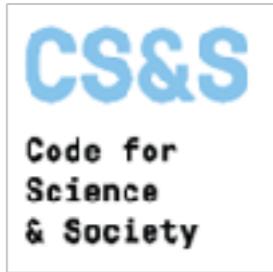




Bringing data science and AI/ML tools to infectious disease research

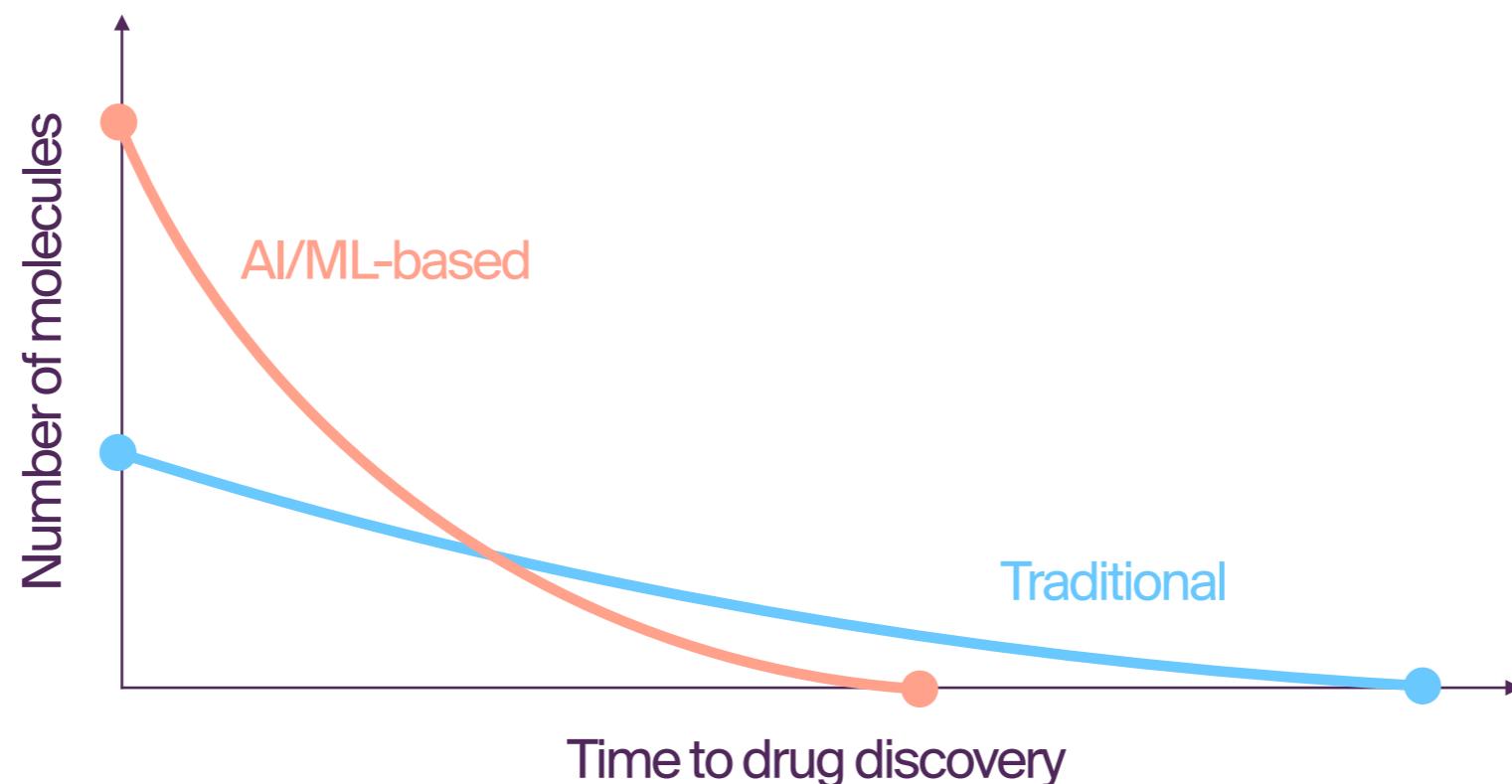
Session 2: Introduction to Supervised Machine Learning

Event Sponsors

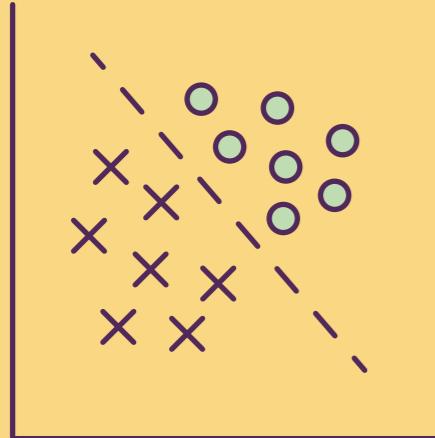


Why AI/ML for drug discovery?

- AI/ML
 - Cost-effective**
 - Fast implementation**
 - Remote**
 - Domain agnostic**
- Small molecule drugs
 - Track record**
 - Neglected diseases**
 - Cheap medications**
 - Inspiring initiatives: MMV, DNDi, OSM...**



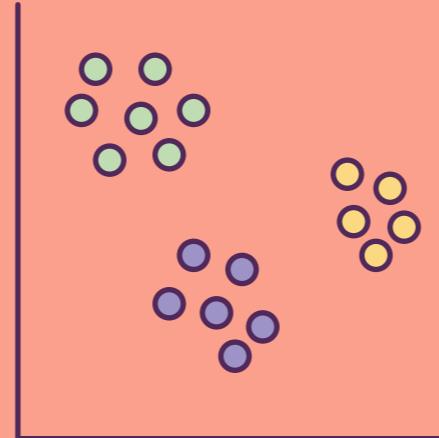
Supervised



Labeled data
Classification
Regression

(Today)

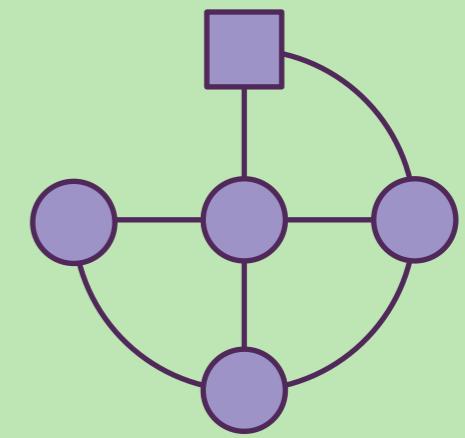
Unsupervised



Unlabeled data
Clustering
2D projection
Similarity search

(Yesterday)

Reinforcement

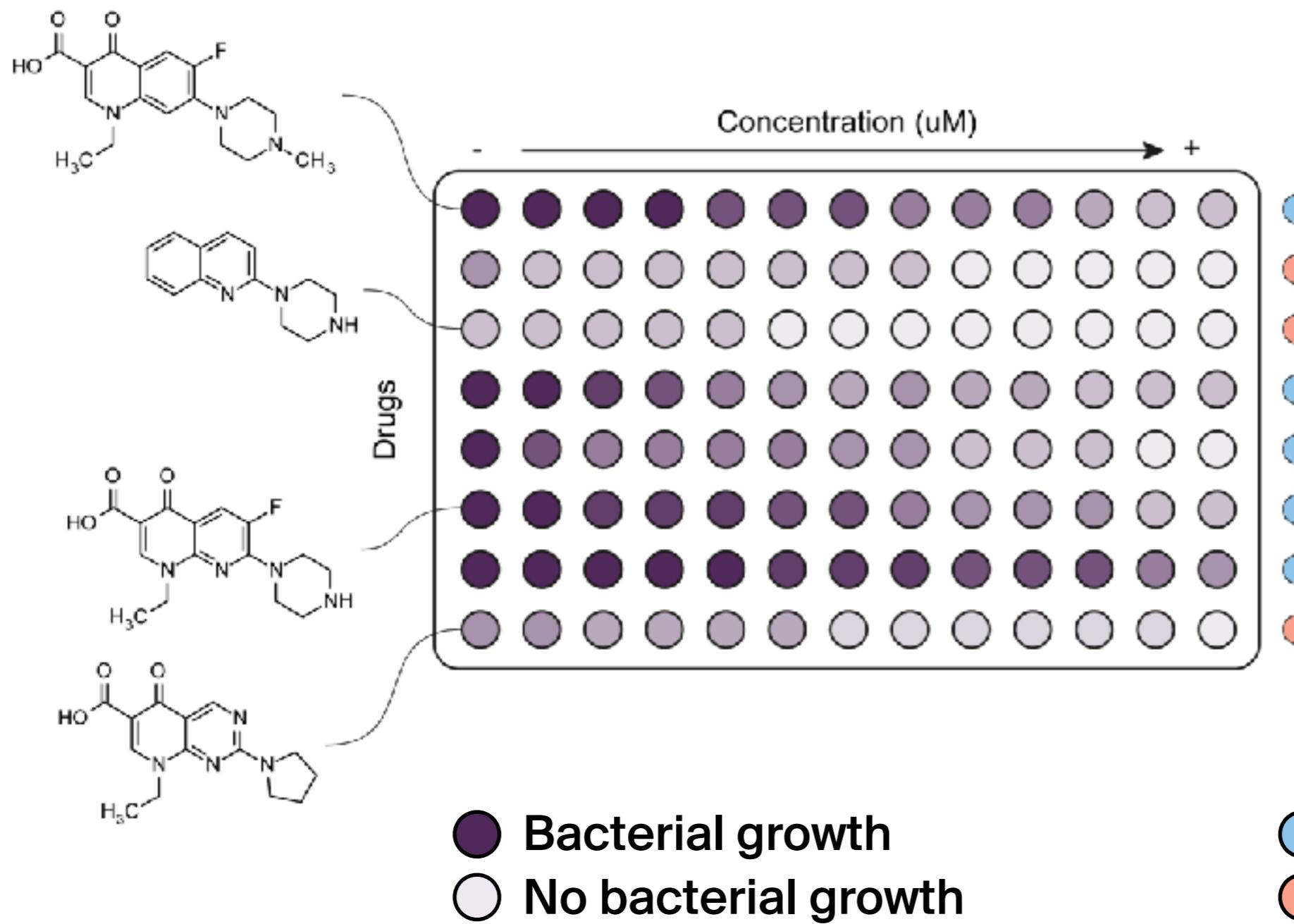


Interaction with
environment / agent
Generative models

(On Friday)

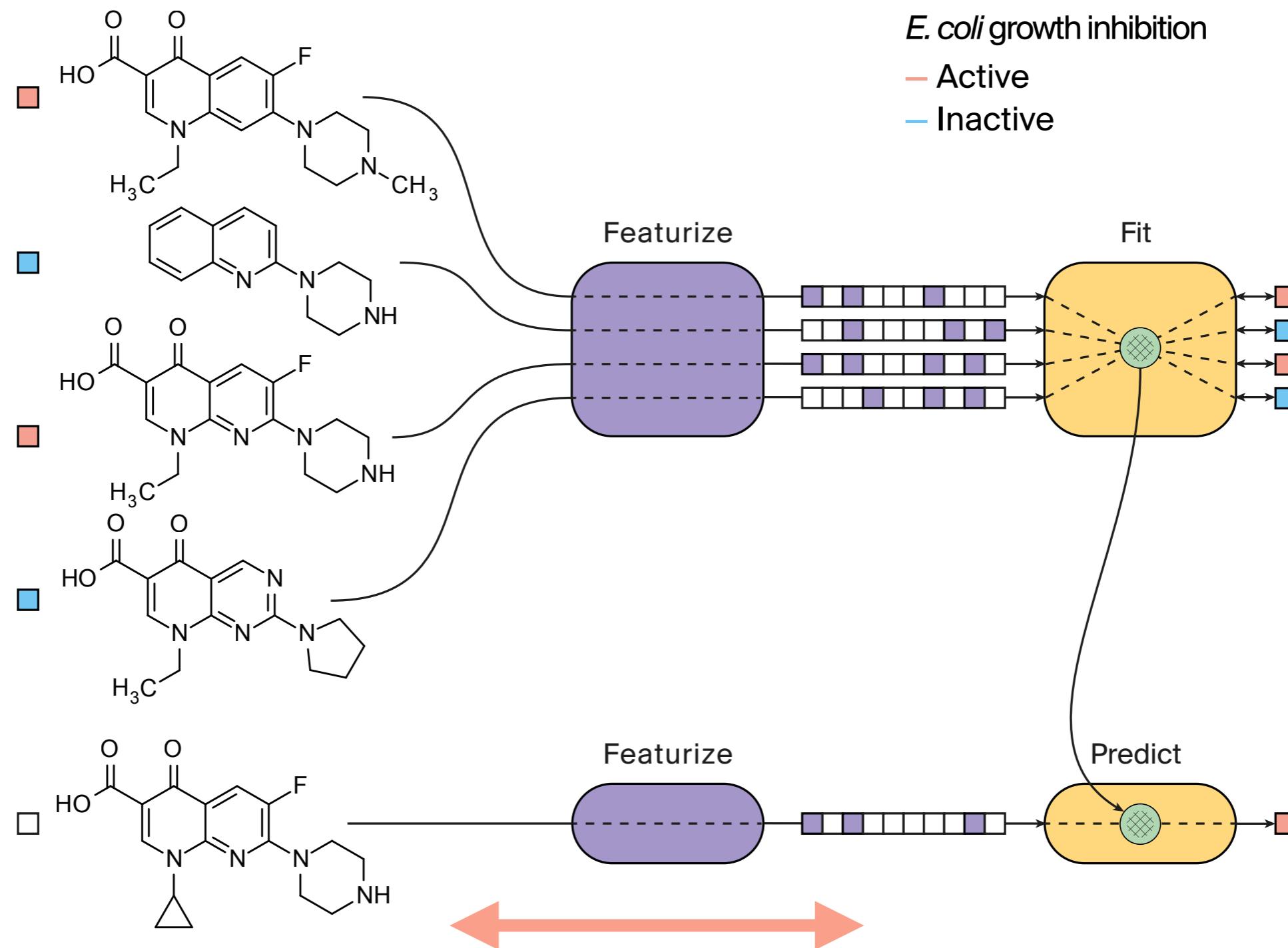
A simple **supervised** task in drug discovery (aka QSAR)

Antibiotic activity experiments (*E. coli*)



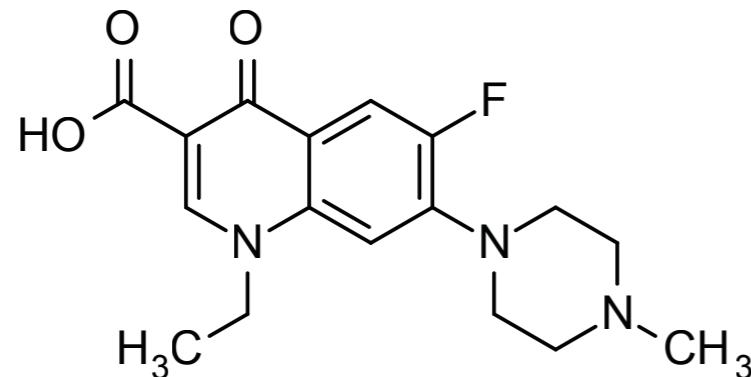
A simple **supervised** task in drug discovery (aka QSAR)

Antibiotic activity experiments (*E. coli*)

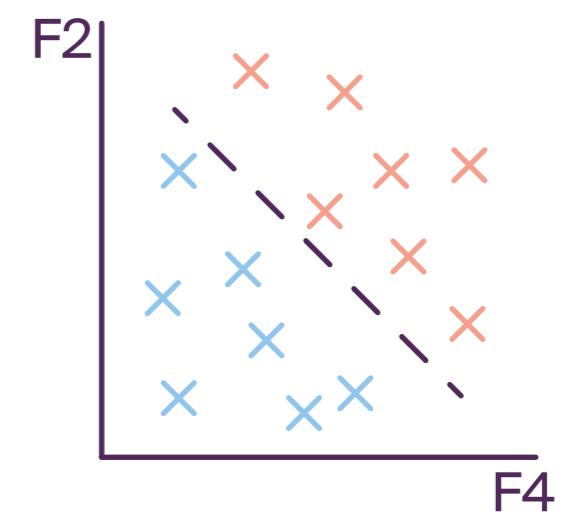
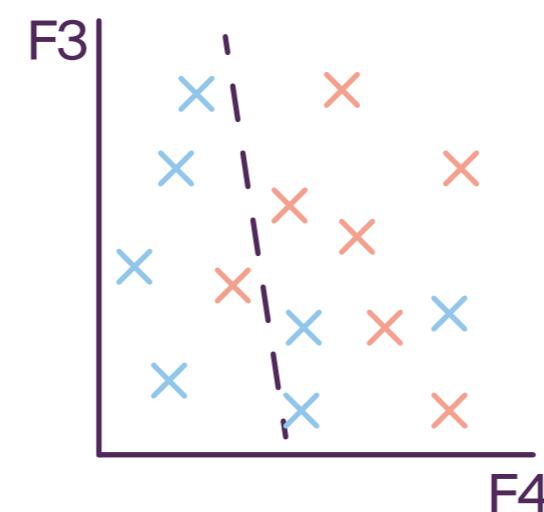
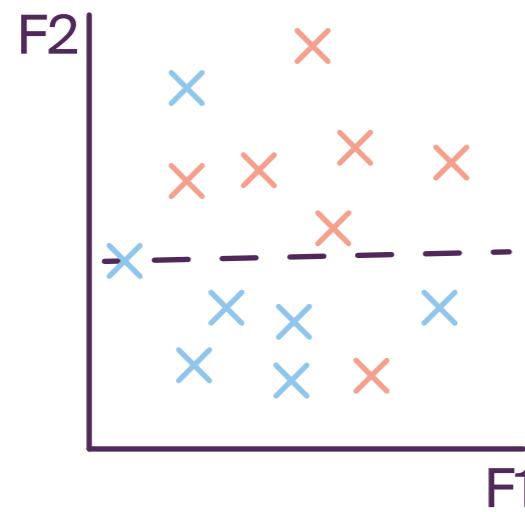


Featurization

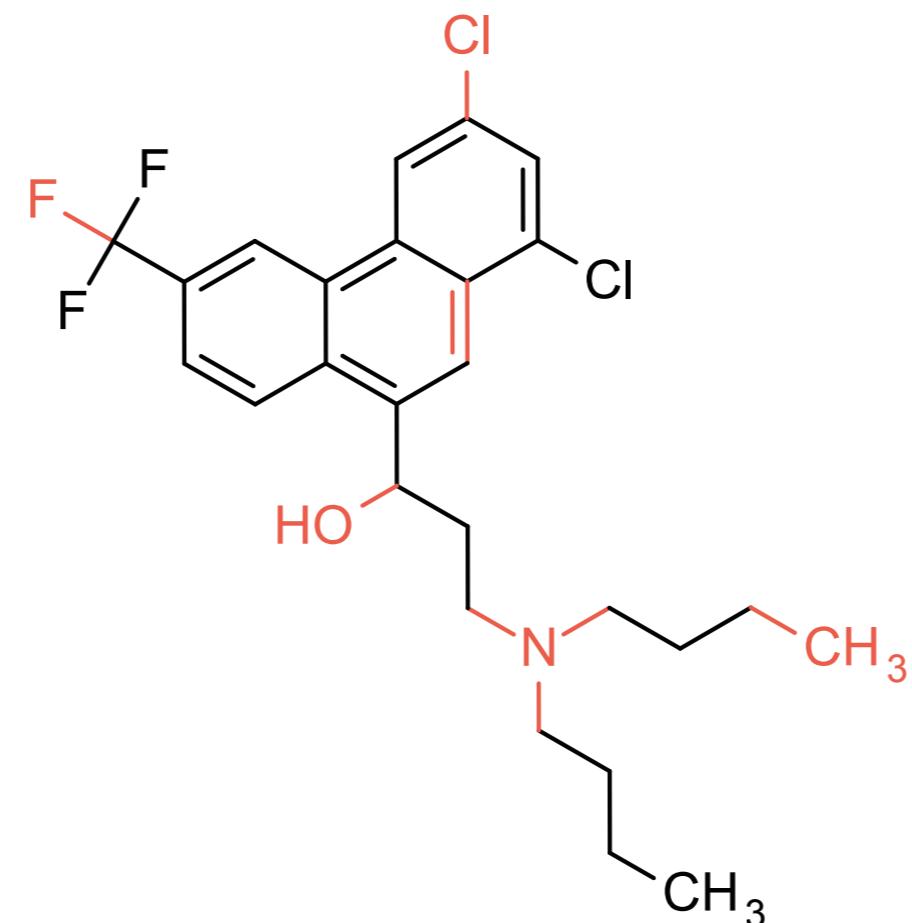
- Choose and design features that are discriminative



- F1: number of rings
- F2: solubility
- F3: number of fluors
- F4: molecular weight



Classical chemical fingerprints



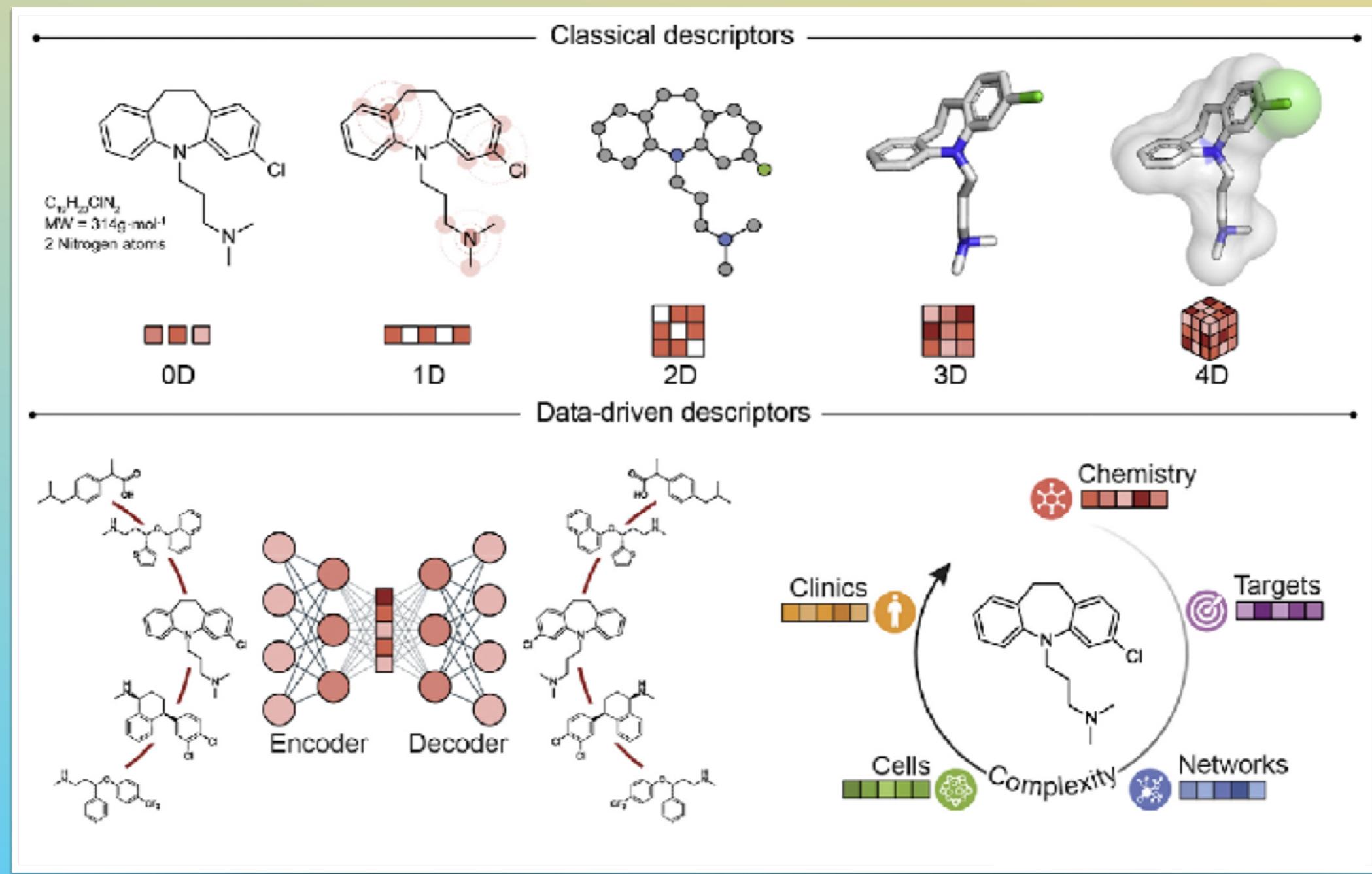
Duran-Frigola et al. Nature Communications, 2014

Chemical Checker signature

Halofantrine belongs to the class of organic compounds known as **phenanthrenes** and derivatives. These are polycyclic compounds containing a phenanthrene moiety, which is a tricyclic aromatic compound with three non-linearly fused benzene. Halofantrine is a synthetic **antimalarial** which acts as a **blood schizonticide**. It is effective against multi drug resistant (including mefloquine resistant) *P. falciparum* malaria. The mechanism of action of Halofantrine may be similar to that of chloroquine, quinine, and mefloquine; by forming toxic **complexes with ferritoporphyrin IX** that damage the membrane of the parasite. It appears to inhibit polymerisation of heme molecules (by the parasite enzyme '**heme polymerase**'), resulting in the parasite being poisoned by its own waste. Halofantrine has been shown to preferentially block open and inactivated **HERG channels** leading to some degree of **cardiotoxicity**. Side effects include coughing noisy, rattling, troubled breathing, loss of appetite, aches and pain in joints, indigestion, and **skin itching** or rash, *et cetera, et cetera*.

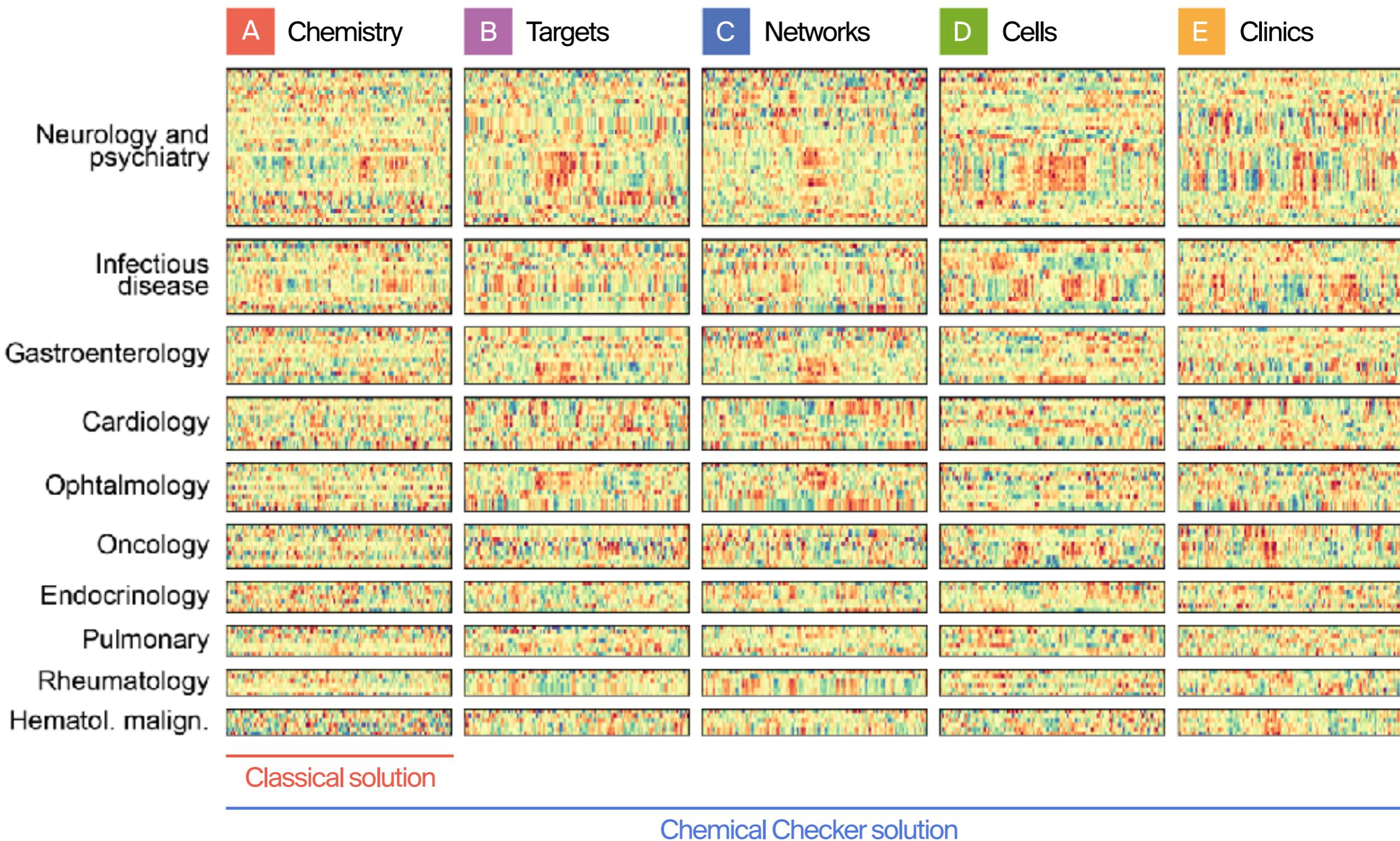


Chemical descriptors

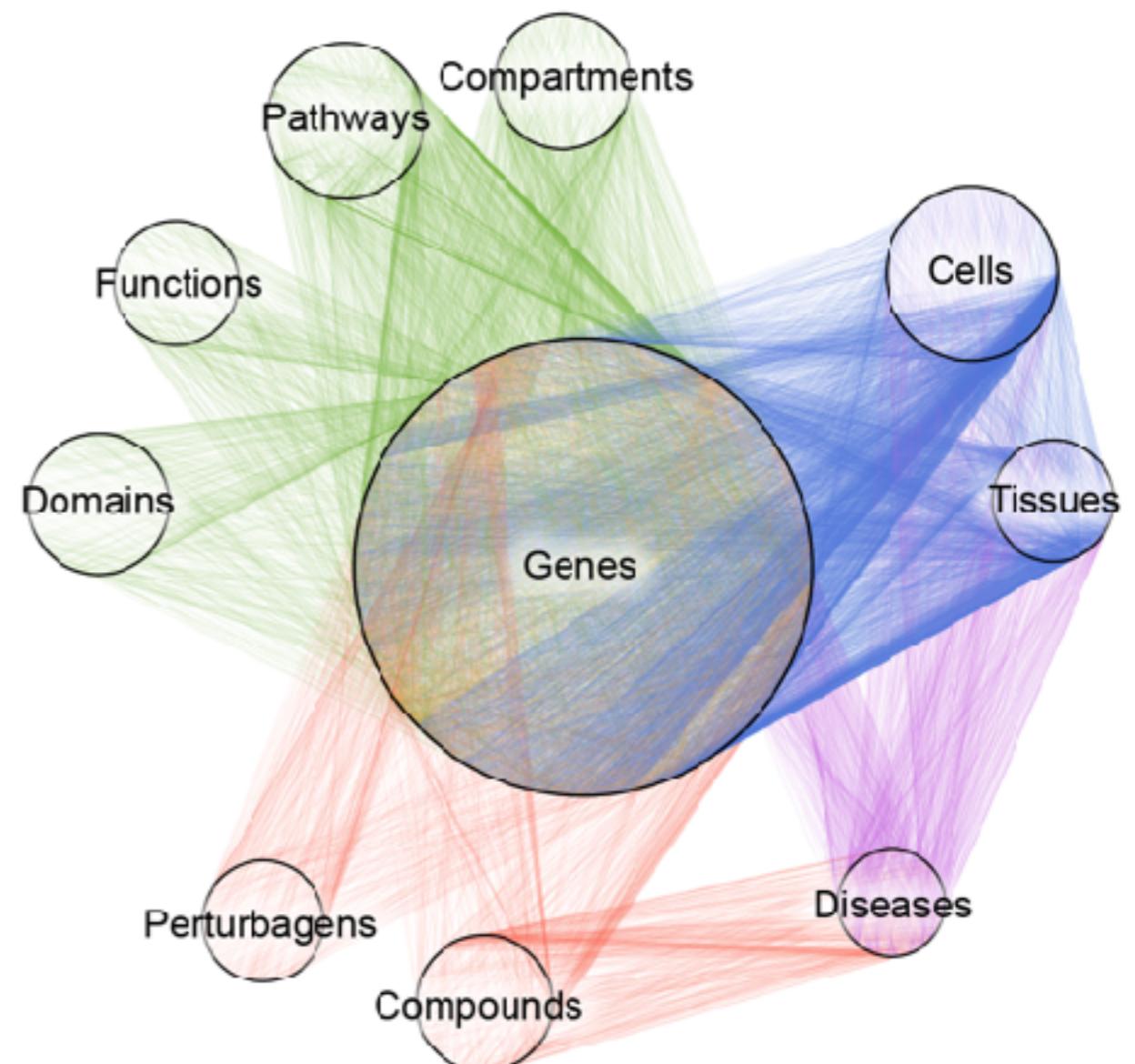
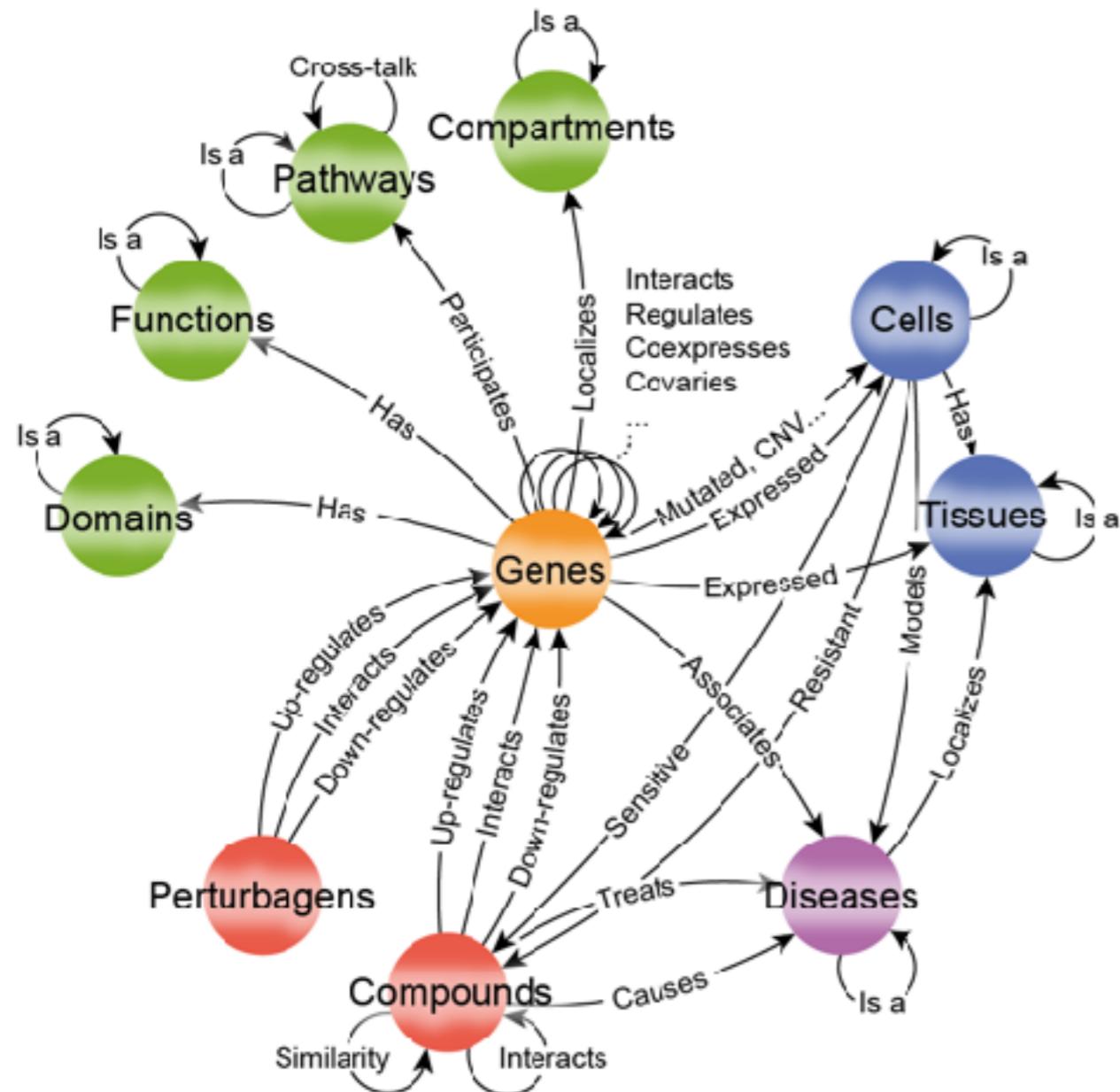


Chemical Checker signatures

Duran-Frigola et al, Nature Biotechnology, 2020

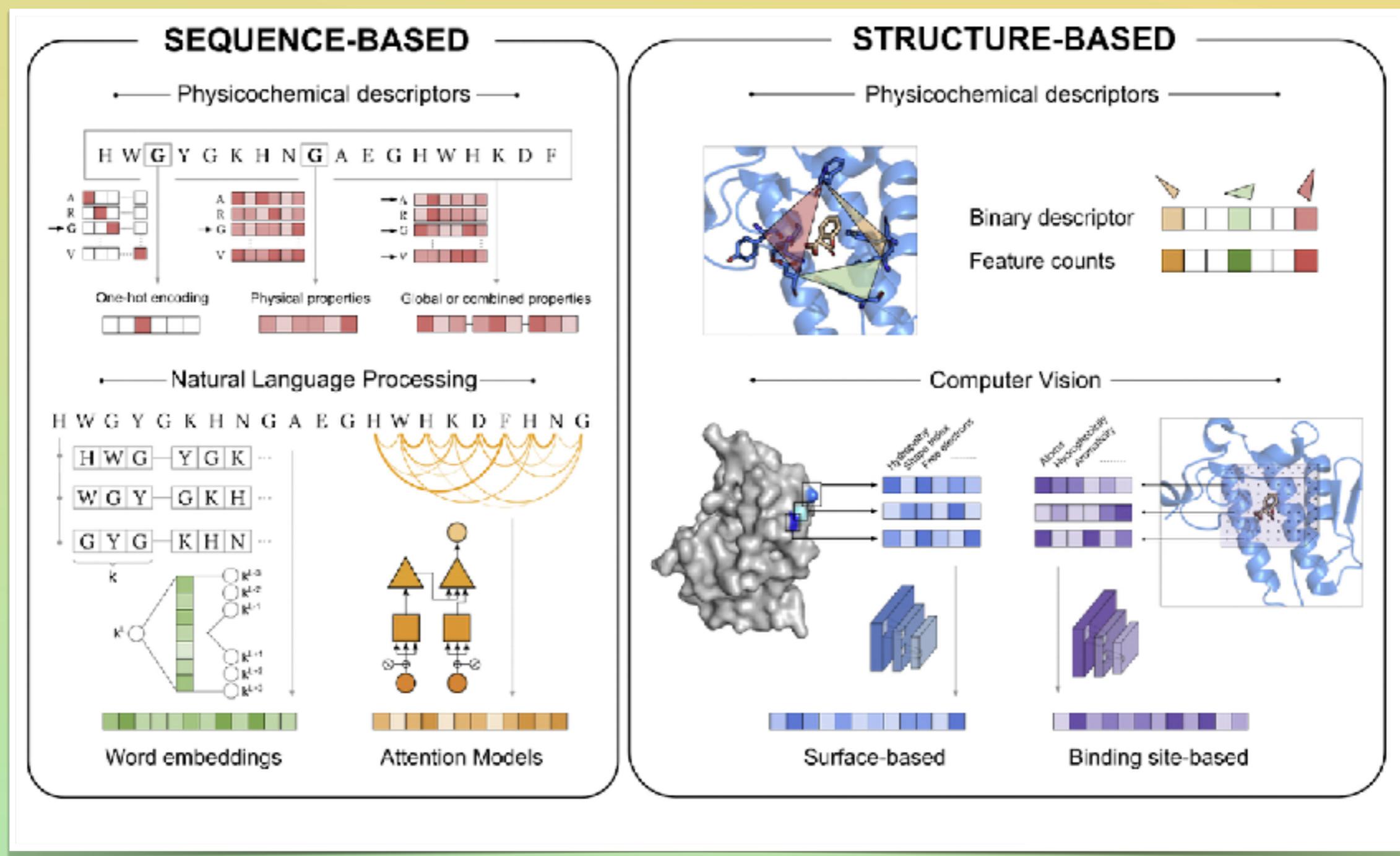


Biomedical knowledge graph

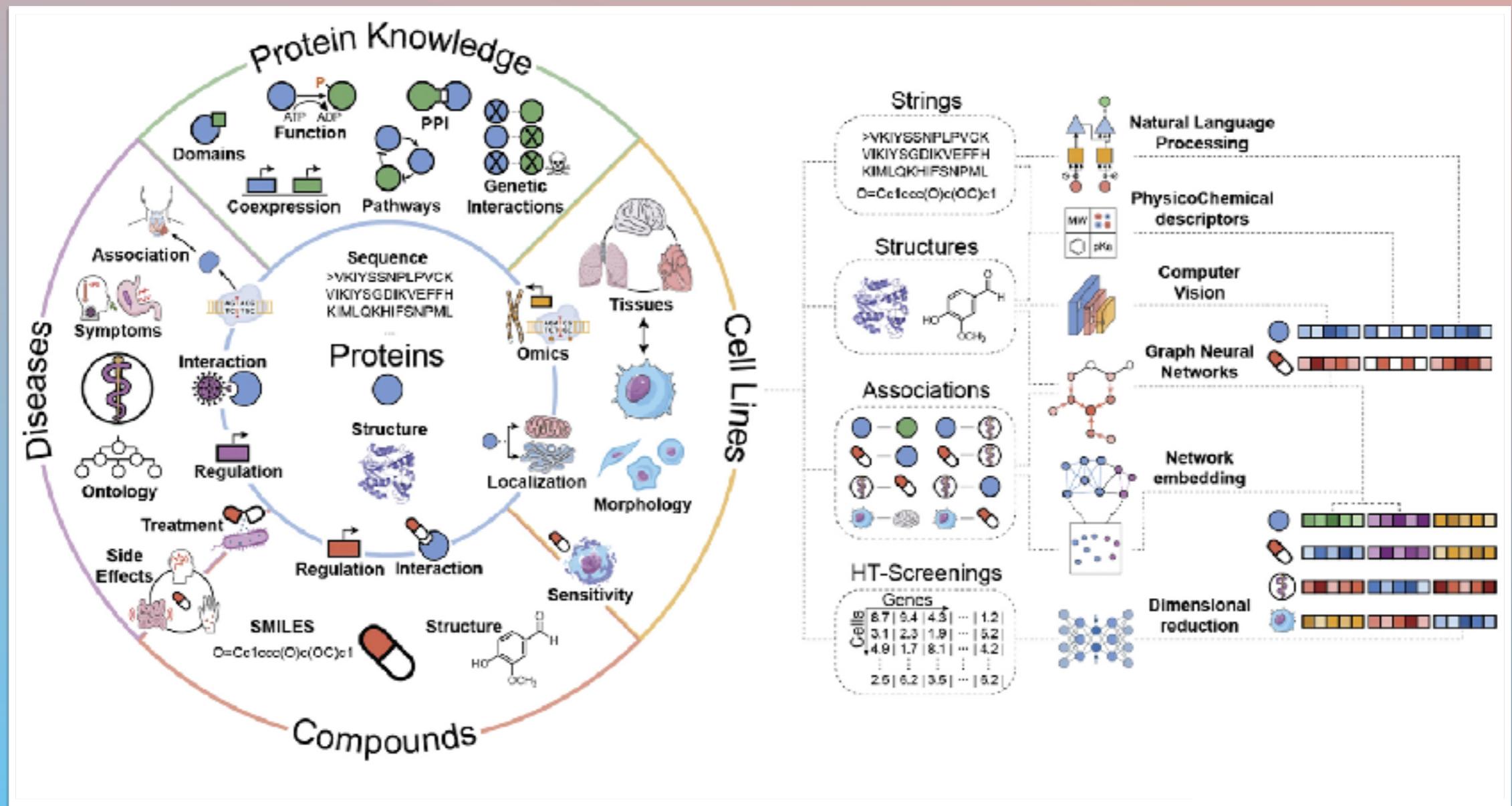


Duran-Frigola et al. WIREs Advanced Reviews, 2019
Fernandez-Torras et al, Nature Communications, 2022

Protein descriptors

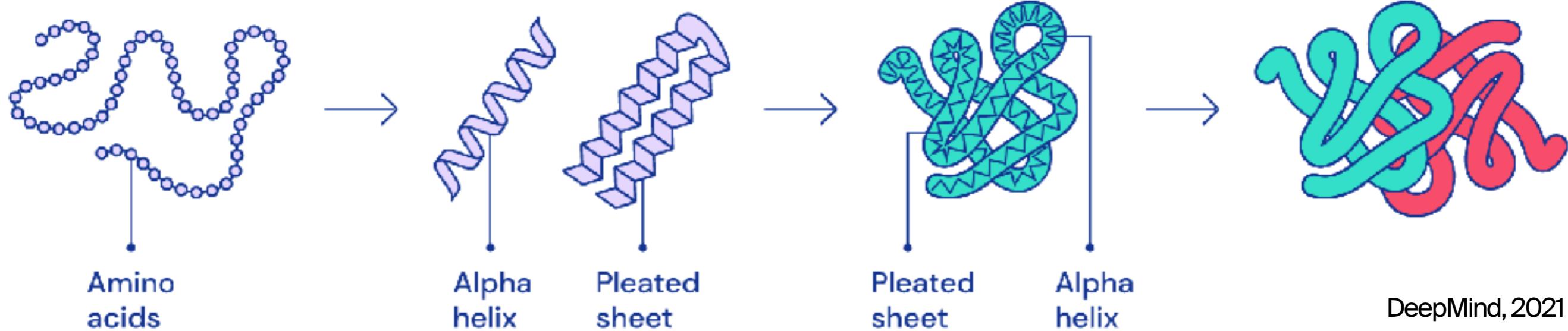


Holistic protein descriptors

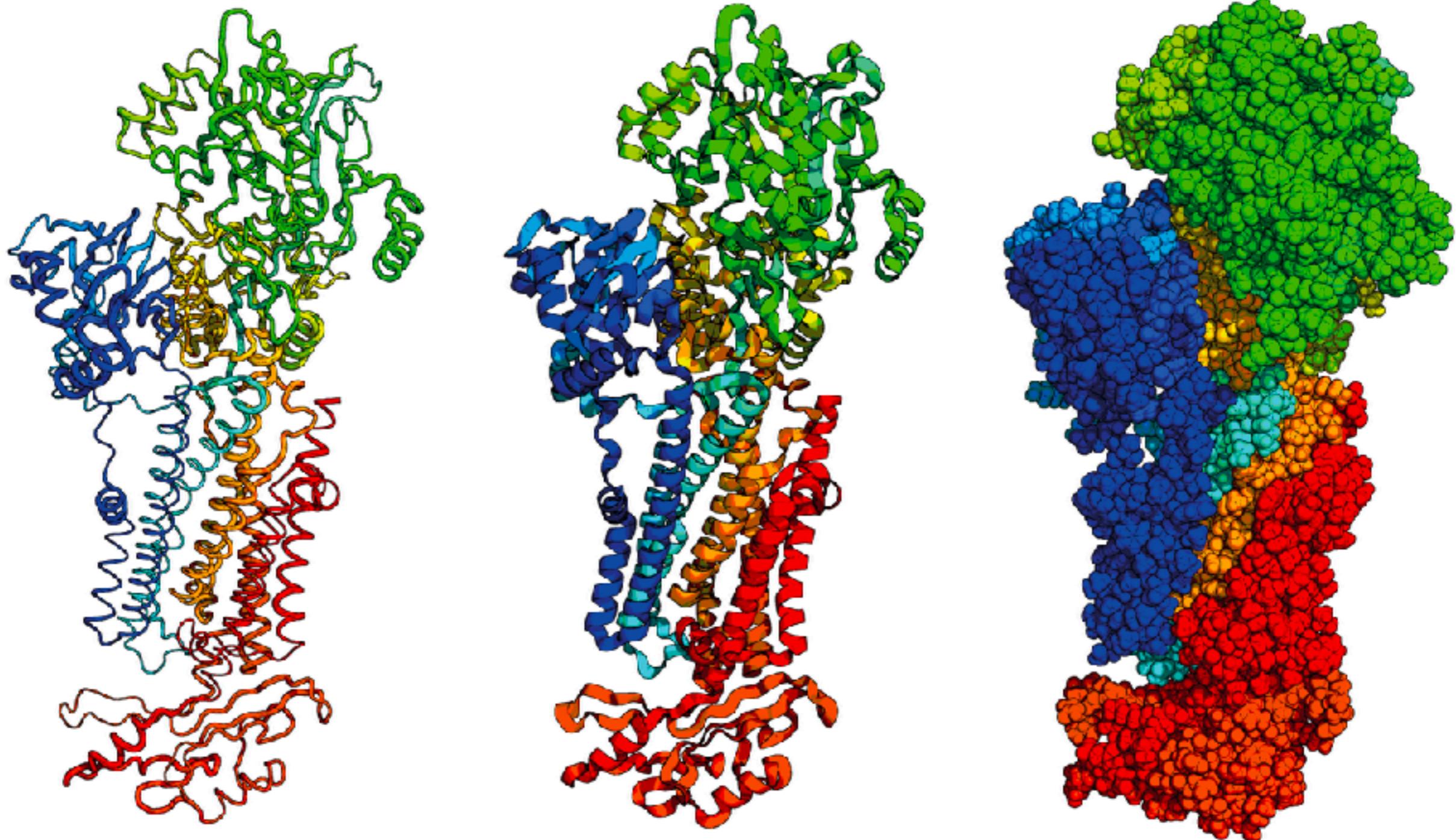




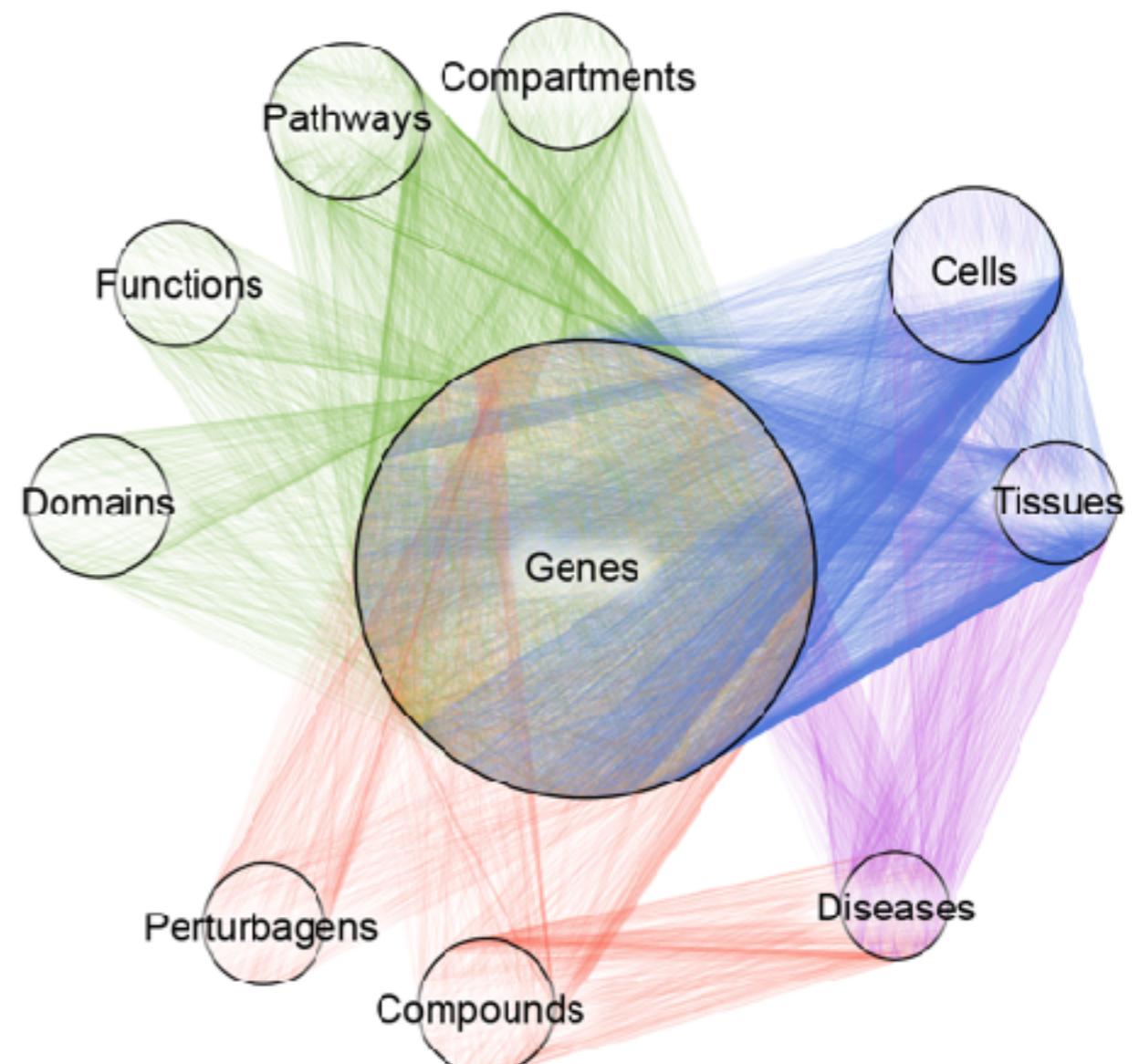
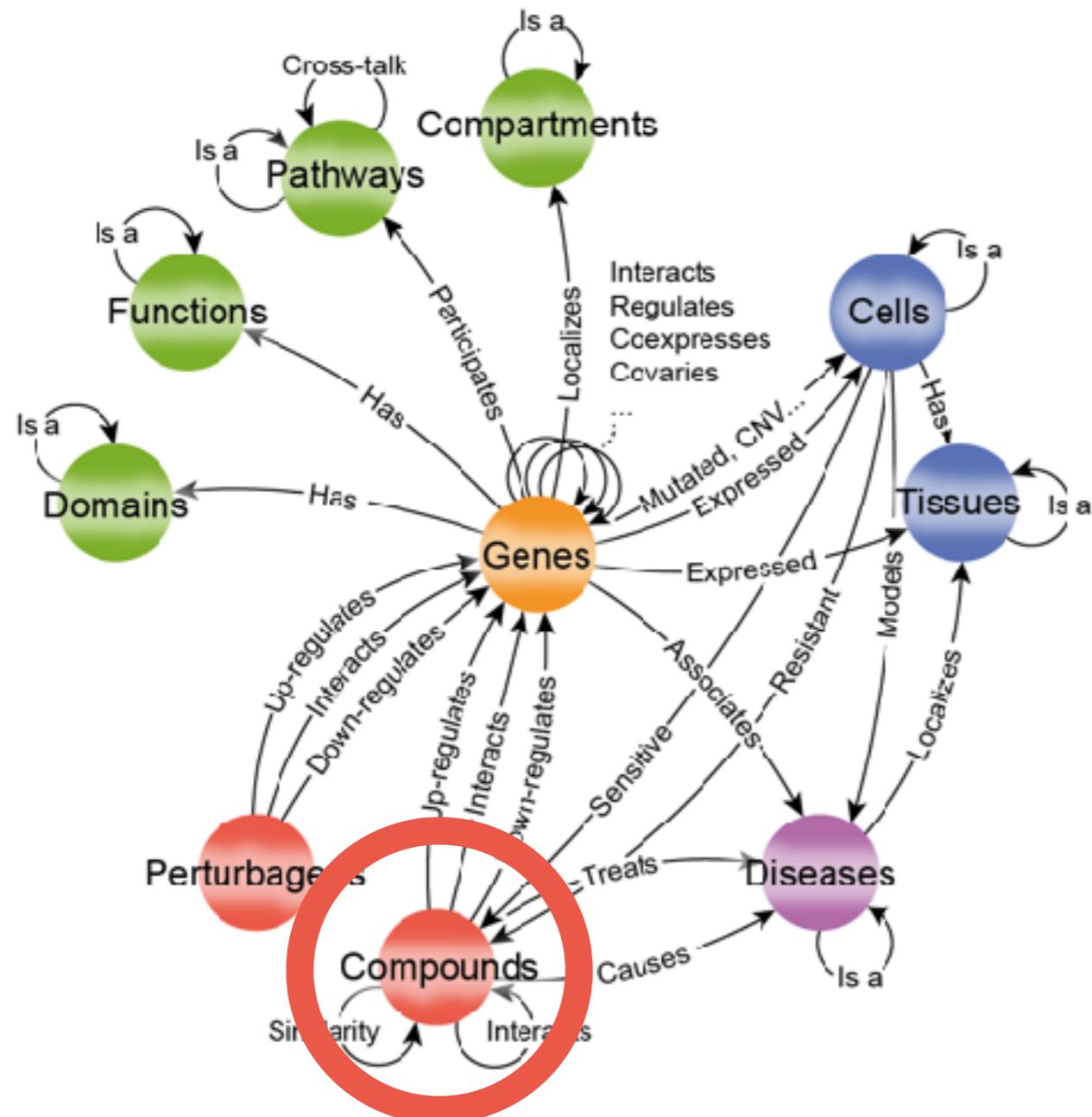
ALPHAGO



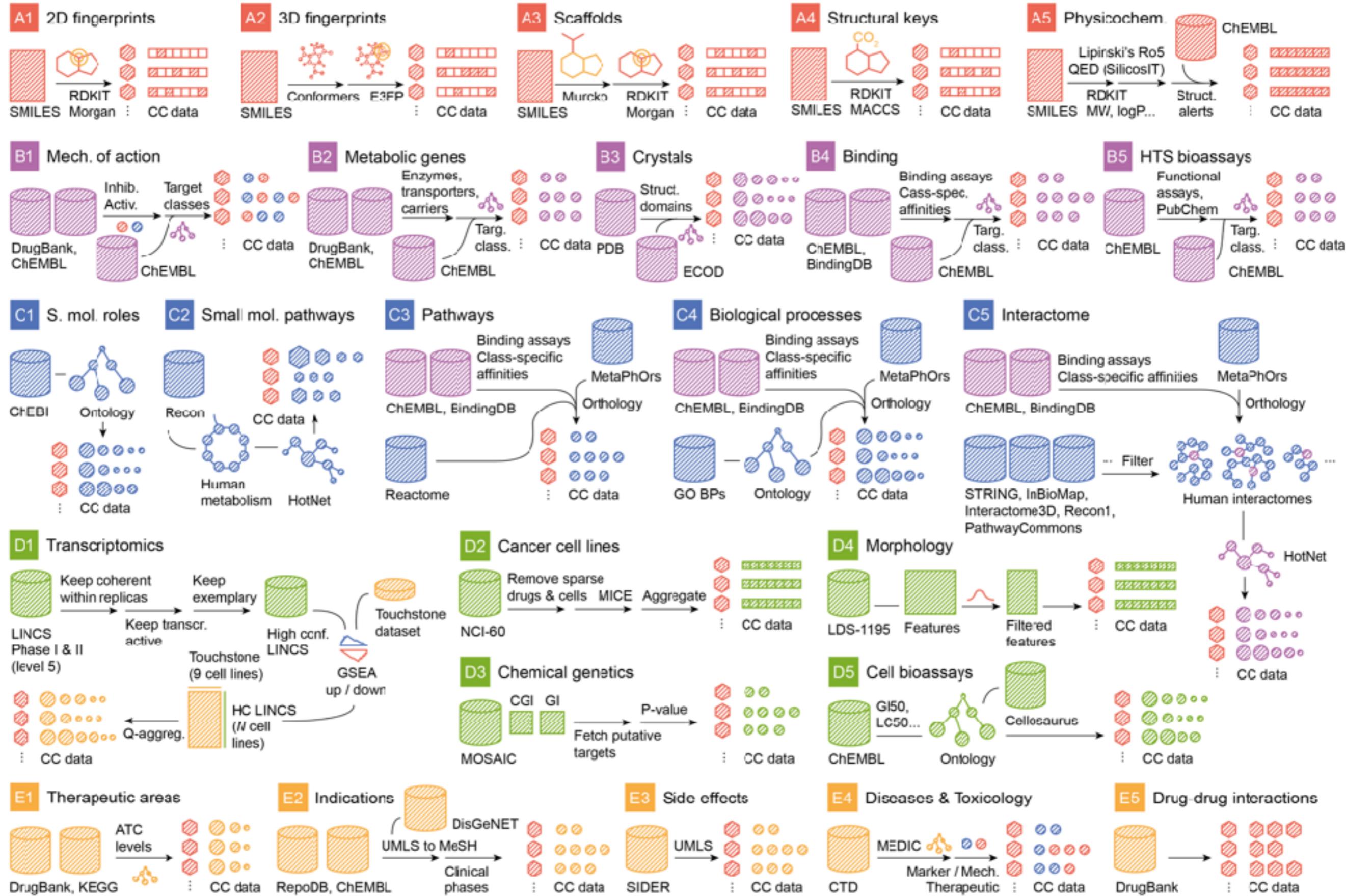
PfATP4 target of *P. falciparum*



Back to compounds!



Duran-Frigola et al. WIREs Advanced Reviews, 2019
 Fernandez-Torras et al, Nature Communications, 2022



The bioactive chemical space, organised

Chemistry	2D fingerprints	3D fingerprints	Scaffolds	Structural keys	Physico-chemistry
Targets	Mechanism of action	Metabolic genes	Crystals	Binding	HTS bioassays
Networks	Small mol. roles	Small mol. pathways	Signaling pathways	Biological processes	Interactome
Cells	Gene expression	Cancer cell lines	Chemical genetics	Morphology	Cell bioassays
Clinics	Therapeutic areas	Indications	Side effects	Diseases and toxicology	Drug-drug interactions

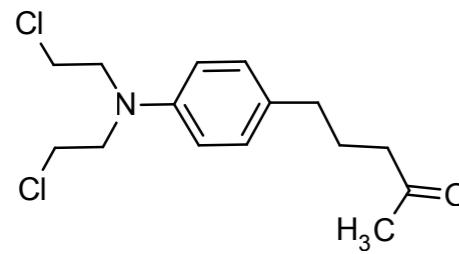
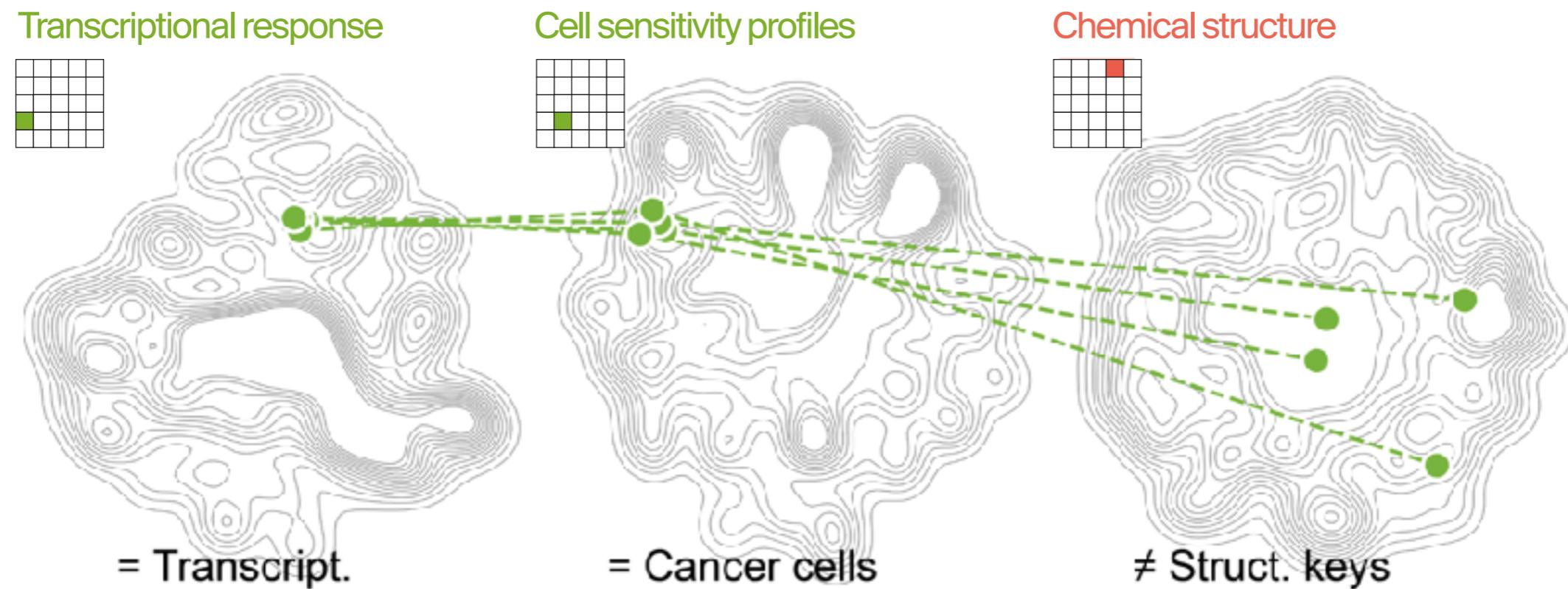
- 800k bioactive molecules
- 25 data types, from chemistry to the clinics
- The major small molecule databases are integrated
- chemicalchecker.org
- bioactivitysignatures.org



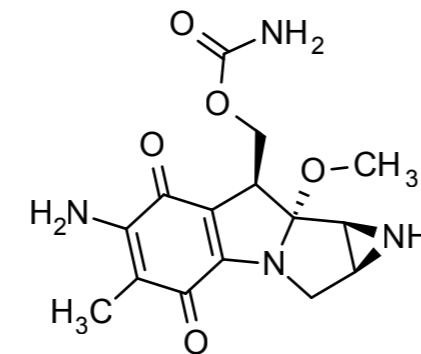
Duran-Frigola et al. Nat Biotech, 2020

 CTD2-pancancer DREAM challenge, 2020

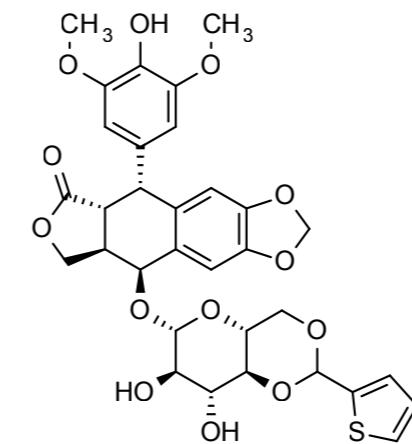
Complex similarity searches



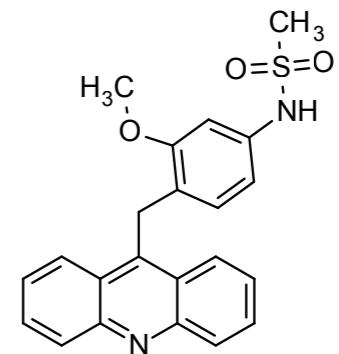
Chlorambucil



Mitomycin

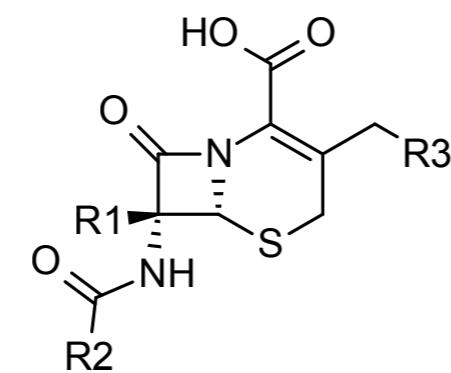
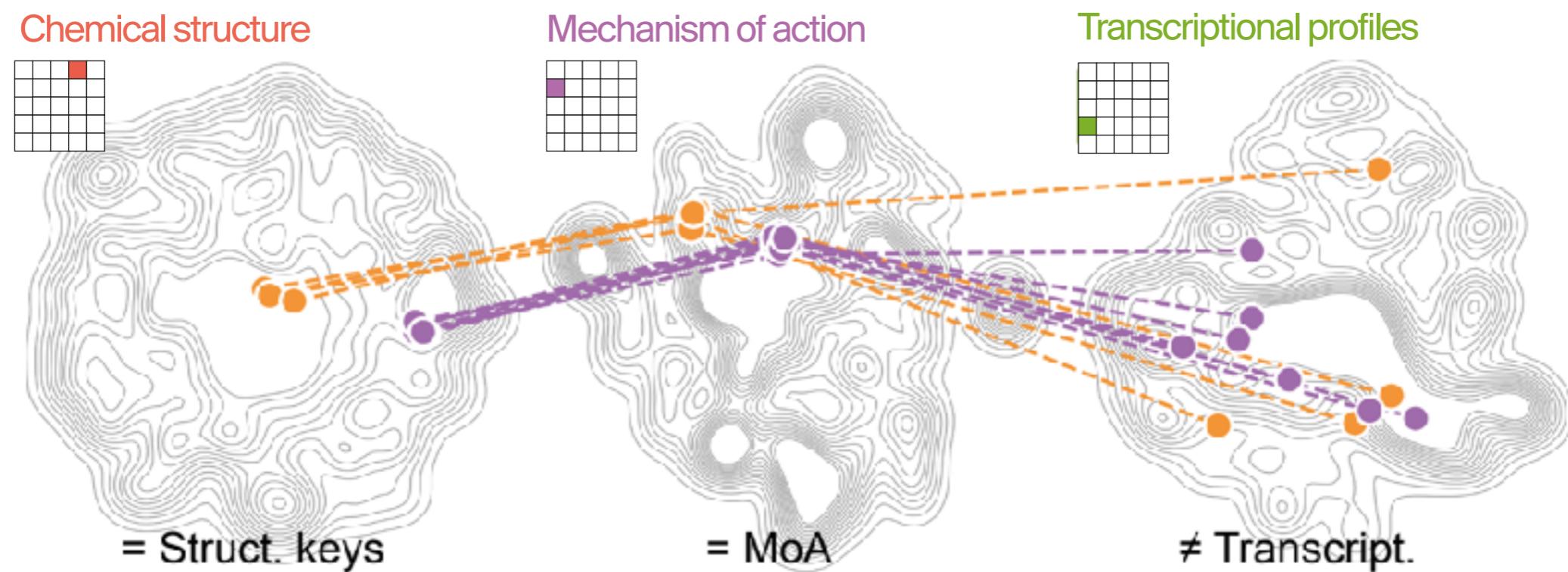


Teniposide

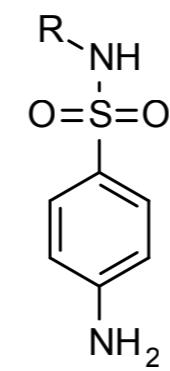


New molecule

Complex similarity searches



Cephalosporins



Sulphonamides

Chemical mimetics of biological drugs

>10,000 compounds considered



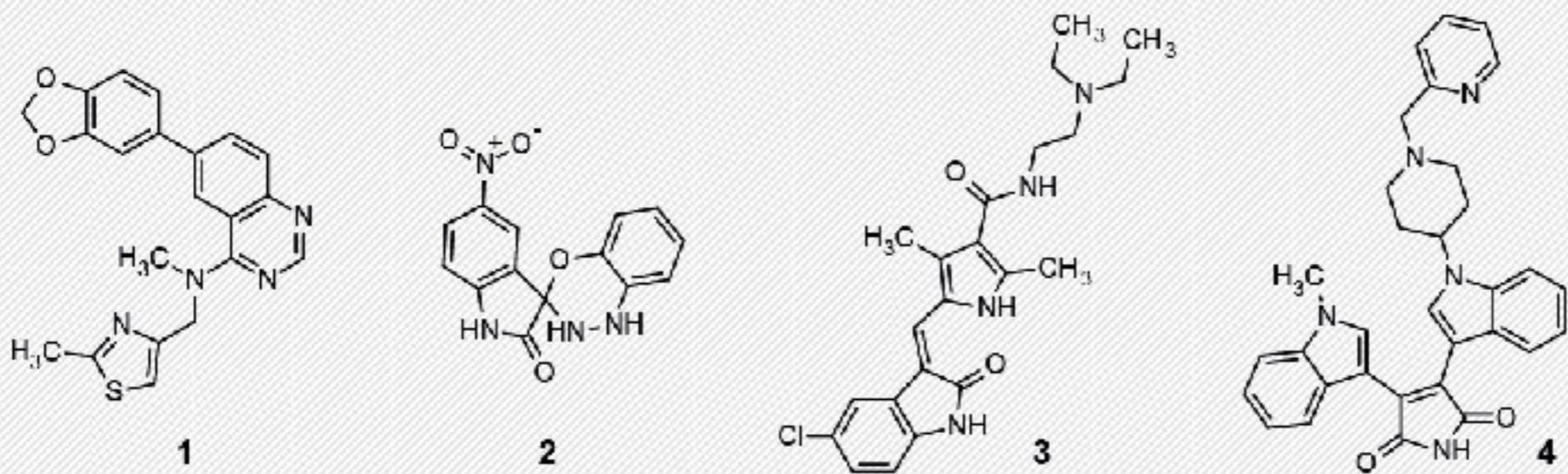
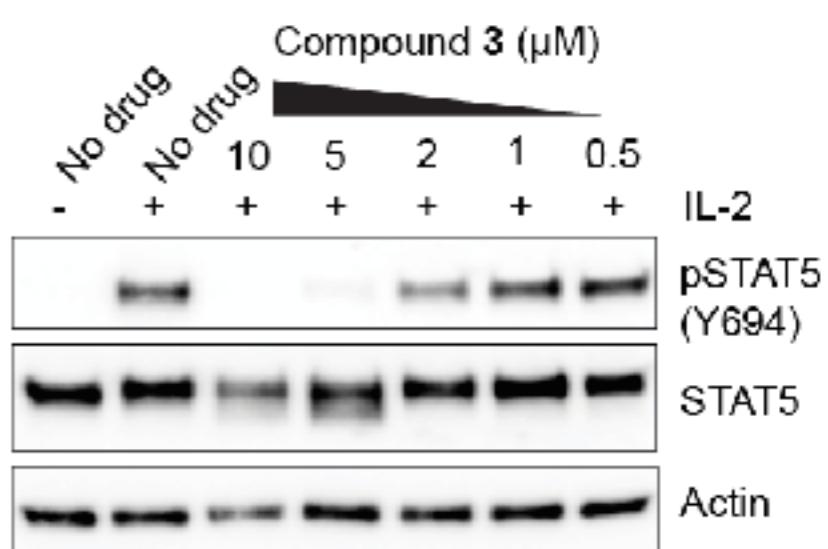
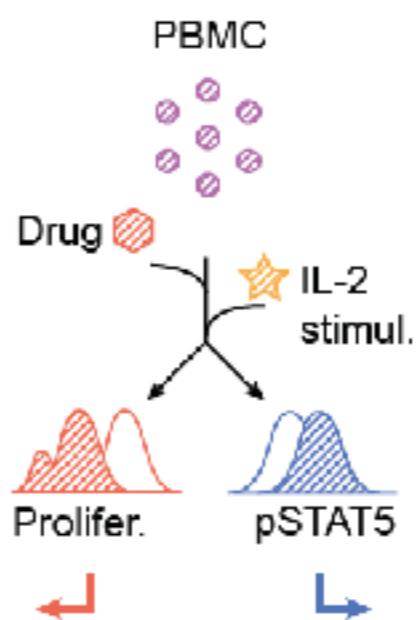
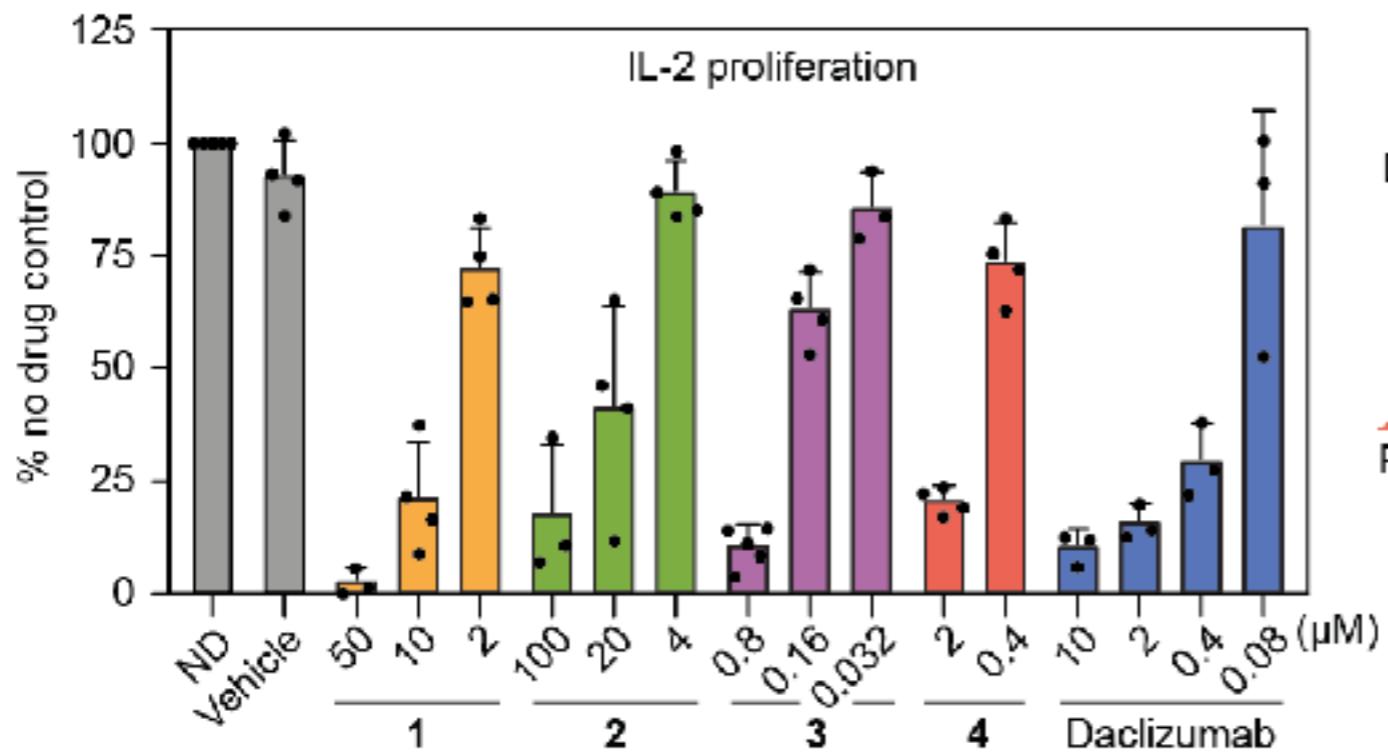
Known and inferred targets of the compounds, together with their signalling pathways.

IL-2 receptor antibody (daclizumab)



Known targets of the antibody drug, together with their signalling pathways.

In vitro experiments



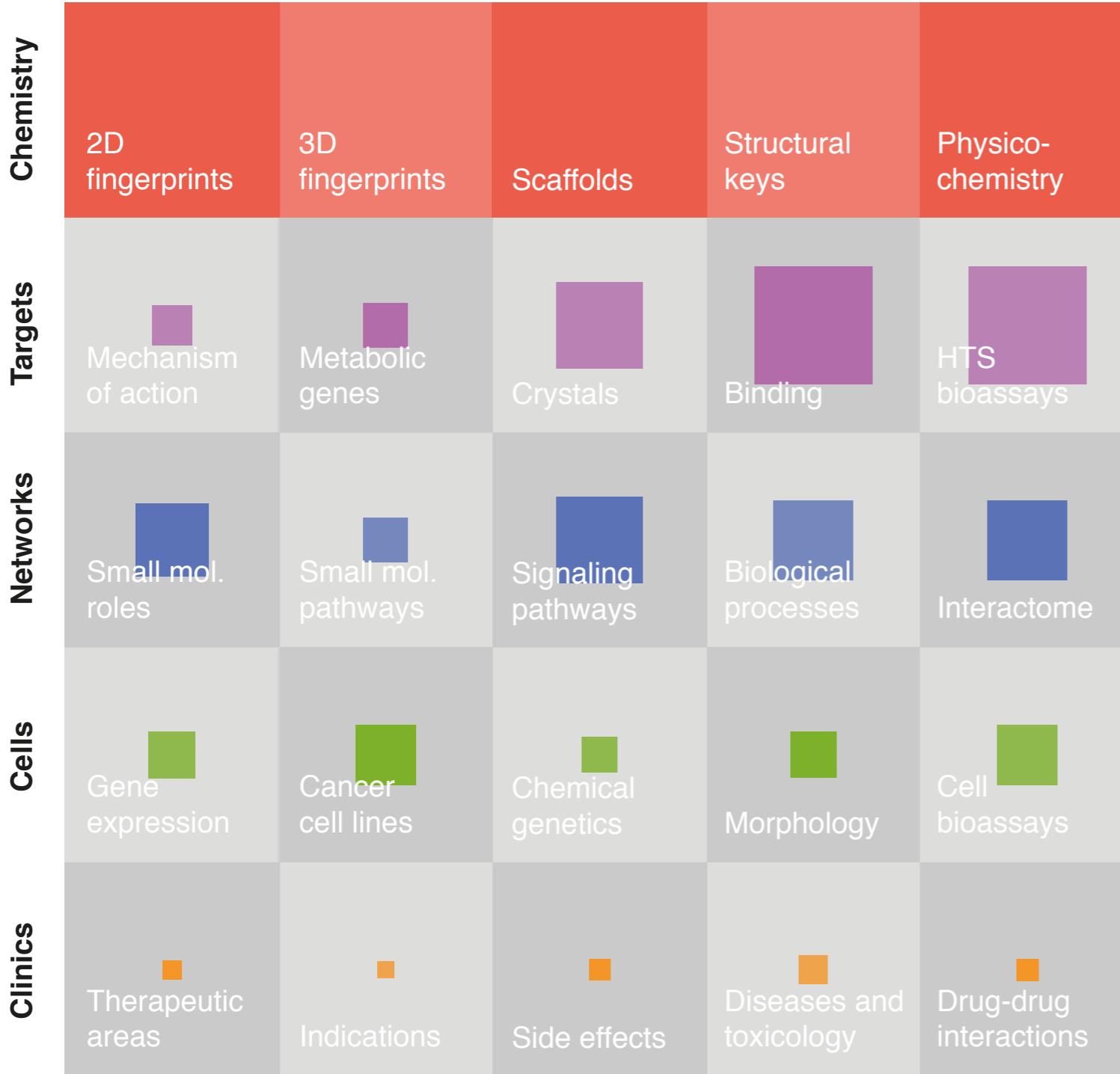
...14 of 19 tested worked*

Chemical space

Chemical space

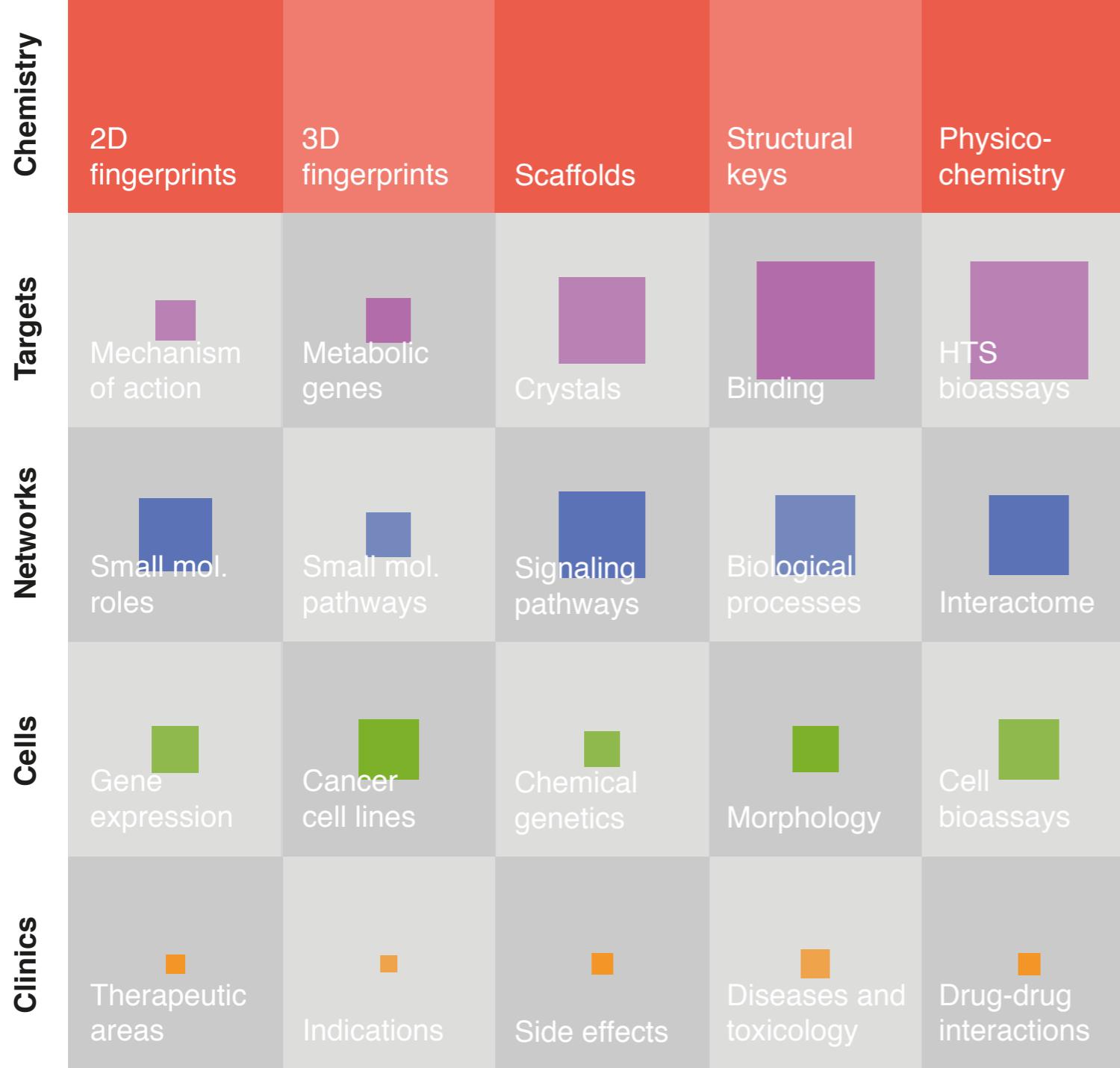


An honest view of the Chemical Checker

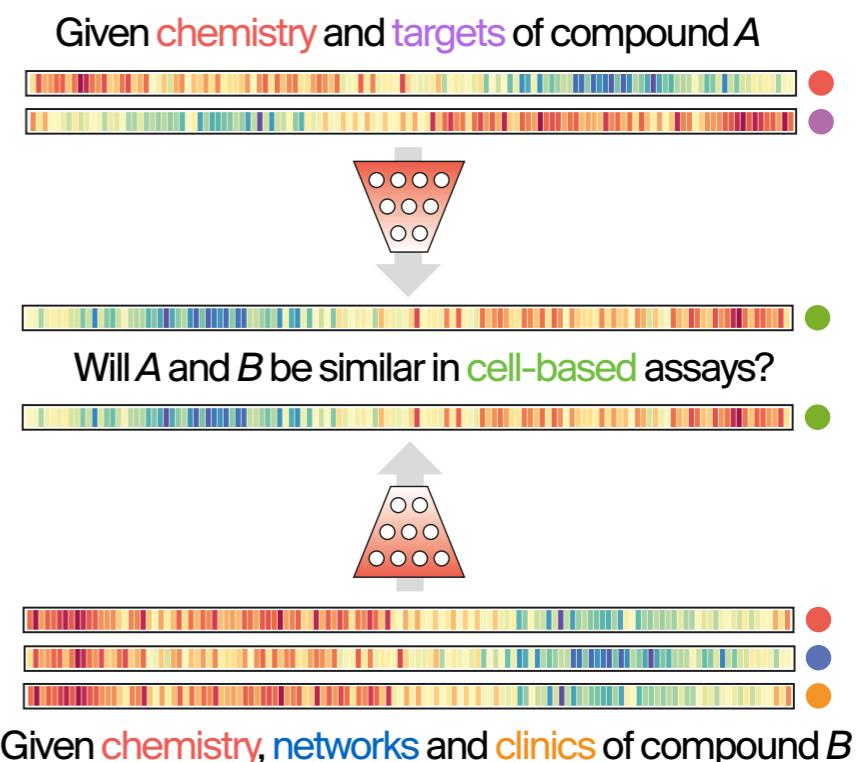


Bertoni*, Duran-Frigola* et al. Nat Commun, 2021

An honest view of the Chemical Checker



– Siamese neural nets 😎



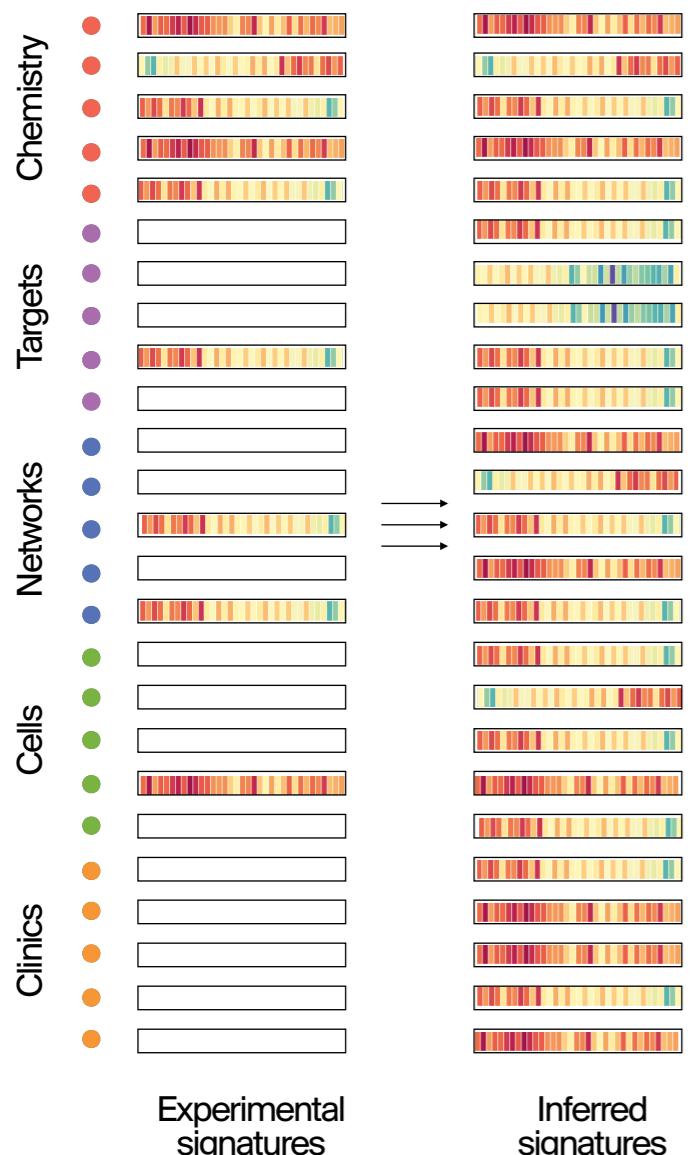
– Trained over a week in a supercomputing center

Bertoni*, Duran-Frigola* et al. Nat Commun, 2021



Large-scale inference of bioactivity signatures

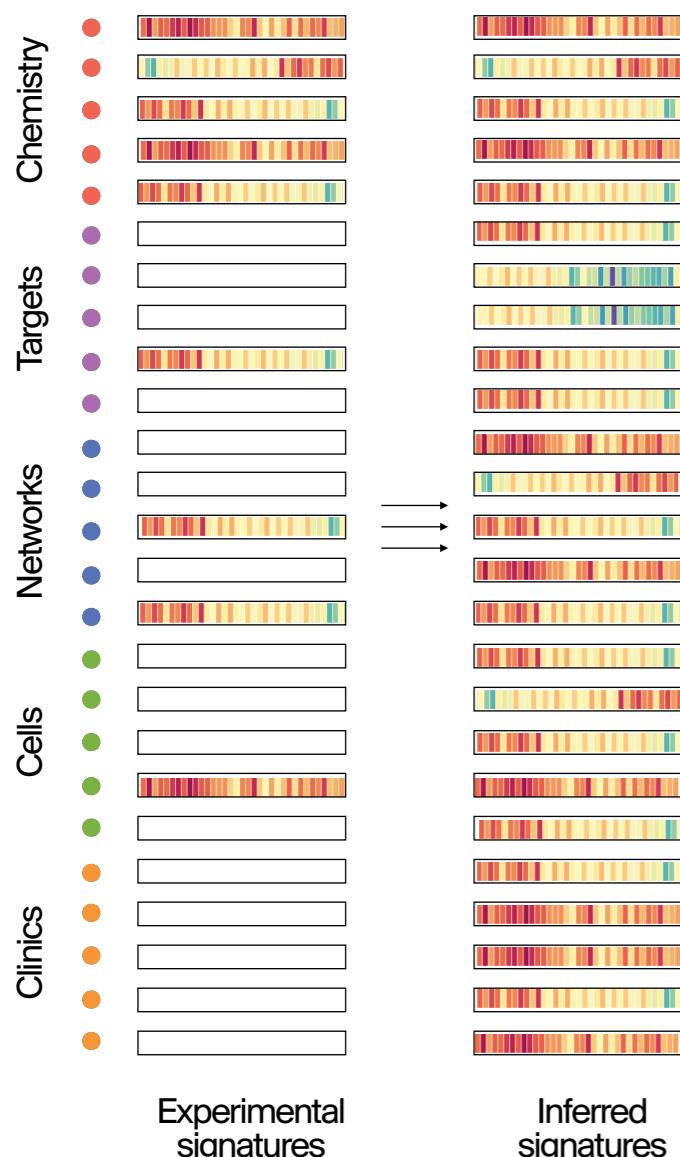
Siamese nets



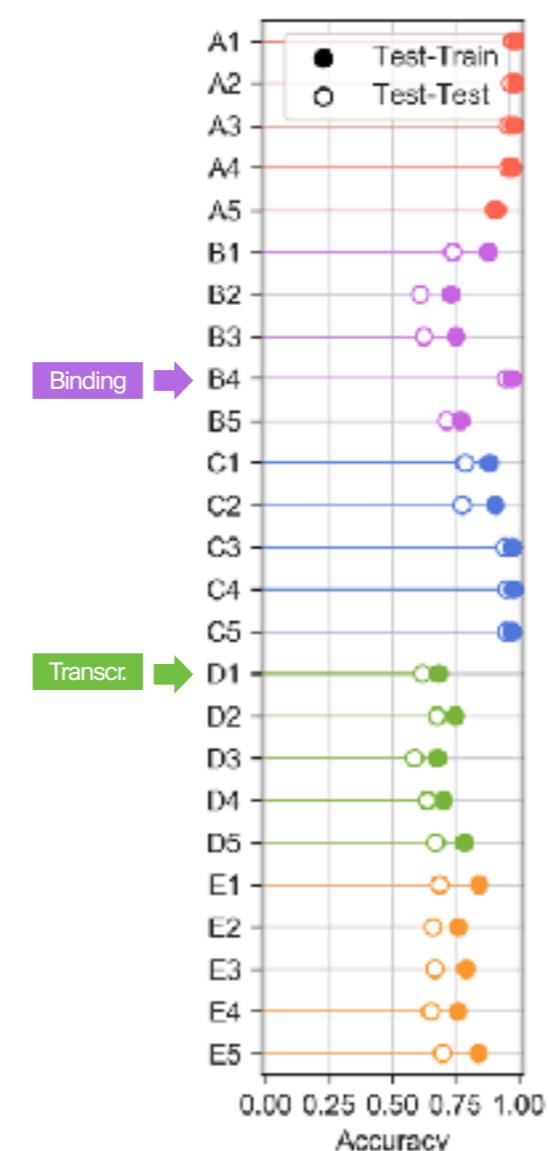
	1	2	3	4	5	
A	A1: 2D fingerprints	A2: 3D fingerprints	A3: Scaffolds	A4: Structural keys	A5: Physicochemistry	
B	B1: Mechanisms of action	B2: Metabolic genes	B3: Crystals	B4: Binding	B5: HTS bioassays	
C	C1: Small molecule roles	C2: Small molecule pathways	C3: Signaling pathways	C4: Biological processes	C5: Interactome	
D	D1: Transcription	D2: Cancer cell lines	D3: Chemical genetics	D4: Morphology	D5: Cell bioassays	
E	E1: Therapeutic areas	E2: Indications	E3: Side effects	E4: Diseases & toxicology	E5: Drug-drug interactions	

Large-scale inference of bioactivity signatures

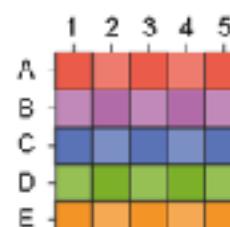
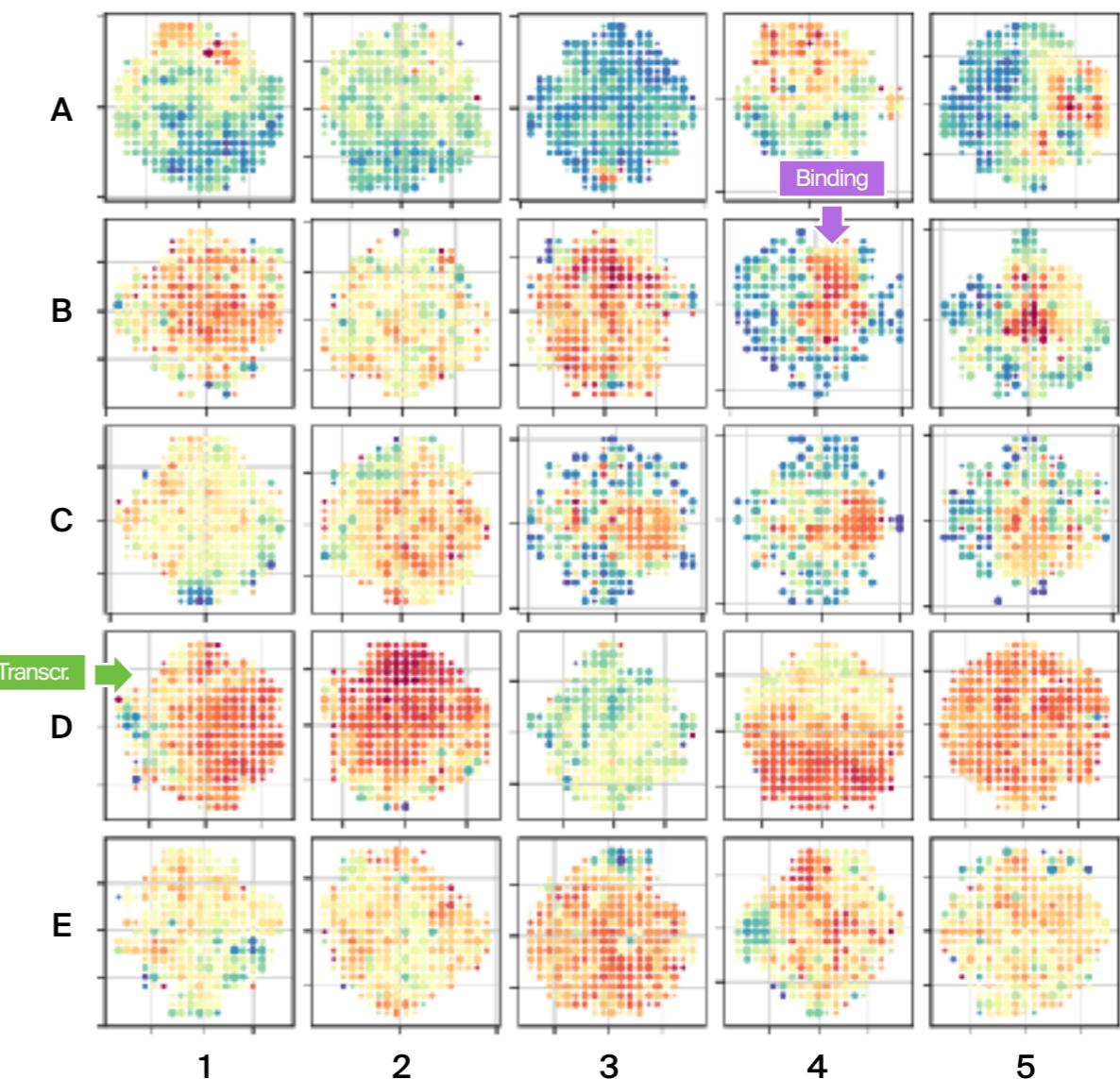
Siamese nets



(a) Accuracy



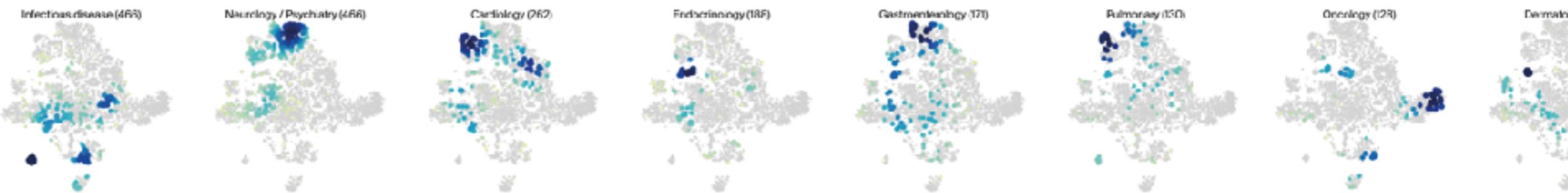
(b) Confidence



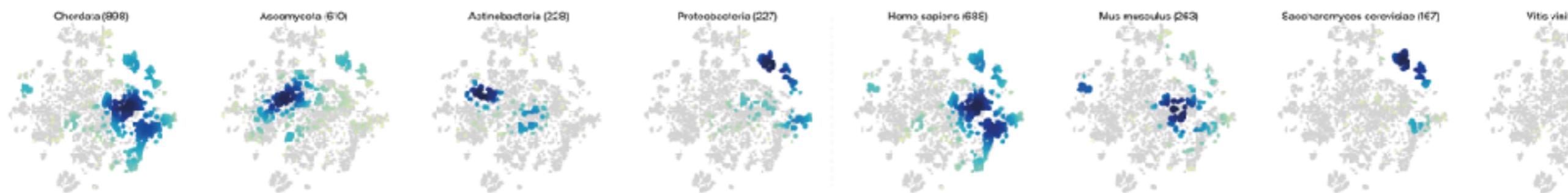
- | | | | | |
|--------------------------|-----------------------------|------------------------|---------------------------|----------------------------|
| A1: 2D fingerprints | A2: 3D fingerprints | A3: Scaffolds | A4: Structural keys | A5: Physicochemistry |
| B1: Mechanisms of action | B2: Metabolic genes | B3: Crystals | B4: Binding | B5: HTS bioassays |
| C1: Small molecule roles | C2: Small molecule pathways | C3: Signaling pathways | C4: Biological processes | C5: Interactome |
| D1: Transcription | D2: Cancer cell lines | D3: Chemical genetics | D4: Morphology | D5: Cell bioassays |
| E1: Therapeutic areas | E2: Indications | E3: Side effects | E4: Diseases & toxicology | E5: Drug-drug interactions |

High
Medium
Low

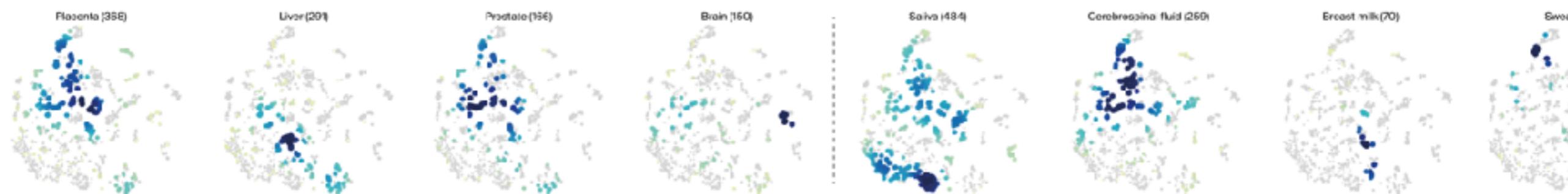
– Drug molecules



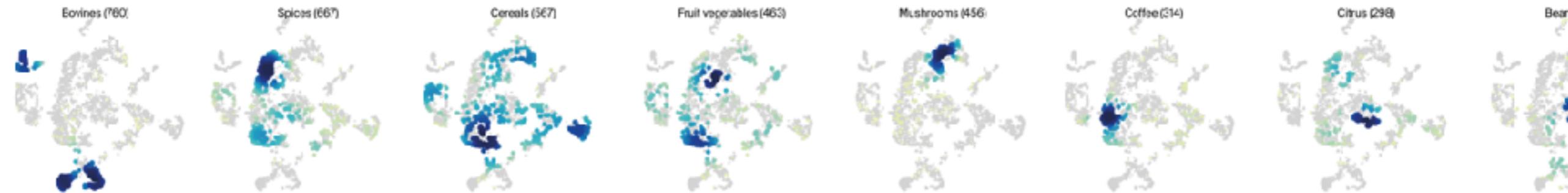
– Species metabolites



– Tissues and fluids

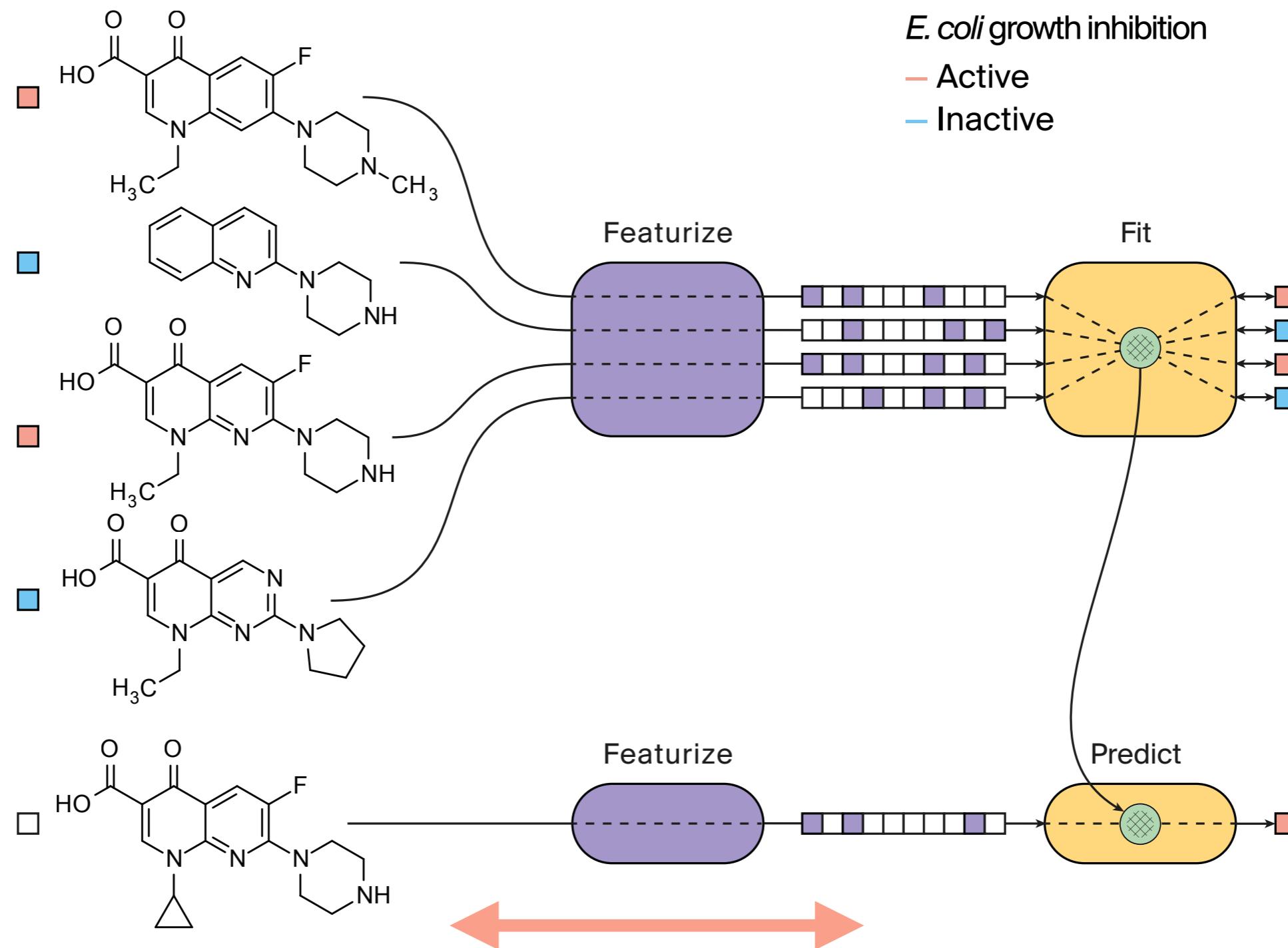


– Food ingredients



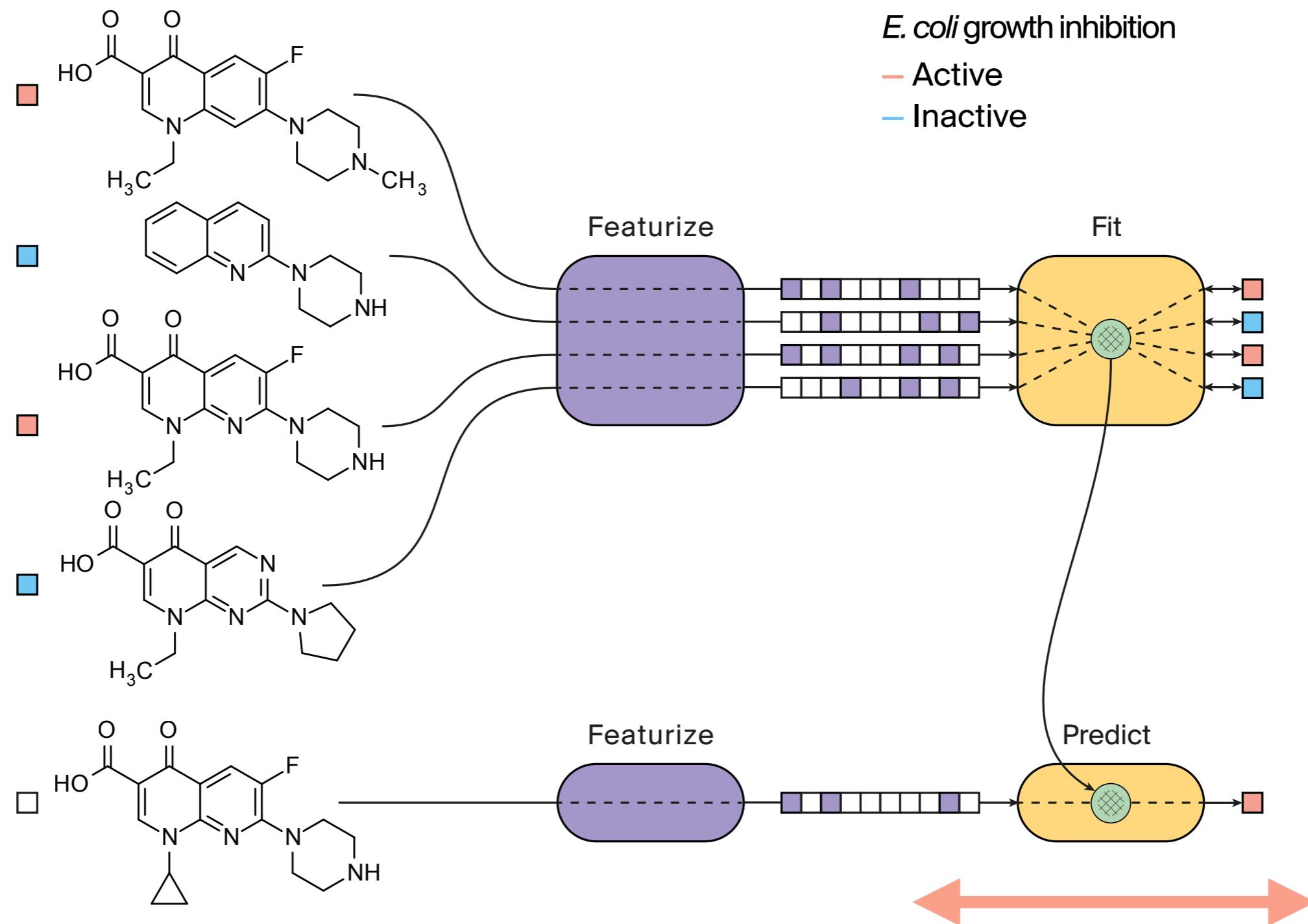
A simple **supervised** task in drug discovery (aka QSAR)

Antibiotic activity experiments (*E. coli*)

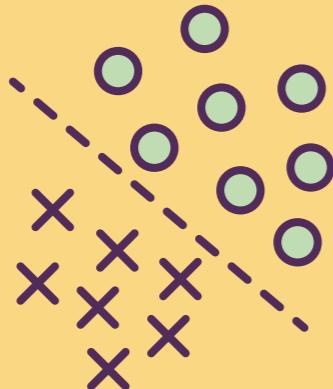


A simple **supervised** task in drug discovery (aka QSAR)

Antibiotic activity experiments (*E. coli*)



Classification

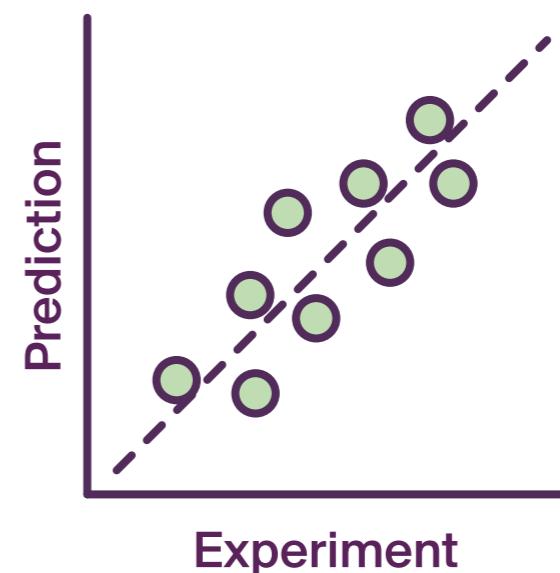


Prediction task:

Active = 1

Inactive = 0

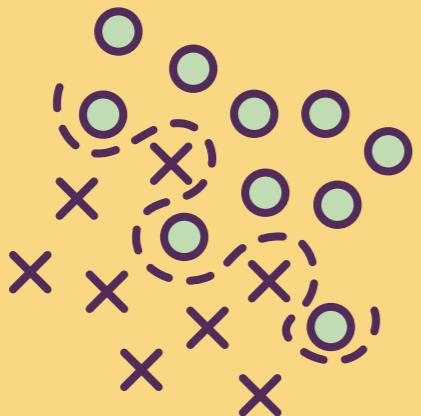
Regression



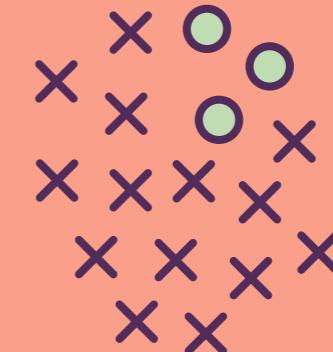
Prediction task:

IC50 value

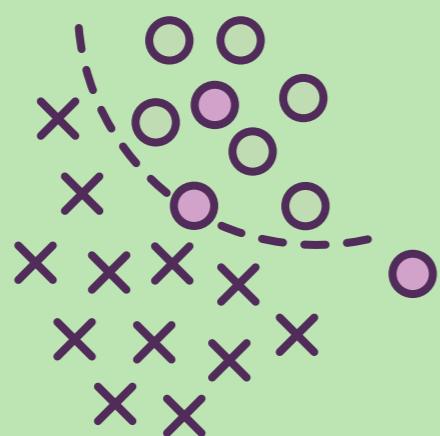
Overfitting



Imbalance



Confidence

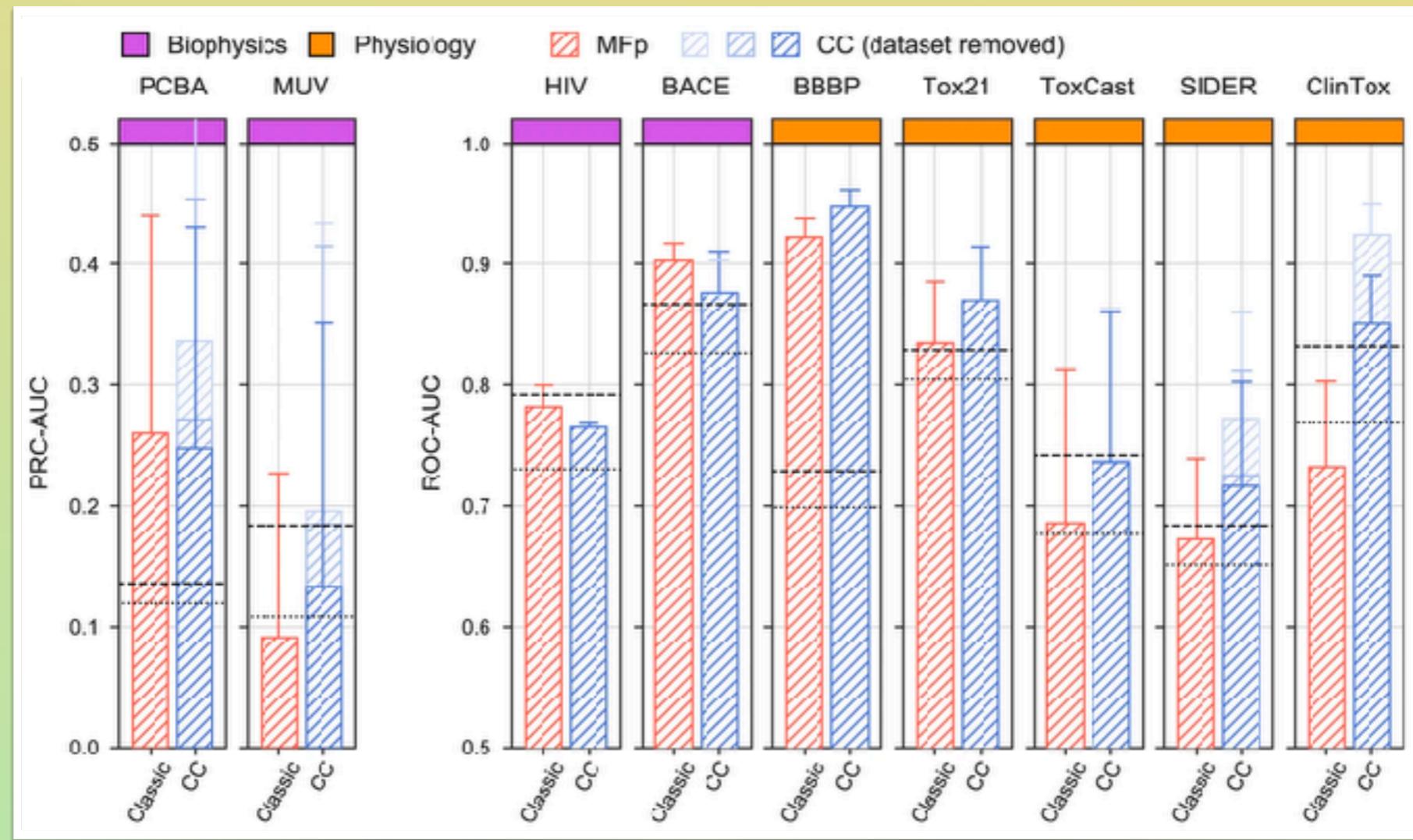


Interpretability



Automated AI/ML pipelines for pharmacology

State-of-the-art performance

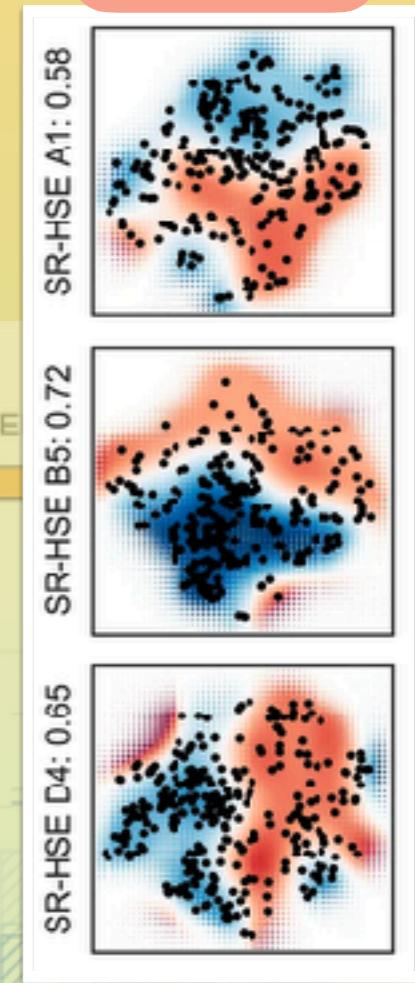
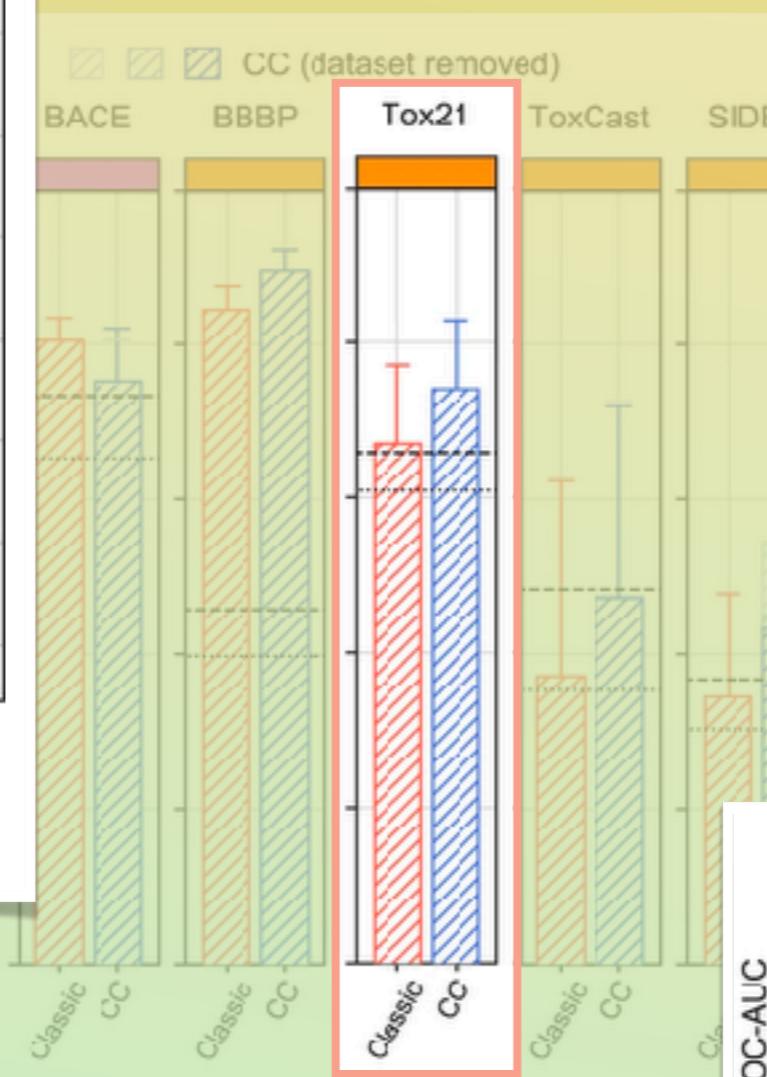
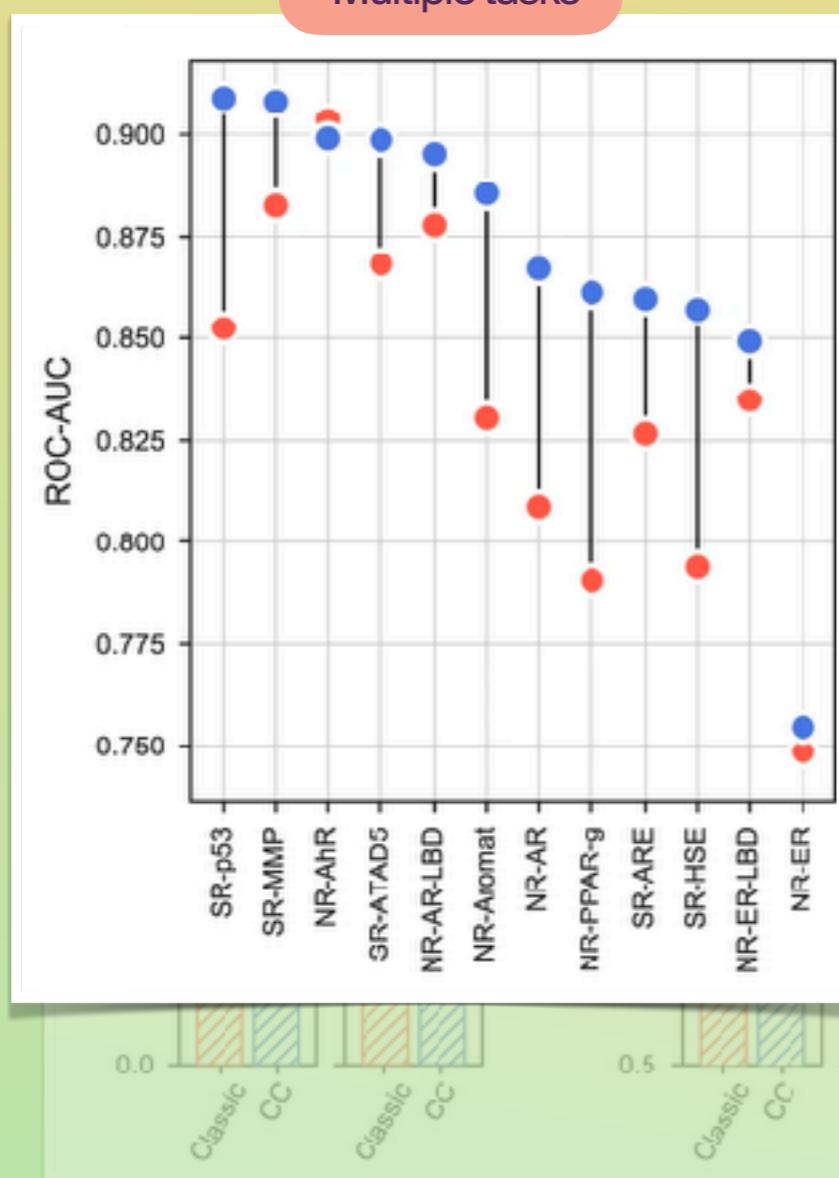


Bertoni*, Duran-Frigola* et al, Nat Commun, 2021
Ersilia's tool: ZairaChem

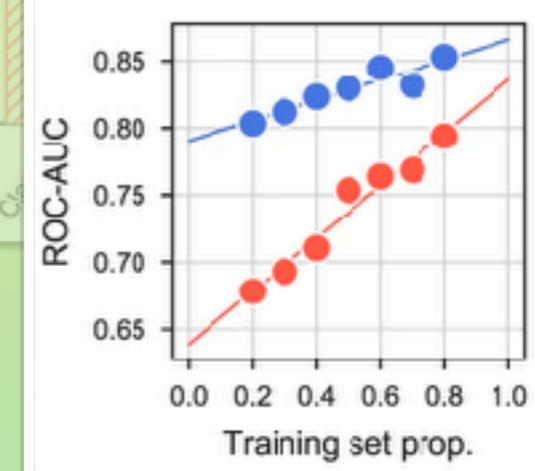
Performance of automated AI/ML

Visualization

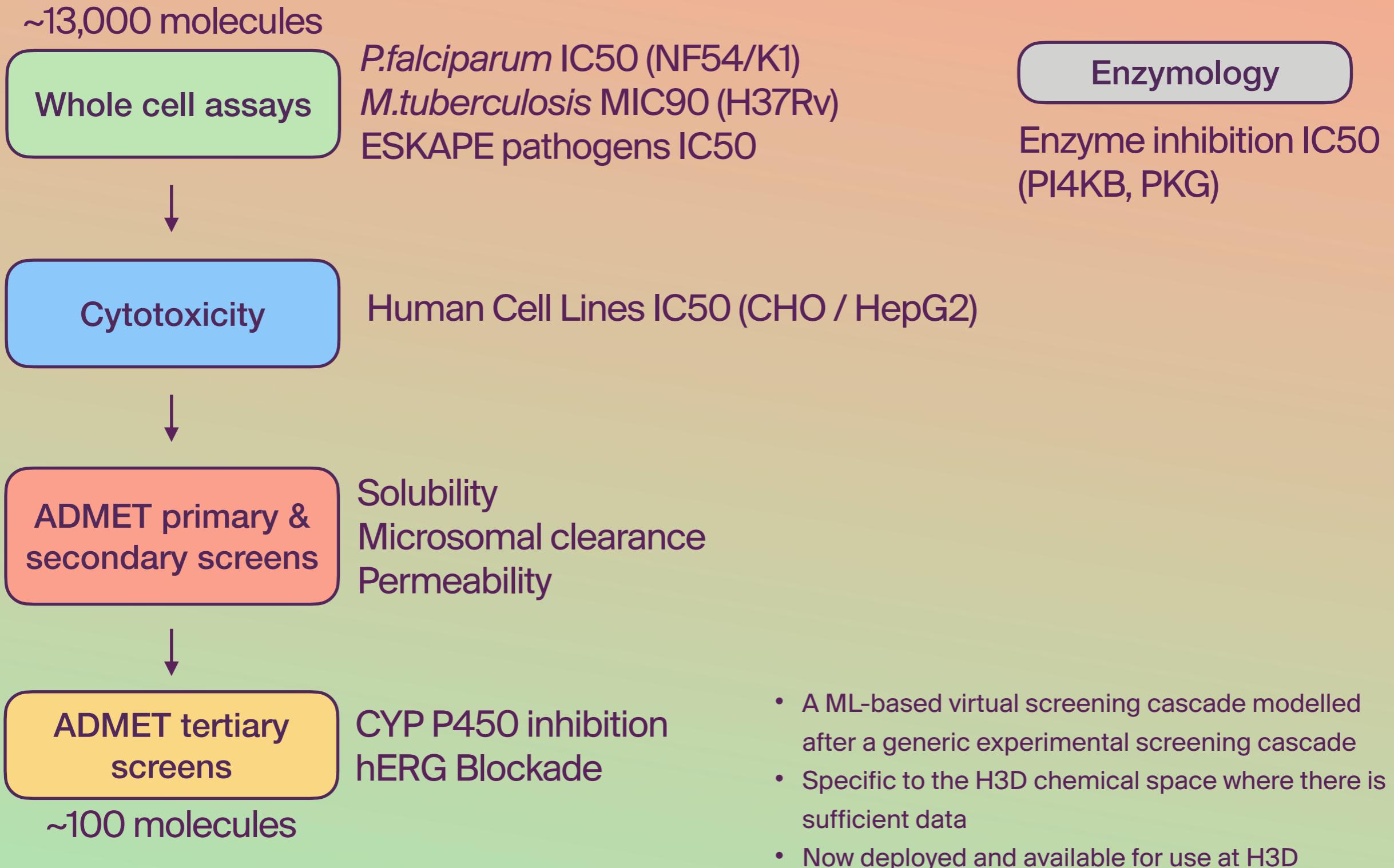
Multiple tasks



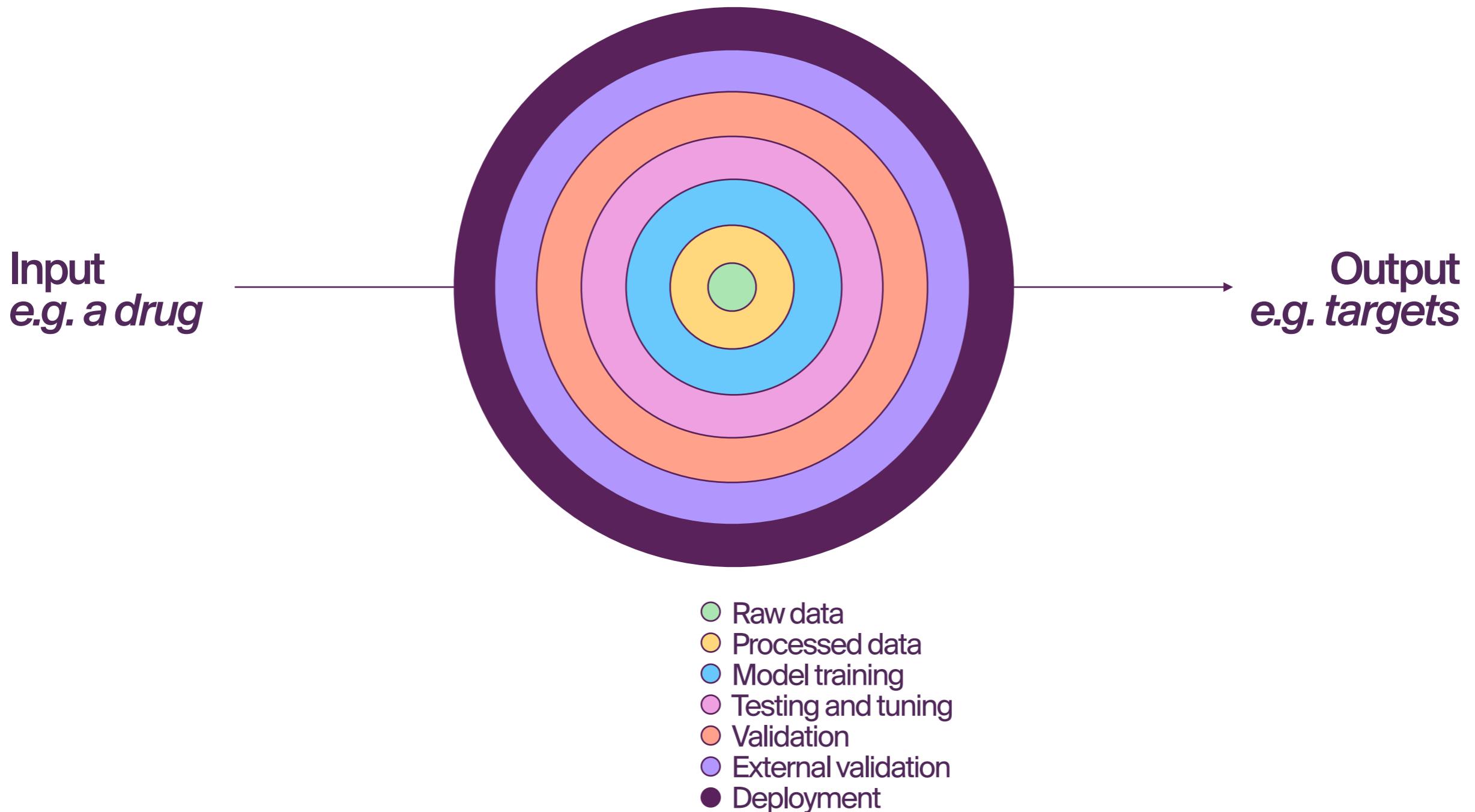
Robustness



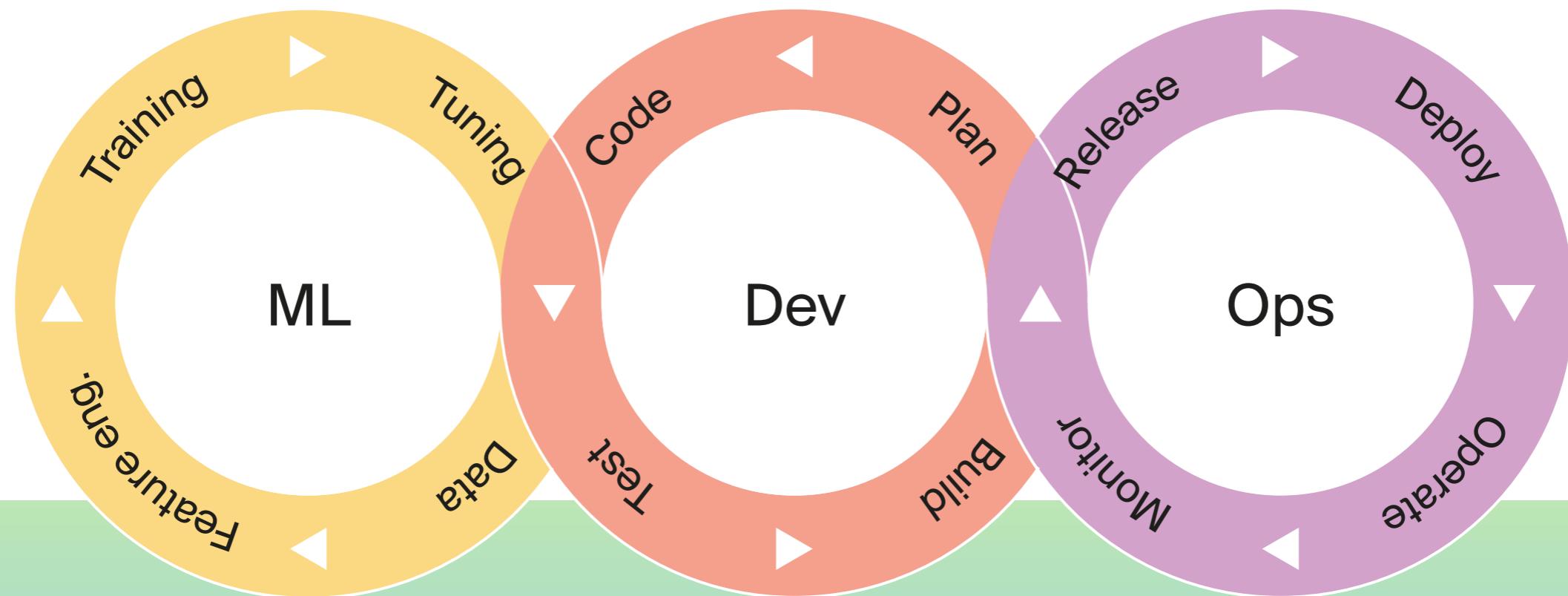
A generic H3D screening cascade



Our goal: to provide ready-to-use AI/ML models



Continuous delivery and automation pipelines



- Training mode: once or continuously updated
- Inference mode: black-box or interpretable
- Privacy: prevent reverse engineering
- Interoperability: easy distribution of AI/ML assets

AI/ML in collaboration

 Ersilia ‘trains’ a model based on partner’s data



Natural product activity
University of Buea



AI/ML in collaboration

 Ersilia ‘trains’ a model based on partner’s data



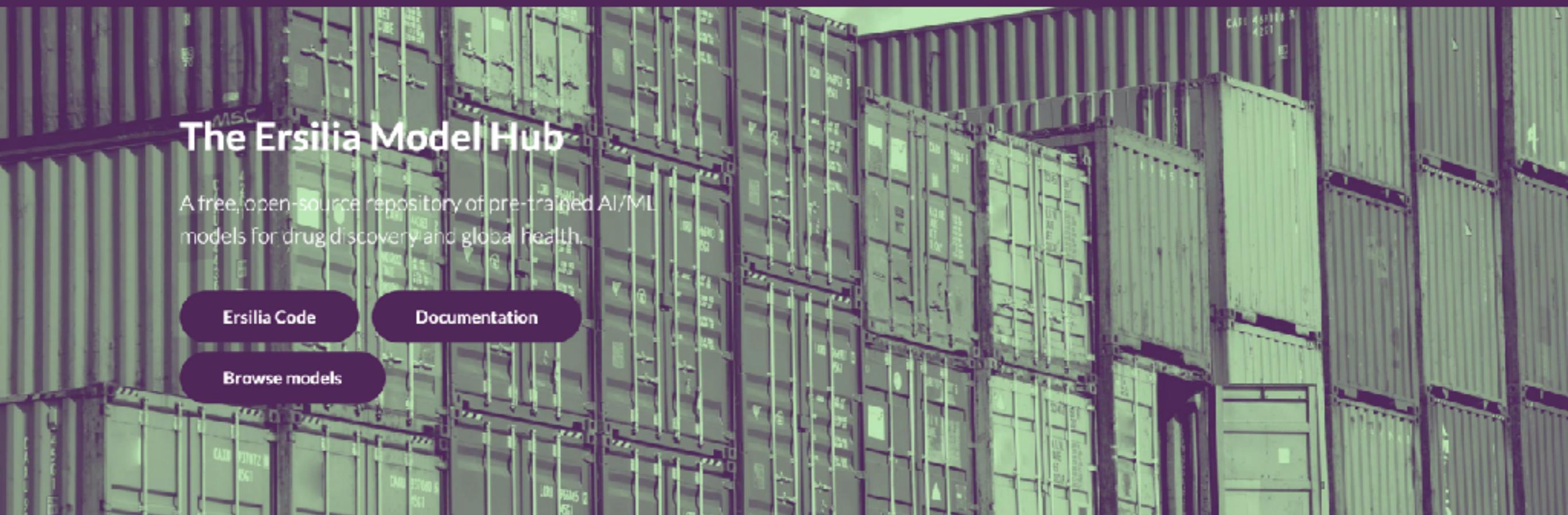
Natural product activity
University of Buea

Humimycin A

Active!

The Ersilia Model Hub

A free, open-source repository of pre-trained AI/ML models for drug discovery and global health.

[Ersilia Code](#)[Documentation](#)[Browse models](#)

A hundred models in the Ersilia Model Hub

Check the latest additions to the Ersilia Model Hub! We systematically look for AI/ML models and datasets in the scientific literature and incorporate them in our platform.

Classification of carcinogenic human metabolites

Predicts six properties using an ensemble machine learning method to predict carcinogenicity

[Toxicity](#)

SMILES transformer descriptor

Molecular fingerprint based on SMILES representations for low-data drug discovery

[Vector](#)

SARS-CoV inhibition

This model was developed to support the early efforts in the identification of novel drugs against SARS-CoV2. It pred...

[Antiviral activity](#)[More models](#)

convert them to a readable CSV file.

chemical libraries as well as generative models

molecule bioactivity signatures.

`bioassay-db`

`chem-sampler`

`chemical-checker`

[Code](#)

[Code](#)

[Code](#)

End-to-end record linkage

A fully automated pipeline for fuzzy matching of medical datasets.

`cldr-e2e-linkage`

[Code](#)

Helper tools for the ESTHER record linkage course

Helper methods to simplify record linkage in the context of the ESTHER course.

`estherlinkage`

[Code](#)

Griddify

Griddify high-dimensional tabular data for easy visualization and deep learning

`griddify`

[Code](#)

A data lake of pre-calculated predictions based on the Ersilia Model Hub

Isaura is a data lake of pre-calculated descriptors based on the Ersilia Model Hub

`isaura`

[Code](#)

On-Disk X and Y matrices

Manage X and Y matrices on disk for scalable machine learning

`ondisk-xy`

[Code](#)

Chemical space visualization

Visualization of the chemical space based on TMAP trees

`tmap-chemistry`

[Code](#)

© 2022 Ersilia Open Source Initiative. Content on this website licensed under CC BY 4.0



[Organisation](#)

About

Code of conduct

Diversity and inclusion

[Help](#)

Request model

Email

FAQ

[Support us](#)

Donate

Volunteer

Code with us

convert them into a readable CSV file.

chemical libraries as well as generative models

molecule bioactivity signatures.

[bioassay-db](#)

[chem-sampler](#)

[chemical-checker](#)

[Code](#)

[Code](#)

[Code](#)

End-to-end record linkage

Helper tools for the ESTHER record linkage course

Griddify

A fully automated pipeline for fuzzy matching of medical datasets.

Helper methods to simplify record linkage in the context of the ESTHER course.

Griddify high-dimensional tabular data for easy visualization and deep learning

[cidrz-e2e-linkage](#)

[estherlinkage](#)

[griddify](#)

[Code](#)

[Code](#)

[Code](#)

A data lake of pre-calculated predictions based on the Ersilia Model Hub

Isaura is a data lake of pre-calculated descriptors based on the Ersilia Model Hub

[isaura](#)

[Code](#)

On-disk X and Y matrices

Manage X and Y matrices on disk for scalable machine learning

[ondisk-xy](#)

[Code](#)

Chemical space visualization

Visualization of the chemical space based on TMAP trees

[tmap-chemistry](#)

[Code](#)

© 2022 Ersilia Open Source Initiative. Content on this website licensed under CC BY 4.0



[Organisation](#)

[Help](#)

[Support us](#)

[About](#)

[Request model](#)

[Donate](#)

[Code of conduct](#)

[Email](#)

[Volunteer](#)

[Diversity and inclusion](#)

[FAQ](#)

[Code with us](#)

Welcome to the Ersilia Model Hub!

Broad spectrum antimicrobial activity

 Type to search model...

Tags

Tox21

Toxicity

MoleculeNet

Grover

Graph Transformer

Output

Antibiotic activity

Toxicity

Synthetic accessibility

Antiviral activity

Target

Mode

Pretrained

Retrained

In-house

Online

License

Carcinogenic potential of metabolites and small molecules

eos1579

metabokiller

Carcinogenicity is a result of several potential effects on cells. This model predicts the carcinogenic potential of a small molecule based on their potential to induce cellular proliferation, genomic instability, oxidative stress, anti-apoptotic responses and epigenetic alterations.

Metabokiller uses the Chemical Checker signaturizer to featurize the molecules, and the Lime package to provide interpretable results.

Using Metabokiller, the authors screened a panel of human metabolites and experimentally demonstrated two of the predicted carcinogenic metabolites induced carcinogenic transformations in yeast and human cells.

Molecular maps based on broadly learned knowledge-based representations

eos6m4j

bidd-molmap

Descriptor-based or fingerprint-based molecular maps (images) are created. Typically, the goal is to use these images as inputs for an image-based deep learning model such as a convolutional neural network.

SMILES transformer descriptor

eos2lm0

smiles-transformer

Molecular fingerprint based on natural language processing. It converts SMILES into fingerprints using an unsupervised model pre-trained on a very large SMILES dataset. The transformer is particularly well-suited for low-data drug discovery.

Similarity search in the ZINC database

eos54c7

zinc-similarity

Look for 100 nearest neighbors, according to ECFP4 Tanimoto similarity, in the ZINC database.

This model posts queries to an external online server.

Similarity search in ChEMBL

cos2a9n

chembl-similarity

Given a molecule, this model looks for its 100 nearest neighbors in the ChEMBL database, according to ECFP4 Tanimoto similarity.

This model is redirected to ChEMBL web server directly, so molecules are submitted to an external server.

[See more](#)

Can't find what you are looking for?

If you are looking for a specific model, please reach out to us. We will be happy to help!

[Contact us](#)



[About](#)

[Request model](#)

[Backlog](#)

[Volunteer](#)

Take-home messages

- Supervised AI/ML for activity prediction
- Importance of chemical descriptors
- Data-driven descriptors with the Chemical Checker
- Challenges of supervised AI/ML methods
- Pre-trained models in the Ersilia Model Hub
- Use Slack!
- Add terms to the Glossary!
- Learn more in the hands-on sessions!



Bringing data science and AI/ML tools to infectious disease research

Session 2: Introduction to Supervised Machine Learning

Event Sponsors

