

# BMI 702: Biomedical Artificial Intelligence

Foundations of Biomedical Informatics II, Spring 2023

Lecture 12: Overview of drug discovery and development, small-molecule generation, molecule optimization, identification and characterization of therapeutic targets, high-throughput genetic and chemical perturbations



Marinka Zitnik  
[marinka@hms.harvard.edu](mailto:marinka@hms.harvard.edu)

# Problem set 3 is released

**BMI 702: Foundations in Biomedical Informatics II** Spring 2023  
**Problem Set 3**  
Due Wednesday, May 3rd, 2023, at 23:59pm EST

## General Instructions

The questions in this assignment require careful consideration but do not necessitate lengthy answers. Please provide succinct responses. Additionally, you will need to write code to perform the analysis and generate the data and insights that you will include in your report.

**Submitting answers:** (1) Prepare your answers to the questions and collect them into a single PDF document, which you submit via Canvas by the due date. (2) In addition to the written report, you need to submit your own copy of the Google Colab, [BMI702Assgn3.ipynb](#), with the required code blocks filled in.

**Academic integrity:** Unless noted otherwise, students are expected to complete the assignments individually, not as teams. Discussions about the problem set are encouraged on Canvas; however, solutions sharing is not allowed. Even though students are encouraged to consult websites for coding problems, they may never just copy code. See the course syllabus at <https://zitniklab.ms.harvard.edu/BMI702/syllabus> for more information on Harvard's academic integrity.

**Questions 3 and 4:** In this problem set, we offer two variants of the coding question. Each is worth the same number of points. You can decide which variant you want to solve. You need to select one variant and solve it fully. Mixing-and-matching is not allowed, meaning that you need to solve either (Question 1, Question 2, and Question 3) or (Question 1, Question 2, and Question 4) to get full credits on the problem set.

- Question 3 [No knowledge of PyTorch is needed]: Generating virtual immunofluorescence images
- Question 4 [Familiarity with PyTorch is needed]: Weakly-supervised learning for cancer diagnosis

## 1 Aneurysm Classification [20 points]

An aneurysm is a localized enlargement of an artery. Your lab has developed a new technique for imaging the circulatory system in horizontal slices that allows you to determine the thickness of a patient's blood vessels at a specific height. You have been tasked to create a *supervised ML algorithm* that can predict whether a new patient might have an aneurysm in their femoral artery.

If you feel that you need additional background on medical imaging for cardiovascular diseases, check out "Cardiovascular Diseases (Chapter 13 of Artificial Intelligence in Medical Imaging)" linked in the [Optional Reading section of Lecture 11](#).

**Questions:**

1. (5 points) Considering the task you are solving, what kind of data would you need to begin with? Specifically, describe a dataset (i.e., dataset elements, such as data instances, labels, dataset size, etc.) that you need to obtain in order to solve this task (i.e., train a supervised ML model for predicting aneurysms in femoral arteries).

**BMI 702: Foundations in Biomedical Informatics II - Problem Set 3**  
Page 3 of 7

**BMI 702: Foundations in Biomedical Informatics II - Problem Set 3**  
Page 6 of 7

### BMI 702: Foundations in Biomedical Informatics II - Problem Set 3

**4 Weakly-Supervised Deep Learning for Cancer Diagnosis in Computational Pathology [40 points]**

**Questions 3 and 4:** In this problem set, we offer two variants of the coding question. Each is worth the same number of points. You can decide which variant you want to solve. You need to select one variant and solve it fully. Mixing-and-matching is not allowed, meaning that you need to solve either (Question 1, Question 2, and Question 3) or (Question 1, Question 2, and Question 4) to get full credits on the problem set.

- Question 3 [No knowledge of PyTorch is needed]: Generating virtual immunofluorescence images
- Question 4 [Familiarity with PyTorch is needed]: Weakly-supervised learning for cancer diagnosis

Tissue phenotyping is a fundamental problem in computational pathology (CPATH) to characterize histopathologic features within gigapixel whole-slide images (WSIs) for cancer diagnosis, prognosis, and prediction of treatment response. Unlike natural images, whole-slide imaging is a challenging computer vision domain in which image resolutions can be as large as  $150,000 \times 150,000$  pixels, with many methods using the following three-stage, weakly-supervised framework based on multiple

**BMI 702: Foundations in Biomedical Informatics II - Problem Set 3**  
Page 6 of 7

instance learning (MIL): 1) tissue patching at a single magnification objective ("zoom"), 2) patch-level feature extraction to construct a sequence of embedding instances, and 3) global pooling of instances to construct a slide-level representation for weak-supervision using slide-level labels (e.g., subtype, grade, stage, survival, origin).

The following exercise is based on the Lung Adenocarcinoma (LUAD) vs. Lung Squamous Cell Carcinoma (LUSC) subtyping task used in experimental setup in Lu *et al.* 2021 ([Data-efficient and weakly supervised computational pathology on whole-slide images](#)) and its open-sourced codebase on GitHub, <https://github.com/mahmoodlab/CLAM>, in which we will:

1. Train and evaluate a "naive" MIL algorithm called **AverageMIL**, which takes the average of patch embeddings (as the global pooling operator). A minimalistic but fully-implemented training setup in PyTorch is provided.
2. Implement a more sophisticated algorithm called Attention-Based Multiple Instance Learning (**ABMIL**), which learns attention weights for computing a weighted average of patch embeddings.
3. Compare and contrast **AverageMIL** and **ABMIL**, discussing which algorithm performs better and potential limitations.

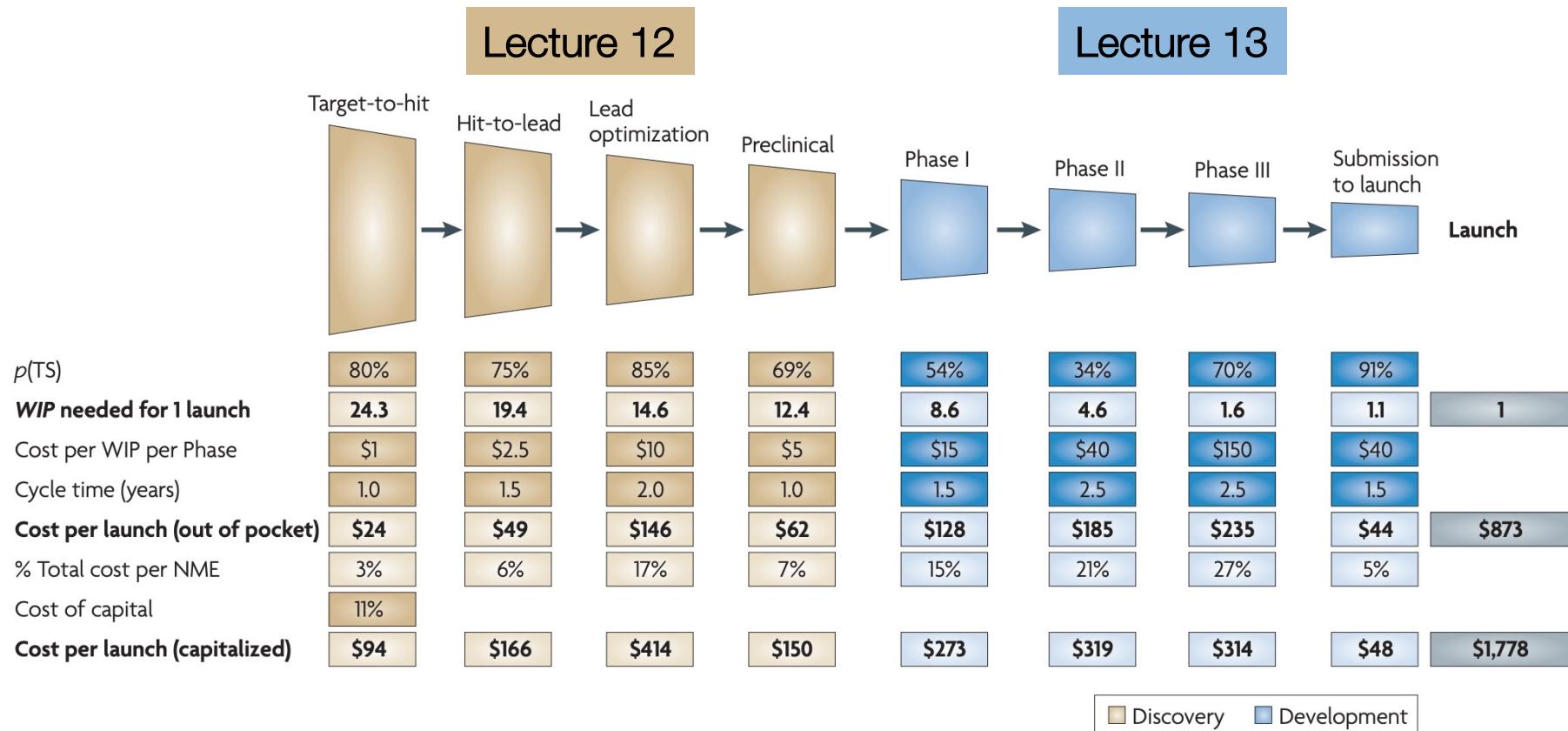
Use the following Google Colab: <https://colab.research.google.com/drive/16i8EI2nHAnYs09H-4Au-ZrvGV7wIVGb?usp=sharing>

**Questions:**

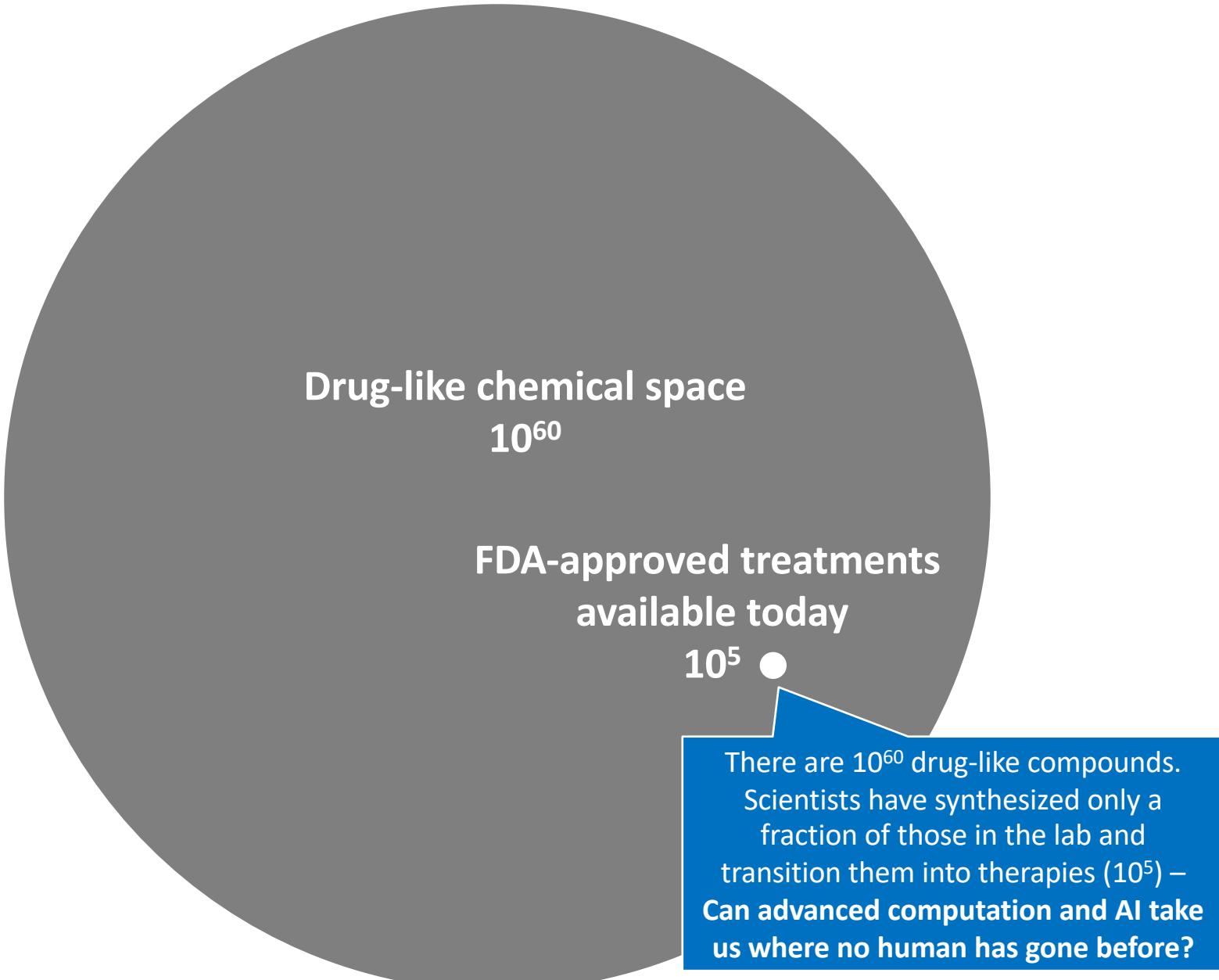
1. (15 points) Implemented in the Google Colab notebook is a minimalistic training setup that performs weakly-supervised learning via **AverageMIL** on LUAD vs. LUSC subtyping using diagnostic H&E tissue slides from the The Cancer Genome Atlas (features already pre-extracted). You can run the cells in the Google Colab Notebook and see how well this algorithm performs in 20 epochs. Though fully implemented, some performance metrics are missing which would help you understand how well this model performs.

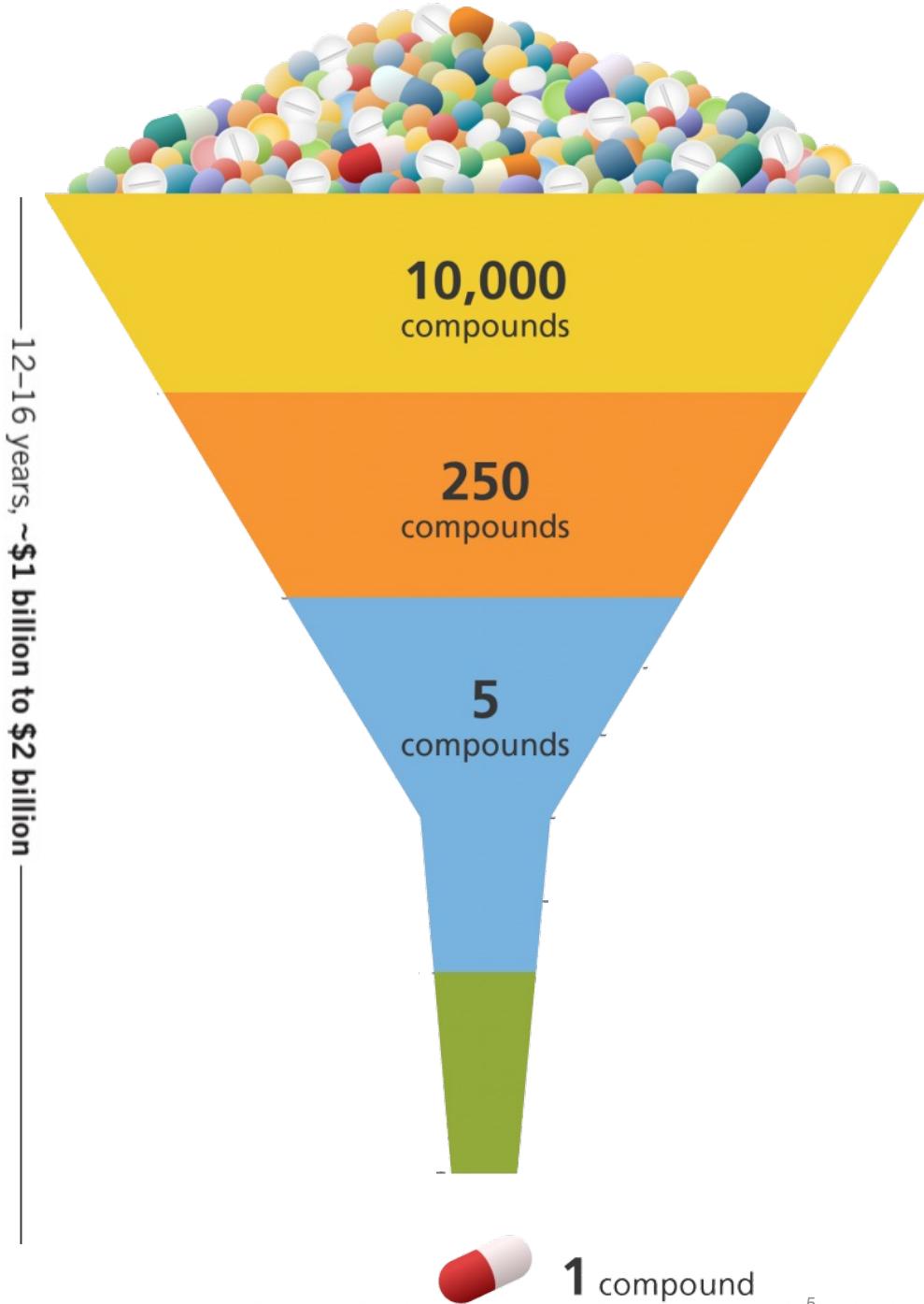
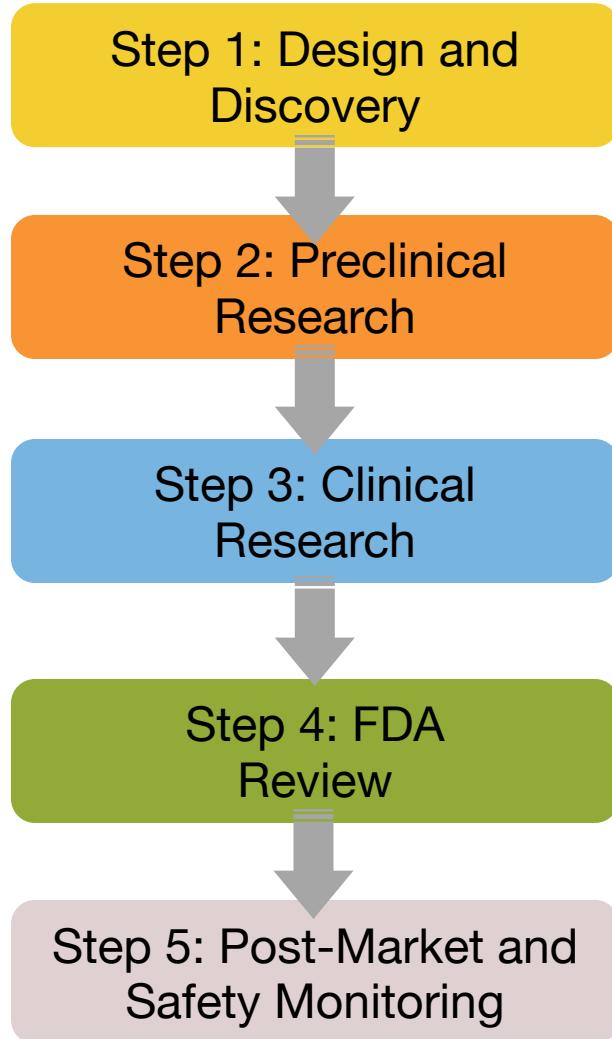
Thematic focus is on clinical language models (LLMs), biomedical imaging, and computational pathology

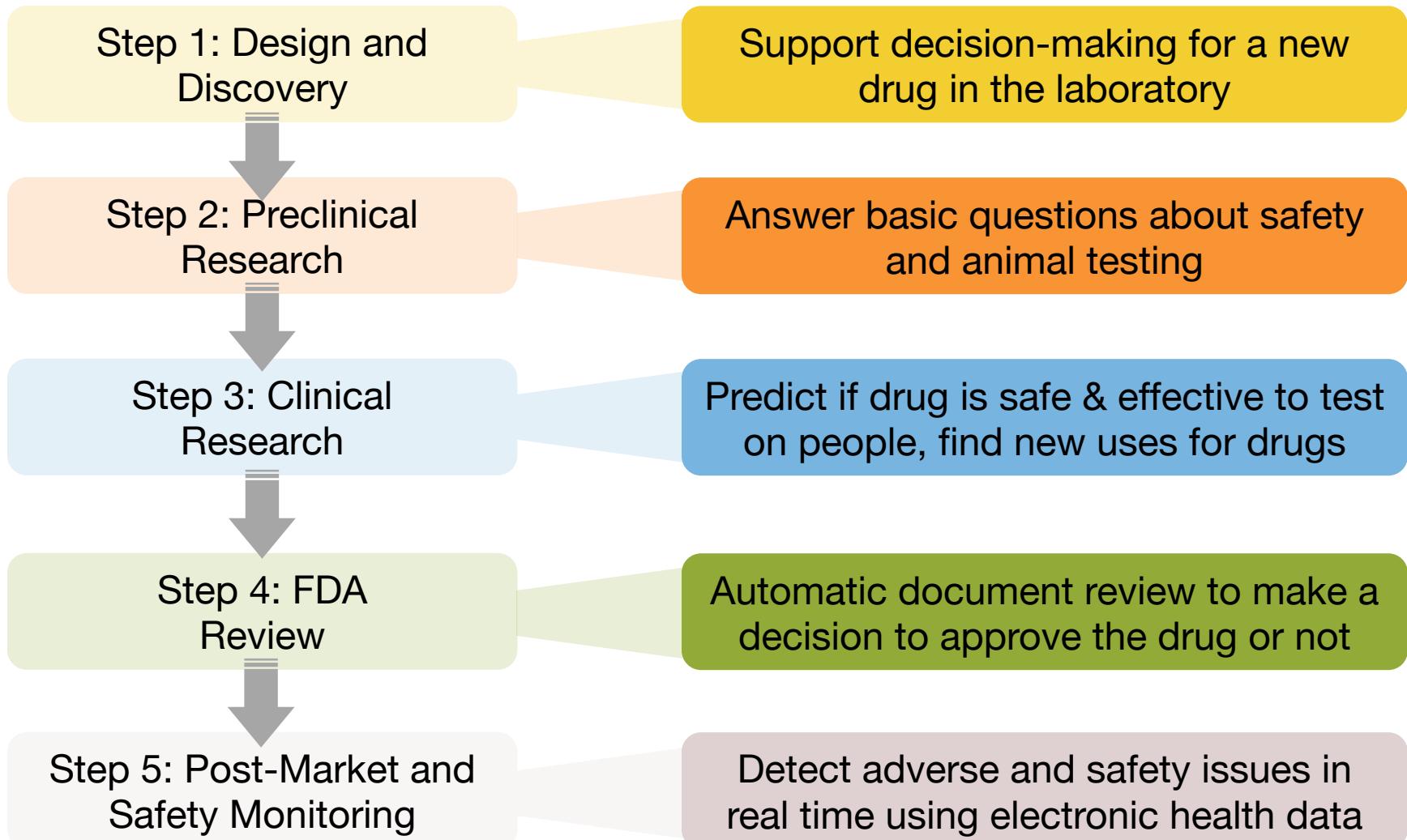
# Distinct phases of drug discovery and development from the initial stage of target-to-hit to the final stage, launch



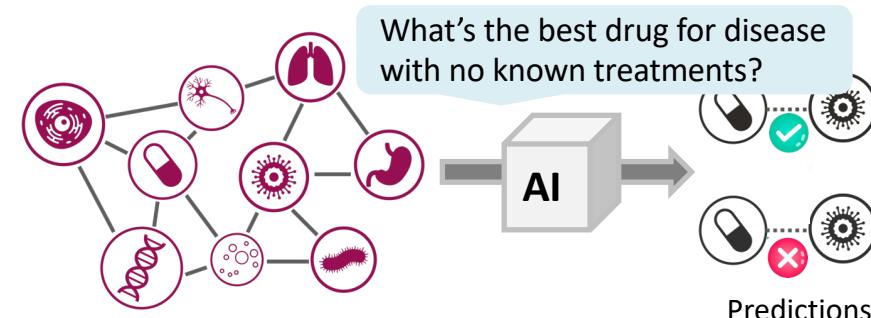
$p(\text{TS})$  – probability of successful transition from one stage to the next; NME – new molecular entity; WIP – work in process



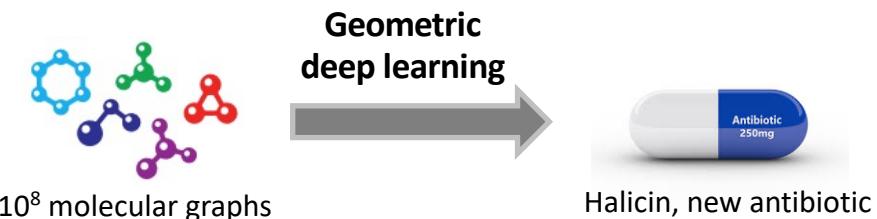




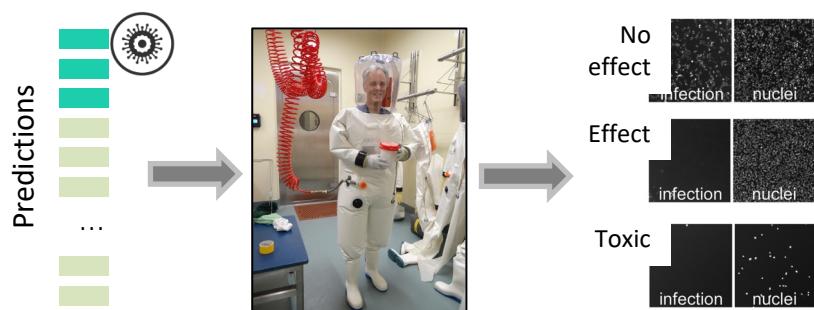
# Examples of success



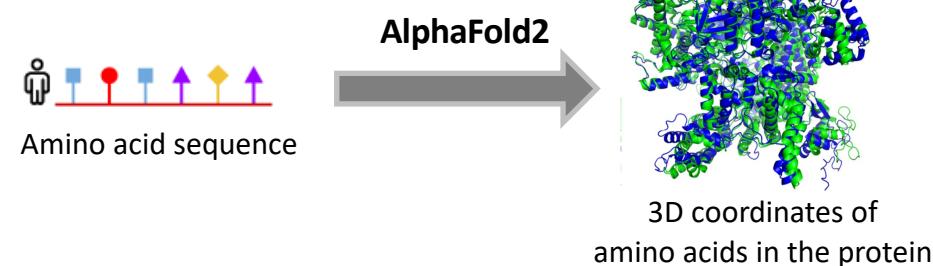
Huang et al., Zero-shot prediction of therapeutic use with geometric deep learning and clinician centered design, 2023



Stokes et al., A Deep Learning Approach to Antibiotic Discovery, 2020



Gysi et al., Network medicine framework for identifying drug-repurposing opportunities for COVID-19, 2021



Jumper et al., Highly accurate protein structure prediction with AlphaFold, 2021

Identify meaningful tasks and datasets

Design AI/ML methods



Biomedical  
scientists



THERAPEUTICS  
DATA COMMONS



AI  
scientists

Data

↔ AI models ↔

Scientific hypotheses



Make drug discovery and development more efficient

Global initiative to access and evaluate AI across therapeutic modalities and stages of drug discovery

210,000 active use cases of AI for therapy design / 90,000 users worldwide



Caltech



BROWN



PRINCETON  
UNIVERSITY



UNIVERSITY OF  
ILLINOIS  
URBANA-CHAMPAIGN



Georgia Institute  
of Technology



DeepMind

Carnegie Mellon



Stanford  
University



Microsoft  
Research



Cornell University

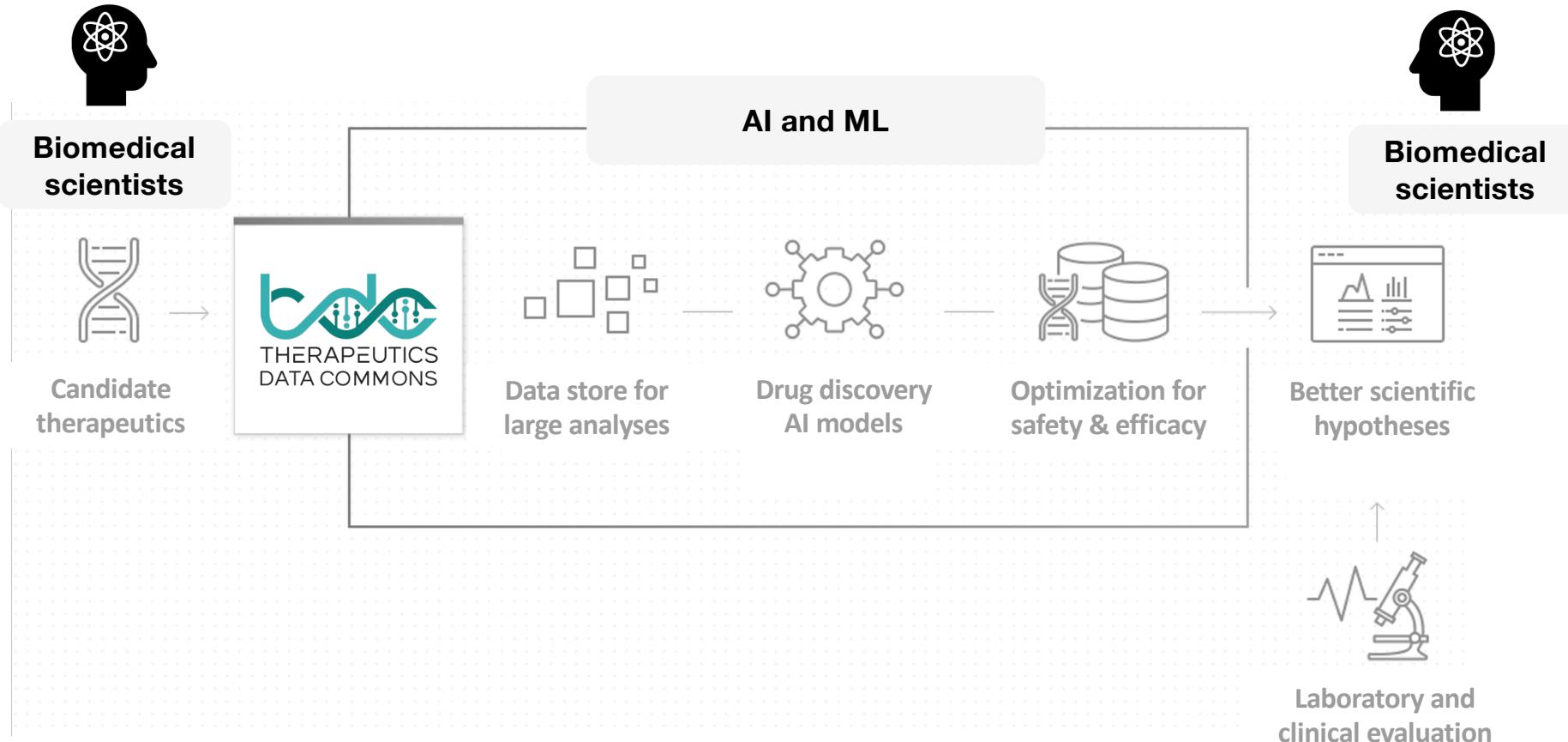


UNIVERSITY OF  
CAMBRIDGE

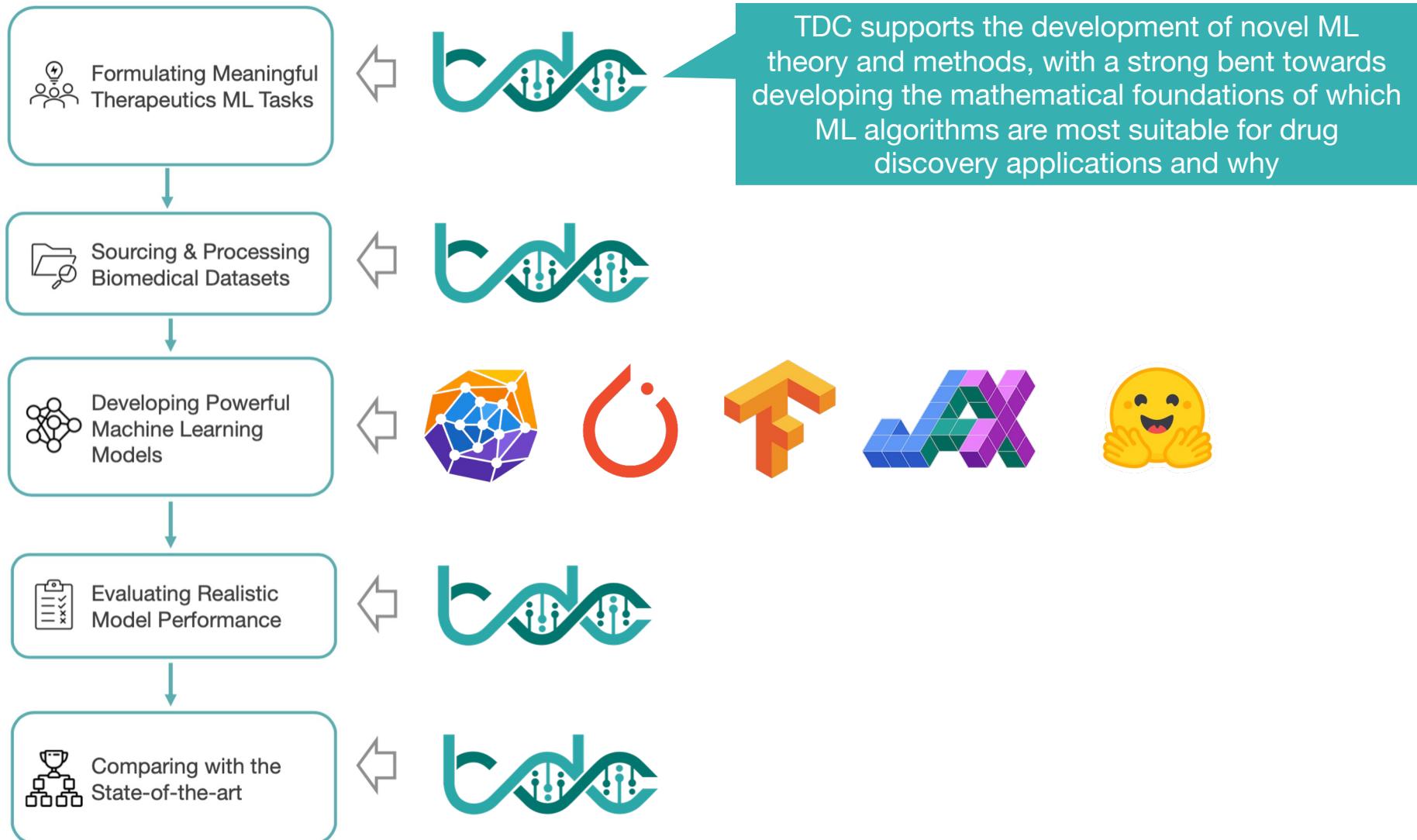
Wang et al., *Nature* '23; Ektefaie et al., *Nature Machine Intelligence* '23; McDermott et al., *Nature Machine Intelligence* '23;

Li et al., *Nature Biomedical Engineering* '22

# Therapeutics Commons: How it works?

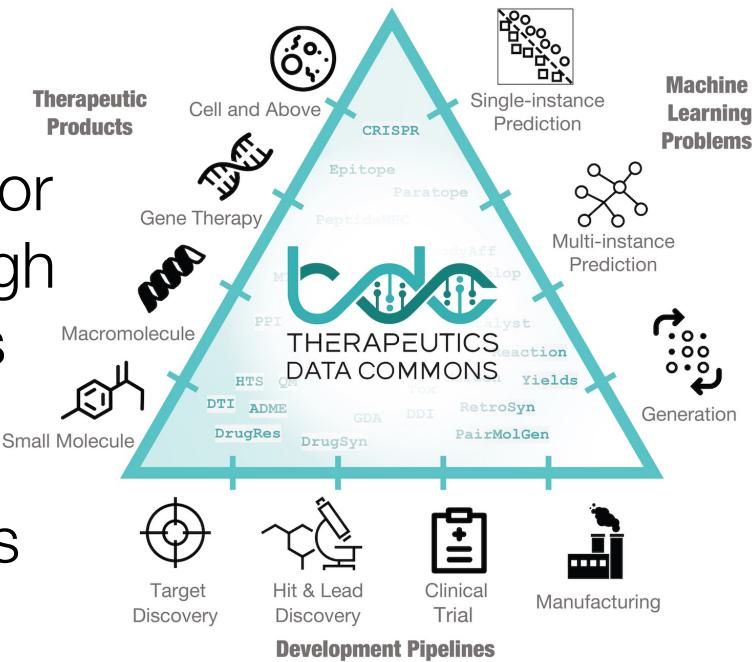


# AI workflows in drug discovery



# What tasks can we address with these workflows?

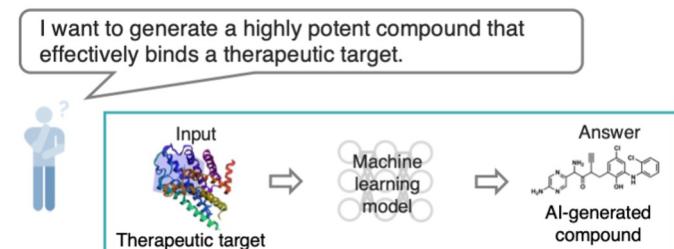
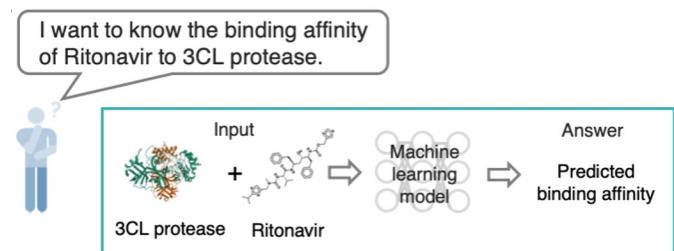
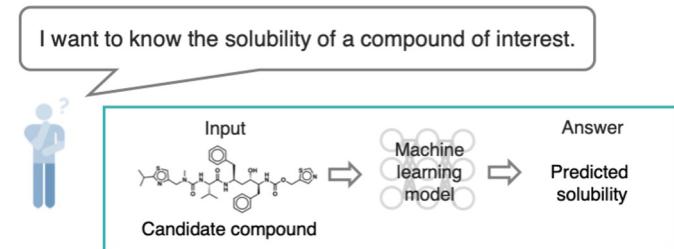
- Target discovery
  - Identify candidate drug targets
- Activity modeling
  - Screen and generate individual or combinatorial therapies with high binding activity towards targets
- Efficacy and safety
  - Optimize therapeutic signatures predictive of safety & efficacy
- Manufacturing
  - Synthesis of therapeutics



# Outline for today's class

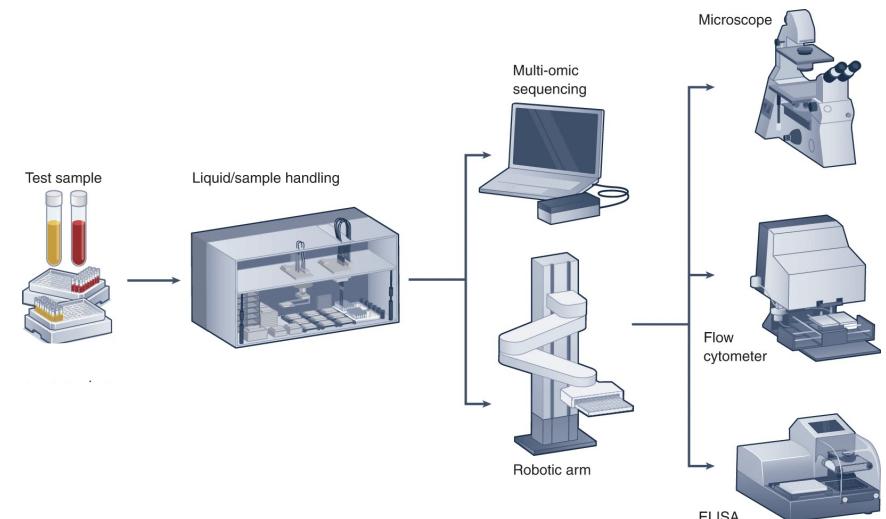


- Optimization & generation of small molecules
- Binding of drugs to therapeutic targets
- High-throughput genetic & chemical perturbations



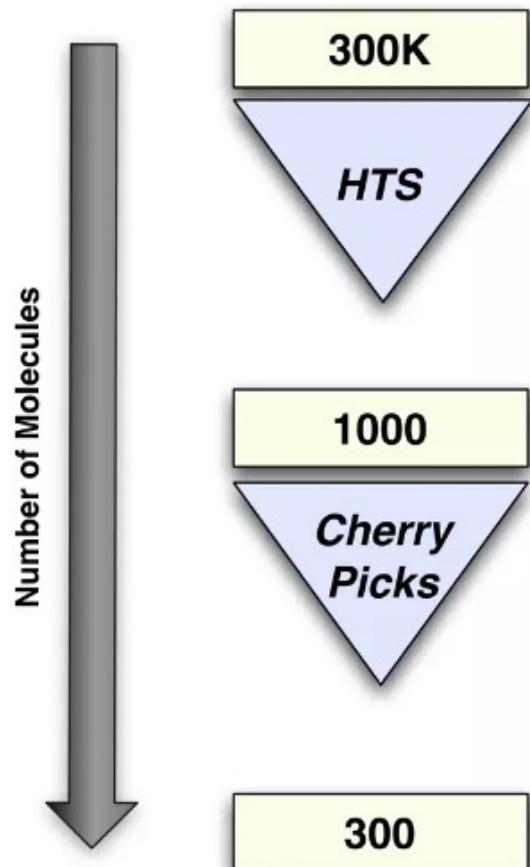
# High throughput screening (HTS)

- Test thousands to hundreds of thousands of compounds in one or more assays
  - Biochemical, genetic, and pharmacological assays
- Integrate with robotics for self-driving lab
- **Goal:** Rapidly identify novel modulators of biological systems
  - Cellular basis of diseases
  - Therapeutic agents



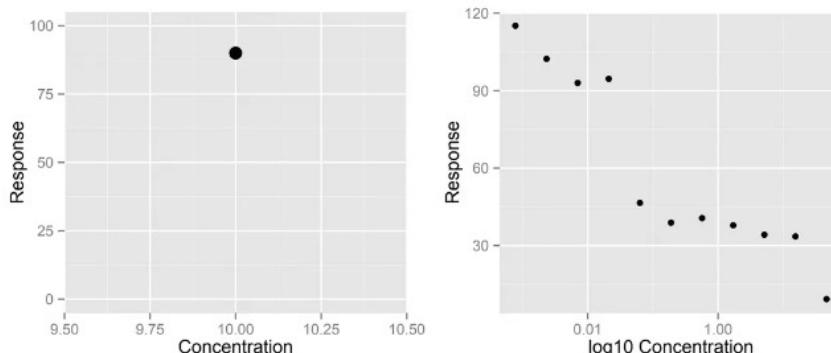
# Goals of high throughput screening

- Rapidly screen large collections of compounds (chemical libraries)
- Efficiently identify active compounds
  - Test them in slower, accurate, expensive screens
- Use the data to learn what types of compounds tend to be active

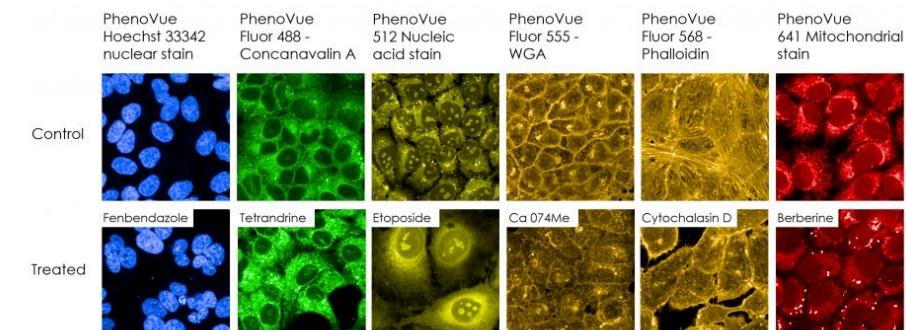


# HTS data types

- Categorical: active/inactive or toxic/nontoxic
- Continuous: single-point or dose-response
- Multiple readouts:
  - Might read at different wavelengths or time points
  - More complex when dealing with images



Single-point vs. dose-response readouts



Cell painting for phenotypic drug discovery

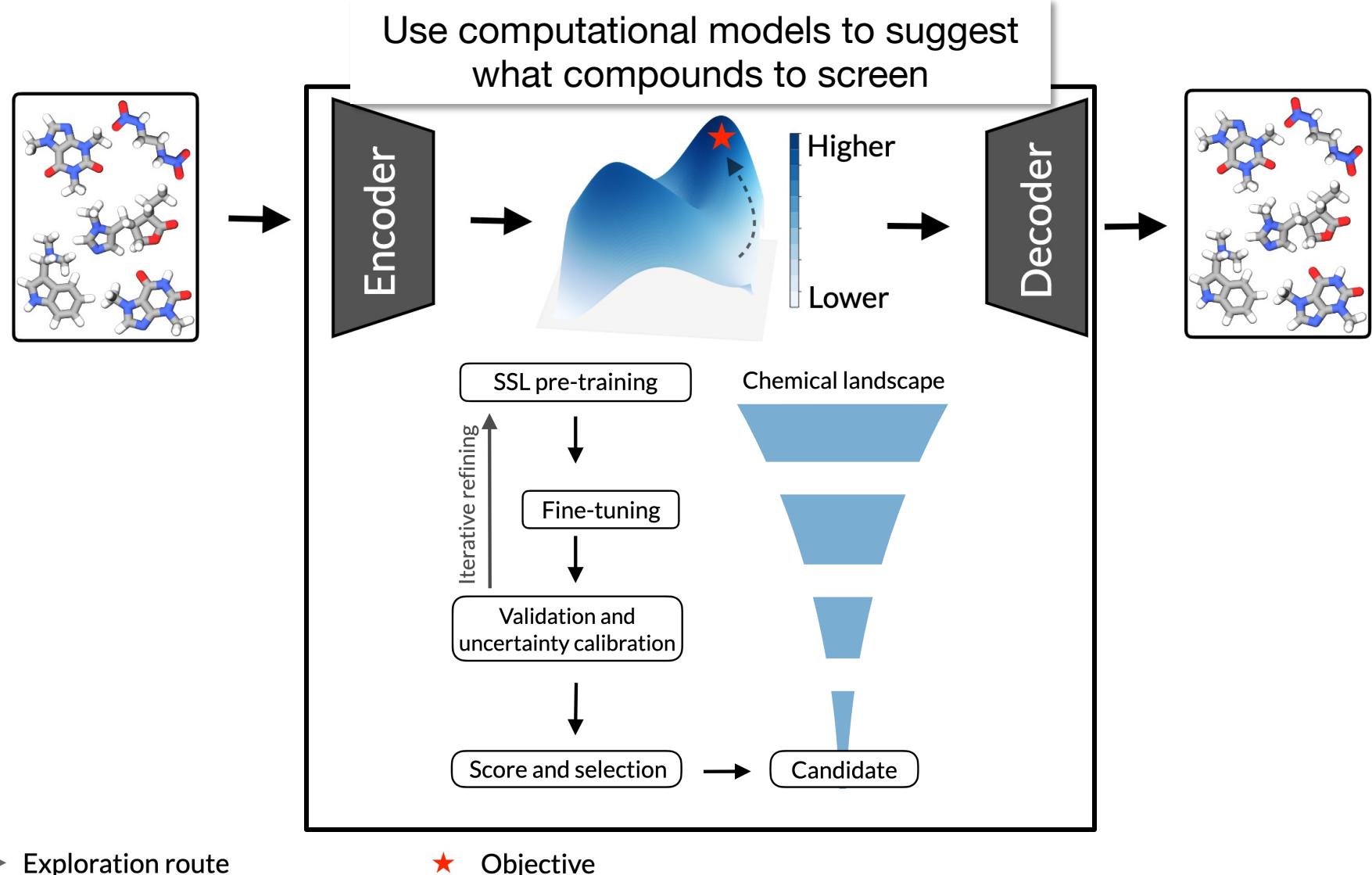
# HTS: Machine learning setup

- HTS tests the activity of molecules:

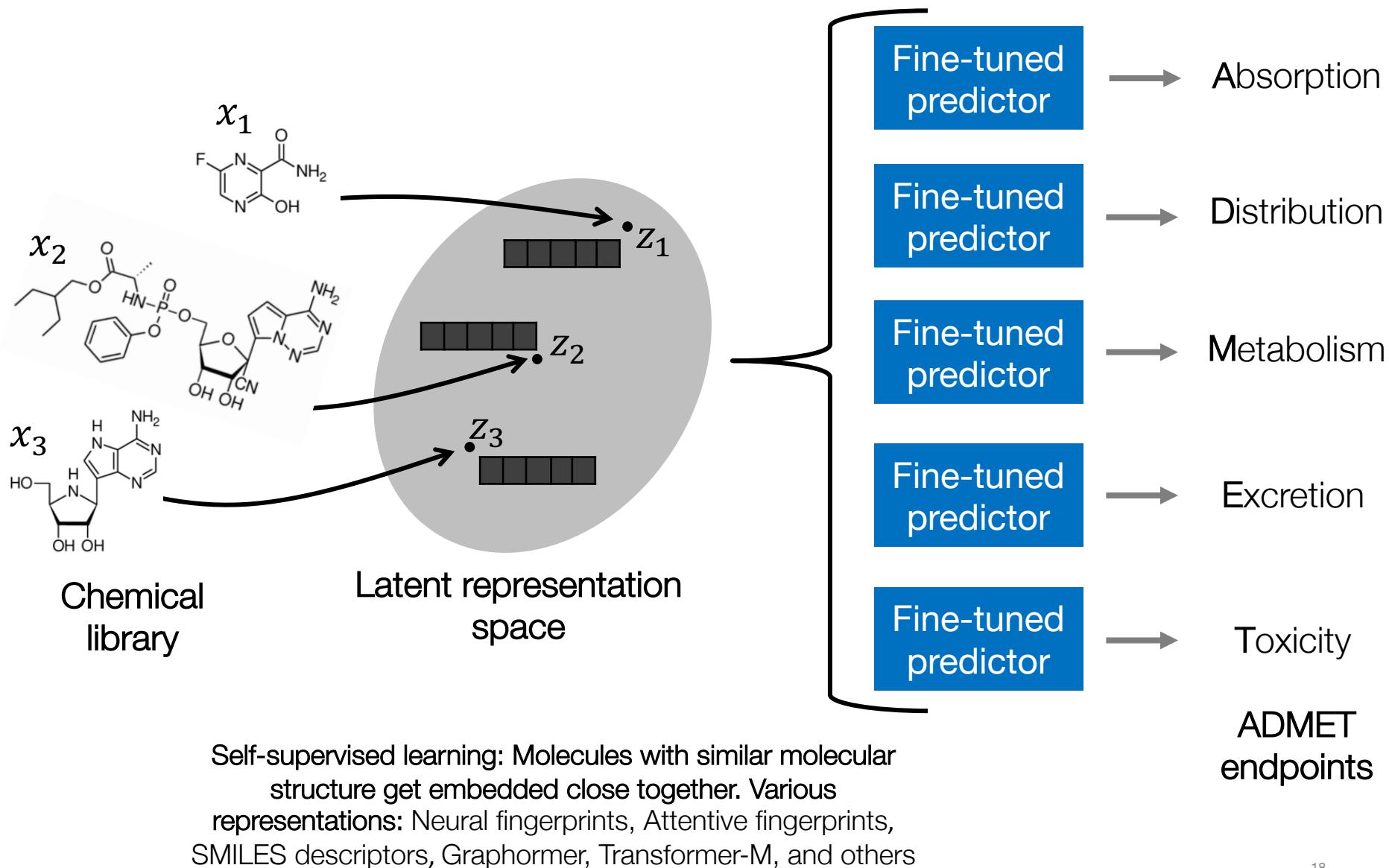
$$\text{Activity} = f(\text{Structure})$$

- We need to describe the molecular structure
  - Various discrete or real-valued descriptors
  - Surfaces (3D)
  - Binary fingerprints
  - Learned molecular embeddings

# In-silico screening and optimization of molecular structure



# Molecular property prediction

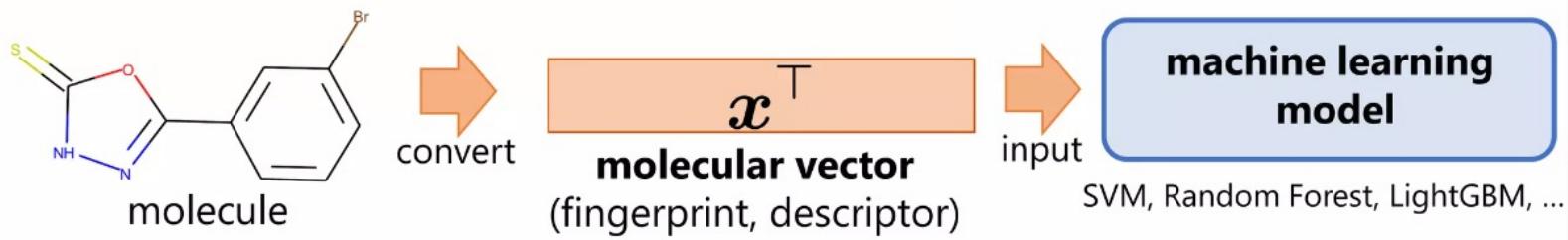


# What can we use molecular representations for?

- **Search**
  - Given a potent active molecule, find similar ones (or dissimilar but also potent)
- **Prediction of various endpoints**
  - Given a set of active and inactive molecules, build a model to predict which members from a chemical library will be active
- **Clustering**
  - Given a set of molecules, do they cluster into structurally different groups?

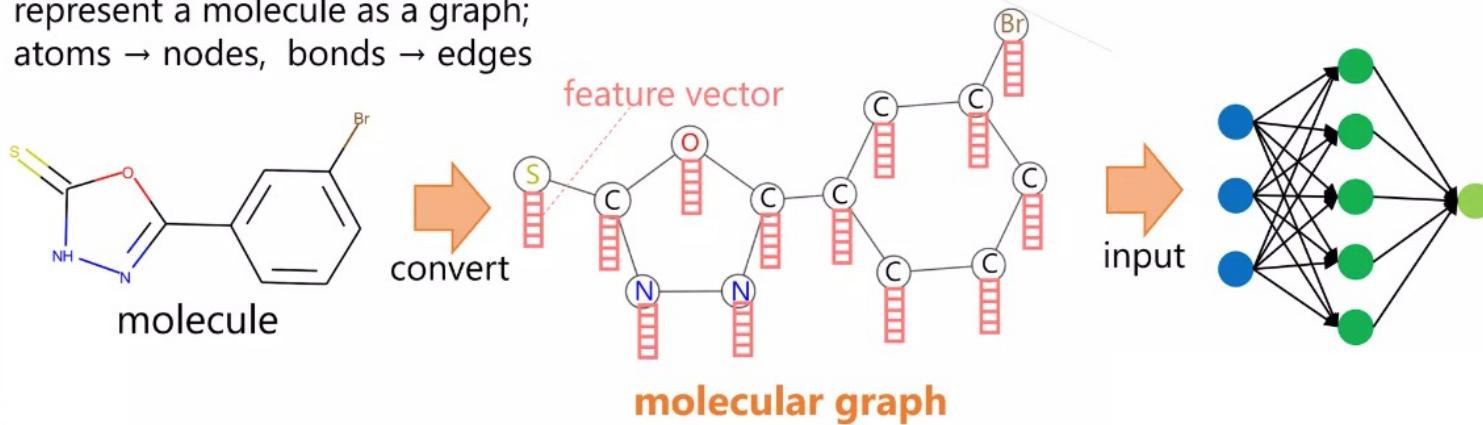
# Two strategies for producing molecular representations

## Traditional approach

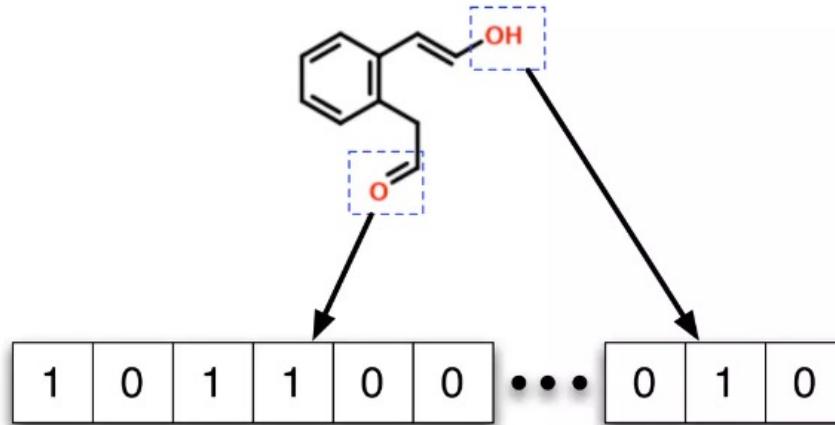


## Graph convolutional network (GCN) approach

represent a molecule as a graph;  
atoms → nodes, bonds → edges



# Fingerprint representations



- Lots of types of fingerprints
- Keyed fingerprints indicate the presence or absence of a structural feature
- Length can vary from 166 to 4096 bits or more
- Fingerprints usually compared to each other using the Tanimoto metric

# Towards neural fingerprints

---

**Algorithm 1** Circular fingerprints
 

---

```

1: Input: molecule, radius  $R$ , fingerprint length  $S$ 
2: Initialize: fingerprint vector  $\mathbf{f} \leftarrow \mathbf{0}_S$ 
3: for each atom  $a$  in molecule
4:    $\mathbf{r}_a \leftarrow g(a)$             $\triangleright$  lookup atom features
5: for  $L = 1$  to  $R$             $\triangleright$  for each layer
6:   for each atom  $a$  in molecule
7:      $\mathbf{r}_1 \dots \mathbf{r}_N = \text{neighbors}(a)$ 
8:      $\mathbf{v} \leftarrow [\mathbf{r}_a, \mathbf{r}_1, \dots, \mathbf{r}_N]$      $\triangleright$  concatenate
9:      $\mathbf{r}_a \leftarrow \text{hash}(\mathbf{v})$             $\triangleright$  hash function
10:     $i \leftarrow \text{mod}(r_a, S)$        $\triangleright$  convert to index
11:     $\mathbf{f}_i \leftarrow 1$                   $\triangleright$  Write 1 at index
12: Return: binary vector  $\mathbf{f}$ 
  
```

---

**Algorithm 2** Neural graph fingerprints
 

---

```

1: Input: molecule, radius  $R$ , hidden weights  $H_1^1 \dots H_R^5$ , output weights  $W_1 \dots W_R$ 
2: Initialize: fingerprint vector  $\mathbf{f} \leftarrow \mathbf{0}_S$ 
3: for each atom  $a$  in molecule
4:    $\mathbf{r}_a \leftarrow g(a)$             $\triangleright$  lookup atom features
5: for  $L = 1$  to  $R$             $\triangleright$  for each layer
6:   for each atom  $a$  in molecule
7:      $\mathbf{r}_1 \dots \mathbf{r}_N = \text{neighbors}(a)$ 
8:      $\mathbf{v} \leftarrow \mathbf{r}_a + \sum_{i=1}^N \mathbf{r}_i$             $\triangleright$  sum
9:      $\mathbf{r}_a \leftarrow \sigma(\mathbf{v} H_L^N)$             $\triangleright$  smooth function
10:     $\mathbf{i} \leftarrow \text{softmax}(\mathbf{r}_a W_L)$        $\triangleright$  sparsify
11:     $\mathbf{f} \leftarrow \mathbf{f} + \mathbf{i}$                   $\triangleright$  add to fingerprint
12: Return: real-valued vector  $\mathbf{f}$ 
  
```

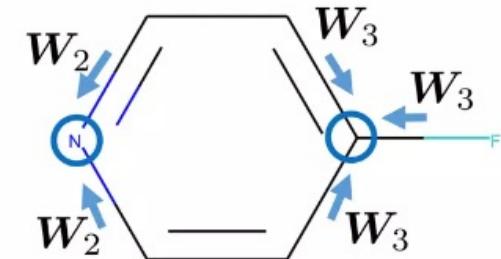
---

Figure 2: Pseudocode of circular fingerprints (*left*) and neural graph fingerprints (*right*). Differences are highlighted in blue. Every non-differentiable operation is replaced with a differentiable analog.

# Neural fingerprint representations

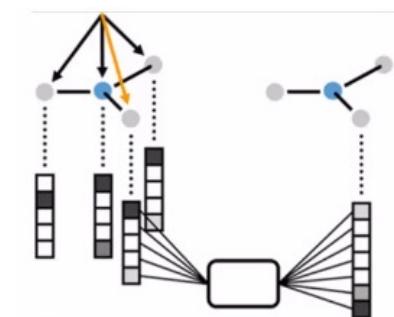
## 1) Neural graph fingerprints

- Generate molecular fingerprints with a neural network
- Update atom features using only adjacent atoms
- Use different weights for node degrees



## 2) Molecular graphs

- Update atom features by convolutional and pooling layers using adjacent atoms

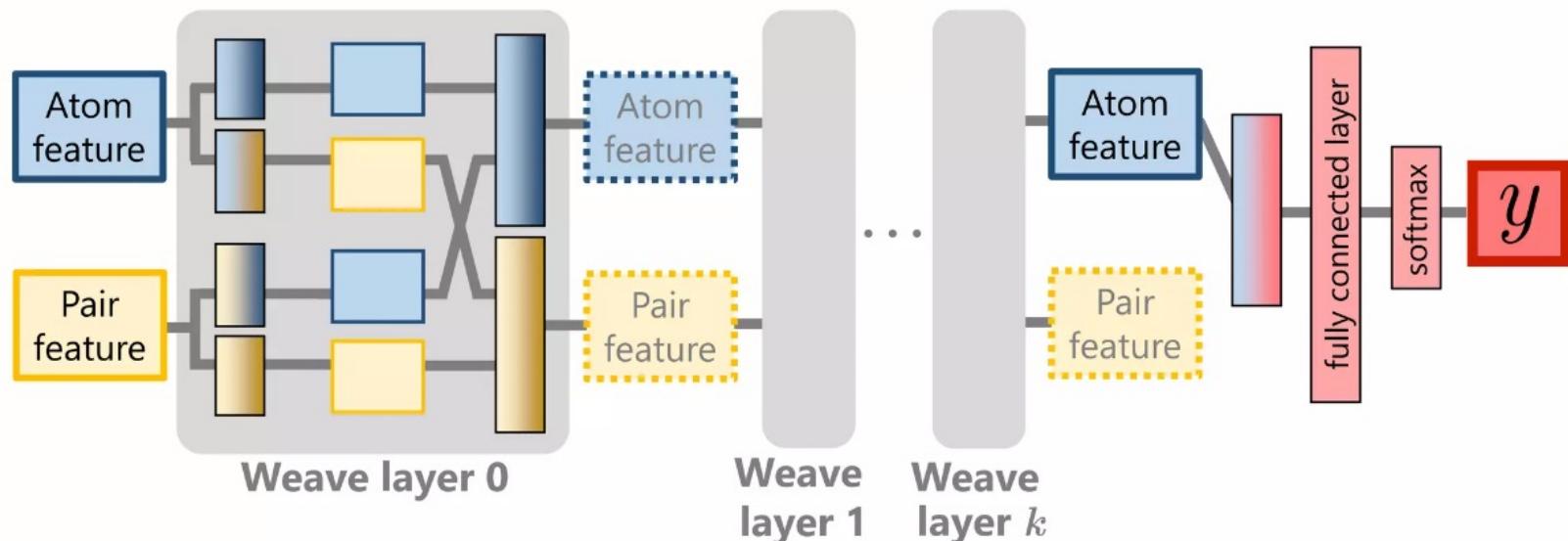


- They did not consider property of edges (bonds)
- They did not consider atoms other than 1-neighbor

# Neural fingerprint representations

## 3) Weave module

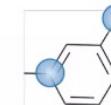
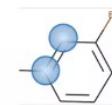
- Weave module considers also atom pairs



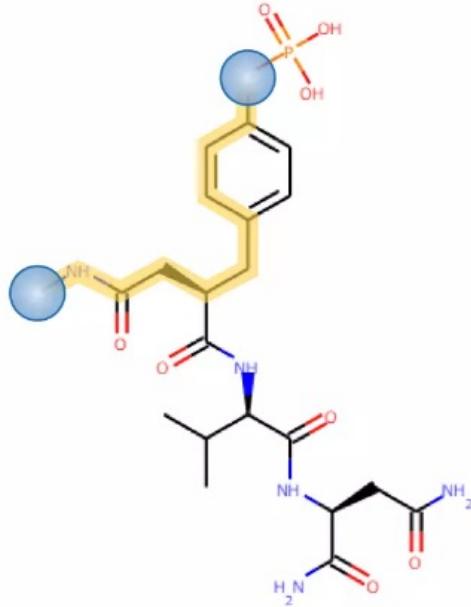
Information of distant atoms can be considered (not just 1-neighbors)



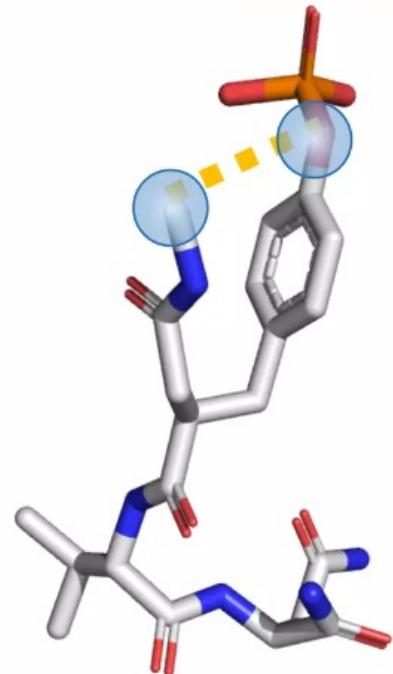
**The difference in distance between atom pairs was not considered**



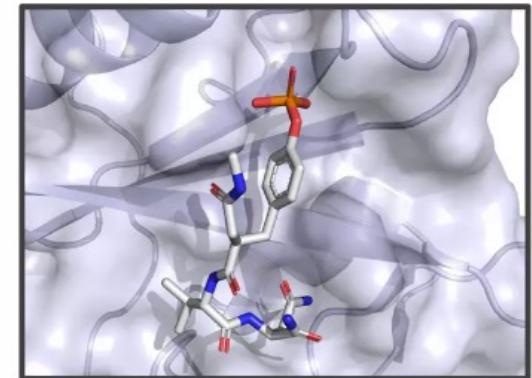
# Graphs vs. 3D structures



Molecular Graph



3D Structure



The distance on the graph does not necessarily correlate with the Euclidean distance between atoms on the 3D structure

Need to consider modifying the definition  
of graph distance

# Datasets

22 datasets with ADMET endpoints

## A: Absorption

Caco2 (Cell Permeability)  
HIA (Intestinal Absorption)  
Pgp (P-glycoprotein)  
Bioavailability  
Lipophilicity  
Solubility

## D: Distribution

BBB (Blood-Brain Barrier)  
PPBR (Plasma Protein Binding)  
VDss (Volume of Distribution)

## M: Metabolism

CYP2C9/2D6/3A4 Inhibition  
CYP2C9/2D6/3A4 Substrate

## E: Excretion

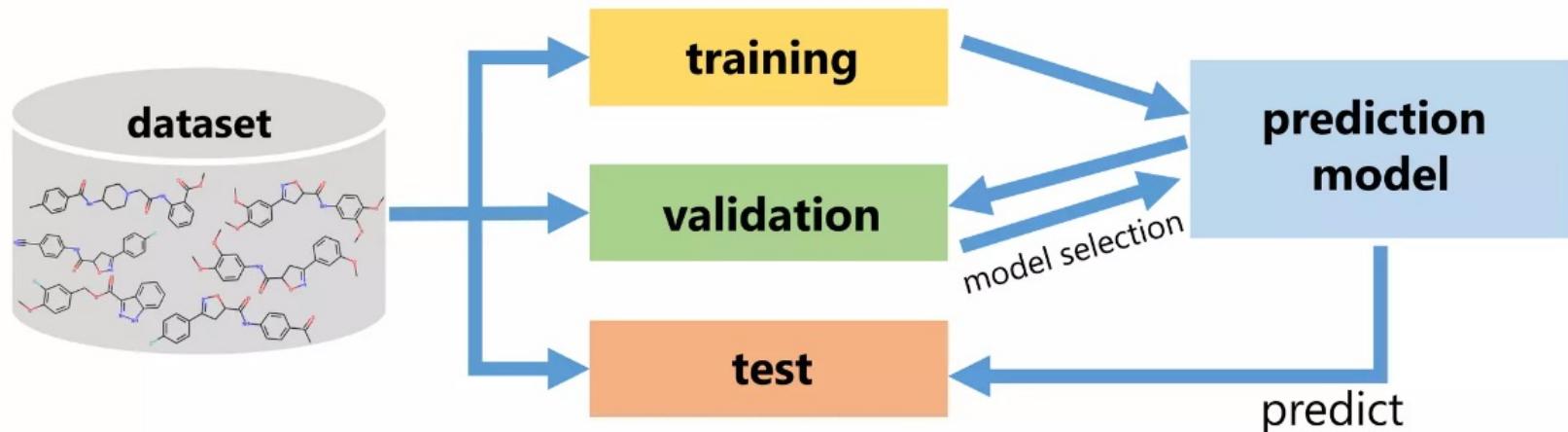
Half Life  
Clearance (Hepatocyte)  
Clearance (Microsome)

## T: Toxicity

LD50 (Acute Toxicity)  
hERG blocker  
Ames Mutagenicity  
Drug Induced Liver Injury



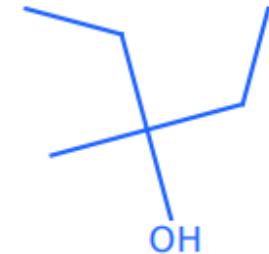
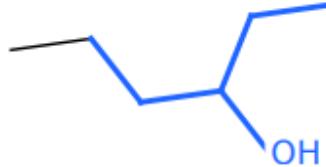
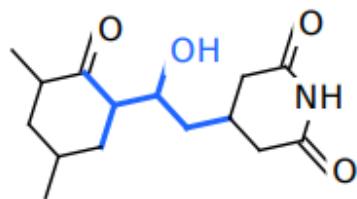
# Experimental setup



- Demonstrate that fingerprints are interpretable
  - Show substructures which most activate individual features in a fingerprint vector
  - **Fingerprint features** can each only be activated by a single fragment of a single radius, except for accidental collisions
  - In contrast, **neural fingerprint features** can be activated by variations of the same structure, making them more interpretable, and allowing shorter feature vectors.

# Results: Examining neural fingerprints

Fragments most activated by pro-solubility feature



---

Fragments most activated by anti-solubility feature

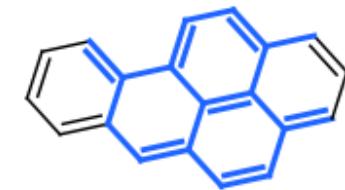
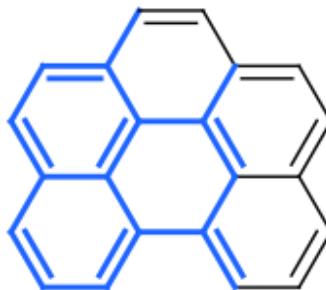
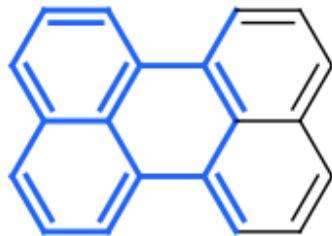
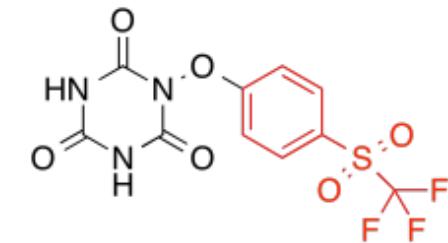
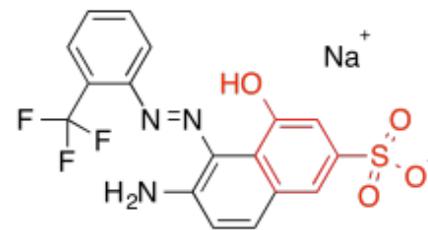
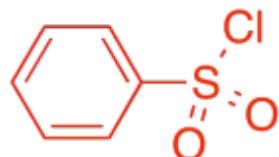


Figure 4: Examining fingerprints optimized for predicting solubility. Shown here are representative examples of molecular fragments (highlighted in blue) which most activate different features of the fingerprint. *Top row:* The feature most predictive of solubility. *Bottom row:* The feature most predictive of insolubility.

# Results: Examining neural fingerprints

Fragments most activated by toxicity feature on SR-MMP dataset



Fragments most activated by toxicity feature on NR-AHR dataset

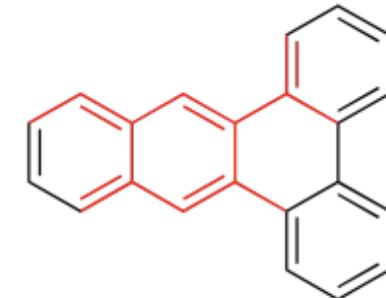
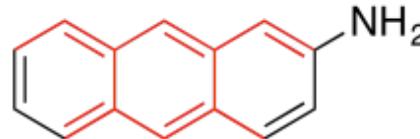
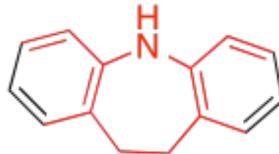


Figure 5: Visualizing fingerprints optimized for predicting toxicity. Shown here are representative samples of molecular fragments (highlighted in red) which most activate the feature most predictive of toxicity. *Top row:* the most predictive feature identifies groups containing a sulphur atom attached to an aromatic ring. *Bottom row:* the most predictive feature identifies fused aromatic rings, also known as polycyclic aromatic hydrocarbons, a well-known carcinogen.

# Results: Molecular property prediction

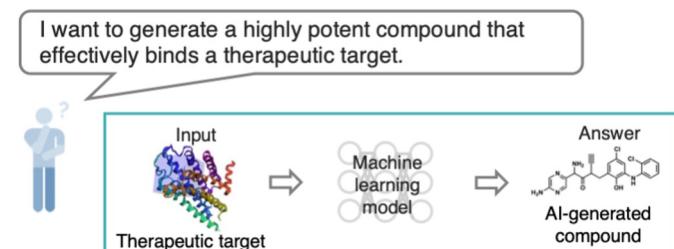
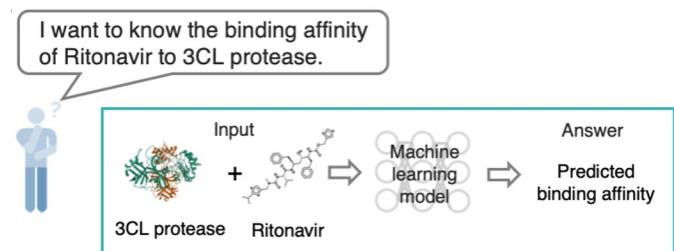
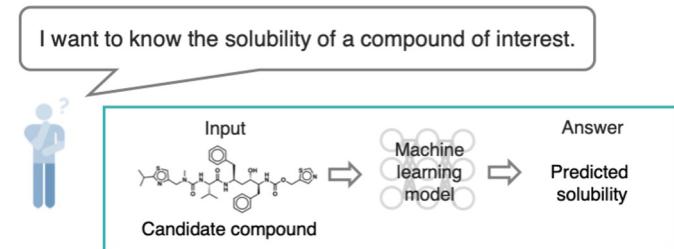
Raw Feature Type		Expert-Curated Methods		SMILES	Molecular Graph-Based Methods (state-of-the-Art in ML)				
Dataset	Metric	Morgan	RDKit2D	CNN	NeuralFP	GCN	AttentiveFP	AttrMasking	ContextPred
	# Params.	1477K	633K	227K	480K	192K	301K	2067K	2067K
<b>TDC.Caco2</b> (↓)	MAE	$0.908 \pm 0.060$	<b><math>0.393 \pm 0.024</math></b>	$0.446 \pm 0.036$	$0.530 \pm 0.102$	$0.599 \pm 0.104$	$0.401 \pm 0.032$	$0.546 \pm 0.052$	$0.502 \pm 0.036$
<b>TDC.HIA</b> (↑)	AUROC	$0.807 \pm 0.072$	$0.972 \pm 0.008$	$0.869 \pm 0.026$	$0.943 \pm 0.014$	$0.936 \pm 0.024$	$0.974 \pm 0.007$	<b><math>0.978 \pm 0.006</math></b>	$0.975 \pm 0.004$
<b>TDC.Pgp</b> (↑)	AUROC	$0.880 \pm 0.006$	$0.918 \pm 0.007$	$0.908 \pm 0.012$	$0.902 \pm 0.020$	$0.895 \pm 0.021$	$0.892 \pm 0.012$	<b><math>0.929 \pm 0.006</math></b>	$0.923 \pm 0.005$
<b>TDC.Bioav</b> (↑)	AUROC	$0.581 \pm 0.086$	<b><math>0.672 \pm 0.021</math></b>	$0.613 \pm 0.013$	$0.632 \pm 0.036$	$0.566 \pm 0.115$	$0.632 \pm 0.039$	$0.577 \pm 0.087$	$0.671 \pm 0.026$
<b>TDC.Lipo</b> (↓)	MAE	$0.701 \pm 0.009$	$0.574 \pm 0.017$	$0.743 \pm 0.020$	$0.563 \pm 0.023$	<u><math>0.541 \pm 0.011</math></u>	$0.572 \pm 0.007$	$0.547 \pm 0.024$	<b><math>0.535 \pm 0.012</math></b>
<b>TDC.AqSol</b> (↓)	MAE	$1.203 \pm 0.019$	$0.827 \pm 0.047$	$1.023 \pm 0.023$	$0.947 \pm 0.016$	$0.907 \pm 0.020$	<b><math>0.776 \pm 0.008</math></b>	$1.026 \pm 0.020$	$1.040 \pm 0.045$
<b>TDC.BBB</b> (↑)	AUROC	$0.823 \pm 0.015$	$0.889 \pm 0.016$	$0.781 \pm 0.030$	$0.836 \pm 0.009$	$0.842 \pm 0.016$	$0.855 \pm 0.011$	$0.892 \pm 0.012$	<b><math>0.897 \pm 0.004</math></b>
<b>TDC.PPBR</b> (↓)	MAE	$12.848 \pm 0.362$	$9.994 \pm 0.319$	$11.106 \pm 0.358$	<b><math>9.292 \pm 0.384</math></b>	$10.194 \pm 0.373$	<u><math>9.373 \pm 0.335</math></u>	$10.075 \pm 0.202$	$9.445 \pm 0.224$
<b>TDC.VD</b> (↑)	Spearman	$0.493 \pm 0.011$	<b><math>0.561 \pm 0.025</math></b>	$0.226 \pm 0.114$	$0.258 \pm 0.162$	$0.457 \pm 0.050$	$0.241 \pm 0.145$	<u><math>0.559 \pm 0.019</math></u>	$0.485 \pm 0.092$
<b>TDC.CYP2D6-I</b> (↑)	AUPRC	$0.587 \pm 0.011$	$0.616 \pm 0.007$	$0.544 \pm 0.053$	$0.627 \pm 0.009$	$0.616 \pm 0.020$	$0.646 \pm 0.014$	<u><math>0.721 \pm 0.009</math></u>	<b><math>0.739 \pm 0.005</math></b>
<b>TDC.CYP3A4-I</b> (↑)	AUPRC	$0.827 \pm 0.009$	$0.829 \pm 0.007$	$0.821 \pm 0.003$	$0.849 \pm 0.004$	$0.840 \pm 0.010$	$0.851 \pm 0.006$	<u><math>0.902 \pm 0.002</math></u>	<b><math>0.904 \pm 0.002</math></b>
<b>TDC.CYP2C9-I</b> (↑)	AUPRC	$0.715 \pm 0.004$	$0.742 \pm 0.006$	$0.713 \pm 0.006$	$0.739 \pm 0.010$	$0.735 \pm 0.004$	$0.749 \pm 0.004$	<u><math>0.829 \pm 0.003</math></u>	<b><math>0.839 \pm 0.003</math></b>
<b>TDC.CYP2D6-S</b> (↑)	AUPRC	$0.671 \pm 0.066$	$0.677 \pm 0.047$	$0.485 \pm 0.037$	$0.572 \pm 0.062$	$0.617 \pm 0.039$	$0.574 \pm 0.030$	<u><math>0.704 \pm 0.028</math></u>	<b><math>0.736 \pm 0.024</math></b>
<b>TDC.CYP3A4-S</b> (↑)	AUROC	$0.633 \pm 0.013$	$0.639 \pm 0.012$	<b><math>0.662 \pm 0.031</math></b>	$0.578 \pm 0.020$	$0.590 \pm 0.023$	$0.576 \pm 0.025$	$0.582 \pm 0.021$	$0.609 \pm 0.025$
<b>TDC.CYP2C9-S</b> (↑)	AUPRC	$0.380 \pm 0.015$	$0.360 \pm 0.040$	$0.367 \pm 0.059$	$0.359 \pm 0.059$	$0.344 \pm 0.051$	$0.375 \pm 0.032$	<u><math>0.381 \pm 0.045</math></u>	<b><math>0.392 \pm 0.026</math></b>
<b>TDC.Half_Life</b> (↑)	Spearman	<b><math>0.329 \pm 0.083</math></b>	$0.184 \pm 0.111$	$0.038 \pm 0.138$	$0.177 \pm 0.165$	<u><math>0.239 \pm 0.100</math></u>	$0.085 \pm 0.068$	$0.151 \pm 0.068$	$0.129 \pm 0.114$
<b>TDC.CL-Micro</b> (↑)	Spearman	$0.492 \pm 0.020$	<b><math>0.586 \pm 0.014</math></b>	$0.252 \pm 0.116$	$0.529 \pm 0.015$	$0.532 \pm 0.033$	$0.365 \pm 0.055$	<u><math>0.585 \pm 0.034</math></u>	$0.578 \pm 0.007$
<b>TDC.CL-Hepa</b> (↑)	Spearman	$0.272 \pm 0.068$	$0.382 \pm 0.007$	$0.235 \pm 0.021$	$0.401 \pm 0.037$	$0.366 \pm 0.063$	$0.289 \pm 0.022$	<u><math>0.413 \pm 0.028</math></u>	<b><math>0.439 \pm 0.026</math></b>
<b>TDC.hERG</b> (↑)	AUROC	$0.736 \pm 0.023$	<b><math>0.841 \pm 0.020</math></b>	$0.754 \pm 0.037$	$0.722 \pm 0.034$	$0.738 \pm 0.038$	<u><math>0.825 \pm 0.007</math></u>	$0.778 \pm 0.046$	$0.756 \pm 0.023$
<b>TDC.AMES</b> (↑)	AUROC	$0.794 \pm 0.008$	$0.823 \pm 0.011$	$0.776 \pm 0.015$	$0.823 \pm 0.006$	$0.818 \pm 0.010$	$0.814 \pm 0.008$	<b><math>0.842 \pm 0.008</math></b>	$0.837 \pm 0.009$
<b>TDC.DILI</b> (↑)	AUROC	$0.832 \pm 0.021$	$0.875 \pm 0.019$	$0.792 \pm 0.016$	$0.851 \pm 0.026$	$0.859 \pm 0.033$	<u><math>0.886 \pm 0.015</math></u>	<b><math>0.919 \pm 0.008</math></b>	$0.861 \pm 0.018$
<b>TDC.LD50</b> (↓)	MAE	$0.649 \pm 0.019$	$0.678 \pm 0.003$	$0.675 \pm 0.011$	$0.667 \pm 0.020$	$0.649 \pm 0.026$	$0.678 \pm 0.012$	<b><math>0.685 \pm 0.025</math></b>	$0.669 \pm 0.030$

- No single method performs the best across all scenarios
- Pre-training boost performance
- Pre-trained graph models yield strongest predictors overall

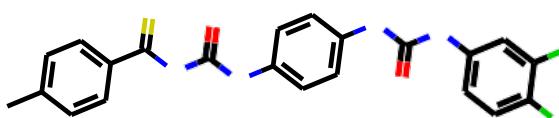
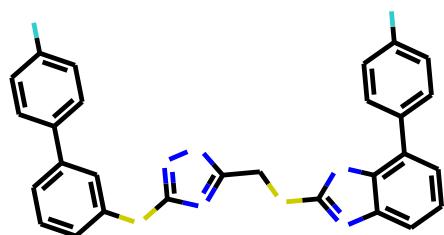
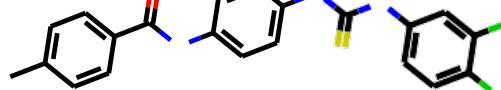
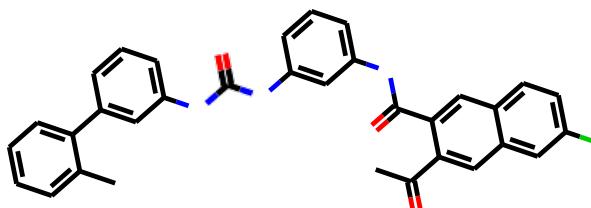
# Outline for today's class



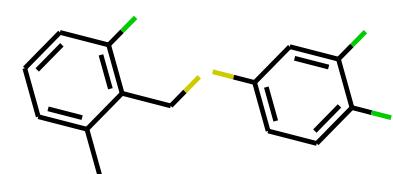
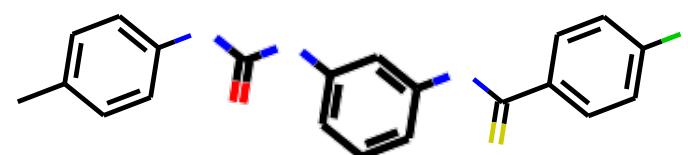
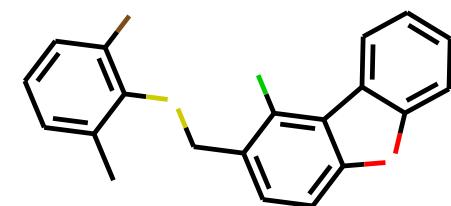
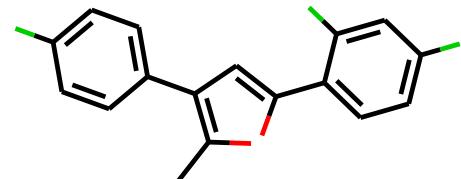
- Optimization & generation of small molecules
- Binding of drugs to therapeutic targets
- High-throughput genetic & chemical perturbations



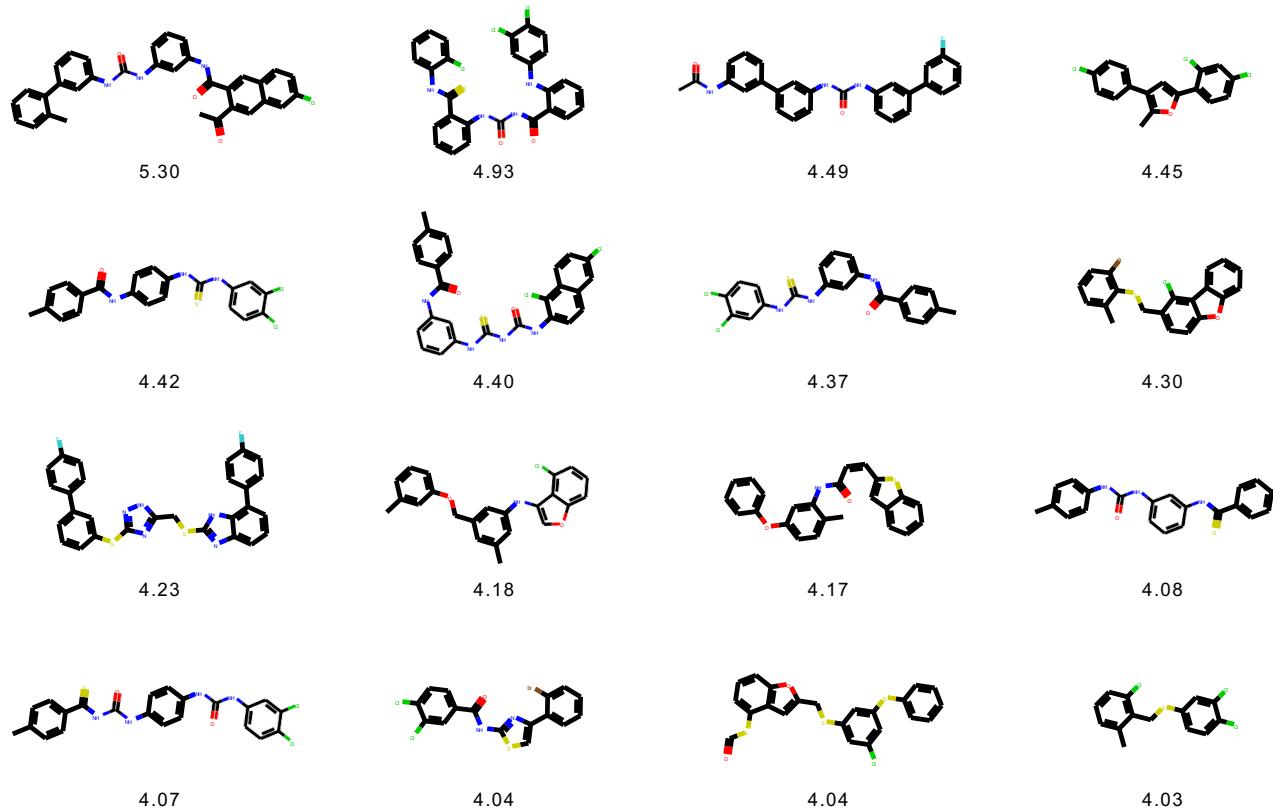
# Molecular graph generation



## Generative model

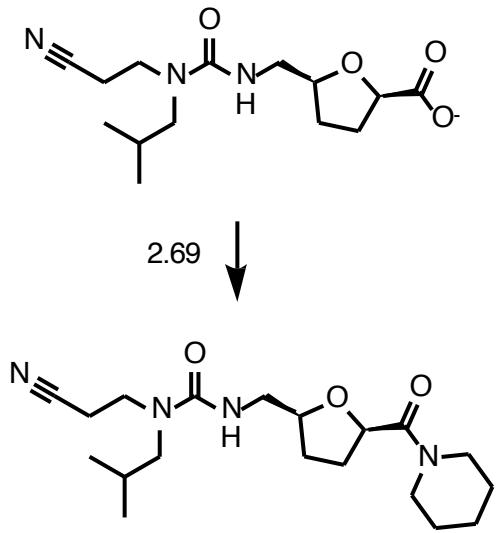
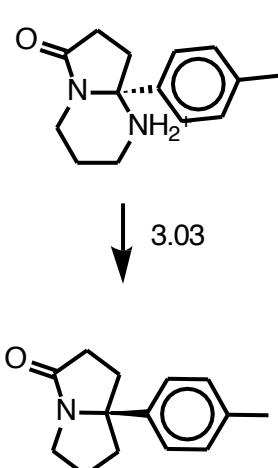
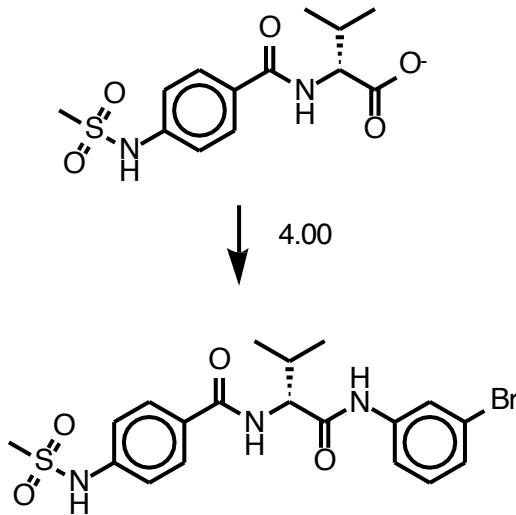


# Molecular graph generation



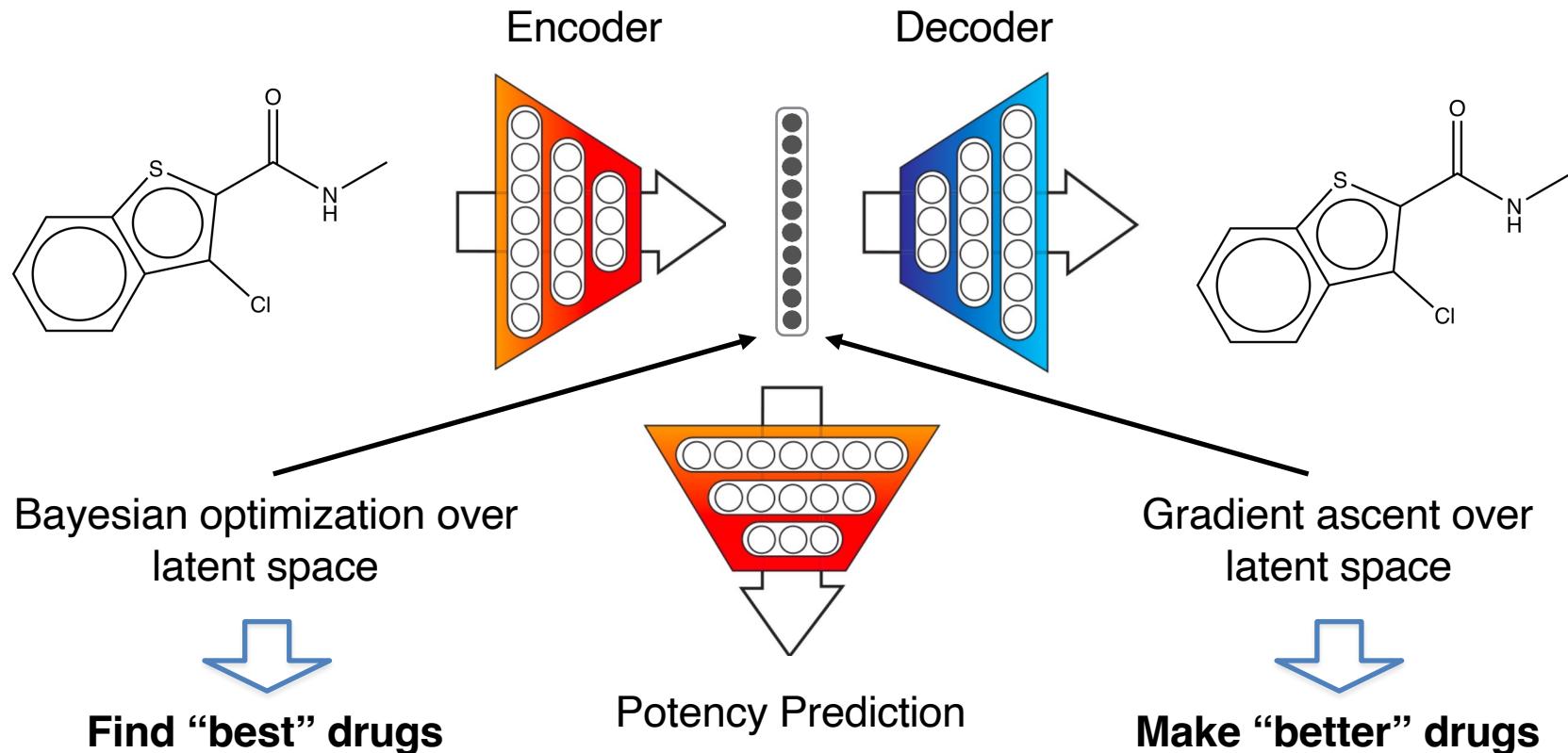
Generate molecules with high potency

# Molecular graph generation



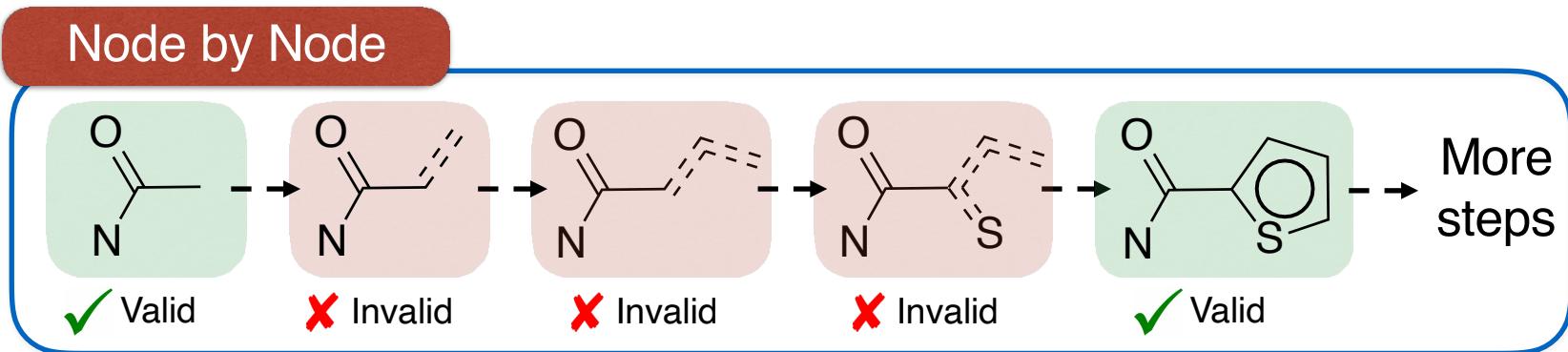
Modify molecules to increase potency

# Molecular variational autoencoder



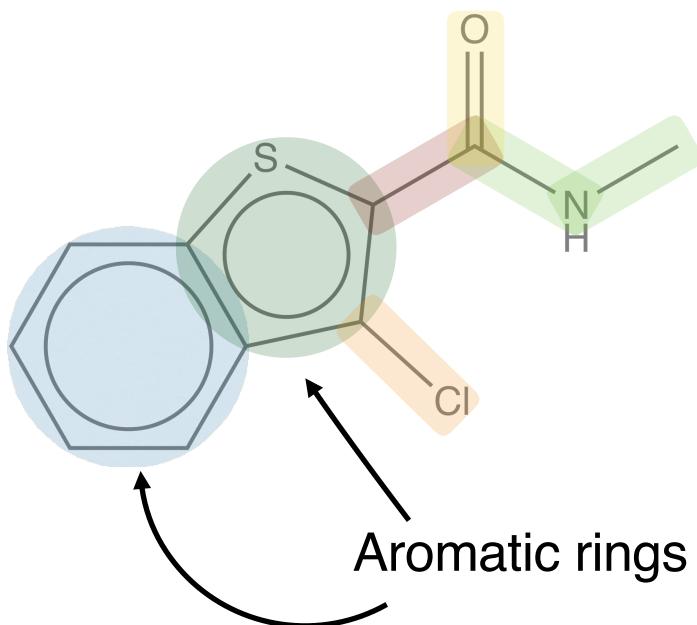
[1] Gomez-Bombarelli et al., Automatic chemical design using a data-driven continuous representation of molecules, 2016

# How to generate graphs?

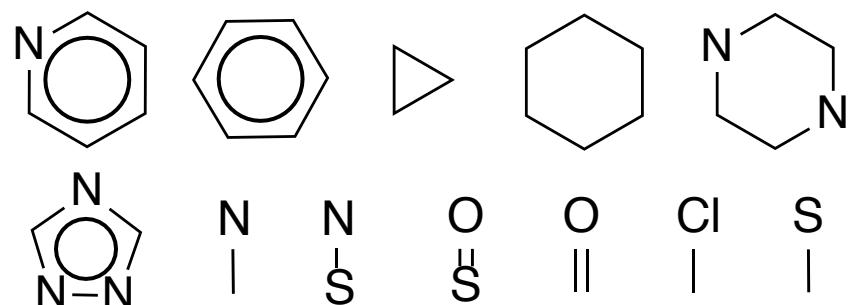


- Not every graphs is chemically valid
- Invalid intermediate states → hard to validate
- Very long intermediate steps → difficult to train (Li et al., 2018)

# Functional groups

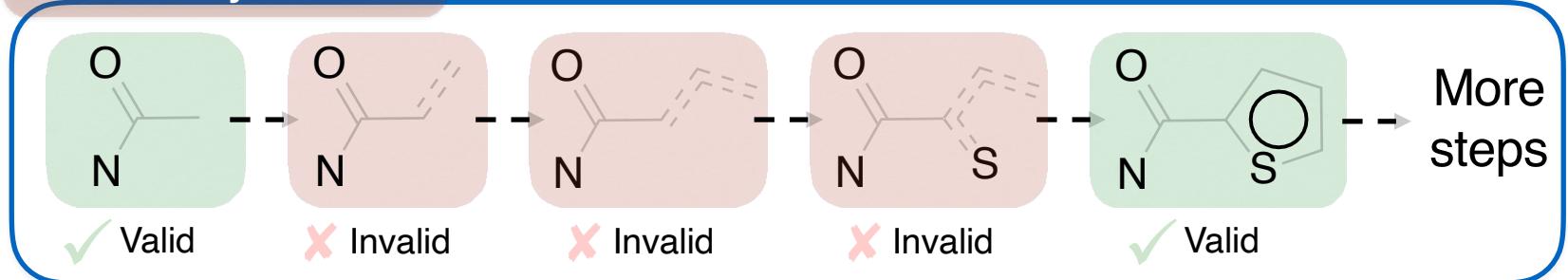


## Functional Groups

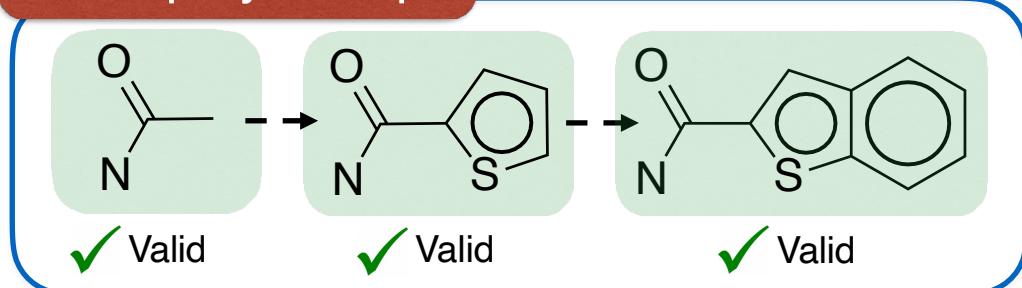


# How to generate graphs?

## Node by Node

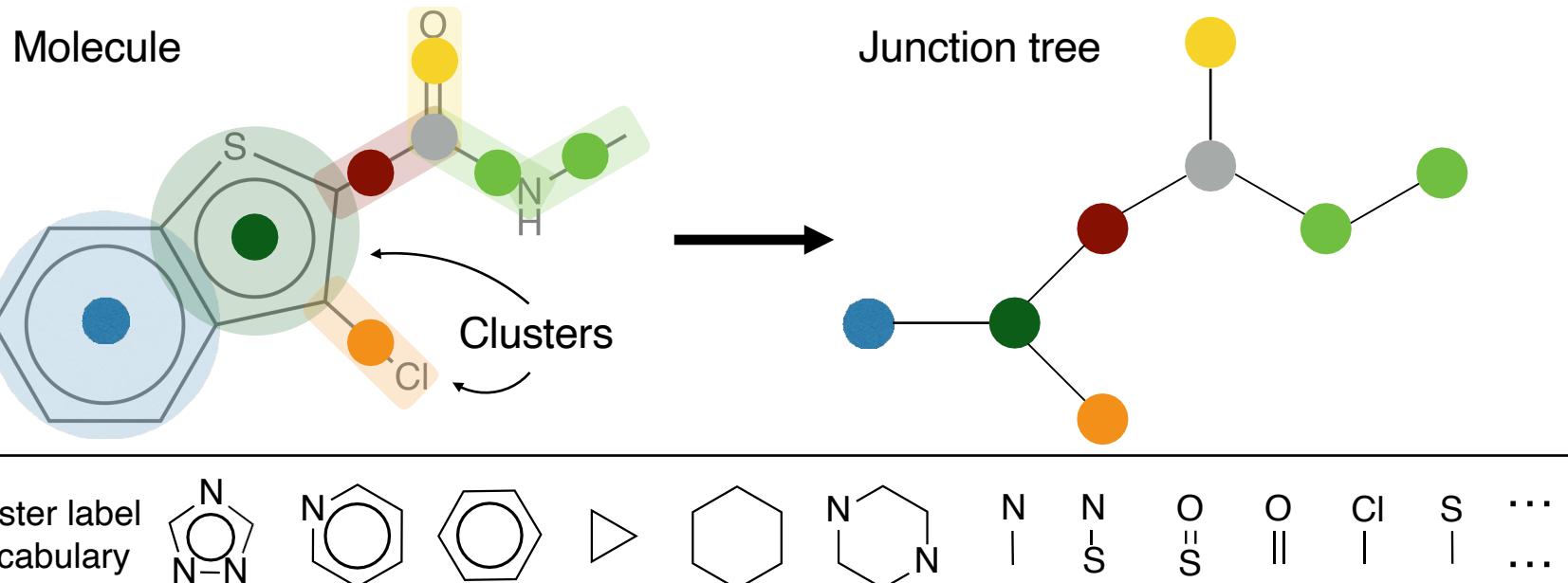


## Group by Group



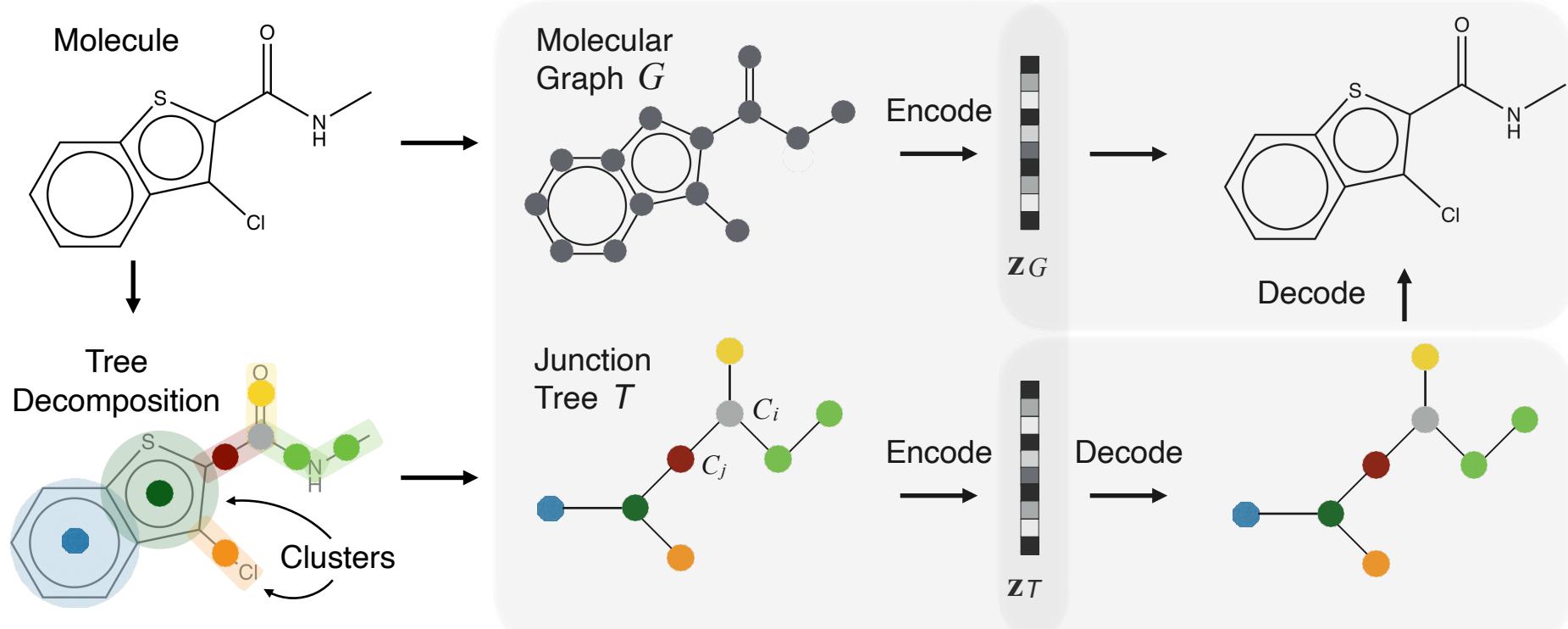
- Shorter action sequence
- Easy to check validity

# Tree decomposition

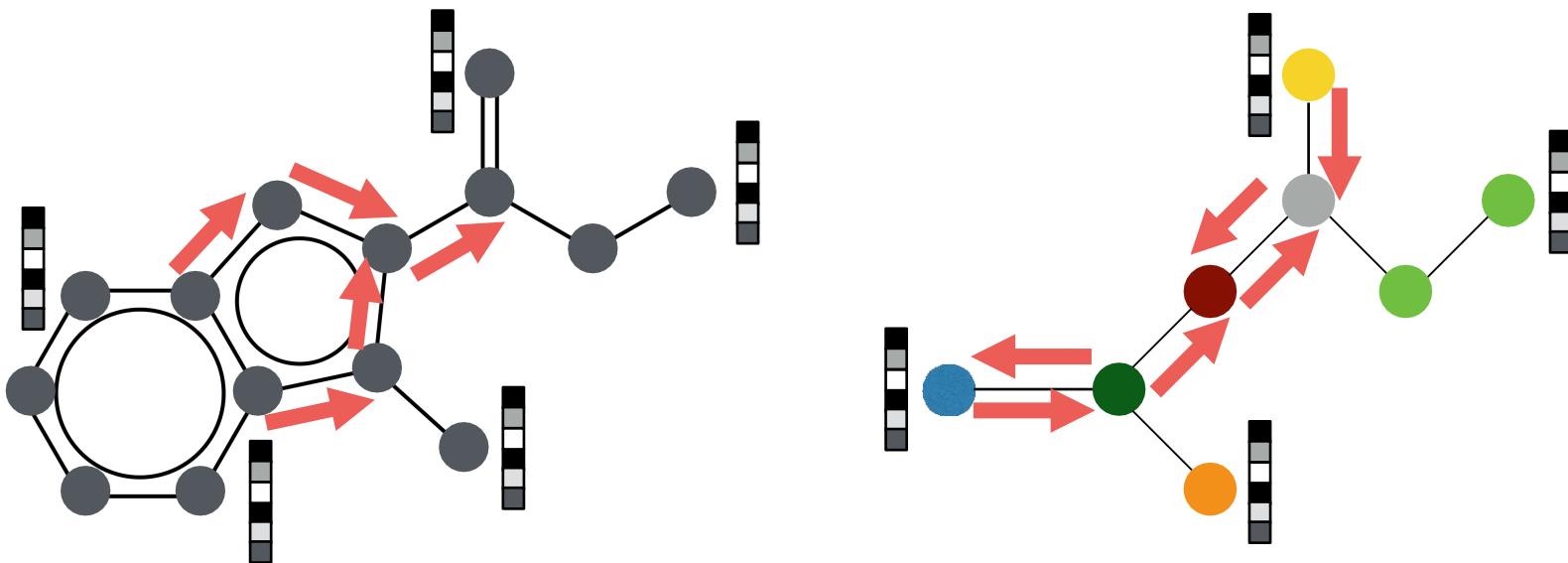


- Generate junction tree → Generate graph group by group
- Vocabulary size: less than 800 given 250K molecules

# Approach: Junction-tree variational autoencoder

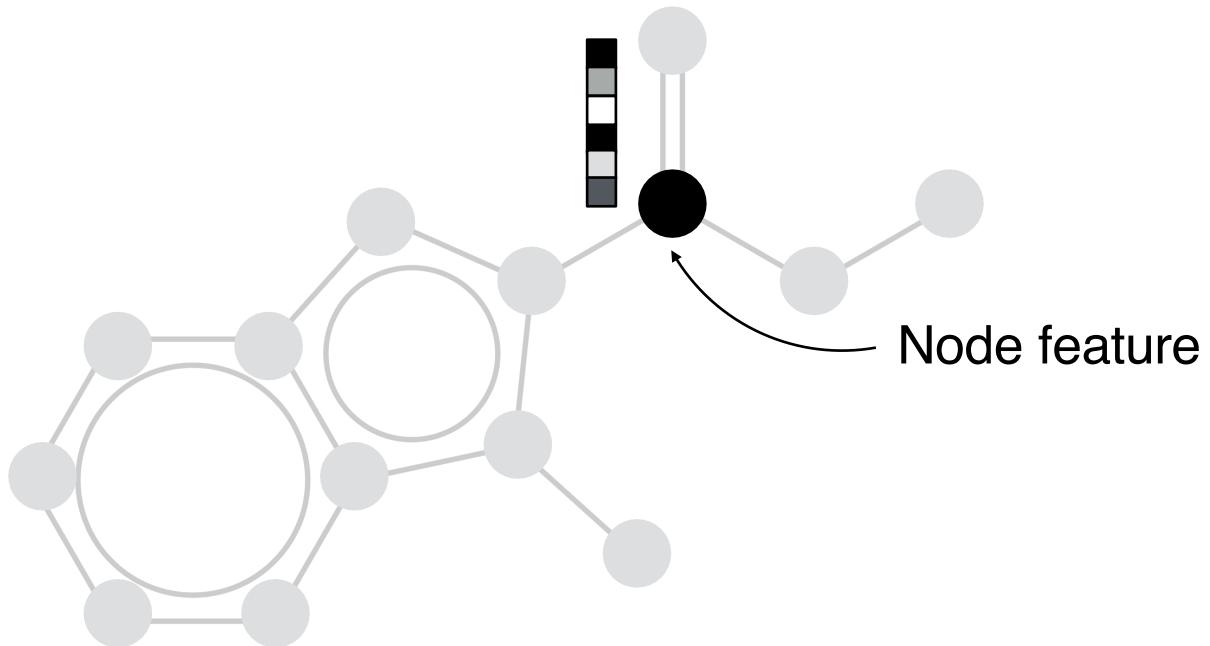


# Graph and tree encoders

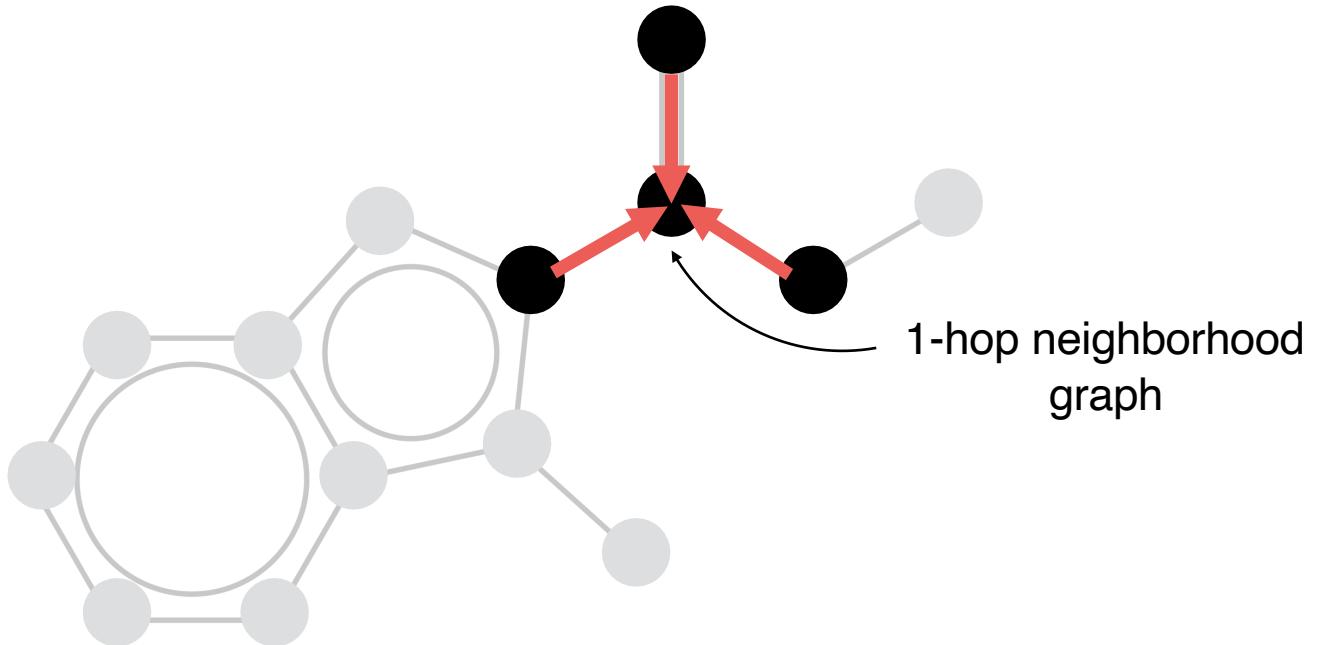


Neural Message Passing Network (MPN)

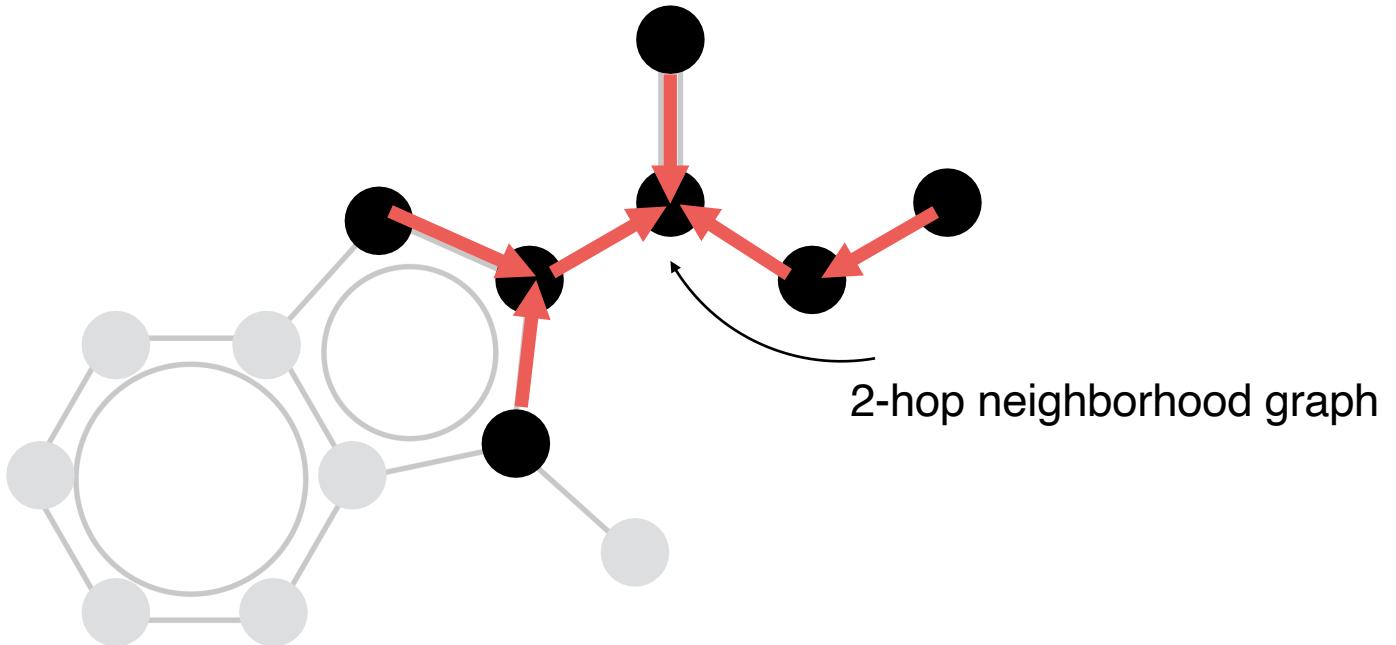
# Graph encoding



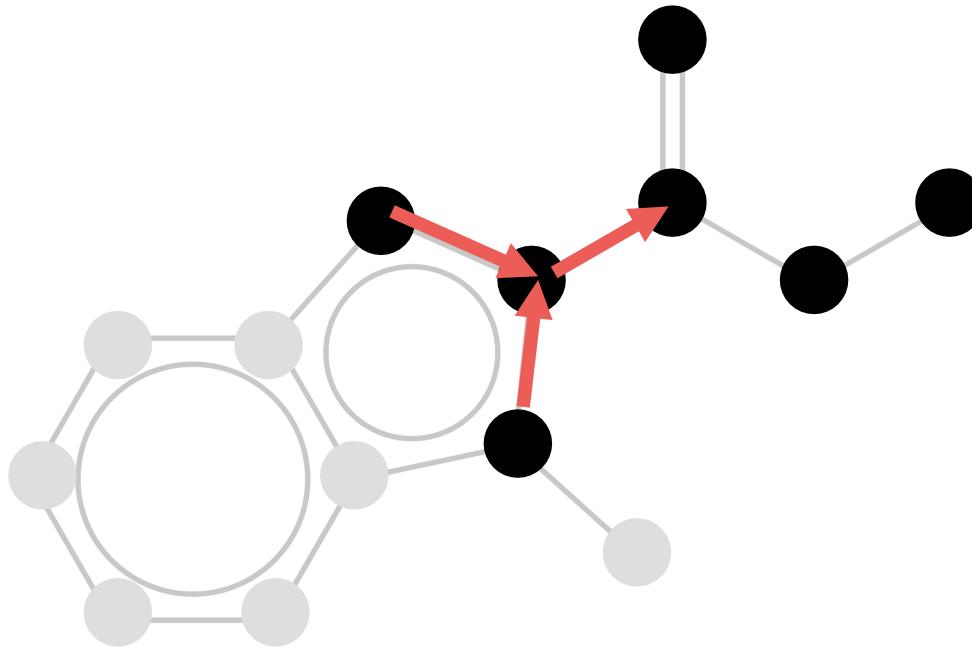
# Graph encoding



# Graph encoding



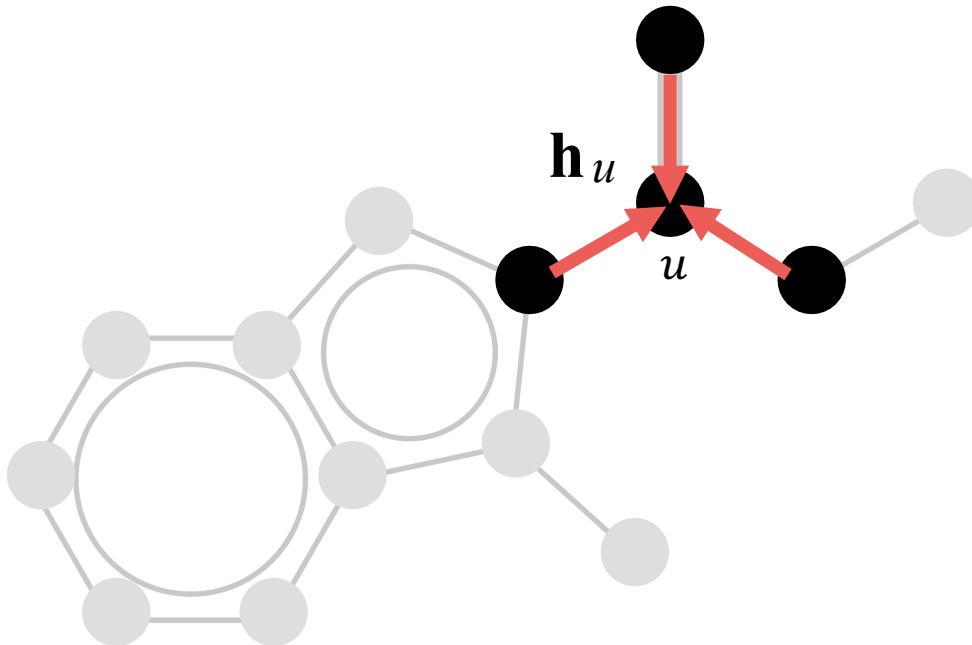
# Graph encoding



$$\nu_{uv}^{(t)} = \tau(\mathbf{W}_1^g \mathbf{x}_u + \mathbf{W}_2^g \mathbf{x}_{uv} + \mathbf{W}_3^g \sum_{w \in N(u) \setminus v} \nu_{wu}^{(t-1)})$$

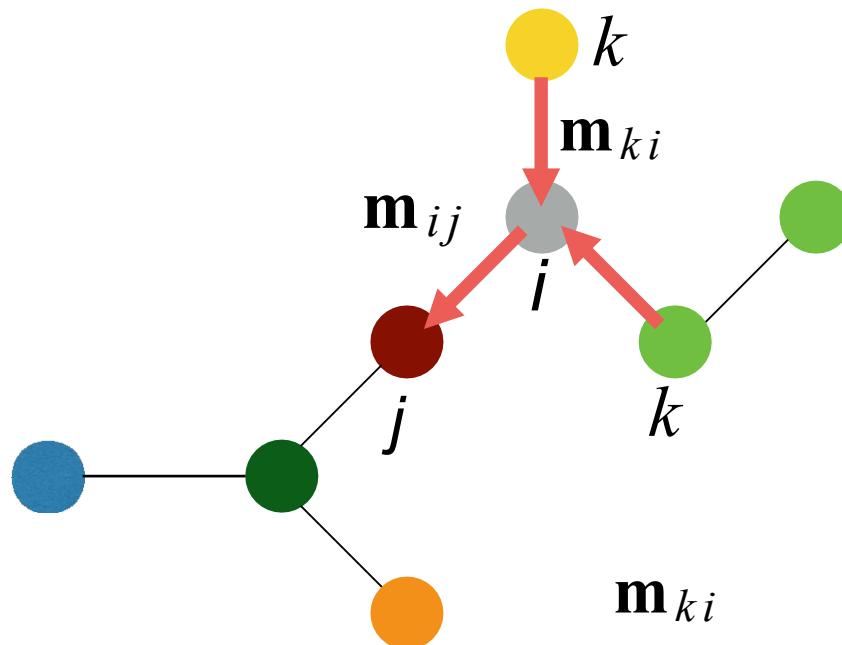
Messages    Node feature    Edge feature     $w \in N(u) \setminus v$

# Graph encoding



$$\mathbf{h}_u = \tau(\mathbf{U}_1^g \mathbf{x}_u + \sum_{v \in N(u)} \mathbf{U}_2^g \boldsymbol{\nu}_{vu}^{(T)})$$

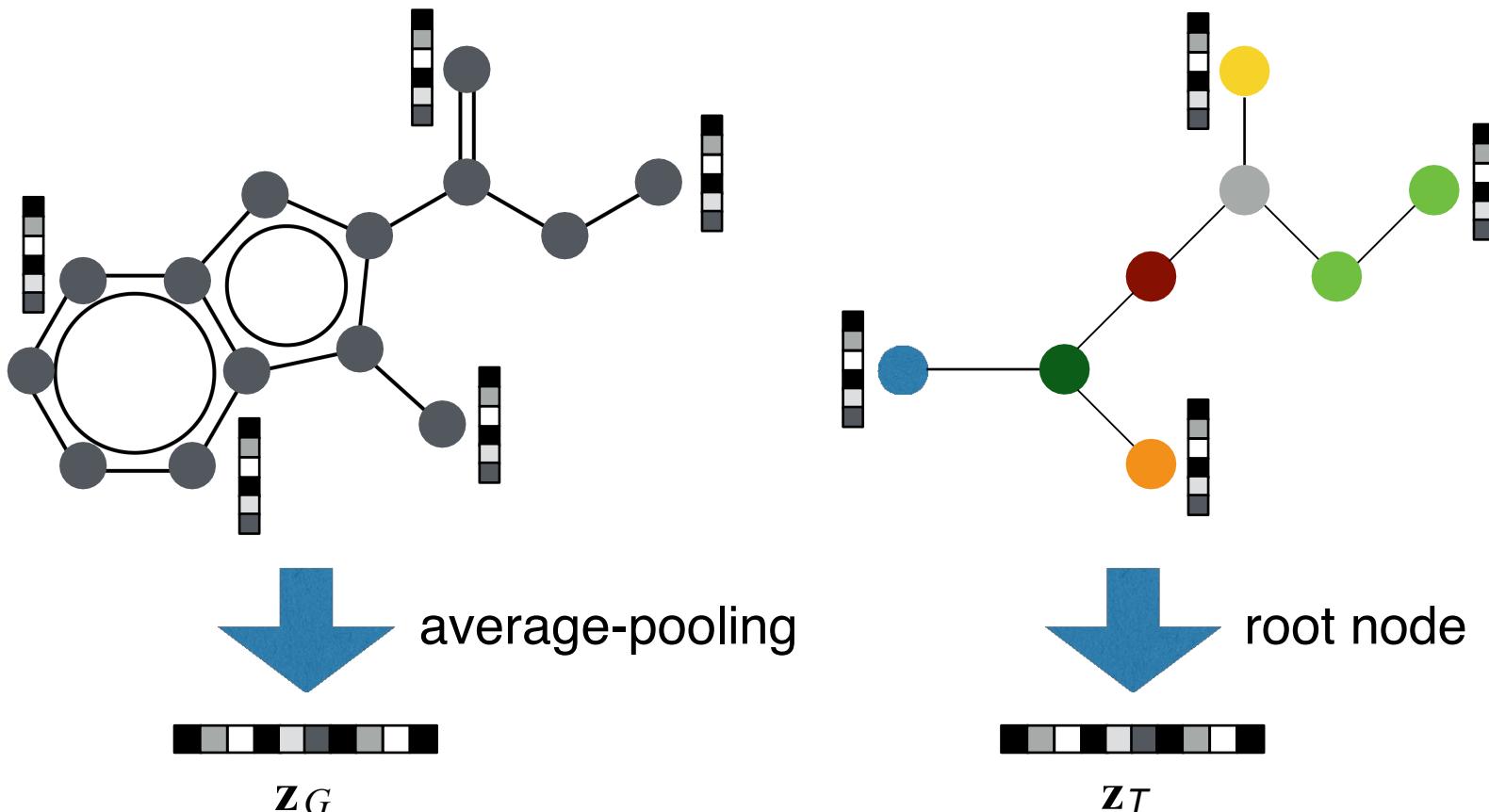
# Tree encoding



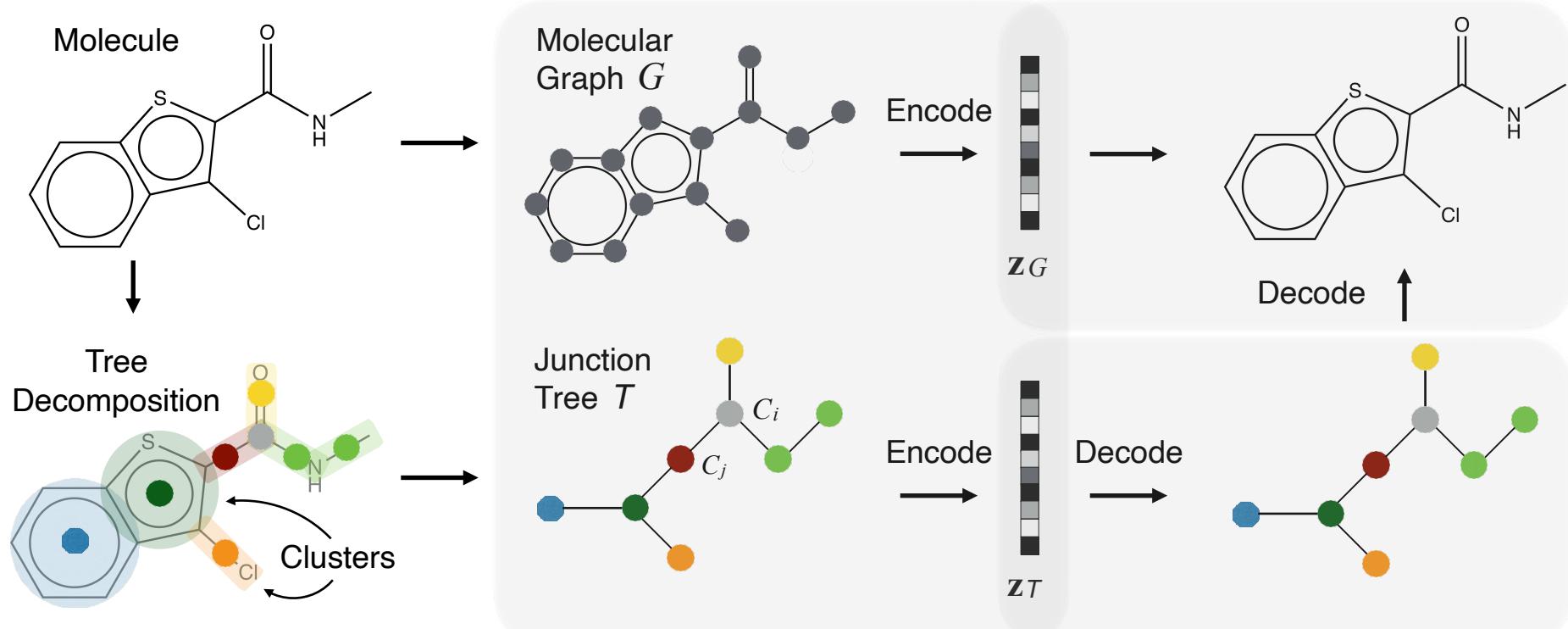
$$\mathbf{m}_{ij} = \text{GRU}(\mathbf{x}_i, \{\mathbf{m}_{ki}\}_{k \in N(i) \setminus j})$$

To capture long range interactions

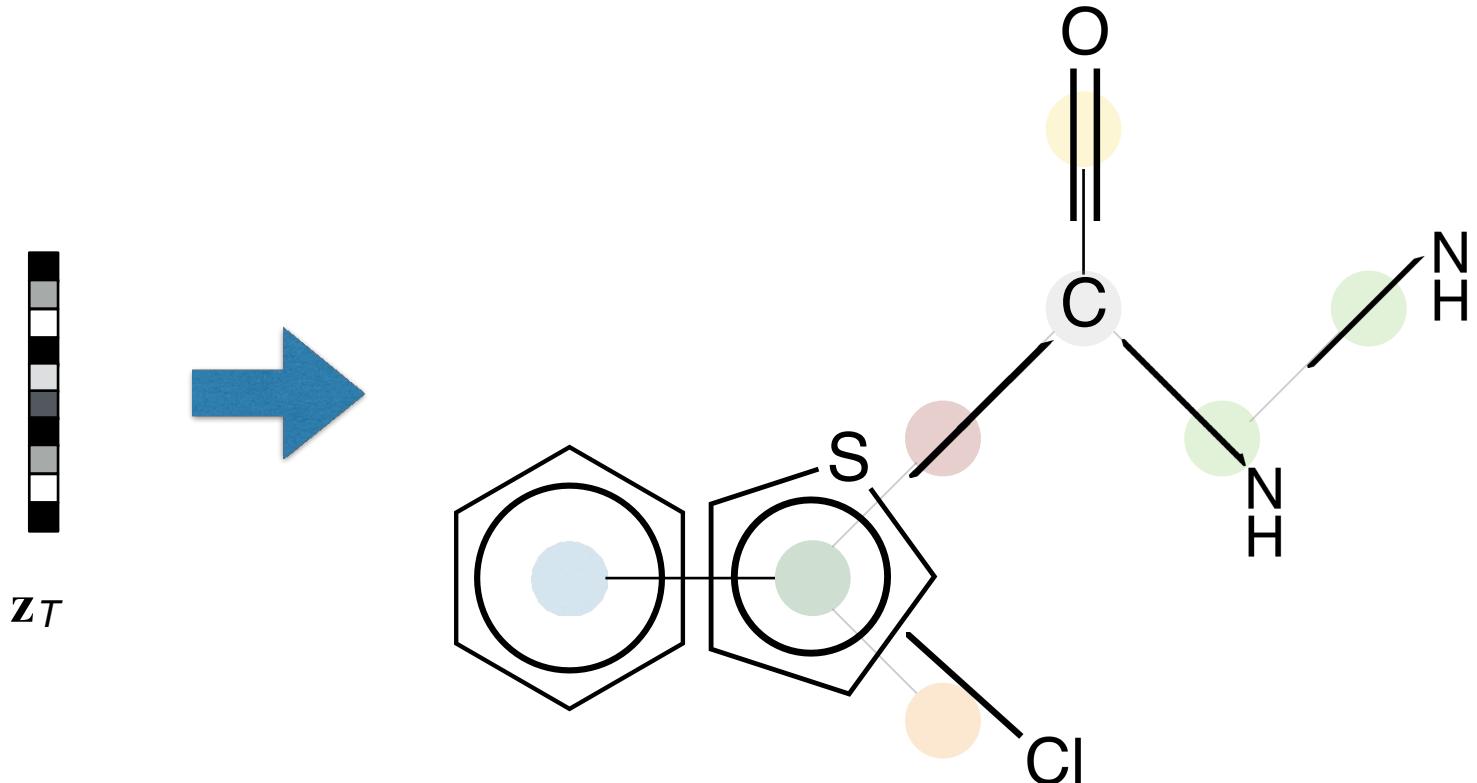
# Graph and tree encoders



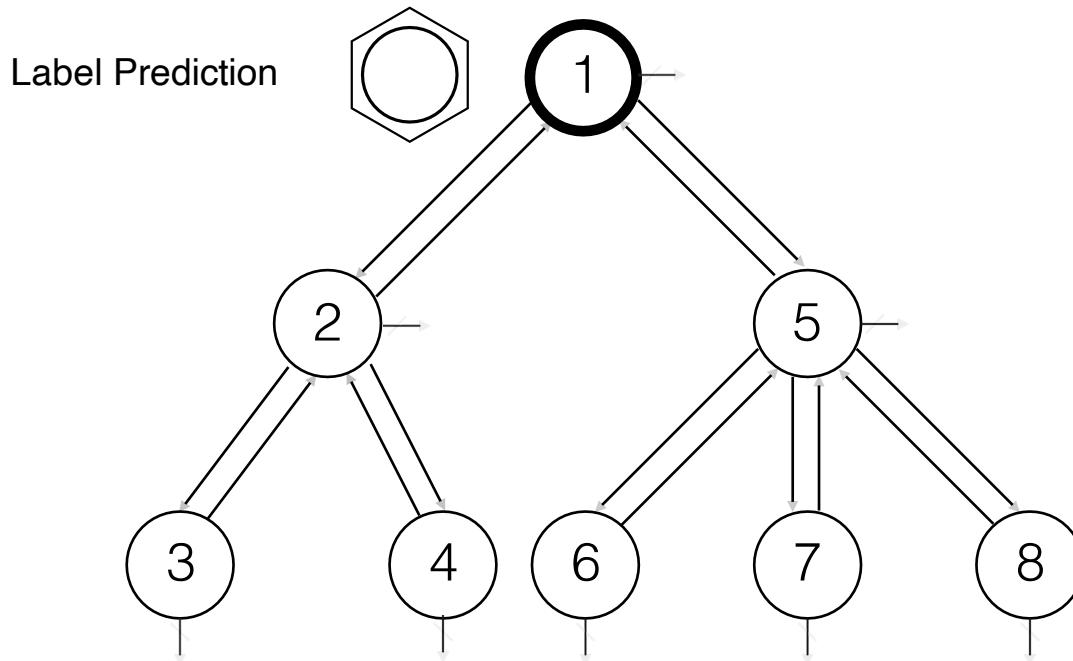
# Approach: Junction-tree variational autoencoder



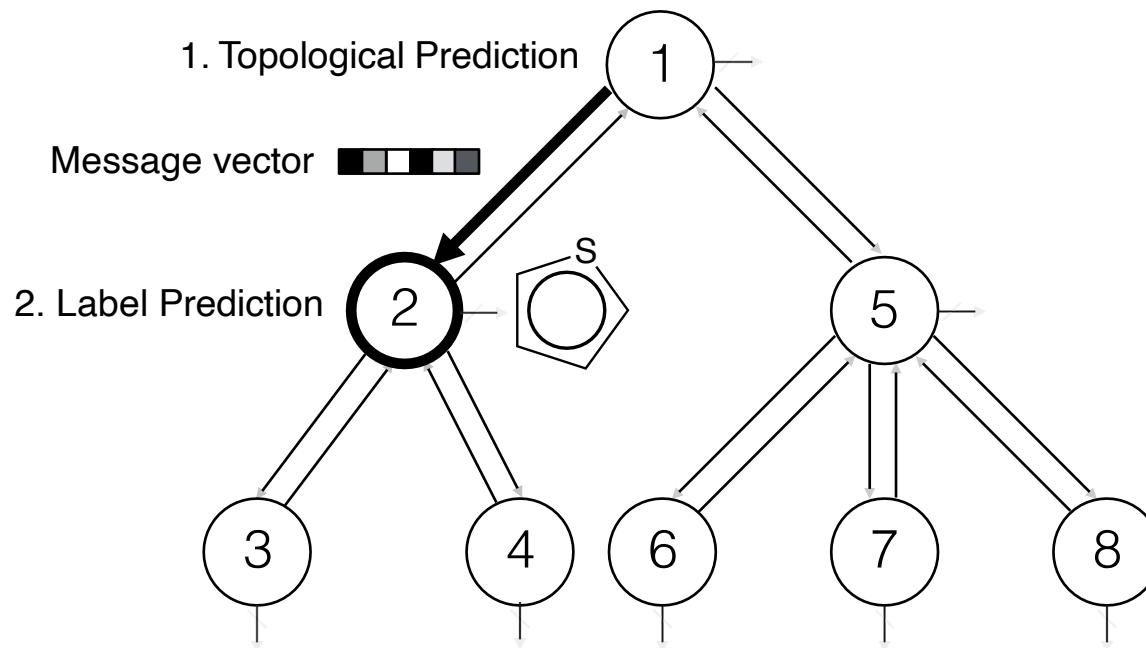
# Tree decoder



# Tree decoder



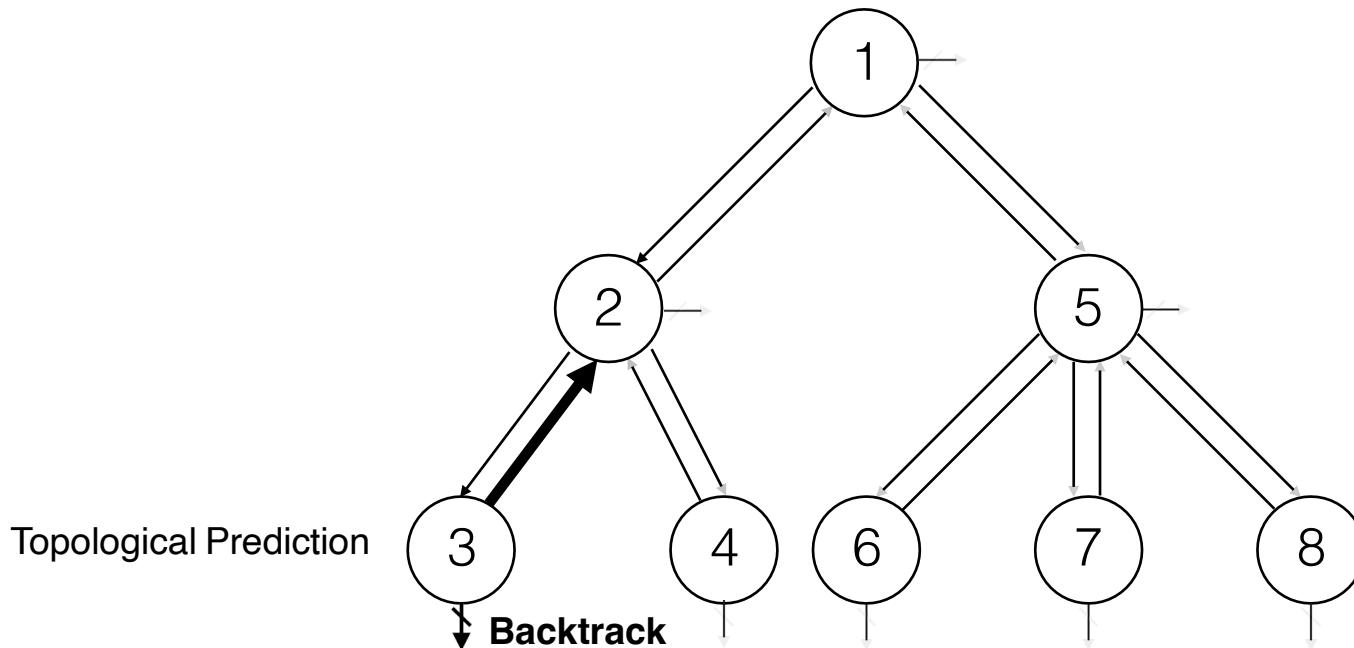
# Tree decoder



**Topological Prediction:** Whether to expand a child or backtrack?

**Label Prediction:** What is the label of a node?

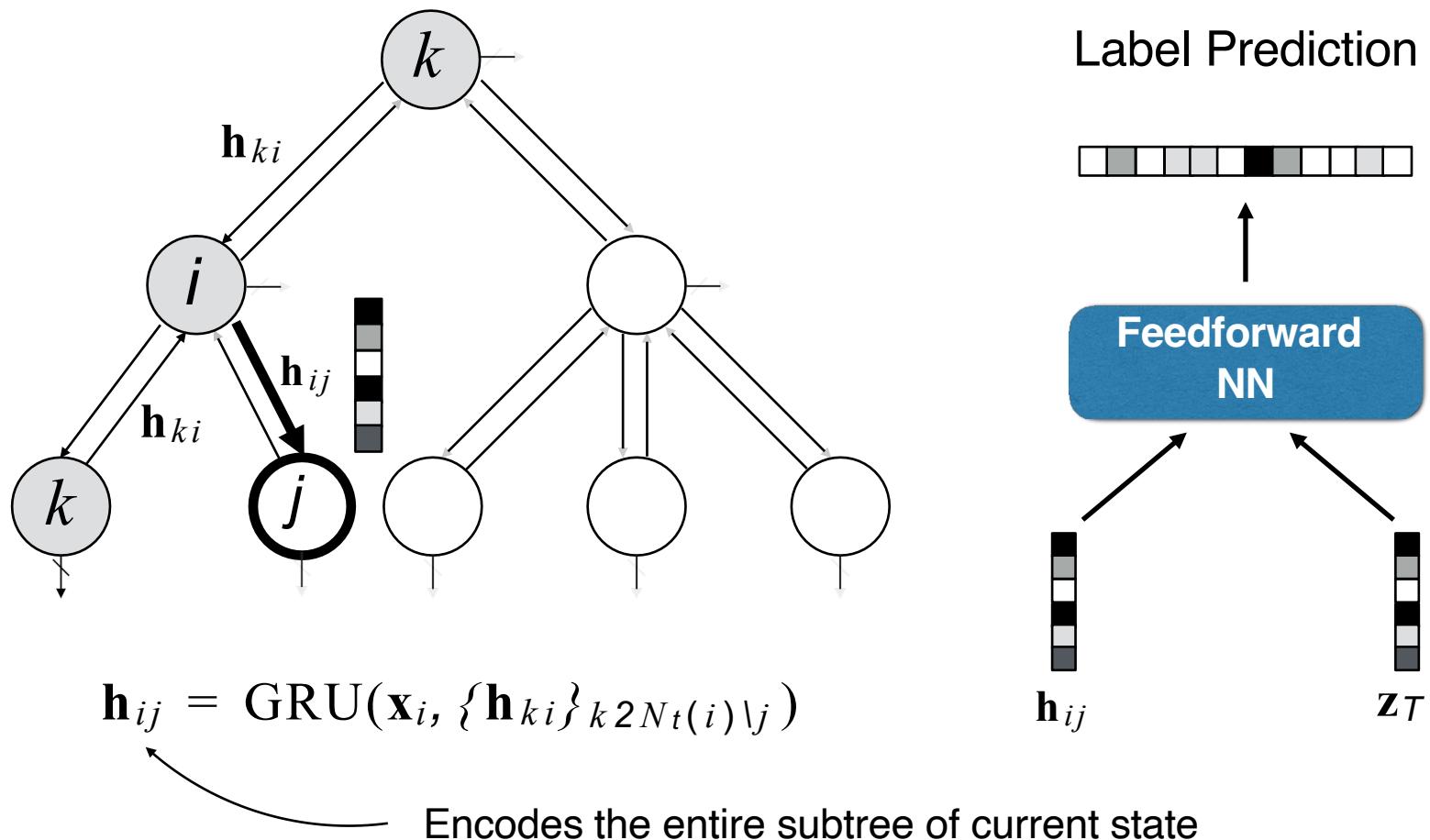
# Tree decoder



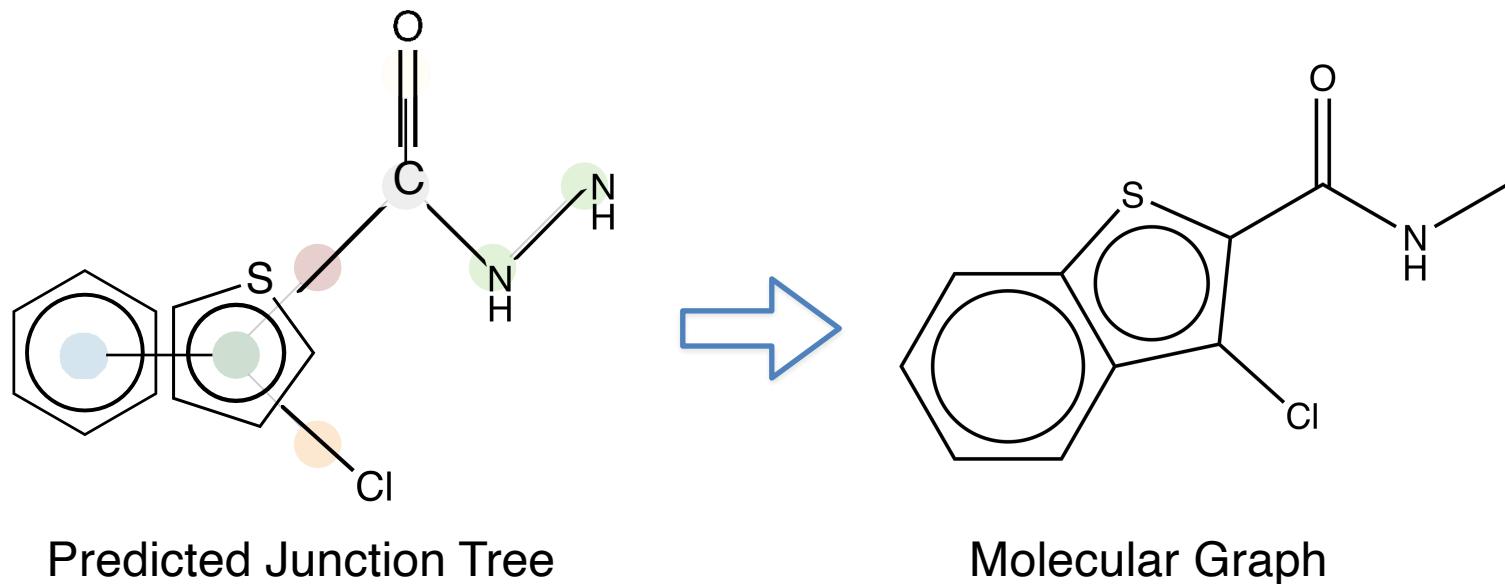
**Topological Prediction:** Whether to expand a node or backtrack?

**Label Prediction:** What is the label of a node?

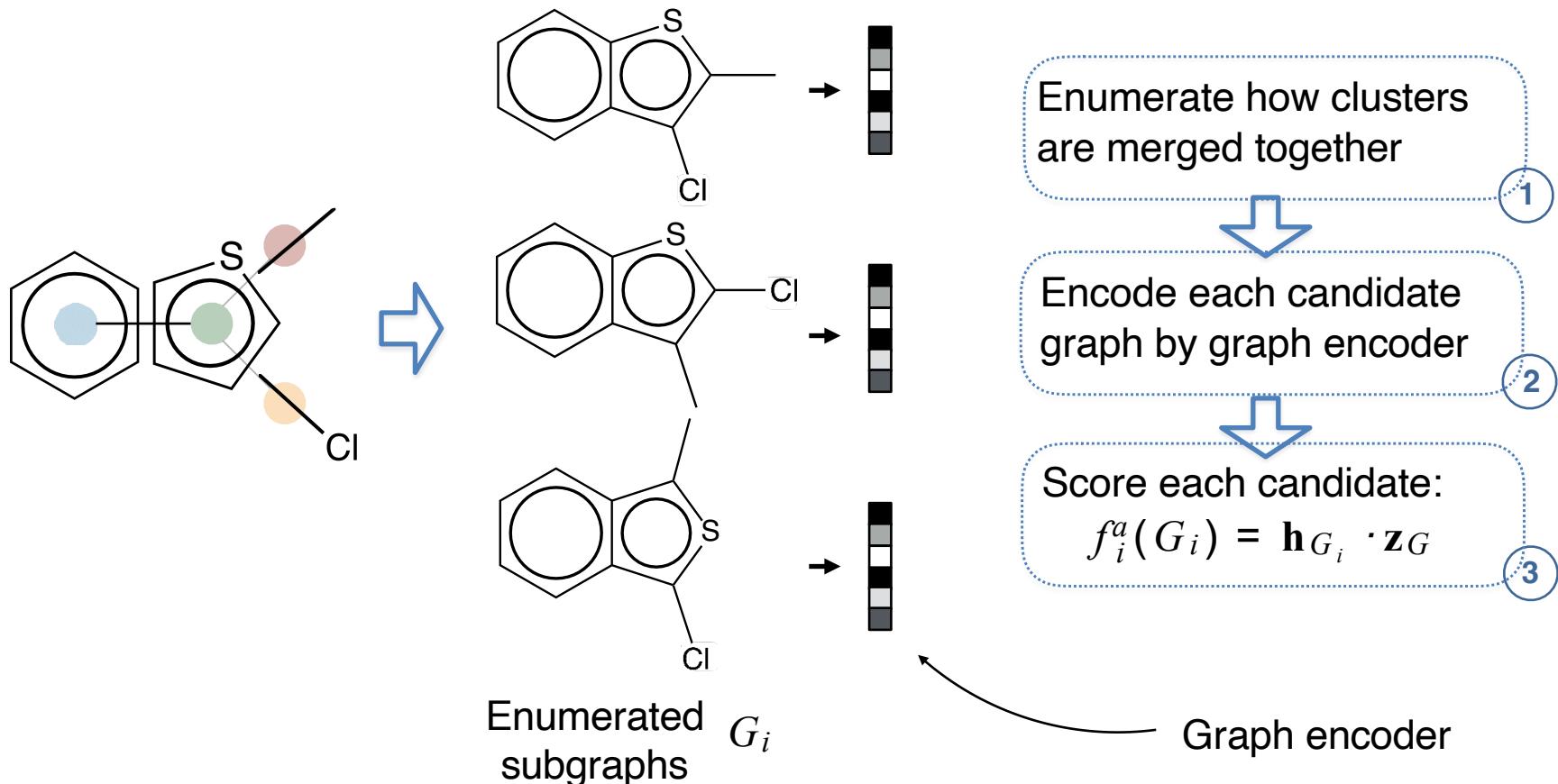
# Tree decoder



# Graph decoder

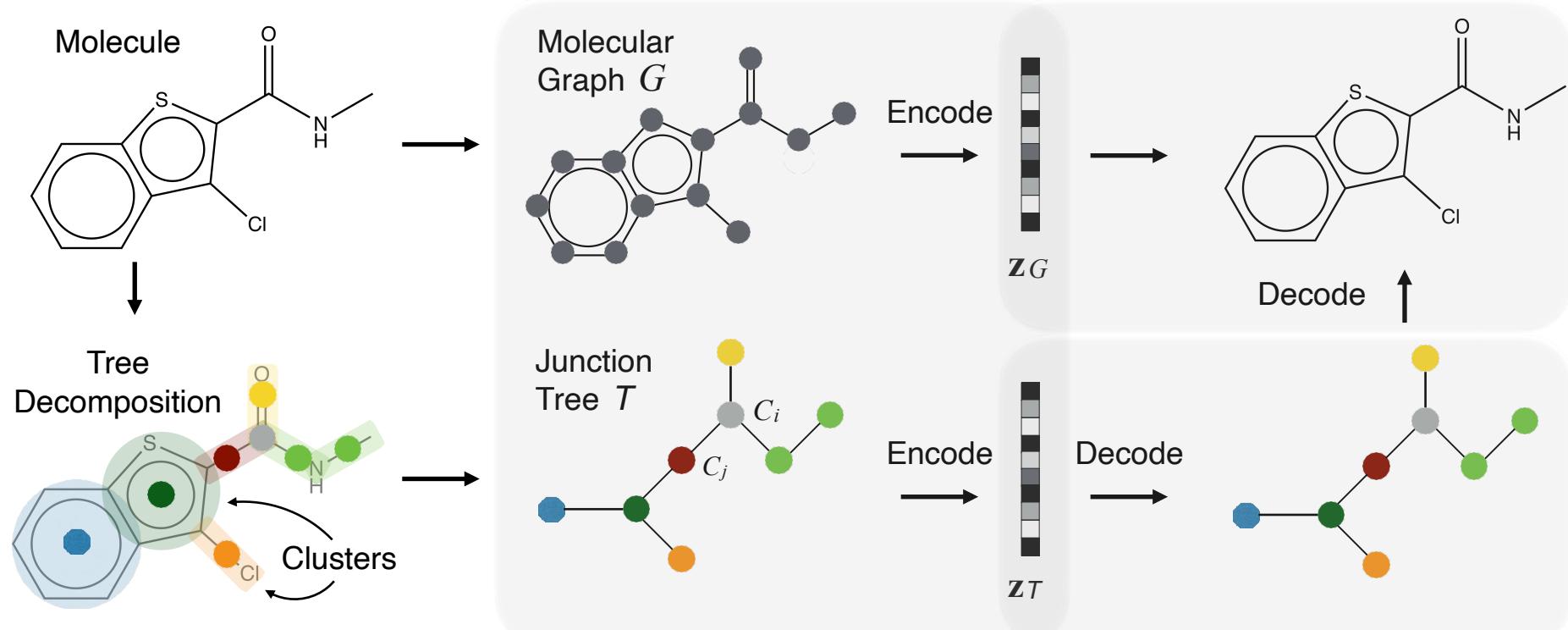


# Graph decoder



$$\mathcal{L}_g(G) = \sum_i \left[ f^a(G_i) - \log \sum_{G'_i \in \mathcal{G}_i} \exp(f^a(G'_i)) \right] \quad (16)$$

# Recap: Junction-tree variational autoencoder



# Experiments

- **Data:** 250K compounds from ZINC dataset
- **Molecule Generation:** How many molecules are valid when sampled from Gaussian prior?
- **Molecule Optimization**
  - **Global:** Find the best molecule in the entire latent space.
  - **Local:** Modify a molecule to increase its potency

# Baselines

## **SMILES string based:**

1. Grammar VAE (GVAE) (Kusner et al., 2017);
2. Syntax-directed VAE (SD-VAE) (Dai et al., 2018)

## **Graph based:**

1. Graph VAE (Simonovsky & Komodakis, 2018)
2. DeepGMG (Li et al., 2018)

---

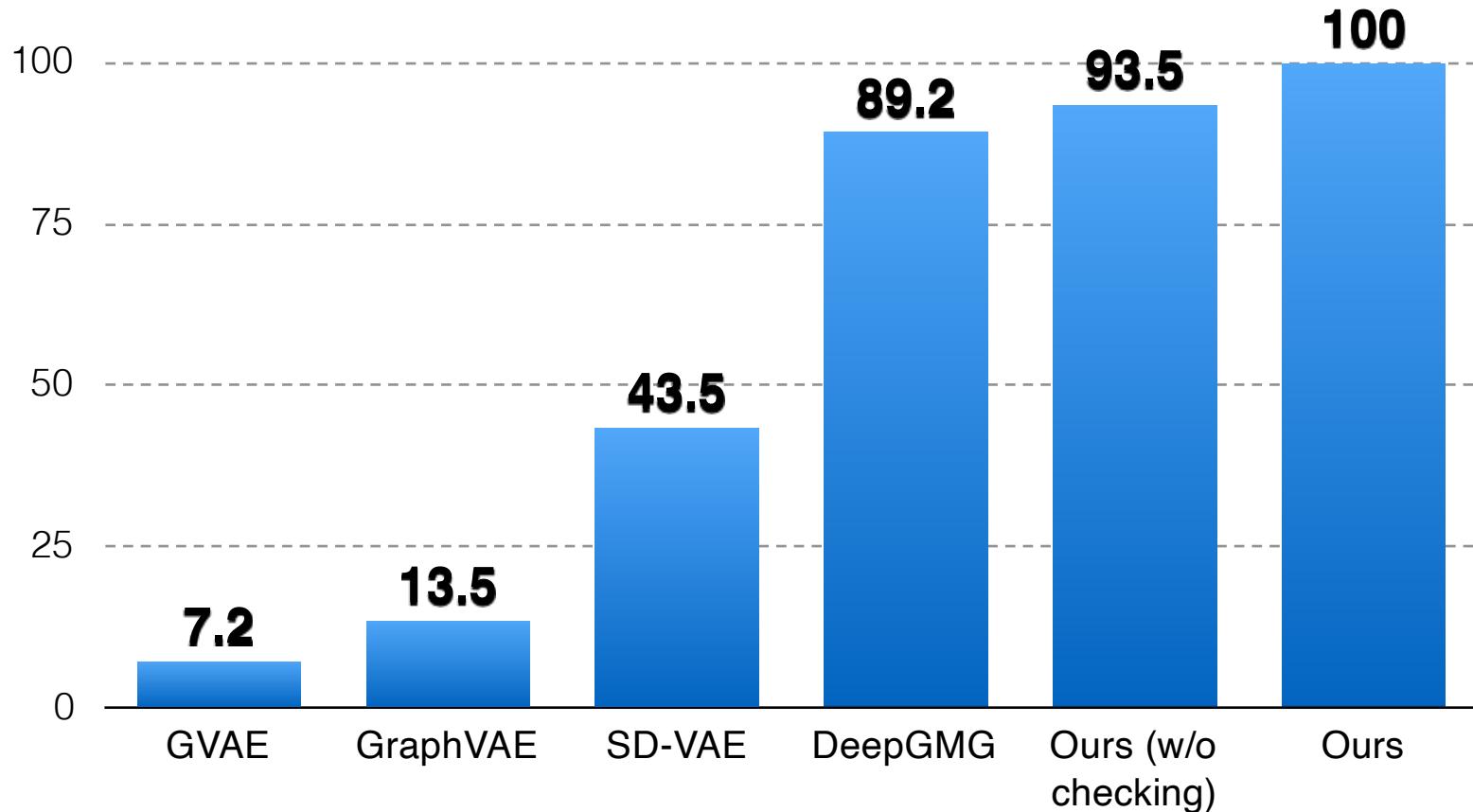
[2] Li et al., Learning Deep Generative Models of Graphs, 2018

5 Kusner et al., Grammar Variational Autoencoder, 2017

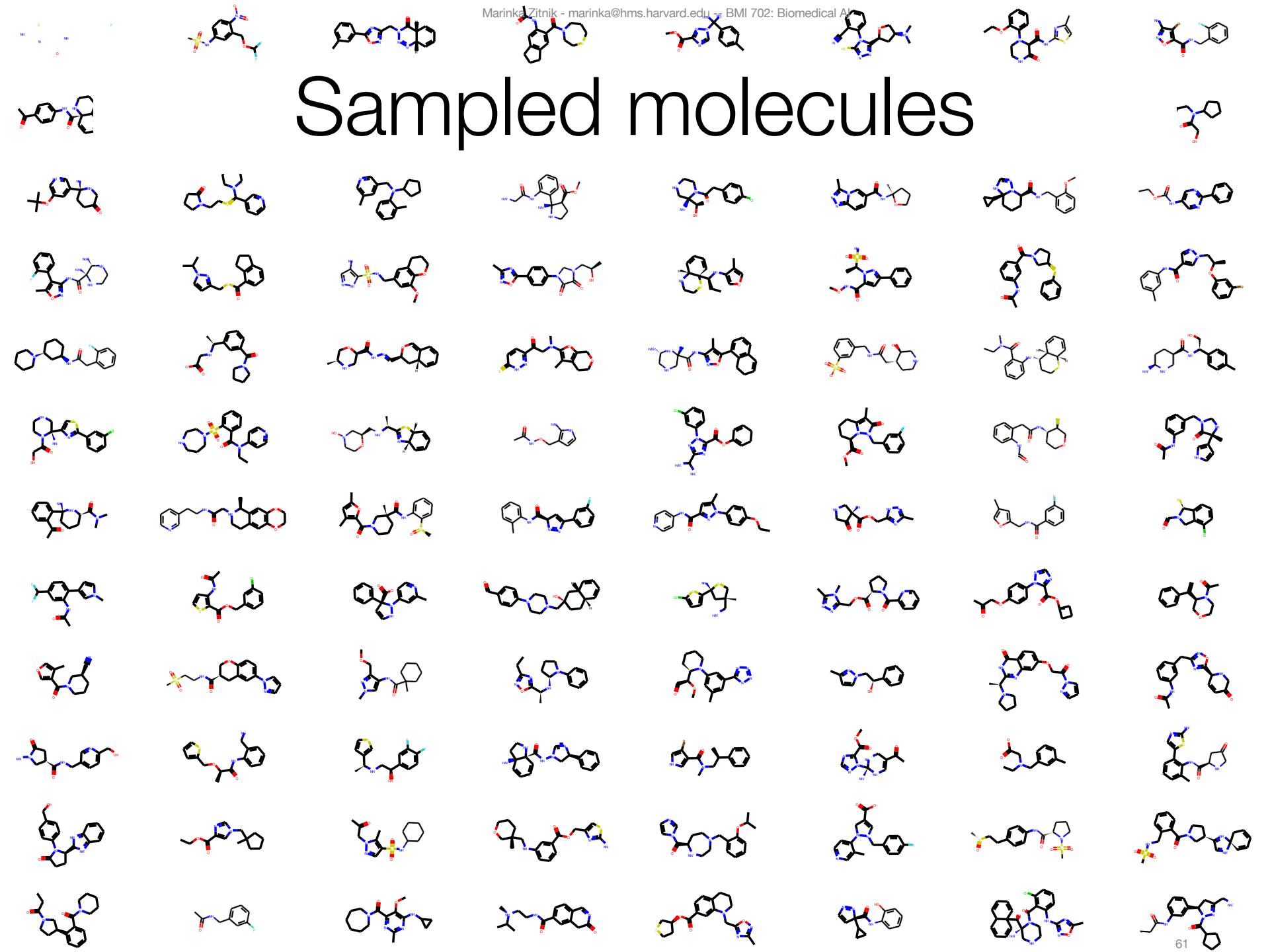
6 Dai et al., Syntax-directed Variational Autoencoder for structured data, 2018

7 Simonovsky & Komodakis, GraphVAE: Towards generation of small graphs using variational autoencoders

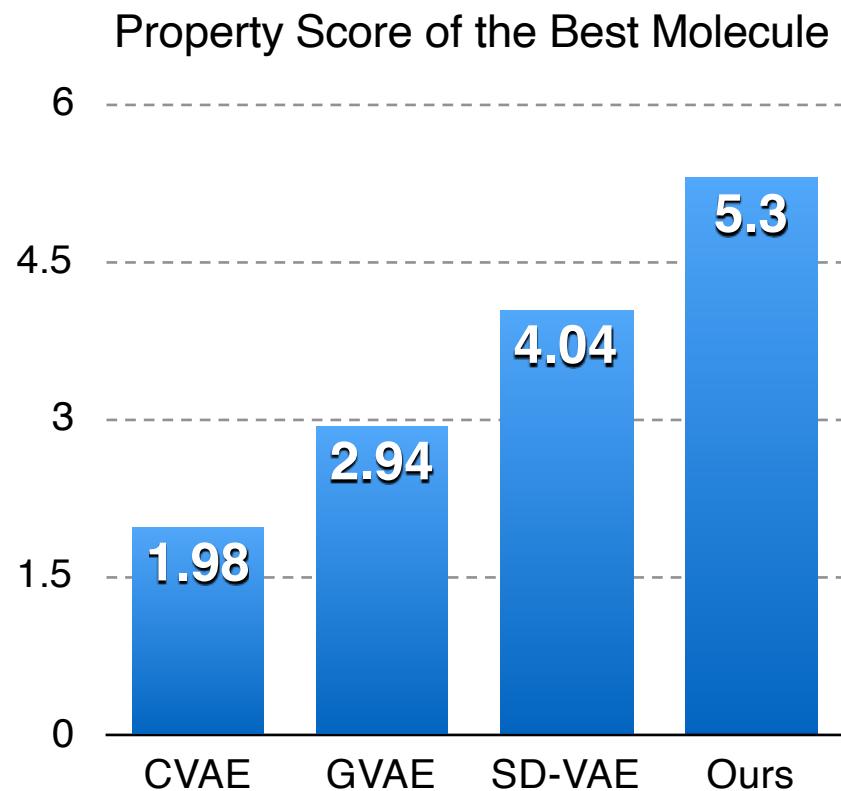
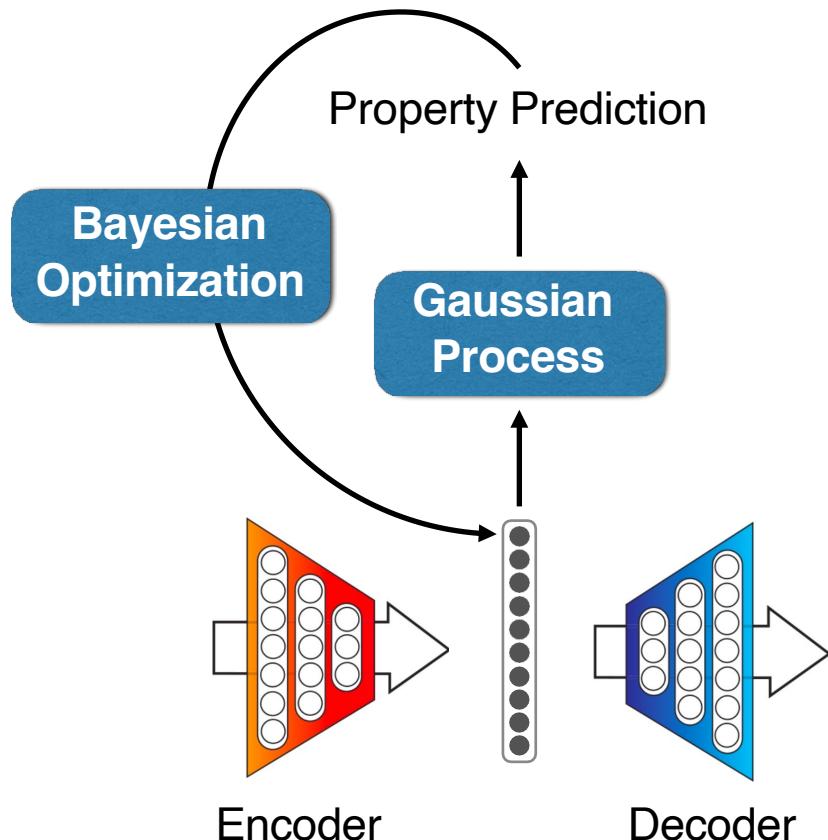
# Molecule generation (Validity)



# Sampled molecules

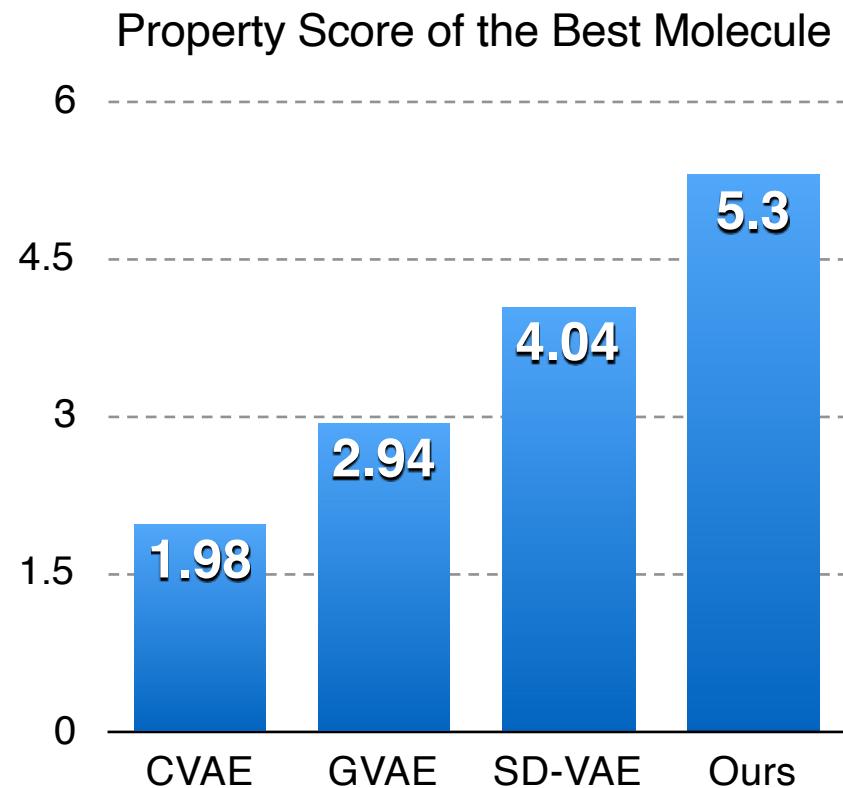
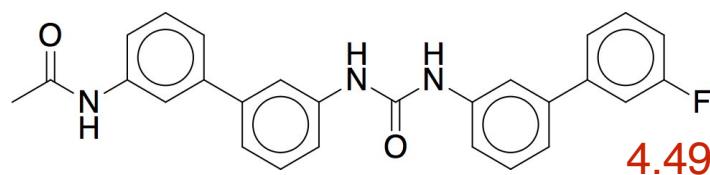
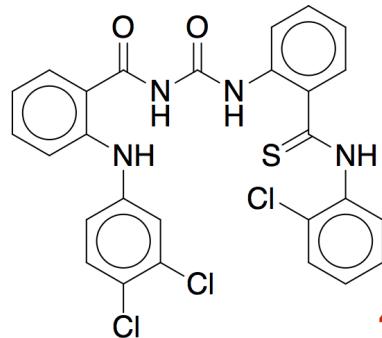
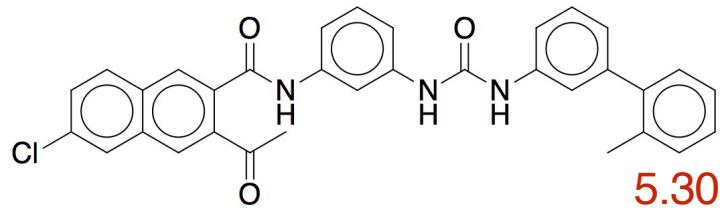


# Molecule optimization (Global)



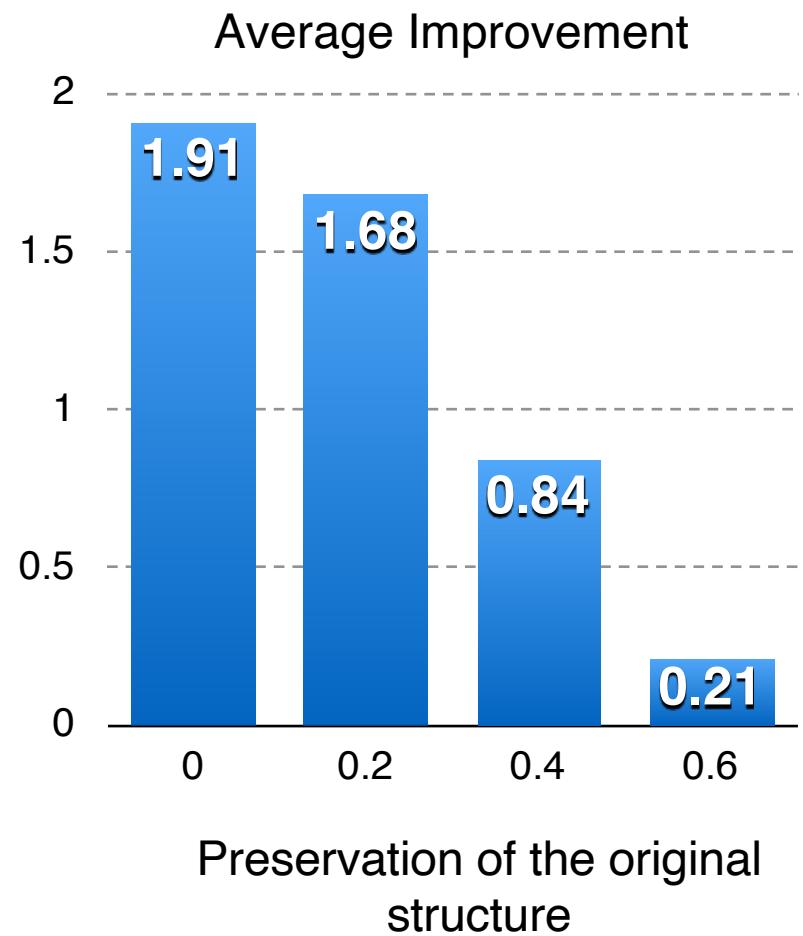
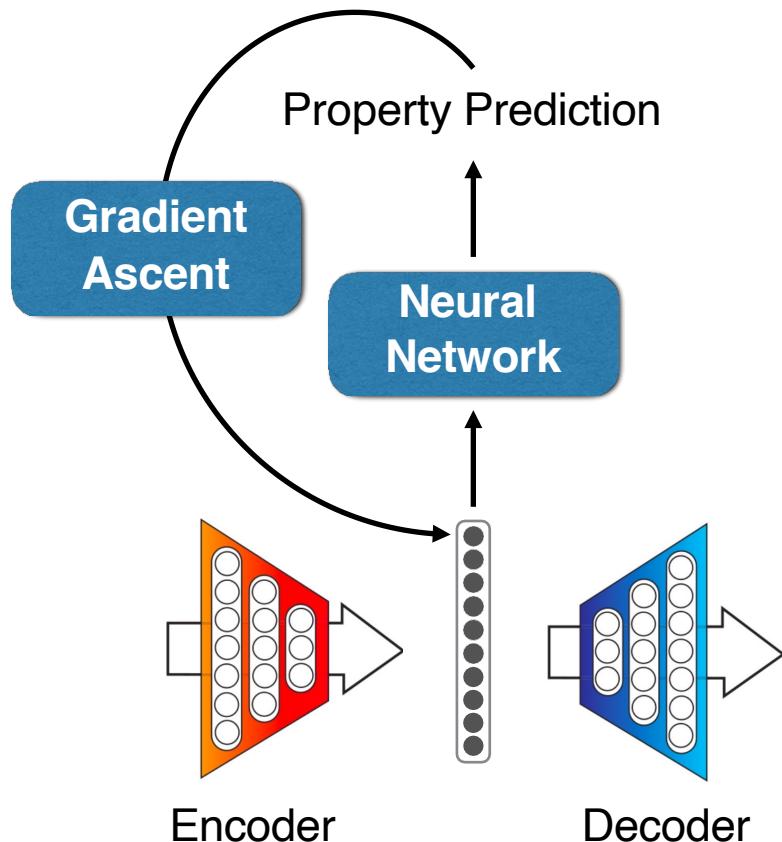
Property: Solubility + Ease of Synthesis

# Molecule optimization (Global)

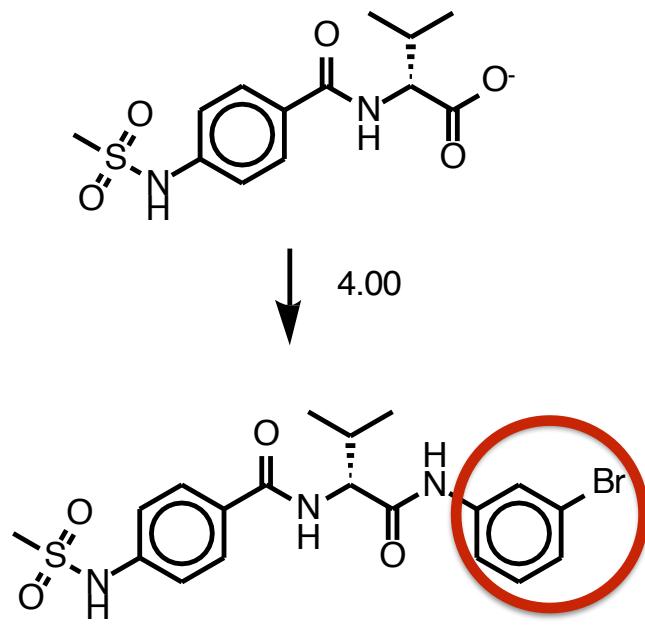


Property: Solubility + Ease of Synthesis

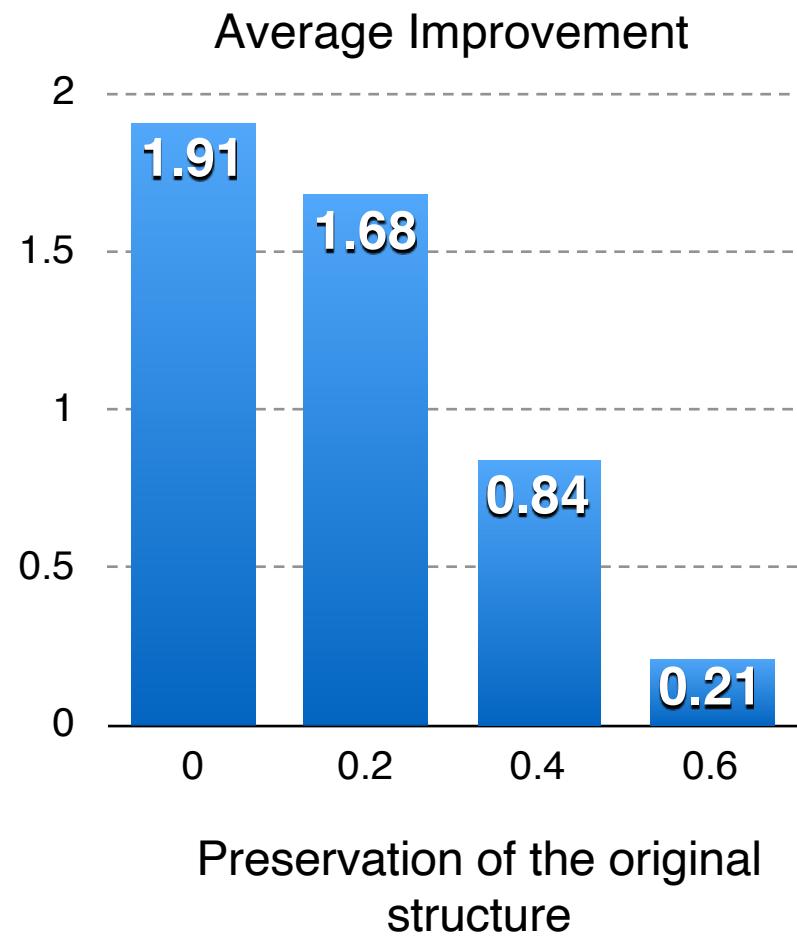
# Molecule optimization (Local)



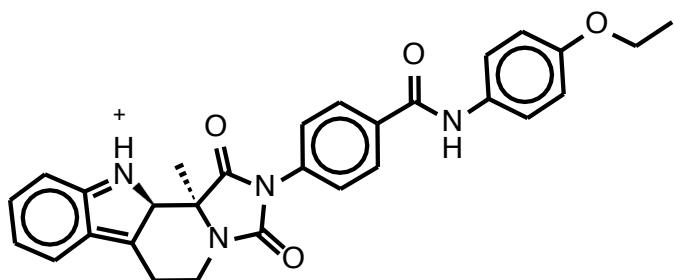
# Molecule optimization (Local)



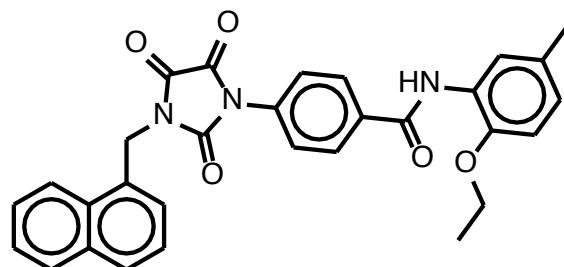
Preservation  $\approx 0.6$



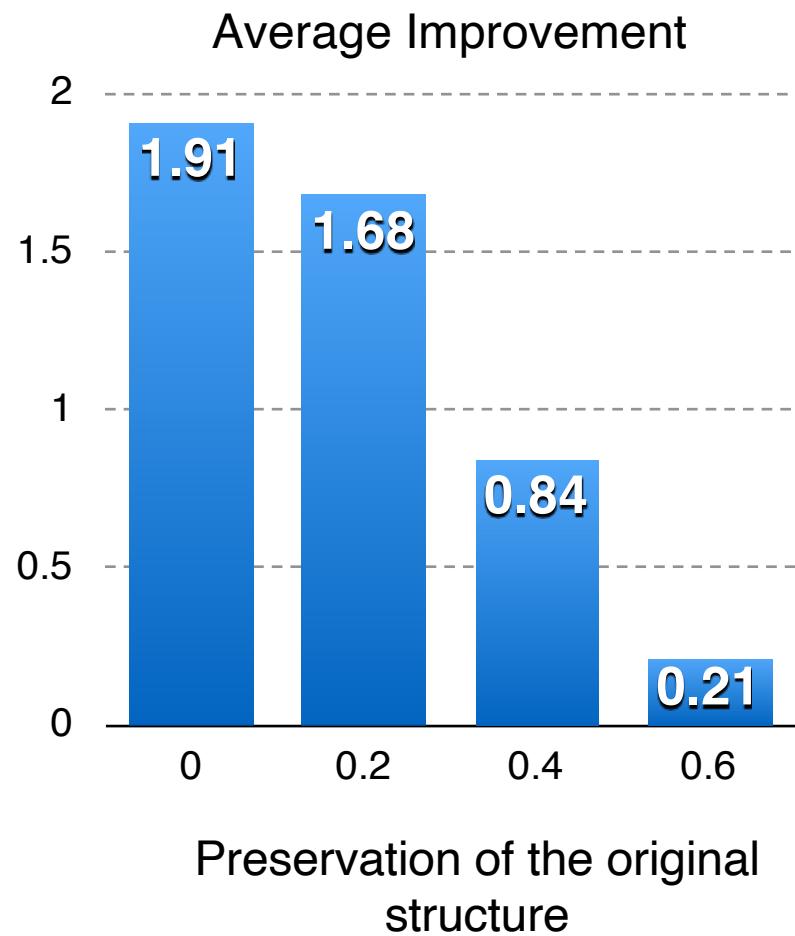
# Molecule optimization (Local)



↓  
5.69



Preservation ≈ 0.4



# Quick Check

<https://forms.gle/4DJYidWsL4KkDau57>

BMI 702: Biomedical Artificial Intelligence  
*Foundations of Biomedical Informatics II, Spring 2023*

Quick check quiz for lecture  
12: Overview of drug discovery and development, small-molecule generation, molecule optimization, identification and characterization of therapeutic targets, high-throughput genetic and chemical perturbations.

Course website and slides: <https://zitniklab.hms.harvard.edu/BMI702>

\* Indicates required question

First and last name \*

Your answer

Harvard email address \*

Your answer

Describe two challenges that models for generating molecular graphs need to address. \*

Your answer

What is the difference between traditional vs. neural fingerprint representations? \*

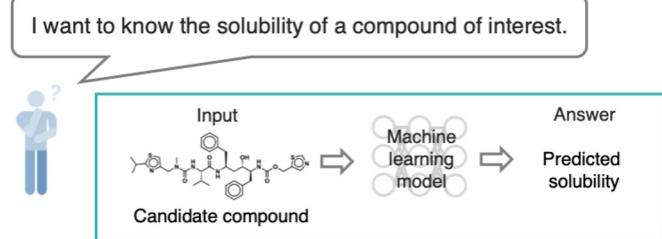
Your answer

**Submit** **Clear form**

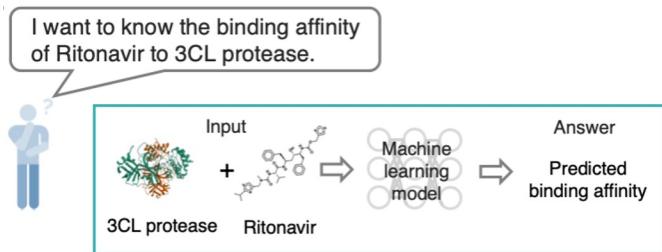
Never submit passwords through Google Forms.

# Outline for today's class

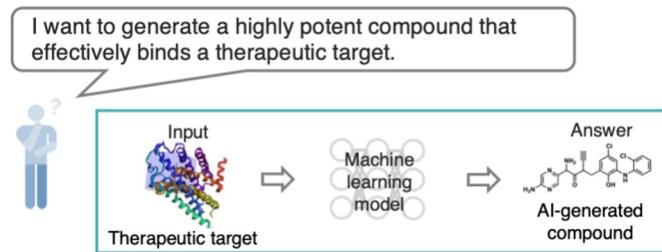
- Optimization & generation of small molecules



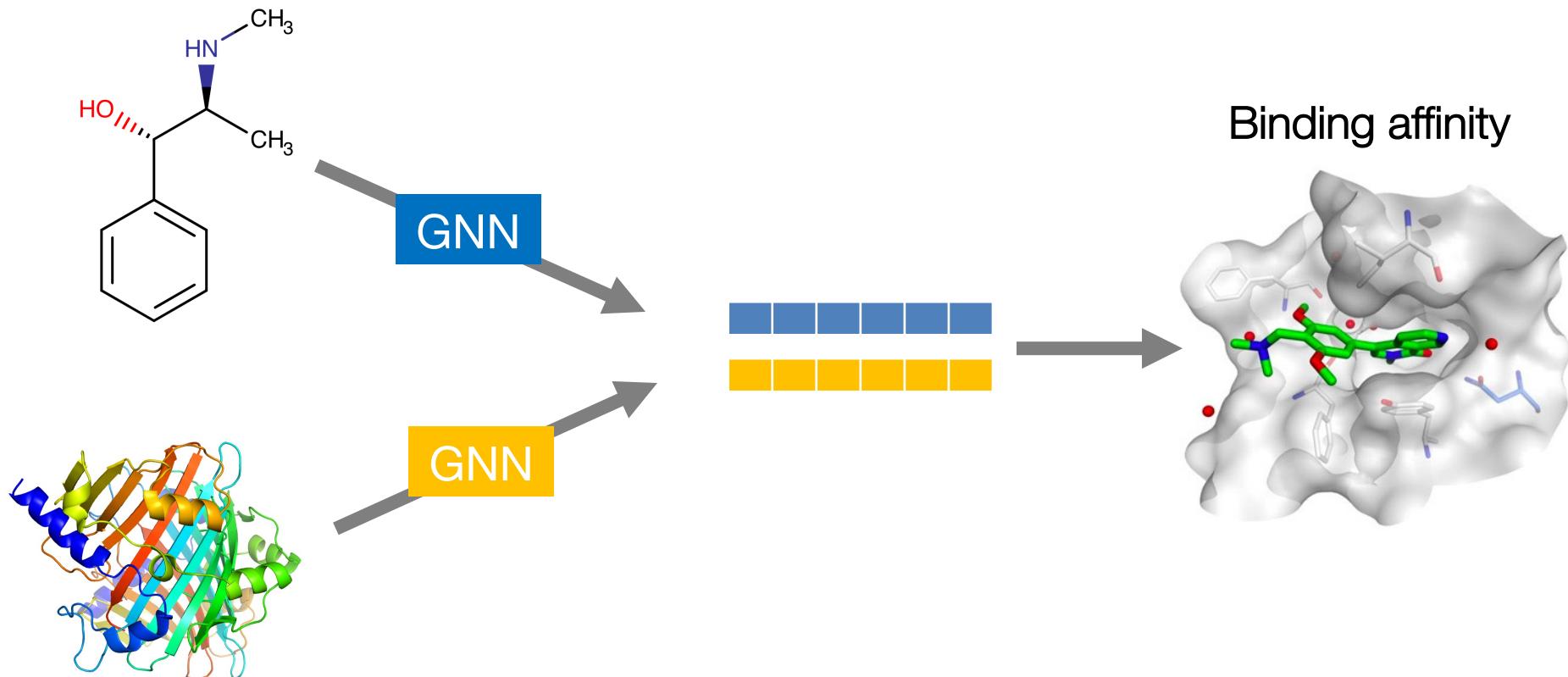
- Binding of drugs to therapeutic targets



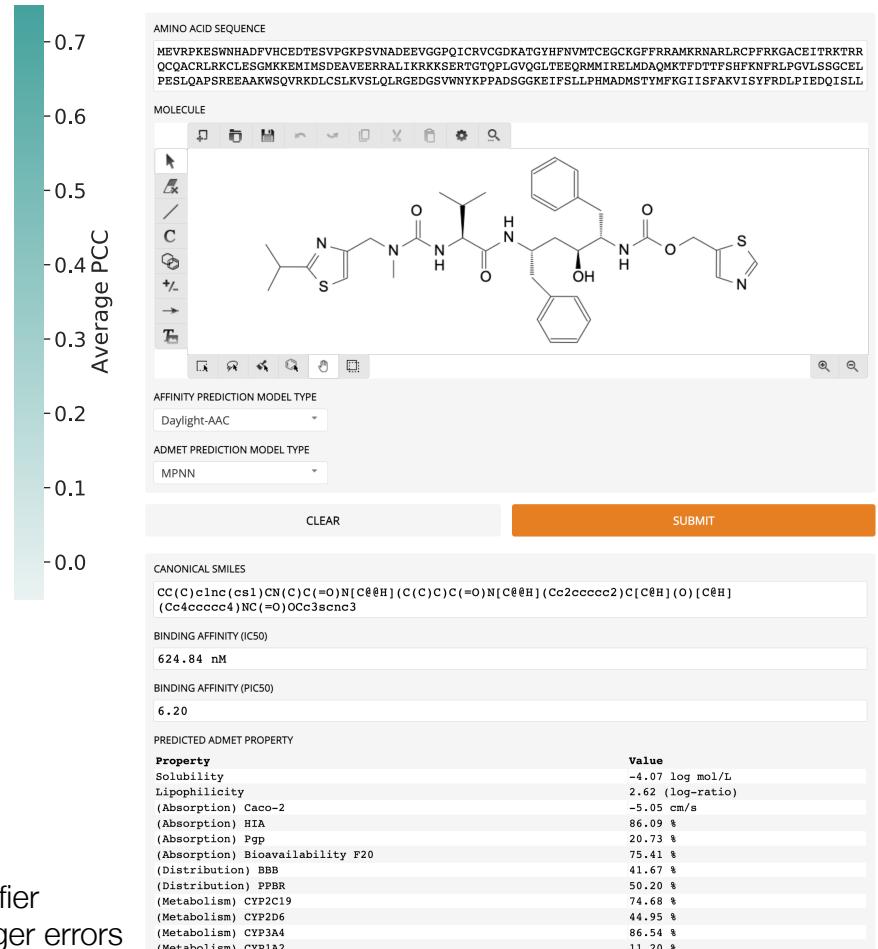
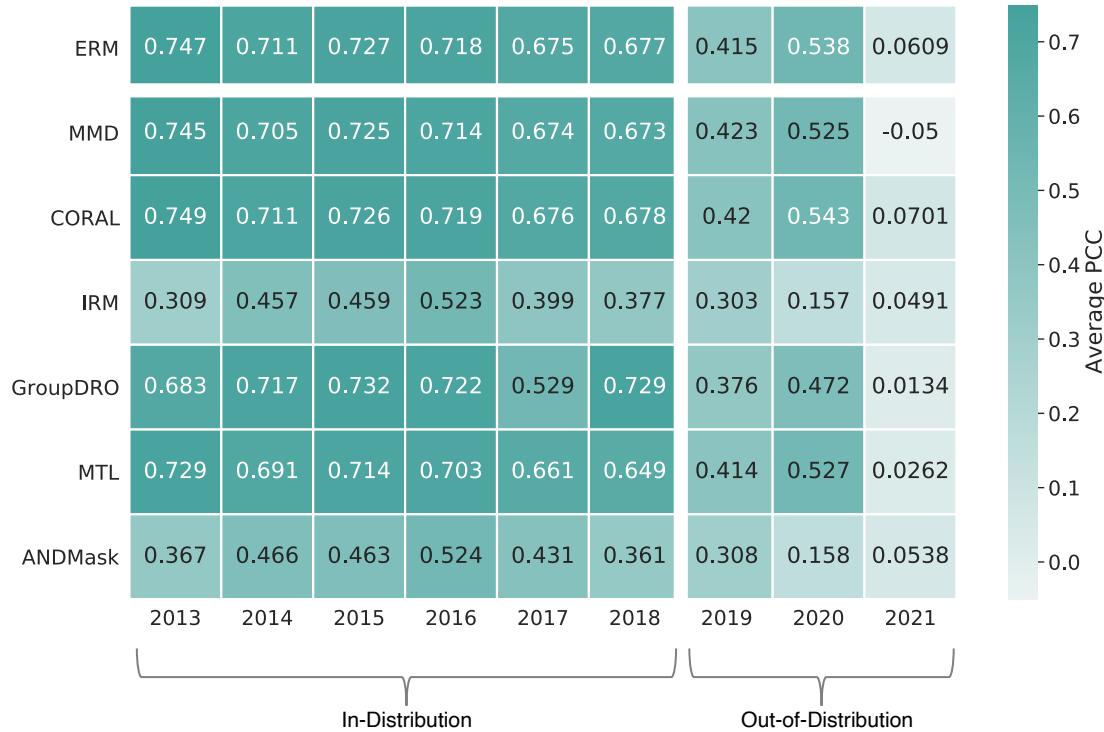
- High-throughput genetic & chemical perturbations



# Geometric modeling of binding



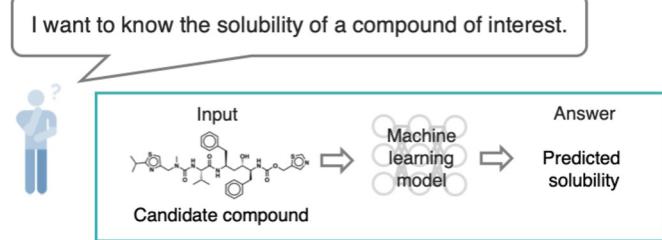
# Results: Binding affinity prediction



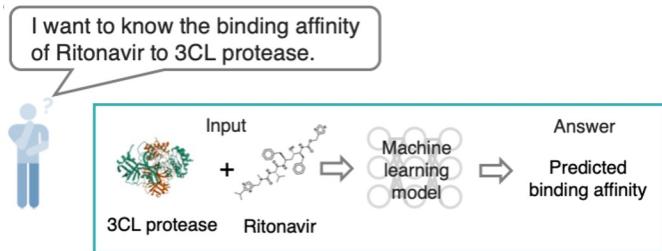
Modern data management  
Human-AI collaboration

# Outline for today's class

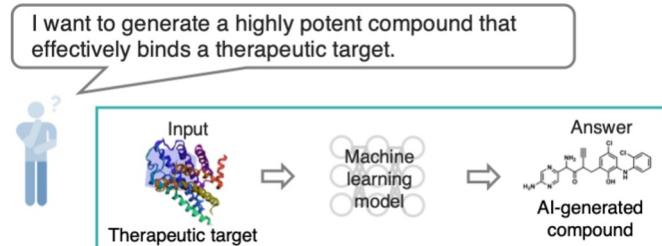
- Optimization & generation of small molecules



- Binding of drugs to therapeutic targets



- High-throughput genetic & chemical perturbations





Words and genes share a correspondence:  
their **meanings** arise from their **context**.

Gene perturbation measurements across diverse cell contexts  
induce **semantics for genes**

(under the right approach)

“apple” is a **polysemic** word...



grow an apple

buy an apple|

... whose **particular meaning** is resolved via **sentence context**.



grow an apple

grow an apple tree

grow an apple tree from seed

grow an apple tree in a pot

grow an apple tree indoors



buy an apple|

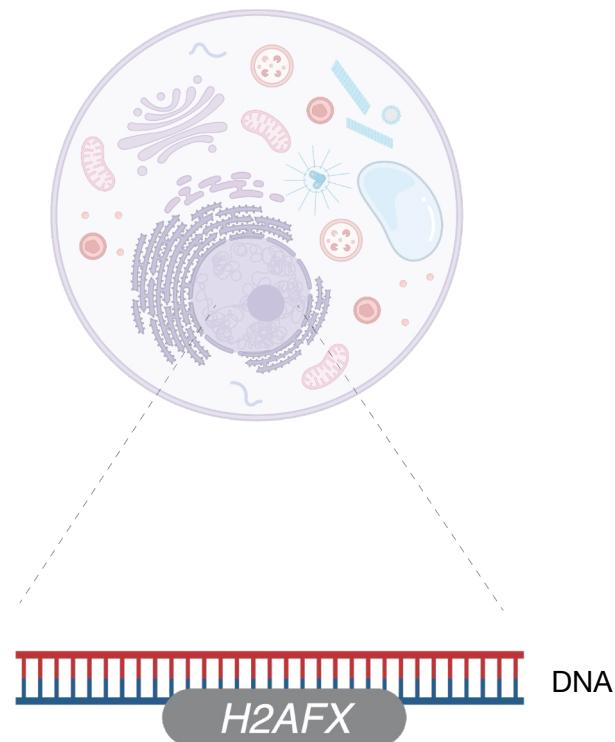
buy an apple watch

buy an apple gift card

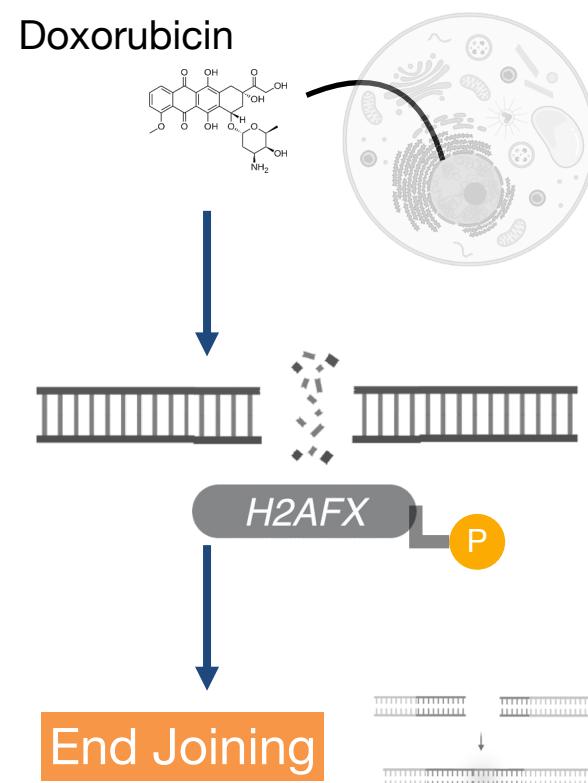
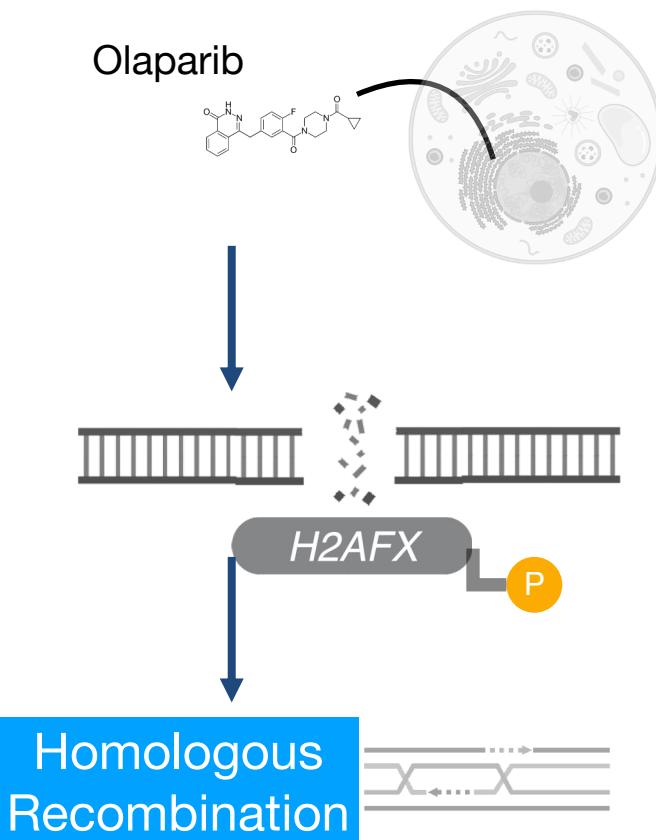
buy an apple tv



H2AFX is a **pleiotropic** gene...



... whose **particular function** is resolved via **cell context**.



While unsupervised learning of word polysemy is **common**...

**Data:** corpus  
of sentence contexts

**Approach:** word embeddings  
w/ linear semantics

$$king - man + woman \approx queen$$

unsupervised learning of gene pleiotropy is **unsolved**

**Data:** ?

**Approach:** ?

$$geneA - func1 + func2 \approx geneB$$

# Our goal for today

Unsupervised learning of gene pleiotropy with applications to therapeutic science

**Data:**

?

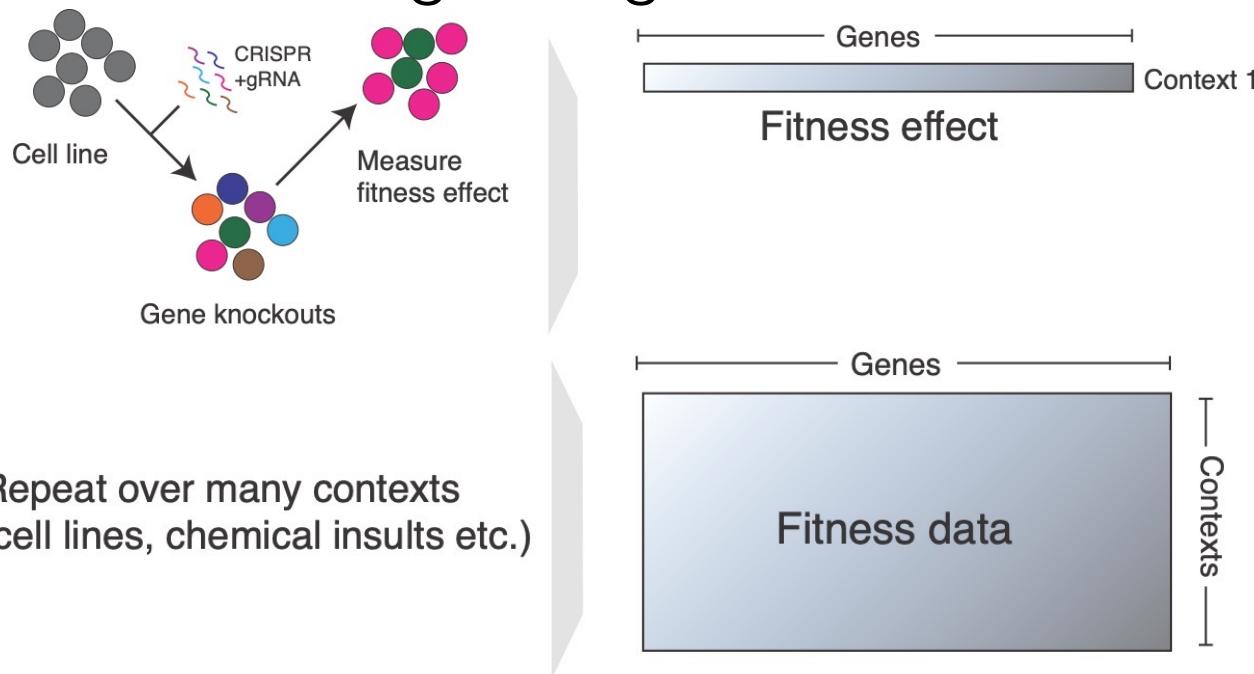
**Approach:**

?

$$geneA - func1 + func2 \approx geneB$$

# Data

Use gene perturbation effect measurements for inferring biological functions

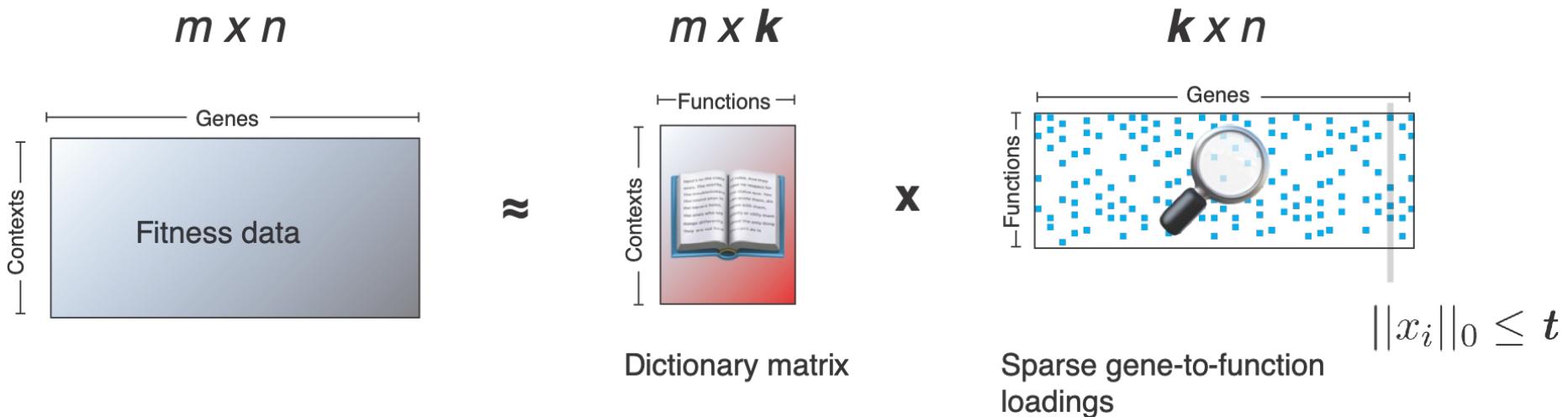


Why perturbation datasets? Alternative data types:

- **Transcriptomics:** gene co-expression is necessary but not sufficient for co-function
- **Protein-protein interactions:** direct interactions are not necessary for co-function

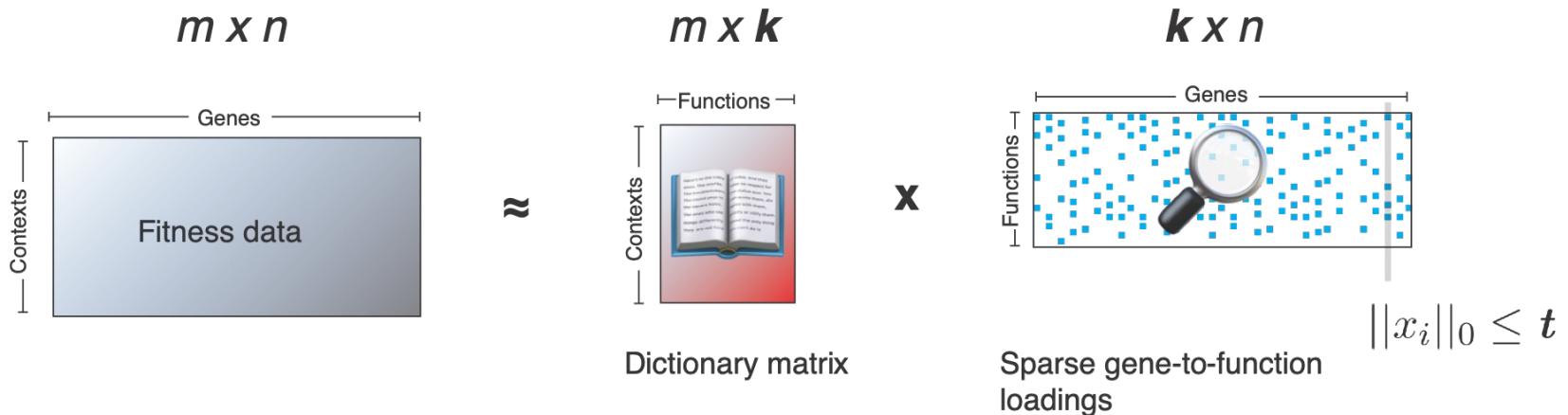
# Approach: Webster

- Low-dimensional vector embeddings that satisfy three criteria:
  - Sparse
  - Latents are biologically meaningful
  - Account for redundancy between cell contexts



# Approach: Webster

Webster learns a dictionary matrix that sparsely approximates gene effects...

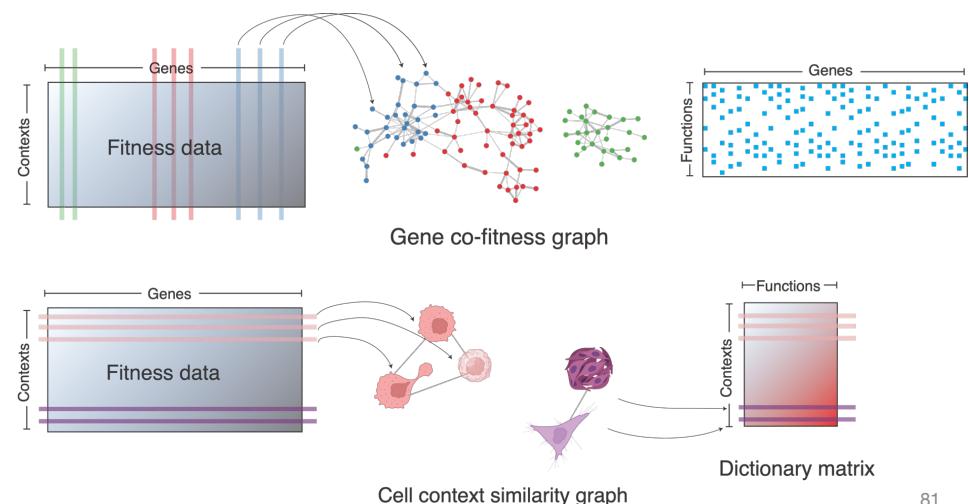


1

... while preserving  
interpretable relationships  
between genes

2

... and accounting for  
redundancies between cell  
contexts

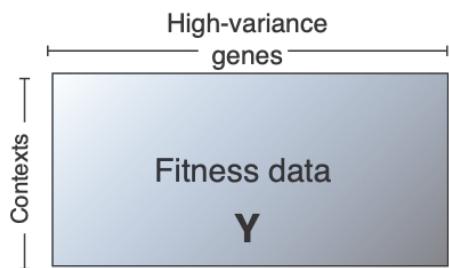


# Overview of Webster

## Preprocessing

Raw fitness data

Standardize cell lines  
Center gene effects  
Filter genes by variance



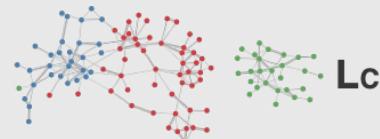
## Graph-regularized dictionary learning *Objectives*

Reduce dimensionality

$$\mathbf{Y} \approx \mathbf{D} \times \mathbf{X}$$

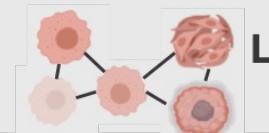
$$\|Y - DX\|_F^2$$

Preserve gene similarity



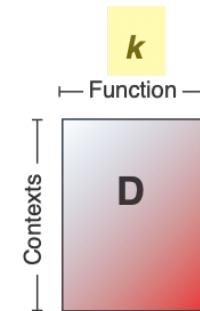
$$Tr(XL_cX^T)$$

Preserve cell context similarity

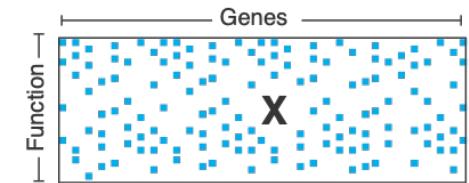


$$Tr(D^T LD)$$

## Output



Dictionary matrix

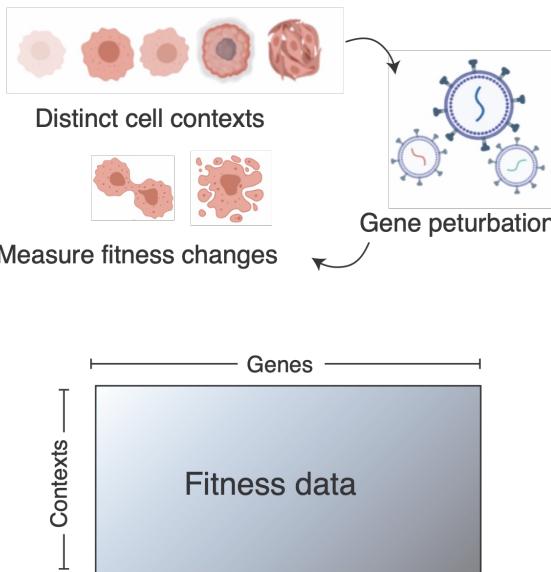
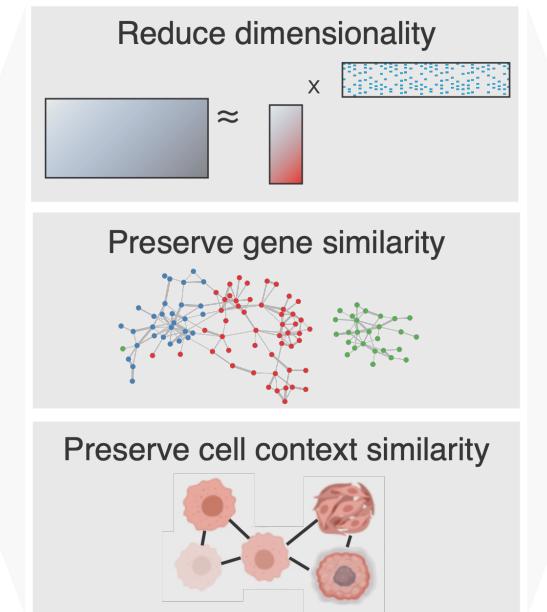
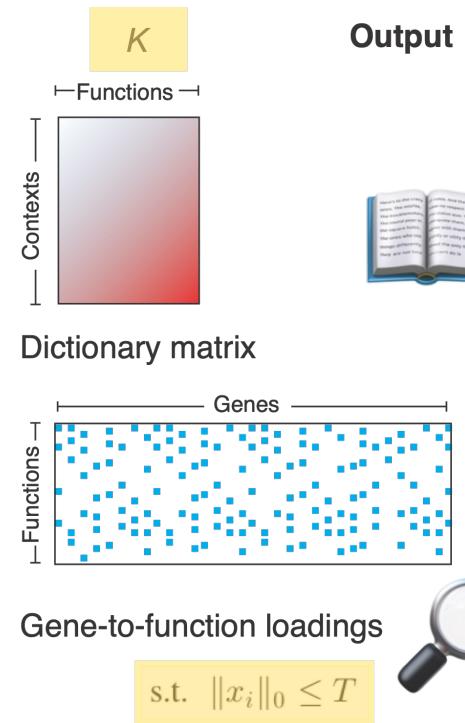


Gene-to-function loadings

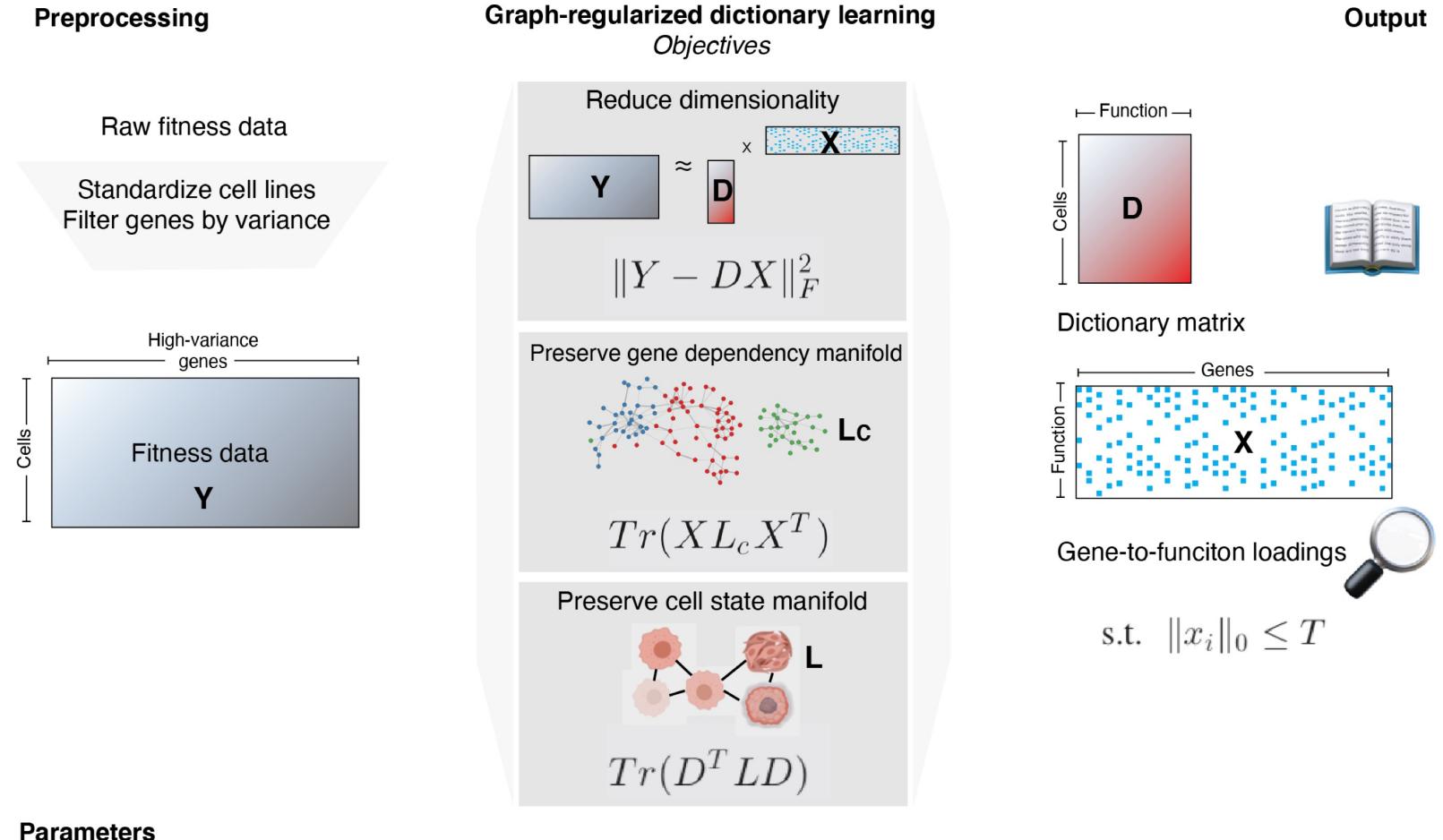
$$\|x_i\|_0 \leq t$$

**t** = key hyperparameters

# Its key parameters are dictionary size ( $K$ ) and sparsity on loadings ( $T$ )

**Input****Graph regularized dictionary learning** $K$ **Output**

# Model optimization



$k$  =latent dimension size

$\alpha$  =weight of cell Laplacian

$L$  =cell Laplacian (num neighbors, metric)

$\beta$  =weight of gene Laplacian

$L_c$  =gene Laplacian (num neighbors, metric)

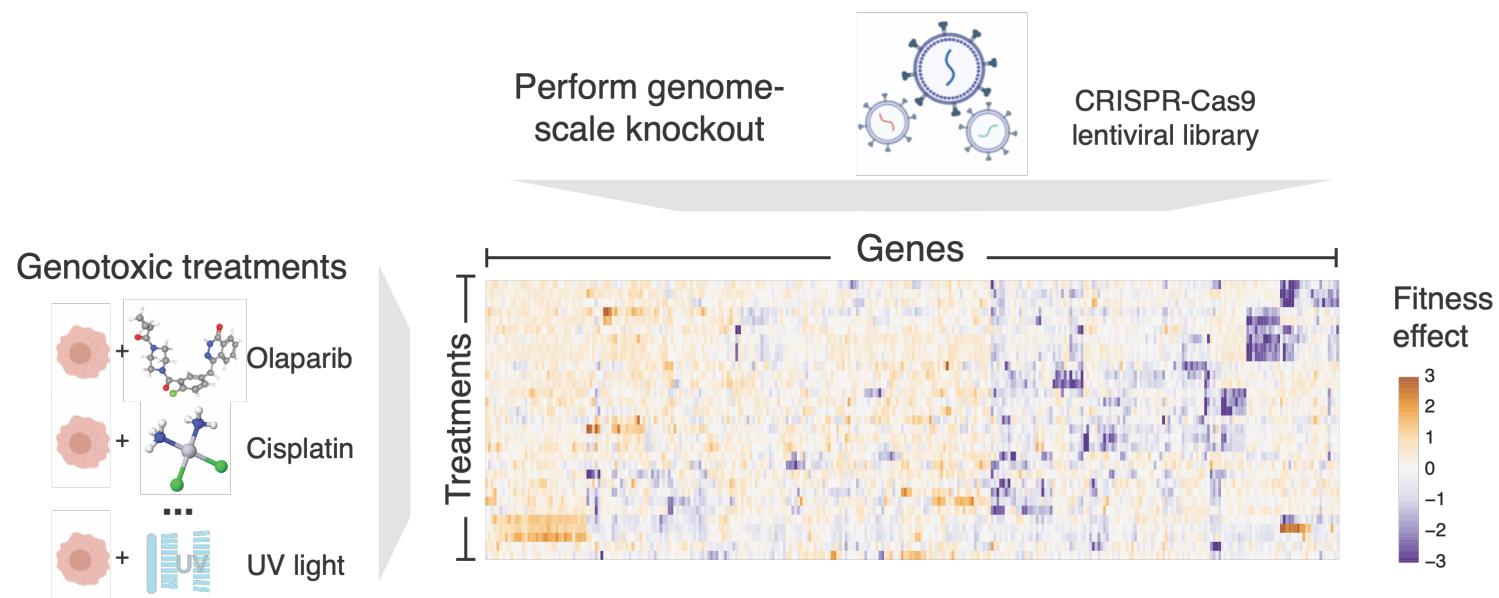
$T$  =sparsity

# Applications to three screens of gene perturbation effects

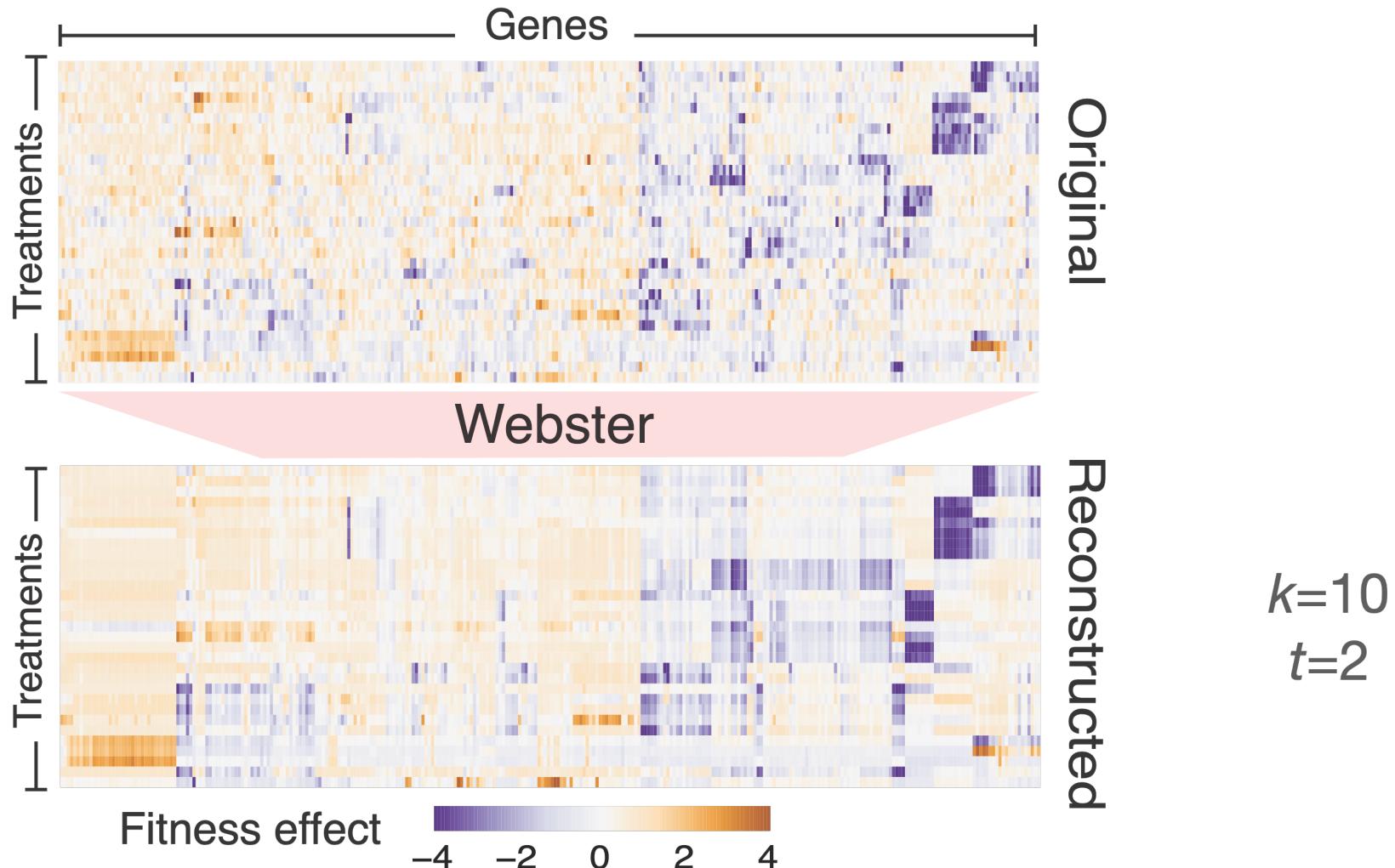
- 1) Genotoxic screens
- 2) Cancer fitness screens
- 3) Compound sensitivity screens

# Part 1: Genotoxic screens

**Olivieri et al. 2020:** fitness effect of gene knockout in presence of genotoxins

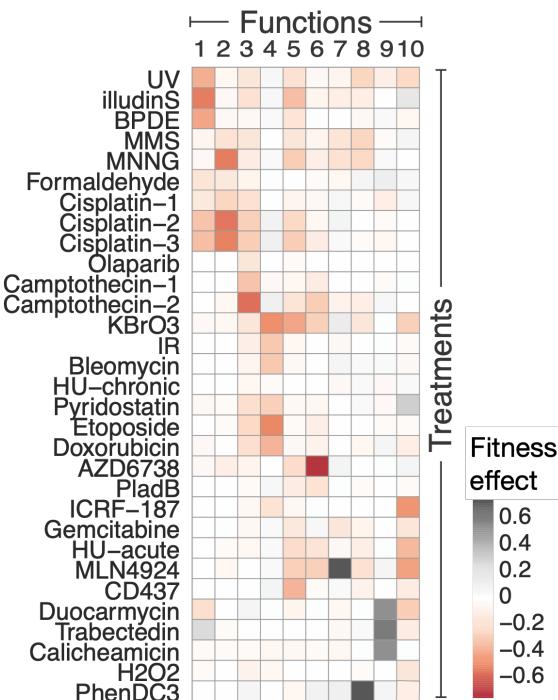


# Webster approximates the input data matrix...

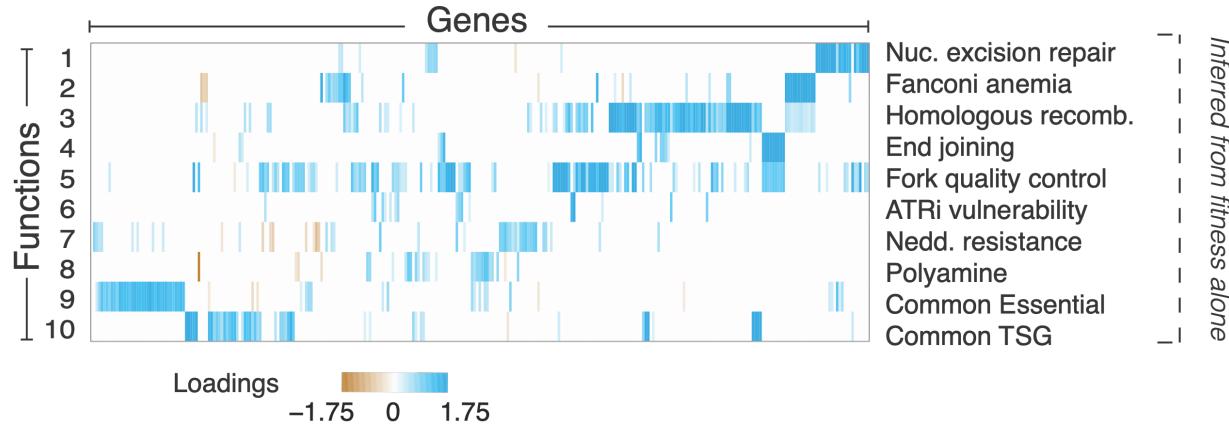


# ... as a product between a dictionary matrix and a loadings matrix

Dictionary matrix



Gene-to-function loadings

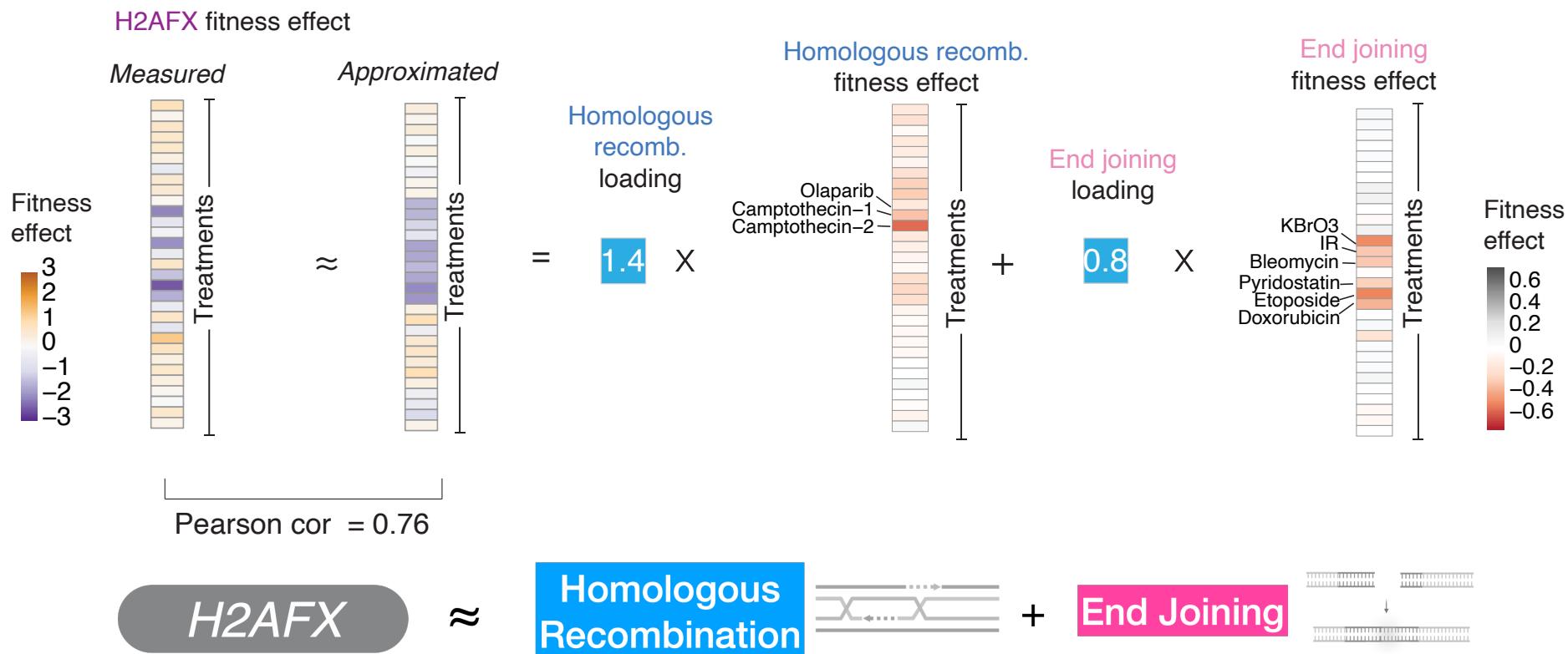


Literature annotations



Learned gene-to-function loadings recover  
biological genesets hidden during model training

# Latents inferred by the model recapitulate pleiotropy *without prior knowledge*

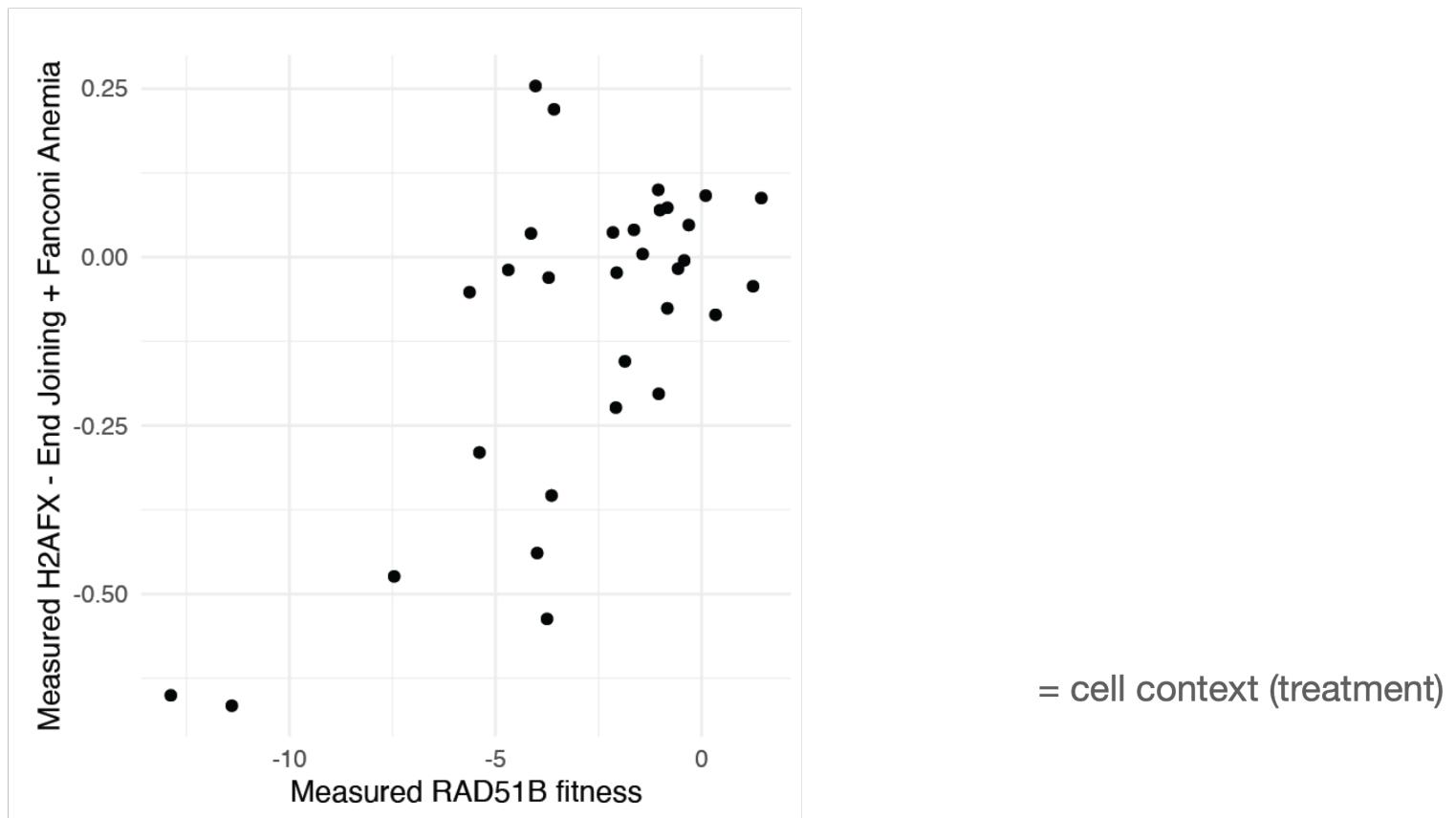


(hidden during model training!)

# Latents are biologically meaningful

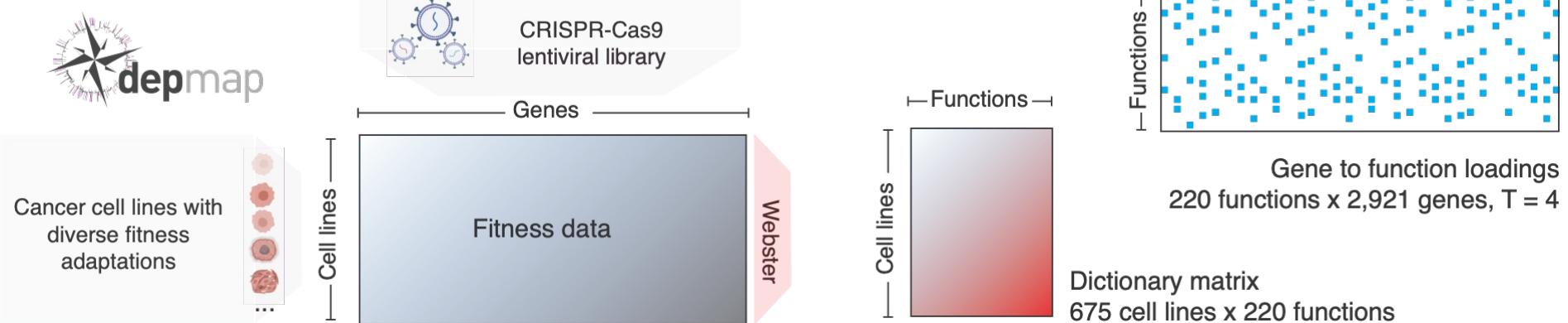
$$geneA - func1 + func2 \approx geneB$$

**H2AFX - End Joining + Fanconi Anemia  $\approx$  RAD51B**



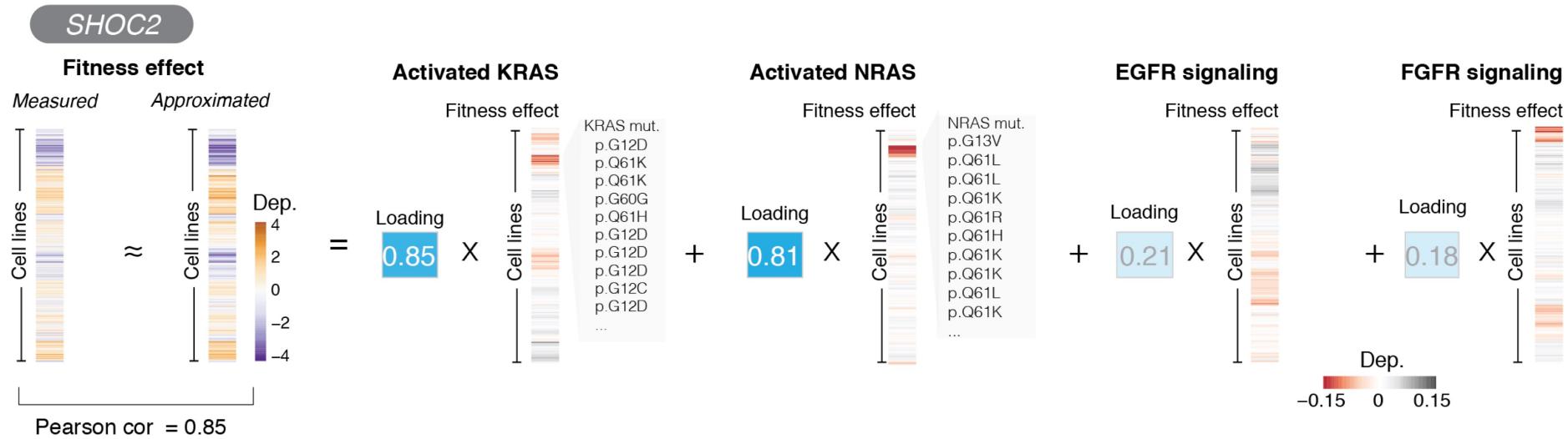
= cell context (treatment)

# Part 2: Cancer fitness screens

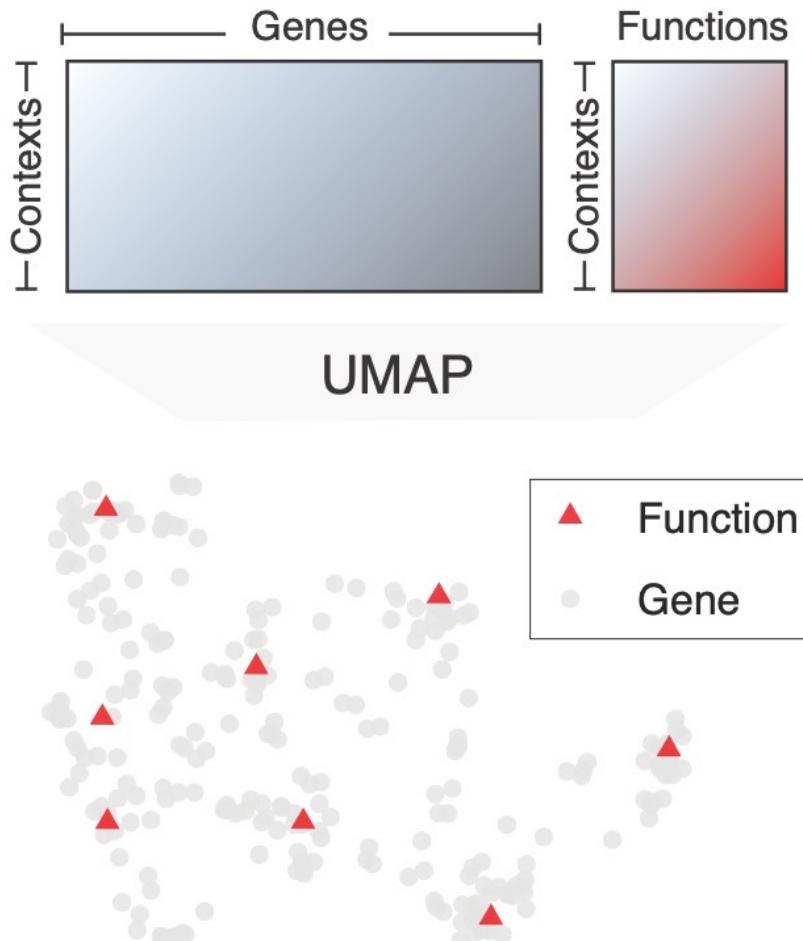


# Pleiotropic genes obey linear semantics in the latent space

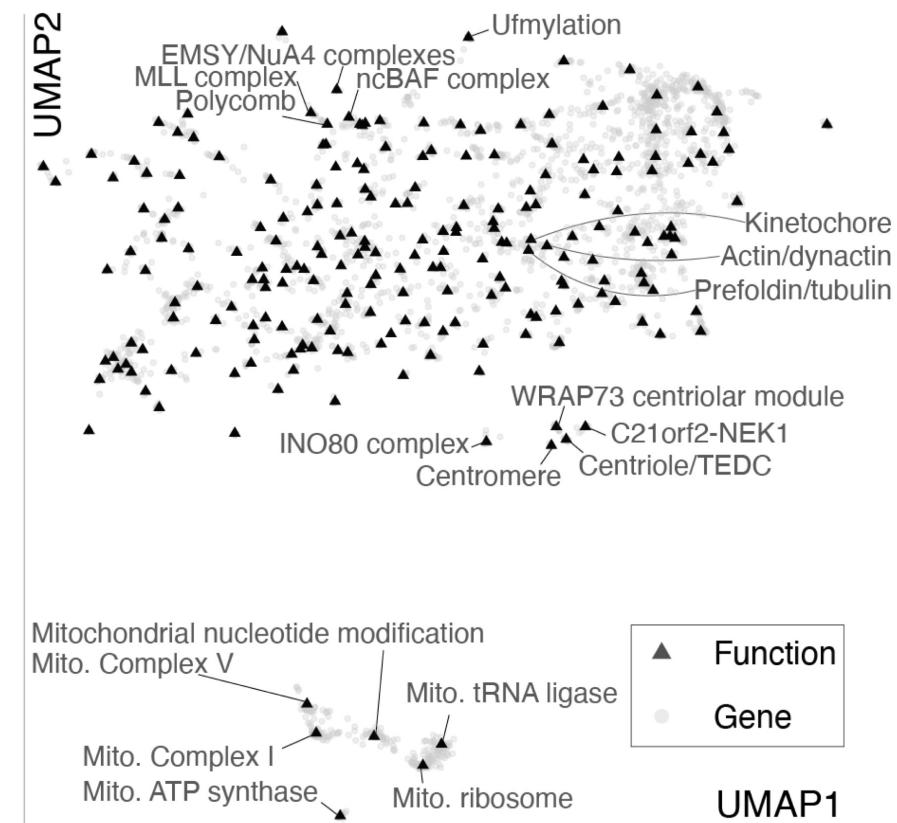
**$SHOC2 \approx \text{Activated KRAS} + \text{Activated NRAS} + \text{EGFR Signaling} + \text{FGFR Signaling}$**



# Joint embedding space of genes and functions

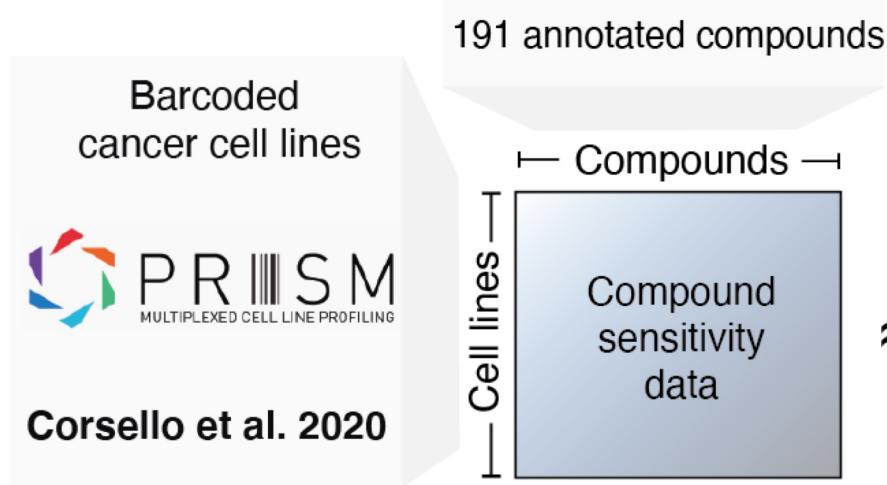


It captures interpretable processes in cancer

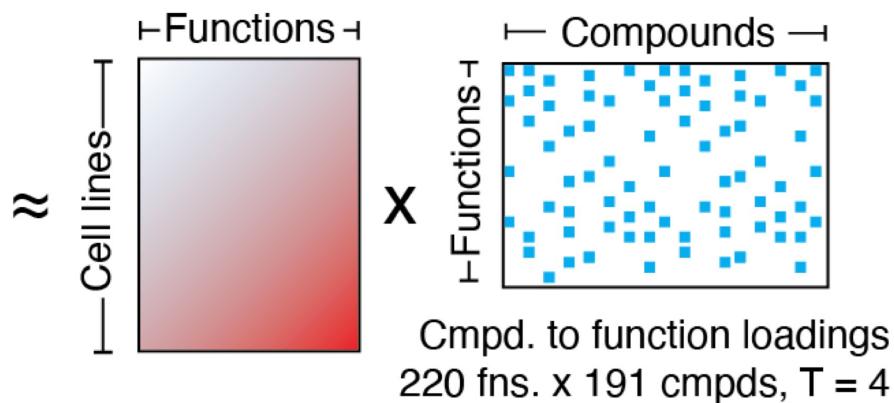


# Part 3: Compound sensitivity screens

**Query:** Drug Repurposing dataset



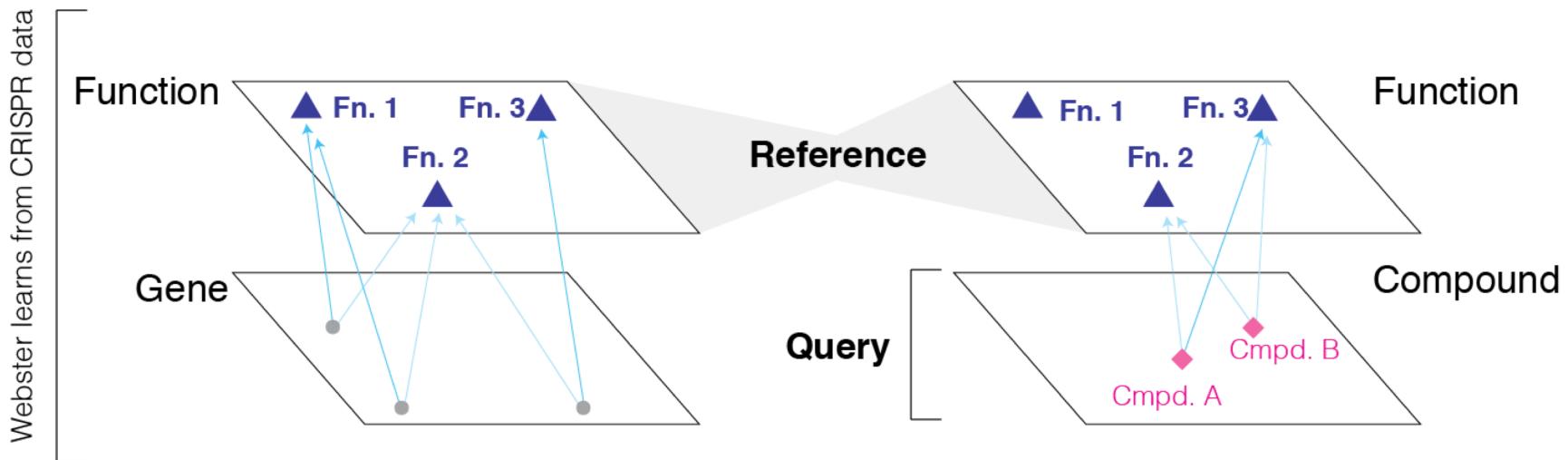
**Reference:**  
CRISPR dictionary



Modeling compound sensitivity profiles as mixtures of functions learned from CRISPR

# Modeling compounds as mixtures of latent functions

## Reference-query projection

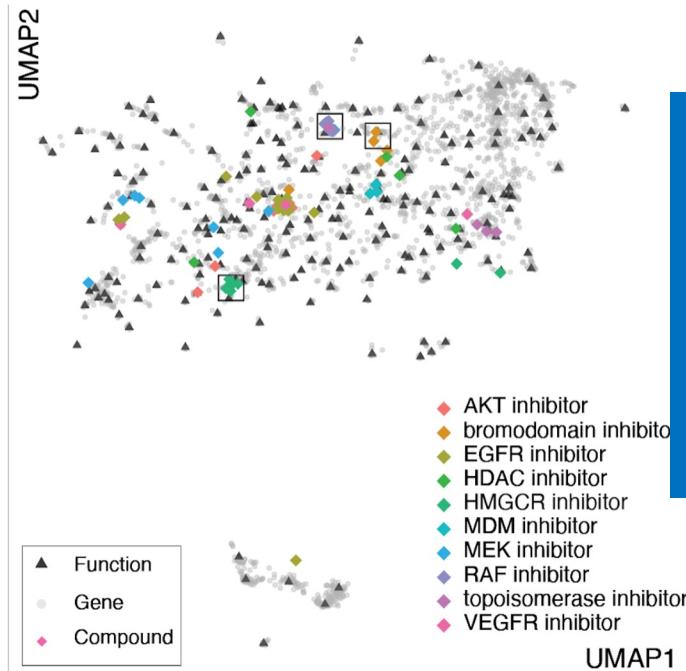


- Modeling compounds as mixtures of functions learned from CRISPR signatures with high similarity represent useful and previously unrecognized connections
  - between two proteins operating in the same pathway
  - between a small-molecule and its protein target
  - between two small-molecules of similar function but structural dissimilarity
- Such a catalog of connections can serve as a functional look-up table of compounds to predict sensitivity and genotoxic profiles and to inform therapeutic use

# Compounds' mechanisms of action

Compounds are embedded nearby gene functions, reflecting their mechanism of action

Projecting compound sensitivity into gene fn. map



BRAF signaling

Loadings

BRAF  
SOX10  
SOX9

H2A.Z maintenance

Loadings

KDM2A  
H2AFZ  
KANSL3

Mevalonate synthesis

Loadings

UBIAD1  
HMGCR  
MVK

Refer

Modeling compounds as mixtures of functions learned from CRISPR signatures with high similarity represent useful and previously unrecognized connections

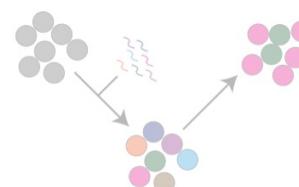
- between two proteins operating in the same pathway
- between a compound and its protein target
- between two compounds of similar function but structural dissimilarity



# Key takeaways

- Analogously to word semantics, genes can be modeled as **distributions over latent bio functions**
  - **Sparse learning** is an effective strategy for learning bio functions from high-dimensional chemical and genetic perturbations
  - New perturbations can be **projected** into learned space

**Data:** high-dimensional gene perturbation measurements



**Approach:** sparse approximation embeddings

$$\begin{matrix} & \times \\ \begin{matrix} \text{red gradient} \end{matrix} & \times \begin{matrix} \text{blue sparse matrix} \end{matrix} \end{matrix}$$
A diagram illustrating the mathematical operation of sparse approximation. It shows a red gradient square being multiplied by a blue sparse matrix, resulting in a final output.

$$geneA - func1 + func2 \approx geneB$$

# https://depmap.org/webster

Webster

Published Paper at Cell Systems    Code for paper    Dictionary learning code    Figshare data  
Design write-up

Explore relationships between genes and biological functions learned from CRISPR fitness screens using Webster.

Read The Paper: "Sparse Dictionary Learning Recovers Pleiotropy From Human Cell Fitness Screens" [For More Details.](#)

+ About this tool

Genotoxic

Select function group

ATRi vulnerability (V3)

Nedd. resistance (V5)

Polyamine (V1)

Search to select a gene or function

Selected function:  
**ATRi vulnerability (V3)**

Pan UMAP w/  
w  
A S D  
OR  
← →  
↑ ↓  
Q Q

2d   3d  

● Functions   ● Genes   ● Gene positive association   ● Gene negative association

Native mouse controls: <=>= pan right left   ^~= zoom

highlighted in plot

Gene  
**DHX35**  
(ex. ### loading, function name)  
**1.08**  
ATRi vulnerability (V3)

1.00  
Fork quality control (V9)

Approximation quality (Pearson)  
0.74

# Outline for today's class

- Optimization & generation of small molecules
- Binding of drugs to therapeutic targets
- High-throughput genetic & chemical perturbations

