# TEACHING FREE ENERGY CALCULATIONS TO LEARN

**John D. Chodera**
MSKCC Computational and Systems Biology Program
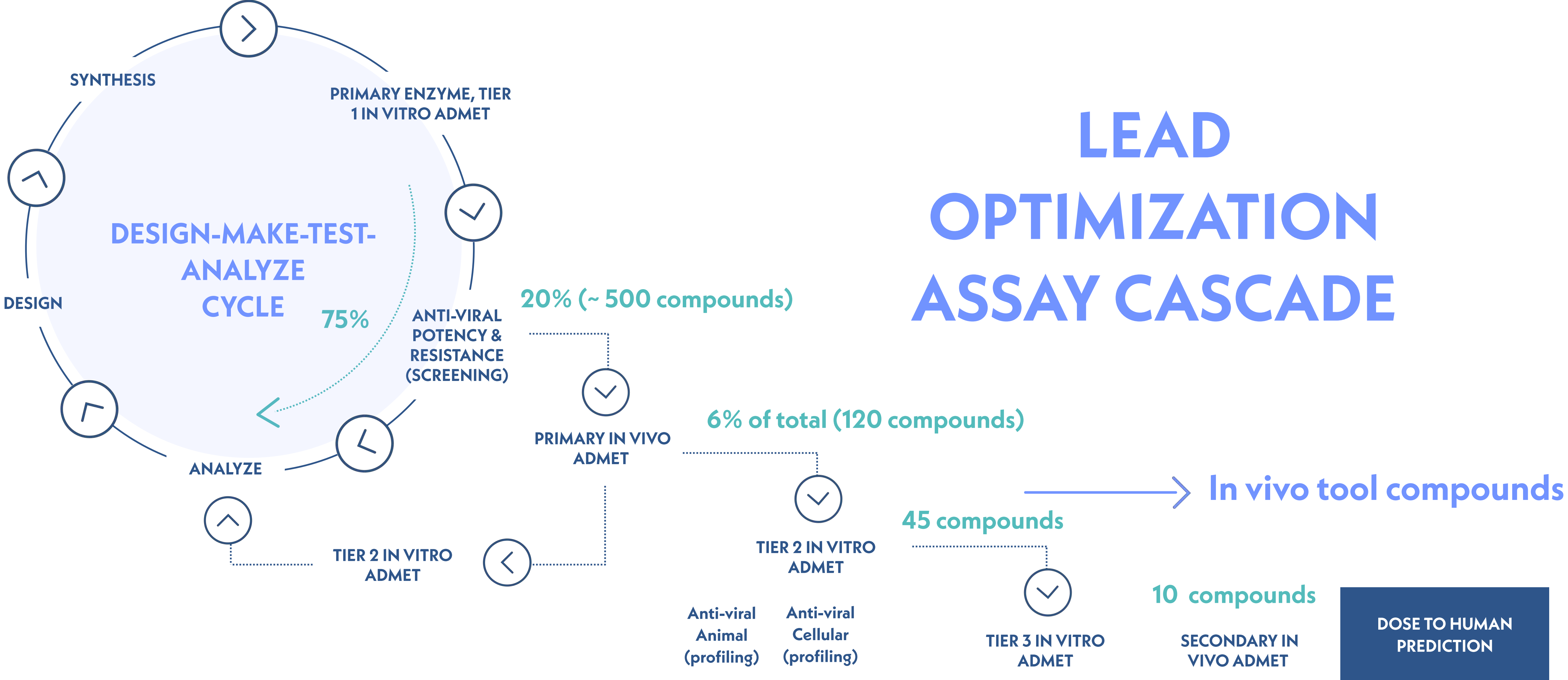Sildes will be posted to http://www.choderalab.org/news

29 Apr 2022 - ICML MLDD - Cyberspace

# MODELS TO STEER DESIGN-MAKE-TEST-ANALYZE CYCLES CAN DIRECTLY IMPACT DISCOVERY PROGRAMS

# STRUCTURAL DATA IS NOW AN ABUNDANT RESOURCE FOR DRUG DISCOVERY



PDB statistics

$16B

AlphaFold2-like methods can generate structural models for many more targets

100,000 new structures
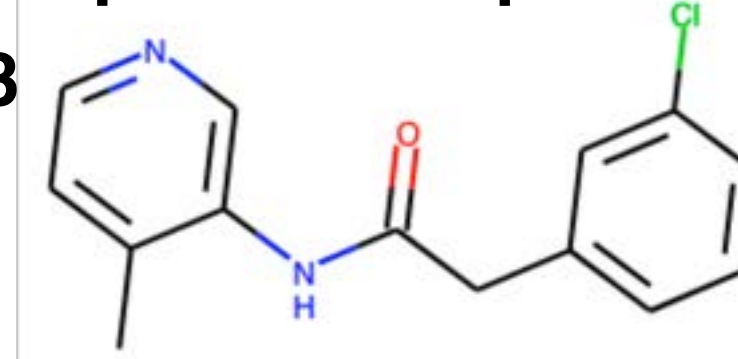
last decade

# WE COMMONLY NEED TO MAKE DECISIONS BETWEEN MANY RELATED SYNTHETICALLY FEASIBLE ANALOGUES

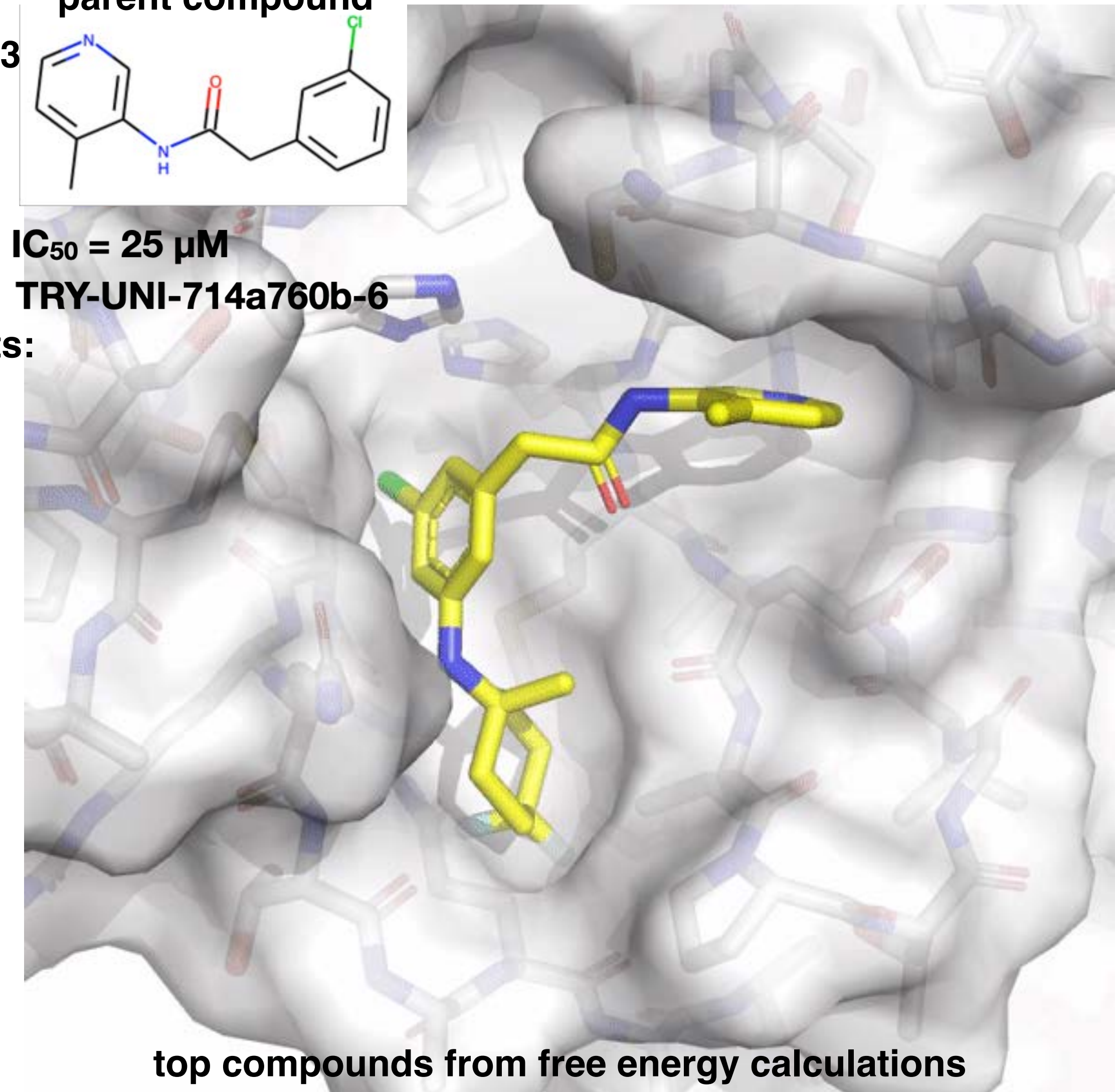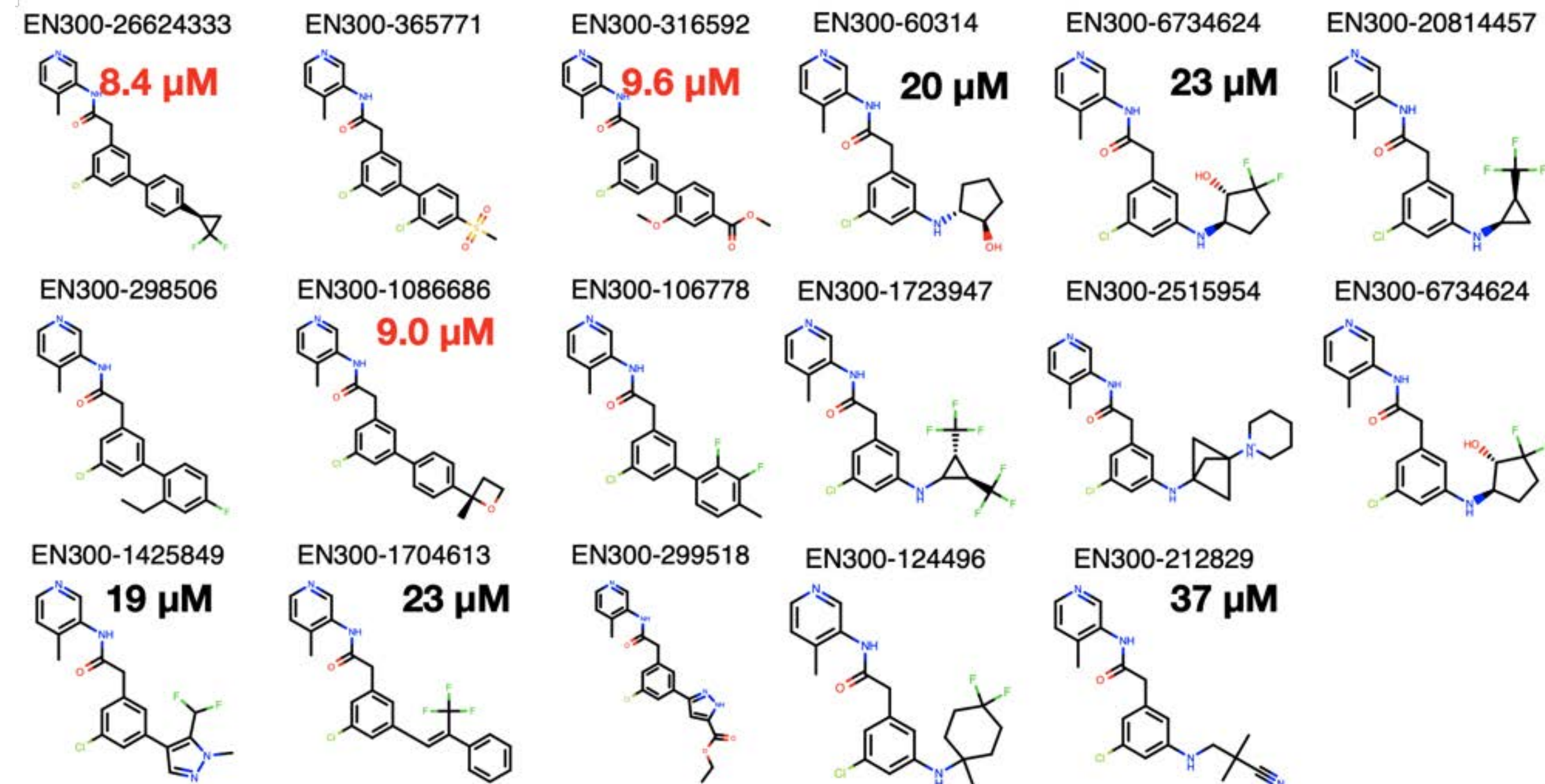Can we engage S4 from this 5,000-compound virtual synthetic library varying R3
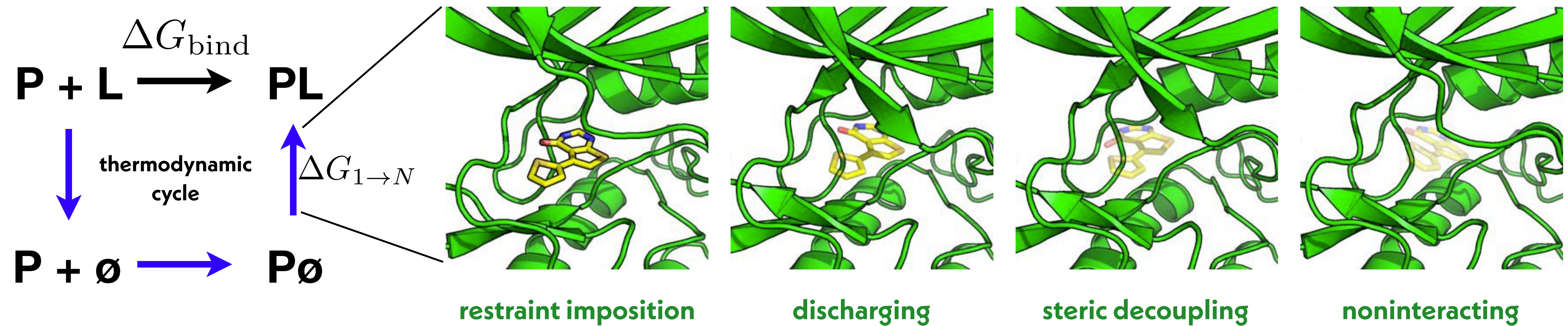


parent compound

IC$_{50}$ = 25 μM
TRY-UNI-714a760b-6

Top free energy calculation compounds and experimental affinity measurements:



EN300-26624333 **8.4 μM**   EN300-365771   EN300-316592 **9.6 μM**   EN300-60314 **20 μM**   EN300-6734624 **23 μM**   EN300-20814457

EN300-298506   EN300-1086686 **9.0 μM**   EN300-106778   EN300-1723947   EN300-2515954   EN300-6734624

EN300-1425849 **19 μM**   EN300-1704613 **23 μM**   EN300-299518   EN300-124496   EN300-212829 **37 μM**

top compounds from free energy calculations

**COVID Moonshot:** [Moonshot] [Fragalysis] [Dashboard]

# ALCHEMICAL FREE ENERGY CALCULATIONS HAVE PROVEN TO BE A USEFUL WAY TO EXPLOIT STRUCTURAL DATA TO PREDICT AFFINITIES

simulations of alchemical intermediates with attenuated interactions

$$\Delta G_{\text{bind}}$$

$$P + L \longrightarrow PL$$

thermodynamic cycle

$$\Delta G_{1 \to N}$$

$$P + \emptyset \longrightarrow P\emptyset$$



restraint imposition      discharging      steric decoupling      noninteracting

## Includes all contributions from enthalpy and entropy of binding to a flexible receptor

$$\Delta G_{1 \to N} = -\beta^{-1} \ln \frac{Z_N}{Z_1} = -\beta^{-1} \ln \frac{Z_2}{Z_1} \cdot \frac{Z_3}{Z_2} \cdots \frac{Z_N}{Z_{N-1}}$$

$$Z_n = \int dx\, e^{-\beta U_n(x)} \quad \text{partition function}$$

# CURRENT ACCURACIES ARE SUFFICIENT TO ACCELERATE DISCOVERY, BUT HOW CAN WE GO FURTHER?
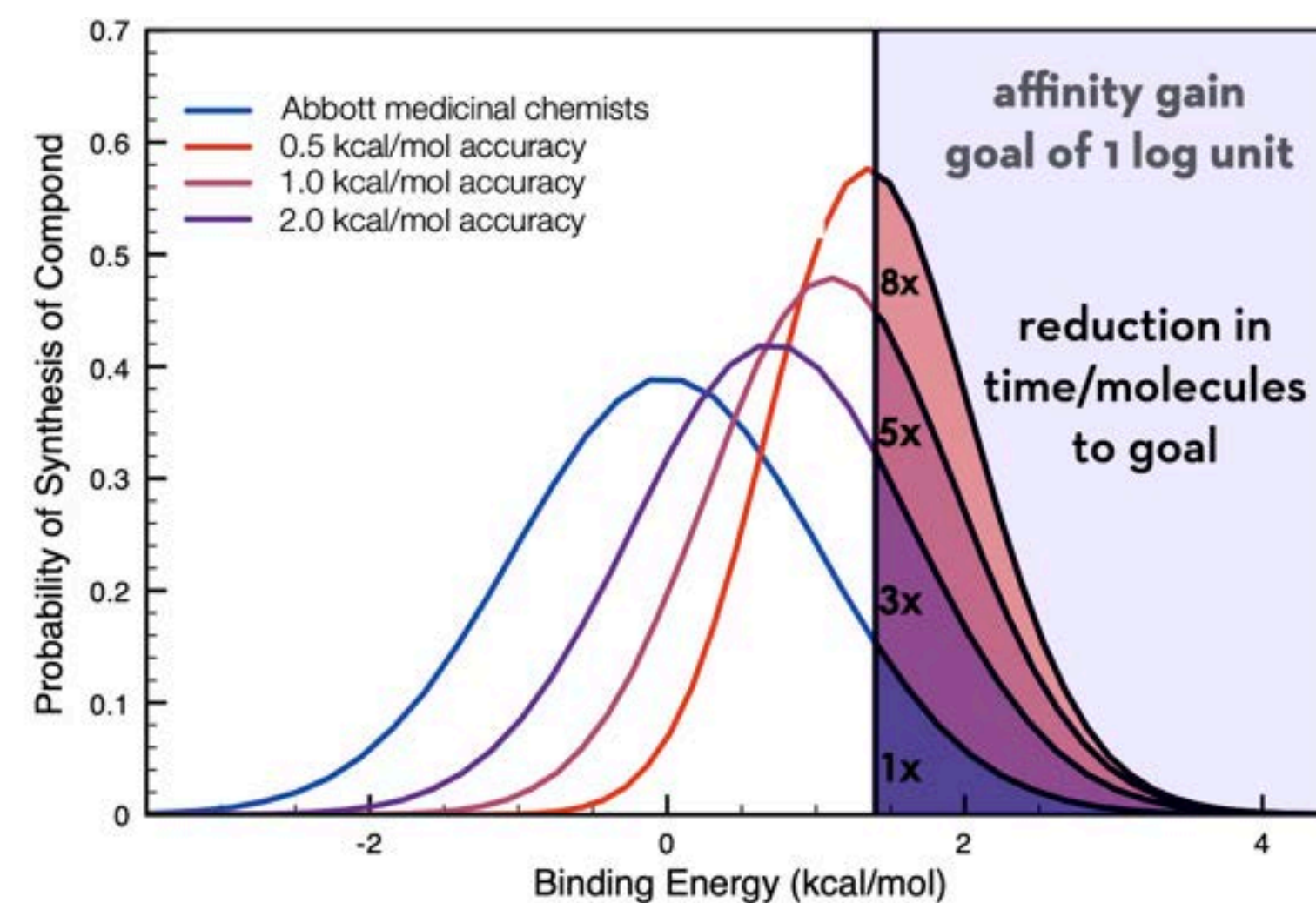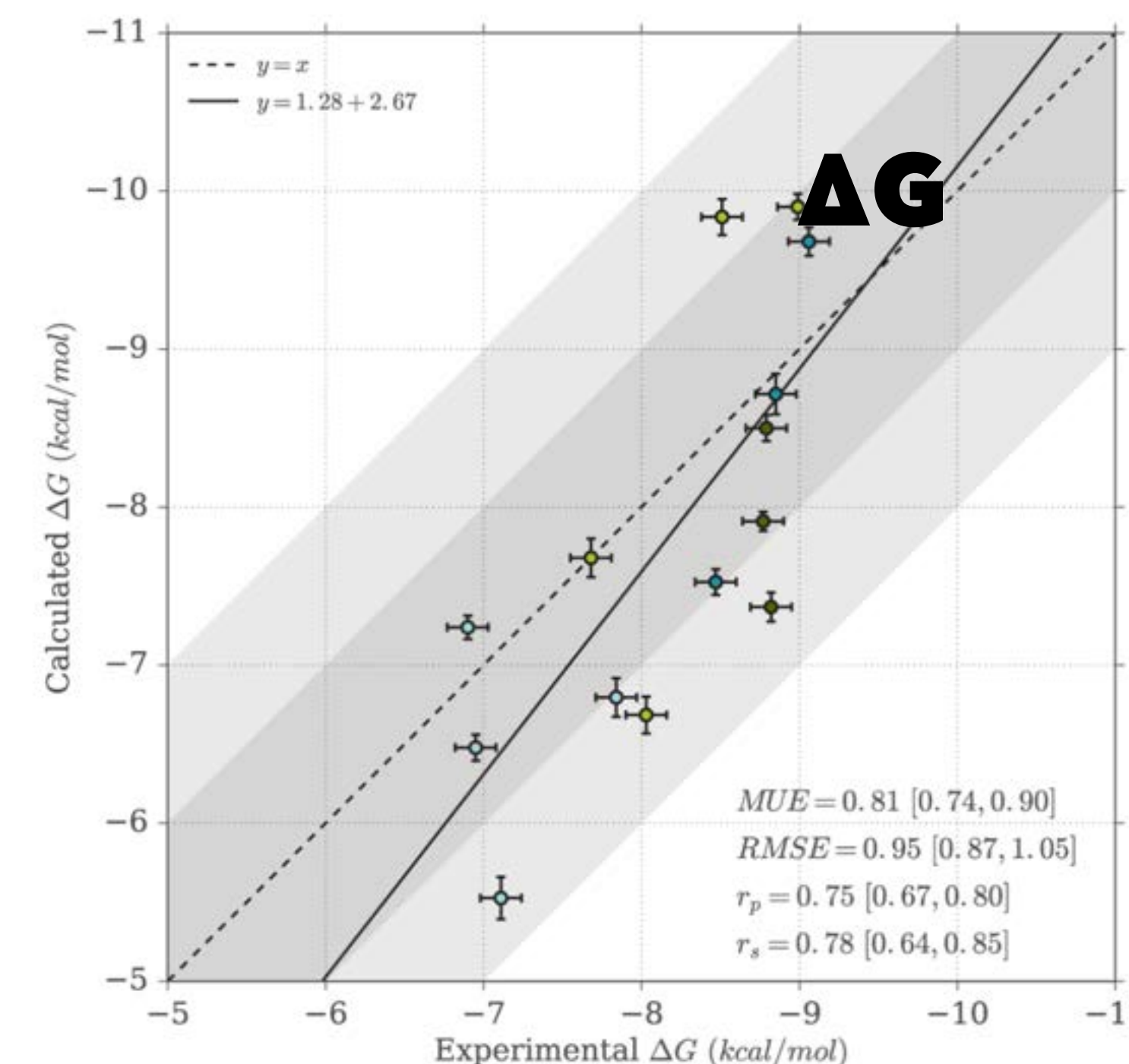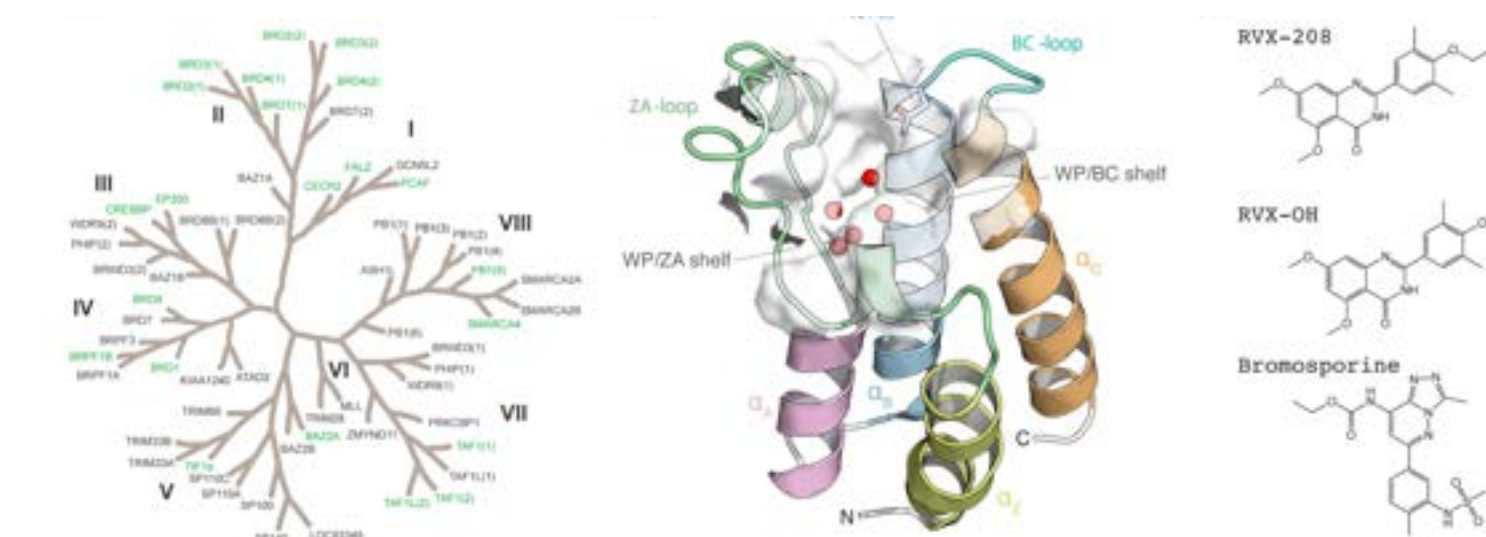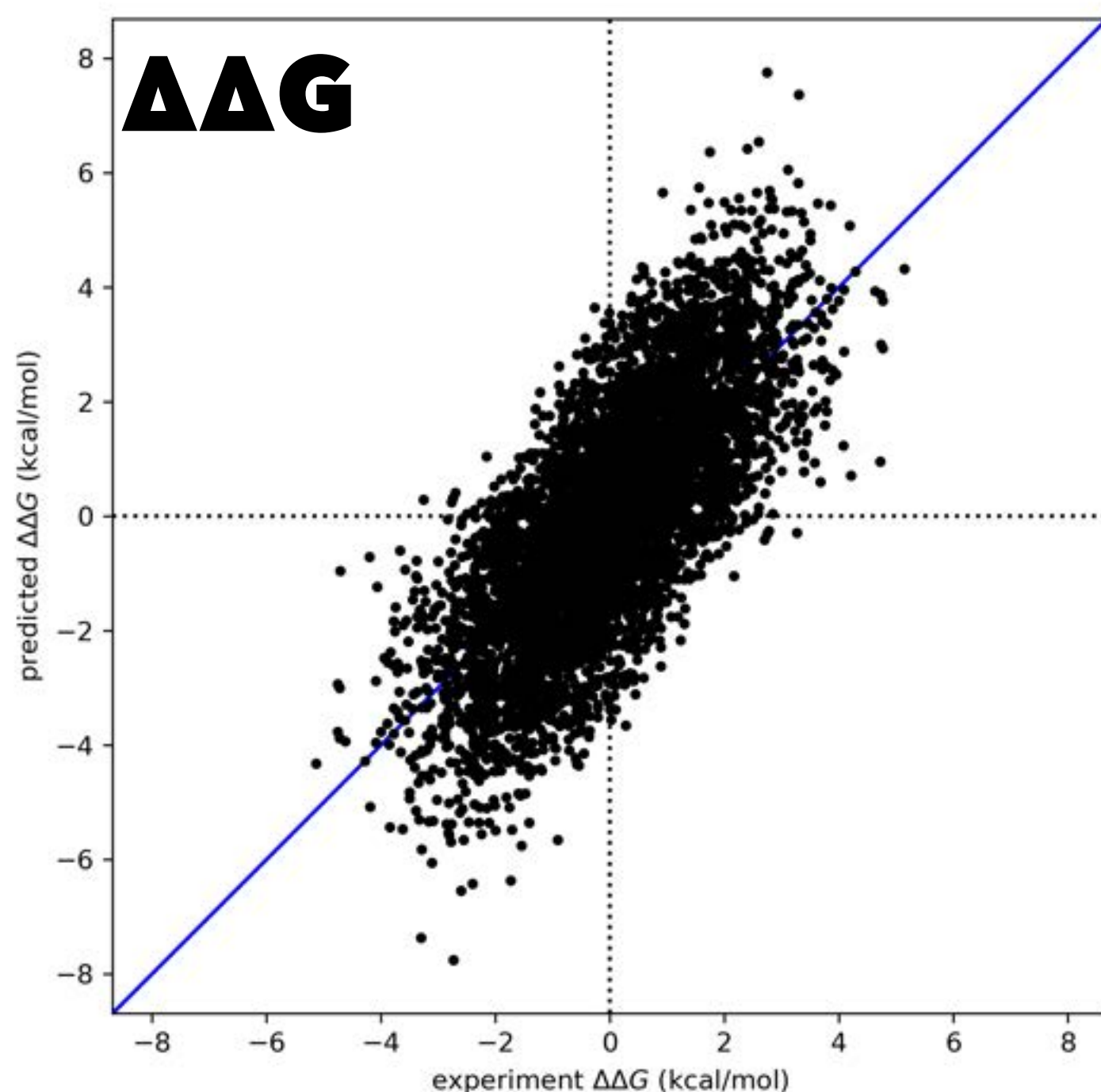
## RELATIVE

## ABSOLUTE



all within-target pairs ΔΔG (N = 5620)

```
RMSE: OPLS    1.37 [95%:  1.34,  1.39] kcal/mol
MUE : OPLS    1.09 [95%:  1.07,  1.11] kcal/mol
R2  : OPLS    0.10 [95%:  0.06,  0.15] kcal/mol
rho : OPLS    0.73 [95%:  0.72,  0.74] kcal/mol
```
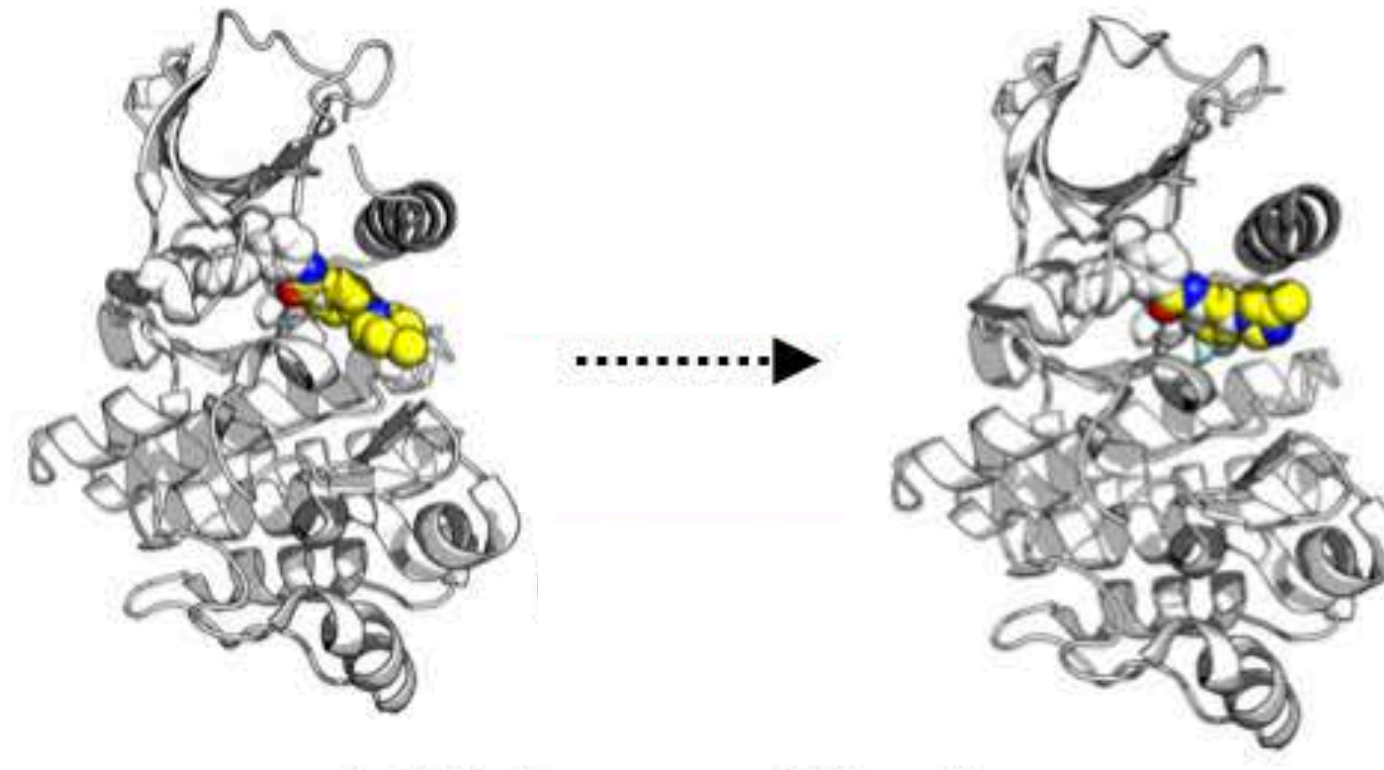
ΔΔG RMSE ~ 1.4 kcal/mol for well-behaved* proteins/chemistries:
**3-5x reduction in molecules synthesized**





**ΔΔG**



**ΔG**

$MUE = 0.81 [0.74, 0.90]$
$RMSE = 0.95 [0.87, 1.05]$
$r_p = 0.75 [0.67, 0.80]$
$r_s = 0.78 [0.64, 0.85]$

**\*best-case scenarios!**

Wang et al. (Schrödinger) JACS 137:2695, 2015
https://doi.org/10.1021/ja512751q
Reanalysis: http://github.com/jchodera/jacs-dataset-analysis

Aldeghi et al. JACS 139:946, 2017.
https://doi.org/10.1021/jacs.6b11467

# ALCHEMICAL FREE ENERGY CALCULATIONS HAVE A BROAD DOMAIN OF APPLICABILITY

## driving affinity / potency
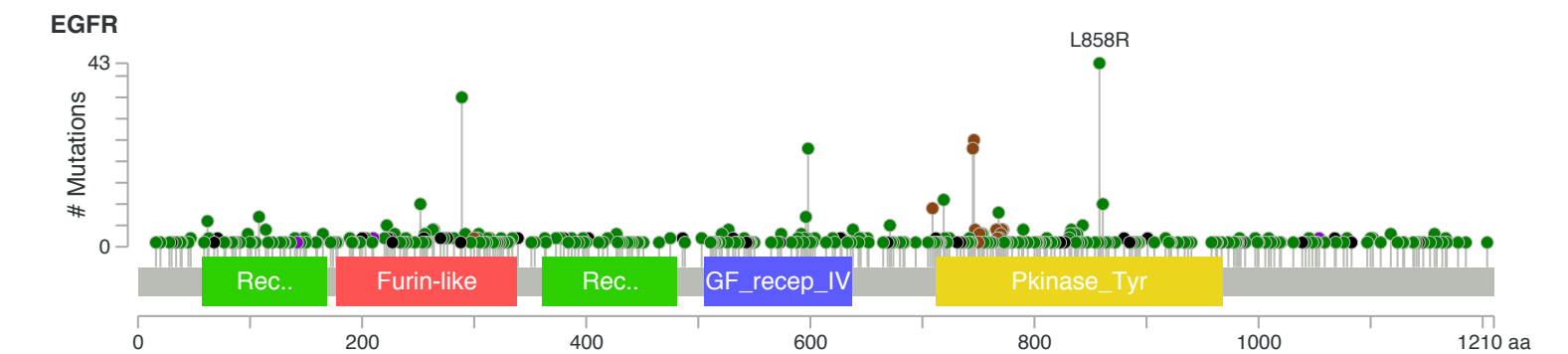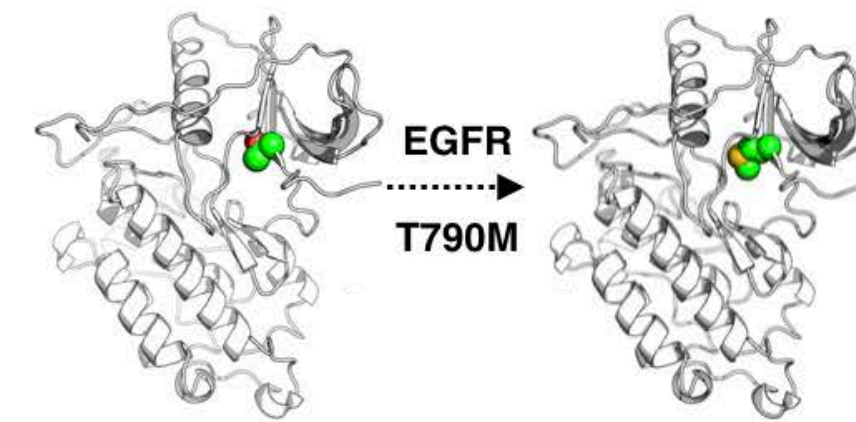Schindler, Baumann, Blum et al. JCIM 11:5457, 2020
https://doi.org/10.1021/acs.jcim.0c00900

## driving selectivity
Moraca, Negri, de Olivera, Abel JCIM 2019
https://doi.org/10.1021/acs.jcim.9b00106
Aldeghi et al. JACS 139:946, 2017.
https://doi.org/10.1021/jacs.6b11467

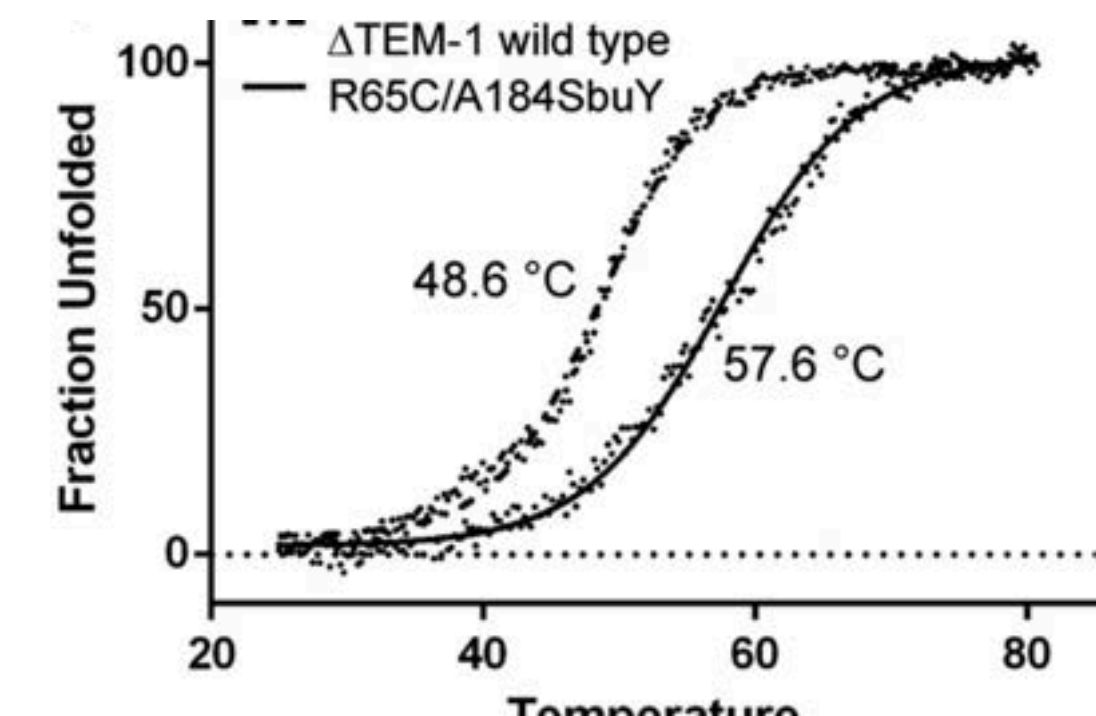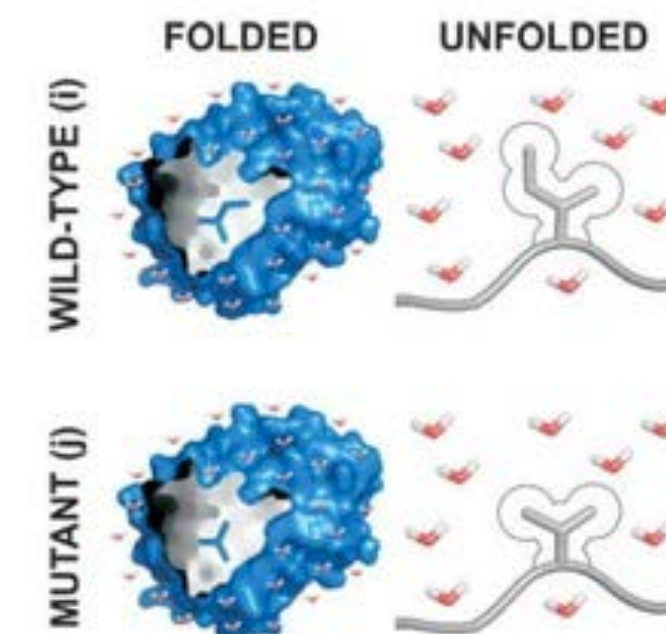## predicting clinical drug resistance/sensitivity
Hauser, Negron, Albanese, Ray, Steinbrecher, Abel, Chodera, Wang.
Communications Biology 1:70, 2018
https://doi.org/10.1038/s42003-018-0075-x
Aldeghi, Gapsys, de Groot. ACS Central Science 4:1708, 2018
https://doi.org/10.1021/acscentsci.8b00717

## optimizing thermostability
Gapsys, Michielssens, Seeliger, and de Groot. Angew Chem 55:7364, 2016
https://doi.org/10.1002/anie.201510054
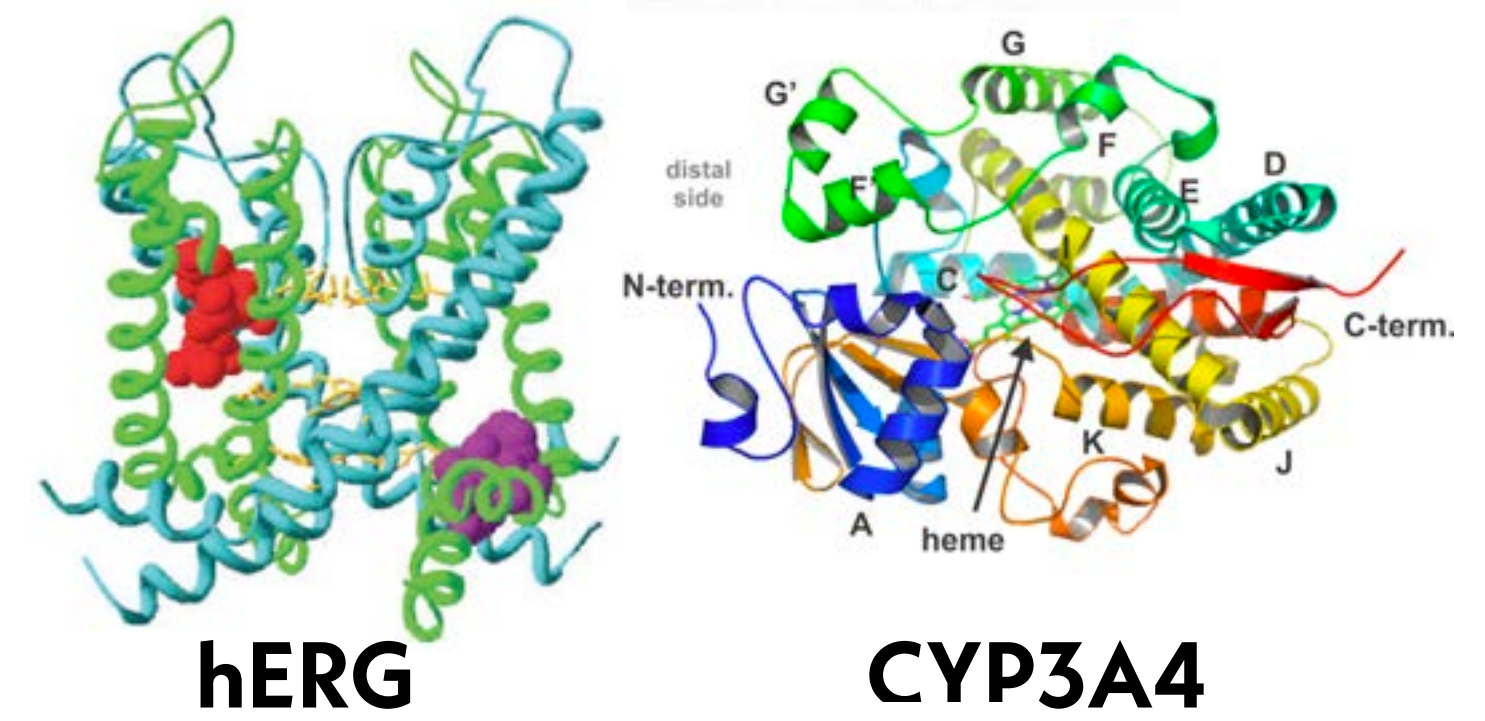


Erlotinib   Lapatinib

# ...AND HOLD THE POTENTIAL FOR EVEN BROADER APPLICABILITY AS MORE STRUCTURAL DATA EMERGES

partition coefficients (logP, logD) and permeabilities

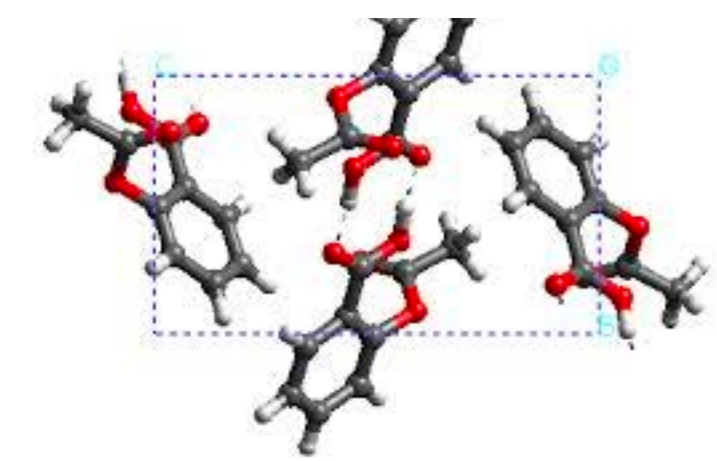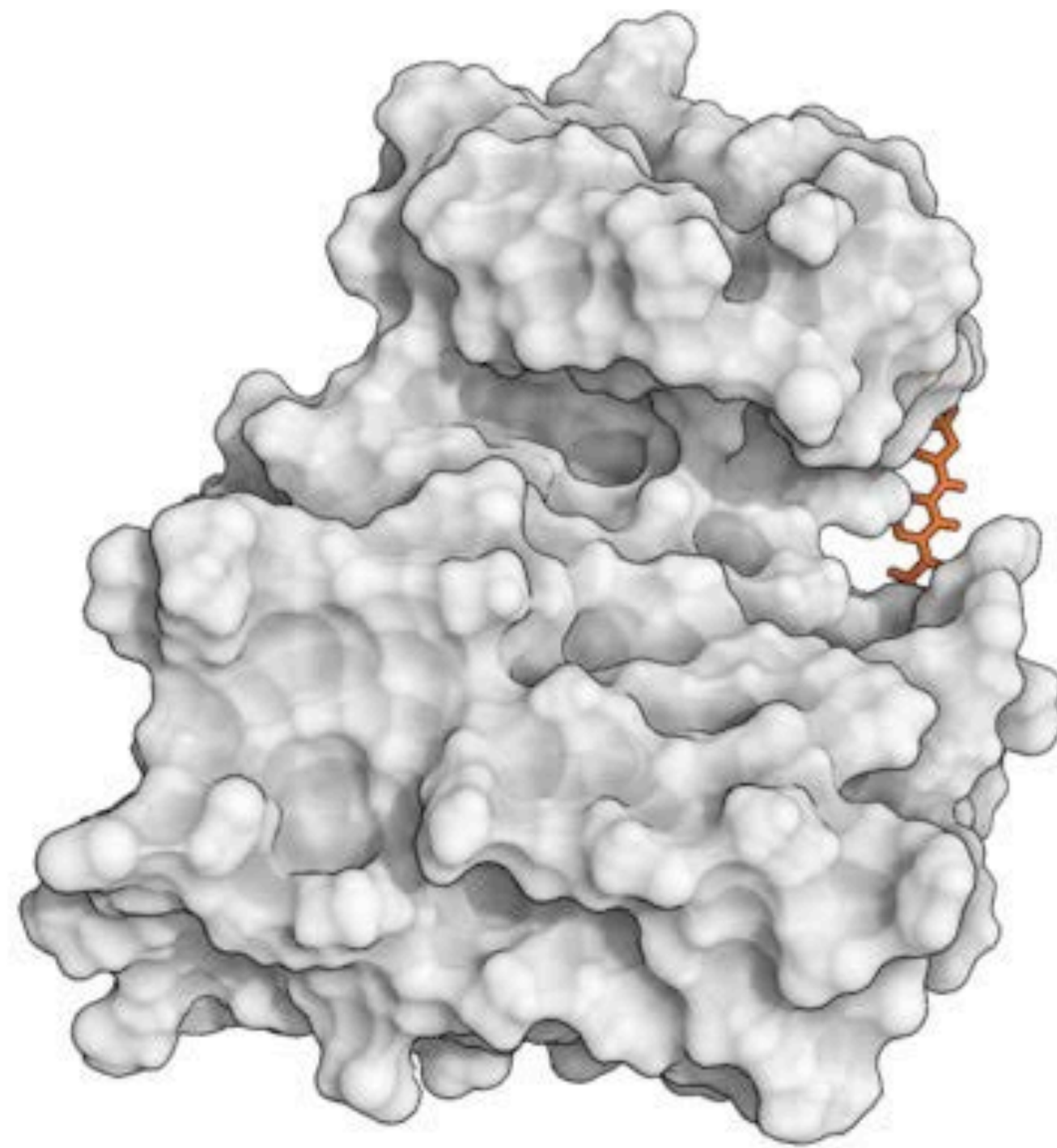structure-enabled ADME/Tox targets
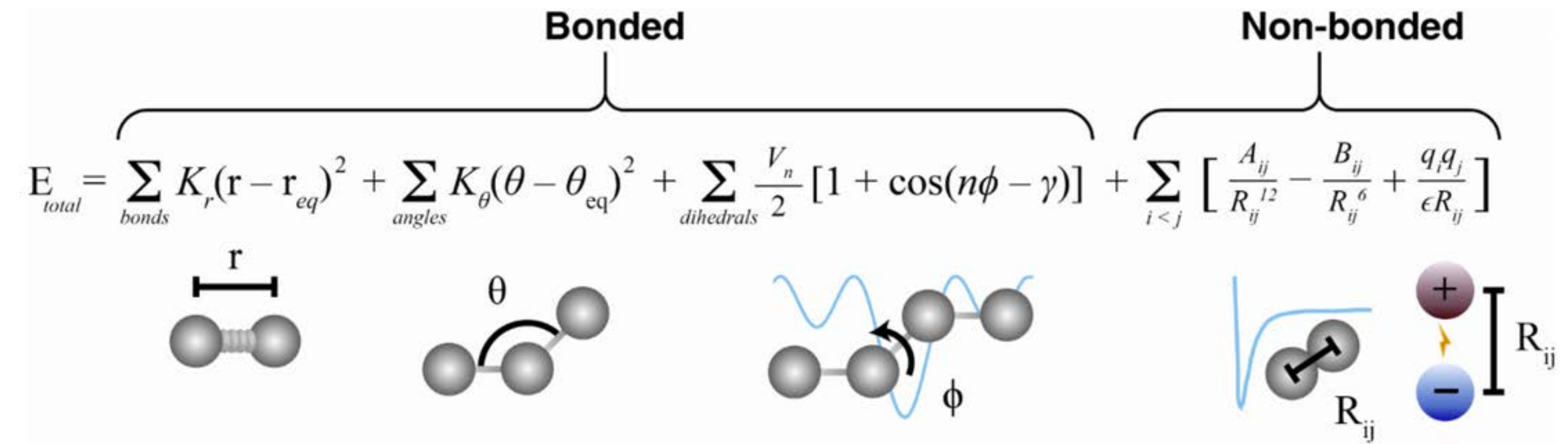
hERG          CYP3A4

porin permeation

crystal polymorphs, etc.

# FREE ENERGY CALCULATIONS (AND MUCH OF COMP CHEM) FUNDAMENTALLY RELIES ON MOLECULAR MECHANICS FORCE FIELDS

## typical class I molecular mechanics force field



$$E_{total} = \underbrace{\sum_{bonds} K_r(r - r_{eq})^2 + \sum_{angles} K_\theta(\theta - \theta_{eq})^2 + \sum_{dihedrals} \frac{V_n}{2}[1 + \cos(n\phi - \gamma)]}_{\text{Bonded}} + \underbrace{\sum_{i < j}\left[\frac{A_{ij}}{R_{ij}^{12}} - \frac{B_{ij}}{R_{ij}^{6}} + \frac{q_i q_j}{\epsilon R_{ij}}\right]}_{\text{Non-bonded}}$$

Shan, Kim, Eastwood, Dror, Seeliger, Shaw. JACS 133:9181, 2011
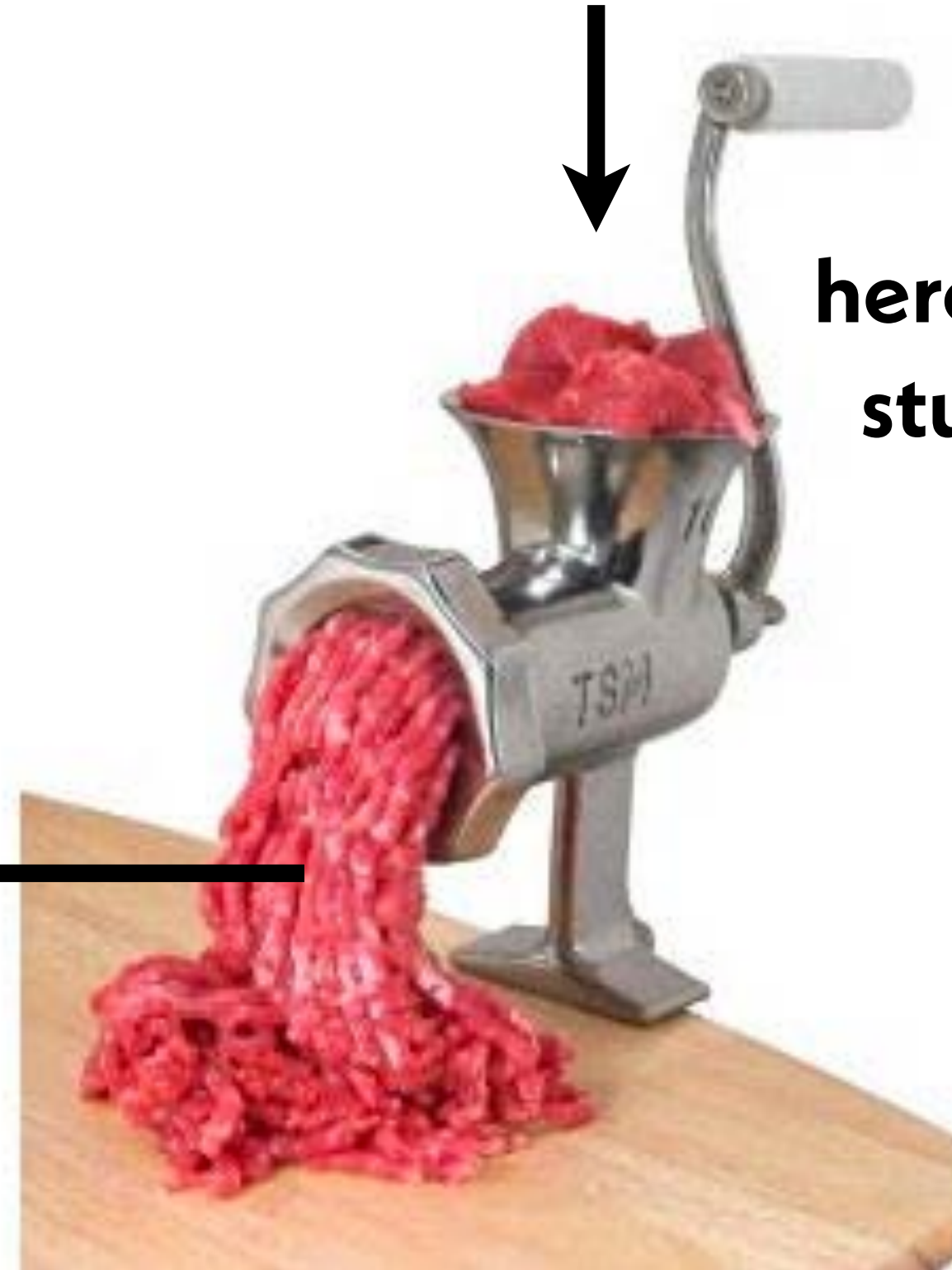Durrant, McCammon. Molecular dynamics simulations and drug discovery. BMC Biology, 2011

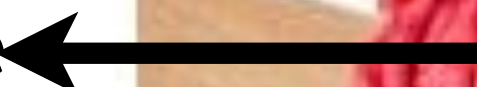# FORCE FIELDS HAVE TRADITIONALLY BEEN HEROIC PRODUCTS OF HUMAN EFFORT
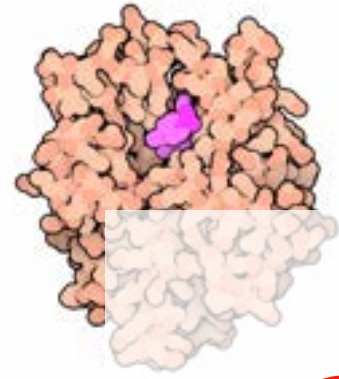
experimental data
quantum chemistry
keen chemical intuition

heroic effort by graduate
students and postdocs

a parameter set we
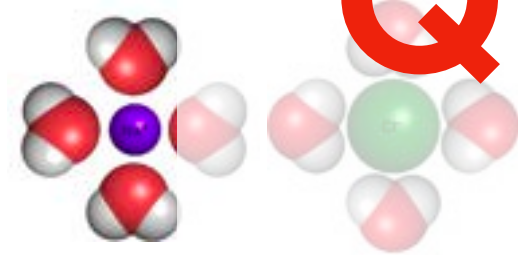desperately hope someone
actually uses

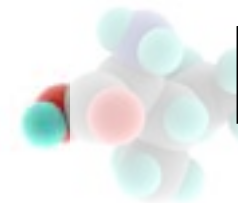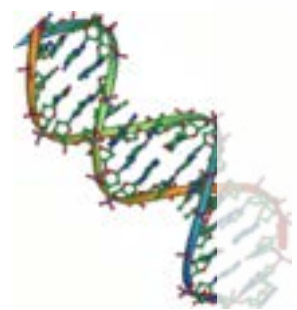# FORCE FIELDS HAVE TRADITIONALLY BEEN HEROIC PRODUCTS OF HUMAN EFFORT
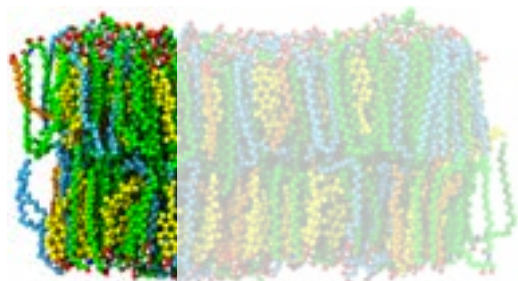
proteins

post-translational modifications

**Amber20 recommendations**

J. A. Maier; C. Martinez; K. Kasavajhala; L. Wickstrom; K. E. Hauser; C. Simmerling. ff14SB: Improving the Accuracy of Protein Side Chain and Backbone Parameters from ff99SB. *J. Chem. Theory Comput.*, **2015**, *11*, 3696–3713.

## Quickly adds up to >100 human-years

W. D. Cornell; P. Cieplak; C. I. Bayly; I. R. Gould; K. M. Merz, Jr.; D. M. Ferguson; D. C. Spellmeyer; ... force field for the simulation of proteins, nucleic ... *J. Am. ...* 1995, 117, 5179–5197.

N. Homeyer; A. H. C. Horn; H. Lanig; H. Sticht. AMBER force-field parameters for phosphorylated amino acids in different protonation states: phosphoserine, phosphothreonine, phosphotyrosine, and phosphohistidine. *J. Mol. Model.*. **2006**. *12*. 281–289.

H. W. Horn; W. C. Swope; J. W. Pitera; J. D. Madura; T. J. Dick; G. L. Hura; T. Head-Gordon. Development of an improved four-site water model for biomolecular simulations: TIP4P-Ew. *J. Chem. Phys.*, **2004**, *120*, 9665–9678.

**Intended to be compatible, but not co-parameterized**

J. S. Joung; T. E. Cheatham III. Molecular dynamics simulations of the dynamic and energetic properties ... fic ion parameters. *J. Phys. Chem. B*, 2009, *113*, 13279–13291.

**Significant effort is required to extend to new areas**

P. Li; B. P. Roberts; D. K. Chakravorty; K. M. Merz, Jr. Rational Design of Particle Mesh Ewald Compatible ... tions in Explicit Solvent. *J. Chem. Theory Comput.*, **2013**, *9*, 2733–2748.

**(e.g. covalent inhibitors, bio-inspired polymers, etc.)**

J. Wang; R. M. Wolf; J. W. Caldwell; P. A. Kollamn; D. A. Case. Development and testing of a general ... 1157–1174.

**Nobody is going to want to refit this based on some new data**

R. Galindo-Murillo; J. C. Robertson; M. Zgarbovic; J. Sponer; M. Otyepka; P. Jureska; T. E. Cheatham. ... State of Amber Force Field Modifications for DNA. *J. Chem. Theory Comput.*, **2016**,

A. Perez; I. Marchan; D. Svozil; J. Sponer; T. E. Cheatham; C. A. Laughton; M. Orozco. Refinement of the AMBER Force Field for Nucleic Acids: Improving the Description of alpha/gamma Conformers. *Biophys. J.*, **2007**, *92*, 3817–3829.

M. Zgarbova; M. Otyepka; J. Sponer; A. Mladek; P. Banas; T. E. Cheatham; P. Jurecka. Refinement of the Cornell et al. Nucleic Acids Force Field Based on Reference Quantum Chemical Calculations of Glycosidic

lipids

## How can we bring this problem into the modern era?

Á. Skjevik; B. D. Madej; R. C. Walker; K. Teigen. Lipid11: A modular framework for lipid simulations using amber. *J. Phys. Chem. B*, **2012**, *116*, 11124–11136.

C. J. Dickson; B. D. Madej; A. A. Skjevik; R. M. Betz; K. Teigen; I. R. Gould; R. C. Walker. Lipid14: The Amber Lipid Force Field. *J. Chem. Theory Comput.*, **2014**, *10*, 865–879.

carbohydrates

K. N. Kirschner; A. B. Yongye; S. M. Tschampel; J. González-Outeiriño; C. R. Daniels; B. L. Foley; R. J. Woods. GLYCAM06: A generalizable biomolecular force field. Carbohydrates. *J. Comput. Chem.*, **2008**, 29, 622–655.

# AS DRUG DISCOVERY EXPLORES NEW PARTS OF CHEMICAL SPACE, HOW CAN FORCEFIELDS KEEP UP?

**The Generalized Amber Forcefield (GAFF) only understands this space of chemistries:**



GAFF 1 was finished in **1999**, still awaiting GAFF 2 completion

Extension to new chemical space is **nontrivial**

Parameter fitting code was **never released**

Atom types have introduced numerous **errors**

Wang J, Wolf RM, Caldwell JW, Kollman PA, and Case DA. J Comput Chem 25:1157, 2004.

# CAN WE MAKE BUILDING BIMOLECULAR FORCE FIELDS AS EASY AS TRAINING A MACHINE LEARNING MODEL?

### training a neural network

```python
import tensorflow as tf
mnist = tf.keras.datasets.mnist

(x_train, y_train),(x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

model = tf.keras.models.Sequential([
  tf.keras.layers.Flatten(input_shape=(28, 28)),
  tf.keras.layers.Dense(128, activation='relu'),
  tf.keras.layers.Dropout(0.2),
  tf.keras.layers.Dense(10, activation='softmax')
])

model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

model.fit(x_train, y_train, epochs=5)
model.evaluate(x_test, y_test)
```

Run code now    Try in Google's interactive notebook

https://www.tensorflow.org/overview

import your tools

grab a standard, curated dataset

define a novel model architecture

declare your objectives in training it

fit it
use it

# CAN WE MAKE BUILDING BIMOLECULAR FORCE FIELDS AS EASY AS TRAINING A MACHINE LEARNING MODEL?

## training a neural network

```python
import tensorflow as tf
mnist = tf.keras.datasets.mnist

(x_train, y_train),(x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

model = tf.keras.models.Sequential([
  tf.keras.layers.Flatten(input_shape=(28, 28)),
  tf.keras.layers.Dense(128, activation='relu'),
  tf.keras.layers.Dropout(0.2),
  tf.keras.layers.Dense(10, activation='softmax')
])

model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

model.fit(x_train, y_train, epochs=5)
model.evaluate(x_test, y_test)
```

Run code now    Try in Google's interactive notebook

https://www.tensorflow.org/overview

## fitting a force field

```python
import openforcefield as off
training_data, benchmark_data = off.datasets.load('2019-Q1')

force_field_model = off.models.ForceFieldModel([
    off.models.forces.HarmonicBond(),
    off.models.forces.HarmonicAngle(),
    off.models.forces.PeriodicTorsion(max_order=6),
    off.models.forces.LennardJones(),
    off.models.forces.BondChargeCorrections(),
])

model.compile(optimizer='L-BFGS',
    loss='error-weighted',
    metrics=['accuracy'])

model.fit(training_data)

model.evaluate(test_data)
```

Run code now    Try in Google's interactive notebook

**open forcefield**

# An open and collaborative approach to better force fields

**OPEN SOURCE**

Software permissively licensed under the MIT License and developed openly on GitHub.

**OPEN SCIENCE**

Scientific reports as blog posts, webinars and preprints

**OPEN DATA**

Curated quantum chemical and experimental datasets used to parameterize and benchmark Open Force Fields.

NEWS        TUTORIALS        ROADMAP

http://openforcefield.org

# THE OPEN FORCE FIELD INITIATIVE AIMS TO BUILD A MODERN INFRASTRUCTURE FOR FORCE FIELD SCIENCE

**Open source <u>Python Toolkit</u>:** use the parameters in most simulation packages

**Open curated QM / physical property datasets:** build your own force fields
**MolSSI QCArchive quantum chemical data: <u>http://qcarchive.molssi.org</u>**

**Open source infrastructure:** for improving force fields with in-house data

**Open science:** everything we do is free, permissively licensed, and online

**<u>http://openforcefield.org</u>**

# WE'VE MADE RAPID AND SIGNIFICANT PROGRESS IN ACCURACY, BUT WE'RE STILL STICK WITH SLOW GENERATIONS



Open Force Field Initiative

GAFF 1 (1999)  OPLS2.1 (2015)  GAFF 2 (2016)  smirnoff99Frosst (2018)  openff 1.0 (2019)

"parsley"

thrombin
PDB101: 1PPB

HANNAH BRUCE MACDONALD
MSKCC

http://github.com/choderalab/perses

DOMINIC RUFA

# NEW GENERATIONS OF MACHINE LEARNING MODELS ARE PARTICULARLY WELL-SUITED TO CHEMISTRY

**molecule**   **bond**   **atom**

predict
properties

Learns **electronegativity** ($e_i$) and **hardness** ($s_i$)

subject to fixed charge sum constraint:

$$\{\hat{q}_i\} = \underset{q_i}{\mathrm{argmin}} \sum_i \hat{e}_i q_i + \frac{1}{2}\hat{s}_i q_i^2$$

$$\sum_i \hat{q}_i = \sum_i q_i = Q$$

Figure adapted from Zhou Z
arXiv:1706.09916

control experiment:
direct prediction of charges: RMSE **0.2800 e**

| | $R^2$ | RMSE(e) | # Samples |
|---|---|---|---|
| Overall | $0.9936_{0.9935}^{0.9937}$ | $0.0223_{0.0221}^{0.0225}$ | 299811 |

DFT charges on ChEMBL dataset from Bleiziffer, Schaller, Riniker JCIM 58:579, 2018

$q_{pred}$    $q_{true}$

C   N   O   S   P   F   Cl   Br   I   H

$$\mathbf{e}_k^{(t+1)} = \phi^e(\mathbf{e}_k^{(t)}, \sum_{i \in \mathcal{N}_k^e} \mathbf{v}_i, \mathbf{u}^{(t)}), \qquad \text{(edge update)}$$

$$\bar{\mathbf{e}}_i^{(t+1)} = \rho^{e \to v}(E_i^{(t+1)}), \qquad \text{(edge to node aggregate)}$$

$$\mathbf{v}_i^{(t+1)} = \phi^v(\bar{\mathbf{e}}_i^{(t+1)}, \mathbf{v}_i^{(t)}, \mathbf{u}^{(t)}), \qquad \text{(node update)}$$
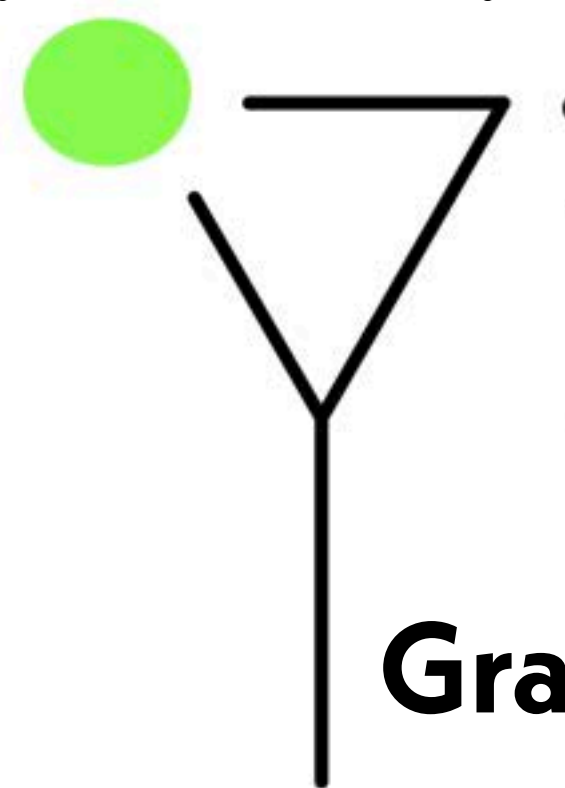
$$\bar{\mathbf{e}}^{(t+1)} = \rho^{e \to u}(E^{(t+1)}), \qquad \text{(edge to global aggregate)}$$

$$\bar{\mathbf{v}}^{(t+1)} = \rho^{v \to u}(V^{(t)}), \qquad \text{(node to global aggregate)}$$

$$\mathbf{u}^{(t+1)} = \phi^u(\bar{\mathbf{e}}^{(t+1)}, \bar{\mathbf{v}}^{(t+1)}, \mathbf{u}^{(t)}), \qquad \text{(global update)}$$

# ∇imlet

## Graph Inference on MoLEcular Topology

**preprint:** https://arxiv.org/abs/1909.07903
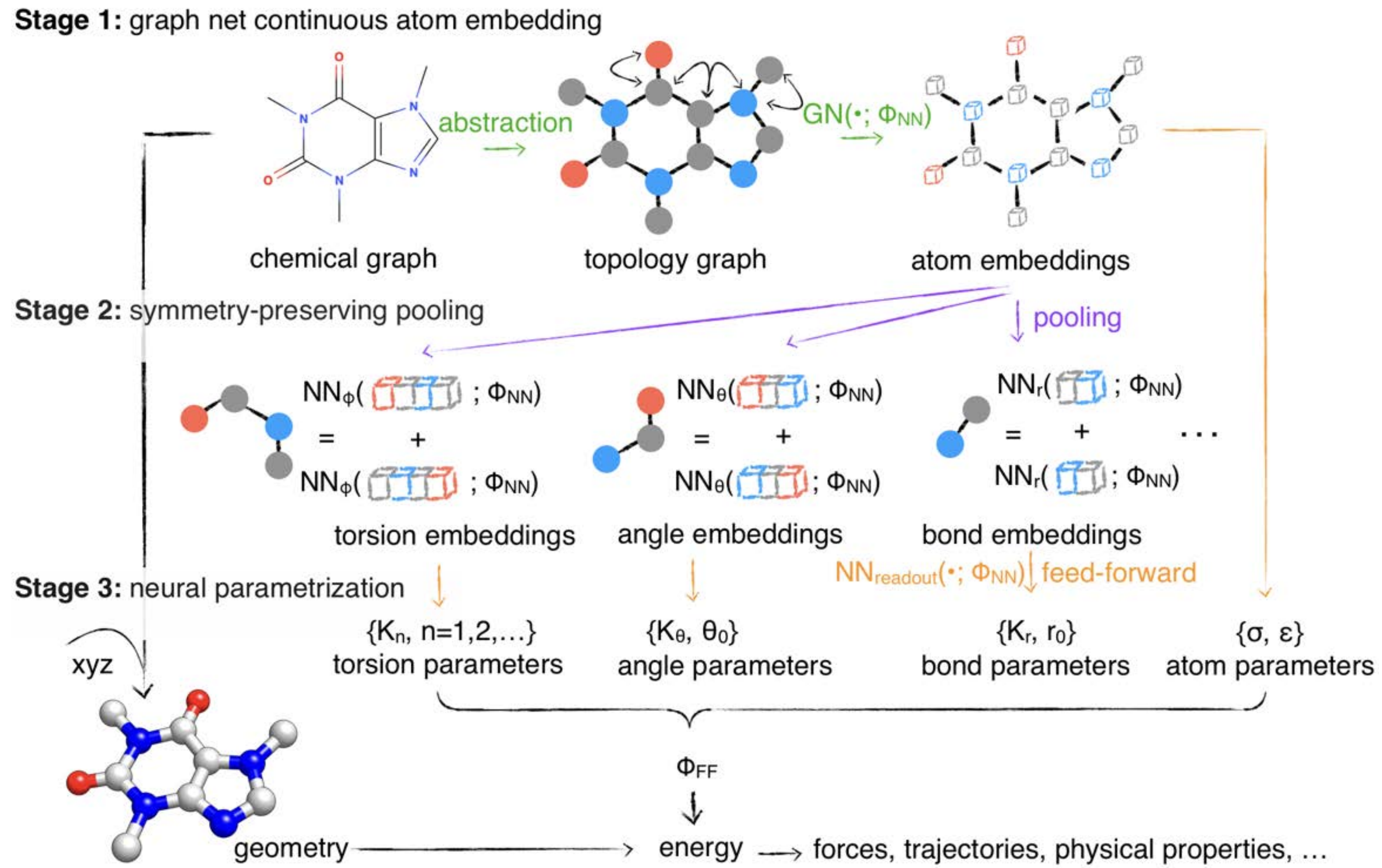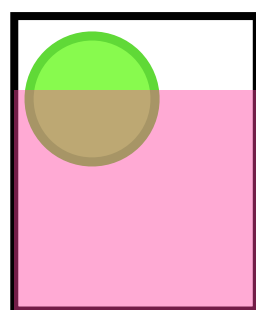**code:** http://github.com/choderalab/gimlet

YUANQING
WANG

# espaloma: extensible surrogate potential of *ab initio* learned and optimized by message-passing algorithm

use of only **chemical graph** means that model can generate parameters for small molecules, proteins, nucleic acids, covalent ligands, carbohydrates, etc.
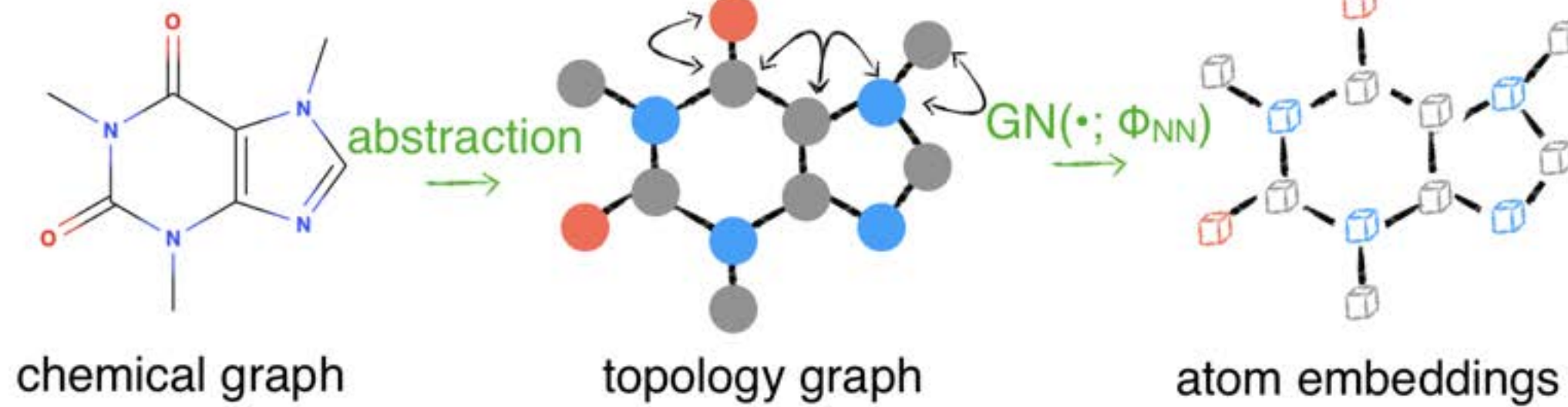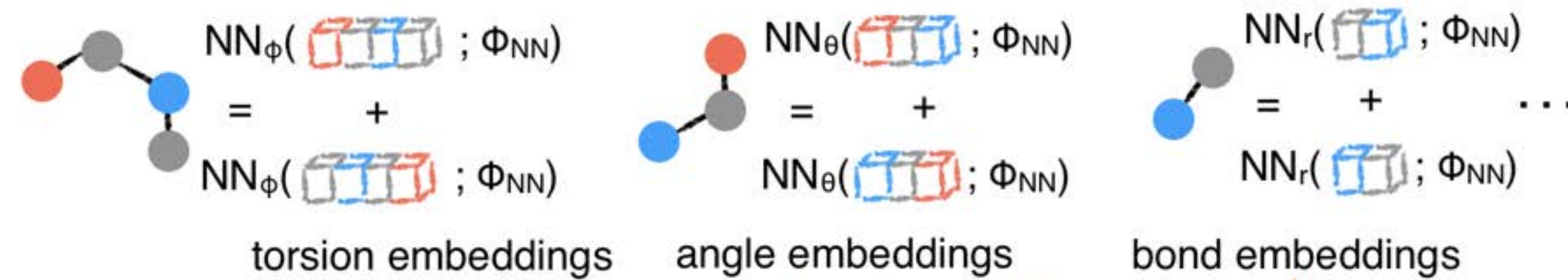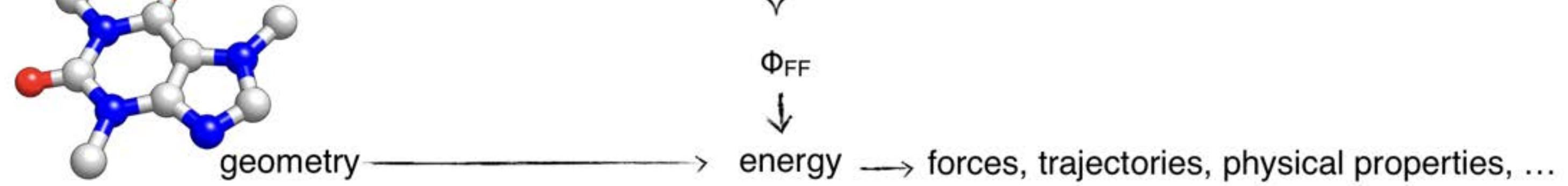


**JOSH FASS**

**YUANQING WANG**

preprint: https://arxiv.org/abs/2010.01196
code: https://github.com/choderalab/espaloma

# **espaloma**: **e**xtensible **s**urrogate **p**otential of **a**b *initio* **l**earned and **o**ptimized by **m**essage-passing **a**lgorithm
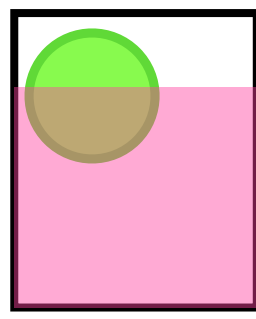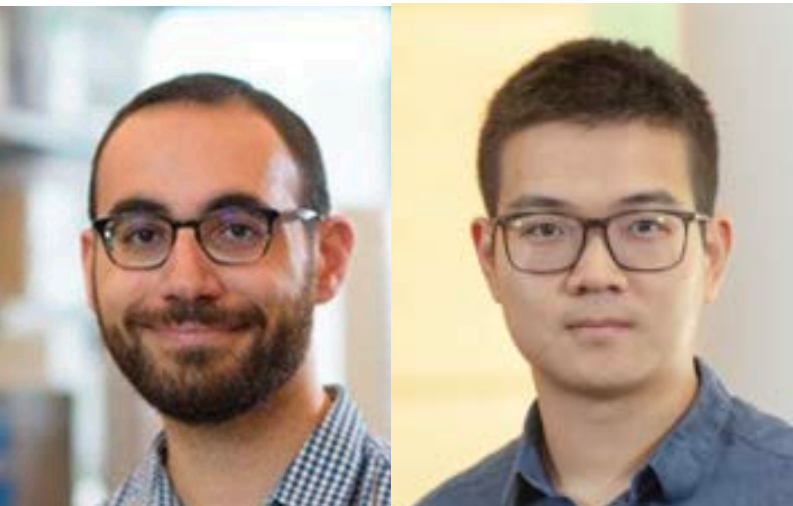


entire model is **end-to-end differentiable** so can be fit to any loss function by standard automatic differentiation machine learning frameworks
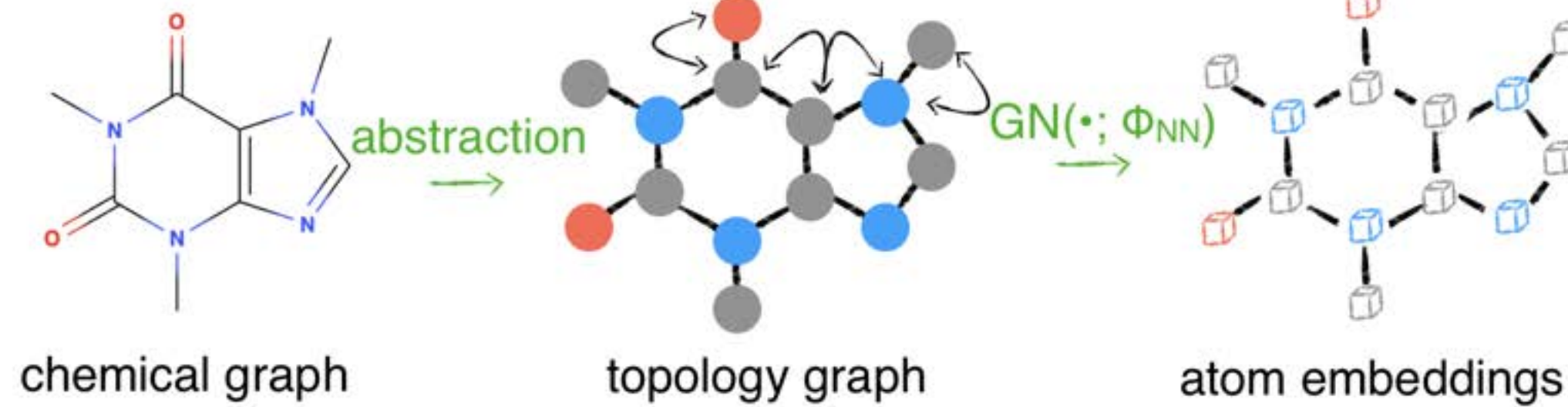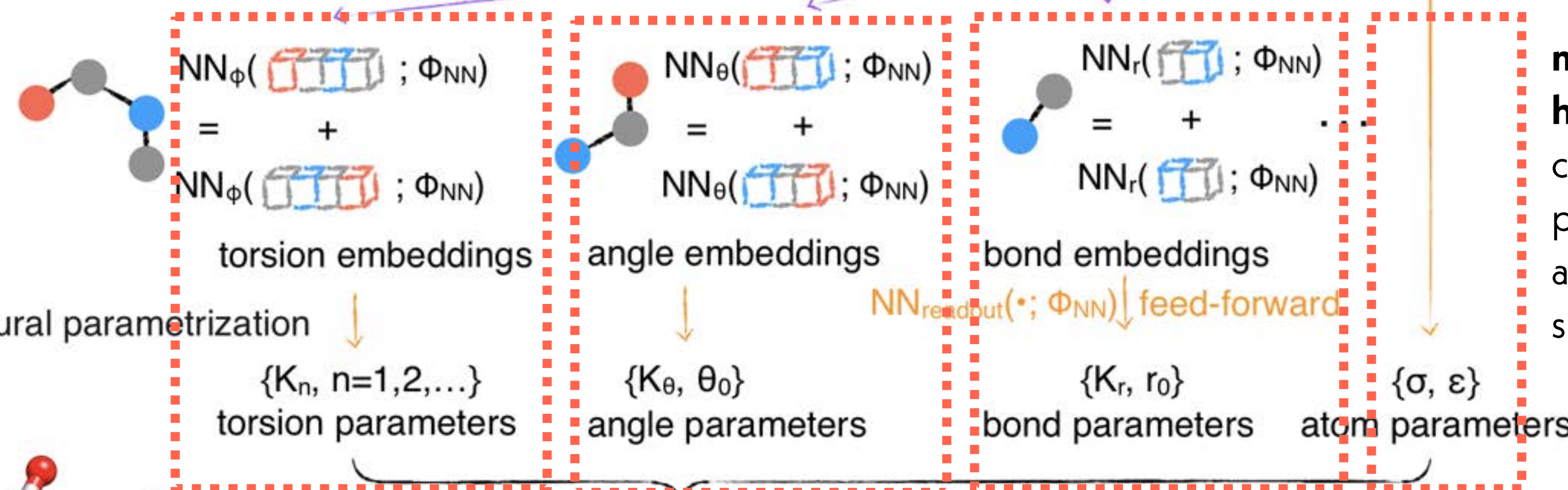
**JOSH FASS**

**YUANQING WANG**

**preprint**: https://arxiv.org/abs/2010.01196
**code**: https://github.com/choderalab/espaloma

# **es**paloma: **e**xtensible **s**urrogate **p**otential of *ab initio* **l**earned and **o**ptimized by **m**essage-passing **a**lgorithm



**Stage 1:** graph net continuous atom embedding

chemical graph → *abstraction* → topology graph → GN(•; $\Phi_{NN}$) → atom embeddings

**Stage 2:** symmetry-preserving pooling

$NN_\phi( \quad ; \Phi_{NN})$ = + $NN_\phi( \quad ; \Phi_{NN})$
torsion embeddings

$NN_\theta( \quad ; \Phi_{NN})$ = + $NN_\theta( \quad ; \Phi_{NN})$
angle embeddings

$NN_r( \quad ; \Phi_{NN})$ = + $NN_r( \quad ; \Phi_{NN})$ · ...
bond embeddings

↓ pooling

$NN_{readout}(•; \Phi_{NN})$ feed-forward

**Stage 3:** neural parametrization

xyz

$\{K_n, n=1,2,\ldots\}$ torsion parameters

$\{K_\theta, \theta_0\}$ angle parameters

$\{K_r, r_0\}$ bond parameters

$\{\sigma, \epsilon\}$ atom parameters

geometry ⟶ $\Phi_{FF}$ ↓ energy ⟶ forces, trajectories, physical properties, …
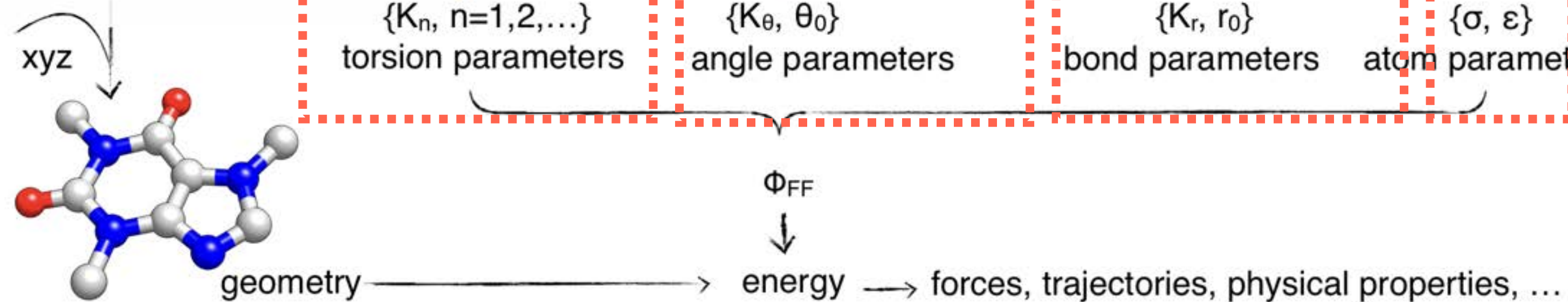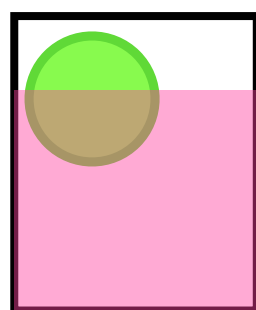
**modular and extensible handling of potential terms:** charge model parameters, point polarizabilities, alternative vdW forms, special 1-4 parameters, etc.

JOSH FASS

YUANQING WANG

**preprint:** https://arxiv.org/abs/2010.01196
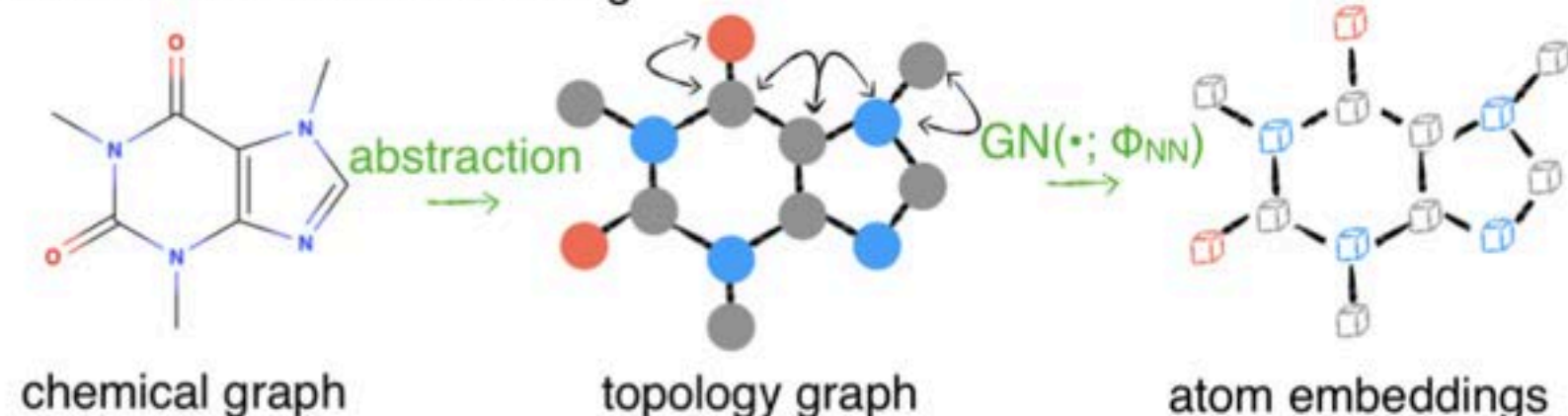**code:** https://github.com/choderalab/espaloma

# ESPALOMA MAKES BUILDING A NEW FORCE FIELD EASY

## espaloma architecture

**building a new force field**



(implemented in pytorch)

http://github.com/choderalab/espaloma

**YUANQING WANG**

```
import torch, dgl, espaloma as esp

# retrieve OpenFF Gen2 Optimization Dataset
dataset = esp.data.dataset.GraphDataset.load("gen2").view(batch_size=128)

# define Espaloma stage I: graph -> atom latent representation
representation = esp.nn.Sequential(
    layer=esp.nn.layers.dgl_legacy.gn("SAGEConv"), # use SAGEConv implementation in DGL
    config=[128, "relu", 128, "relu", 128, "relu"], # 3 layers, 128 units, ReLU activation
)

# define Espaloma stage II and III:
# atom latent representation -> bond, angle, and torsion representation and parameters
readout = esp.nn.readout.janossy.JanossyPooling(
    in_features=128, config=[128, "relu", 128, "relu", 128, "relu"],
    out_features={              # define modular MM parameters Espaloma will assign
        1: {"e": 1, "s": 1}, # atom hardness and electronegativity
        2: {"coefficients": 2}, # bond linear combination
        3: {"coefficients": 3}, # angle linear combination
        4: {"k": 6}, # torsion barrier heights (can be positive or negative)
    },
)

# compose all three Espaloma stages into an end-to-end model
espaloma_model = torch.nn.Sequential(
                representation, readout,
                esp.mm.geometry.GeometryInGraph(), esp.mm.energy.EnergyInGraph(),
                esp.nn.readout.charge_equilibrium.ChargeEquilibrium(),
)

# define training metric
metrics = [
    esp.metrics.GraphMetric(
            base_metric=torch.nn.MSELoss(), # use mean-squared error loss
            between=['u', "u_ref"],         # between predicted and QM energies
            level="g", # compare on graph level
    )
    esp.metrics.GraphMetric(
            base_metric=torch.nn.MSELoss(), # use mean-squared error loss
            between=['q', "q_hat"],         # between predicted and reference charges
            level="n1", # compare on node level
    )
]

# fit Espaloma model to training data
results = esp.Train(
    ds_tr=dataset, net=espaloma_model, metrics=metrics,
    device=torch.device('cuda:0'), n_epochs=5000,
    optimizer=lambda net: torch.optim.Adam(net.parameters(), 1e-3), # use Adam optimizer
).run()

torch.save(espaloma_model, "espaloma_model.pt") # save model
```
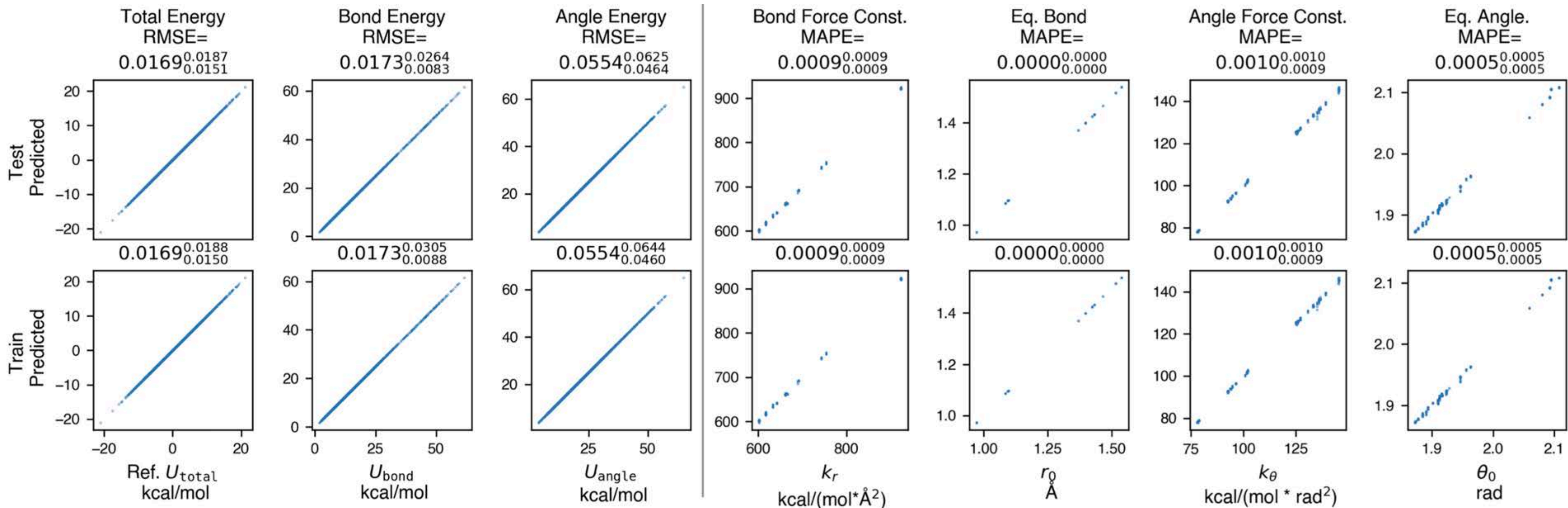
**Listing 1.** Defining and training a modular Espaloma model.

# ESPALOMA CAN LEARN TO REPRODUCE LEGACY MM FORCE FIELDS WITH LOW RMSE ERROR IN CONFORMATIONAL ENERGIES

## conformer energies

## force field parameters

**YUANQING WANG**

# ESPALOMA OUTPERFORMS CURRENT FORCE FIELDS IN QM ACCURACY AND CAN BE EASILY TRAINED FOR HETEROGENEOUS SYSTEMS
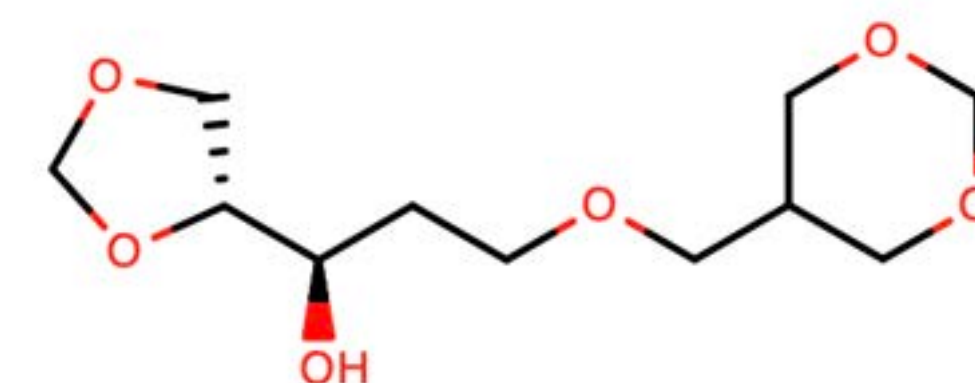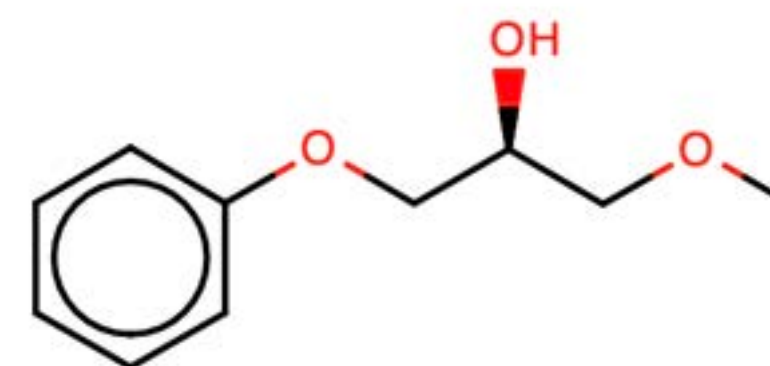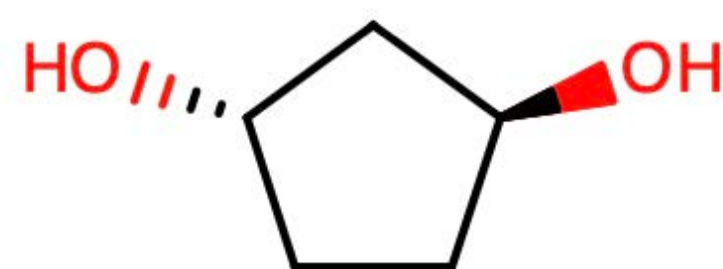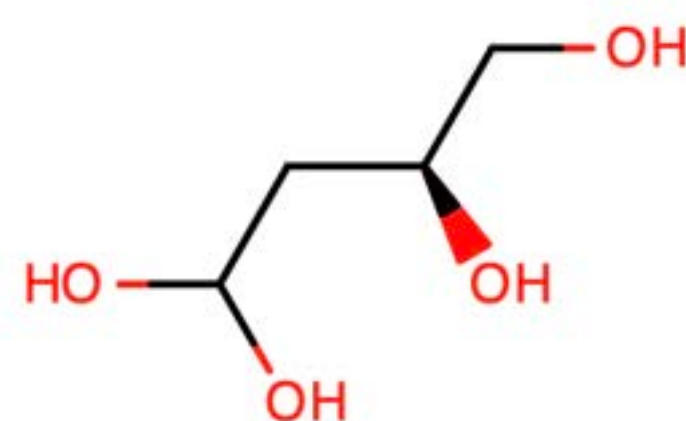
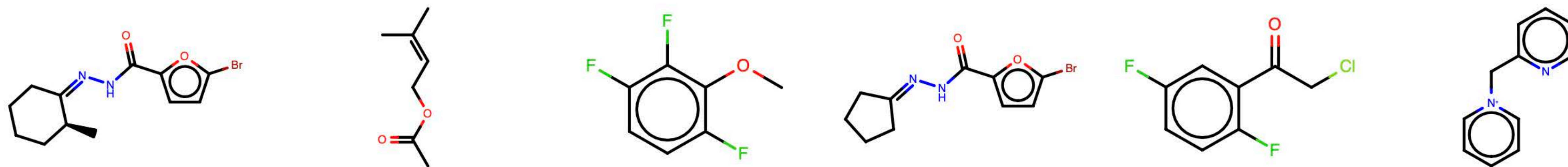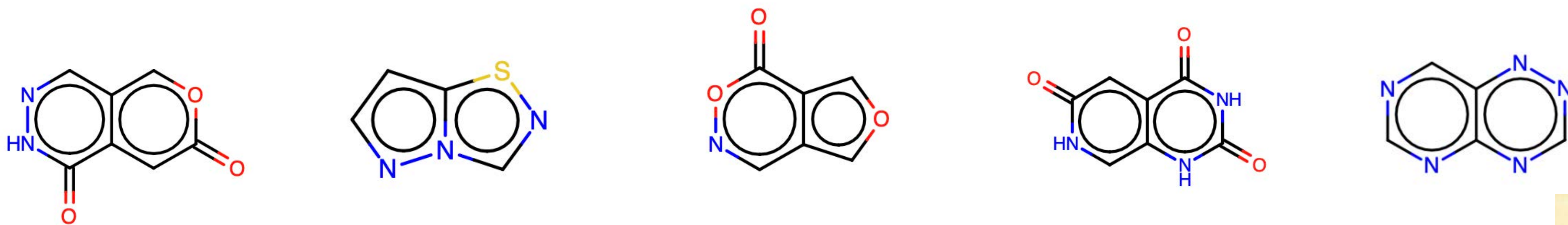| (a) dataset | # mols | # trajs | # snapshots | Espaloma RMSE | | Legacy FF RMSE (kcal/mol) (Test molecules) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Train | Test | OpenFF 1.2.0 | GAFF-1.81 | GAFF-2.11 | Amber ff14SB |
| **PhAlkEthOH** (simple CHO) | 7408 | 12592 | 244036 | $0.8656^{0.9131}_{0.8225}$ | $1.1398^{1.2332}_{1.0715}$ | $1.6071^{1.6915}_{1.5197}$ | $1.7267^{1.7935}_{1.6543}$ | $1.7406^{1.8148}_{1.6679}$ | |
| **OpenFF Gen2 Optimization** (druglike) | 792 | 3977 | 23748 | $0.7413^{0.7920}_{0.6914}$ | $0.7600^{0.8805}_{0.6644}$ | $2.1768^{2.3388}_{2.0380}$ | $2.4274^{2.5207}_{2.3300}$ | $2.5386^{2.6640}_{2.4370}$ | |
| **VEHICLe** (heterocyclic) | 24867 | 24867 | 234326 | $0.4476^{0.4690}_{0.4273}$ | $0.4233^{0.4414}_{0.4053}$ | $8.0247^{8.2456}_{7.8271}$ | $8.0077^{8.2313}_{7.7647}$ | $9.4014^{9.6434}_{9.2135}$ | |
| **PepConf** (peptides) | 736 | 7560 | 22154 | $1.2714^{1.3616}_{1.1899}$ | $1.8727^{1.9749}_{1.7309}$ | $3.6143^{3.7288}_{3.4870}$ | $4.4446^{4.5738}_{4.3386}$ | $4.3356^{4.4641}_{4.1965}$ | $3.1502^{3.1859,*}_{3.1117}$ |
| **joint**   OpenFF Gen2 Optimization | 1528 | 11537 | 45902 | $0.8264^{0.9007}_{0.7682}$ | $1.8764^{1.9947}_{1.7827}$ | $2.1768^{2.3388}_{2.0380}$ | $2.4274^{2.5207}_{2.3300}$ | $2.5386^{2.6640}_{2.4370}$ | |
|   PepConf | | | | $1.2038^{1.3056}_{1.1178}$ | $1.7307^{1.8439}_{1.6053}$ | $3.6143^{3.7288}_{3.4870}$ | $4.4446^{4.5738}_{4.3386}$ | $4.3356^{4.4641}_{4.1965}$ | $3.1502^{3.1859,*}_{3.1117}$ |

preprint: https://arxiv.org/abs/2010.01196
code: http://github.com/choderalab/espaloma

**YUANQING WANG**

# ESPALOMA OUTPERFORMS CURRENT FORCE FIELDS IN QM ACCURACY AND CAN BE EASILY TRAINED FOR HETEROGENEOUS SYSTEMS

| (a) dataset | # mols | # trajs | # snapshots | Espaloma RMSE | | Legacy FF RMSE (kcal/mol) (Test molecules) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Train | Test | OpenFF 1.2.0 | GAFF-1.81 | GAFF-2.11 | Amber ff14SB |
| **PhAlkEthOH** (simple CHO) | 7408 | 12592 | 244036 | $0.8656_{0.8225}^{0.9131}$ | $1.1398_{1.0715}^{1.2332}$ | $1.6071_{1.5197}^{1.6915}$ | $1.7267_{1.6543}^{1.7935}$ | $1.7406_{1.6679}^{1.8148}$ | |

## PhAlkEthOh: Phenyls, Alkanes, Ethers, and alcohols (OH)
### (a low-complexity chemical space)

**YUANQING WANG**

# ESPALOMA OUTPERFORMS CURRENT FORCE FIELDS IN QM ACCURACY AND CAN BE EASILY TRAINED FOR HETEROGENEOUS SYSTEMS

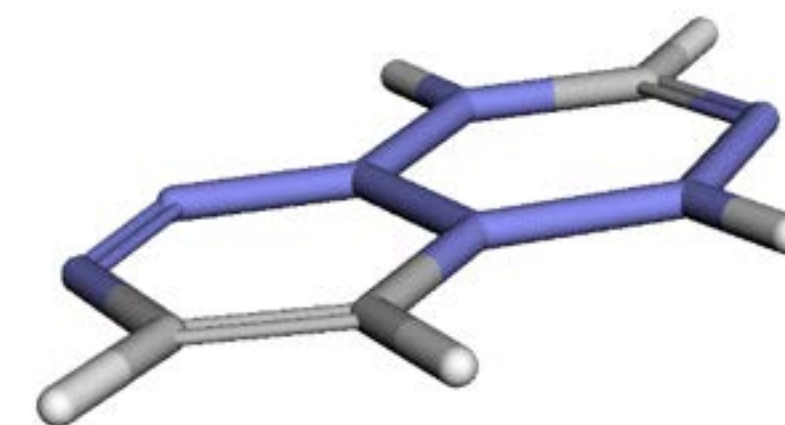| (a) dataset | # mols | # trajs | # snapshots | Espaloma RMSE | | Legacy FF RMSE (kcal/mol) (Test molecules) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Train | Test | OpenFF 1.2.0 | GAFF-1.81 | GAFF-2.11 | Amber ff14SB |
| **PhAlkEthOH** (simple CHO) | 7408 | 12592 | 244036 | $0.8656_{0.8225}^{0.9131}$ | $1.1398_{1.0715}^{1.2332}$ | $1.6071_{1.5197}^{1.6915}$ | $1.7267_{1.6543}^{1.7935}$ | $1.7406_{1.6679}^{1.8148}$ | |
| **OpenFF Gen2 Optimization** (druglike) | 792 | 3977 | 23748 | $0.7413_{0.6914}^{0.7920}$ | $0.7600_{0.6644}^{0.8805}$ | $2.1768_{2.0380}^{2.3388}$ | $2.4274_{2.3300}^{2.5207}$ | $2.5386_{2.4370}^{2.6640}$ | |

## OpenFF Gen2 Optimization set: Diverse druglike fragments challenging for force fields
### (a moderate-complexity chemical space)

**YUANQING WANG**

# ESPALOMA OUTPERFORMS CURRENT FORCE FIELDS IN QM ACCURACY AND CAN BE EASILY TRAINED FOR HETEROGENEOUS SYSTEMS

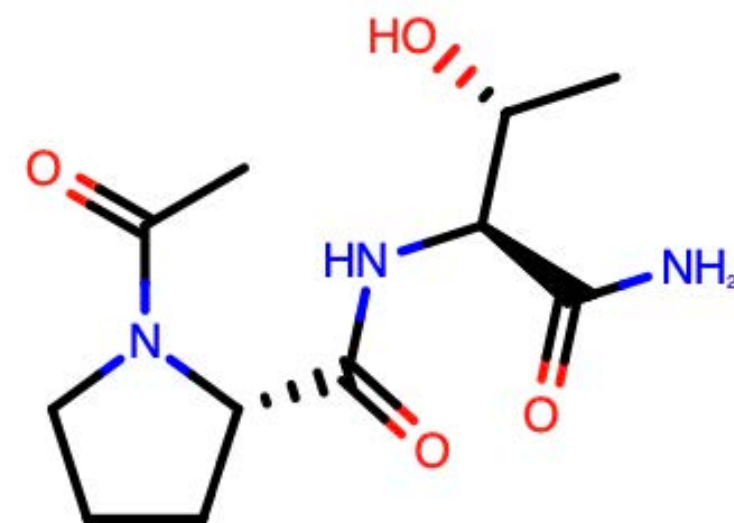| (a) dataset | # mols | # trajs | # snapshots | Espaloma RMSE | | Legacy FF RMSE (kcal/mol) (Test molecules) | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | Train | Test | OpenFF 1.2.0 | GAFF-1.81 | GAFF-2.11 | Amber ff14SB |
| **PhAlkEthOH** (simple CHO) | 7408 | 12592 | 244036 | $0.8656_{0.8225}^{0.9131}$ | $1.1398_{1.0715}^{1.2332}$ | $1.6071_{1.5197}^{1.6915}$ | $1.7267_{1.6543}^{1.7935}$ | $1.7406_{1.6679}^{1.8148}$ | |
| **OpenFF Gen2 Optimization** (druglike) | 792 | 3977 | 23748 | $0.7413_{0.6914}^{0.7920}$ | $0.7600_{0.6644}^{0.8805}$ | $2.1768_{2.0380}^{2.3388}$ | $2.4274_{2.3300}^{2.5207}$ | $2.5386_{2.4370}^{2.6640}$ | |
| **VEHICLe** (heterocyclic) | 24867 | 24867 | 234326 | $0.4476_{0.4273}^{0.4690}$ | $0.4233_{0.4053}^{0.4414}$ | $8.0247_{7.8271}^{8.2456}$ | $8.0077_{7.7647}^{8.2313}$ | $9.4014_{9.2135}^{9.6434}$ | |

**VEHICLe**: Virtual exploratory heterocyclic drug scaffold library
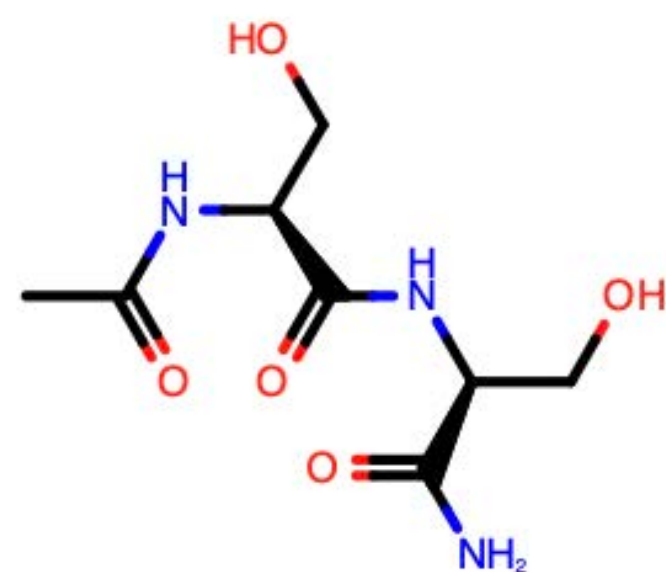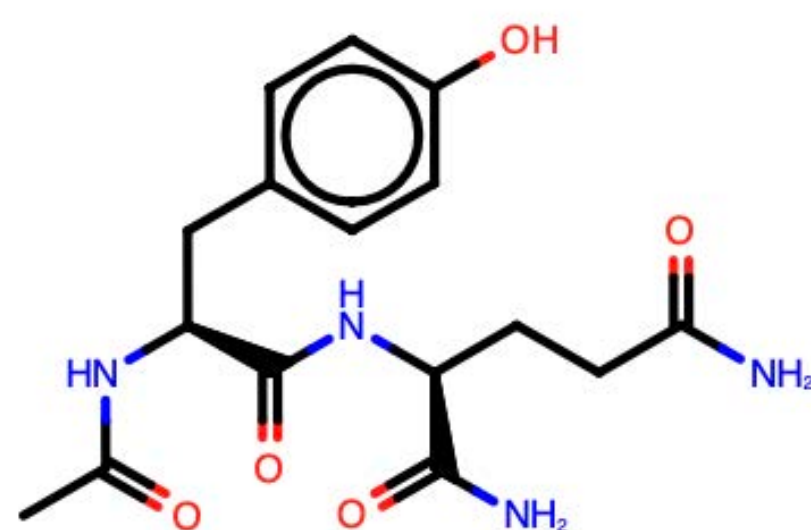(aromatic bicyclic heterocyclic compounds containing C, N, O, S, H)

**YUANQING WANG**

# ESPALOMA OUTPERFORMS CURRENT FORCE FIELDS IN QM ACCURACY AND CAN BE EASILY TRAINED FOR HETEROGENEOUS SYSTEMS

| (a) dataset | # mols | # trajs | # snapshots | Espaloma RMSE | | Legacy FF RMSE (kcal/mol) (Test molecules) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Train | Test | OpenFF 1.2.0 | GAFF-1.81 | GAFF-2.11 | Amber ff14SB |
| **PhAlkEthOH** (simple CHO) | 7408 | 12592 | 244036 | $0.8656^{0.9131}_{0.8225}$ | $1.1398^{1.2332}_{1.0715}$ | $1.6071^{1.6915}_{1.5197}$ | $1.7267^{1.7935}_{1.6543}$ | $1.7406^{1.8148}_{1.6679}$ | |
| **OpenFF Gen2 Optimization** (druglike) | 792 | 3977 | 23748 | $0.7413^{0.7920}_{0.6914}$ | $0.7600^{0.8805}_{0.6644}$ | $2.1768^{2.3388}_{2.0380}$ | $2.4274^{2.5207}_{2.3300}$ | $2.5386^{2.6640}_{2.4370}$ | |
| **VEHICLe** (heterocyclic) | 24867 | 24867 | 234326 | $0.4476^{0.4690}_{0.4273}$ | $0.4233^{0.4414}_{0.4053}$ | $8.0247^{8.2456}_{7.8271}$ | $8.0077^{8.2313}_{7.7647}$ | $9.4014^{9.6434}_{9.2135}$ | |

## Comparison with QCArchive data



initial

QM minimized

DFT B3LYP-D3(BJ) / DZVP

**YUANQING WANG**

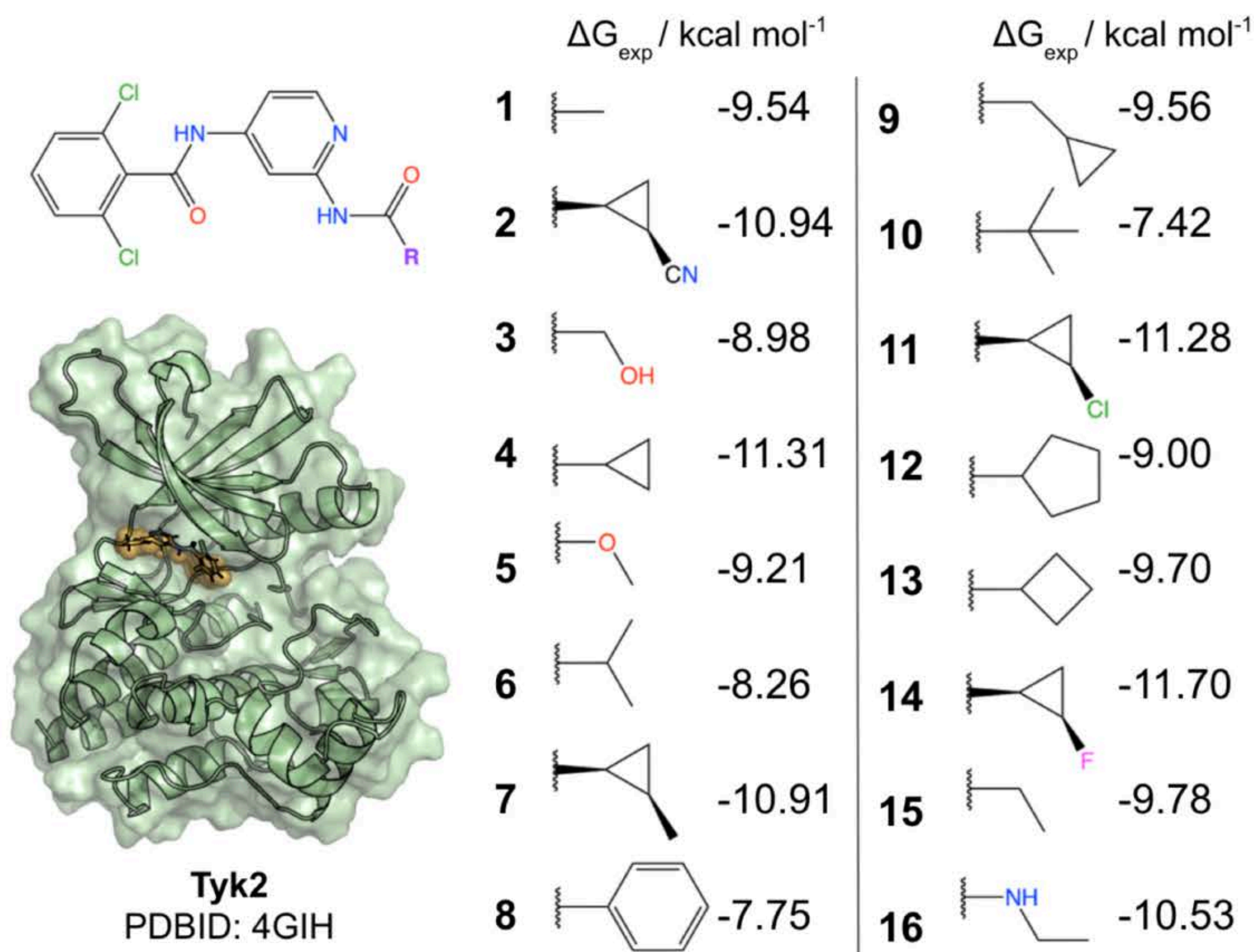# ESPALOMA OUTPERFORMS CURRENT FORCE FIELDS IN QM ACCURACY AND CAN BE EASILY TRAINED FOR HETEROGENEOUS SYSTEMS

| (a) dataset | # mols | # trajs | # snapshots | Espaloma RMSE | | Legacy FF RMSE (kcal/mol) (Test molecules) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Train | Test | OpenFF 1.2.0 | GAFF-1.81 | GAFF-2.11 | Amber ff14SB |
| **PhAlkEthOH** (simple CHO) | 7408 | 12592 | 244036 | $0.8656_{0.8225}^{0.9131}$ | $1.1398_{1.0715}^{1.2332}$ | $1.6071_{1.5197}^{1.6915}$ | $1.7267_{1.6543}^{1.7935}$ | $1.7406_{1.6679}^{1.8148}$ | |
| **OpenFF Gen2 Optimization** (druglike) | 792 | 3977 | 23748 | $0.7413_{0.6914}^{0.7920}$ | $0.7600_{0.6644}^{0.8805}$ | $2.1768_{2.0380}^{2.3388}$ | $2.4274_{2.3300}^{2.5207}$ | $2.5386_{2.4370}^{2.6640}$ | |
| **VEHICLe** (heterocyclic) | 24867 | 24867 | 234326 | $0.4476_{0.4273}^{0.4690}$ | $0.4233_{0.4053}^{0.4414}$ | $8.0247_{7.8271}^{8.2456}$ | $8.0077_{7.7647}^{8.2313}$ | $9.4014_{9.2135}^{9.6434}$ | |
| **PepConf** (peptides) | 736 | 7560 | 22154 | $1.2714_{1.1899}^{1.3616}$ | $1.8727_{1.7309}^{1.9749}$ | $3.6143_{3.4870}^{3.7288}$ | $4.4446_{4.3386}^{4.5738}$ | $4.3356_{4.1965}^{4.4641}$ | $3.1502_{3.1117}^{3.1859,*}$ |

**PepConf**: Short peptides, including disulfides and cyclic peptides

**YUANQING WANG**

# ESPALOMA OUTPERFORMS CURRENT FORCE FIELDS IN QM ACCURACY AND CAN BE EASILY TRAINED FOR HETEROGENEOUS SYSTEMS

| (a) dataset | | # mols | # trajs | # snapshots | Espaloma RMSE | | Legacy FF RMSE (kcal/mol) (Test molecules) | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Train | Test | OpenFF 1.2.0 | GAFF-1.81 | GAFF-2.11 | Amber ff14SB |
| **PhAlkEthOH** (simple CHO) | | 7408 | 12592 | 244036 | $0.8656^{0.9131}_{0.8225}$ | $1.1398^{1.2332}_{1.0715}$ | $1.6071^{1.6915}_{1.5197}$ | $1.7267^{1.7935}_{1.6543}$ | $1.7406^{1.8148}_{1.6679}$ | |
| **OpenFF Gen2 Optimization** (druglike) | | 792 | 3977 | 23748 | $0.7413^{0.7920}_{0.6914}$ | $0.7600^{0.8805}_{0.6644}$ | $2.1768^{2.3388}_{2.0380}$ | $2.4274^{2.5207}_{2.3300}$ | $2.5386^{2.6640}_{2.4370}$ | |
| **VEHICLe** (heterocyclic) | | 24867 | 24867 | 234326 | $0.4476^{0.4690}_{0.4273}$ | $0.4233^{0.4414}_{0.4053}$ | $8.0247^{8.2456}_{7.8271}$ | $8.0077^{8.2313}_{7.7647}$ | $9.4014^{9.6434}_{9.2135}$ | |
| **PepConf** (peptides) | | 736 | 7560 | 22154 | $1.2714^{1.3616}_{1.1899}$ | $1.8727^{1.9749}_{1.7309}$ | $3.6143^{3.7288}_{3.4870}$ | $4.4446^{4.5738}_{4.3386}$ | $4.3356^{4.4641}_{4.1965}$ | $3.1502^{3.1859,*}_{3.1117}$ |
| joint | OpenFF Gen2 Optimization | 1528 | 11537 | 45902 | $0.8264^{0.9007}_{0.7682}$ | $1.8764^{1.9947}_{1.7827}$ | $2.1768^{2.3388}_{2.0380}$ | $2.4274^{2.5207}_{2.3300}$ | $2.5386^{2.6640}_{2.4370}$ | |
| | PepConf | | | | $1.2038^{1.3056}_{1.1178}$ | $1.7307^{1.8439}_{1.6053}$ | $3.6143^{3.7288}_{3.4870}$ | $4.4446^{4.5738}_{4.3386}$ | $4.3356^{4.4641}_{4.1965}$ | $3.1502^{3.1859,*}_{3.1117}$ |





**Tyk2 from OpenFF benchmark set**
espaloma **joint** model
+ TIP3P water

**YUANQING WANG**

# ESPALOMA SMALL MOLECULE PARAMETERS PERFORM AS WELL OR BETTER THAN MODERN BIOMOLECULAR FORCE FIELDS



MIKE HENRY

IVÁN PULIDO

IVY ZHANG

DOMINIC RUFA

HANNAH BRUCE CDONALD

YUANQING WANG

Tyk2
PDBID: 4GIH

**OpenFF 1.2.0** small molecule
Amber ff14SB protein
TIP3P water

**espaloma "joint" 0.2.2** small molecule
Amber ff14SB protein
TIP3P water

# ESPALOMA CAN ALSO FIT EXPERIMENTAL FREE ENERGIES

experimental hydration
free energies from **FreeSolv**
https://github.com/MobleyLab/FreeSolv

loss function:

$$L(\Phi_{NN}) = \sum_{n=1}^{N} \frac{[\Delta G_n(\Phi_{NN}) - \Delta G_n^{\exp}]^2}{\sigma_n^2}$$

Here, ΔG estimated via one-step free energy perturbation, but can easily differentiate properties through MBAR

**YUANQING WANG**

**JOSH FASS**

### OBC2 GBSA FreeSolv RMSE



- training
- validation
- --- FreeSolv reference calculations

**preprint**: https://arxiv.org/abs/2010.01196
**code**: https://github.com/choderalab/espaloma

# A NEW GENERATION OF QUANTUM MACHINE LEARNING (QML) POTENTIALS PROVIDE SIGNIFICANTLY MORE FLEXIBILITY IN FUNCTIONAL FORM, THOUGH AT MUCH GREATER COST

**ANI** family of quantum machine learning (QML) potentials

**radial** and **angular** features

deep neural network for each atom

excellent agreement with DFT



OLEXANDR ISAYEV    ADRIAN ROITBERG

# HYBRID QUANTUM MACHINE LEARNING / MOLECULAR MECHANICS (QML/MM) FREE ENERGY CALCULATIONS CUT ERROR IN HALF



Rufa, Bruce Macdonald, Fass, Wieder, Grinaway, Roitberg, Isayev, and **Chodera**.
**preprint:** https://doi.org/10.1101/2020.07.29.227959
**code:** https://github.com/choderalab/qmlify

# HYBRID QUANTUM MACHINE LEARNING / MOLECULAR MECHANICS (QML/MM) POST-PROCESSING CAN IMPROVE ACCURACY



**A**      ML/MM AUGMENTED THERMODYNAMIC CYCLE

# HYBRID QUANTUM MACHINE LEARNING / MOLECULAR MECHANICS (QML/MM) FREE ENERGY CALCULATIONS CUT ERROR IN HALF

**MM** (OPLS2.1 + CM1A-BCC charges)
Missing torsions from LMP2/cc-pVTZ(-f) QM calculations
SPC water

**MM** (OpenFF 1.0.0 "Parsley")
AMBER14SB protein force field
TIP3P; Joung and Cheatham ions

**QML/MM** (OpenFF 1.0.0 + ANI2x)
AMBER14SB protein force field
TIP3P; Joung and Cheatham ions



|  | Tyk2 |
| --- | --- |
| no. of compds | 16 |
| binding affinity range (kcal/mol) | 4.3 |
| crystal structure | 4GIH |
| series ref | 52,53 |
| no. of perturbations | 24 |
| MUE FEP | 0.75 ± 0.11 |
| RMSE FEP | 0.93 ± 0.12 |

Free energies are in units of kilocalories per mole.

Tyk2 benchmark system from Wang et al. JACS 137:2695, 2015
replica-exchange free energy calculations with solute tempering (FEP/REST)

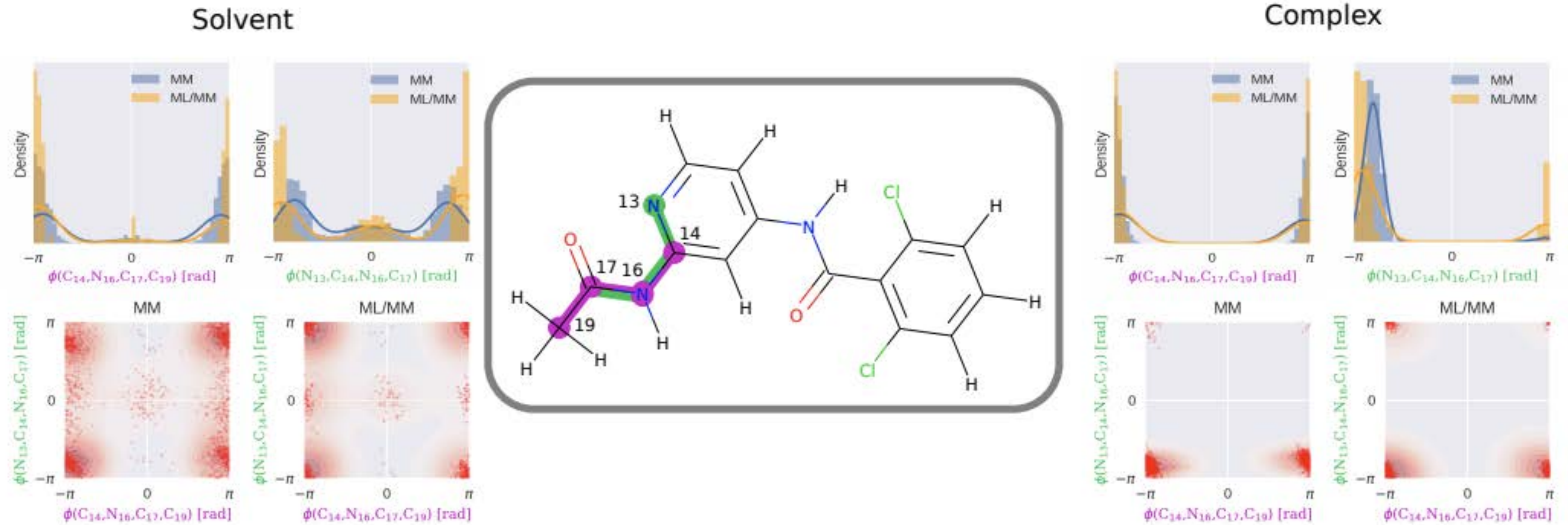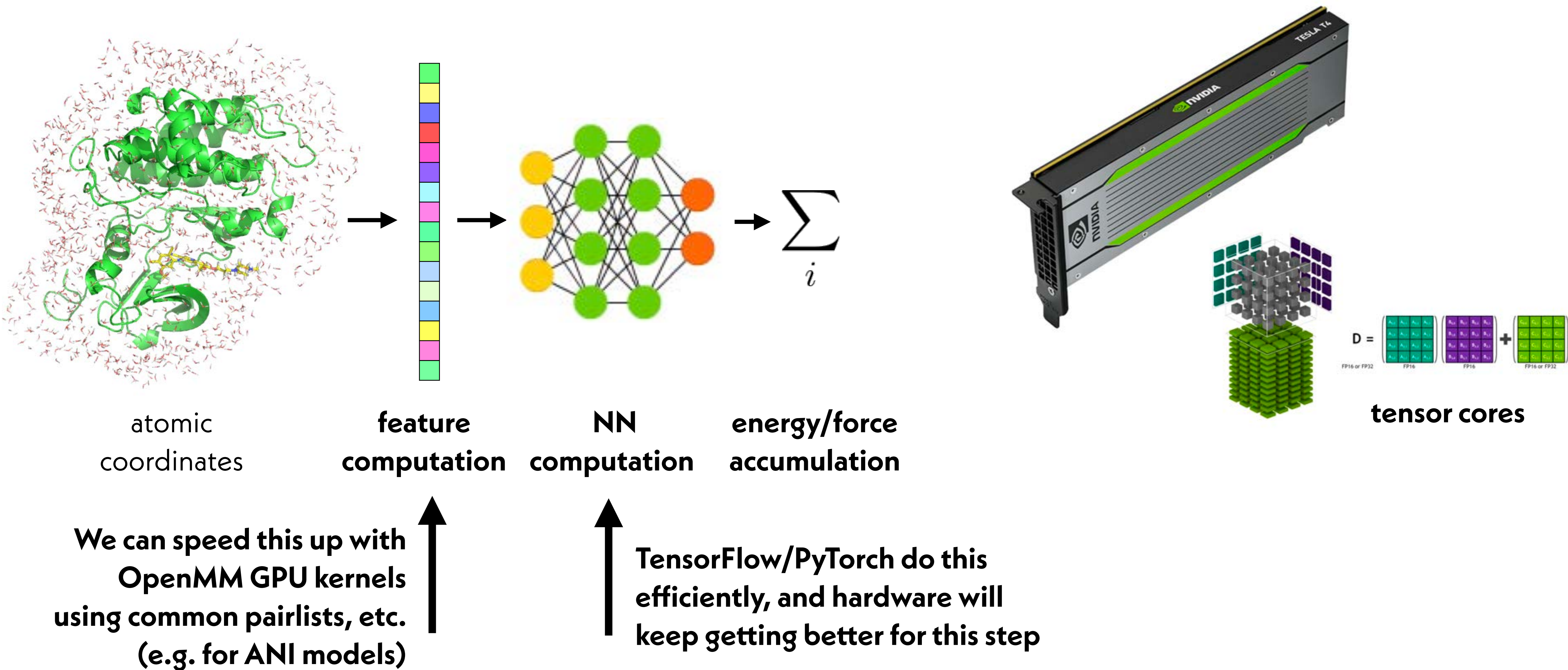replica-exchange free energy calculations with perses
**preprint:** https://doi.org/10.1101/2020.07.29.227959
**code:** https://github.com/choderalab/perses
https://github.com/choderalab/qmlify

# HYBRID QUANTUM MACHINE LEARNING / MOLECULAR MECHANICS (QML/MM) POST-PROCESSING CAN IMPROVE ACCURACY

# COMPUTATIONAL BOTTLENECKS IN CURRENT QML MODELS CAN BE SPED UP WITH CUSTOM GPU KERNELS



atomic coordinates     feature computation     NN computation     energy/force accumulation

                                   tensor cores

We can speed this up with OpenMM GPU kernels using common pairlists, etc. (e.g. for ANI models)

TensorFlow/PyTorch do this efficiently, and hardware will keep getting better for this step

# COMPUTATIONAL BOTTLENECKS IN CURRENT QML MODELS CAN BE SPED UP WITH CUSTOM GPU KERNELS

| PDB ID | # res | # heavy atoms | OpenMM ns/day (4 fs timestep) | TorchANI QML/MM ns/day (2 fs timestep) | OpenMM QML/MM* ns/day (2 fs timestep) |
|--------|-------|---------------|-------------------------------|----------------------------------------|----------------------------------------|
| 3BE9 | 328 | 48 | 436 | 10.4 | 96.5 / 50.8 |
| 2P95 | 286 | 50 | 430 | 7.93 | 96.8 / 49.8 |
| 1HPO | 198 | 64 | 547 | 9.12 | 101 / 44.6 |
| 1AJV | 198 | 75 | 666 | 9.19 | 101 / 40.7 |

* ANI ensemble size:  1 / 8

**NNPOps** library
https://github.com/openmm/nnpops
*   CUDA/CPU accelerated kernels
*   API for inclusion in MD engines
*   Ops wrappers for ML frameworks (PyTorch, TensorFlow, JAX)
*   Community-driven, package agnostic

(~2.5x slower than GPU MD right now, but need 2x smaller timestep) **model distillation** will become important in building single models that are efficient on hardware

# WE WANT TO MAKE IT EASY TO RUN QML/MM SIMULATIONS WITH OPENMM

```python
# Use Amber 14SB and TIP3P-FB for the protein and solvent
forcefield = ForceField('amber14-all.xml', 'amber14/tip3pfb.xml')
# Use OpenFF for the ligand
from openmmforcefields.generators import SMIRNOFFTemplateGenerator
smirnoff = SMIRNOFFTemplateGenerator(molecules=molecules)
# Create an OpenMM MM system
mm_system = forcefield.createSystem(topology)
# Replace ligand intramolecular energetics with ANI-2x
potential = MLPotential('ani2x')
ml_system = potential.createMixedSystem(topology, mm_system, ligand_atoms)
```

https://github.com/openmm/openmm-ml

# PURE QUANTUM MACHINE LEARNING (QML) POTENTIALS CAN BE USED TO COMPUTE FREE ENERGY DIFFERENCES BETWEEN CHEMICAL SPECIES

Potentials are free of singularities, so **simple linear alchemical potentials** can robustly compute alchemical free energies

$$U(x;\lambda) = (1-\lambda)U_{\lambda=0}(x) + \lambda U_{\lambda=1}(x)$$



$U_{\lambda=0}$     $\longleftrightarrow$   $\lambda$     $U_{\lambda=1}$

Simple atomic restraints can be used to improve efficiency by preventing atoms from flying away

**JOSH FASS**    **MARCUS WIEDER**



**ANI-2x**

preprint: https://doi.org/10.1101/2020.10.24.353318
code: https://github.com/choderalab/neutromeratio

# QML POTENTIALS CAN LEARN FROM EXPERIMENTAL DATA TO IMPROVE PHYSICAL MODELS

physical models are data-efficient: retraining on small number of experimental measurements improves accuracy and generalizes well

$\Delta G$

**train:** 221 tautomer pairs

**validate:** 57 tautomer pairs

**test:** 72 tautomer pairs



JOSH FASS   MARCUS WIEDER

**preprint**: https://doi.org/10.1101/2020.10.24.353318
**code**: https://github.com/choderalab/neutromeratio

**OpenMM and the Open Force Field Initiative are working closely with MoISSI to expand the QCArchive to support the construction of next-generation machine learning force fields**

| | | | |
|---|---|---|---|
| SPICE DES Monomers Single Points Dataset v1.1 | 2021-11-15-QMDataset-DES-monomers-single-points | Single point energy calculation of DES monomers. | I, C, Br, P, Cl, H, S, O, F, N |
| SPICE Solvated Amino Acids Single Points Dataset v1.1 | 2021-11-08-QMDataset-Solvated-Amino-Acids-single-points | Single point energy calculation of solvated amino acids. | N, S, O, C, H |
| SPICE DES370K Single Points Dataset v1.0 | 2021-11-08-QMDataset-DES370K-single-points | SPICE single point dataset for ML applications. | 'N', 'O', 'Mg', 'H', 'F', 'K', 'Br', 'Na', 'P', 'Cl', 'I', 'Ca', 'S', 'Li', 'C' |
| SPICE DES370K Single Points Dataset Supplement v1.0 | 2022-02-18-QMDataset-DES370K-single-points-supplement | SPICE single point dataset for ML applications. | F, H, Cl, S, I, Br, N, Li, O, C, Na |
| SPICE Dipeptides Single Points Dataset v1.2 | 2021-11-08-QMDataset-Dipeptide-single-points | SPICE single point dataset for ML applications. | C ,N ,O ,H ,S |
| SPICE PubChem Set 1 Single Points Dataset v1.2 | 2021-11-08-QMDataset-pubchem-set1-single-points | SPICE single point dataset for ML applications. | 'O', 'Cl', 'N', 'C', 'P', 'Br', 'S', 'F', 'I', 'H' |
| SPICE PubChem Set 2 Single Points Dataset v1.2 | 2021-11-09-QMDataset-pubchem-set2-single-points | SPICE single point dataset for ML applications. | 'H', 'P', 'C', 'Cl', 'Br', 'N', 'F', 'S', 'O', 'I' |
| SPICE PubChem Set 3 Single Points Dataset v1.2 | 2021-11-09-QMDataset-pubchem-set3-single-points | SPICE single point dataset for ML applications. | 'N', 'C', 'S', 'Cl', 'Br', 'F', 'P', 'I', 'H', 'O' |
| SPICE PubChem Set 4 Single Points Dataset v1.2 | 2021-11-09-QMDataset-pubchem-set4-single-points | SPICE single point dataset for ML applications. | 'N', 'S', 'Br', 'O', 'C', 'F', 'H', 'I', 'Cl', 'P' |
| SPICE PubChem Set 5 Single Points Dataset v1.2 | 2021-11-09-QMDataset-pubchem-set5-single-points | SPICE single point dataset for ML applications. | 'F', 'H', 'S', 'Br', 'Cl', 'N', 'P', 'C', 'I', 'O' |
| SPICE PubChem Set 6 Single Points Dataset v1.2 | 2021-11-09-QMDataset-pubchem-set6-single-points | SPICE single point dataset for ML applications. | 'Cl', 'O', 'N', 'H', 'C', 'P', 'S', 'F', 'Br', 'I' |

http://qcarchive.molssi.org

https://github.com/openmm/spice-dataset

# CAN WE CHANGE PRACTICE IN STRUCTURE-ENABLED DRUG DISCOVERY BY LEVERAGING DATA WE GENERATE?



week 1

week 2

**2021**

| MON | TUE | WED | THU | FRI | SAT | SUN |
|-----|-----|-----|-----|-----|-----|-----|
| designs/ predictions | synthesis | | | new data | | |

| MON | TUE | WED | THU | FRI | SAT | SUN |
|-----|-----|-----|-----|-----|-----|-----|
| designs/ predictions | synthesis | | | new data | | |

using published force field model

using the same published force field model!
we haven't learned anything from the data

week 1

week 2

**2025**

| MON | TUE | WED | THU | FRI | SAT | SUN |
|-----|-----|-----|-----|-----|-----|-----|
| designs/ predictions 1.0 | synthesis | | | new data | build model 2.0! | |

| MON | TUE | WED | THU | FRI | SAT | SUN |
|-----|-----|-----|-----|-----|-----|-----|
| designs/ predictions 2.0 | synthesis | | | | | |

using force field model
built from public + private data

using new model tuned to target
from first week's data

# PREPRINTS AND CODE

**gimlet:** graph convolutional networks for partial charge assignment
**preprint:** https://arxiv.org/abs/1909.07903
**code**: http://github.com/choderalab/gimlet

**espaloma:** end-to-end differentiable assignment of force field parameters
**preprint**: https://arxiv.org/abs/2010.01196
**code**: https://github.com/choderalab/espaloma

**qmlify:** hybrid QML/MM alchemical free energy calculations for protein-ligand binding
**preprint:** https://doi.org/10.1101/2020.07.29.227959
**code:** https://github.com/choderalab/qmlify

**neutromeratio:** alchemical free energy calculations with fully QML potentials for tautomer ratio prediction
**preprint:** https://doi.org/10.1101/2020.10.24.353318
**code:** https://github.com/choderalab/neutromeratio

# CHODERA LAB

# MM WILL MOVE TOWARD POTENTIALS THAT BLEND SHORT-RANGE ML AND LONG-RANGE PHYSICS

## PhysNet

## 4D-HGNNP



$$E = \sum_{i=1}^{N} E_i + \sum_{i=1}^{N} \sum_{j>i}^{N} k_e \frac{q_i q_j}{r_{ij}}$$

Energy

$$\text{Forces} \quad F_i = -\frac{\partial E}{\partial r_i}$$

$$\text{Dipole} \quad p = \sum_{i=1}^{N} q_i r_i$$

MD codes need to interoperate with ML frameworks and implement optimized ML potentials using common atomic featurizations

# ALCHEMICAL FREE ENERGY CALCULATIONS CAN PREDICT SELECTIVITIES BETTER THAN AFFINITIES

# HOW WELL CAN WE PREDICT SELECTIVITY?



CDK9/cyclin T
(4BCI)

CDK2/cyclin A
(4BCK)

inhibition reinstates apoptosis in cancer cells

essential for S-phase progression

**STEVEN ALBANESE**

# HOW MUCH DOES CANCELLATION OF ERROR HELP SELECTIVITY PREDICTION?



CDK2

ERK2

- 🟠 Charged (negative)
- 🔵 Charged (positive)
- ⚪ Glycine
- 🟢 Hydrophobic
- 🔵 Polar
- —• Pi-cation
- → H-bond
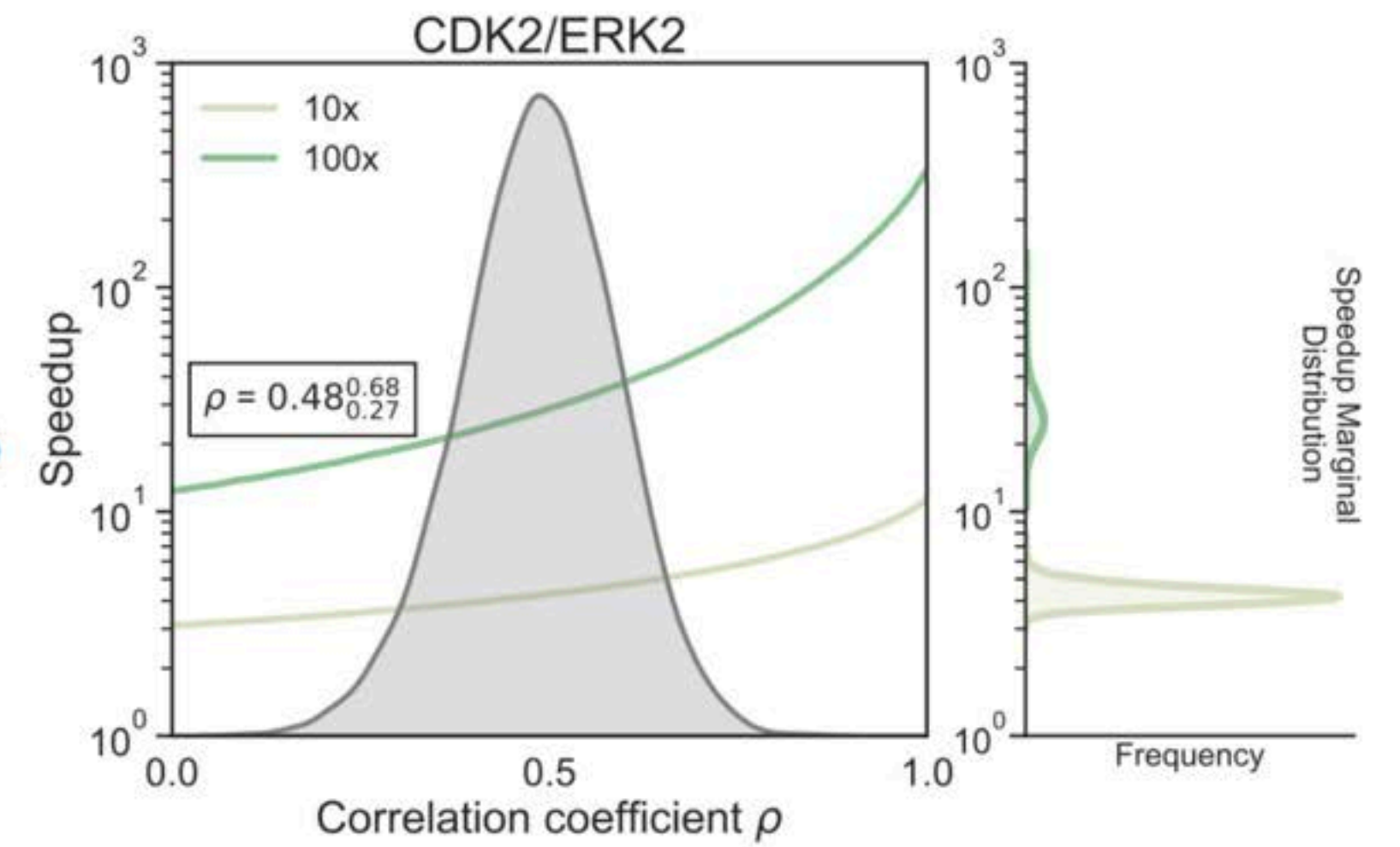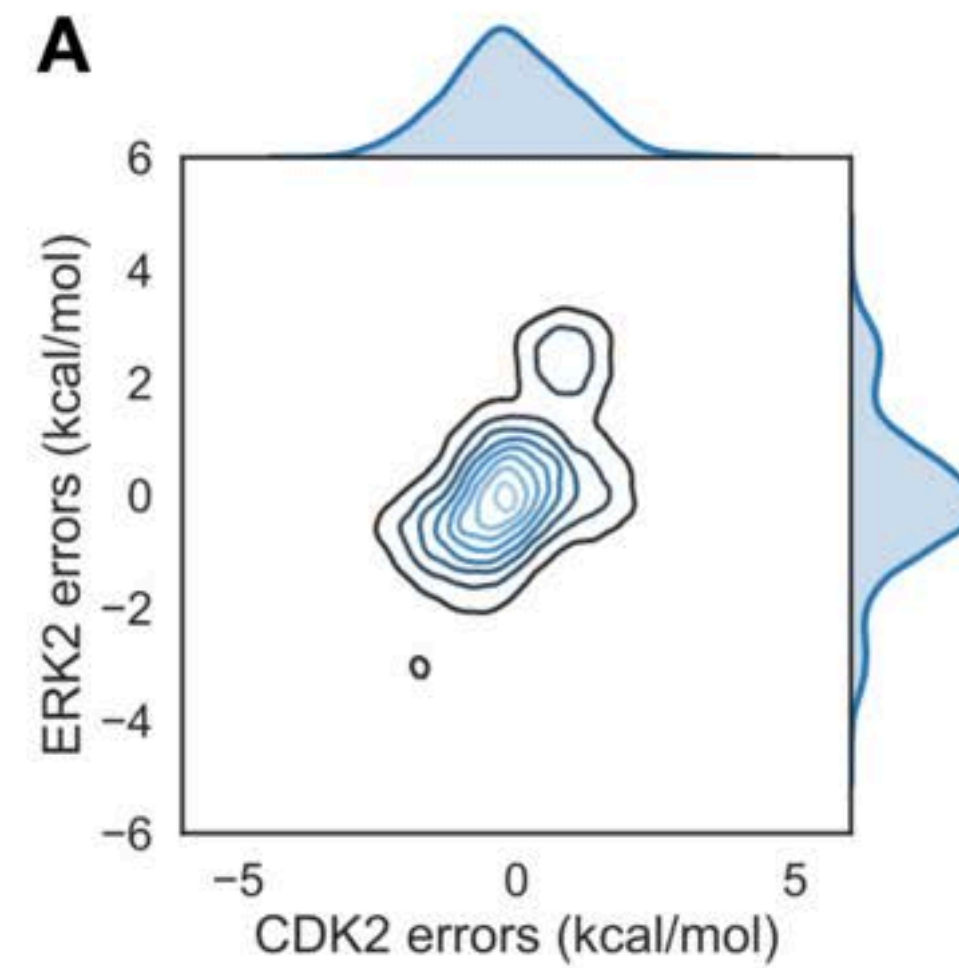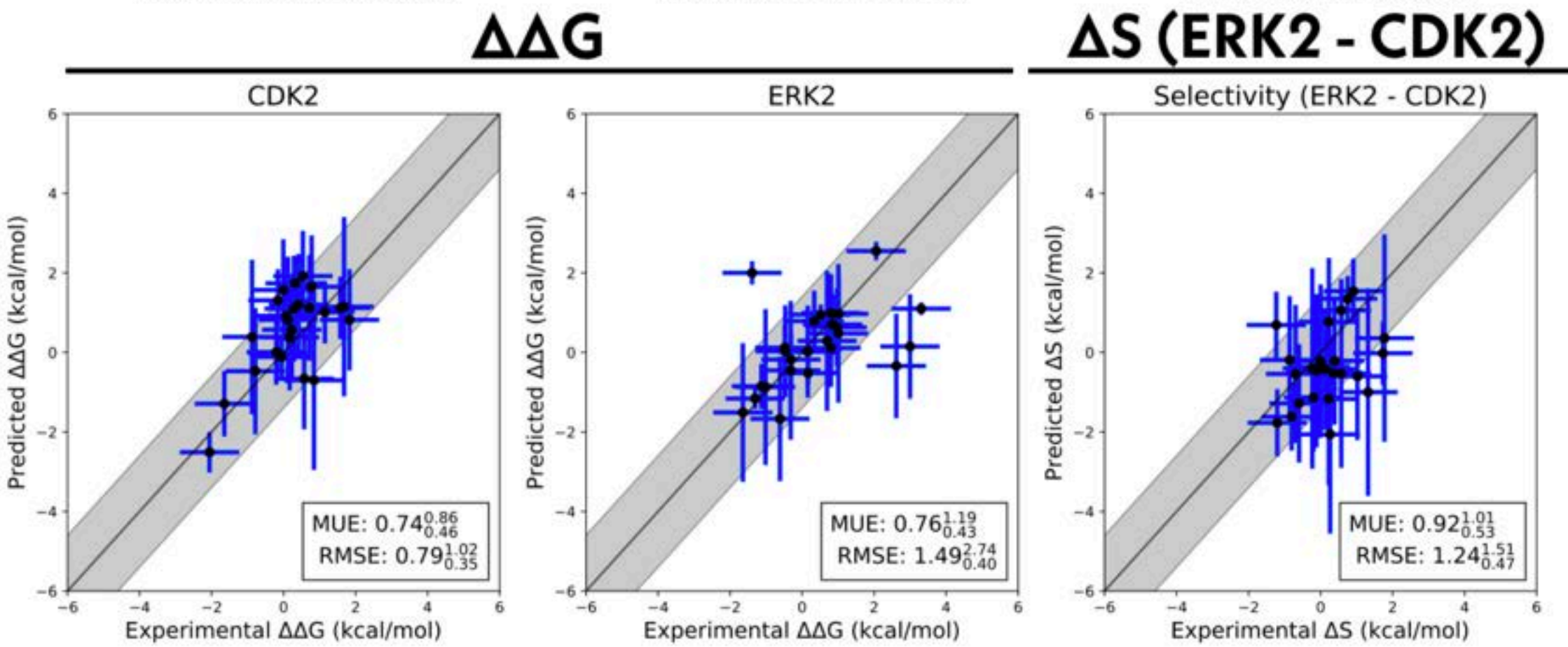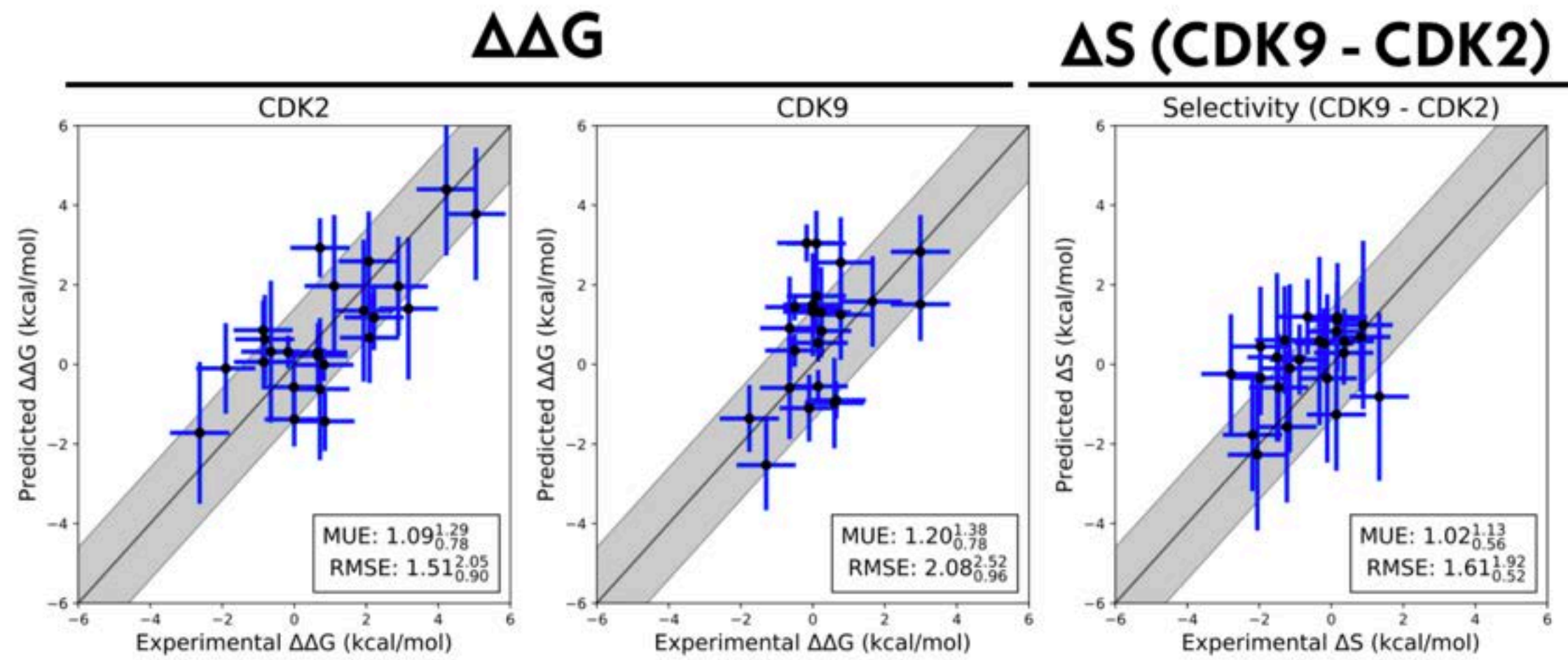- ⚪ Solvent exposure

Quantify via the **correlation coefficient**

$$\rho \equiv \frac{\mathrm{cov}(\epsilon_1, \epsilon_2)}{\sqrt{\mathrm{var}(\epsilon_1)\mathrm{var}(\epsilon_2)}}$$

of the **error**

$$\epsilon_* \equiv \Delta\Delta G_*^{\mathrm{FEP}} - \Delta\Delta G_*^{\mathrm{exp}}$$

# DIFFERENT SELECTIVITY PROBLEMS SHOW DIFFERENT DEGREES OF CANCELLATION
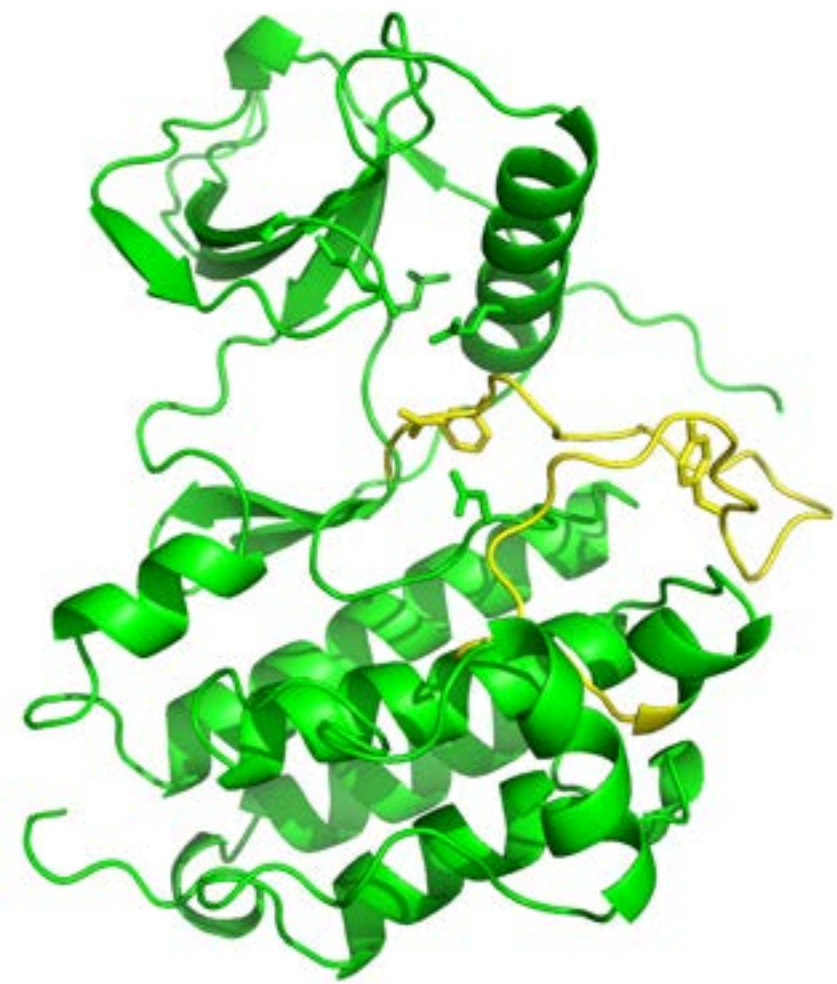


**FEP+/OPLS3**

**STEVEN ALBANESE**

**LINGLE WANG**
**SCHRÖDINGER**

# INTERLINE WILL PURSUE A NUMBER OF SELECTIVITY-FOCUSED DESIGN PROBLEMS



target
(promotes downstream activity)
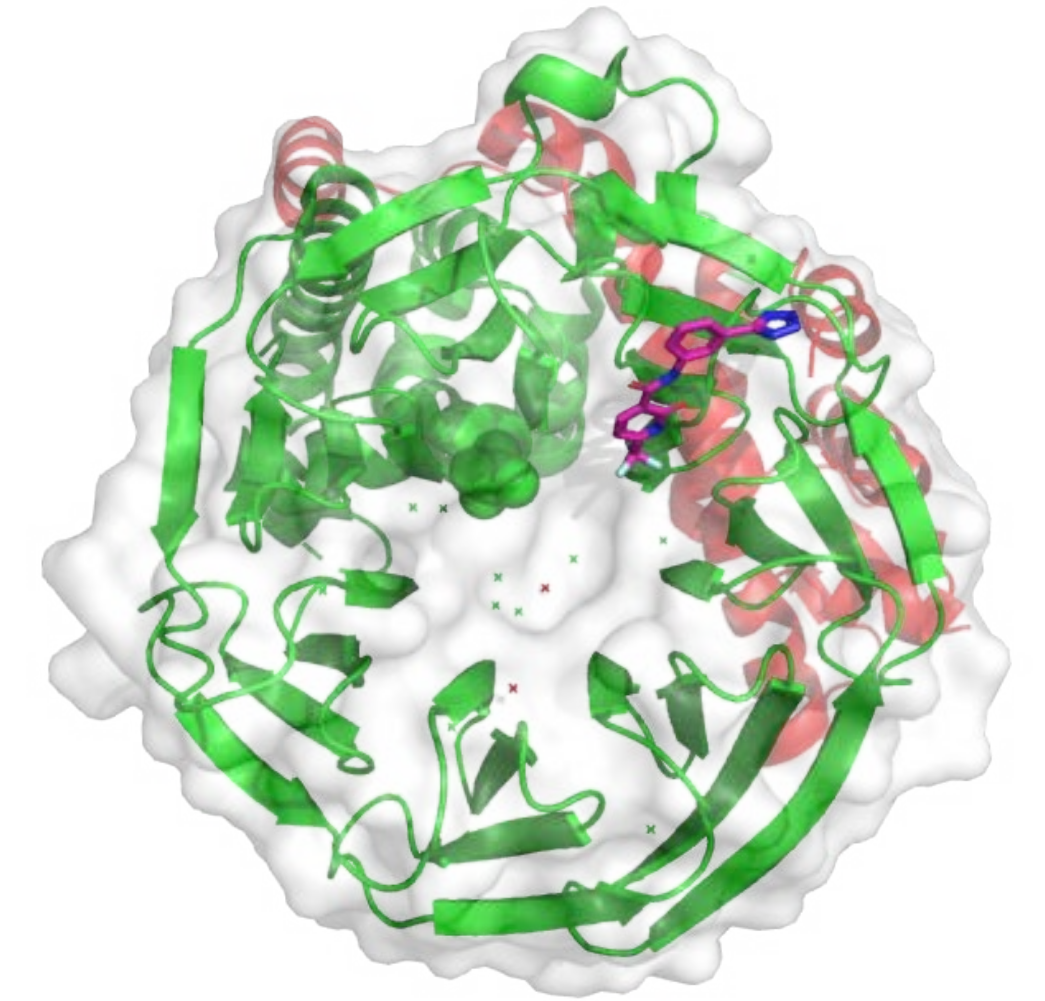
antitarget
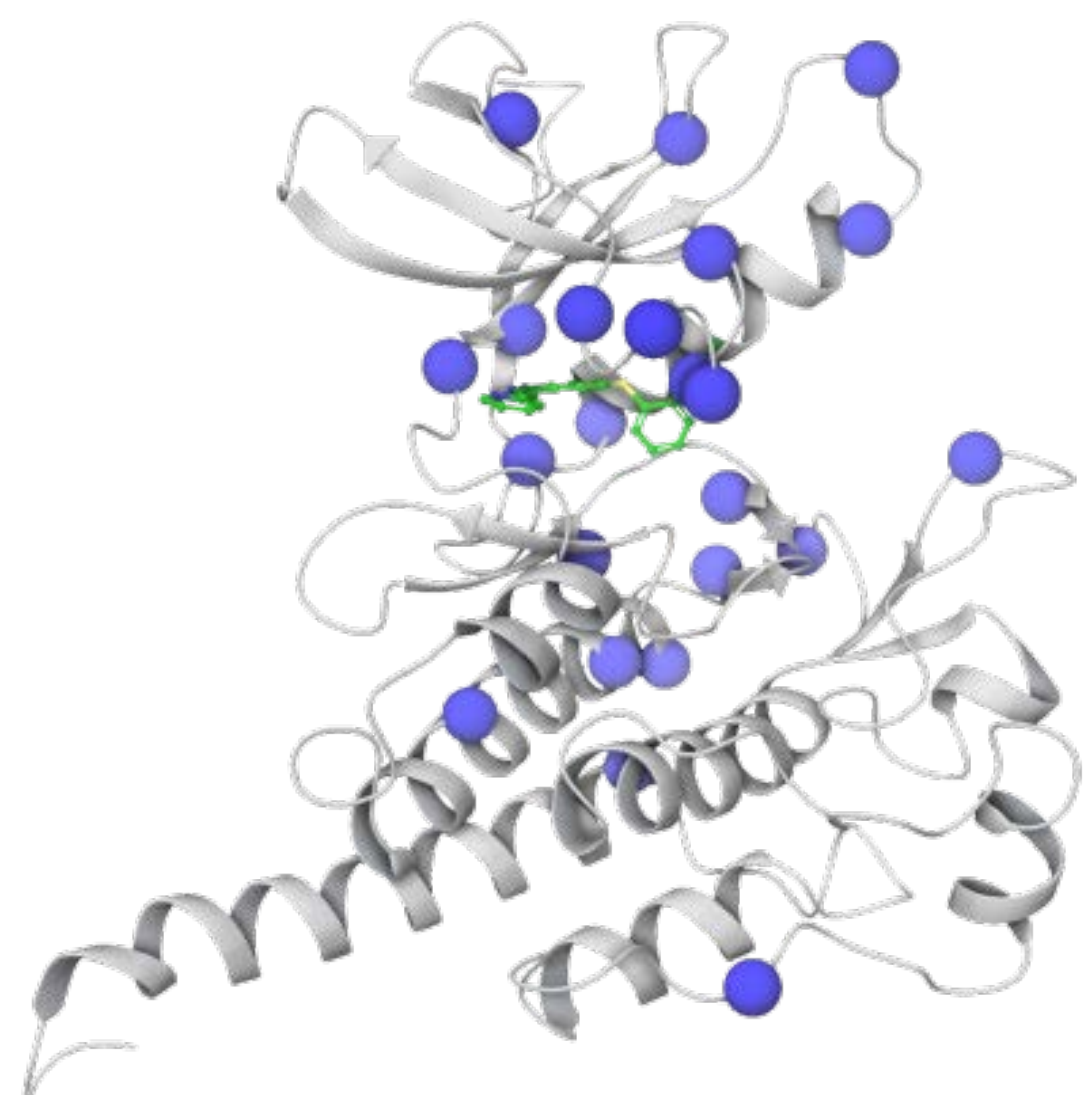(inhibits downstream activity)

**selective (de)stabilization
of target conformations**

target
(complex to be stabilized)
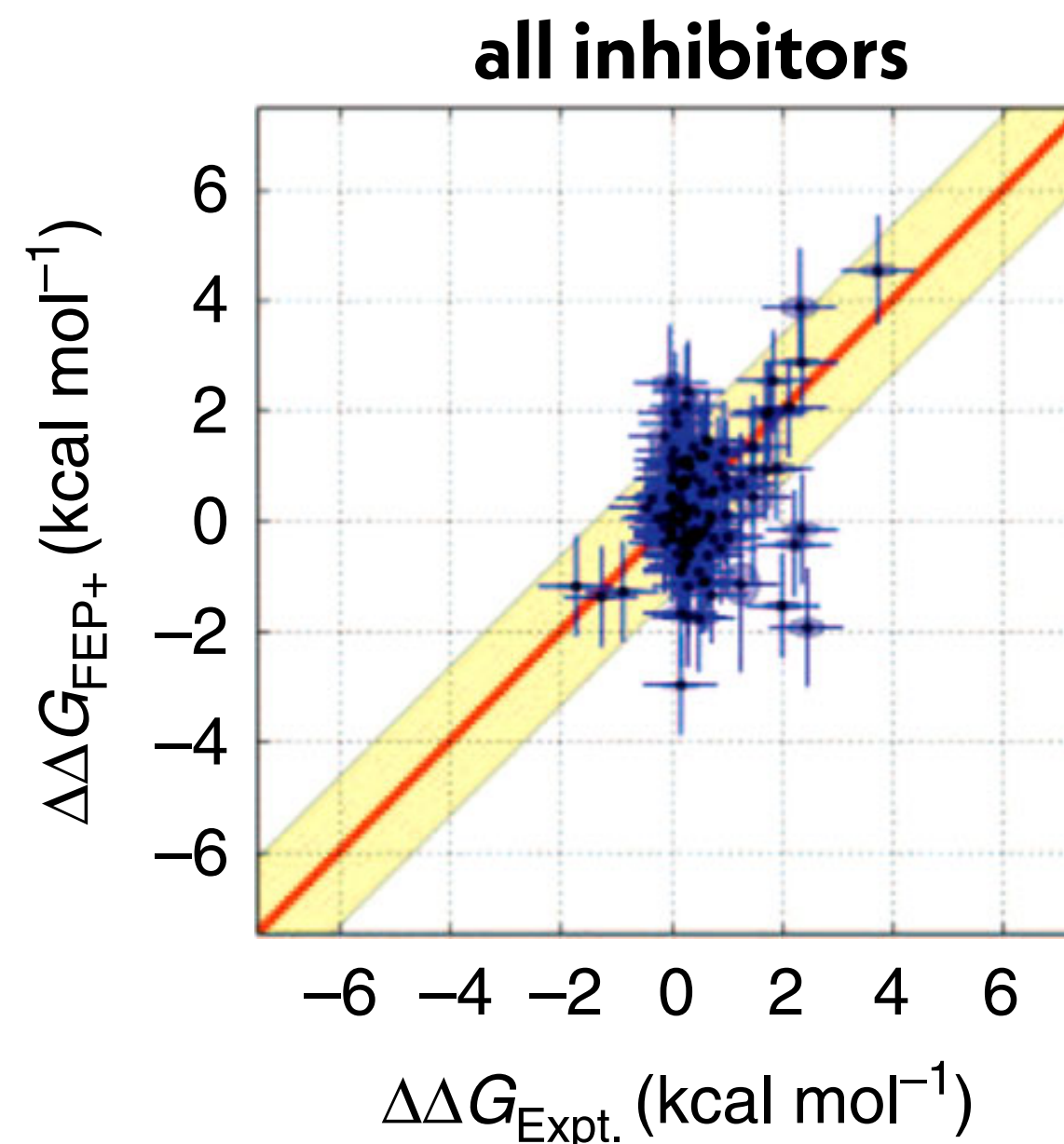
antitarget
(individual binding partners)

**selective (de)stabilization
of complexes**

# ALCHEMICAL FREE ENERGY CALCULATIONS CAN PREDICT THE IMPACT OF MUTATIONS ON LIGAND BINDING OR PROTEIN-PROTEIN INTERACTIONS



**all inhibitors**

| TKI | $N_{mut}$ | R | S |
|---|---|---|---|
| Axitinib | 26 | 0 | 26 |
| Bosutinib | 21 | 4 | 17 |
| Dasatinib | 21 | 5 | 16 |
| Imatinib | 21 | 5 | 16 |
| Nilotinib | 21 | 4 | 17 |
| Ponatinib | 21 | 0 | 21 |
| Subtotal | 131 | 18 | 113 |
| Erlotinib | 7 | 1 | 6 |
| Gefitinib | 6 | 0 | 6 |
| Total | 144 | 19 | 125 |

$N_{mut}$ Total number of mutants for which $\Delta pIC_{50}$ data was available
Number of **R**esistant, **S**usceptible mutants using 10-fold affinity change threshold

| RMSE (kcal mol$^{-1}$) | 0.99 | 1.15 / 0.85 |
|---|---|---|

**Prediction**

| | | S | r |
|---|---|---|---|
| **Experiment** | S | 105 | 8 |
| | r | 9 | 9 |

| Accuracy | 0.89 | 0.92 / 0.86 |
|---|---|---|
| Specificity | 0.91 | 0.94 / 0.89 |
| Sensitivity | 0.69 | 1.00 / 0.46 |

Hauser, Negron, Albanese, Ray, Steinbrecher, Abel, **Chodera**, Wang. Co____ic___ Biolog_

**KEVIN HAUSER**
**SCHRÖDINGER (NOW AT RUBRYC)**