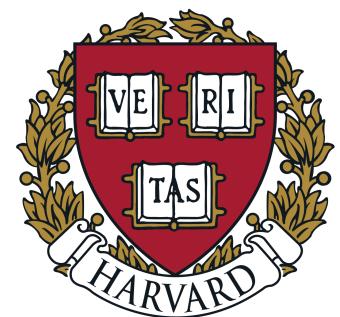# Machine Learning for Drug Development

## Marinka Zitnik

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

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

# Outline

✓ Overview and introduction

✓ **Part 1:** Virtual drug screening and drug repurposing

✓ **Part 2:** Adverse drug effects, drug-drug interactions

✓ **Part 3:** Clinical trial site identification, patient recruitment

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

✓ **Part 5:** Molecular property prediction and transformers

Demos, resources, wrap-up & future directions 👉

# Datasets to facilitate algorithmic innovation

# Therapeutics are one of most exciting areas for computational scientists. However,

Retrieving, curating, and processing datasets is time-consuming and requires extensive domain expertise

Datasets are scattered around the bio repositories and there is no centralized data repository for a variety of therapeutics
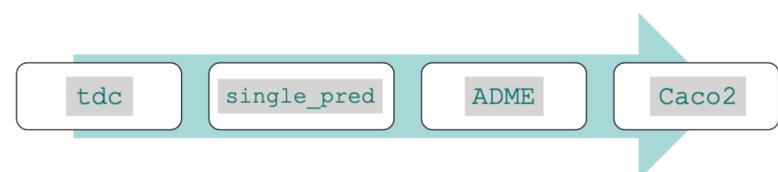
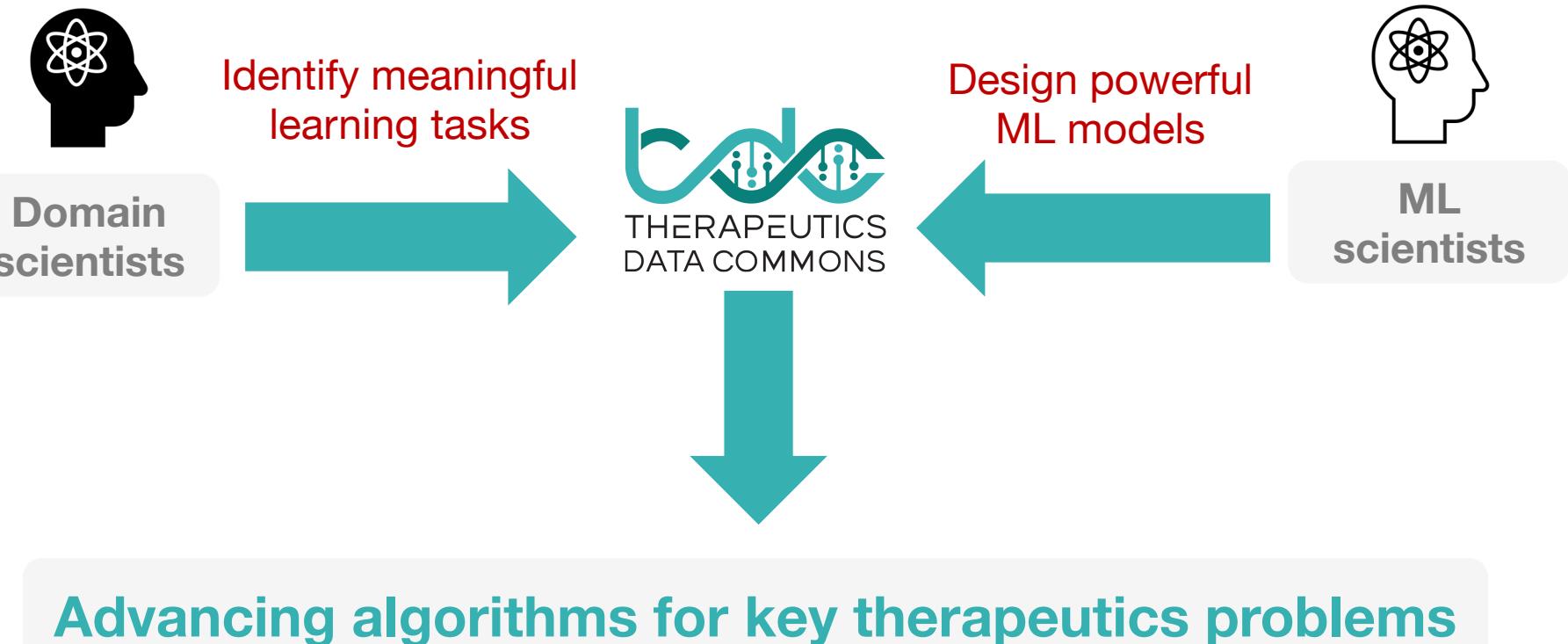Many tasks are under-explored in AI/ML community because of the lack of data access

# THERAPEUTICS DATA COMMONS

- **Open-Source ML Datasets for Therapeutics:**
  - Wide range of tasks: target discovery, activity screening, efficacy, safety, manufacturing
  - Wide range of products: small molecules, antibodies, vaccine, miRNA
- **Numerous Data Functions:**
  - Extensive data functions
  - Model evaluation, data processing and splits, molecule generation oracles, and much more
- **3 Lines of Code:**
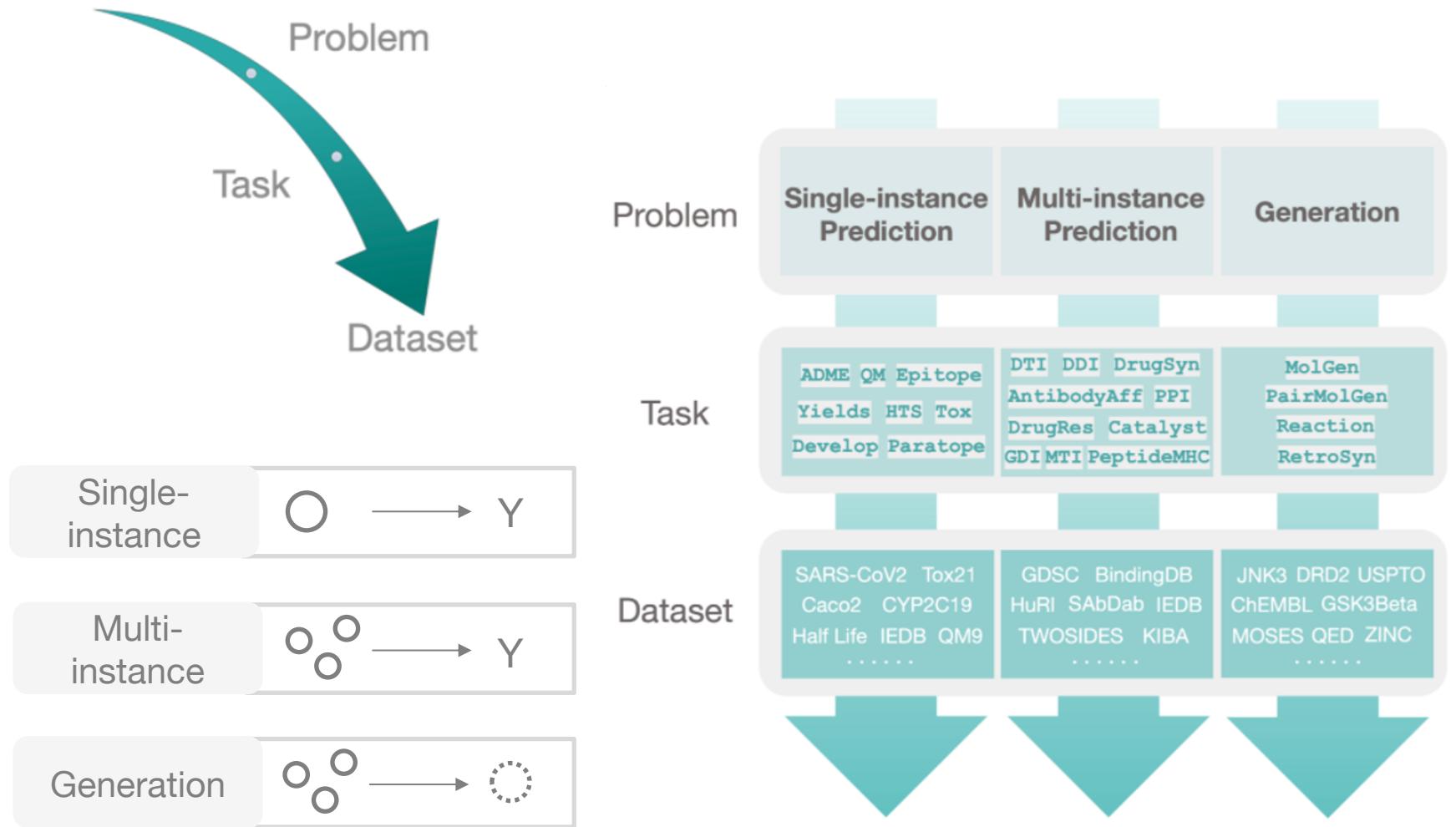  - Minimum package dependency, lightweight loaders

```
from tdc.single_pred import ADME
data = ADME(name = 'Caco2_Wang')
splits = data.split()
```

| tdc | single_pred | ADME | Caco2 |

# Our Vision for TDC



Identify meaningful learning tasks

Design powerful ML models

Domain scientists

THERAPEUTICS DATA COMMONS

ML scientists

**Advancing algorithms for key therapeutics problems**

# Modular Structure of TDC

DATASET INDEX

Absorption

Caco-2 (Cell Effective Permeability), Wang et al.

HIA (Human Intestinal Absorption), Hou et al.

Pgp (P-glycoprotein) Inhibition, Broccatelli et al.

Bioavailability, Ma et al.

Bioavailability F20/F30, eDrug3D

Lipophilicity, AstraZeneca

Solubility, AqSolDB

Solubility, ESOL

Hydration Free Energy, FreeSolv

Distribution

**ADME**

BBB (Blood-Brain Barrier), Adenot et al.

BBB (Blood-Brain Barrier), Martins et al.

PPBR (Plasma Protein Binding Rate), Ma et al.

PPBR (Plasma Protein Binding Rate), eDrug3D

VD (Volumn of Distribution), eDrug3D

Metabolism

CYP P450 2C19 Inhibition, Veith et al.

CYP P450 2D6 Inhibition, Veith et al.

CYP P450 3A4 Inhibition, Veith et al.

CYP P450 1A2 Inhibition, Veith et al.

CYP P450 2C9 Inhibition, Veith et al.

Excretion

Half Life, eDrug3D

Clearance, eDrug3D

---

DATASET INDEX

BindingDB

DAVIS

KIBA

**DTI**

---

DATASET INDEX

SARS-CoV-2 In Vitro, Touret et al.

SARS-CoV-2 3CL Protease, Diamond.

HIV

**HTS**

---

DATASET INDEX

IEDB, Jespersen et al.

PDB, Jespersen et al.

**Epitope**

---

DATASET INDEX

TAP

SAbDab, Chen et al.

**Develop**

---

DATASET INDEX

DisGeNET

**GDA**

---

DATASET INDEX

GDSC1

GDSC2

**DrugRes**

---

DATASET INDEX

OncoPolyPharmacology

**DrugSyn**

**MHC**

**yAff**

miRTarBase

**MTI**

---

DATASET INDEX

USPTO

**Catalyst**

---

DATASET INDEX

DrugBank Multi-Typed DDI

TWOSIDES Polypharmacy Side Effects

**DDI**

---

DATASET INDEX

Tox21

ToxCast

ClinTox

**Tox**

---

DATASET INDEX

USPTO

**Reaction**

---

DATASET INDEX

MOSES

ZINC

ChEMBL

**MolGen**

---

DATASET INDEX

DRD2

QED

LogP

**PairMolGen**

---

DATASET INDEX

USPTO-50K

USPTO

**RetroSyn**

---

DATASET INDEX

HuRI

**PPI**
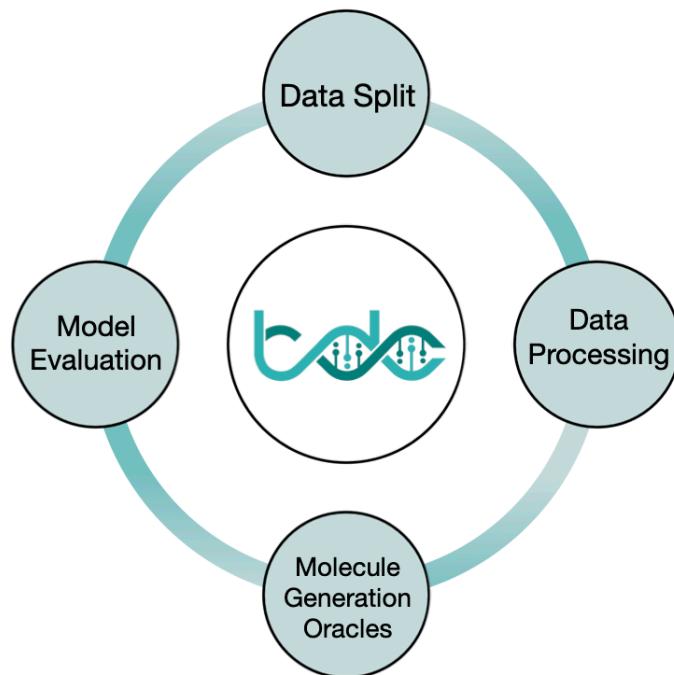
---

DATASET INDEX

Buchwald-Hartwig

USPTO

**Yields**

---

# 67 datasets spread over 22 learning tasks

# Data Functions to Support Your Research



## Model performance evaluators

FUNCTION INDEX

Regression Metric

  Mean Squared Error (MSE)

  Mean Absolute Error (MAE)

  Coefficient of Determination ($R^2$)

Binary Classification Metric

  Area Under the Receiver Operating Characteristic Curve (ROC-AUC)

  Area Under the Precision-Recall Curve (PR-AUC)

  Accuracy Metric

  Precision

  Recall

  F1 Score

Multi-class Classification Metric

  Micro-F1, Micro-Precision, Micro-Recall, Accuracy

  Macro-F1

  Cohen's Kappa (Kappa)

Token-level Classification Metric

  Average ROC-AUC

## A variety of data splits

FUNCTION INDEX

Data Split Overview

  Random Split

  Scaffold Split

  Cold-Start Split

## Data processing helpers

FUNCTION INDEX

Label Distribution Visualization

Label Binarization

Label Units Conversion

Label Meaning

Basic Statistics

Data Balancing

Graph Transformation for Pair Data

Negative Samples for Pair Data

From PubChem CID to SMILES

From Uniprot ID to Amino Acid Sequence

# Molecule Generation Oracles

### Molecule Generation



Generated Molecules → Oracle → Score → Optimize

GuacaMol

MOSES

Literature

```
In [ ]: |
```

**FUNCTION INDEX**

**Goal-oriented Oracles**

- Glycogen Synthase Kinase 3 Beta (GSK3β)
- c-Jun N-terminal Kinases-3 (JNK3)
- Dopamine Receptor D2 (DRD2)
- Synthetic Accessibility (SA)
- IBM RXN Synthetic Accessibility (IBM_RXN)
- Quantitative Estimate of Drug-likeness (QED)
- Octanol-water Partition Coefficient (LogP)
- Rediscovery
- Similarity/Dissimilarity
- Median Molecules
- Isomers
- Multi-Property Objective (MPO)
- Valsartan SMARTS
- Hop

**Distribution Learning Oracles**

- Diversity
- KL divergence
- Frechet ChemNet Distance (FCD)
- Novelty
- Validity
- Uniqueness

GuacaMol: Benchmarking Models for de Novo Molecular Design, J. Chem. Inf. Model., 2019
MOSES: A Benchmarking Platform for Molecular Generation Models, Frontiers in Pharmacology, 2020

# Leaderboards: Submit your Models



zitniklab.hms.harvard.edu/TDC

# You Are Invited to Join TDC! TDC is an Open-Source, Community Effort



**zitniklab.hms.harvard.edu/TDC**

**github.com/mims-harvard/TDC**

## Tutorials

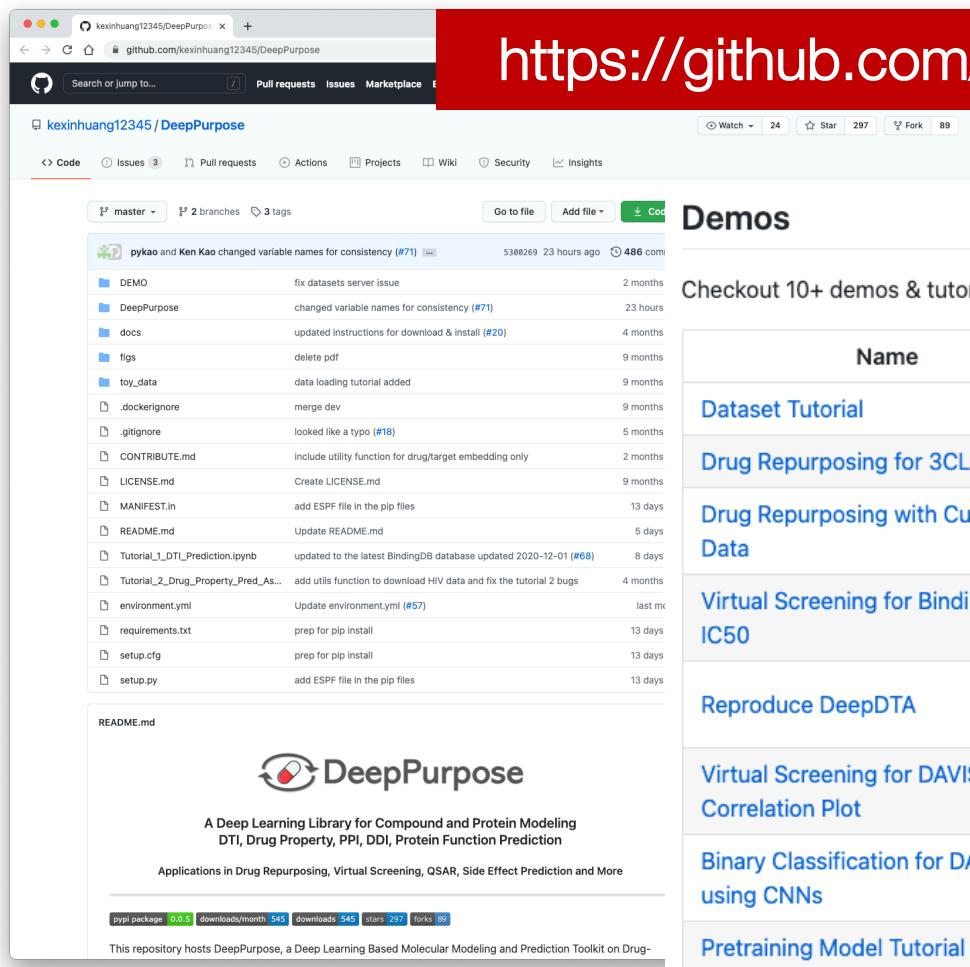We provide a series of tutorials for you to get started using TDC:

| Name | Description |
|------|-------------|
| 101 | Introduce TDC Data Loaders |
| 102 | Introduce TDC Data Functions |
| 103.1 | Walk through TDC Small Molecule Datasets |
| 103.2 | Walk through TDC Biologics Datasets |
| 104 | Generate 21 ADME ML Predictors with 15 Lines of Code |
| 105 | Molecule Generation Oracles |

`pip install PyTDC`

# Demos, tools, and implementations

# DeepPurpose: Deep Learning Library for Compound and Protein Modeling
# DTI, Drug Property, PPI, DDI, Protein Function Prediction

https://github.com/kexinhuang12345/DeepPurpose



## Demos

Checkout 10+ demos & tutorials to start:

| Name | Description |
|---|---|
| Dataset Tutorial | Tutorial on how to use the dataset loader and read customized data |
| Drug Repurposing for 3CLPro | Example of one-liner repurposing for 3CLPro |
| Drug Repurposing with Customized Data | Example of one-liner repurposing with AID1706 Bioassay Data, training from scratch |
| Virtual Screening for BindingDB IC50 | Example of one-liner virtual screening |
| Reproduce DeepDTA | Reproduce DeepDTA with DAVIS dataset and show how to use the 10 lines framework |
| Virtual Screening for DAVIS and Correlation Plot | Example of one-liner virtual screening and evaluate on unseen dataset by plotting correlation |
| Binary Classification for DAVIS using CNNs | Binary Classification for DAVIS dataset using CNN encodings by using the 10 lines framework. |
| Pretraining Model Tutorial | Tutorial on how to load pretraining models |

and more in the DEMO folder!

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction, *Bioinformatics* 2020

# How can domain scientists interact with AI systems?

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction, *Bioinformatics* 2020
MolDesigner: Interactive Design of Efficacious Drugs with Deep Learning, *NeurIPS* 2020

# MolDesigner: Interactive Design of Drugs with Deep Learning



http://deeppurpose.sunlab.org

# DEMO: DRUG-TARGET INTERACTION PREDICTION

Drug: [Remdesivir](#) Remdesivir is indicated for the treatment of adult and pediatric patients aged 12 years and over weighing at least 40 kg for coronavirus disease 2019 (COVID-19) infection requiring hospitalization.

Target protein: [Replicase polyprotein 1ab](#). Multifunctional protein involved in the transcription and replication of viral RNAs



```
>lcl|BSEQ0052511|Replicase polyprotein 1ab
MESLVPGFNEKTHVQLSLPVLQVRDVLVRGFGDSVEEVLSEARQHLKDGTCGLVEVEKGV
LPQLEQPYVFIKRSDARTAPHGHVMVELVAELEGIQYGRSGETLGVLVPHVGEIPVAYRK
VLLRKNGNKGAGGHSYGADLKSFDLGDELGTDPYEDFQENWNTKHSSGVTRELMRELNGG
AYTRYVDNNFCGPDGYPLECIKDLLARAGKASCTLSEQLDFIDTKRGVYCCREHEHEIAW
YTERSEKSYELQTPFEIKLAKKFDTFNGECPNFVFPLNSIIKTIQPRVEKKKLDGFMGRI
RSVYPVASPNECNQMCLSTLMKCDHCGETSWQTGDFVKATCEFCGTENLTKEGATTCGYL
PQNAVVKIYCPACHNSEVGPEHSLAEYHNESGLKTILRKGGRTIAFGGCVFSYVGCHNKC
AYWVPRASANIGCNHTGVVGEGSEGLNDNLLEILQKEKVNINIVGDFKLNEEIAIILASF
SASTSAFVETVKGLDYKAFKQIVESCGNFKVTKGKAKKGAWNIGEQKSILSPLYAFASEA
ARVVRSIFSRTLETAQNSVRVLQKAAITILDGISQYSLRLIDAMMFTSDLATNNLVVMAY
ITGGVVQLTSQWLTNIFGTVYEKLKPVLDWLEEKFKEGVEFLRDGWEIVKFISTCACEIV
GGQIVTCAKEIKESVQTFFKLVNKFLALCADSIIIGGAKLKALNLGETFVTHSKGLYRKC
VKSREETGLLMPLKAPKEIIFLEGETLPTEVLTEEVVLKTGDLQPLEQPTSEAVEAPLVG
TPVCINGLMLLEIKDTEKYCALAPNMMVTNNTFTLKGGAPTKVTFGDDTVIEVQGYKSVN
ITFELDERIDKVLNEKCSAYTVELGTEVNEFACVVADAVIKTLQPVSELLTPLGIDLDEW
SMATYYLFDESGEFKLASHMYCSFYPPDEDEEEGDCEEEEFEPSTQYEYGTEDDYQGKPL
EFGATSAALQPEEEQEEDWLDDDSQQTVGQQDGSEDNQTTTIQTIVEVQPQLEMELTPVV
QTIEVNSFSGYLKLTDNVYIKNADIVEEAKKVKPTVVVNAANVYLKHGGGVAGALNKATN
```

Molecular structure of Remdesivir

Amino acid sequence of Replicase polyprotein 1ab

# How can domain scientists interact with AI systems?

DeepPurpose: a Deep Learning Library for Drug-Target Interaction Prediction, *Bioinformatics* 2020
MolDesigner: Interactive Design of Efficacious Drugs with Deep Learning, *NeurIPS* 2020

# Automating Science

# Automating Science

# How to explain predictions?

## Key idea:

- Summarize where in the data the model "looks" for evidence for its prediction

- Find a small subgraph most influential for the prediction



GNN model training and predictions

Explaning GNN's predictions

GNNExplainer

Approach to generate explanations for graph neural networks based on counterfactual reasoning

# GNNExplainer: Key Idea

- Input: Given prediction $f(x)$ for node/link $x$

- Output: Explanation, a small subgraph $M_x$ together with a small subset of node features:

  - $M_x$ is most influential for prediction $f(x)$

- **Approach:** Learn $M_x$ via counterfactual reasoning

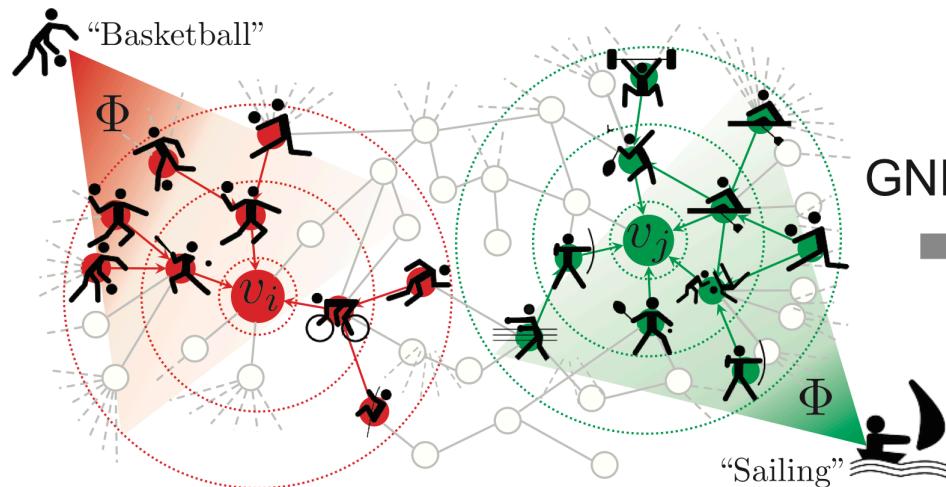  - Intuition: If removing $v$ from the graph strongly decreases the probability of prediction $\Rightarrow v$ is a good counterfactual explanation for the prediction



☐ Node feature vector     ✖ Feature excluded from explanation

GNN Explainer: Generating Explanations for Graph Neural Networks, *NeurIPS* 2019

# Examples of Explanations



"Will rosuvastatin treat hyperlipidemia? What is the disease treatment mechanism?"

New Algorithms: GNNExplainer: Generating Explanations for Graph Neural Networks, *NeurIPS* 2019
New Insights: Discovery of Disease Treatment Mechanisms through the Multiscale Interactome, *Nature Communications 2021* (in press)

# Open challenges and future directions

# Learn about Therapeutics ML!



[https://www.drugsymposium.org](https://www.drugsymposium.org)

Videos from the presentations are now publicly available to everyone through the Symposium Video Channel

# Open Challenges

- **Disconnected, uncoupled biomedical knowledge:**
  - <u>Challenge:</u> Need to combine data in their broadest sense to close the gap between research and patient data

- **Diverse mechanisms of drug action:**
  - <u>Challenge:</u> Need to consider diverse mechanisms through which a drug can treat a disease

- **Novel drugs in development, emerging diseases:**
  - <u>Challenge:</u> Need to learn and reason about never-before-seen phenomena

- **Datasets for a variety of therapeutics tasks:**
  - <u>Challenge:</u> Need datasets and benchmarks to accelerate ML model development, validation and transition into production and clinical implementation

# Outline

✓ Overview and introduction

✓ Part 1: Virtual drug screening and drug repurposing

✓ Part 2: Adverse drug effects, drug-drug interactions

✓ Part 3: Clinical trial site identification, patient recruitment

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

✓ Part 5: Molecular property prediction and transformers

✓ Demos, resources, wrap-up & future directions

https://zitniklab.hms.harvard.edu/drugml