

# Molecular Machine Learning

Bharath Ramsundar

# What is Molecular Machine Learning?

- The biggest applied fields of machine learning today are computer vision, natural language processing, speech processing, and video processing.
- **Molecular ML** is the field attempting to build machine-learning architectures and code-bases for learning to predict properties of molecules
- There are many challenges unique to molecular ML (as opposed to other machine learning fields).

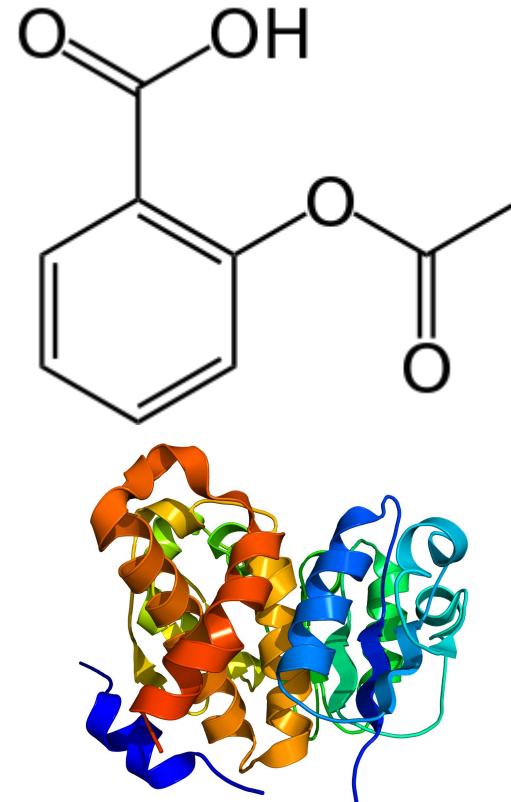
# Molecular ML Challenges: Data Scarcity

- Limited data available publicly.
- Not easy to gather data. Most experiments require highly-trained scientists or technicians to perform. (No Mechanical Turk of data).
- Pharmaceutical companies and biotechs have data, but tough to get due to IP concerns.



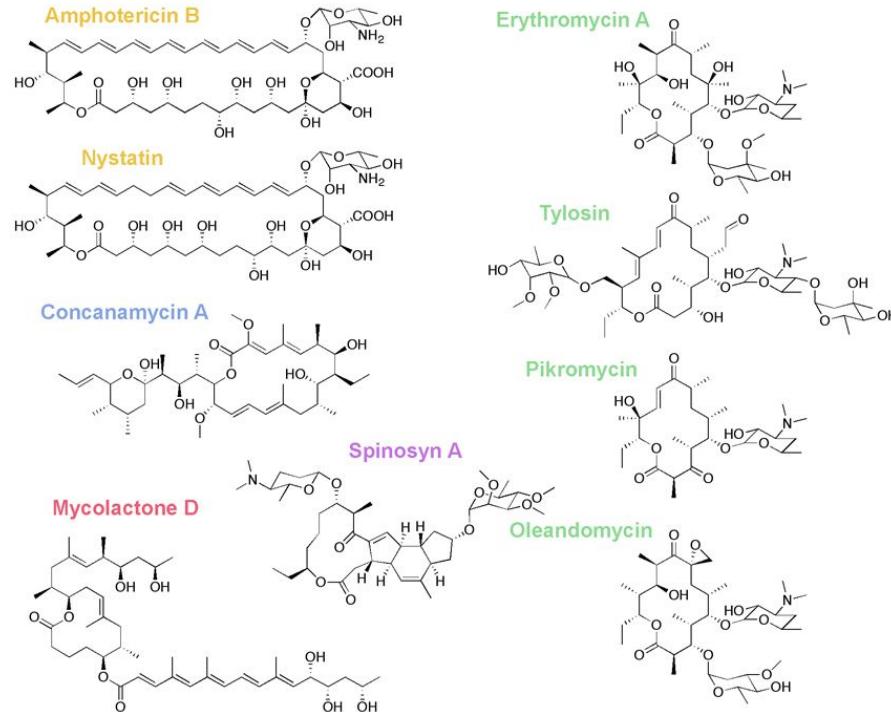
# Molecular ML Challenge: Featurization

- Molecules come in many sizes and shapes.
- How can a molecule be transformed into a vector/matrix for machine learning?
- Turns out different representations needed for different problems.



# Molecular ML Challenge: Generalization

- Chemical space is vast ( $10^{60}$  drug-like molecules) and only a small fraction (100 M compounds) actually made.
- So test set for models has to be different from training set! We want to know which new molecules are useful.

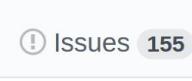
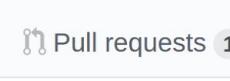
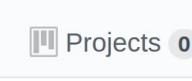


# DeepChem: <https://deepchem.io>

- DeepChem is a molecular machine learning toolchain built on Tensorflow, being used for a variety of projects, including drug-discovery, quantum chemistry, solar cell design, and structural biology.

 [deepchem / deepchem](#)

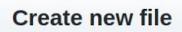
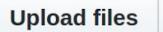
  

Democratizing Deep-Learning for Drug Discovery, Quantum Chemistry, Materials Science and Biology <https://deepchem.io/> 



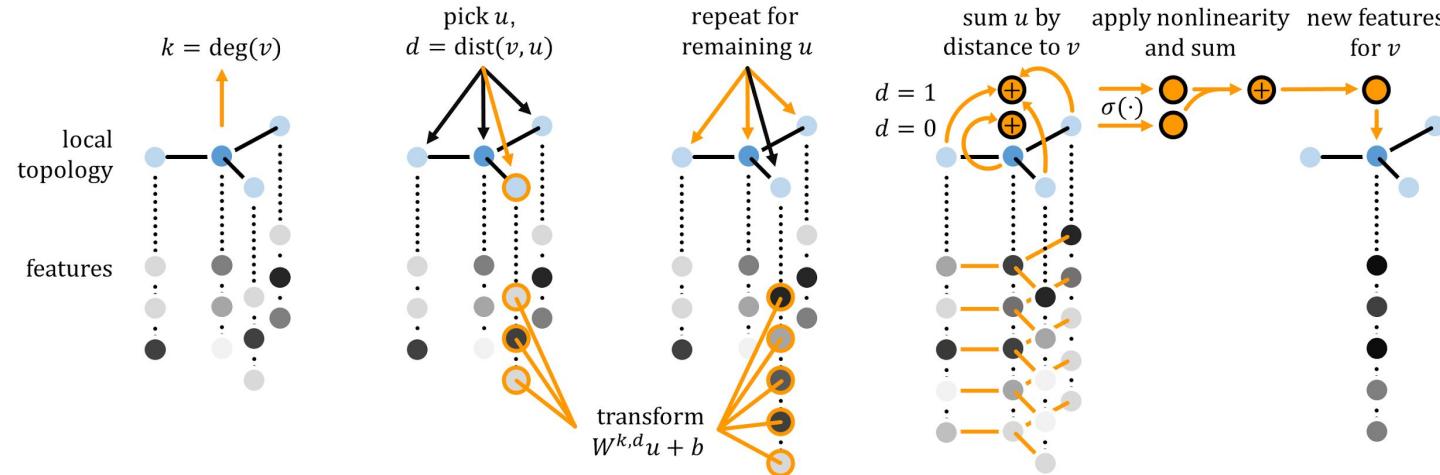
 2,942 commits  5 branches  9 releases  36 contributors  MIT

 Branch: master     

 **rbharath** Merge pull request #924 from deepchem/rbharath-patch-3  Latest commit 401d669 22 hours ago

# **Graph Convolutions**

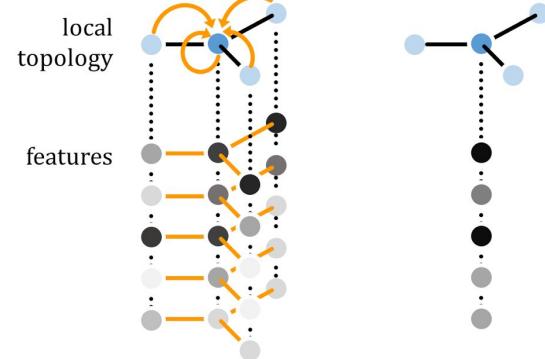
# Graph Convolution



# Graph Pool

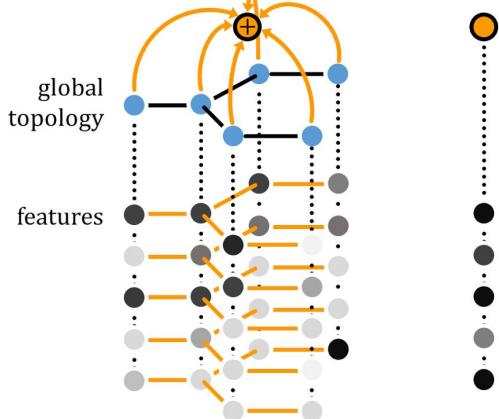
max over  
neighbors and self

new features  
for  $v$

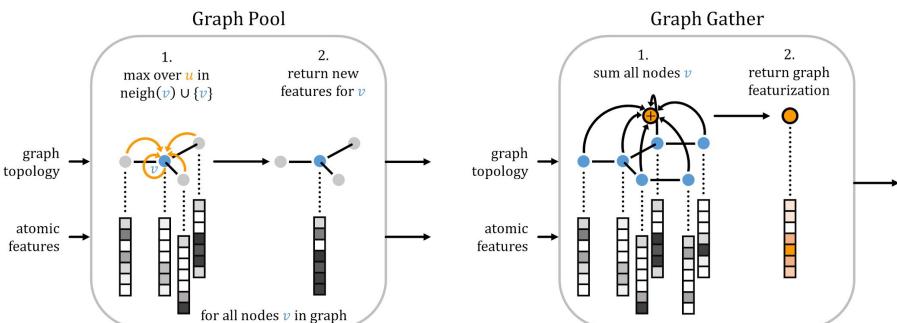
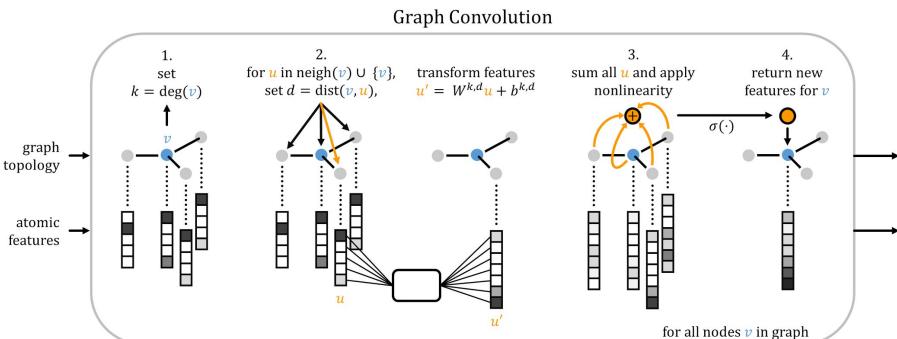
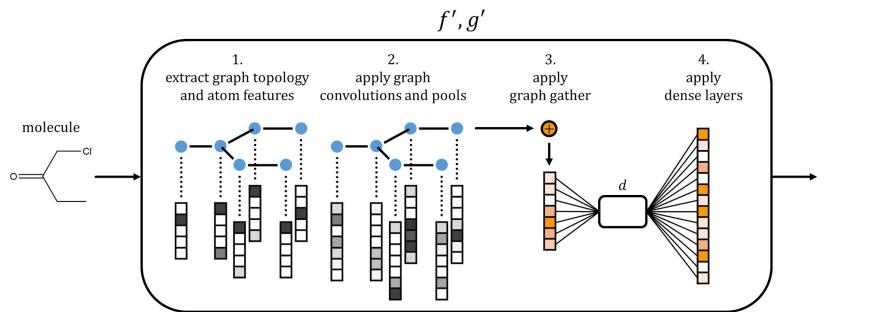


# Graph Gather

sum all nodes  
molecular featurization



# Graph Convolutions Explained Diagrammatically



# Low Data Drug Discovery with One Shot Learning

Han Altae-Tran<sup>\*‡</sup>, Bharath Ramsundar<sup>\*†</sup>, Aneesh S. Pappu<sup>†</sup>, Vijay Pandet<sup>†</sup>

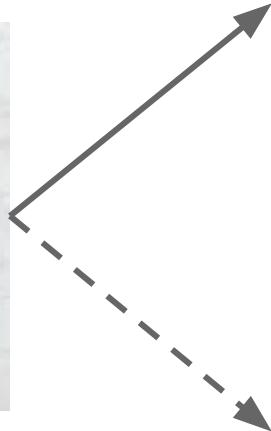
<sup>\*</sup> : equal contribution

<sup>†</sup>: Stanford University

<sup>‡</sup>: Massachusetts Institute of Technology

Published in ACS Central Science.

# One-Shot: Borrowing Tricks from Babies



Giraffe



# Google DeepMind's Matching Networks

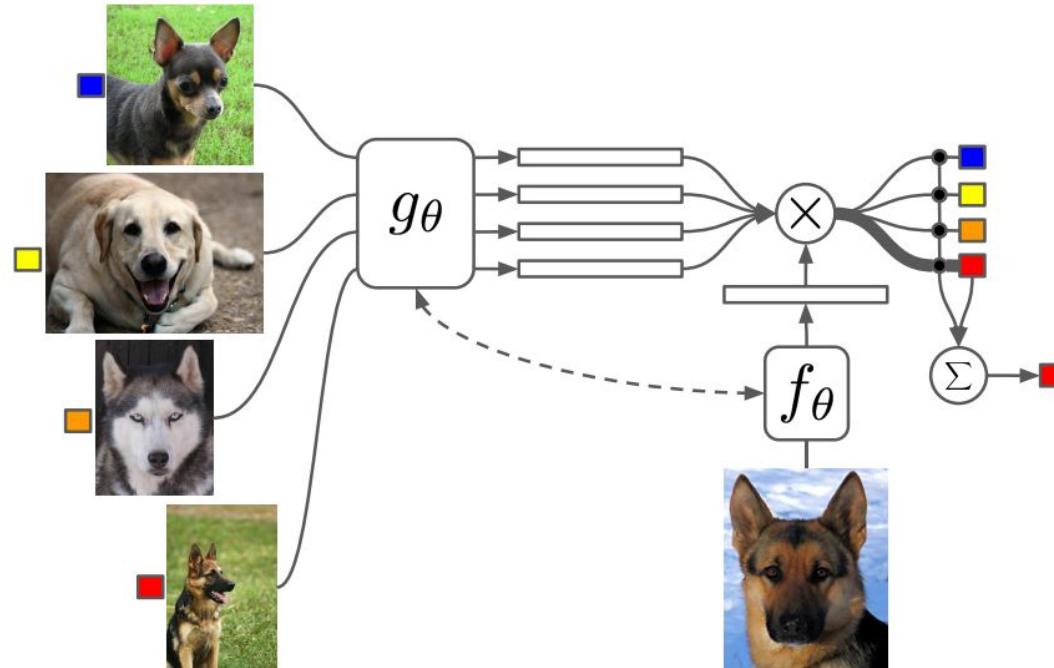


Figure 1: Matching Networks architecture

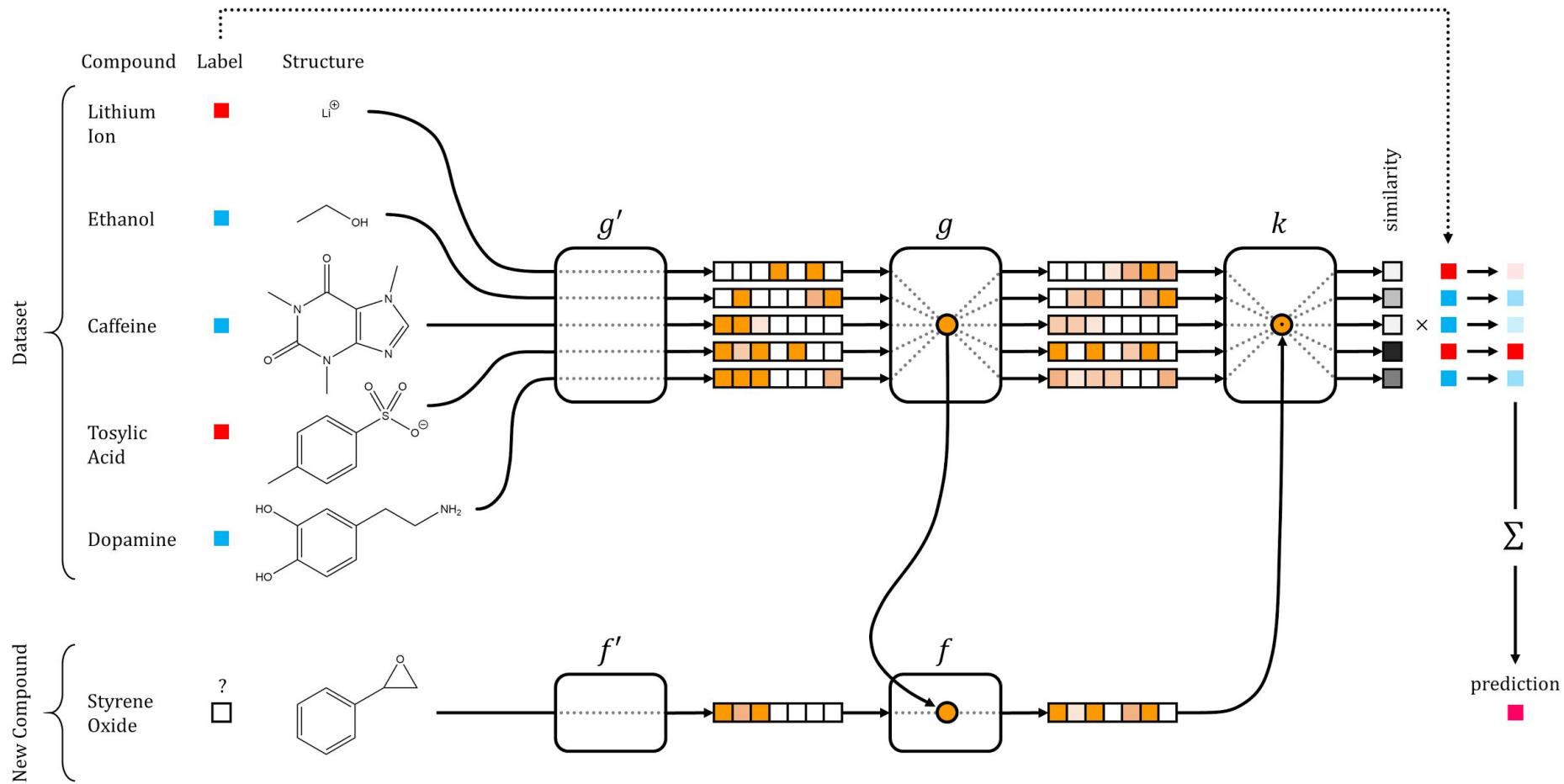


Table 1: ROC-AUC scores of models on median held-out task for each model on Tox21. Numbers reported are means and standard deviations. Randomness is over the choice of support set; experiment is repeated with 20 support sets. The appendix contains results for all held-out Tox21 tasks. The result with highest mean in each row is highlighted. The notation 10+/10- indicates supports with 10 positive examples and 10 negative examples.

| Tox21   | RF (100 trees)    | Graph Conv        | Siamese                             | AttnLSTM          | IterRefLSTM                         |
|---------|-------------------|-------------------|-------------------------------------|-------------------|-------------------------------------|
| 10+/10- | $0.586 \pm 0.056$ | $0.648 \pm 0.029$ | $0.820 \pm 0.003$                   | $0.801 \pm 0.001$ | <b><math>0.823 \pm 0.002</math></b> |
| 5+/10-  | $0.573 \pm 0.060$ | $0.637 \pm 0.061$ | $0.823 \pm 0.004$                   | $0.753 \pm 0.173$ | <b><math>0.830 \pm 0.001</math></b> |
| 1+/10-  | $0.551 \pm 0.067$ | $0.541 \pm 0.093$ | <b><math>0.726 \pm 0.173</math></b> | $0.549 \pm 0.088$ | $0.724 \pm 0.008$                   |
| 1+/5-   | $0.559 \pm 0.063$ | $0.595 \pm 0.086$ | $0.687 \pm 0.210$                   | $0.593 \pm 0.153$ | <b><math>0.795 \pm 0.005</math></b> |
| 1+/1-   | $0.535 \pm 0.056$ | $0.589 \pm 0.068$ | $0.657 \pm 0.222$                   | $0.507 \pm 0.079$ | <b><math>0.827 \pm 0.001</math></b> |

# MoleculeNet: <https://moleculenet.ai>

Zhenqin Wu<sup>\*†</sup>, Bharath Ramsundar<sup>\*†</sup>, Evan N. Feinberg<sup>\*\*†</sup>, Joseph Gomes<sup>\*\*†</sup>, Caleb Geniesse<sup>†</sup>, Evan N. Feinberg<sup>†</sup>, Aneesh S. Pappu<sup>†</sup>, Evan N. Feinberg<sup>†</sup>, Karl Leswing<sup>‡</sup>, Vijay Pandet<sup>†</sup>

<sup>\*</sup> : equal first authors

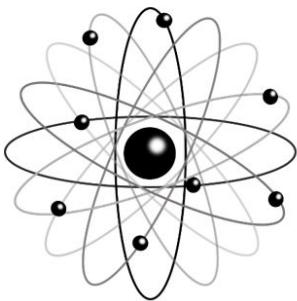
<sup>\*\*</sup>: equal second authors

<sup>†</sup>: Stanford University

<sup>‡</sup>: Schrodinger Inc.

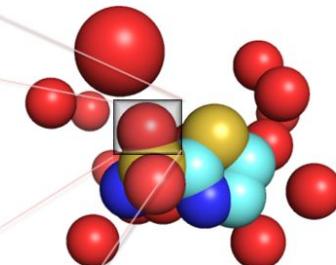
Accepted by Chemical Science.

# MoleculeNet Datasets



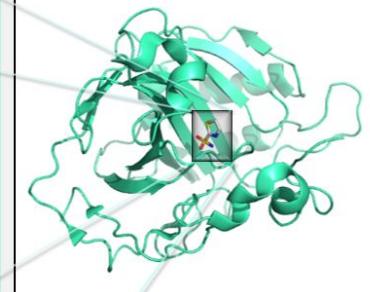
## Quantum Mechanics

- QM7
- QM8
- QM7b
- QM9



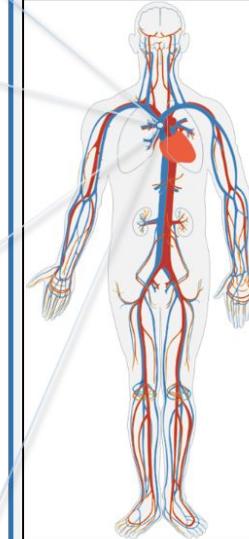
## Physical Chemistry

- ESOL
- FreeSolv
- Lipophilicity



## Biophysics

- HIV
- PDBbind
- BACE
- PCBA
- MUV

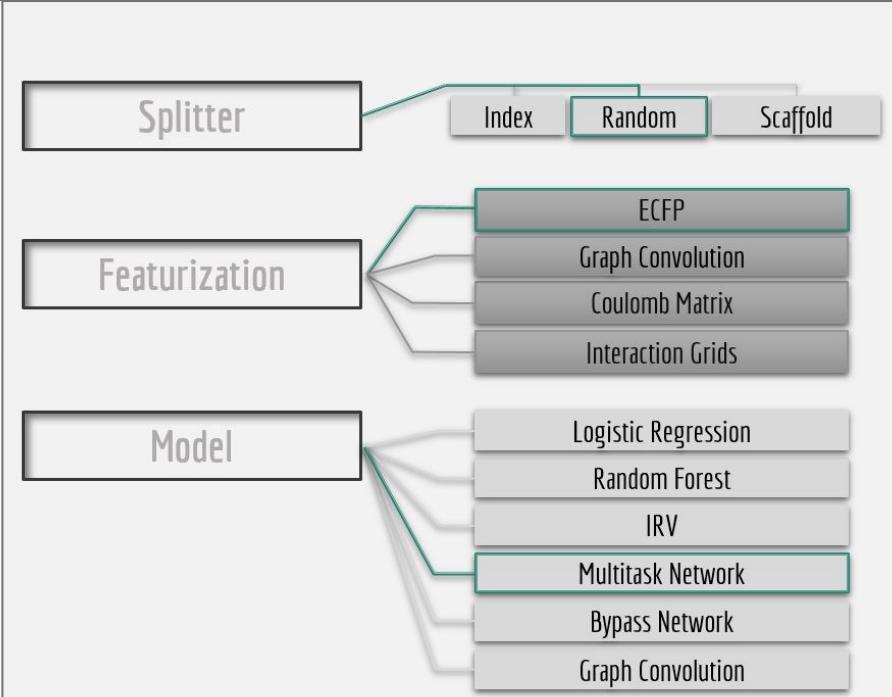


## Physiology

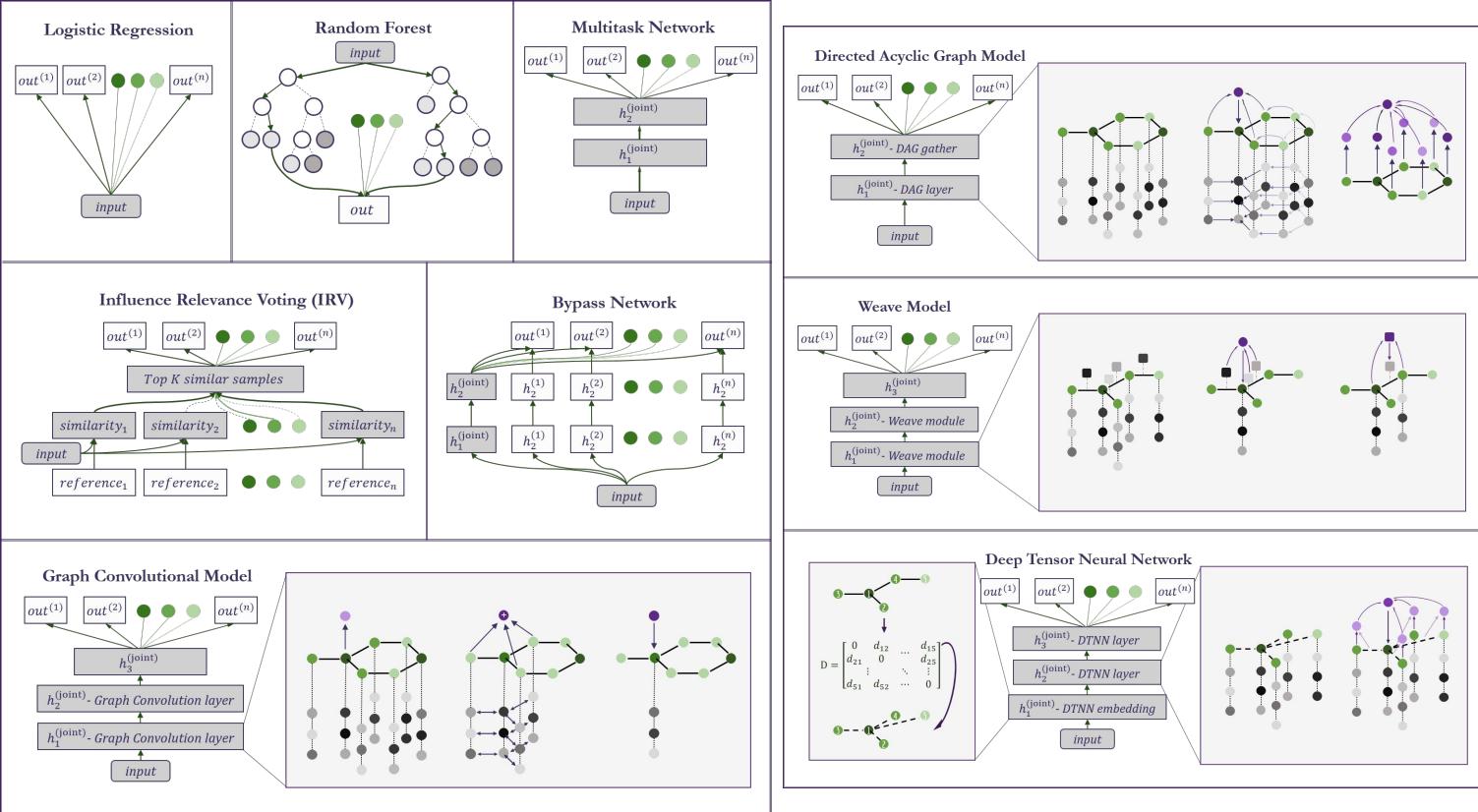
- BBBP
- Tox21
- ToxCast
- SIDER
- ClinTox

# MoleculeNet uses DeepChem to Benchmark

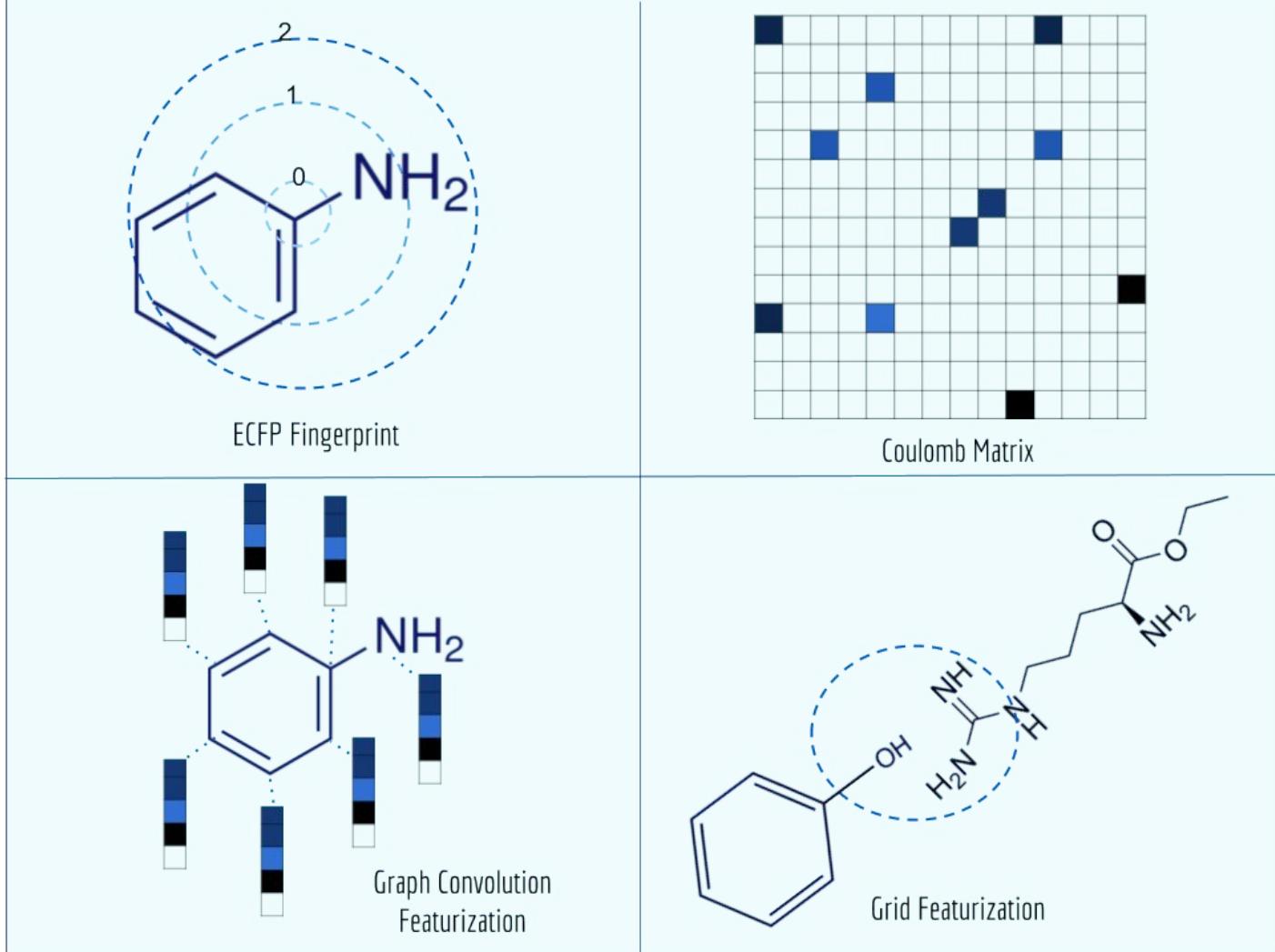
```
1 import numpy as np
2 import deepchem as dc
3 from clintox_datasets import load_clintox
4
5 # Load clintox dataset
6 n_features = 1024
7 clintox_tasks, datasets, transformers = load_clintox(split='random',
8 featurizer='ECFP')
9 train_dataset, valid_dataset, test_dataset = datasets
10 # Defining metrics
11 metric = dc.metrics.Metric(dc.metrics.roc_auc_score, np.mean)
12 # Defining models
13 model = dc.models.TensorflowMultiTaskClassifier(len(clintox_tasks), n_features,
14 layer_sizes=[1000], dropouts=[.25], learning_rate=0.001, batch_size=50)
15 # Fit trained model
16 model.fit(train_dataset)
17 model.save()
18 print("Evaluating model")
19 train_scores = model.evaluate(train_dataset, [metric], transformers)
20 valid_scores = model.evaluate(valid_dataset, [metric], transformers)
21 test_scores = model.evaluate(test_dataset, [metric], transformers)
22 print("Train scores")
23 print(train_scores)
24 print("Validation scores")
25 print(valid_scores)
26 print("Test scores")
27 print(test_scores)
```



# Machine Learning Models Tested for MoleculeNet Benchmarks



# MoleculeNet Featurizations



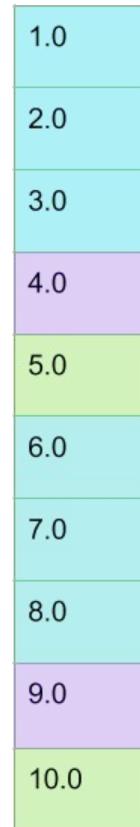
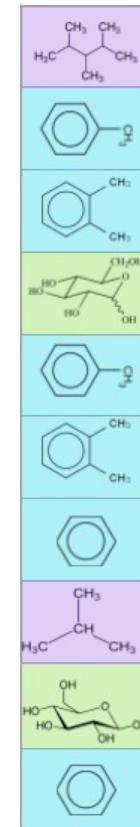
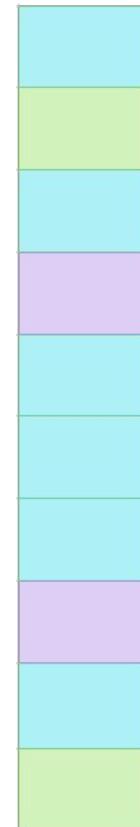
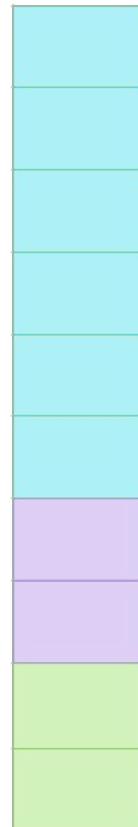
# MoleculeNet

## Tests

### Generalization with Different Splits

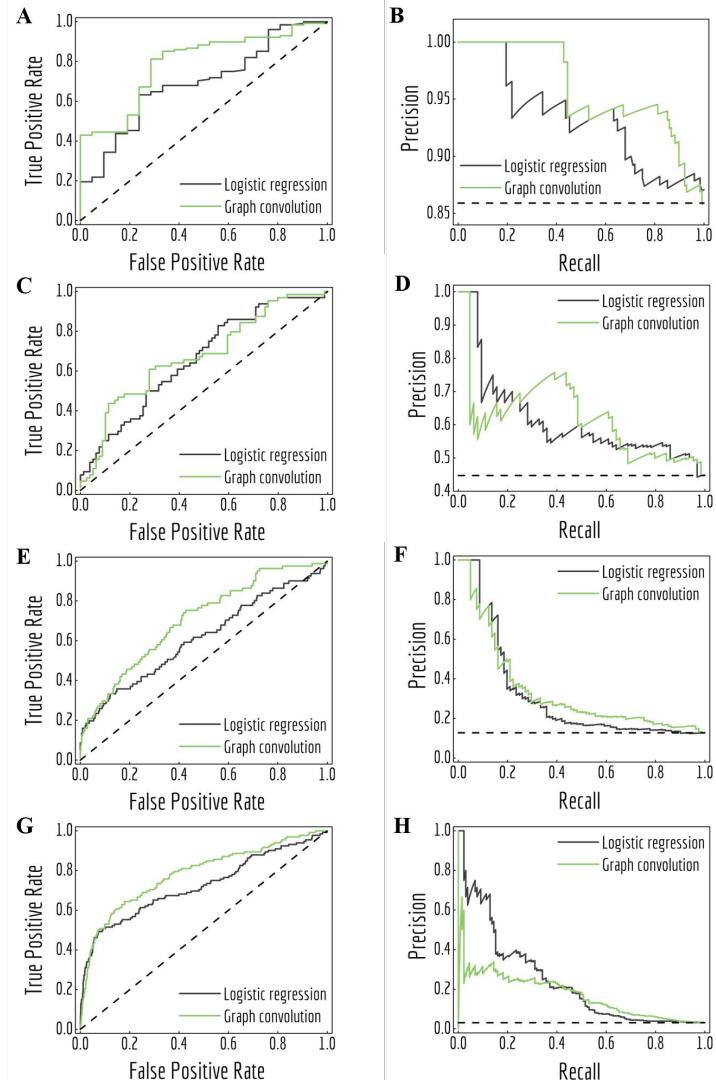
Legend:

- Train
- Valid
- Test



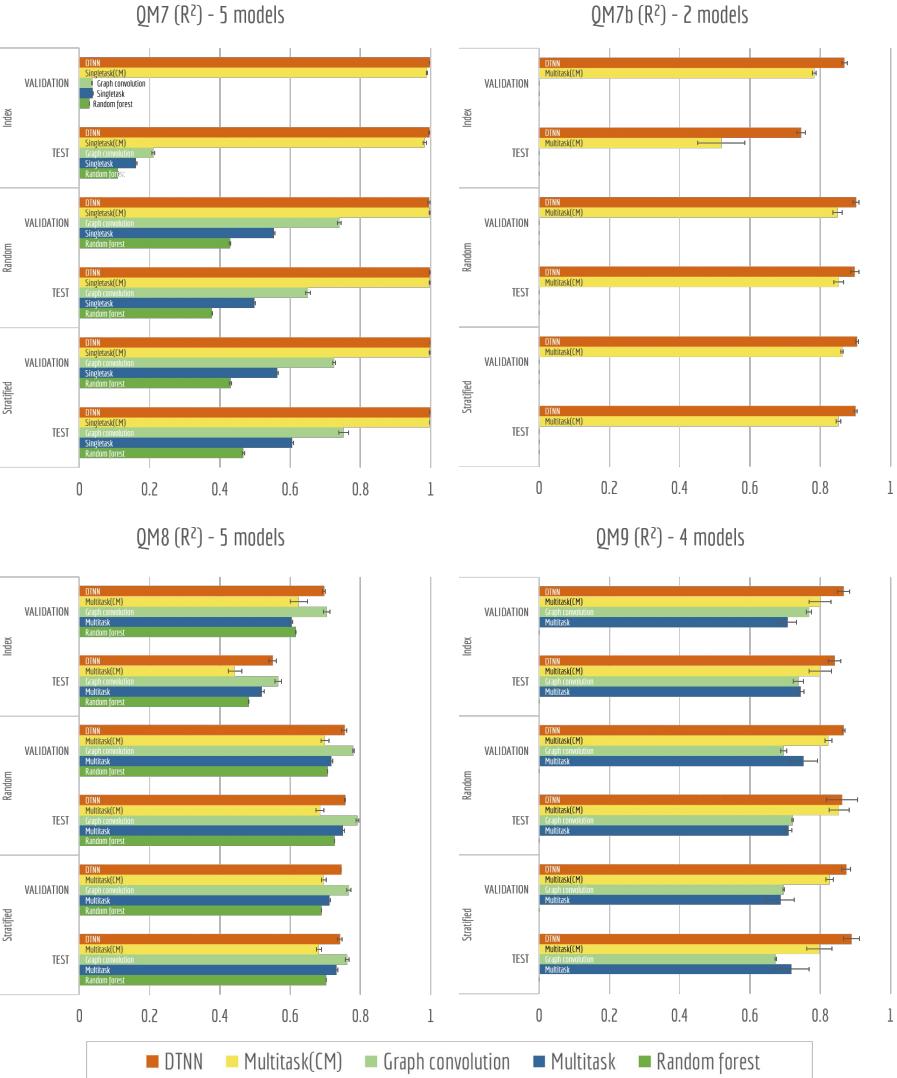
# ROC vs PRC

- Most deep drug discovery benchmarked with ROC-AUC
- However, analysis shows that for very imbalanced datasets, ROC and PRC behave differently.
- Deep methods do not currently win on very imbalanced datasets!



# Quantum Datasets

- Deep Learning methods dominate benchmark.
- End-to-end DTNN wins overall, with hand engineered DNN+Coulomb a close second.



# PubChem BioAssay

PCBA (AUC-ROC) - 4 models



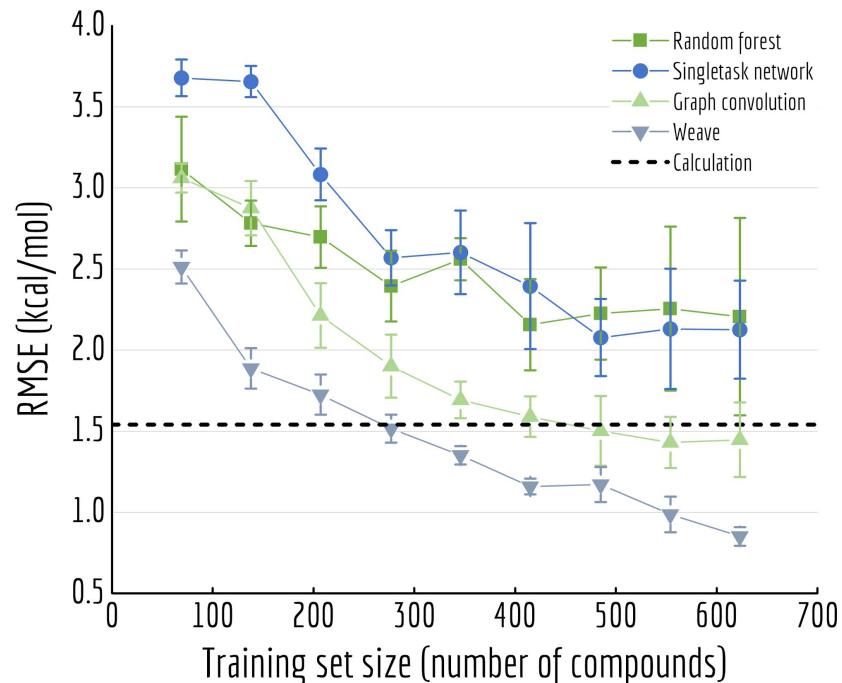
PCBA (AUC-PRC) - 4 models



■ Graph convolution ■ Bypass ■ Multitask ■ Logistic regression

# Molecular ML vs. Physical Computations

- Solubility Data
- Dashed Line measures error of physics-based simulation algorithm.
- 1000 datapoints in set.  
Plots show that Molecular ML can out-perform physics-based with enough data.



# Atomic Convolutional Networks for Predicting Protein-Ligand Binding Affinity

Joseph Gomes\*†, Bharath Ramsundar\*, Evan N. Feinberg†, Vijay Pandet

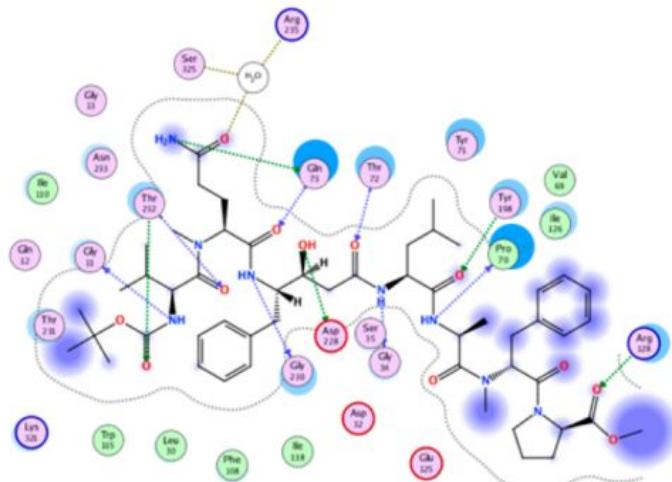
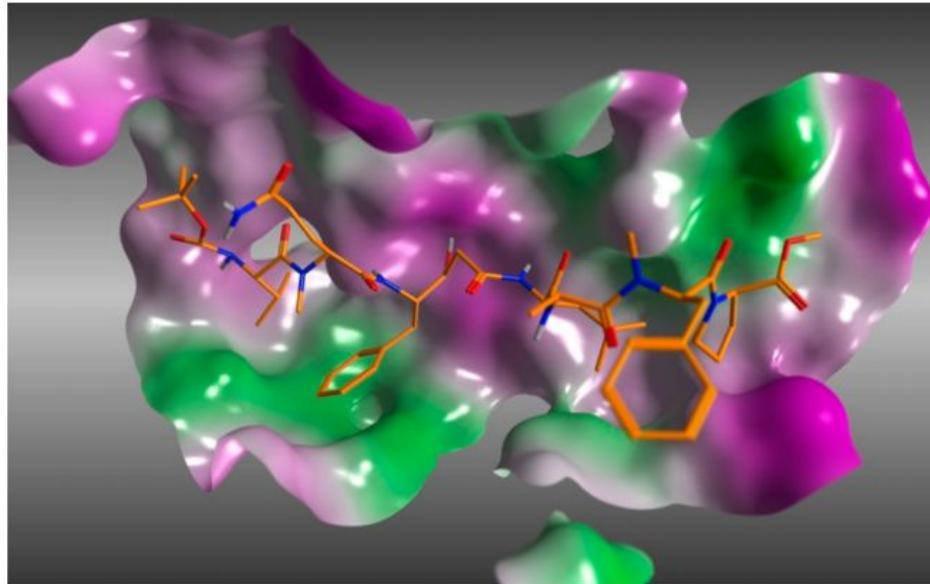
\* : equal contribution

†: Stanford University

Currently on Arxiv. Tentatively plan to submit to ACS Central Science.

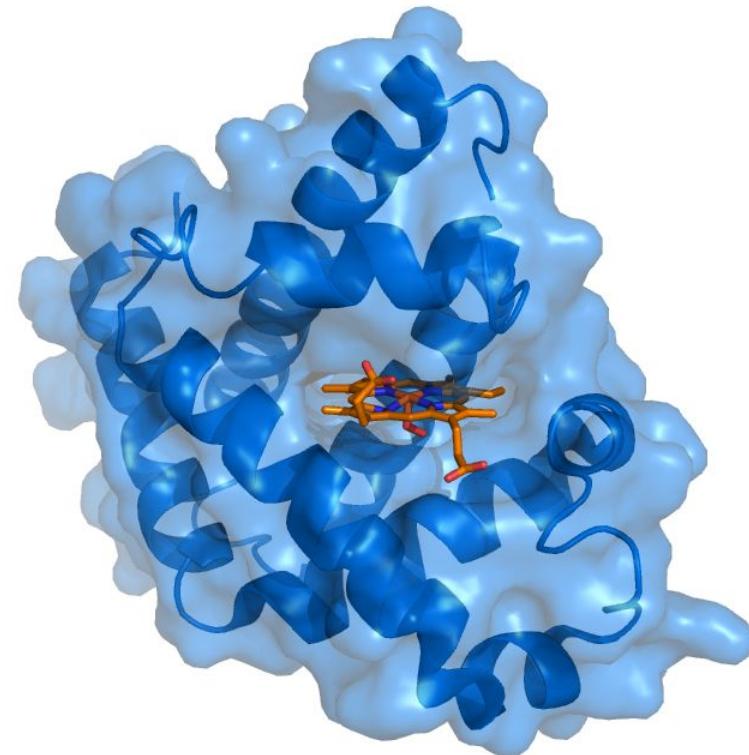
# Example of BACE-1 Binding Pocket Interactions

Scheme 1. Depiction of BACE-1 Binding Site (left) Using the Ligand from PDB Code 3UQP along with the Protein–Ligand Interaction (right)

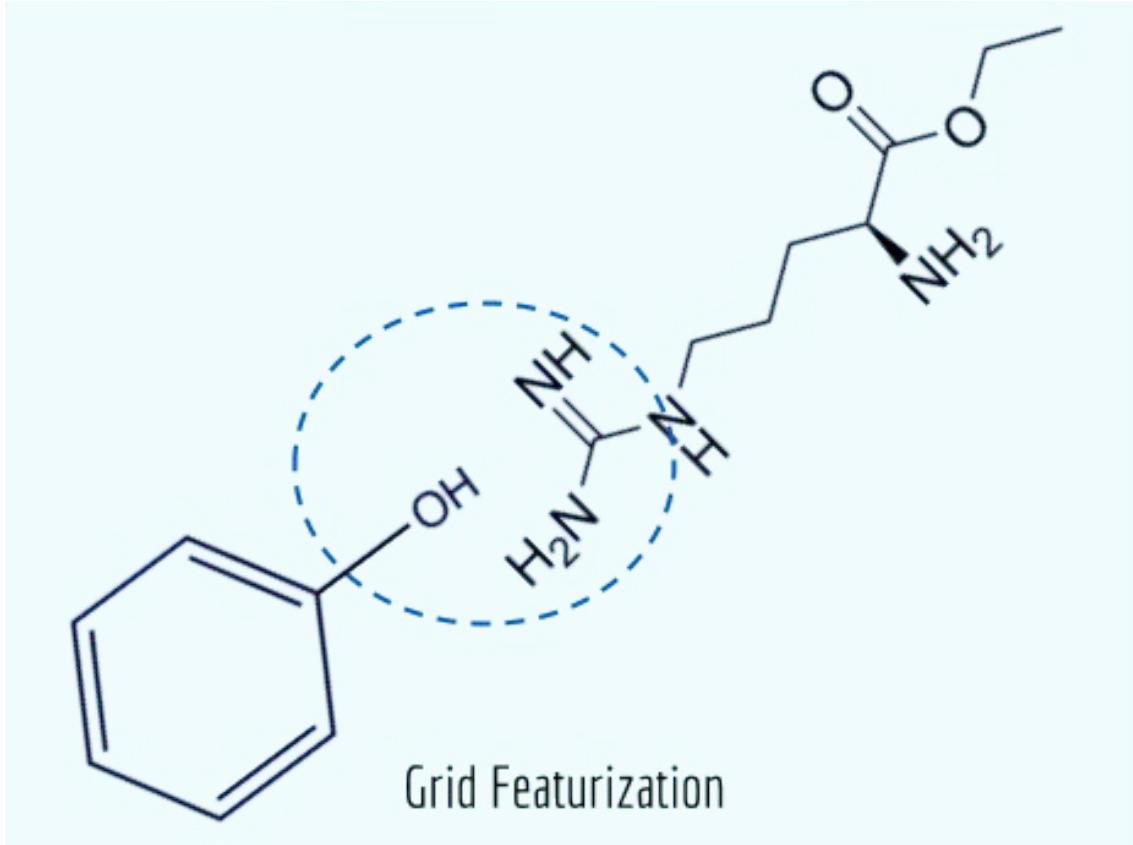


# Exploiting 3D Geometry for Drug Discovery

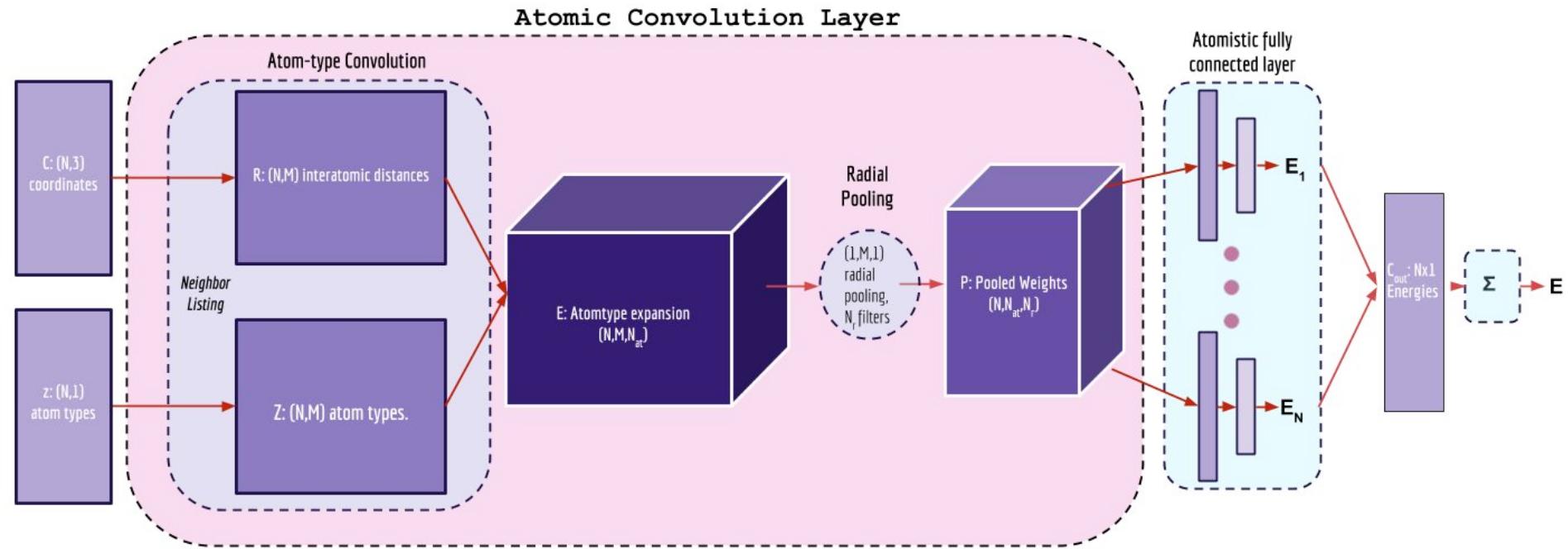
- Learning to predict interaction energies requires modeling the 3D geometry of binding interactions.
- Two approaches: hand tuned and learned featurizations.



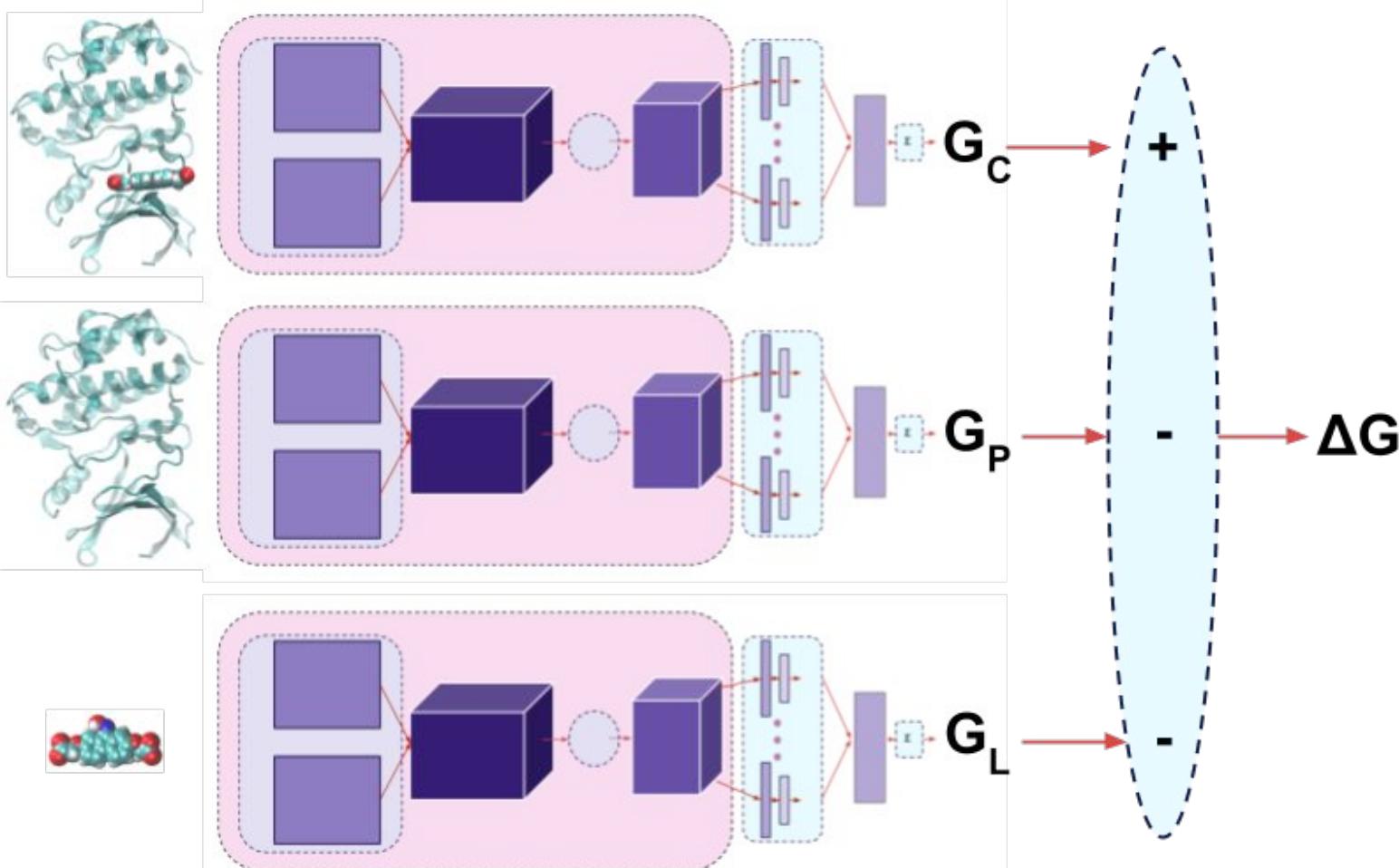
# Hand-Designed Grid Featurization



# Learnable Atomic Convolutions



# Thermodynamic Cycles with End-to-End Training



# Results on PDBBind Benchmark

|            | ACNN  |             | GRID-RF |             | GRID-NN |      | GCNN  |      | ECFP-RF |      | ECFP-NN |      |
|------------|-------|-------------|---------|-------------|---------|------|-------|------|---------|------|---------|------|
| Split      | Train | Test        | Train   | Test        | Train   | Test | Train | Test | Train   | Test | Train   | Test |
| Random     | .705  | .508        | .962    | <b>.546</b> | .976    | .539 | .581  | .403 | .883    | .466 | .850    | .386 |
| Stratified | .793  | .491        | .963    | <b>.562</b> | .983    | .494 | .635  | .410 | .883    | .477 | .840    | .462 |
| Scaffold   | .752  | .267        | .964    | <b>.349</b> | .977    | .334 | .652  | .135 | .861    | .229 | .824    | .146 |
| Temporal   | .704  | <b>.529</b> | .963    | .486        | .978    | .455 | .596  | .315 | .875    | .401 | .827    | .362 |

**Table 3.** Performance (Pearson  $R^2$ ) on PDBBind refined train/test sets.

|            | ACNN  |       | GRID-RF |              | GRID-NN |       | GCNN  |       | ECFP-RF |       | ECFP-NN |       |
|------------|-------|-------|---------|--------------|---------|-------|-------|-------|---------|-------|---------|-------|
| Split      | Train | Test  | Train   | Test         | Train   | Test  | Train | Test  | Train   | Test  | Train   | Test  |
| Random     | 0.518 | 0.653 | 0.240   | <b>0.630</b> | 0.468   | 1.208 | 0.598 | 0.726 | 0.319   | 0.661 | 0.335   | 0.740 |
| Stratified | 0.540 | 0.686 | 0.237   | <b>0.650</b> | 0.439   | 1.109 | 0.559 | 0.710 | 0.319   | 0.660 | 0.344   | 0.678 |
| Scaffold   | 0.555 | 0.579 | 0.245   | <b>0.537</b> | 0.406   | 1.177 | 0.557 | 0.805 | 0.339   | 0.749 | 0.373   | 0.856 |
| Temporal   | 0.526 | 0.668 | 0.235   | <b>0.666</b> | 0.410   | 1.184 | 0.610 | 0.865 | 0.321   | 0.770 | 0.371   | 0.798 |

**Table 4.** Performance (MUE [kcal/mol]) on PDBBind refined train/test sets.