# ALCHEMICAL ACADEMY:
# TEACHING FREE ENERGY CALCULATIONS TO LEARN

**John D. Chodera**

MSKCC Computational and Systems Biology Program

http://choderalab.org

![Memorial Sloan Kettering Cancer Center 1884]

# Sloan-Kettering Institute

In more than 100 laboratories, our scientists are conducting innovative research to advance understanding in the biological sciences and improve human health.

Cornell Weill

MSKCC

Rockefeller

Dana Pe'er

Quaid Morris

Christina Leslie

Joao Xavier

Kushal Dey

John Chodera

Thomas Norman

csbio@MSKCC

# CHODERA LAB

## HOW CAN COMPUTATIONAL BIOPHYSICS AND MACHINE LEARNING ADVANCE DISCOVERY AND TREATMENT IN THE ERA OF CANCER GENOMICS?
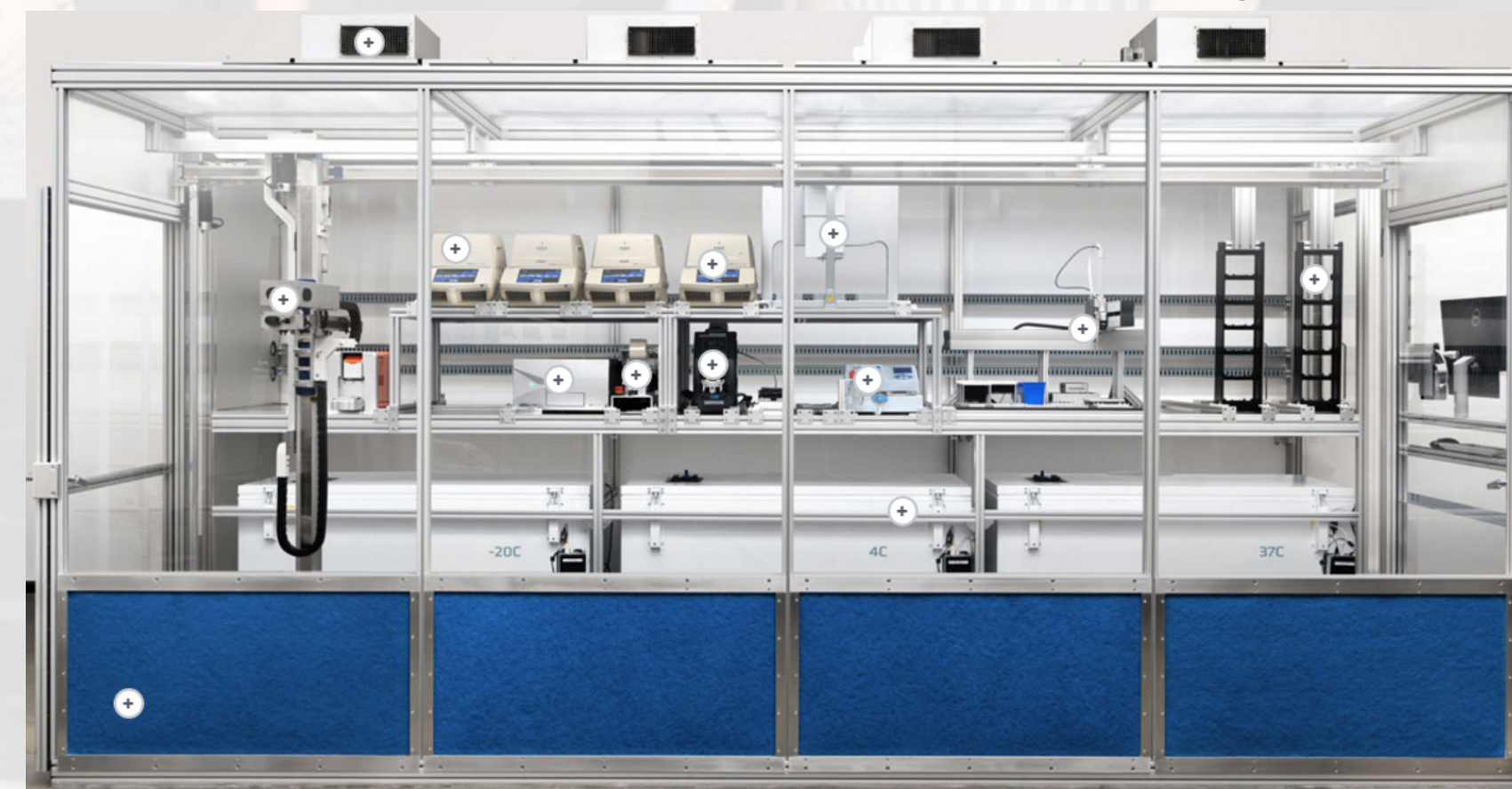
**MODELING**

**AUTOMATION**

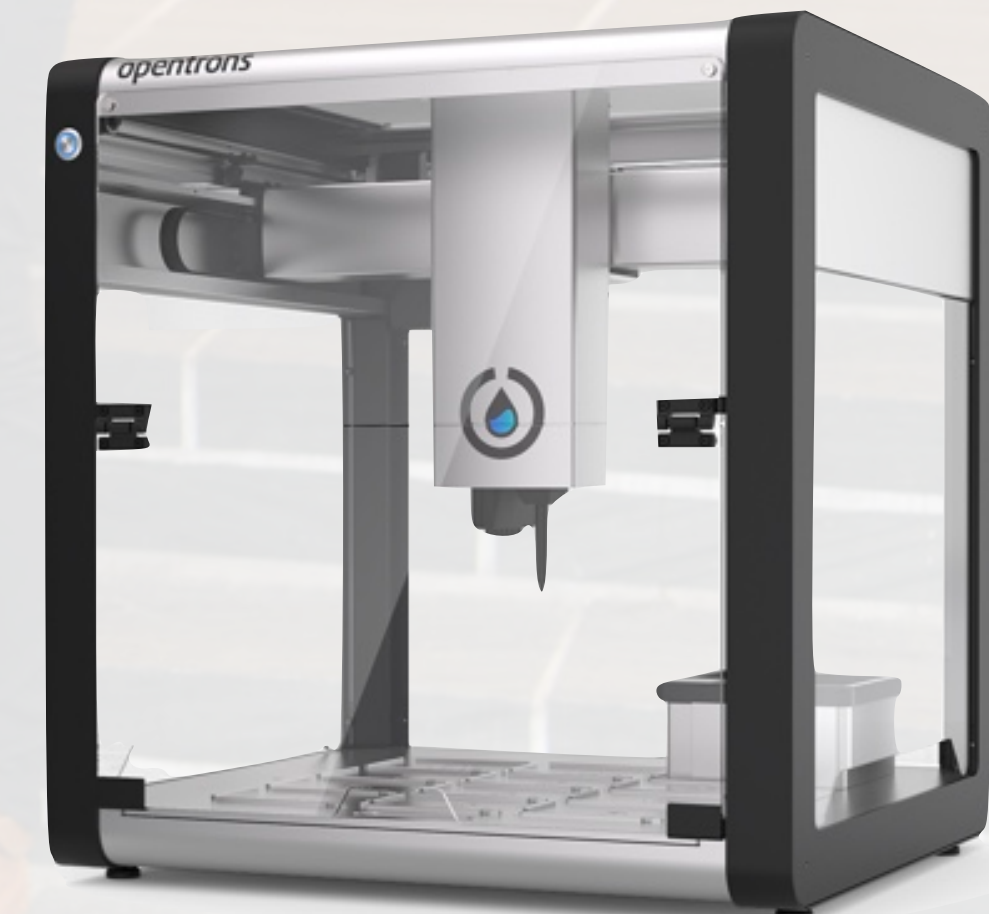FOLDING @HOME

amazon web services™ | EC2

CHODERA LAB, Z17

$$V(\mathbf{q}) = \sum_{\text{bonds}} K_r (r - r_{eq})^2 + \sum_{\text{angles}} K_\theta (\theta - \theta_{eq})^2$$

$$+ \sum_{\text{dihedrals}} \frac{V_n}{2} [1 + \cos(n\phi - \gamma)] + \sum_{i<j} \left[ \frac{A_{ij}}{R_{ij}^{12}} - \frac{B_{ij}}{R_{ij}^{6}} + \frac{q_i q_j}{\epsilon R_{ij}} \right]$$
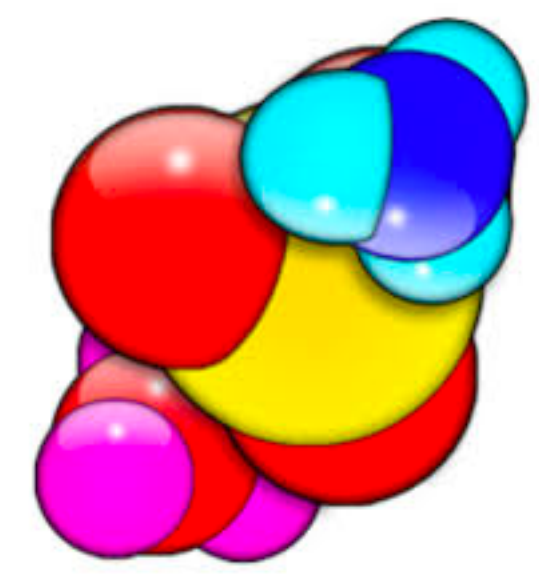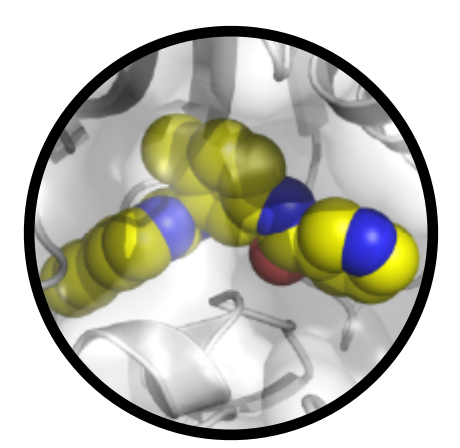
STRATEOS CLOUD WETLAB

OPENTRONS

# WE COLLABORATE BROADLY TO ADVANCE THE STATE OF DRUG DISCOVERY



open source software development initiatives

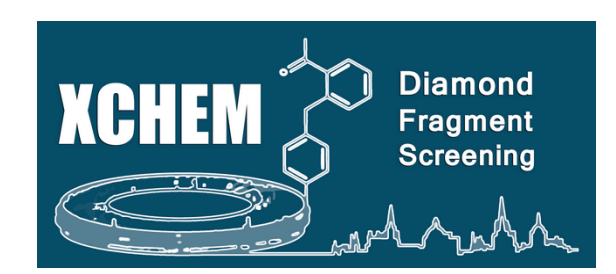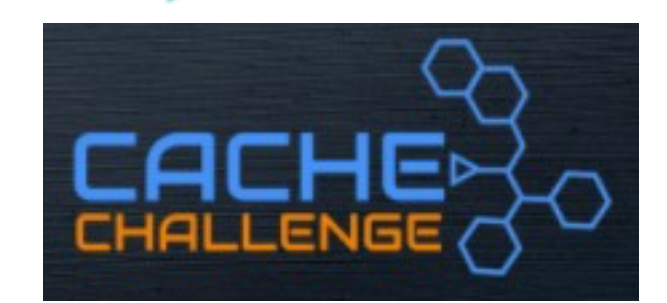data generators, community challenges, and resources

industry collaborations

choderalab
(algorithms and open source software)

academia

IP-generating collaborations

open science / open source software

# OpenMM

A high performance toolkit for molecular simulation. Use it as a library, or as an application. We include extensive language bindings for Python, C, C++, and even Fortran. The code is open source and actively maintained on Github, licensed under MIT and LGPL. Part of the Omnia suite of tools for predictive biomolecular simulation.

ABOUT  FORUM  GITHUB

# Extreme Flexibility. Extreme Speed.

Extreme flexibility through custom forces and integrators. Extreme performance through GPU Acceleration, with optimizations for AMD, NVIDIA, and Intel Integrated GPUs. It's fast on CPUs too. See the benchmarks.

## Install

Install using the conda Python package manager that powers the Omnia ecosystem.

## Docs

For more information about the science, the code base, and the API behind OpenMM.

## Support

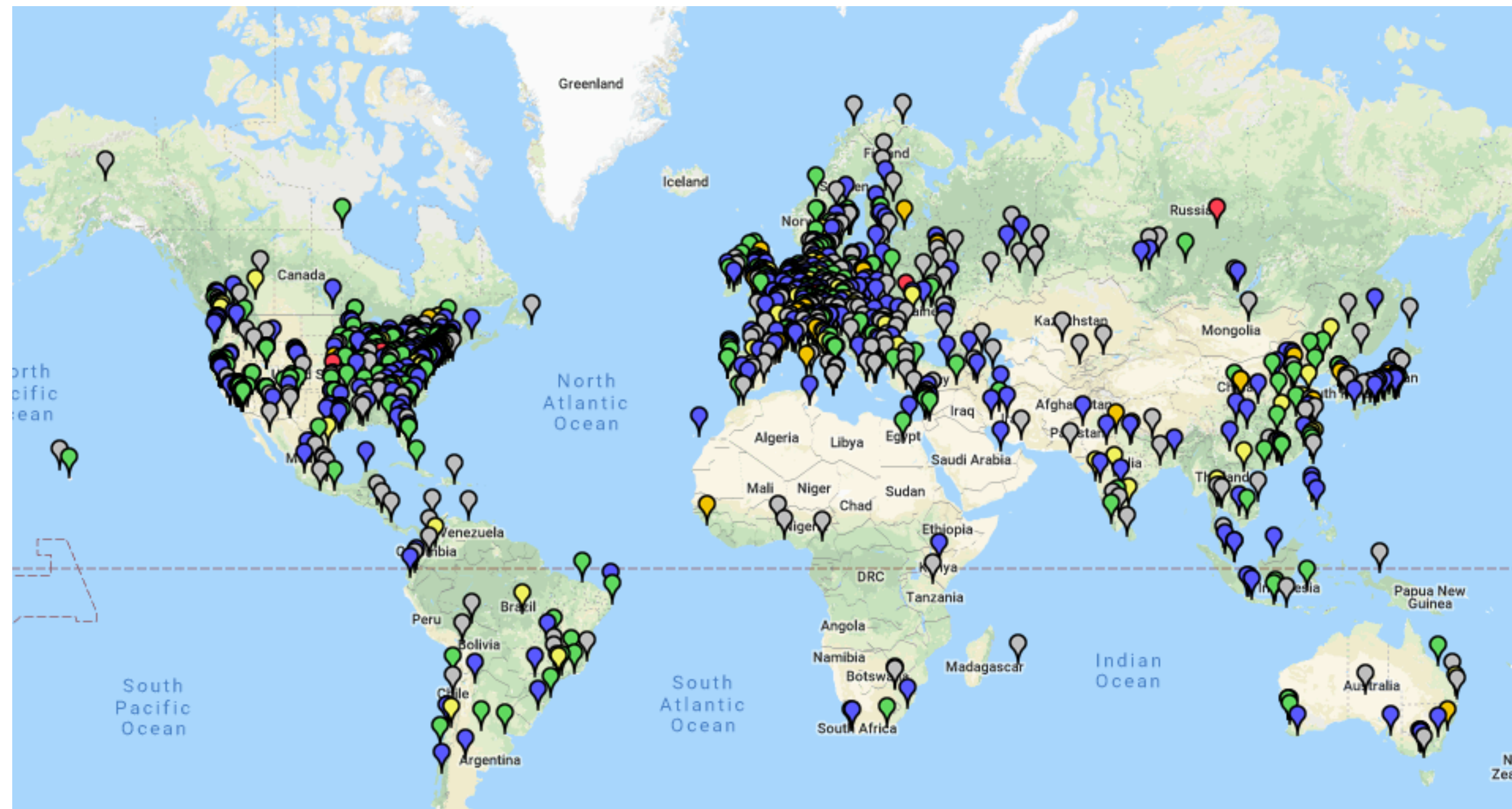For more information about filing bug reports, requesting new features, and other issues.

## Resources

Explore additional libraries and third-party tools built around OpenMM.

## Tutorials

Get started right away with OpenMM tutorials.

http://openmm.org/

# OPENMM ARCHITECTURE MAKES DEVELOPMENT SIMPLE



OpenMM also has bindings for **C++**, **C**, and **FORTRAN**

# OPENMM IS USED BY RESEARCHERS ALL OVER THE WORLD



Geographic statistics from http://simtk.org

**OpenMM**
http://openmm.org

downloads | 1M total

**OpenMMTools**
http://github.com/choderalab/openmmtools
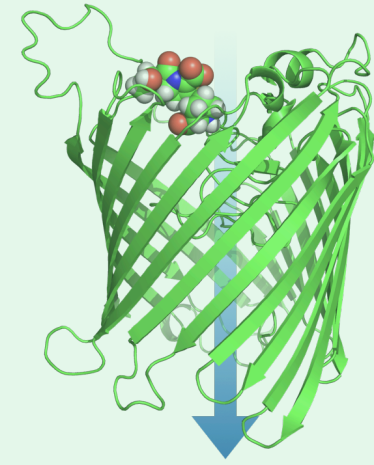
downloads | 174k total

# OPENMM CAN BE USED AS A LIBRARY TO ENABLE APPLICATIONS TO INTEGRATE PHYSICAL MODELING



**isolde**  **perses**  **iapetus**

targeted domain-specific applications
(Python, C++, C, or Fortran)

**APPLICATIONS**

**openmmtools**

high-level simulation algorithms, alchemical tools
(Python to enable rapid development)

**ALGORITHMS**

OpenMM

general GPU-accelerated MD simulation engine
(C++/CUDA/OpenCL with Python API)

**CORE**

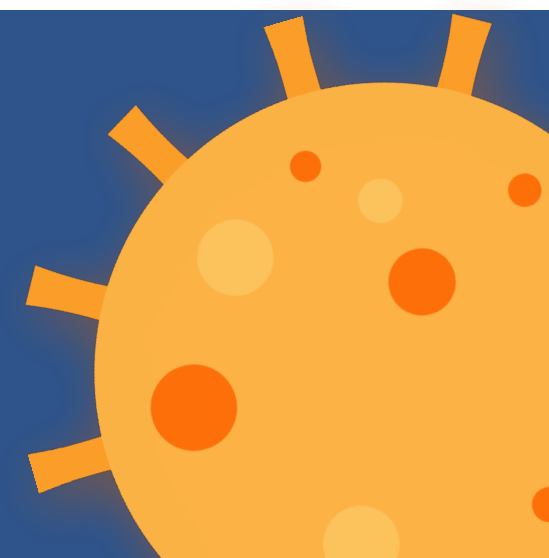# DESIGNING REAL PRECLINICAL DRUG CANDIDATES IS CHALLENGING

**Ed Griffen**

Medchemica

## Target Candidate Profile (TCP) for oral SARS-CoV-2 main viral protease (Mpro) inhibitor
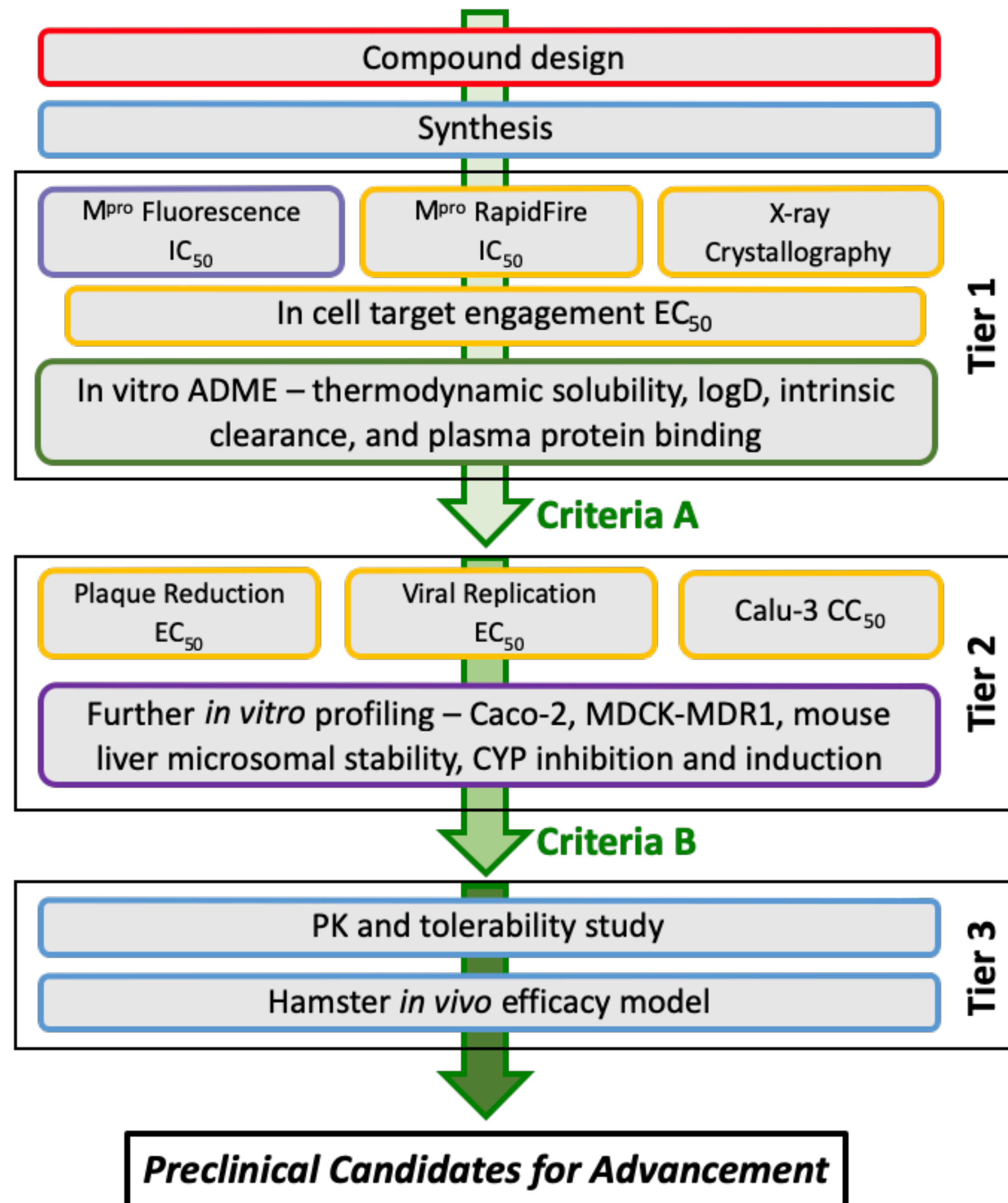
| Property | Target range | Rationale |
|---|---|---|
| protease assay | $IC_{50} < 10$ nM | Extrapolation from other anti-viral programs |
| viral replication assay | $EC_{50} < 5$ μM | Suppression of virus at achievable blood levels |
| plaque reduction assay | $EC_{50} < 5$ μM | Suppression of virus at achievable blood levels |
| route of administration | oral | bid/tid - compromise PK for potency if pharmacodynamic effect achieved |
| solubility | > 5 mg/mL | Aim for biopharmaceutical class 1 assuming <= 750 mg dose |
| half-life | > 8 h (human) est from rat and dog | Assume PK/PD requires continuous cover over plaque inhibition for 24 h max bid dosing |
| safety | Only reversible and monitorable toxicities<br>No significant DDI - clean in 5 CYP450 isoforms<br>hERG and NaV1.5 $IC_{50} > 50$ μM<br>No significant change in QTc<br>Ames negative<br>No mutagenicity or teratogenicity risk | No significant toxicological delays to development<br>DDI aims to deal with co-morbidities / therapies,<br>cardiac safety for COVID-19 risk profile<br>cardiac safety for COVID-19 risk profile<br>Low carcinogenicity risk reduces delays in manufacturing<br>Patient group will include significant proportion of women of childbearing age |

COVID Moonshot

An international effort to
DISCOVER A COVID ANTIVIRAL

**https://covid.postera.ai/covid**

# TO GET THERE, DRUG DESIGN INVOLVES MAKING A LOT OF DECISIONS ABOUT WHICH MOLECULES WILL ACHIEVE CERTAIN OBJECTIVES



**assay purpose**

Does it inhibit the target? How does it bind?

Does it work in cells?

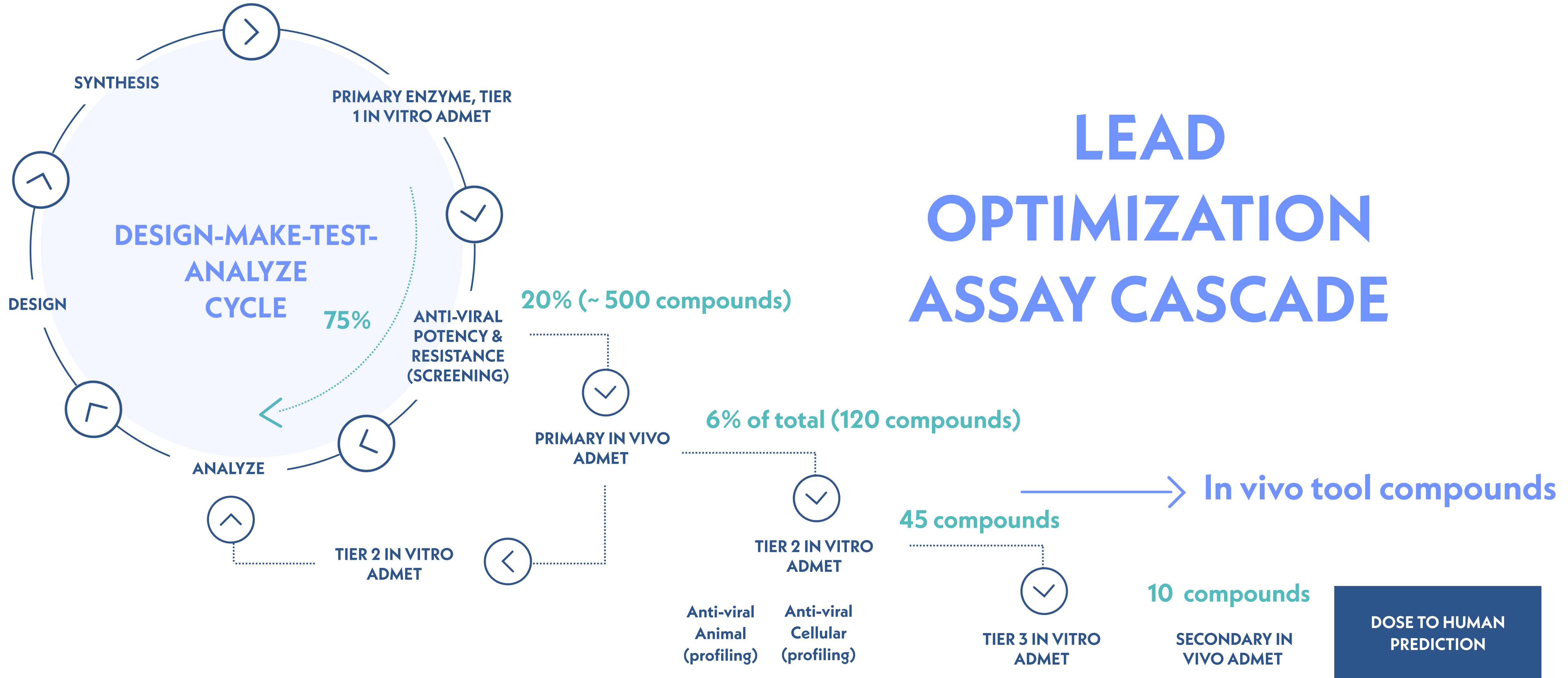Does it have a chance of working in humans?

Does it kill the virus in cells?

Could it cause bad side effects?
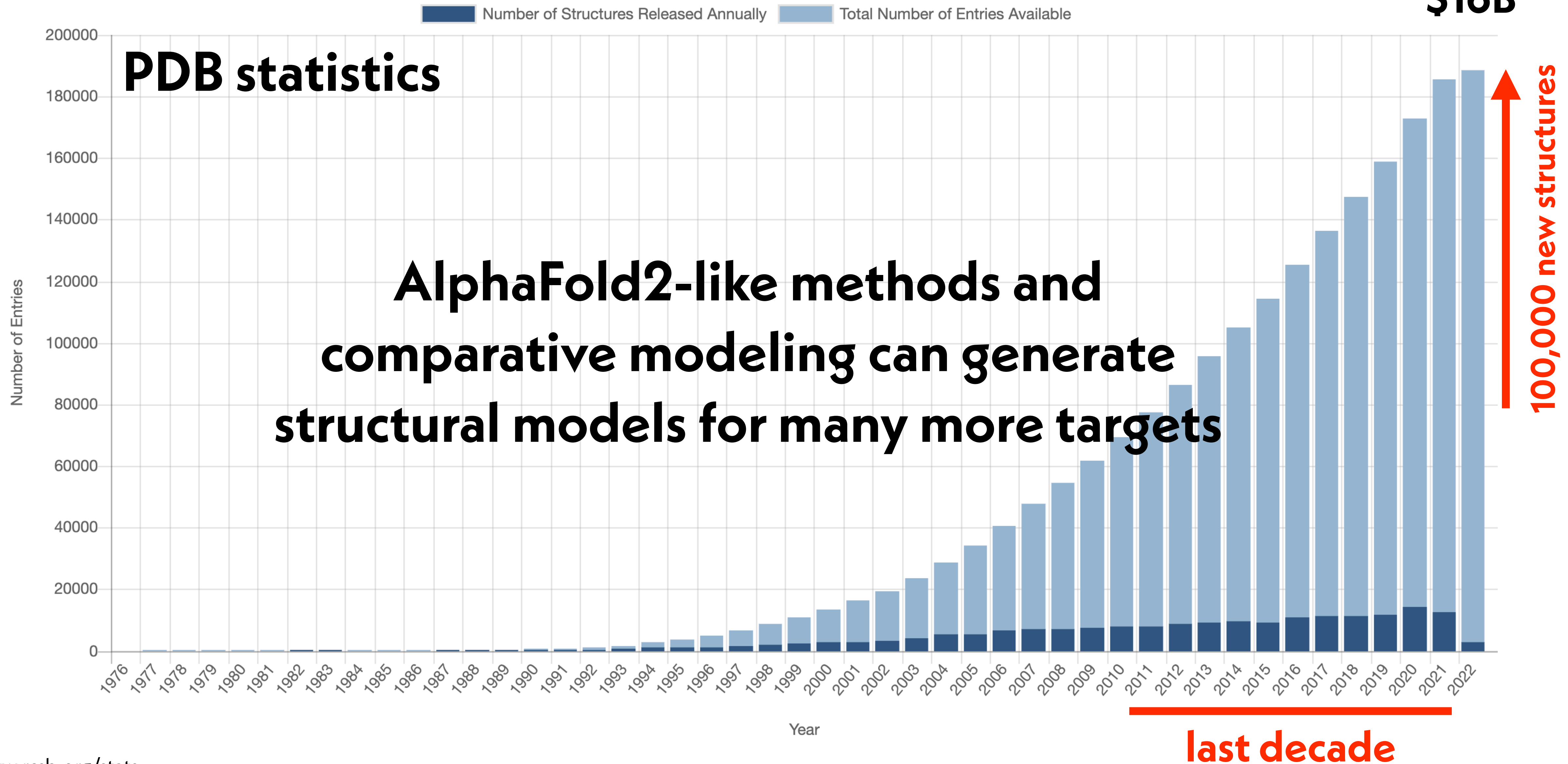
Can oral dosing deliver sufficient drug?

Does it actually work against the disease?

# MUCH OF THE TIME IS SPENT IN PREDICTING COMPOUNDS THAT WILL IMPROVE OR MAINTAIN POTENCY
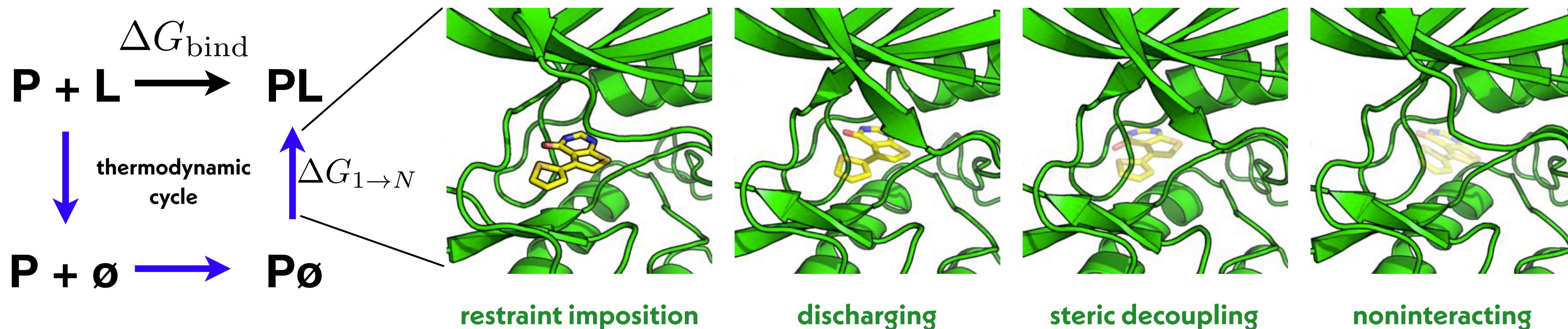


**LEAD OPTIMIZATION ASSAY CASCADE**

SYNTHESIS

PRIMARY ENZYME, TIER 1 IN VITRO ADMET

DESIGN-MAKE-TEST-ANALYZE CYCLE

DESIGN

75%

ANTI-VIRAL POTENCY & RESISTANCE (SCREENING)

ANALYZE

20% (~ 500 compounds)

PRIMARY IN VIVO ADMET

6% of total (120 compounds)

TIER 2 IN VITRO ADMET

In vivo tool compounds

45 compounds

TIER 2 IN VITRO ADMET

Anti-viral Animal (profiling)

Anti-viral Cellular (profiling)

TIER 3 IN VITRO ADMET

10 compounds

SECONDARY IN VIVO ADMET

DOSE TO HUMAN PREDICTION

# STRUCTURAL DATA IS NOW AN ABUNDANT RESOURCE FOR DRUG DISCOVERY



$16B

PDB statistics

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

100,000 new structures

last decade

http://www.rcsb.org/stats

# 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 + ø $\longrightarrow$ Pø

restraint imposition    discharging    steric decoupling    noninteracting

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

$$\Delta G_{0 \to 1} = -k_B T \ln \frac{Z_1}{Z_0} = -k_B T \ln \frac{Z_{\lambda_2}}{Z_{\lambda_1}} \frac{Z_{\lambda_3}}{Z_{\lambda_2}} \cdots \frac{Z_{\lambda_N}}{Z_{\lambda_{N-1}}}$$

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

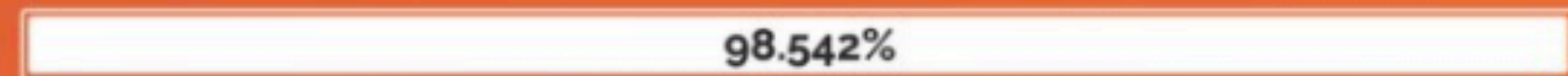# WE'VE RUN LOTS OF FREE ENERGY CALCULATIONS



Folding@home
@foldingathome

Replying to @foldingathome @covid_moonshot and @EnamineLtd

The first @covid_moonshot sprint was a huge success!
Your GPUs worked through 2,353,512 work units of small molecules binding to the #COVID19 main protease.
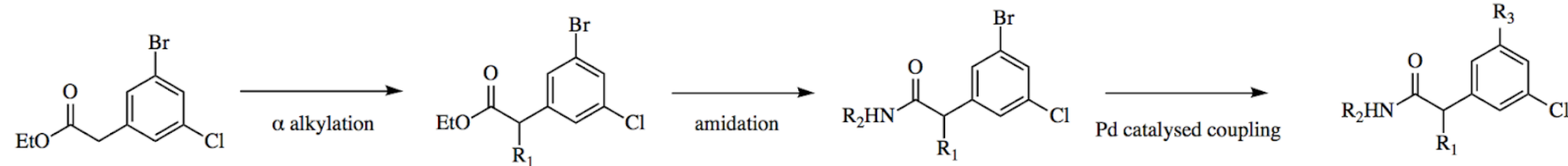That's nearly 10 milliseconds of simulation time!

Progress on the current Sprint 1 to evaluate a batch of potential drugs Started Sun Jul 26 06:31:13 UTC 2020
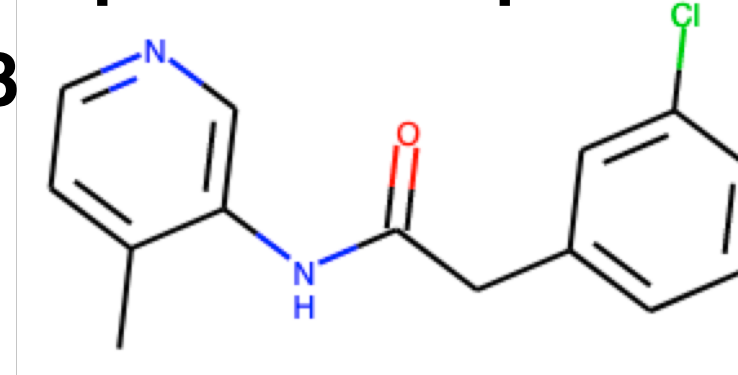
98.542%

8:52 AM · Aug 17, 2020 · TweetDeck

# WE CAN LEVERAGE STRUCTURE TO MAKE DECISIONS BETWEEN MANY RELATED SYNTHETICALLY FEASIBLE ANALOGUES



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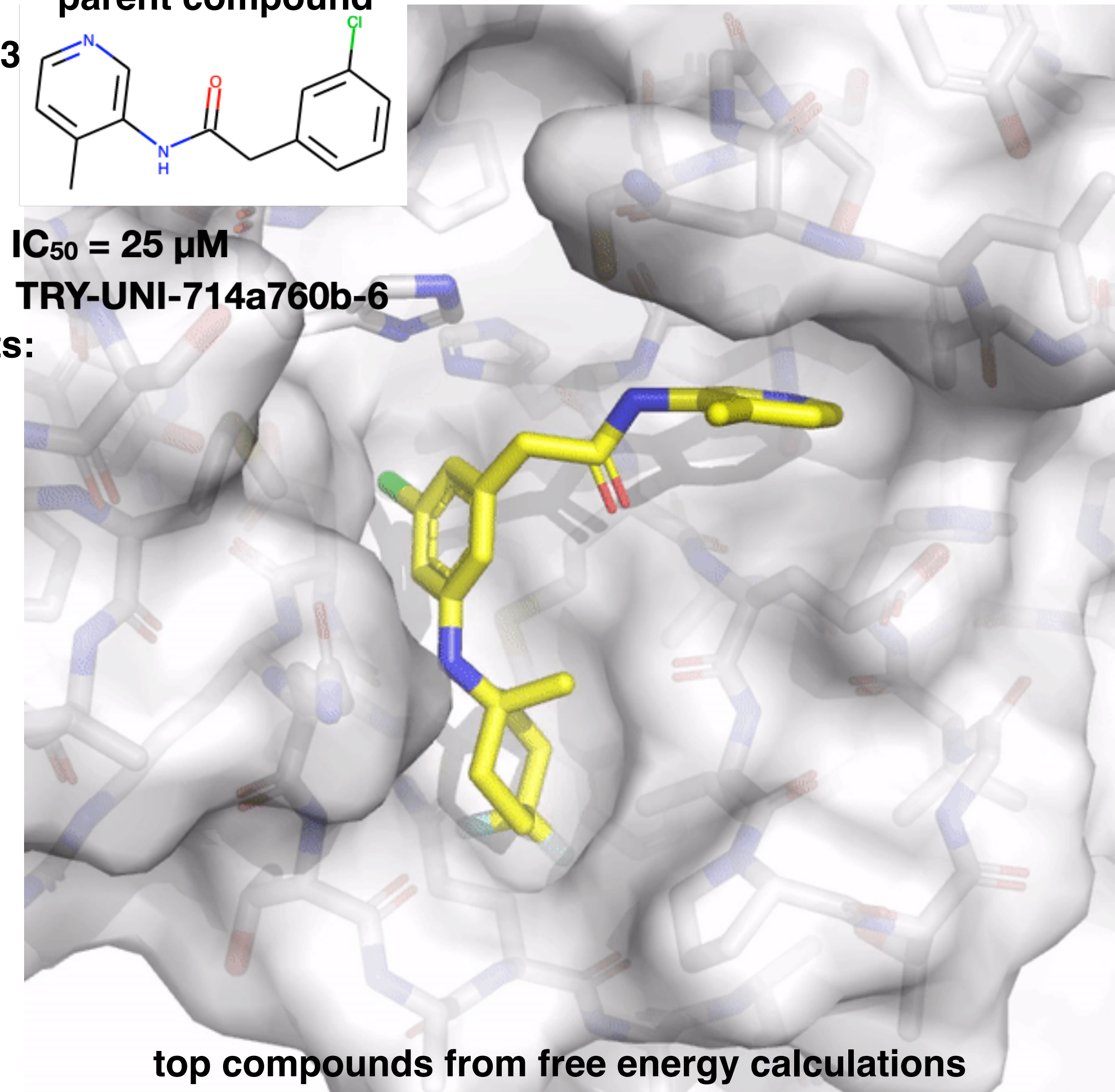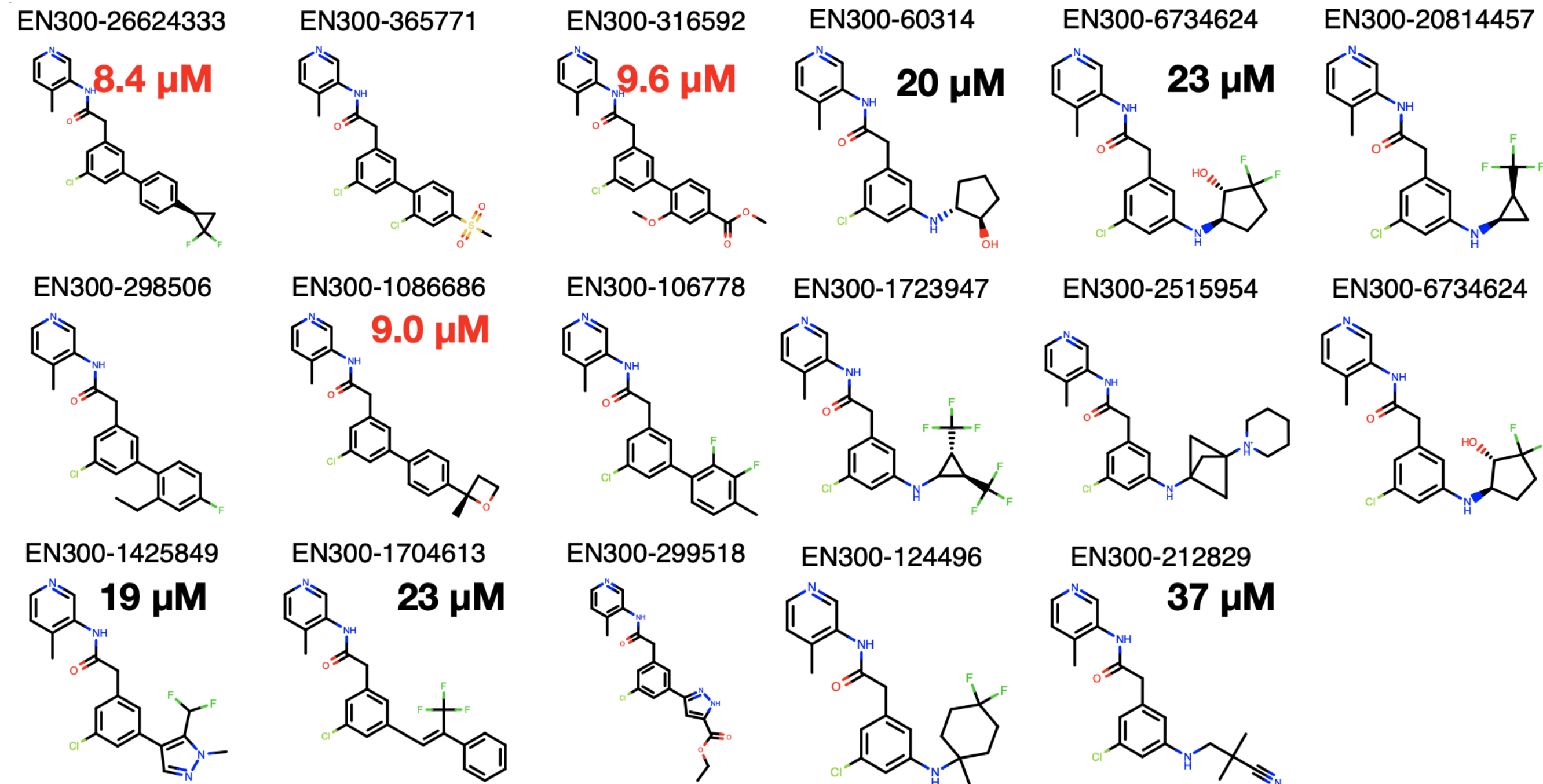
parent compound

$IC_{50}$ = 25 μM
TRY-UNI-714a760b-6

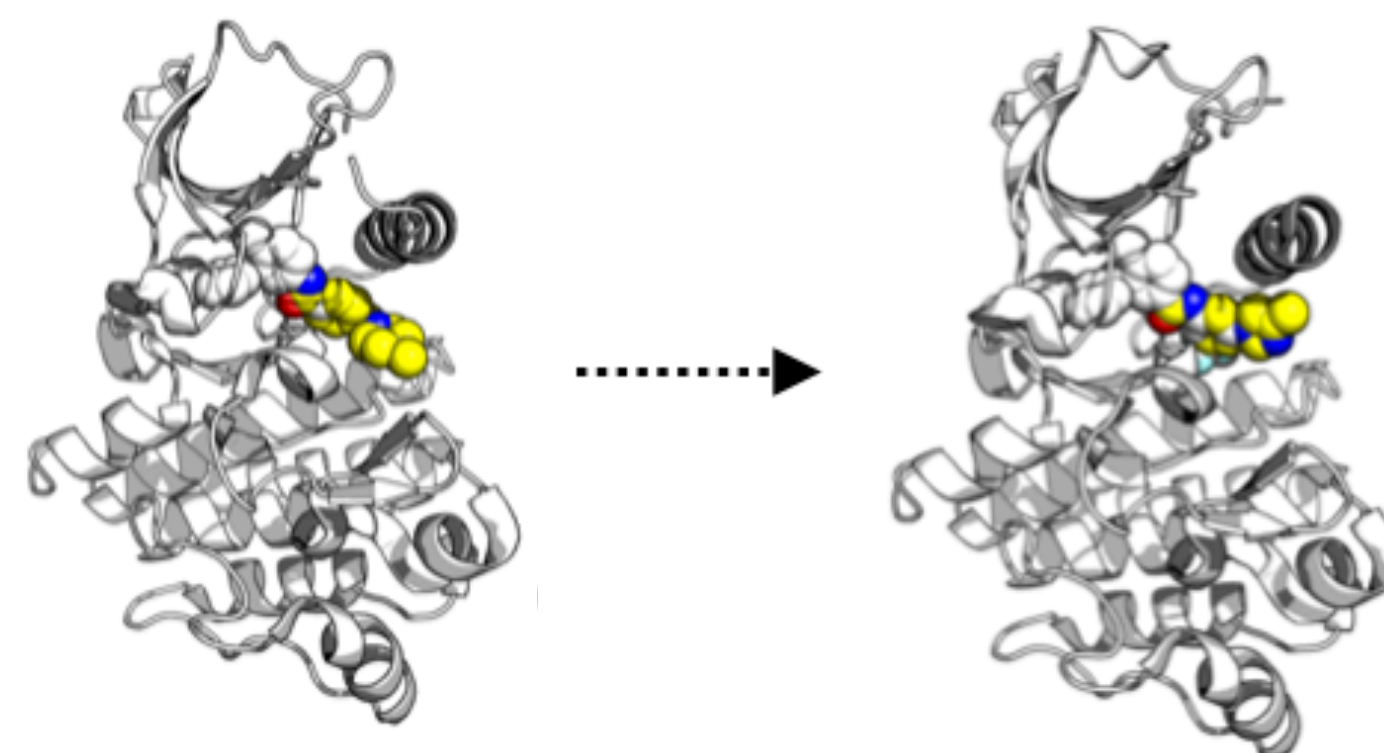**Top free energy calculation compounds and experimental affinity measurements:**

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

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

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

top compounds from free energy calculations

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

# ALCHEMICAL FREE ENERGY CALCULATIONS HAVE A BROAD DOMAIN OF APPLICABILITY IN DRUG DISCOVERY

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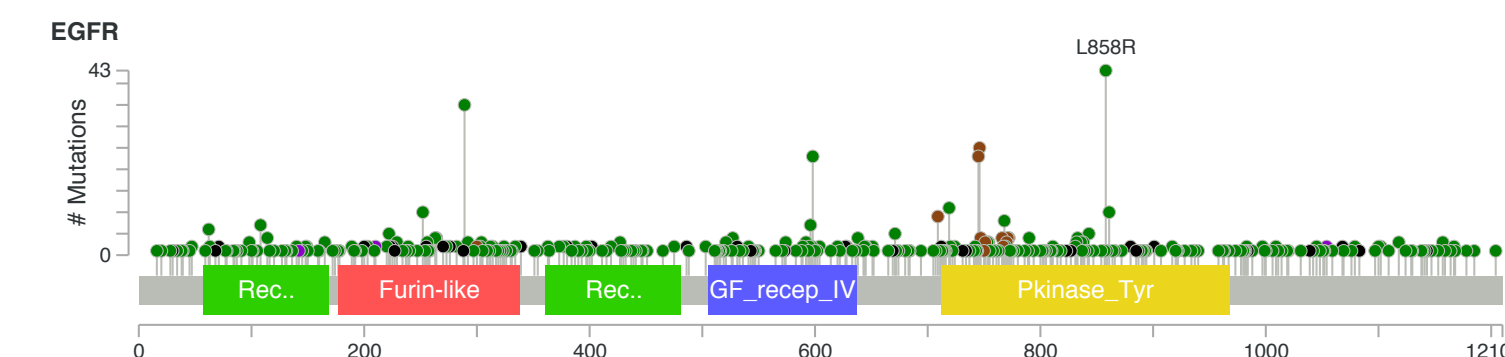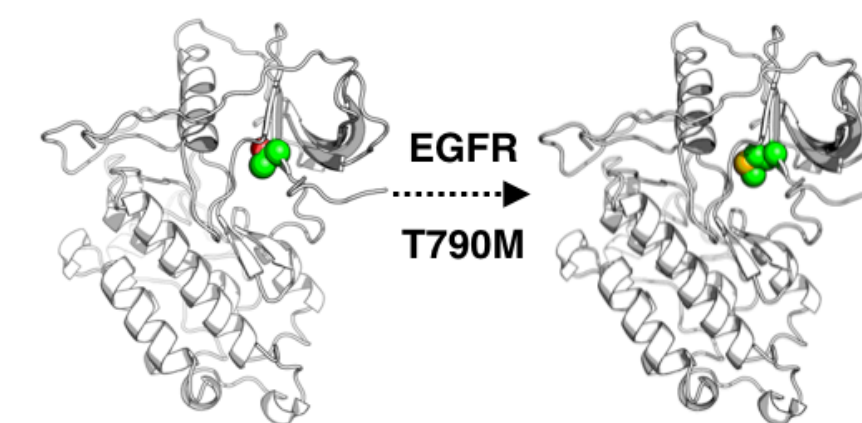Aldeghi et al. JACS 139:946, 2017.
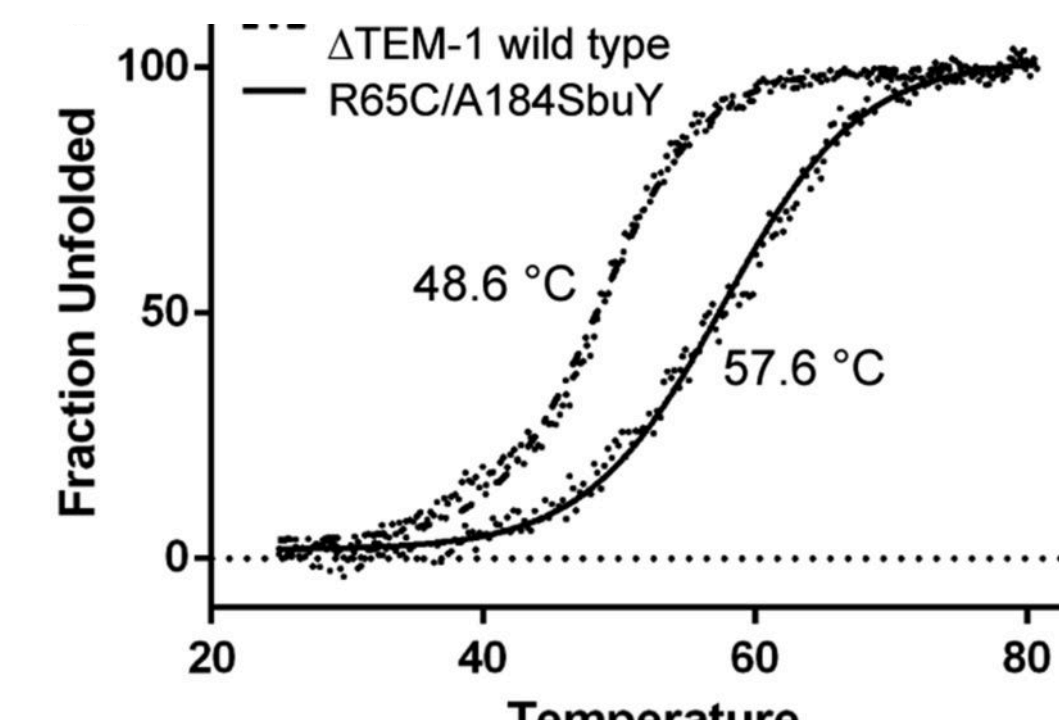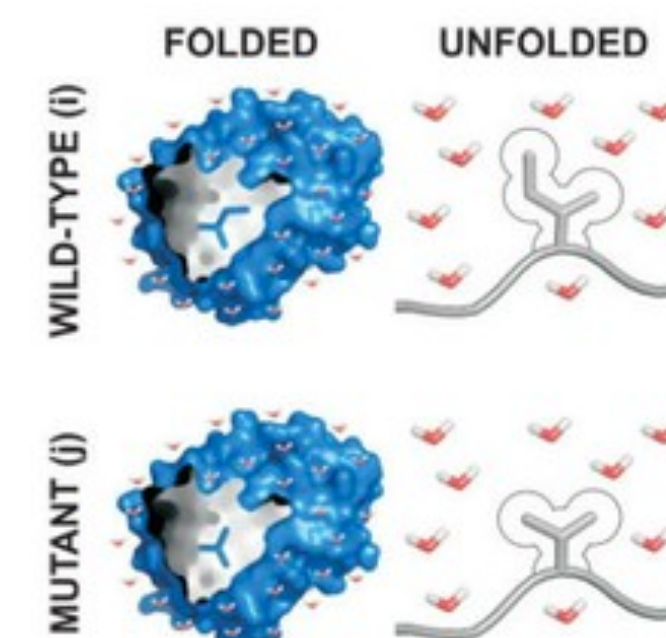https://doi.org/10.1021/jacs.6b11467

**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
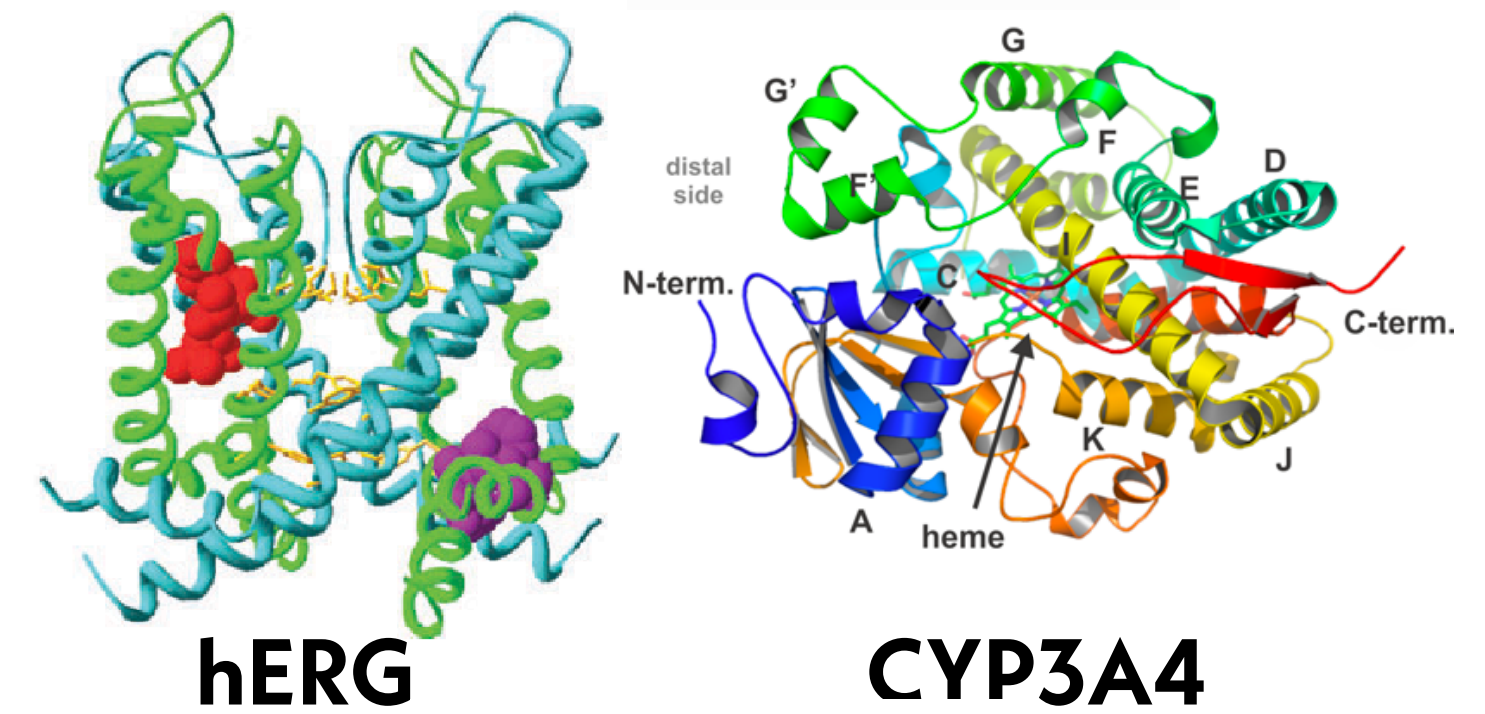https://doi.org/10.1002/anie.201510054

# ...AND HOLD THE POTENTIAL FOR EVEN BROADER APPLICABILITY AS MORE STRUCTURAL DATA EMERGES
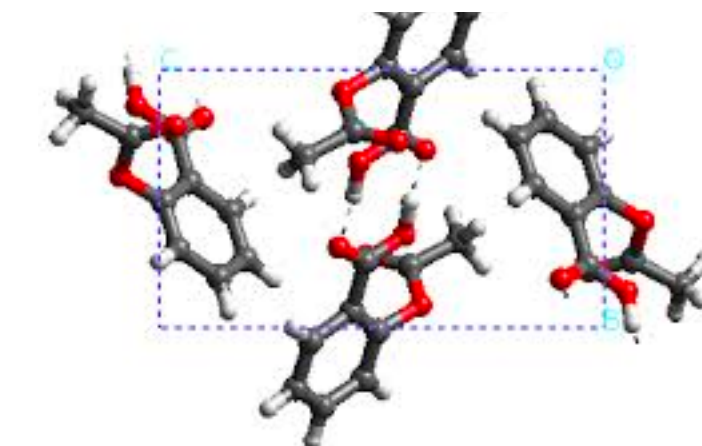
partition coefficients (logP, logD) and permeabilities

structure-enabled ADME/Tox targets

hERG    CYP3A4

porin permeation

crystal polymorphs, etc.

# Obligatory dead person quote

**"Knowledge is a treasure, but practice is the key to it."**
- Emil Fischer, the father of synthetic organic chemistry
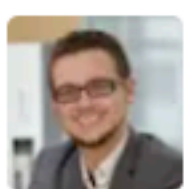
Please provide me with a relevant quote from Emil Fischer about the importance of learning.

Emil Fischer, a German chemist and one of the pioneers in the field of biochemistry, once said:

"Knowledge is a treasure, but practice is the key to it."

**Now provide a relevant quote from Ludwig Boltzmann.**

Ludwig Boltzmann, an Austrian physicist and one of the founders of statistical mechanics, once said:

"The general laws of nature are to be expressed by equations which hold good for all phenomena, whether physical or moral; and these equations are to be derivable from the principle of the conservation of energy."
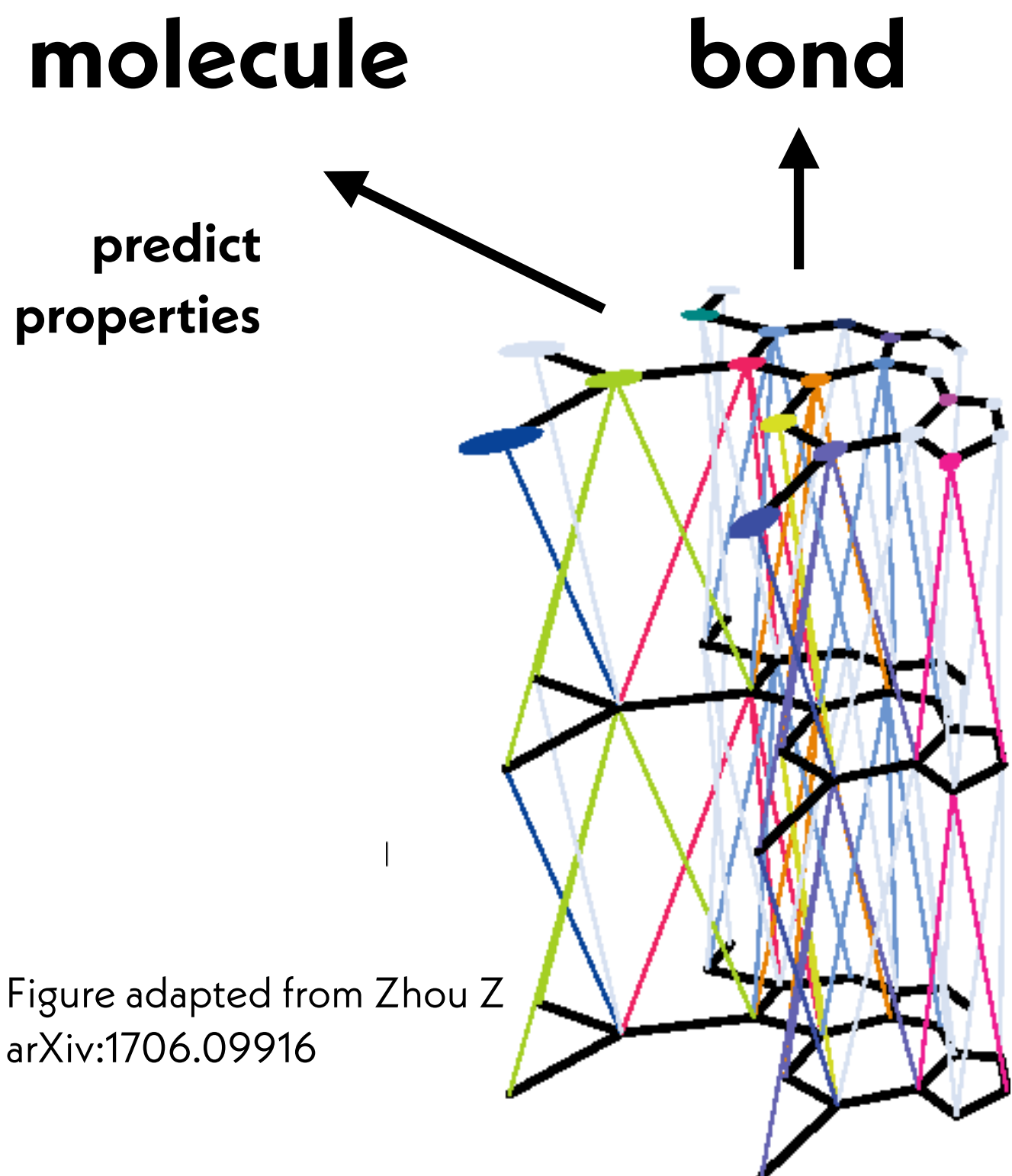
**That quote isn't about learning.**

I apologize for misunderstanding your request. Here is a relevant quote from Ludwig Boltzmann about learning:

"The most important thing is to have a deep understanding of the fundamental principles. Without them, one can learn nothing."
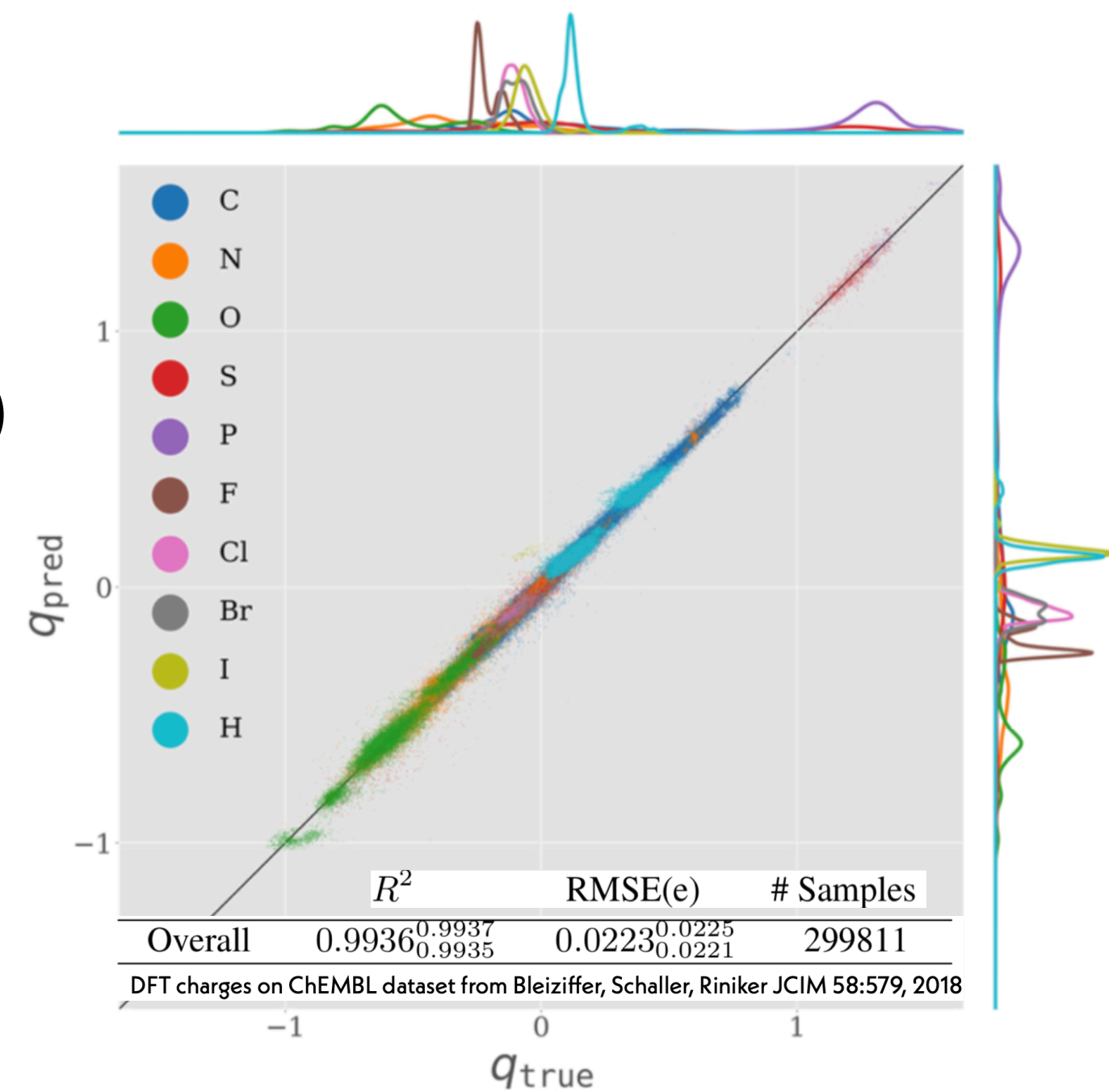
# A DEEP UNDERSTANDING OF THE FUNDAMENTAL PRINCIPLES MAKES LEARNING A HELL OF A LOT EASIER

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

predict
properties



Figure adapted from Zhou Z
arXiv:1706.09916

Learns **electronegativity** ($e_i$) and **hardness** ($s_i$) subject to fixed charge sum constraint:

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

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



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

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

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

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

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

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

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

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

## ∇imlet

**Graph Inference on MoLEcular Topology**

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

**YUANQING WANG**

**Where else can we apply this principle?**

* start with a fundamental physical or statistical mechanical model
* identify areas where a poor approximation has been inserted
* introduce a flexible, learnable model
* train with lots of (potentially synthetic) data

# DRUG DISCOVERY IS NOT A BIG DATA PROBLEM

**DALL-E 2** was trained on a dataset of **650 million** images
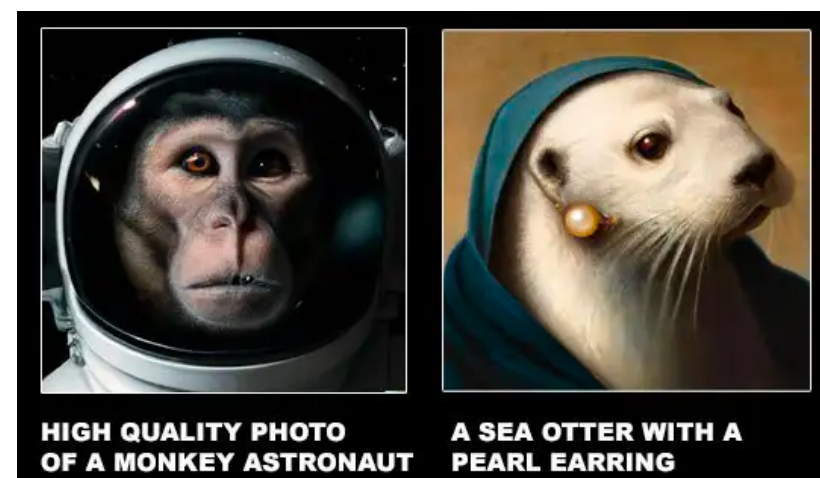
**GPT-3** was trained on a corpus of **22.5 billion pages of text** (45 TB)

Typical drug discovery programs make and test **~2000 compounds**
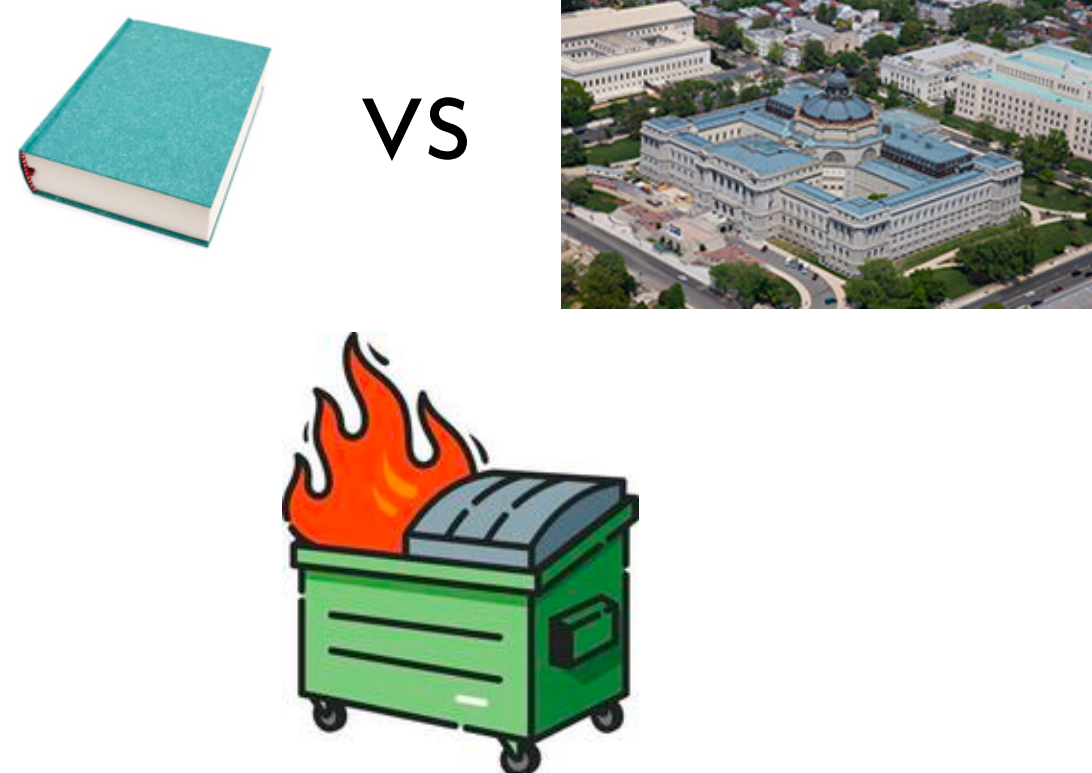
Trying to use public datasets ingested from publications with heterogeneous methods is like "dumpster diving for data"

We need to:
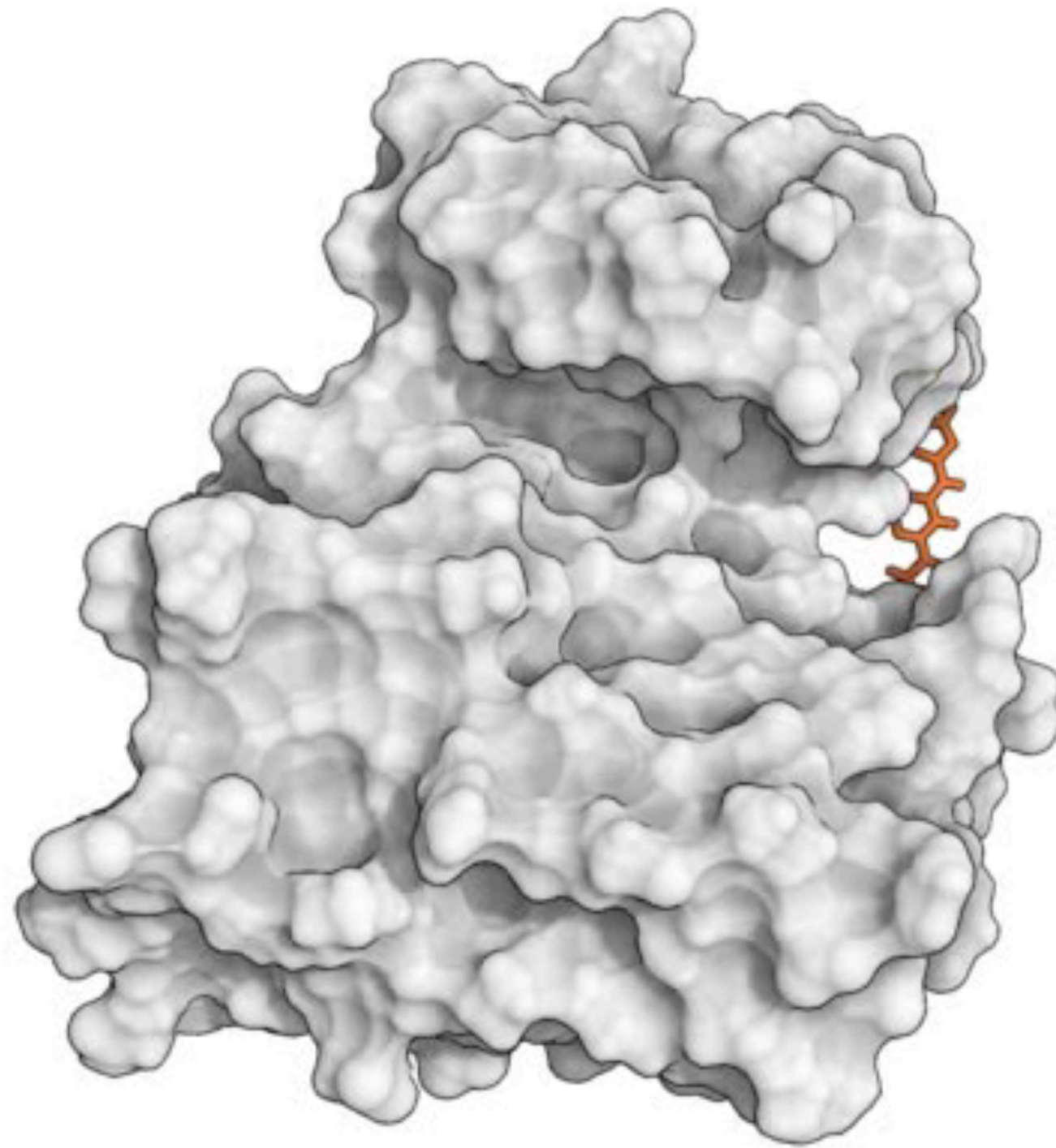* Develop **extremely data efficient** machine learning methods or leverage **synthetic data** (e.g. quantum chemistry) where possible
* Find a way to make data from different discovery programs **fit into the same model** (pool all data together)

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

## typical class I molecular mechanics force field



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

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

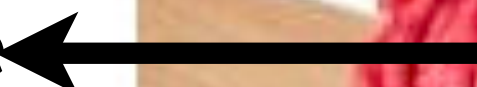# FORCE FIELDS HAVE TRADITIONALLY BEEN HEROIC PRODUCTS OF HUMAN EFFORT

experimental data
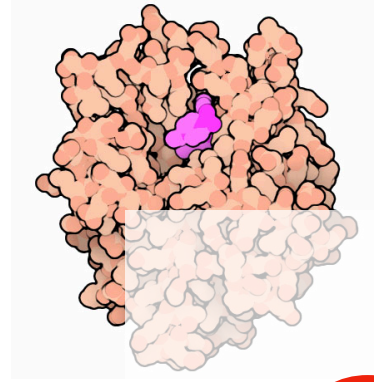quantum chemistry
keen chemical intuition

heroic effort by graduate students and postdocs

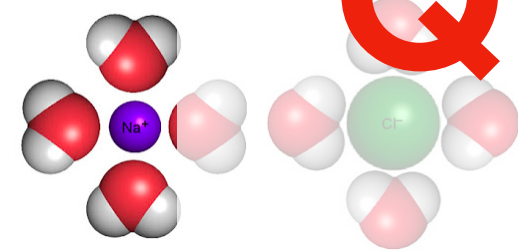a parameter set we desperately hope someone actually uses

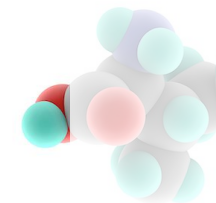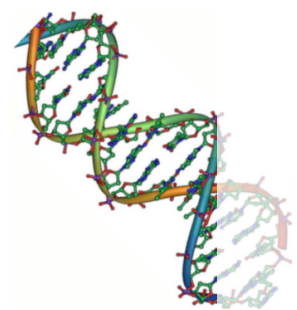# FORCE FIELDS HAVE TRADITIONALLY BEEN HEROIC PRODUCTS OF HUMAN EFFORT
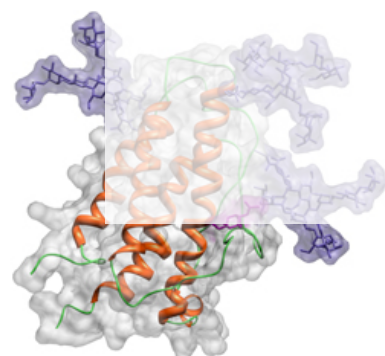
proteins

post-translational modifications

Amber20 recommendations

**Quickly adds up to >100 human-years**

water

ions

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

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

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

nucleic acids

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

lipids

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

carbohydrates

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.

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. Chem. ...*, 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 phosphohisti-dine. *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.

I. S. Joung; T. E. Cheatham, III. Molecular dynamics simulations of the dynamic and energetic properties ... fic ion parameters. *J. Phys. Chem. B*, 2009, 113, 13279–13290.

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

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

R. Galindo-Murillo; J. C. Robertson; M. Zgarbovic; J. Sponer; M. Otyepka; P. Jureska; T. E. Cheatham. ... the Current 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 ...

Å. 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.

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.

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



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

chemical graph →(abstraction)→ topology graph →(GN(•; Φ_NN))→ atom embeddings

**Stage 2:** symmetry-preserving pooling

$NN_\phi(\;;\Phi_{NN}) = NN_\phi(\;;\Phi_{NN}) + NN_\phi(\;;\Phi_{NN})$
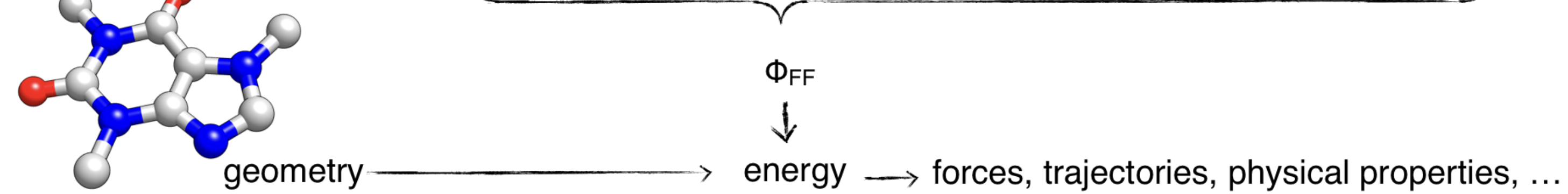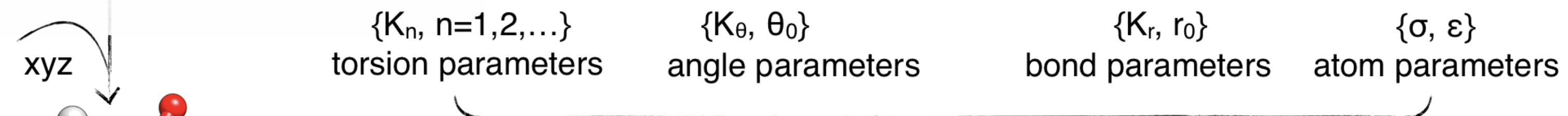
$NN_\theta(\;;\Phi_{NN}) = NN_\theta(\;;\Phi_{NN}) + NN_\theta(\;;\Phi_{NN})$

$NN_r(\;;\Phi_{NN}) = NN_r(\;;\Phi_{NN}) + NN_r(\;;\Phi_{NN}) \ldots$

torsion embeddings   angle embeddings   bond embeddings

$NN_{readout}(•; \Phi_{NN})$ 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   $\{\sigma, \varepsilon\}$ atom parameters

$\Phi_{FF}$

geometry ⟶ energy ⟶ forces, trajectories, physical properties, …

JOSH FASS   YUANQING WANG

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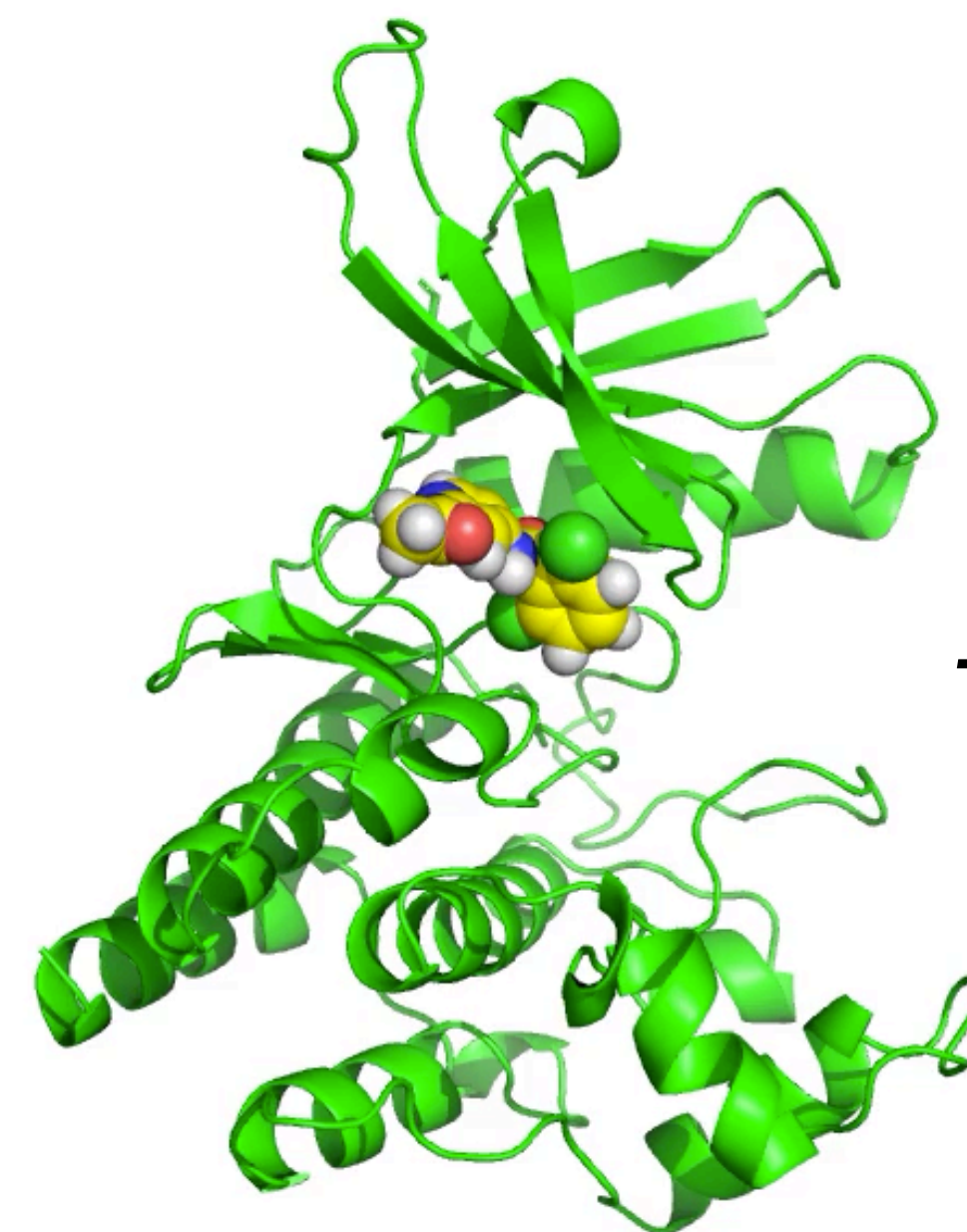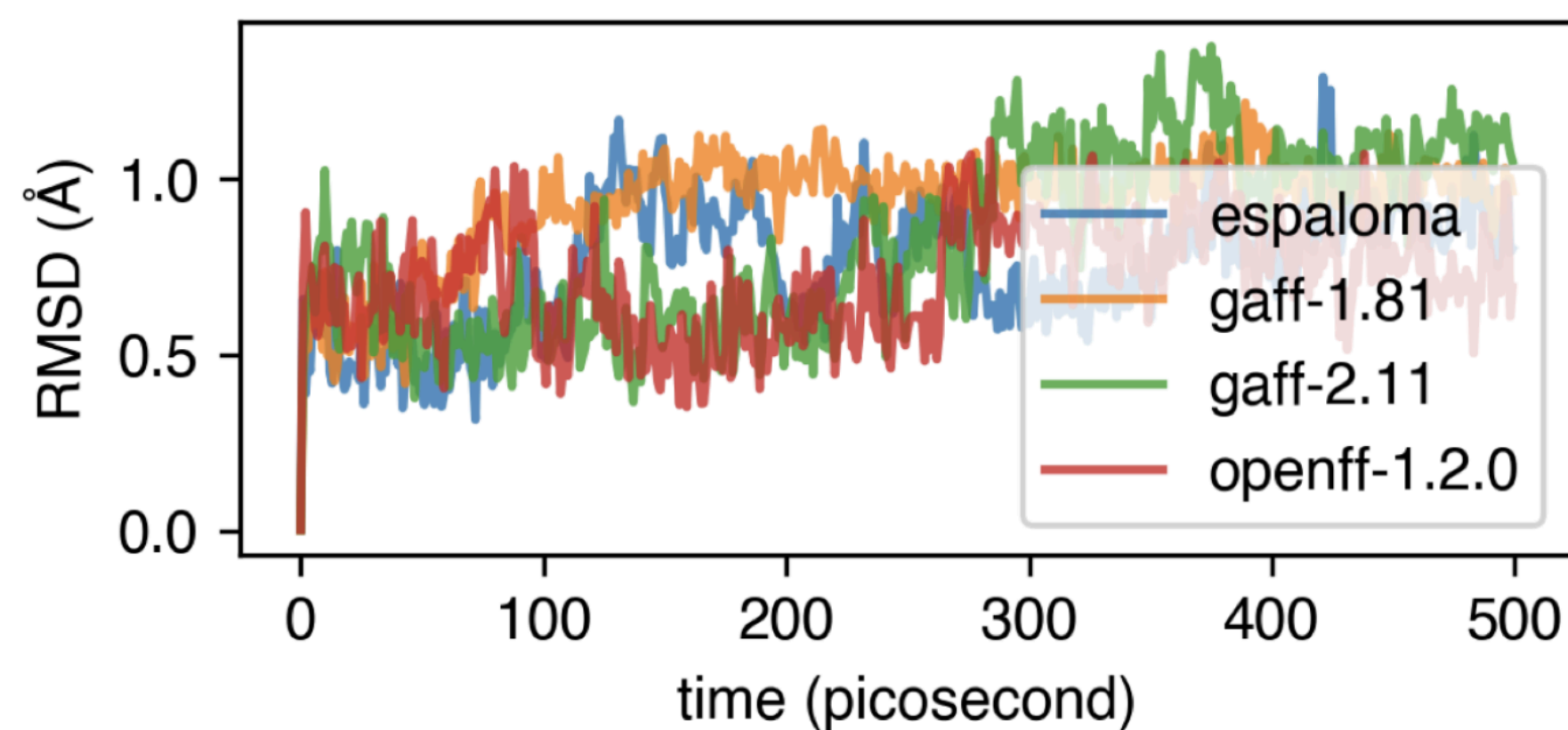**code**: https://github.com/choderalab/espaloma

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

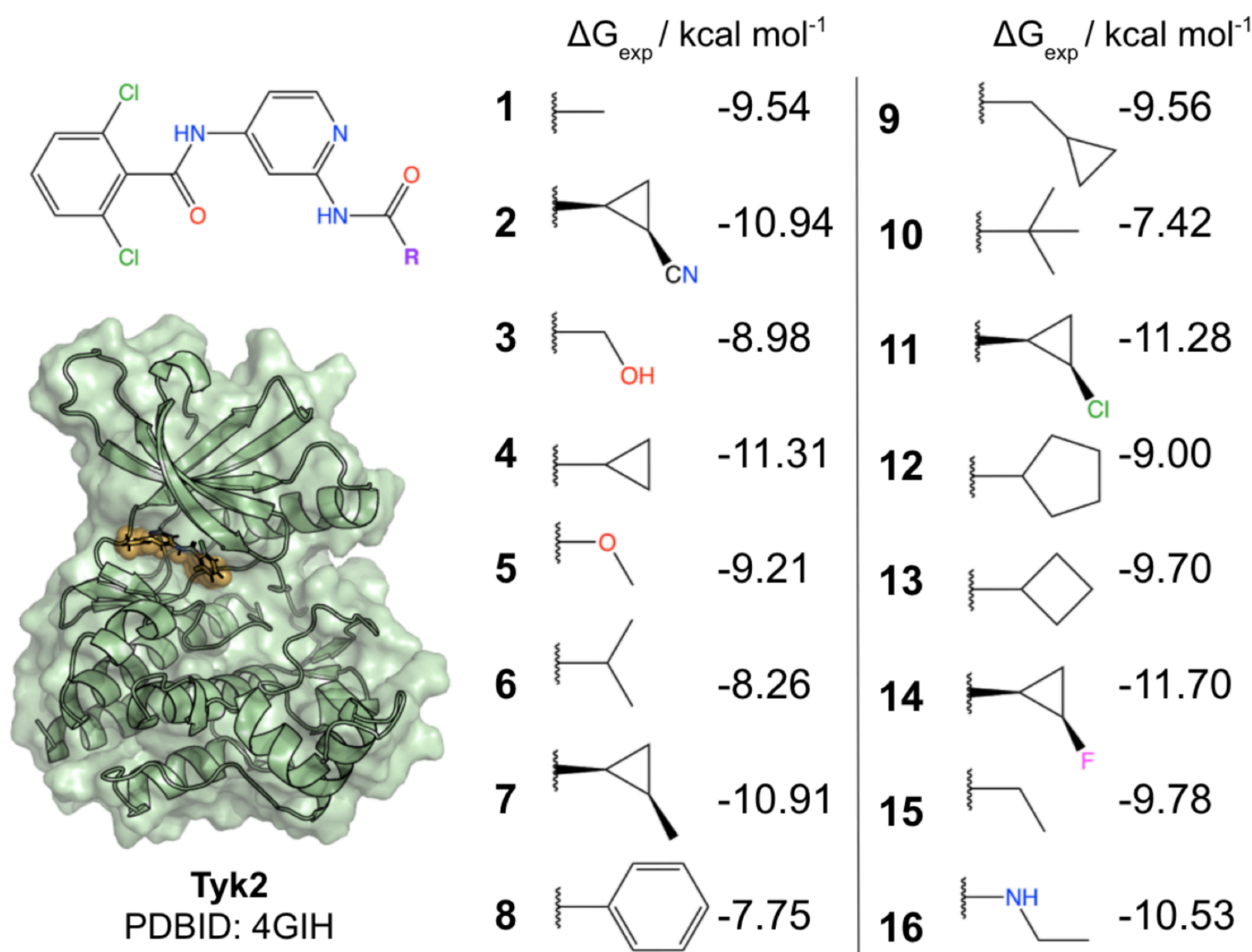| (a) dataset | # mols | # trajs | # snapshots | Espaloma RMSE | | Legacy FF RMSE (kcal/mol) (Test molecules) | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | Train | Test | OpenFF 1.2.0 | GAFF-1.81 | GAFF-2.11 | Amber ff14SB |
| **PhAlkEthOH** (simple CHO) | 7408 | 12592 | 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 | $0.8264^{0.9007}_{0.7682}$ | $1.8764^{1.9947}_{1.7827}$ | $2.1768^{2.3388}_{2.0380}$ | $2.4274^{2.5207}_{2.3300}$ | $2.5386^{2.6640}_{2.4370}$ | |
| **joint** — 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**

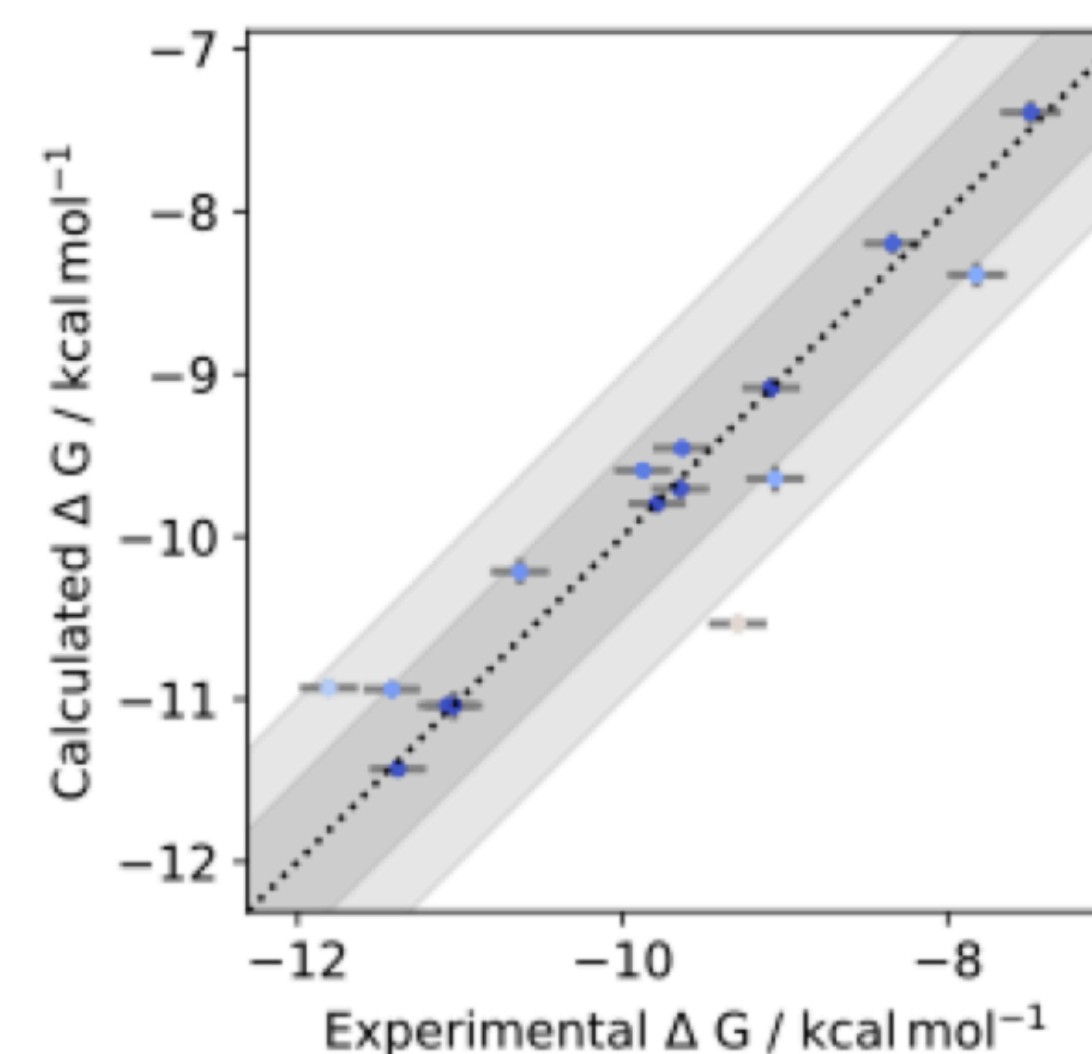# ELIMINATING DISCRETE TYPES APPEARS TO SIGNIFICANTLY IMPROVE ACCURACY IN FREE ENERGY CALCULATIONS

preprint: https://arxiv.org/abs/2010.01196
code: http://github.com/choderalab/espaloma
free energy calculations with http://github.com/choderalab/perses

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

## week 1

**2023**

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

using published force field model

## week 2

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

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

## week 1

**2025**

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

using force field model
built from public + private data

## week 2

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

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

# CAN WE LEARN TO FIT EXPERIMENTAL DATA AS WELL?

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

loss function:

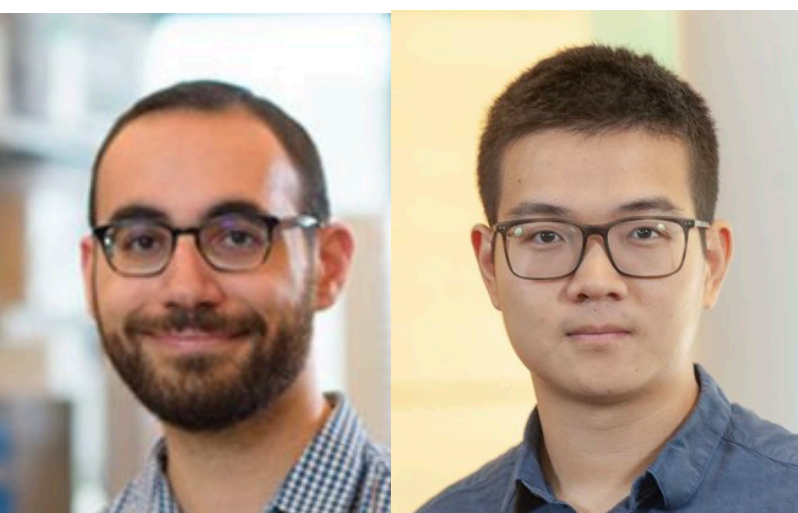$$L(\Phi_{NN}) = \sum_{n=1}^{N} \frac{[\Delta G_n(\Phi_{NN}) - \Delta G_n^{\exp}]^2}{\sigma_n^2}$$

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

**JOSH FASS**

**YUANQING WANG**

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

## OBC2 GBSA FreeSolv RMSE

# WHY SHOULD WE BE STUCK WITH A PHYSICAL MODEL THAT CATERED TO THE CAPABILITIES OF A PDP-11?



**DEC PDP-11**

~45 years old

typical class I molecular mechanics force field



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

shitty Taylor series
truncated at lowest order

crappy Fourier series
truncated at n=6

don't even get me
started on this fucker

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

# WE COULD GO TO CLASS II FORCE FIELDS...
# IT'S CERTAINLY EASY TO DO NOW

$$E = \sum_b [{}^2K_b(b - b_0)^2 + {}^3K_b(b - b_0)^3 + {}^4K_b(b - b_0)^4]$$

$$+ \sum_\theta [{}^2K_\theta(\theta - \theta_0)^2 + {}^3K_\theta(\theta - \theta_0)^3 + {}^4K_\theta(\theta - \theta_0)^4]$$

$$+ \sum_\phi [{}^1K_\phi(1 - \cos \phi) + {}^2K_\phi(1 - \cos 2\phi) + {}^3K_\phi(1 - \cos 3\phi)]$$

$$+ \sum_\chi K_\chi \chi^2 + \sum_{i>j} \frac{q_i q_j}{r_{ij}} + \sum_{i>j} \epsilon \left[ 2\left(\frac{r^*}{r_{ij}}\right)^9 - 3\left(\frac{r^*}{r_{ij}}\right)^6 \right]$$

$$+ \sum_b \sum_{b'} K_{bb'}(b - b_0)(b' - b'_0) + \sum_\theta \sum_{\theta'} K_{\theta\theta'}(\theta - \theta_0) \times$$

$$(\theta' - \theta'_0)$$

$$+ \sum_b \sum_\theta K_{b\theta}(b - b_0)(\theta - \theta_0)$$

$$+ \sum_\phi \sum_b (b - b_0)[{}^1K_{\phi b} \cos \phi + {}^2K_{\phi b} \cos 2\phi + {}^3K_{\phi b} \cos 3\phi]$$

$$+ \sum_\phi \sum_{b'} (b' - b'_0)[{}^1K_{\phi b'} \cos \phi + {}^2K_{\phi b'} \cos 2\phi +$$

$$\quad {}^3K_{\phi b'} \cos 3\phi]$$

$$+ \sum_\phi \sum_\theta (\theta - \theta_0)[{}^1K_{\phi\theta} \cos \phi + {}^2K_{\phi\theta} \cos 2\phi + {}^3K_{\phi\theta} \cos 3\phi]$$

$$+ \sum_\phi \sum_\theta \sum_{\theta'} K_{\phi\theta\theta'} (\theta - \theta_0)(\theta' - \theta'_0) \cos \phi \qquad (1)$$

bond-bond: angle node
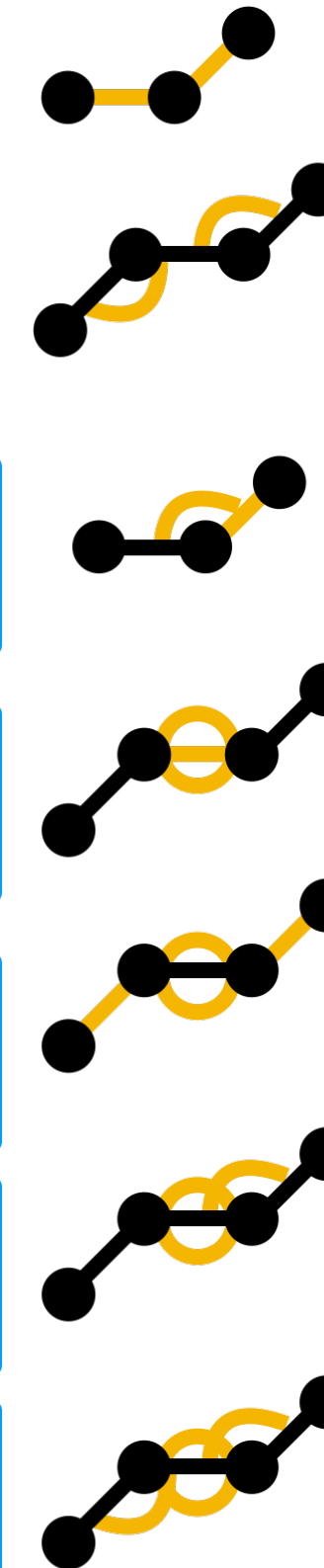
angle-angle: torsion node

bond-angle: angle node

torsion-(center) bond: torsion

torsion-(side) bond: torsion

torsion-angle: torsion

torsion-angle-angle: torsion

But can we do a better job of modeling true many-body local valence terms?

Hwang et al. (1994) http://doi.org/10.1021/ja00085a036

# A NEW GENERATION OF QUANTUM MACHINE LEARNING (QML) POTENTIALS PROVIDE SIGNIFICANTLY MORE FLEXIBILITY IN FUNCTIONAL FORM, THOUGH AT MUCH GREATER COST
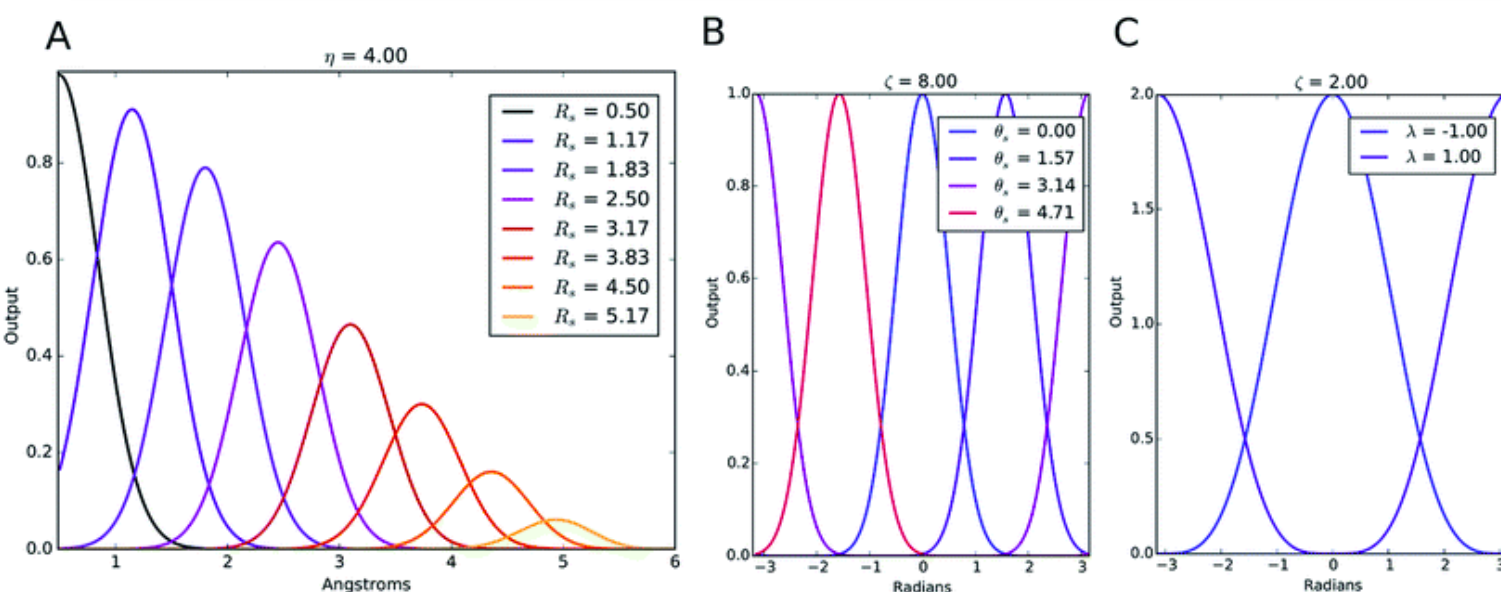
**ANI** family of quantum machine learning (QML) potentials
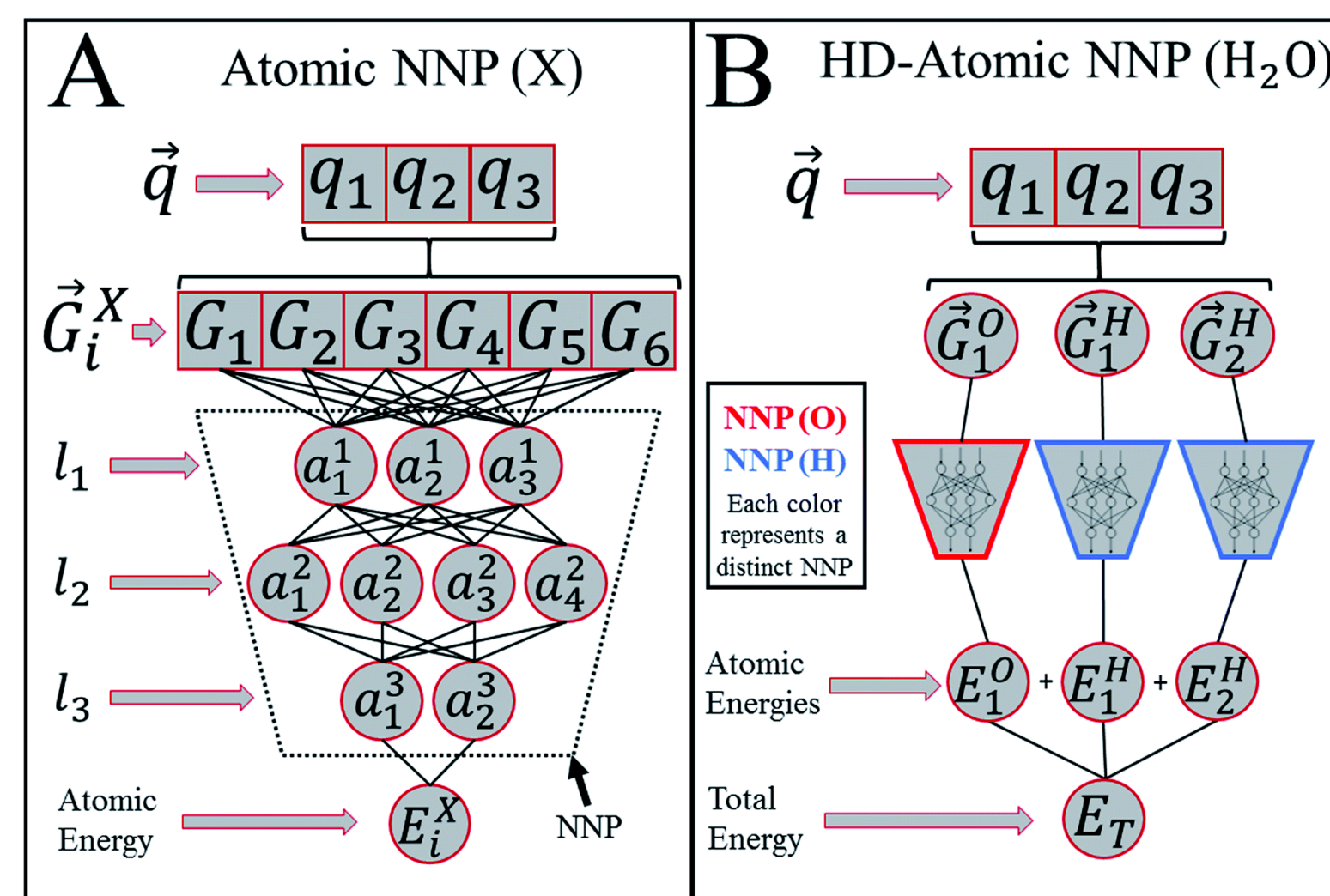
**radial** and **angular** features

deep neural network for each atom

excellent agreement with DFT

# QML POTENTIALS ARE SEEING RAPID EVOLUTION IN ARCHITECTURES

**ANI**     **Deep Tensor Networks**     **Tensor Field Networks**     **PotentialNet**



The **ANI** class of models uses distance- and angle-based features [http://doi.org/10.1039/c6sc05720a].
**Deep Tensor Networks** and **SchNet** use distance-based features for continuous convolutions [https://doi.org/10.1038/ncomms13890].
**Tensor Field Networks** and Clebsch-Gordon nets use spherical harmonics [https://arxiv.org/abs/1802.08219; https://bit.ly/2SRVS67].
**PotentialNet** uses a graph convolutional network augmented by distance-dependent edges [https://doi.org/10.1021/acscentsci.8b00507].

# HYBRID QUANTUM MACHINE LEARNING / MOLECULAR MECHANICS (QML/MM) SIMULATIONS ARE MOST FEASIBLE IN THE NEAR TERM



many QML/MM formulations possible, including those that use QML for protein-ligand interactions

$$U_{QML/MM}(X_P, X_L) = U_{MM}(X_P, X_L) - U_{MM}^{vacuum}(X_L) + U_{QML}^{vacuum}(X_L)$$

| MM | openforcefield 1.0.0 |
| QML | ANI2x |

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

# WE CAN ASSESS HOW WELL QML/MM FREE ENERGY CALCULATIONS MIGHT PERFORM THROUGH A PERTURBATIVE CORRECTION



**A**      ML/MM AUGMENTED THERMODYNAMIC CYCLE

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



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

|  | 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.

**Figure 7. ML/MM corrections to MM binding free energies can be up to 4 kcal mol⁻¹ in magnitude.** The signed $\Delta G^{MM \to ML/MM}$ corrections for each ligand (with R-group shown) are shown, ordered from least positive (slightly disfavoring binding) to most positive (strongly disfavoring binding).

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

Rufa, Bruce Macdonald, Fass, Wieder, Grinaway, Roitberg, Isayev, and **Chodera**.

# 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)

MM: openff-1.0.0 (N = 16)

| RMSE: | 0.97 [95%: 0.68, 1.22] |
| MUE: | 0.77 [95%: 0.51, 1.08] |
| R2: | 0.42 [95%: 0.08, 0.75] |
| rho: | 0.65 [95%: 0.25, 0.88] |

ML/MM: openff-1.0.0 with ANI2x (N = 16)

| RMSE: | 0.47 [95%: 0.32, 0.68] |
| MUE: | 0.35 [95%: 0.24, 0.56] |
| R2: | 0.86 [95%: 0.66, 0.95] |
| rho: | 0.93 [95%: 0.79, 0.97] |

replica-exchange free energy calculations with perses

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

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



atomic coordinates

feature computation

NN computation

energy/force accumulation

$$\sum_i$$

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

**paper:** https://arxiv.org/abs/2201.08110
**code:** https://github.com/openmm/nnpops

# OPENMM 8 WILL MAKE QML/MM SIMULATIONS INCREDIBLY EASY

```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)
```

**OpenMM 8** was just released!

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

# WHY DO WE NEED MM AT ALL?



Can we just use ML force fields for everything?
We can finally be free of the hegemony of bonds!

ANI2x

# 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

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



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

Simple restraints can be used when we need to enforce specific chemical species

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

**ANI-2x**

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

# STATISTICAL MECHANICS IS ESSENTIAL IN TAUTOMER RATIOS.
# EVEN IN VACUUM, ONLY SUMMING OVER MINIMA INTRODUCES HUGE ERRORS.

ANI1ccx alchemical dG vs mining minima dG

$\rho = 0.96; [0.94, 0.97]$

$RMSE = 2.5; [2.1, 2.8]$

alchemical dG [kcal/mol]

mining minima dG [kcal/mol]

JOSH FASS

MARCUS WIEDER

OLEXANDR ISAYEV

ADRIAN ROITBERG

# PURE QUANTUM MACHINE LEARNING (QML) POTENTIALS CAN BE TUNED/RETRAINED BY FREE ENERGIES, REGULARIZED BY QM DATA



test set performance

training / validation optimization

Regularization by QM data

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

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

Fast on-the-fly reweighting enables inexpensive loss/gradient computation without repeating expensive free energy calculation

# The MoISSI Quantum Chemistry Archive

A central source to compile, aggregate, query, and share quantum chemistry data.

**GET STARTED!**

## QCArchive
### A MoISSI Project

**FAIR Data**

MoISSI hosts the QCArchive server, the largest publicly available collection of quantum chemistry data. So far, it stores over ten million computations for the molecular sciences community.

**Interactive Visualization**

Not only for computing and storing quantum chemistry computations at scale, but also for visualizing and understanding results as well.

**Private Instances**

The infrastructure behind QCArchive is fully open-souce. Spin up your own instance to compute private data and share only with collaborators.

| 102,477,973 MOLECULES | 108,469,316 RESULTS | 212 COLLECTIONS |

http://qcarchive.molssi.org

**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**

| Subset | Molecules | Conformations | Atoms | Elements |
|---|---|---|---|---|
| Dipeptides | 677 | 33850 | 26–60 | H, C, N, O, S |
| Solvated Amino Acids | 26 | 1300 | 79–96 | H, C, N, O, S |
| DES370K Dimers | 3490 | 345676 | 2–34 | H, Li, C, N, O, F, Na, Mg, P, S, Cl, K, Ca, Br, I |
| DES370K Monomers | 374 | 18700 | 3–22 | H, C, N, O, F, P, S, Cl, Br, I |
| PubChem | 14643 | 731856 | 3–50 | H, C, N, O, F, P, S, Cl, Br, I |
| Ion Pairs | 28 | 1426 | 2 | Li, F, Na, Cl, K, Br, I |
| Total | 19238 | 1132808 | 2–96 | H, Li, C, N, O, F, Na, Mg, P, S, Cl, K, Ca, Br, I |

**DFT ωB97M-D3(BJ)/def2-TZVPPD level of theory**
>4M core-hours computed on QCFractal academic clusters



**SPICE QML model: 0.7 kcal/mol**
median absolute error

**Figure 5.** Absolute error as a function of the total number of atoms in a molecule. The line indicates the median for all molecules of a certain size, and the gray region contains the central 50% of samples.

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

# SPICE IS OUR FIRST STEP TOWARD BUILDING "FOUNDATION MODELS" THAT CAN BE RAPIDLY TAILORED TO DIFFERENT APPLICATIONS

## Dataset

| QC specification | Dataset | Category | # Mols | # Conformations |
|---|---|---|---|---|
| Openff-default | OpenFF Gen2 Optimization | Small molecules | 1022 | 244944 |
| | OpenFF PepConf Optimization | Di-, Tri-peptides | 522 | 228582 |
| | RNA-BGSU Diverse Dataset | RNA trinucleotides | 64 | 3649 |
| | RNA-BGSU Trinucleotide Dataset | RNA trinucleotides | 64 | 35134 |
| | SPICE Pubchem | Small molecules | 14110 | 601719 |
| | SPICE Dipeptide | Dipeptides | 677 | 25098 |
| | SPICE DES Monomers | Small molecules | 369 | 18450 |

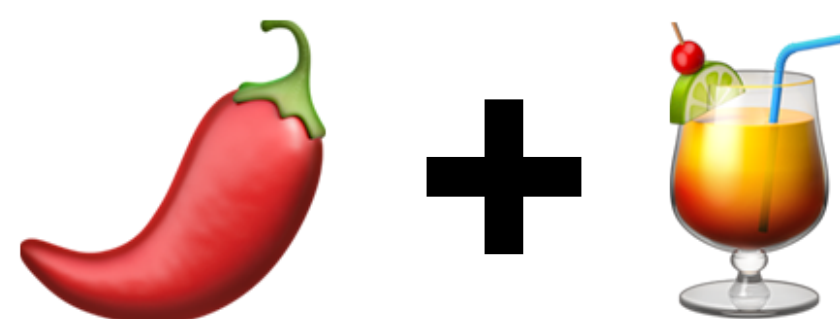| Name | Dataset | Mols | Conformations | Energy RMSE (kcal/mol) Force RMSE (kcal/mol·Å⁻¹) | | Baseline FF Energy RMSE (kcal/mol) (Test molecules) Baseline FF Force RMSE (kcal/mol·Å⁻¹) (Test molecules) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Train | Test | GAFF-1.81 | GAFF-2.11 | OpenFF-1.2.0 | OpenFF-2.0.0 |
| **Joint** | Gen2 | 1022 | 244944 | $1.13^{1.16}_{1.11}$ $90.41^{94.92}_{86.35}$ | $\mathbf{1.77^{1.85}_{1.70}}$ $35.66^{44.16}_{26.51}$ | $2.97^{3.03}_{2.91}$ $10.68^{10.70}_{10.65}$ | $2.96^{3.03}_{2.90}$ $10.65^{10.68}_{10.62}$ | $2.82^{2.90}_{2.76}$ $11.11^{11.14}_{11.08}$ | $2.69^{2.77}_{2.61}$ $\mathbf{10.33^{10.36}_{10.30}}$ |
| | PepConf | 522 | 228582 | $1.61^{1.63}_{1.58}$ $6.17^{6.17}_{6.17}$ | $\mathbf{1.97^{2.03}_{1.92}}$ $\mathbf{6.27^{6.27}_{6.26}}$ | $3.64^{3.68}_{3.60}$ $192.61^{300.66}_{65.86}$ | $4.61^{4.66}_{4.55}$ $78.83^{124.76}_{33.44}$ | $3.09^{3.14}_{3.04}$ $13.67^{17.91}_{10.49}$ | $3.23^{3.28}_{3.18}$ $45.45^{69.14}_{18.05}$ |
| | SPICE-Pubchem | 14110 | 601719 | $2.39^{2.40}_{2.38}$ $9.10^{9.16}_{9.06}$ | $\mathbf{2.68^{2.71}_{2.65}}$ $\mathbf{9.40^{9.46}_{9.35}}$ | $4.44^{4.48}_{4.40}$ $14.62^{14.65}_{14.59}$ | $4.62^{4.66}_{4.58}$ $15.16^{15.21}_{15.11}$ | $4.28^{4.32}_{4.24}$ $14.82^{14.85}_{14.80}$ | $4.32^{4.36}_{4.28}$ $14.66^{14.69}_{14.64}$ |
| | SPICE-Dipeptide | 677 | 25098 | $2.48^{2.51}_{2.45}$ $6.64^{6.65}_{6.63}$ | $\mathbf{2.46^{2.55}_{2.38}}$ $\mathbf{6.39^{6.43}_{6.36}}$ | $4.16^{4.29}_{4.03}$ $12.76^{12.81}_{12.70}$ | $4.15^{4.27}_{4.02}$ $12.67^{12.73}_{12.62}$ | $4.04^{4.17}_{3.91}$ $12.41^{12.46}_{12.36}$ | $3.83^{3.96}_{3.70}$ $12.72^{12.77}_{12.67}$ |
| | SPICE-DES-Monomers | 369 | 18450 | $1.37^{1.40}_{1.35}$ $7.12^{7.17}_{7.07}$ | $\mathbf{1.70^{1.84}_{1.58}}$ $\mathbf{8.69^{8.82}_{8.57}}$ | $2.84^{3.08}_{2.61}$ $13.77^{14.07}_{13.51}$ | $2.75^{2.97}_{2.52}$ $13.35^{13.58}_{13.11}$ | $2.99^{3.20}_{2.80}$ $14.31^{14.63}_{14.02}$ | $3.06^{3.29}_{2.84}$ $14.88^{15.19}_{14.59}$ |
| | RNA-Diverse | 64 | 3649 | $3.35^{3.46}_{3.23}$ $17.71^{17.77}_{17.66}$ | $\mathbf{3.76^{4.09}_{3.40}}$ $17.90^{18.05}_{17.76}$ | $6.67^{7.16}_{6.32}$ $\mathbf{16.32^{16.42}_{16.23}}$ | $6.77^{7.26}_{6.32}$ $17.27^{17.37}_{17.17}$ | $6.47^{6.96}_{5.99}$ $18.52^{18.67}_{18.38}$ | $6.50^{7.04}_{5.97}$ $18.58^{18.73}_{18.43}$ |
| | RNA-Trinucleotide | 64 | 35134 | $3.04^{3.08}_{3.01}$ $17.72^{17.73}_{17.70}$ | $\mathbf{3.53^{3.64}_{3.42}}$ $17.93^{17.98}_{17.88}$ | $6.10^{6.26}_{5.93}$ $\mathbf{16.44^{16.47}_{16.41}}$ | $6.12^{6.30}_{5.95}$ $17.42^{17.45}_{17.39}$ | $6.21^{6.38}_{6.04}$ $18.68^{18.73}_{18.63}$ | $6.19^{6.36}_{6.03}$ $18.75^{18.80}_{18.70}$ |

KEN TAKABA

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



**week 1**

**2023**

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

using published force field model

**week 2**

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

using the <span style="color:red">same</span> published force field model!
we haven't learned anything from the data

**week 1**

**2025**

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

using force field model
built from public + private data

**week 2**

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

using <span style="color:red">new</span> model tuned to target
from first week's data

# HYBRID PHYSICAL / MACHINE LEARNING MODELS COULD DRIVE A NEW ERA OF PRODUCTIVITY IN COMPUTATIONAL CHEMISTRY



- Fast, structure-based **machine learning surrogates** assess designs over vast synthetic chemical spaces prioritize useful calculations

- Adaptive allocation of effort to alchemical free energy calculations guided by **machine learning cost prediction**s

- **Machine learned optimal alchemical transformations** produce faster estimates of free energy differences more cheaply

- **Learnable machine learning potentials** fit to experimental free energy and quantum chemical data produce higher accuracy predictions

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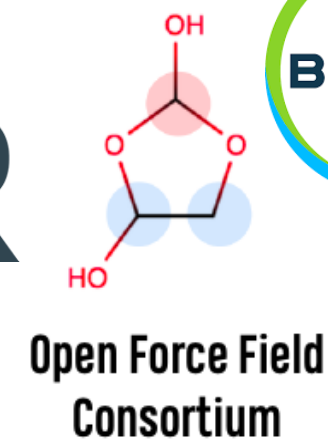**preprint:** https://doi.org/10.1101/2020.10.24.353318
**code:** https://github.com/choderalab/neutromeratio

# CHODERA LAB



National Institutes of Health · NSF · OpenEye SCIENTIFIC · SILICON Therapeutics · PARKER INSTITUTE for CANCER IMMUNOTHERAPY · STIFTUNG CHARITÉ · OpenEye SCIENTIFIC · EINSTEIN Foundation.de · Gerstner FAMILY FOUNDATION · STARR CANCER CONSORTIUM · VIR · Open Force Field Consortium · BAYER · XtalPi · CYCLE FOR SURVIVAL

- All funding: http://choderalab.org/funding