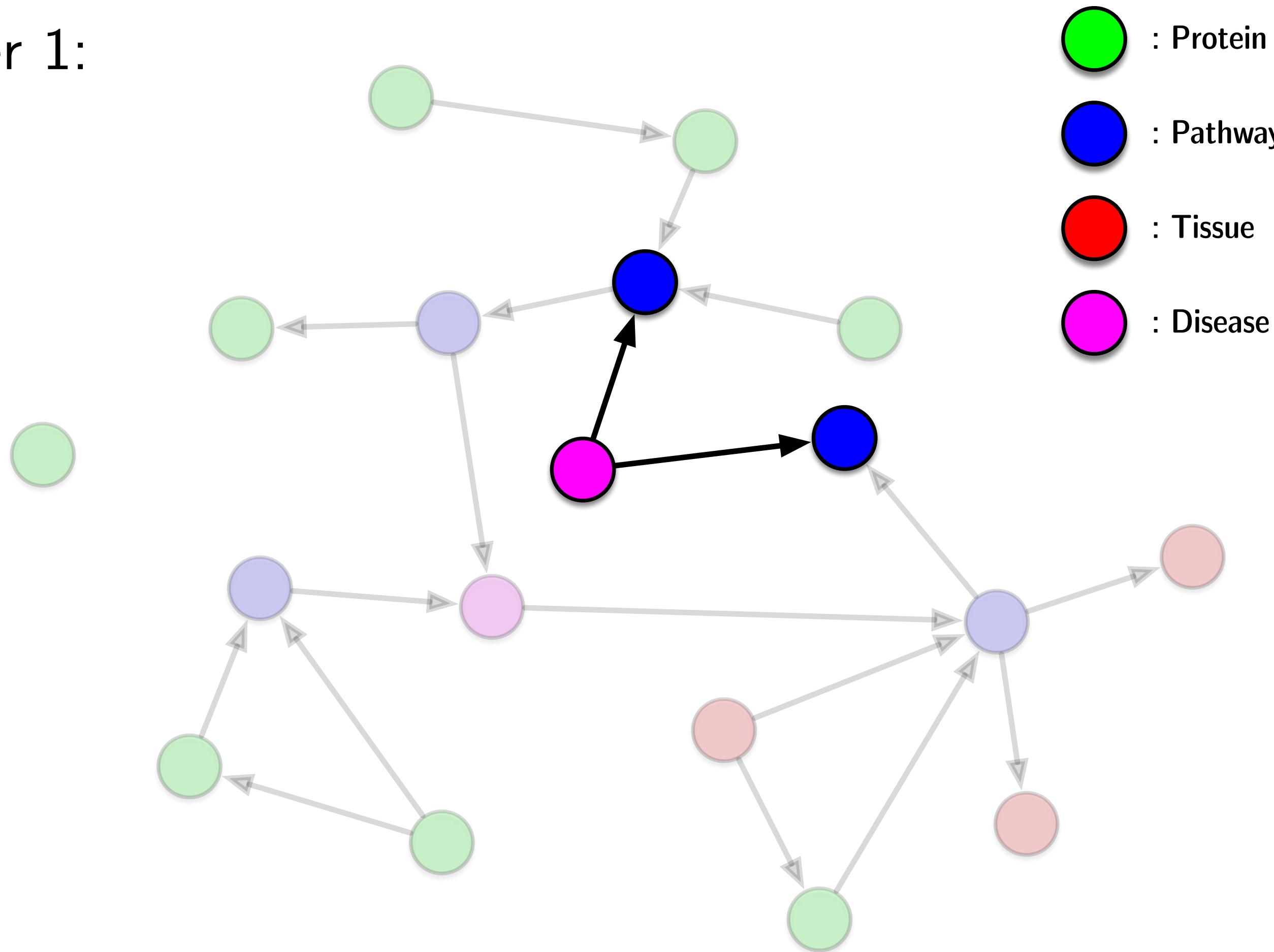
# Graph Neural Networks

- Input:



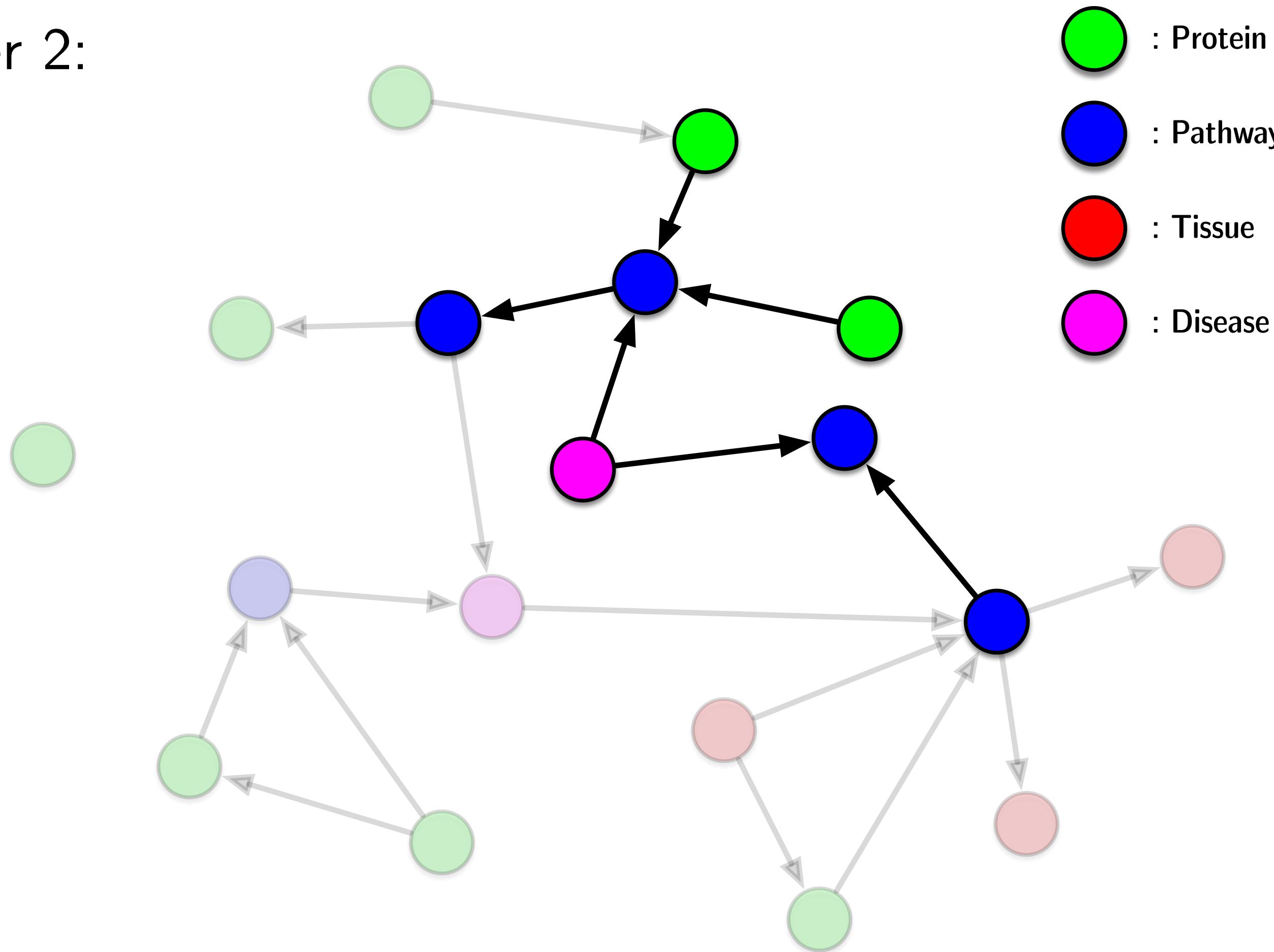
# Graph Neural Networks

- Layer 1:



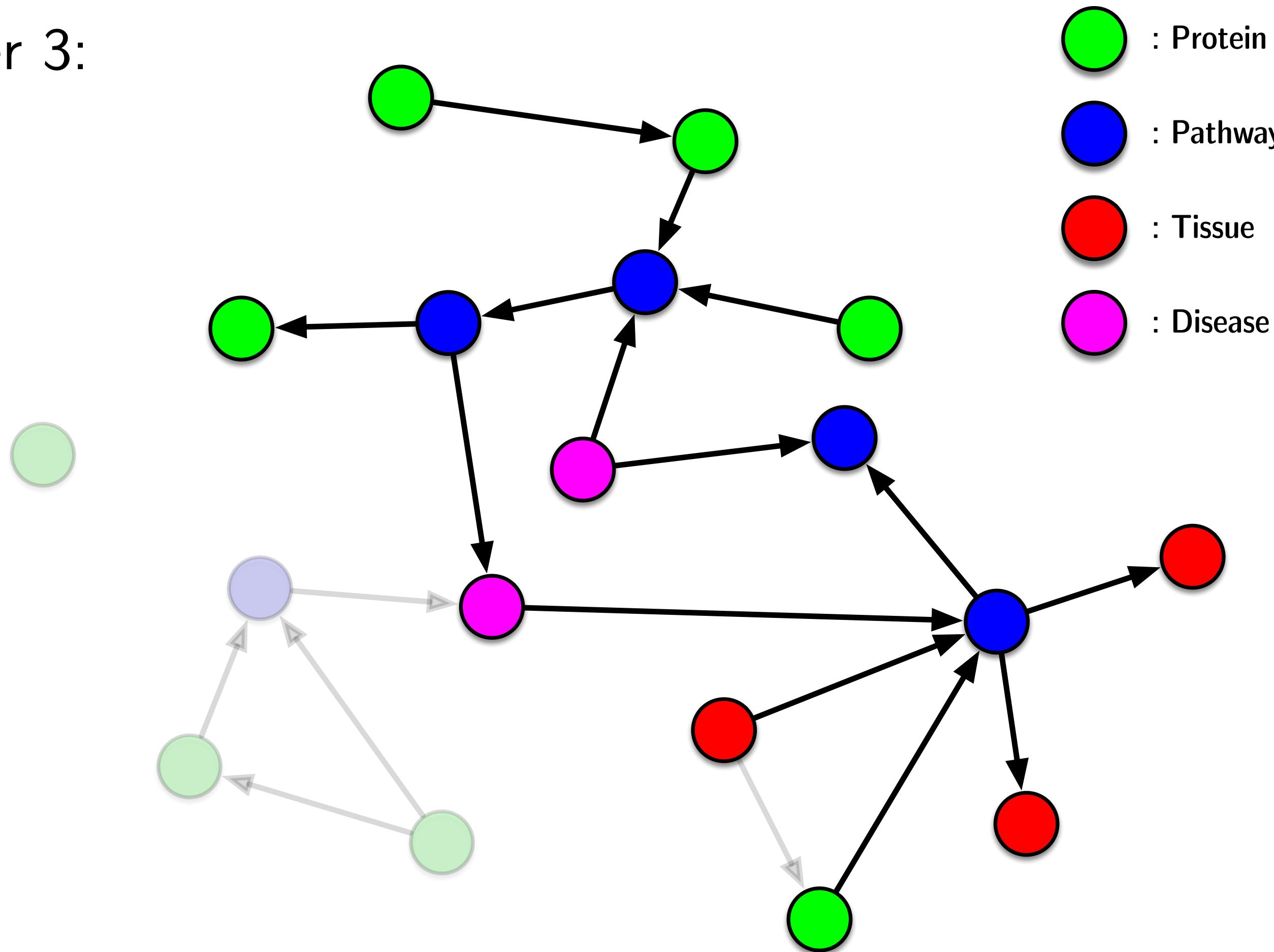
# Graph Neural Networks

- Layer 2:

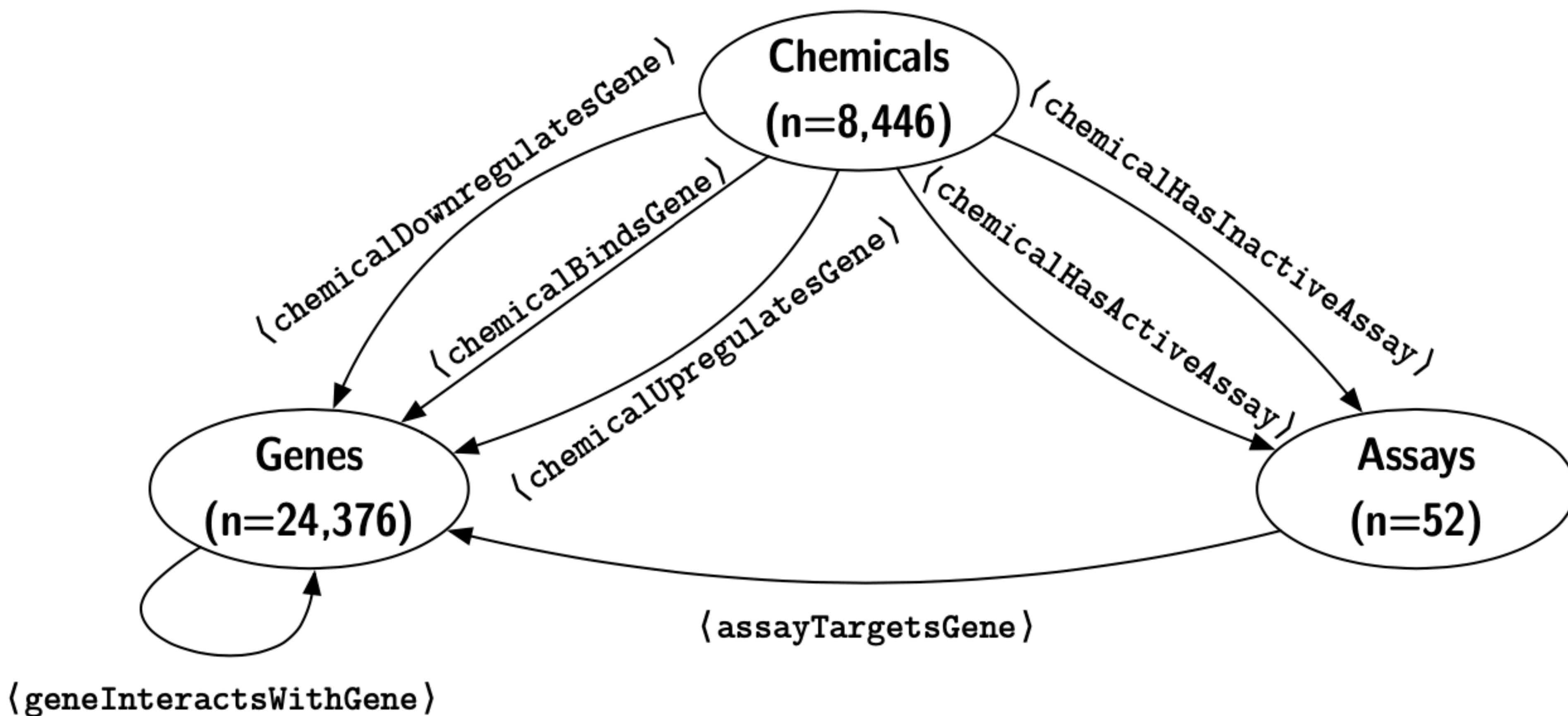


# Graph Neural Networks

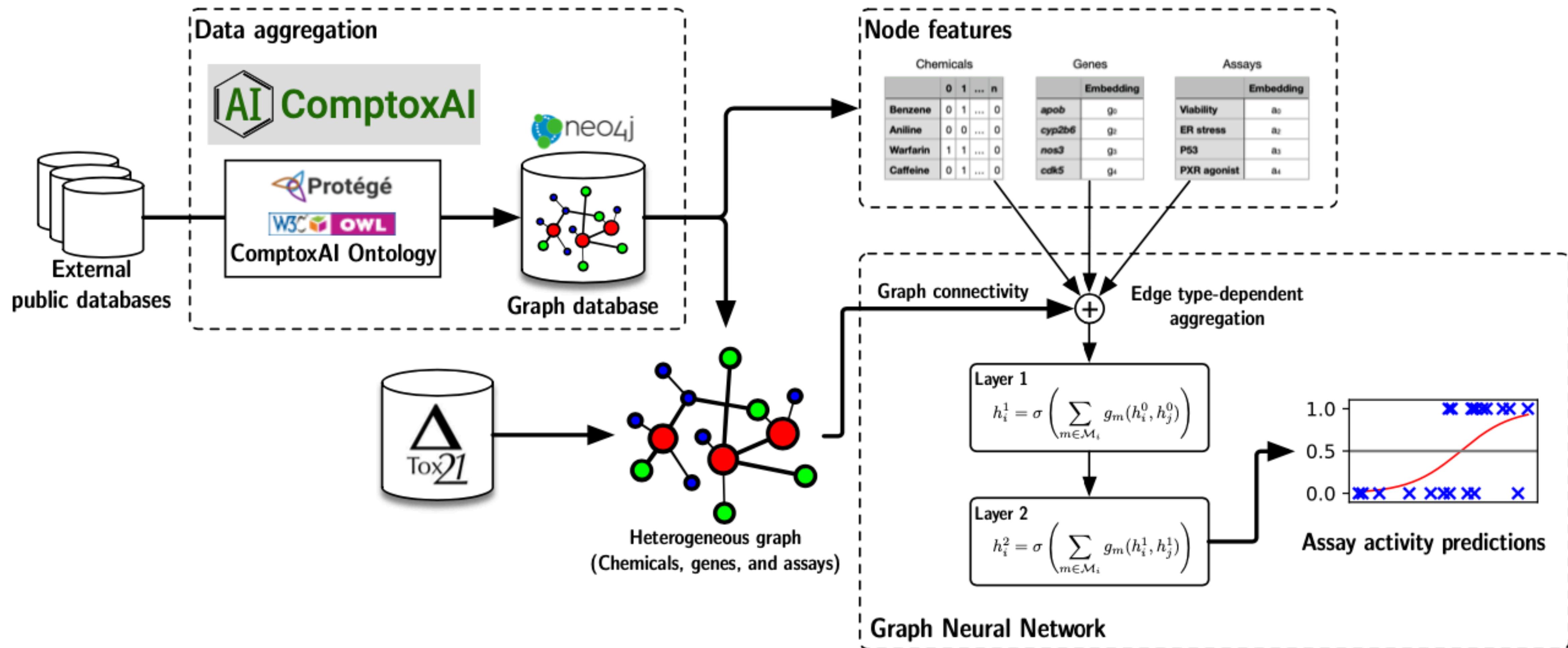
- Layer 3:

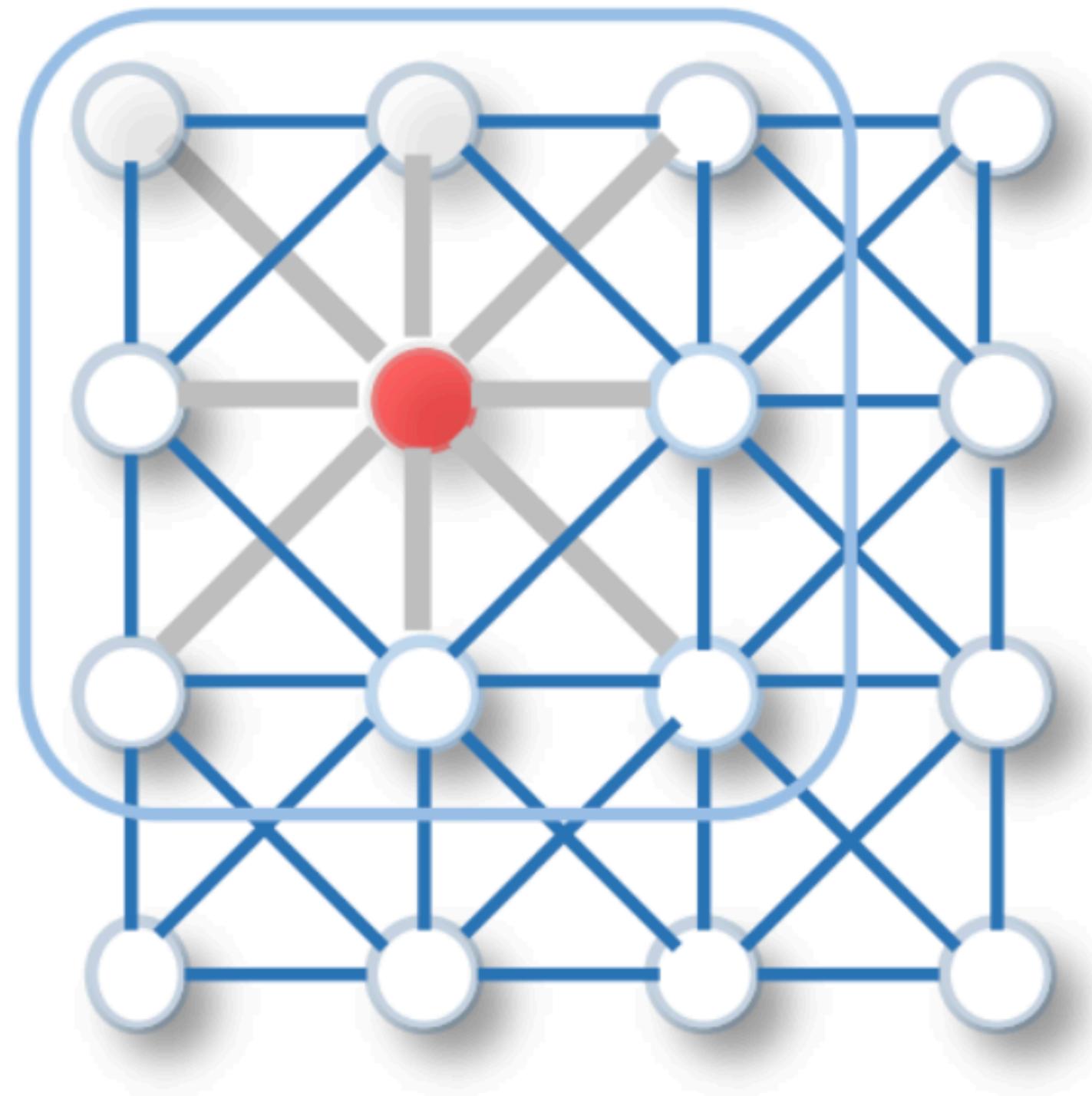


# QSAR Subgraph

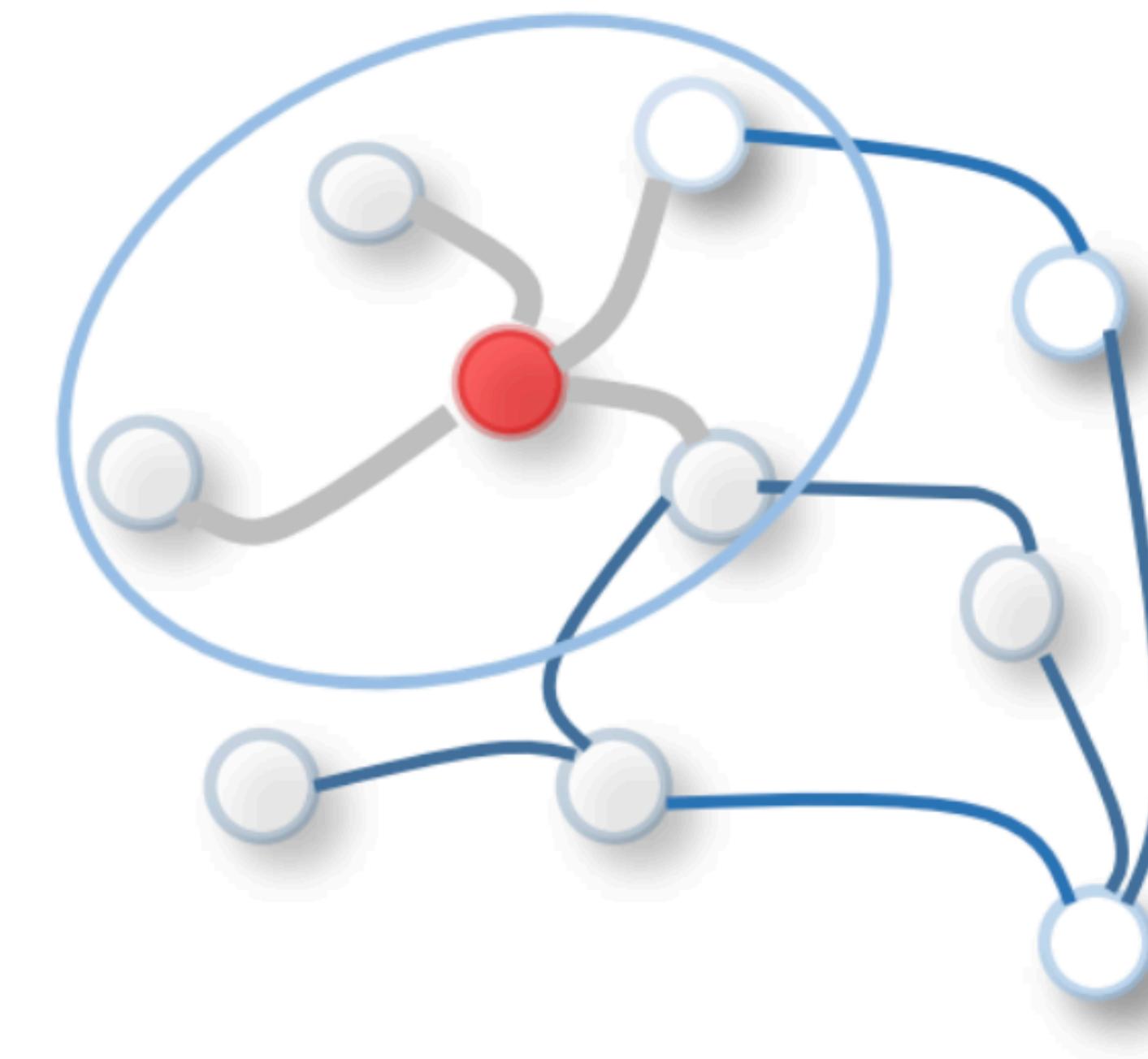


# Predicting chemical toxicity with graph neural networks



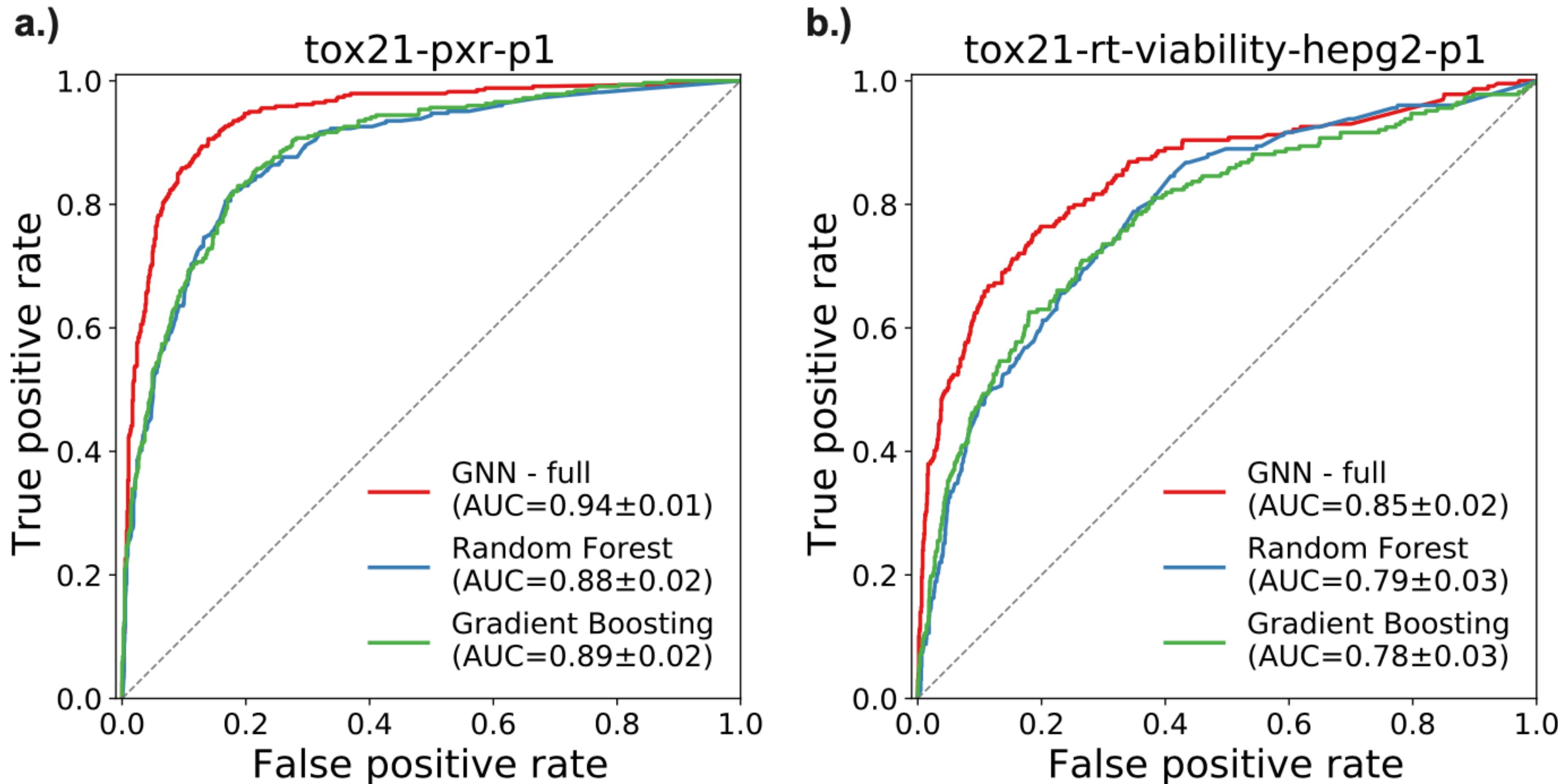


Convolutions on pixels (images)

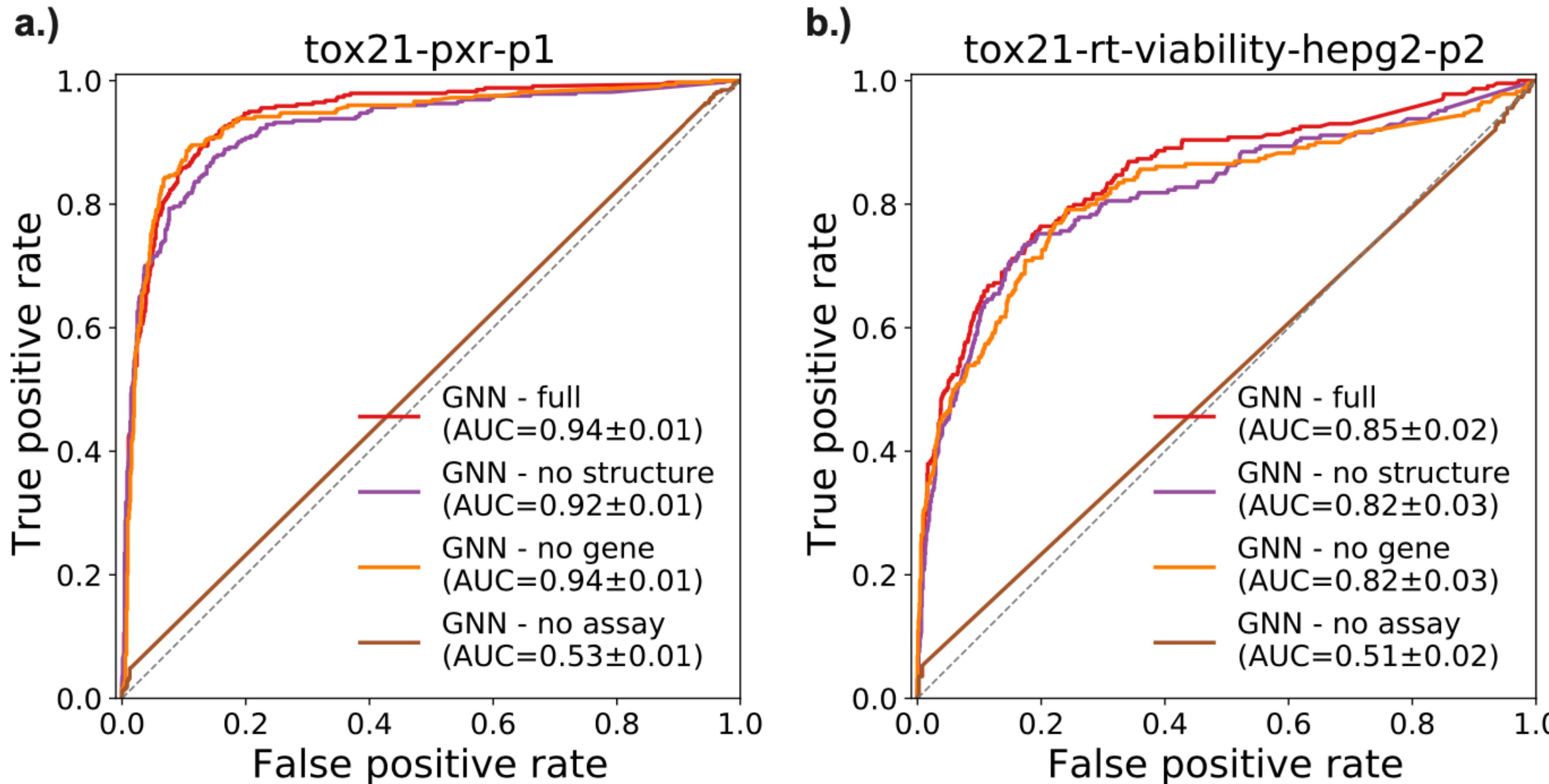


Convolutions on nodes (graphs)

# Graph neural network performance

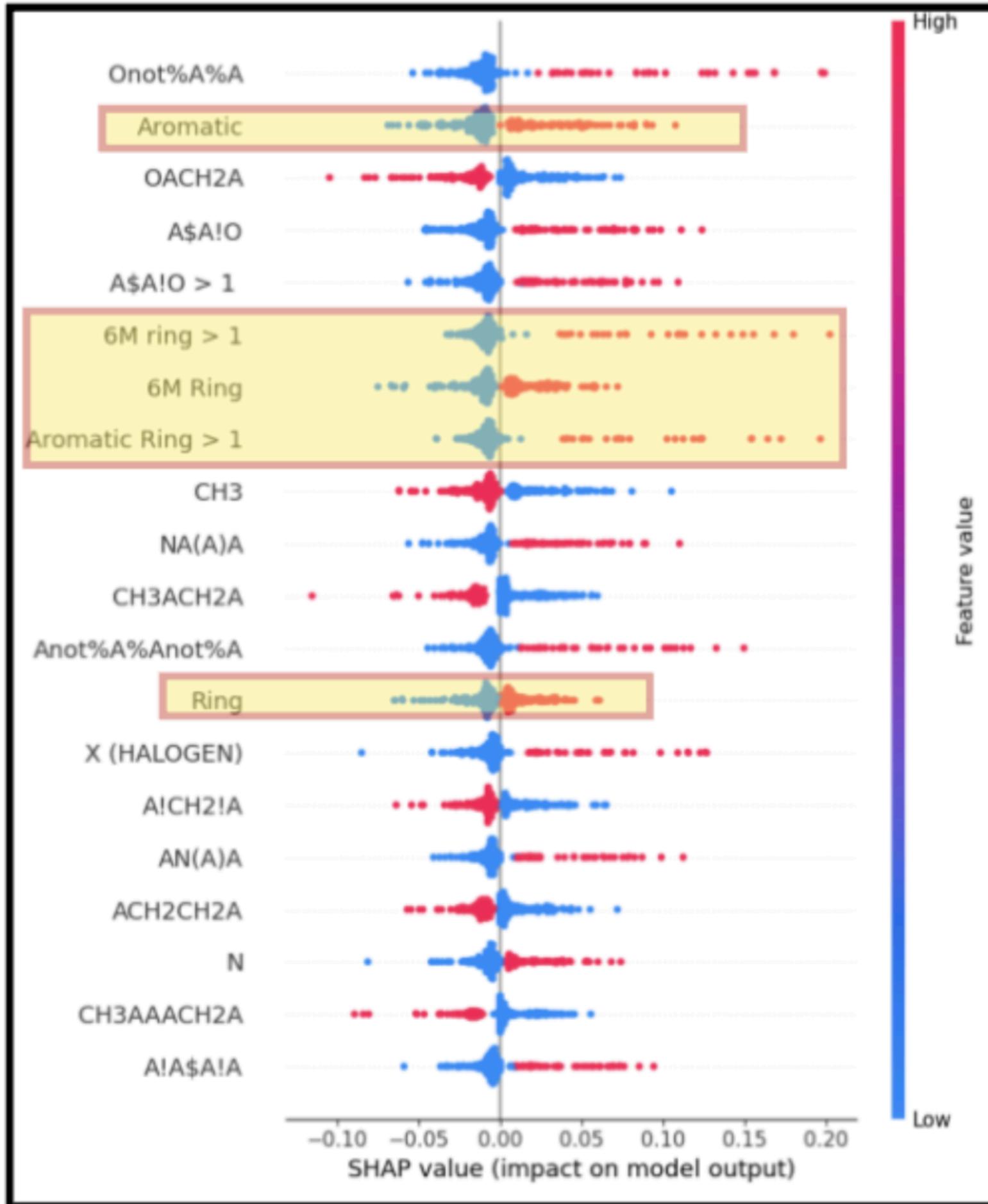


# Why do the GNNs perform so much better?

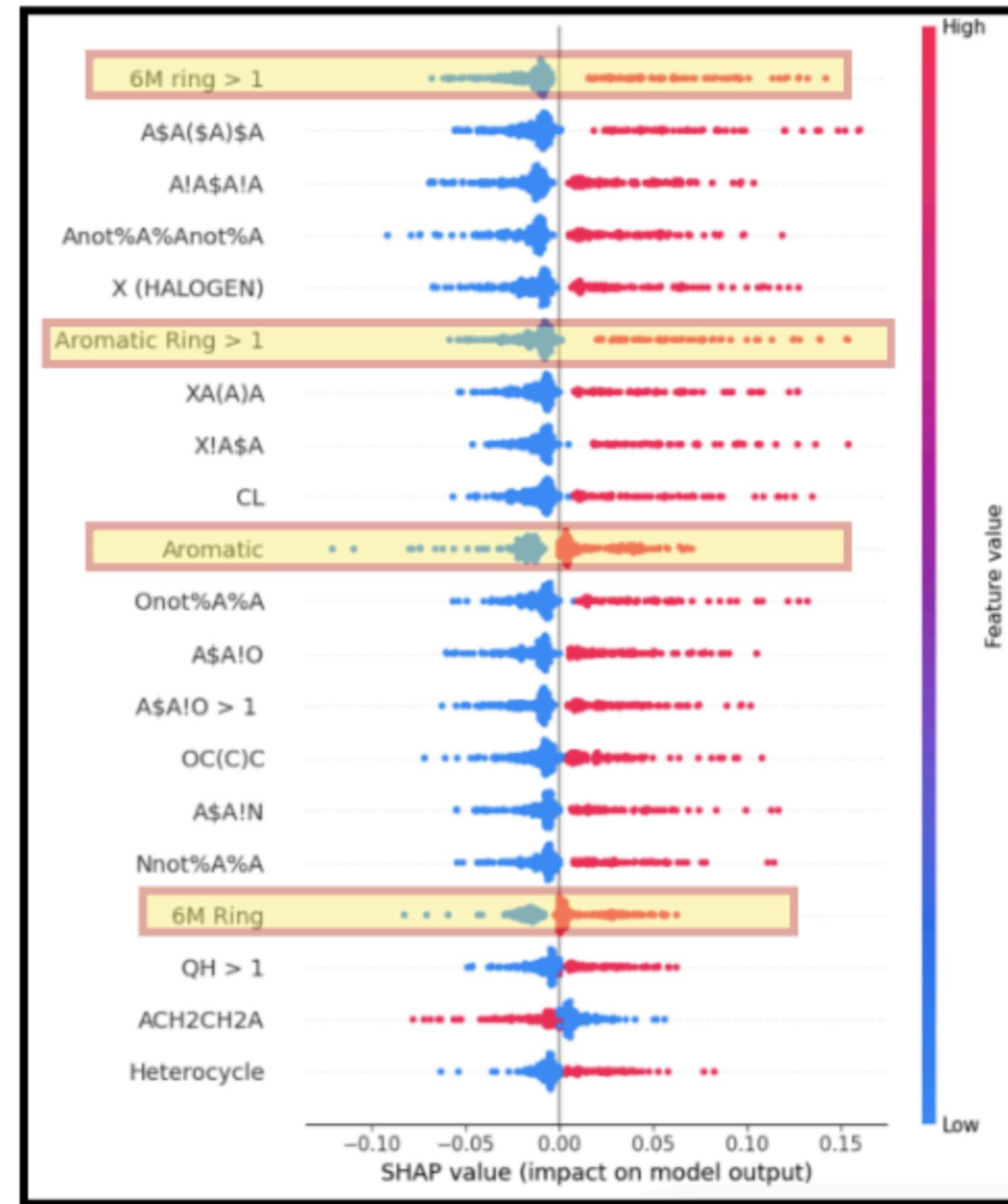


# Interpretable Predictive Toxicology

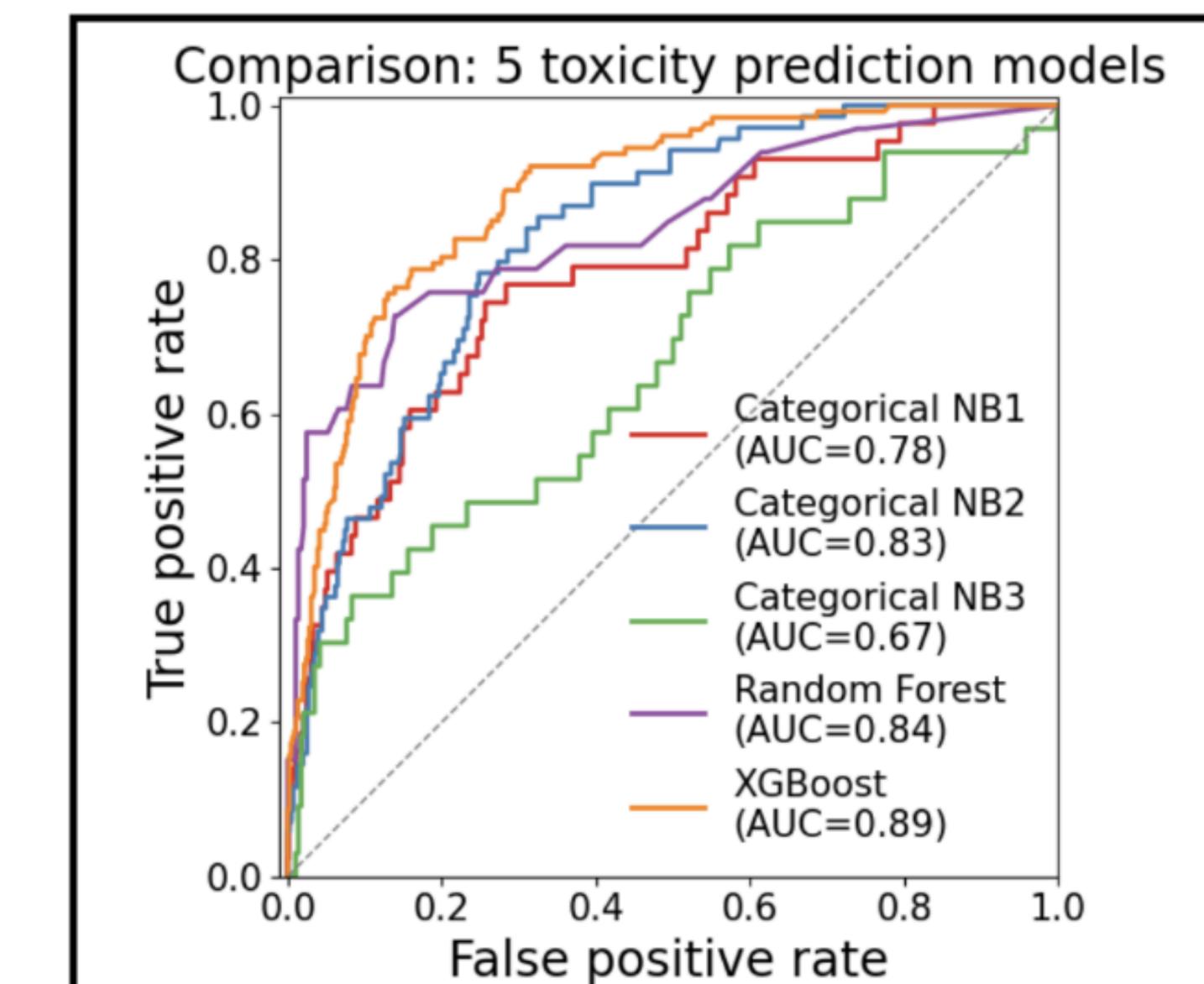
## Food Contact Chemicals Database with Liver Assay



## Hazardous Substances with Colon Cancer Assay

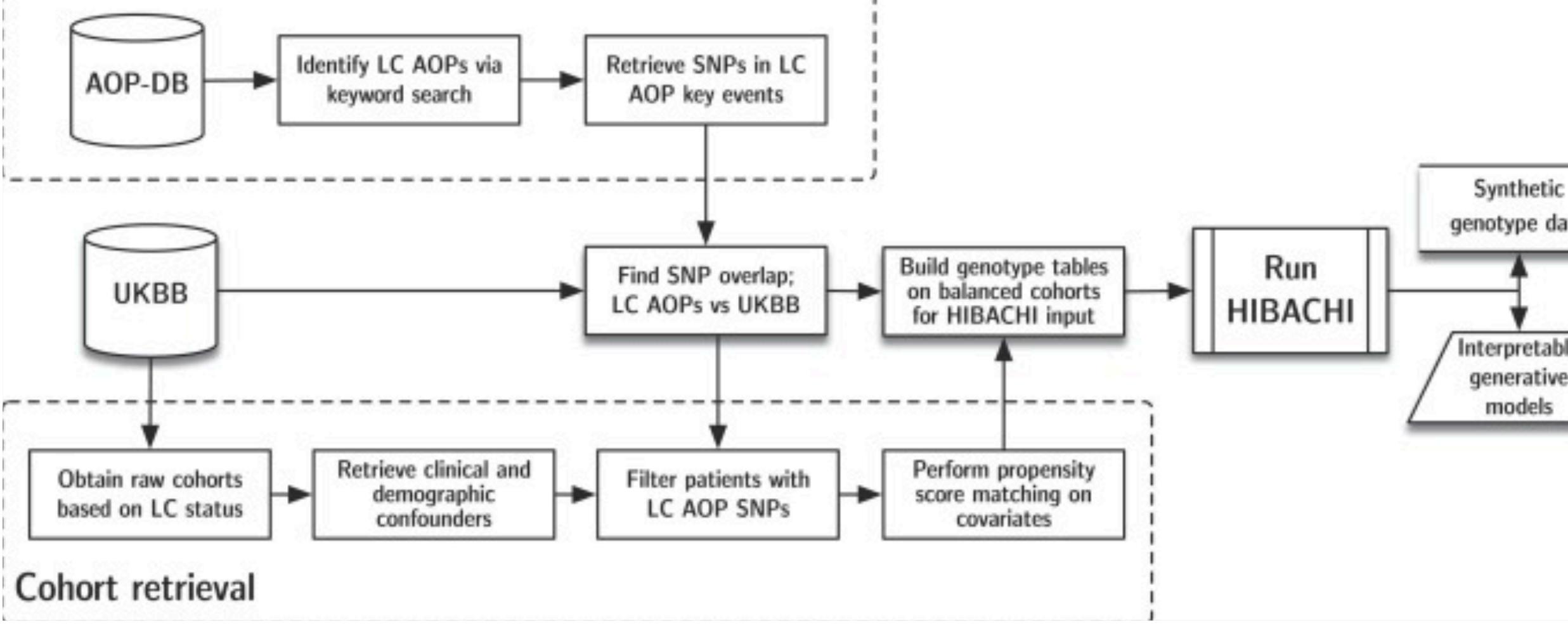


Estefania Morales  
Summer Undergraduate Internship Program (SUIP) Fellow

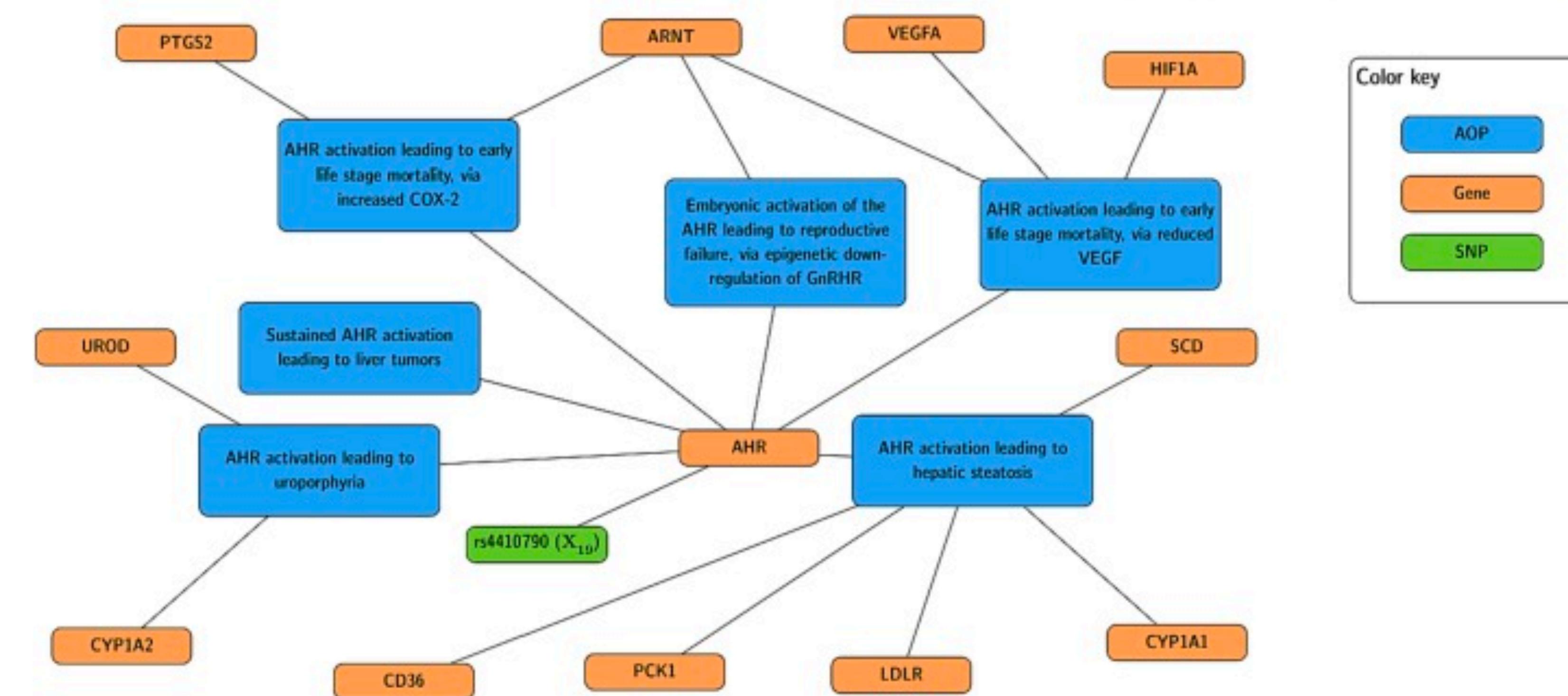


# **Observational data in exposomics research**

## SNP selection



## Cohort retrieval



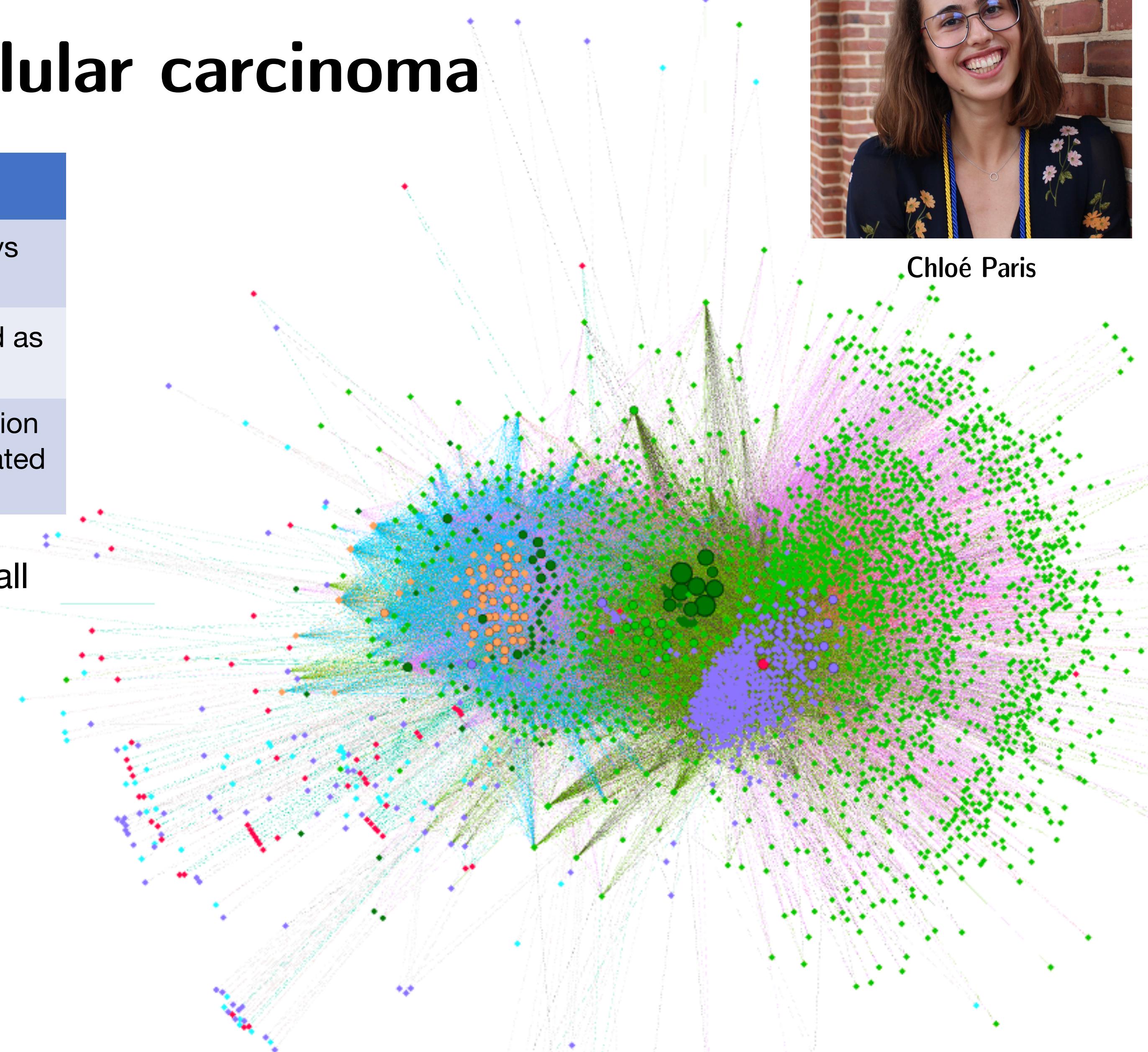
Individual	Model	Fitness Score
1	$((X_{10} \text{ XOR } X_7) \text{ AND } (X_{19} - (X_7 \neq ((\log_{X_{10}}(X_2) + X_{19}) \bmod 2))))$	2.730
2	$((X_{10} \text{ XOR } X_7) \text{ AND } (X_{19} - (X_7 \neq X_9)))$	2.280
3	$((X_{10} \text{ XOR } X_7) \text{ AND } (X_{19} - (X_7 \neq X_{24})))$	2.244
4	$((X_{10} \text{ XOR } X_{19}) \text{ AND } (X_{19} - X_7))$	2.204
5	$(X_{19} + (\neg X_6)) \bmod 2$	1.824
6	$X_9 \neq X_{19}$	1.254
7	$X_{10}!$	0.232

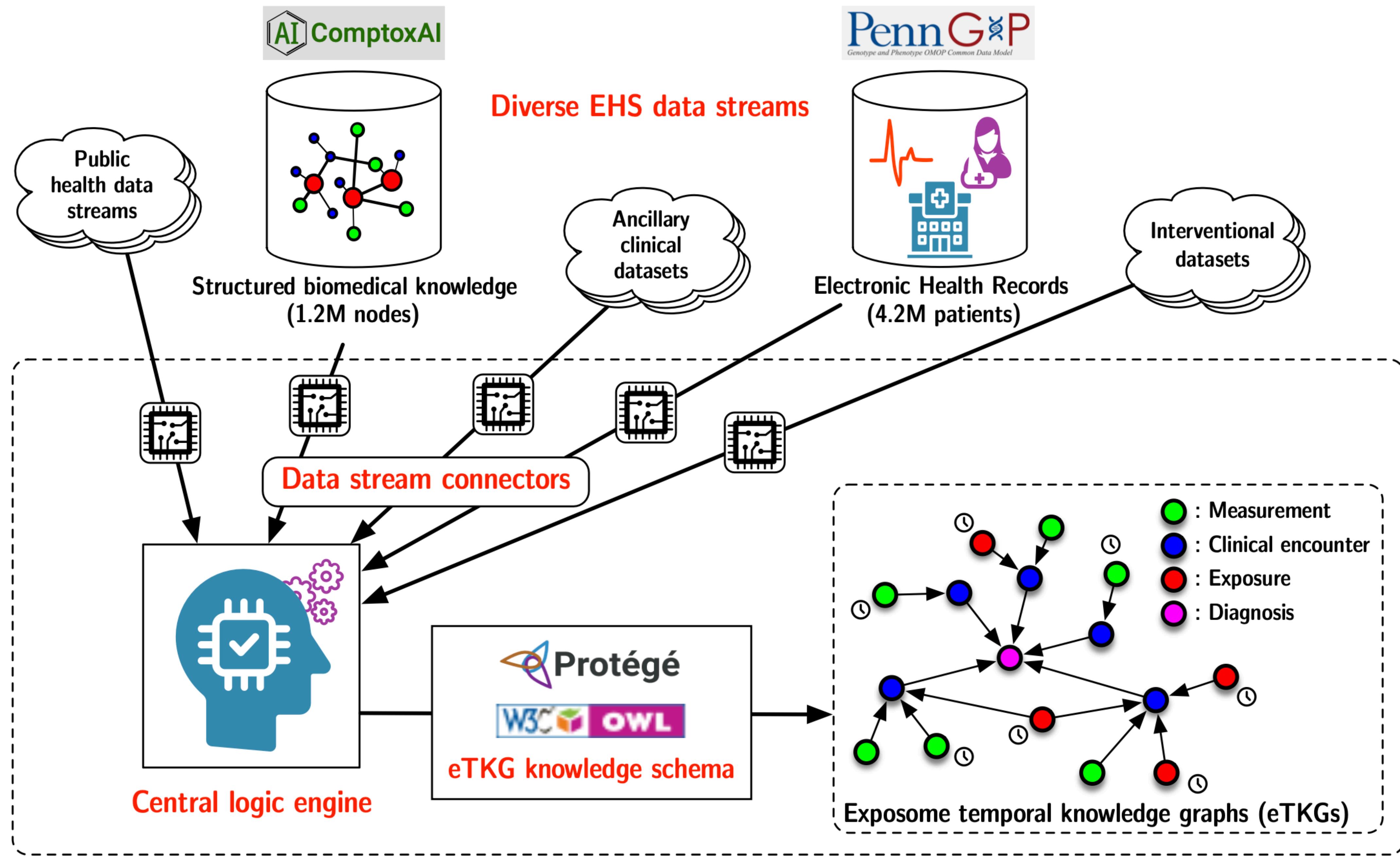
# GxE interactions in PFAS-induced hepatocellular carcinoma



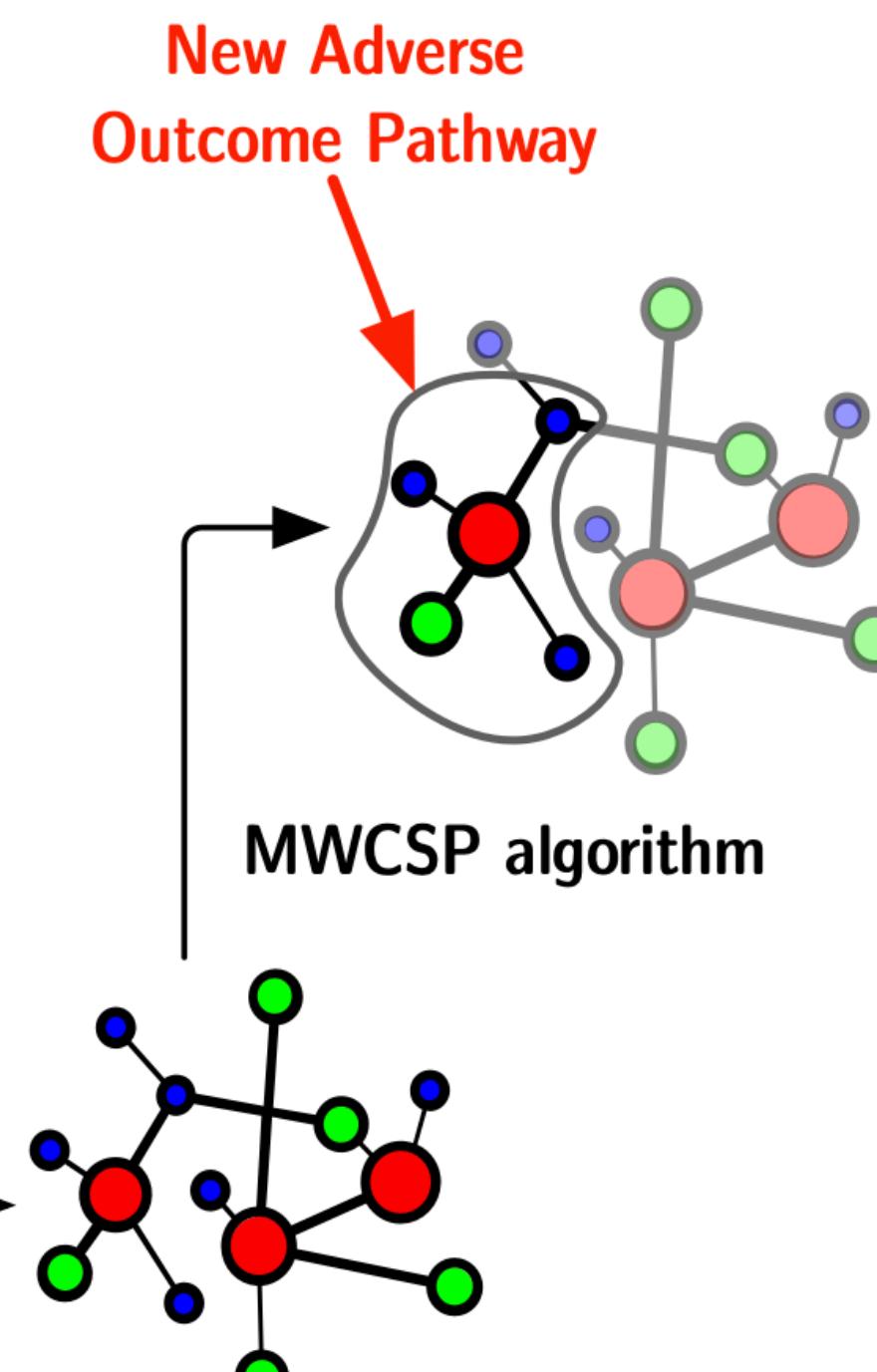
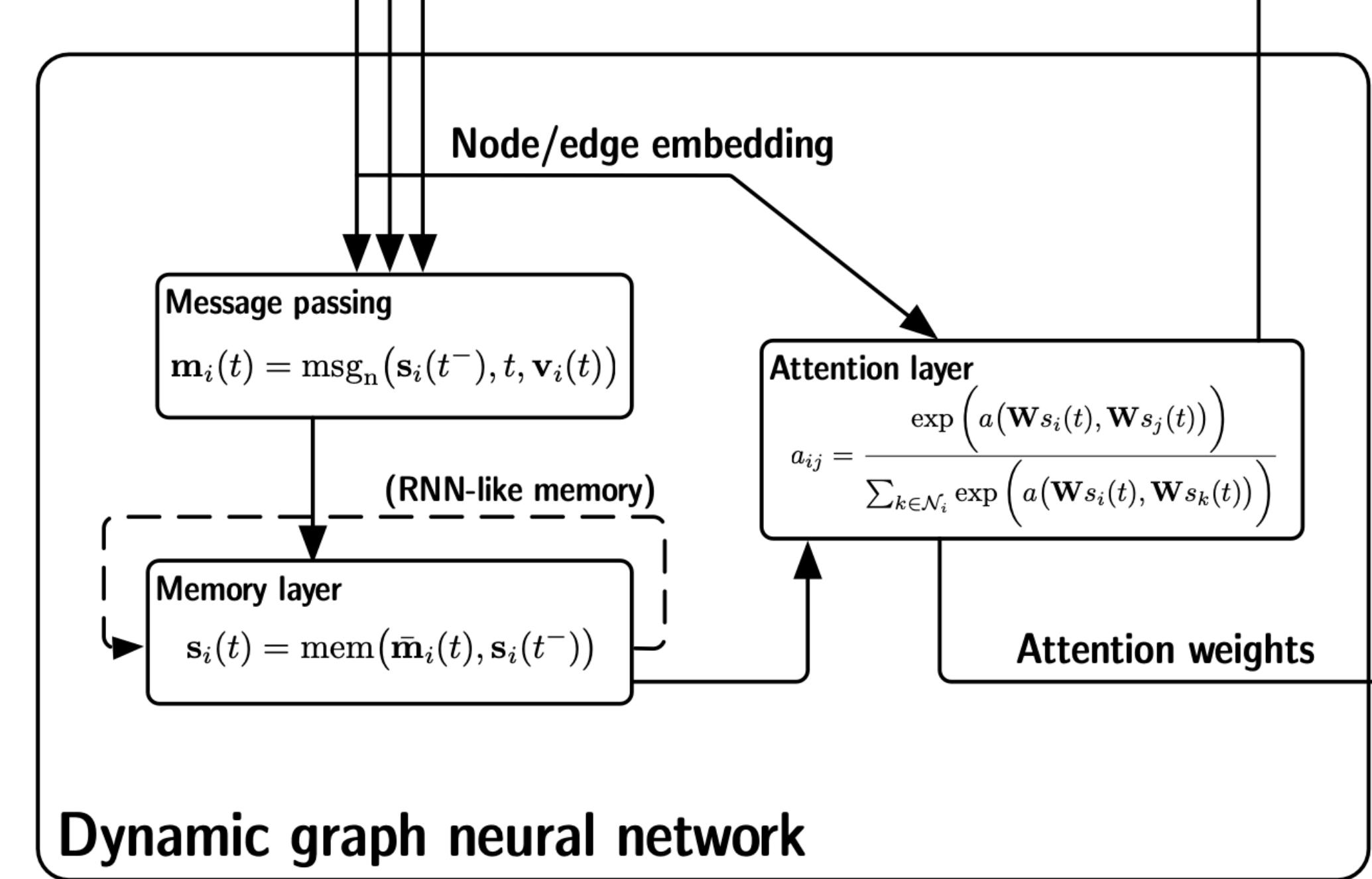
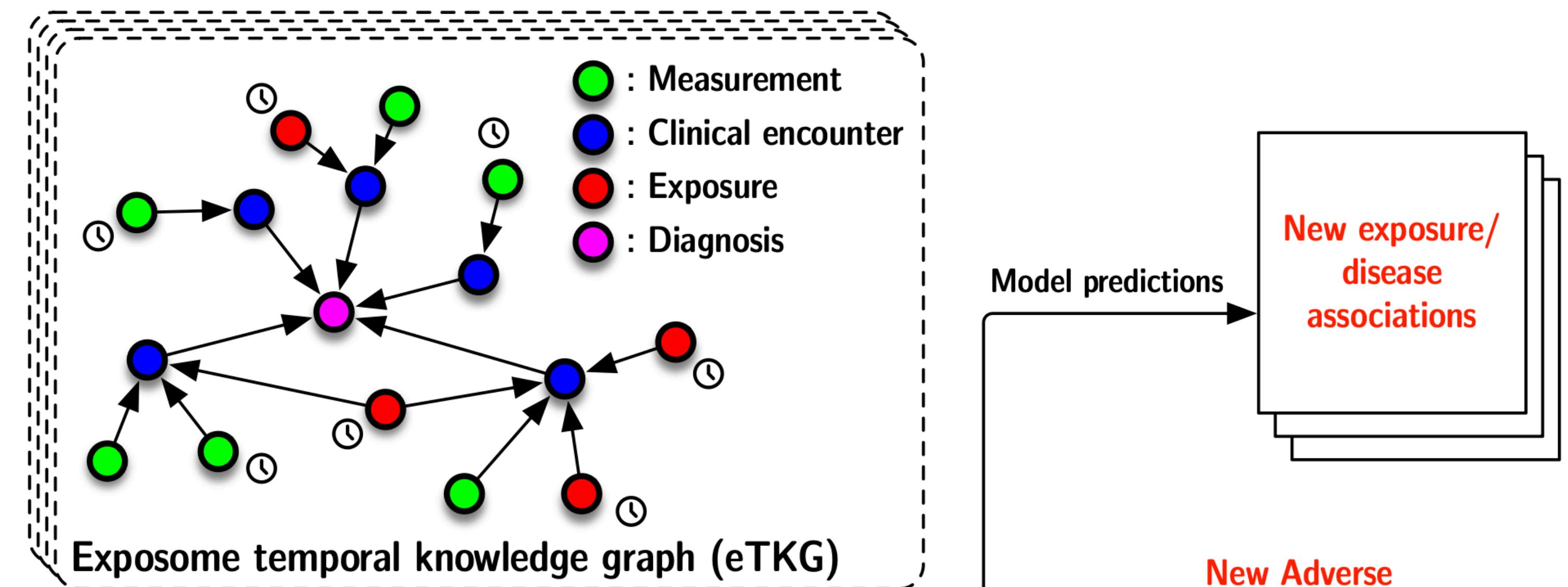
Gene Symbol	Function
<b>TP53</b>	Most common mutation found in HCC (plays important role in genomic stability)
<b>DNAJB1</b>	Fusion of DNAJB1 and PRKACA suggested as a possible marker for FLC
<b>HMOX1</b>	Suggested to play a role in HCC cell migration + its <b>metabolite bilirubin</b> has been associated with HCC

- Calculate the base **communicability score** between all PFAS-HCC node pairs in subgraph
  - Defined as a weighted sum of all possible “walks” linking a pair of nodes
- For each gene in the network:
  - Remove the gene from the network
  - Recalculate communicability scores for all PFAS-HCC pairs
  - Determine the change in communicability score
- Identify genes yielding the greatest mean change in communicability





# Patient-level eTKGs



# Acknowledgements

- The Romano Lab:
  - Chloé Paris
  - Estefania Morales
  - Kevin Shen
  - Tom Pan
  - Jesse Goodrich, PhD (USC)
  - Danielle Mowery, PhD
  - Trevor Penning, PhD
- Li Shen, PhD
- Jason Moore, PhD (Cedars-Sinai)
- Yun Hao, PhD (Simons foundation)
- Grant support:
  - K99/R00-LM013646 (PI: Romano)
  - Colton Informatics Fellowship
  - P30-ES013508 (PI: Penning)

DEPARTMENT of  
**BIOSTATISTICS**  
**EPIDEMIOLOGY &**  
**INFORMATICS**

Center of Excellence in  
Environmental Toxicology



Penn Institute for  
Biomedical Informatics

See more at: [romanolab.org](http://romanolab.org)