# INTEGRATING HIGH-PERFORMANCE SIMULATIONS AND LEARNING TOWARD IMPROVED CANCER THERAPY

**AUSTIN CLYDE**

*Ph.D. Student, University of Chicago*
*Computational science, Argonne National Laboratory*

Austin Clyde[1*], Dave Wright[2*], Katya Ahmed[2], John D. Chodera[3], John Karanicolas[4], Palani Kirubakaran[4], Hyungro Lee[5], Matteo Turilli[5], Shunzhou Wan[2], Peter Coveney[2**], Shantenu Jha[5,6**], Rick Stevens[1**]

[1]*University of Chicago and Argonne National Laboratory*
[2]*University College London*
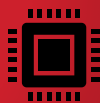[3]*Memorial Sloan Kettering Cancer Center*
[4]*Molecular Therapeutics group, Fox Chase Cancer Center*
[5]*RADICAL, ECE, Rutgers University, Piscataway,NJ 08854, USA*
[6]*Brookhaven National Laboratory, Upton, NY, USA* *Joint First Authors* **Senior Authors*

**DRUG DISCOVERY**

~$10^8$ products
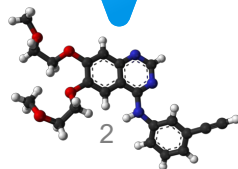
**PRE CLINICAL**

11,000 products

**CLINICAL TRIALS**

6,300 products

**FDA APPROVAL**

111 products

Argonne
NATIONAL LABORATORY

# Target based compound screening

$10^{60}$ estimated drug-like compounds

## COMPOUND DISCOVERY

Mining massive building block or de-novo generated libraries
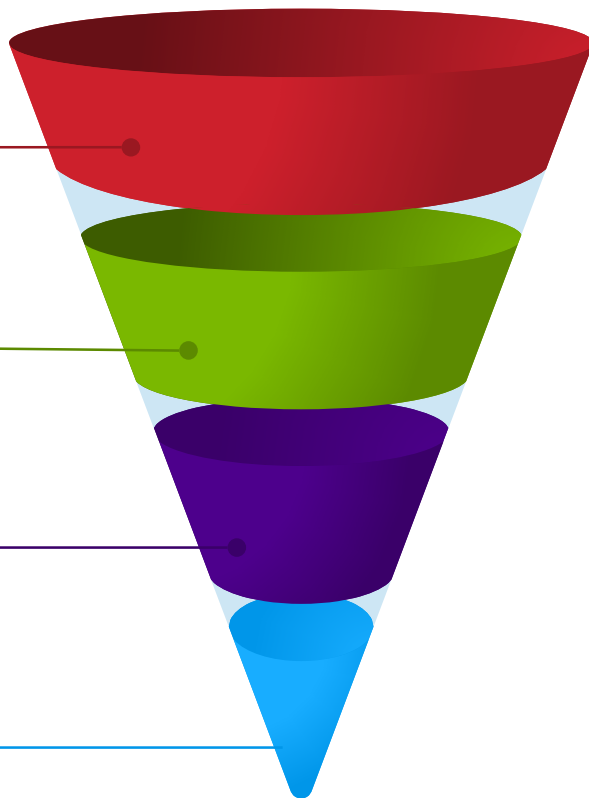
## INTERESTING?

Does this compound inhibit or interact with the target?

## TOXICOLOGY

Is this compound reasonably safe?

## SYNTHESIS

Can we buy it, is it from available building blocks, or do we need to hire a medicinal chemist?

Argonne
NATIONAL LABORATORY

# GOAL:

Design an intelligent system to screen a space of drugs efficiently and intelligently.

## LETTER

### Anthropogenic biases in chemical reaction data hinder exploratory inorganic synthesis

Xiwen Jia[1], Allyson Lynch[1], Yuheng Huang[1], Matthew Danielson[1], Immaculate Lang'at[1], Alexander Milder[1], Aaron E. Ruby[1], Hao Wang[1], Sorelle A. Friedler[2]*, Alexander J. Norquist[1]* & Joshua Schrier[1,3]*
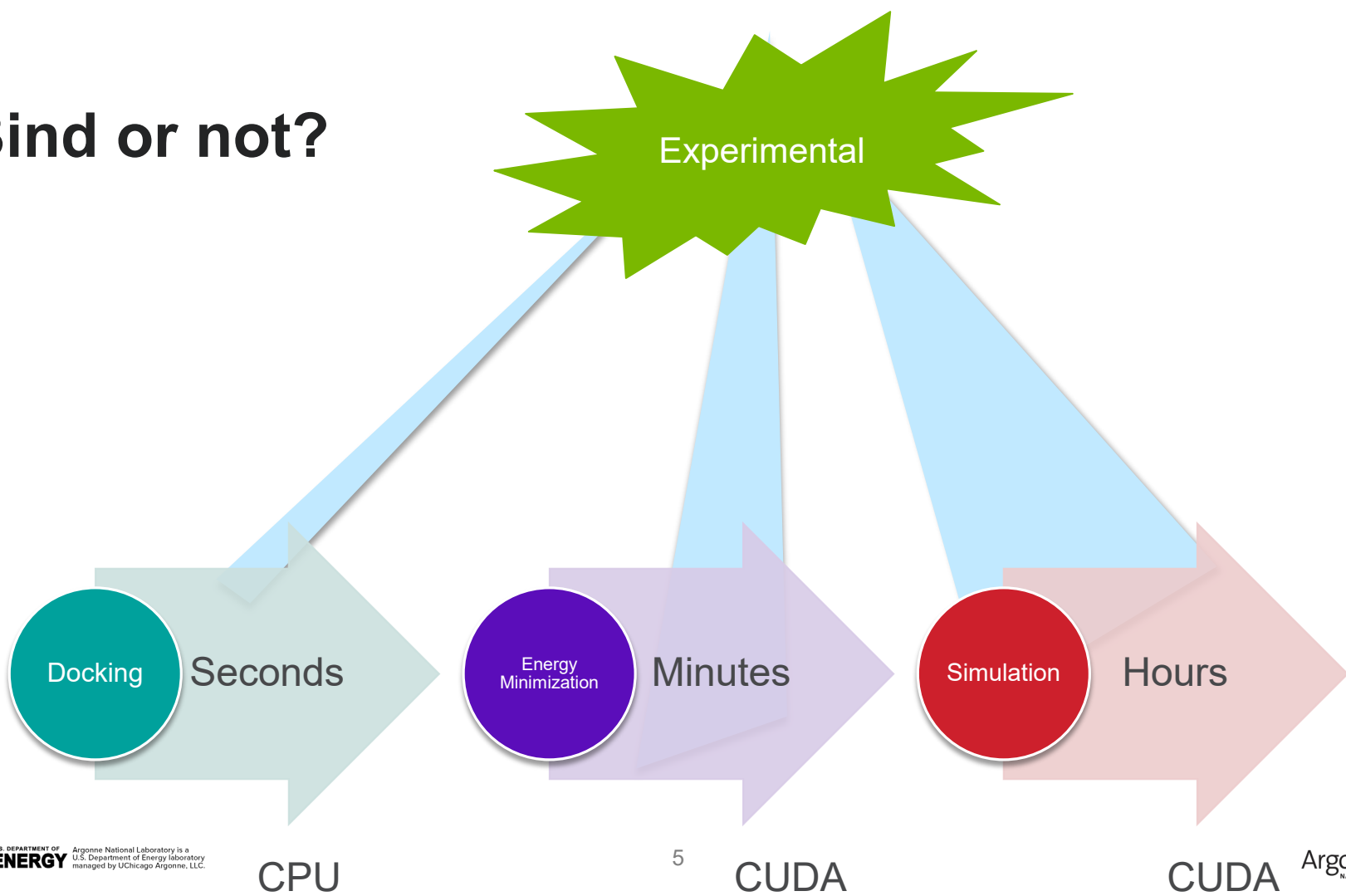
## ARTICLE

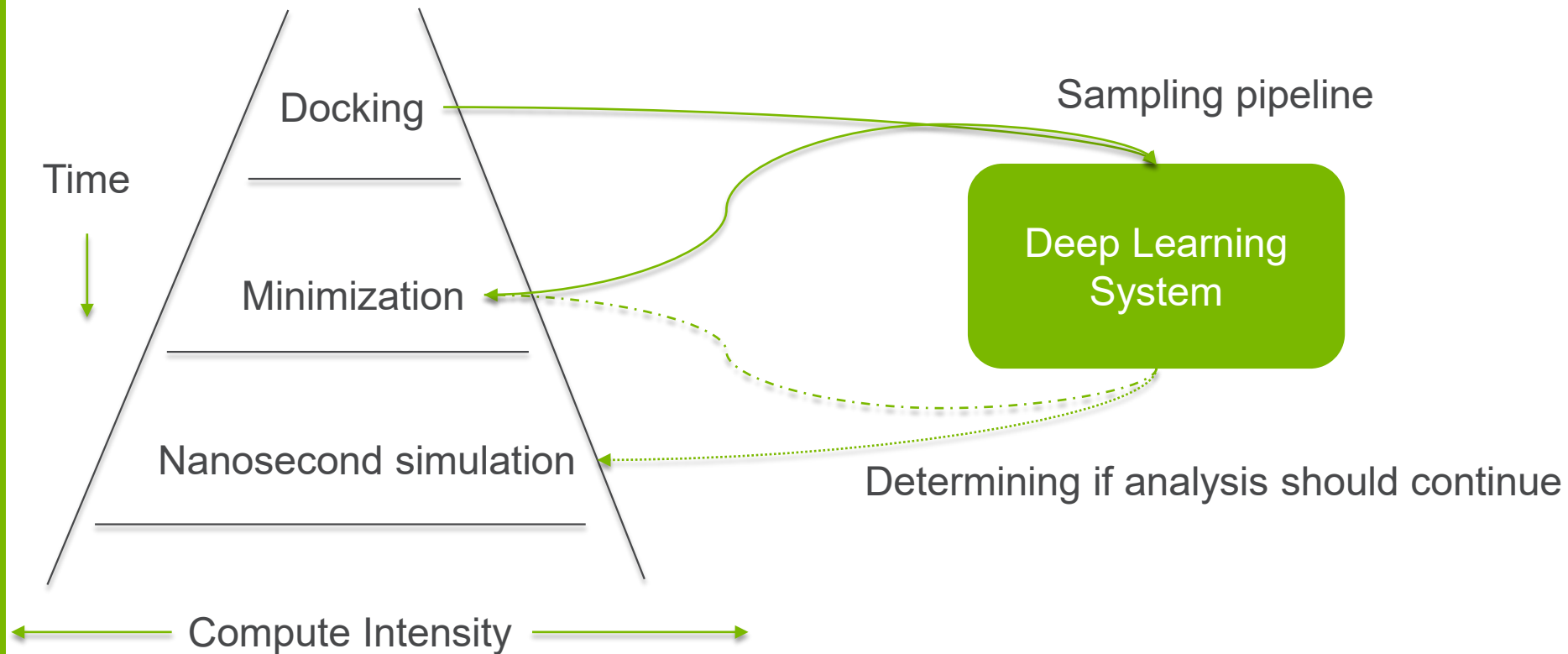### Ultra-large library docking for discovering new chemotypes

Jiankun Lyu[1,2,10], Sheng Wang[3,4,10], Trent E. Balius[1,10], Isha Singh[1,10], Anat Levit[1], Yurii S. Moroz[5,6], Matthew J. O'Meara[1], Tao Che[4], Enkhjargal Algaa[1], Kateryna Tolmachova[7], Andrey A. Tolmachev[7], Brian K. Shoichet[1]*, Bryan L. Roth[4,8,9]* & John J. Irwin[1]*
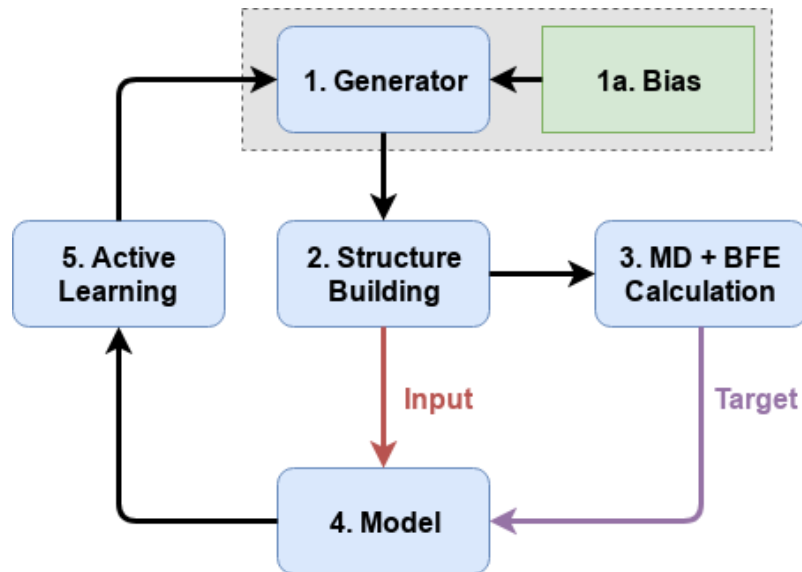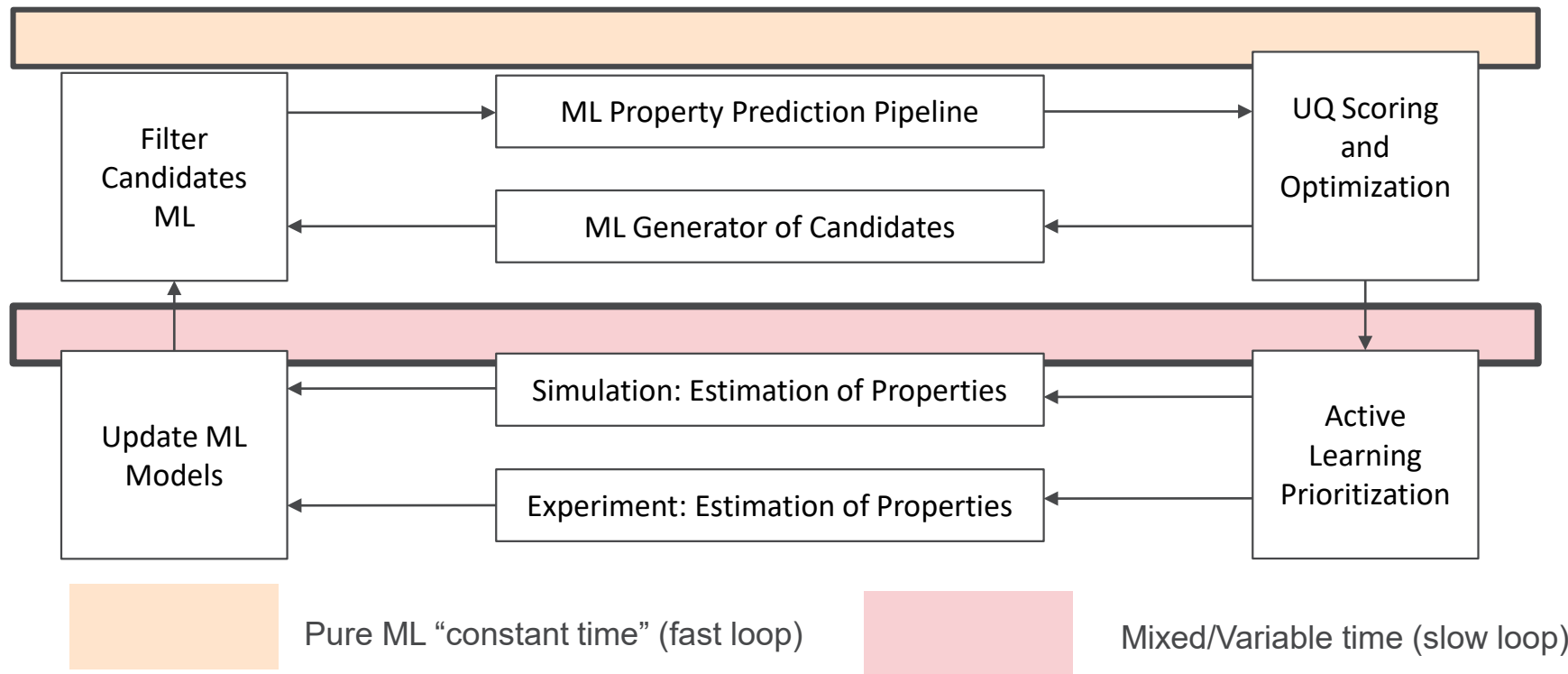
# Bind or not?

Experimental

Docking
Seconds

Energy Minimization
Minutes

Simulation
Hours

CPU

CUDA

CUDA

Argonne
NATIONAL LABORATORY

A pipeline unit

Docking

Sampling pipeline

Time

Deep Learning
System

Minimization

Nanosecond simulation

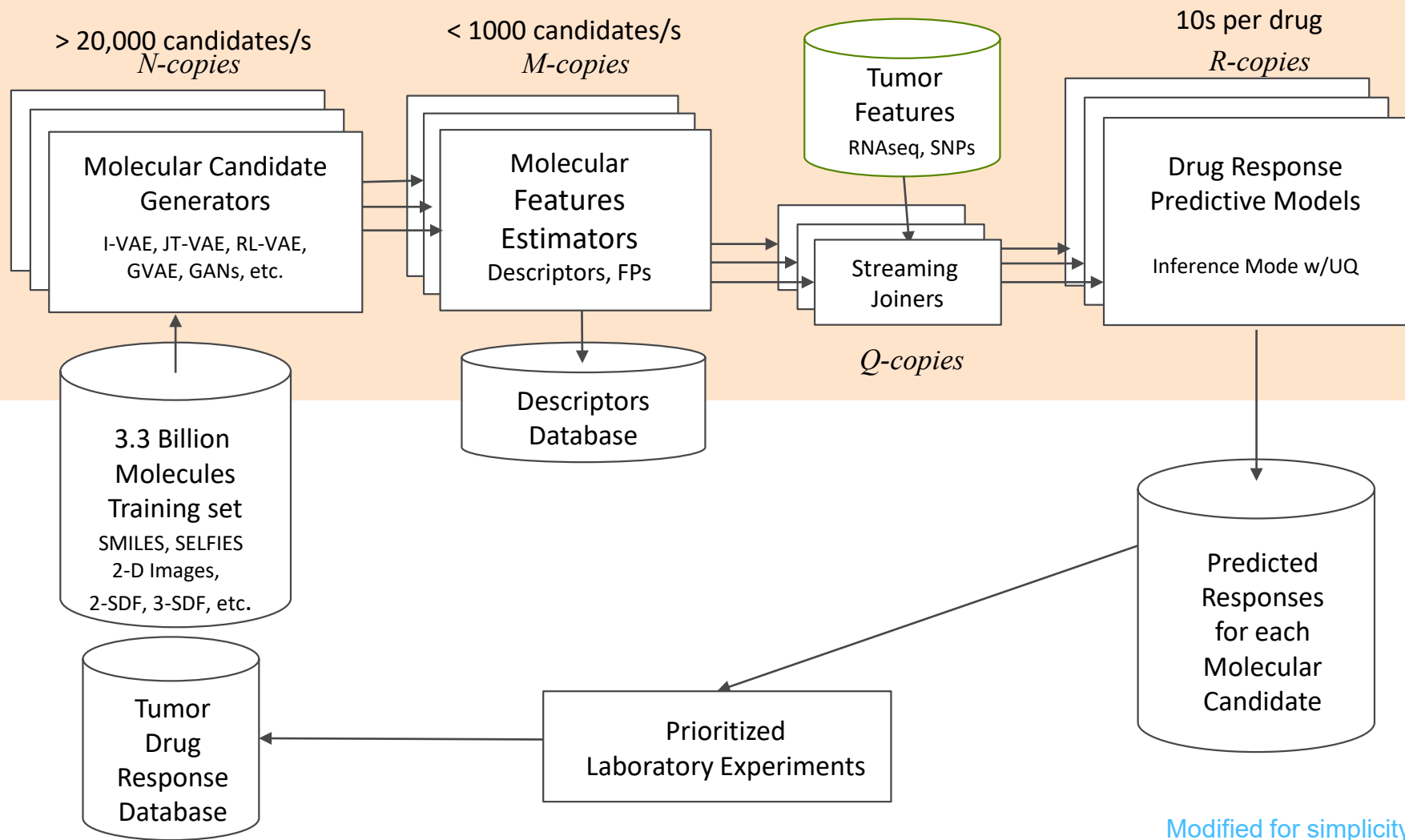Determining if analysis should continue

Compute Intensity

Argonne
NATIONAL LABORATORY
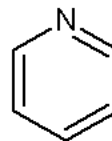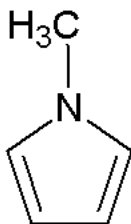
# Pipelining discovery and screening
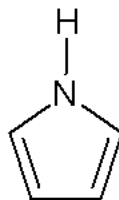
# LAYERED WORKFLOW
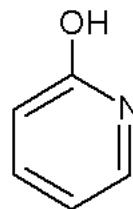
Modified for simplicity*

c1ccco1

O=C1C=CN=CN1

c1ccccn1

Cn1cccc1

c1cccn1

Oc1ccccn1

# PROPERTY PREDICTIONS

## Images, 3D surfaces

Feinberg, Evan N., et al. "Potentialnet for molecular property prediction." *ACS central science* 4.11 (2018): 1520-1530.
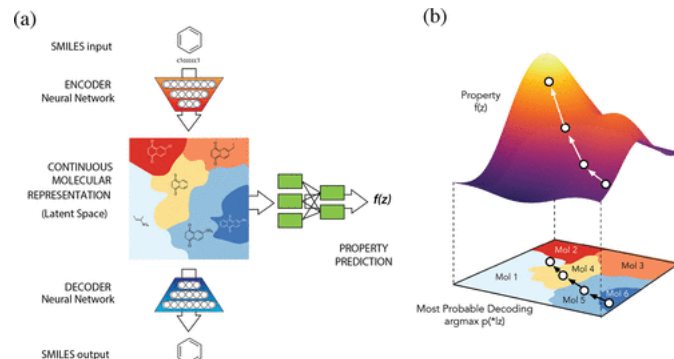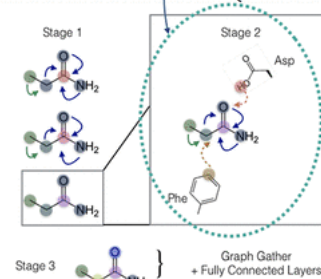


Skalic, Miha, et al. "Shape-Based Generative Modeling for de Novo Drug Design." *Journal of chemical information and modeling* 59.3 (2019): 1205-1214.
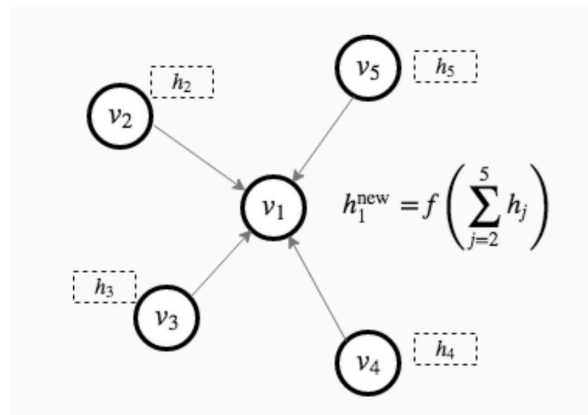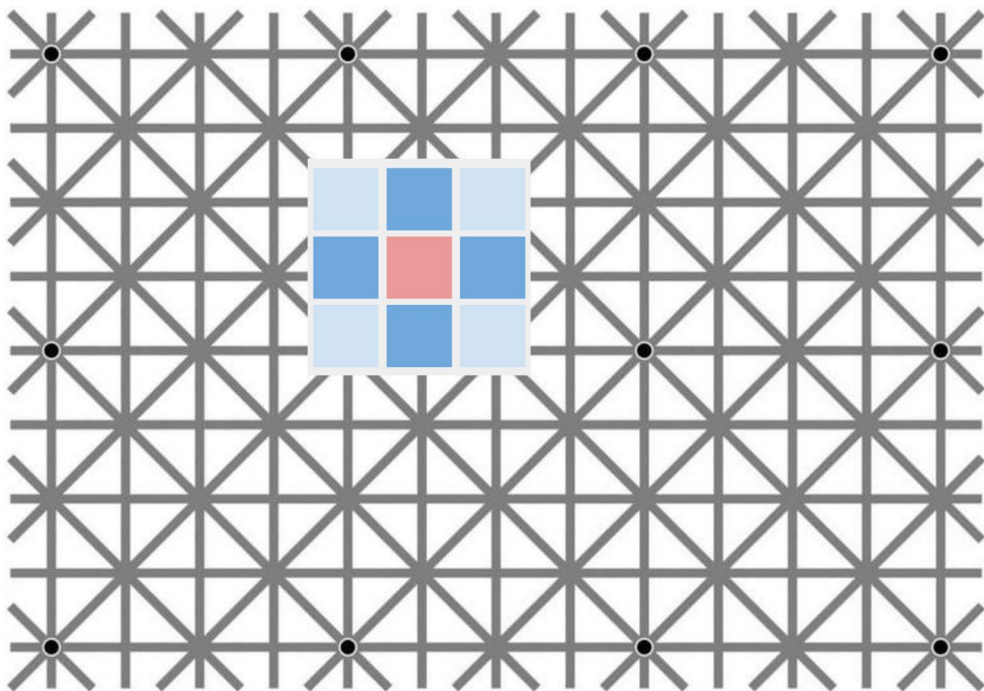
Gómez-Bombarelli, Rafael, et al. "Automatic chemical design using a data-driven continuous representation of molecules." ACS central science 4.2 (2018): 268-276.

$$h_1^{\text{new}} = f\left(\sum_{j=2}^{5} h_j\right)$$

Learning Cuve for SMILES AURORA-1 Kinase MMGBSA Minimization Score
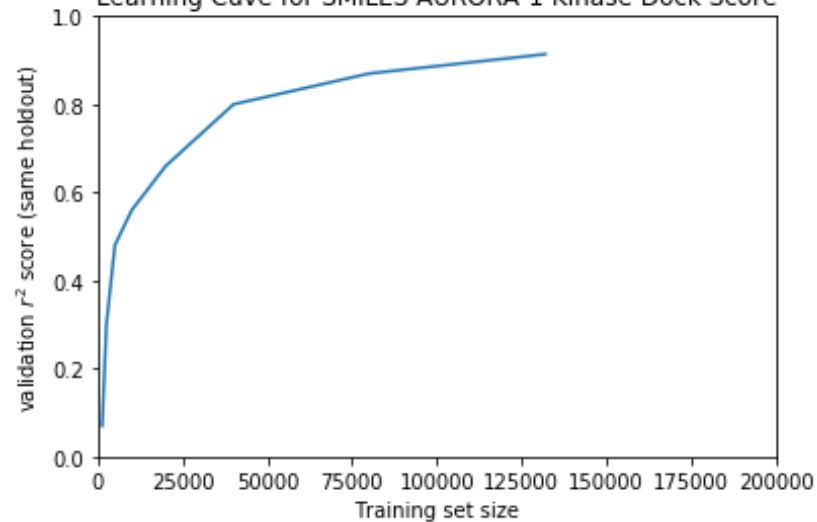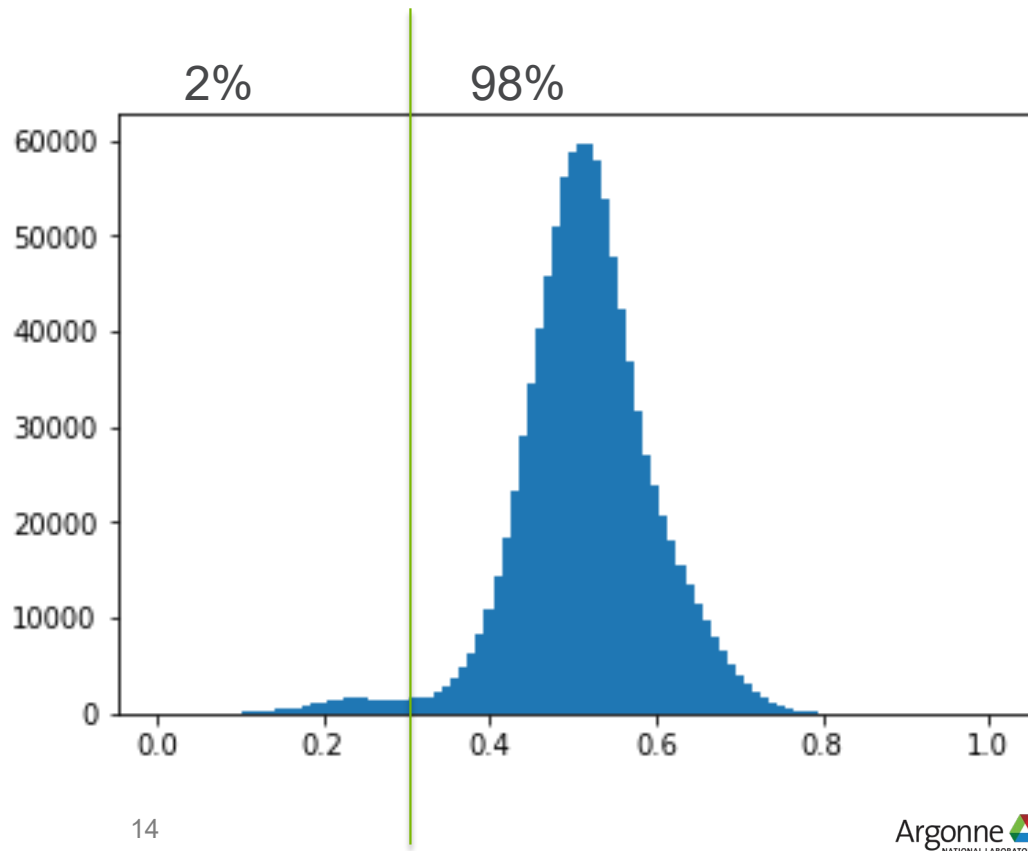
Learning Cuve for SMILES AURORA-1 Kinase Dock Score

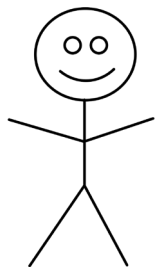# EXAMPLE: ML FOR DOCKING SCORING

## Interested in the left tail

What is r2 score if we just guess everything in that right tail is clipped at the normal distribution? 0.75

Your balanced accuracy? 50%

- Each experiment cost $1,000
- Your boss wants to find leads at the very early stages.

Here is $100,000, find five interesting beta lactamase inhibitors

10 experimental data points, randomly

10 experiments he should run

Alice, the ML hacker

Bob runs 20 experiments, cost $20,000 –but he found 5 leads!

The r$^2$ value was 0.2

Alice, the ML hacker

Metrics measure distance in spaces, not real life goals, objectives
Dreams desires, etc! Especially, not on skewed distributions

- Each experiment cost $1,000
- Your boss wants to find leads at the very early stages.

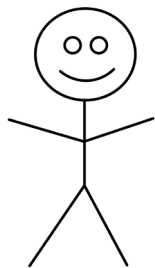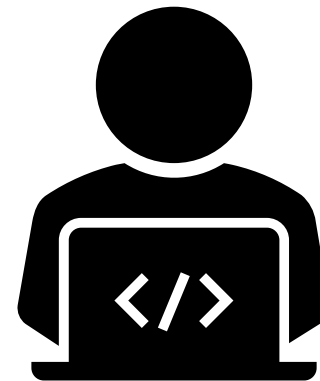Here is $100,000, find some interesting beta lactamase inhibitors
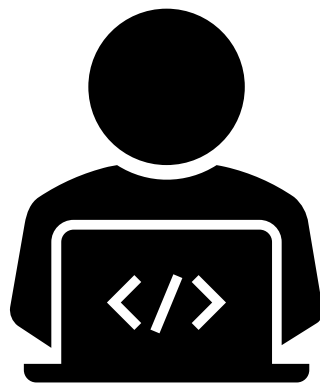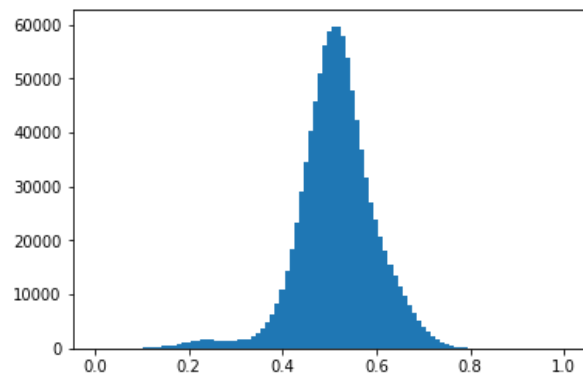
10 experimental data points, randomly

10 experiments he should run

Best x% of of your experimental values

Best x% of of your predicted values

$$EF_{x\%}^{(\text{COUNT})} = \frac{| \text{TopR}(y, x) \cap \text{TopR}(\hat{y}, x) |}{xN}$$

How many values?

What if we replace the need to simulate every molecule?

# Replicating Lyu et al. Giga-Docking with 200x less CPU compute

Trained message-passing network with 500K ampC



(ampC) Regression Detection Surface

Screen 1% of molecules, you'll
Have 50% of the true top 1%

Screen 1% of molecules, you'll
Have 70% of the true top 0.05%

Screen 10% of molecules,
get all of the top 0.1%

*preliminary work, first approximation of a good model

Argonne
NATIONAL LABORATORY

# GPU

10,000 per second

100 s/s  100 s/s  100 s/s  100 s/s

**Super fast, modern generative algorithms**

**Single threaded algorithms for CPU post-processing**

**Even slower simulations**

IBM AC922, 6 GPU node. Balanced Heavily towards GPU, not CPU
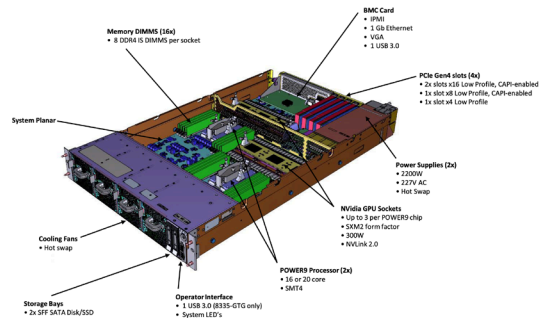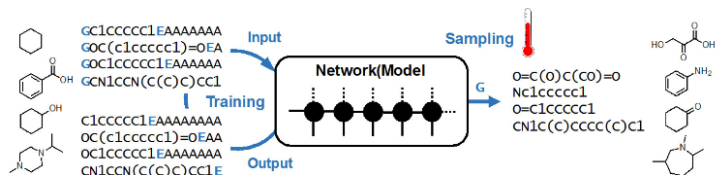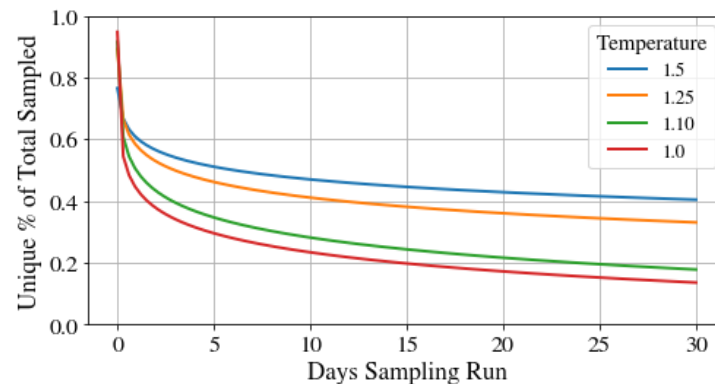
5000 Seconds per smiles

1 SMILE per second
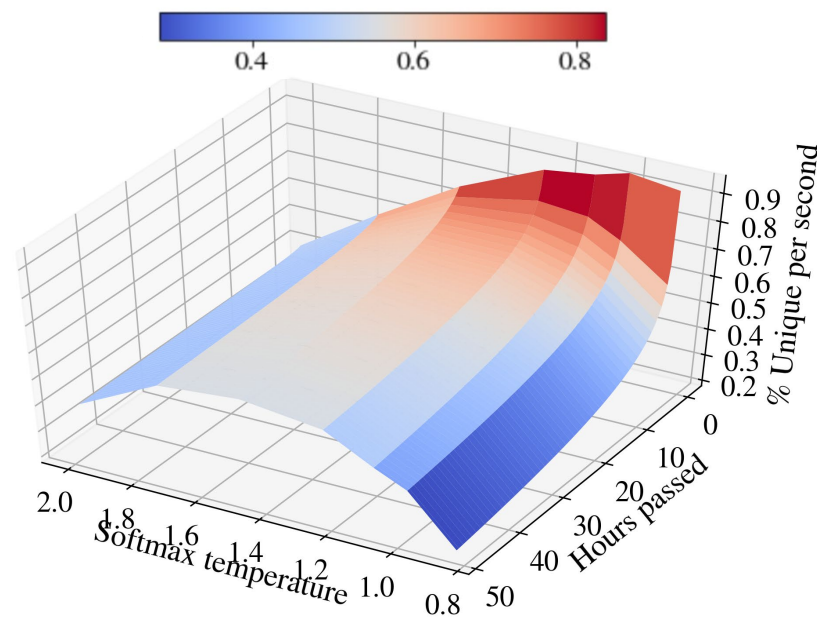
Argonne
NATIONAL LABORATORY

# RNN SMILES Modeling



Gupta, Anvita, et al. "Generative recurrent networks for de novo drug design." Molecular informatics 37.1-2 (2018): 1700111.



Sampling RNN Generator on 4 V100 GPU (first approx., T=1.1)



Samples from RNN on single GPU (<6 minutes)



(Predicted) Unique Molecules as a % of Sample Rate

Argonne NATIONAL LABORATORY

# In order to keep GPUs and CPUs hot, unique stream of molecules needs to stay constant

Database, Experiments

Top 0.001%

Estimated Unique Molecules 27,600 V100 GPUs



$$p(s)_i = \frac{e^{-\beta s_i}}{\sum_{j=0}^{K} e^{-\beta s_j}}$$

Argonne
NATIONAL LABORATORY

# DRUG DISCOVERY

# HIGH THROUGHPUT SCREENING



Generating Drug Leads

Database of Leads

- Generative Neural Networks
- Language modeling
- Graphical models

- Simulation surrogate models
- Uncertainty calibrated
- Ranking Neural networks

| High-Throughput Lab (HTL) | High Performance Computing (HPC) |
|---|---|
| Data Generation | Data Analysis |
| Biological Experiments | *In silico* Experiments |
| Hypothesis Testing | Novel Hypotheses |

Argonne
NATIONAL LABORATORY

# THANKS!

Argonne
NATIONAL LABORATORY