UGM
Boston
Oct. 4th 23

Chemaxon

**The World's Largest Protein-Ligand Complex and Binding Affinity Dataset for Data Driven Methods in Drug Design**

Dr. Arun Subramaniyan
Vice President Cloud & AI,
Strategy & Execution, Intel

# Thanks to the team that made the magic happen

**Alex Iankoulski**
Yusong Wang
Stephen Litster
Srinivas Tadepalli

**AWS**

Rakesh Srivastava
Prathit Chatterjee
Divya Korlepara
Vasavi C.S.
Suyash Gupta
Vishal Kumar
Pradeep Kumar Pal
Aathira Nair

Shivangi Verma
Harshini Anand
Saswati Mallick
Kavita Thakran
Parth Kanani
Indhu Ramachandran
Divya Nayar
**U. Deva Priyakumar**

**IITH**

Faird El Chalouhi
Chetan Rao
Arun Karthi Subramaniyan
Akanksha R. Bilani
Varma S. Konala
**Ramanathan Sethuraman**

**Intel**

**Vladimir Aladinskiy**
Evgeny Kirilin
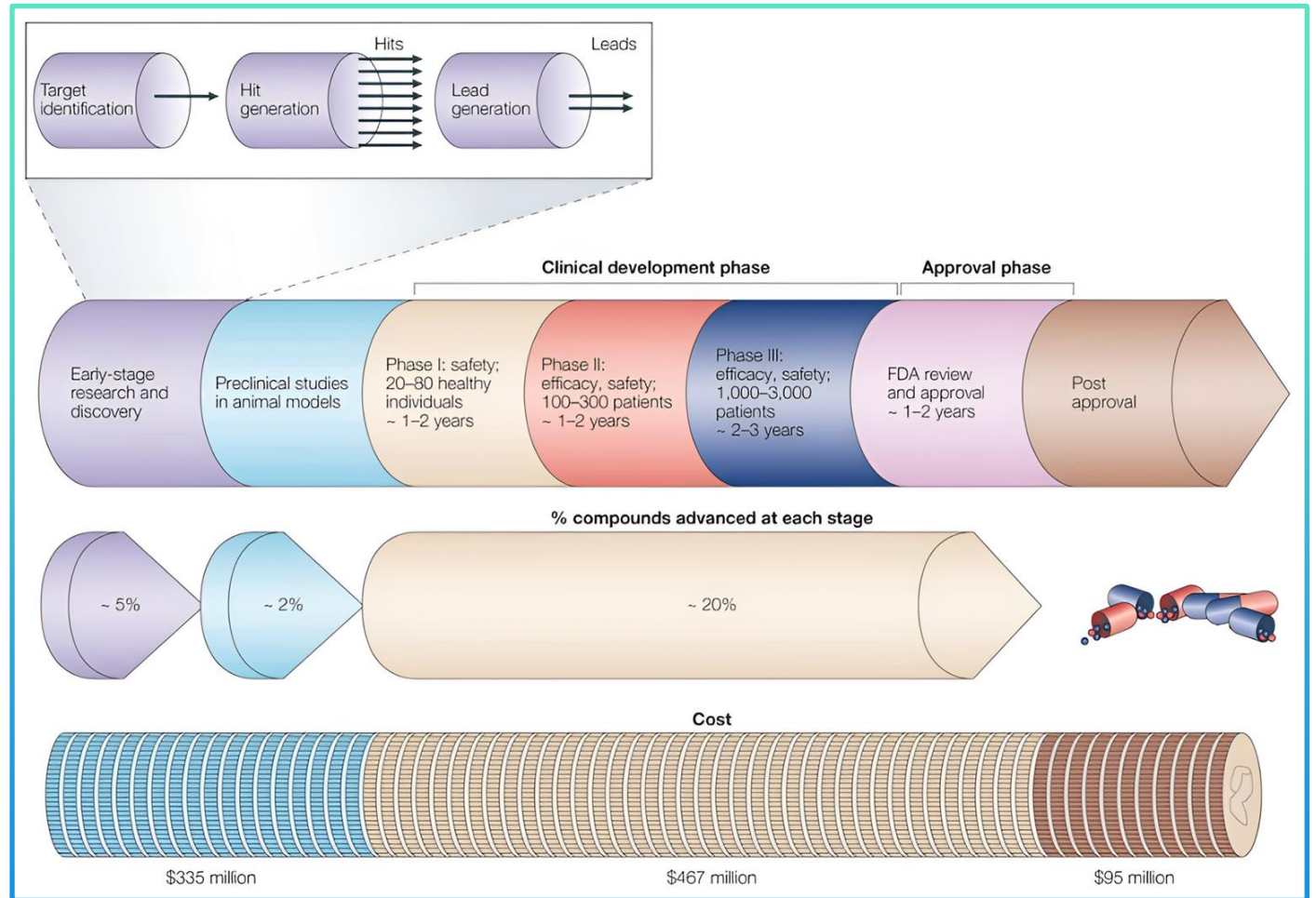Georgiy Andreev
Arkadii Lin
Eugene Babin

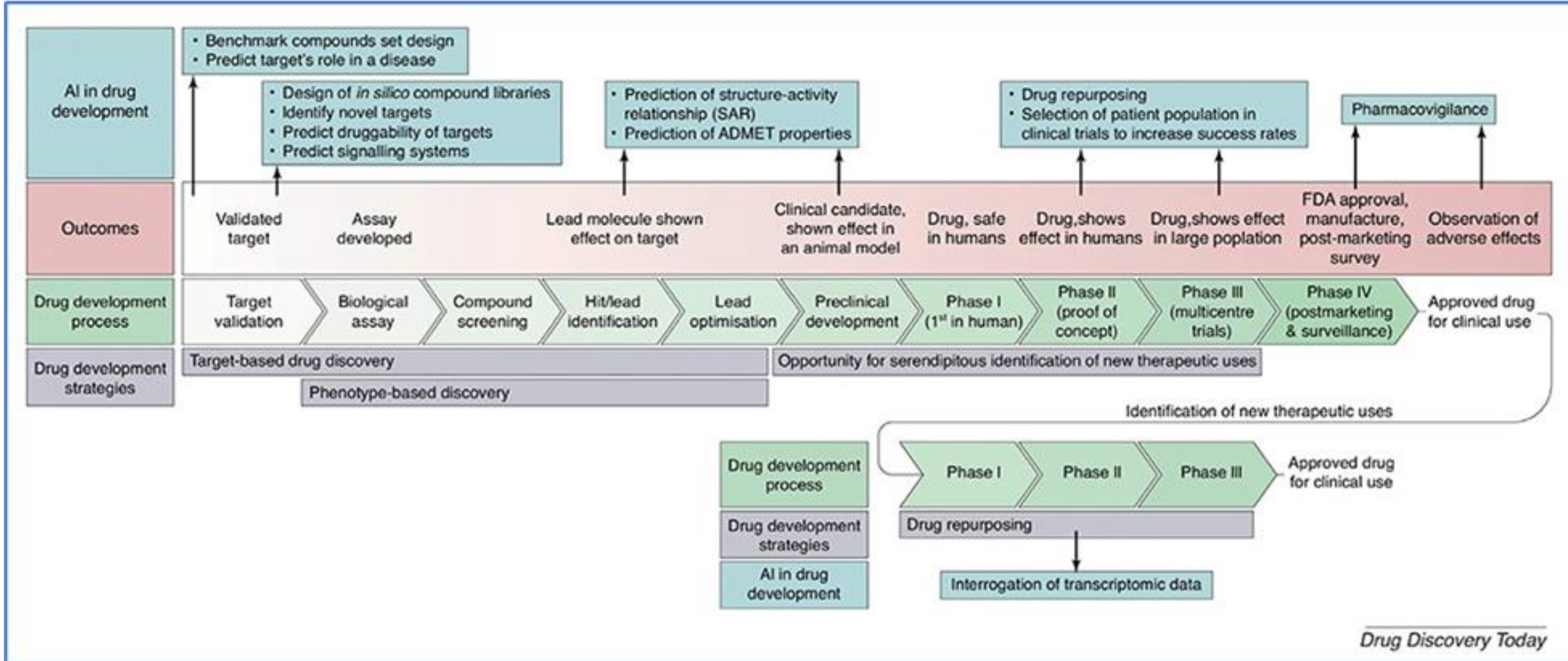**Insilico Medicine**

AI³

# Introduction

# Inspiration

- Traditional drug discovery is costly, time-consuming, and has a low success rate.

- Computational techniques are crucial for revolutionizing drug discovery workflows.

- Recent advances in cloud computing and AI/ML could help accelerate drug discovery.



Traditional drug discovery workflow (dLab, 2023)

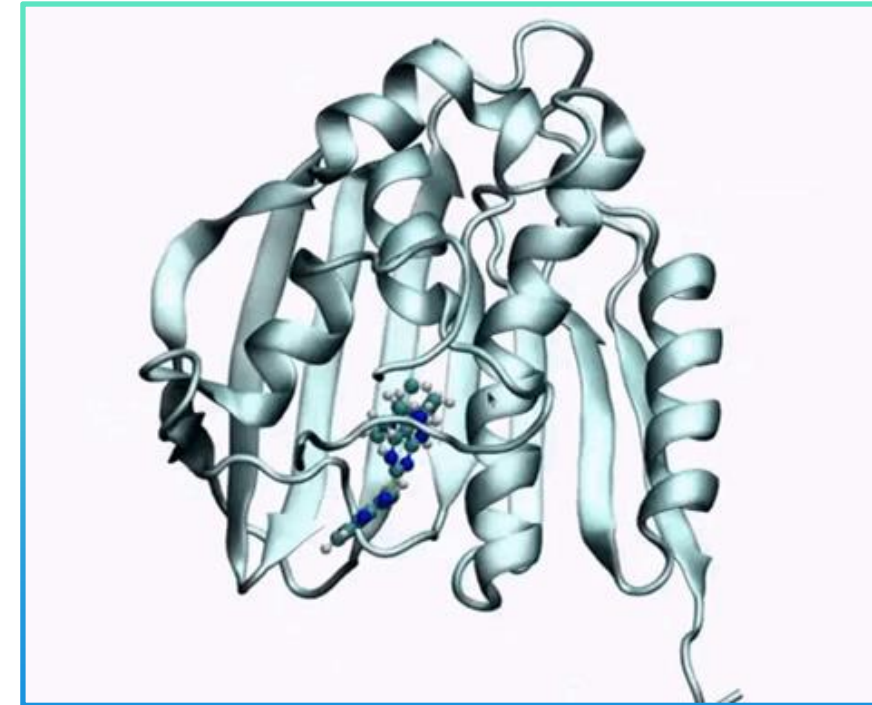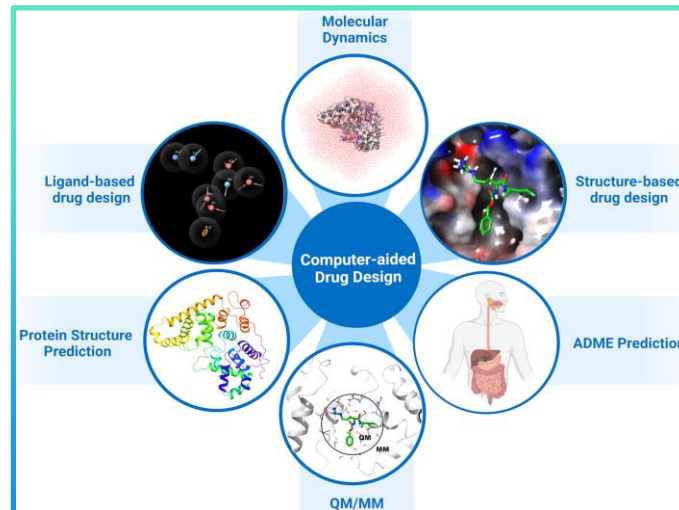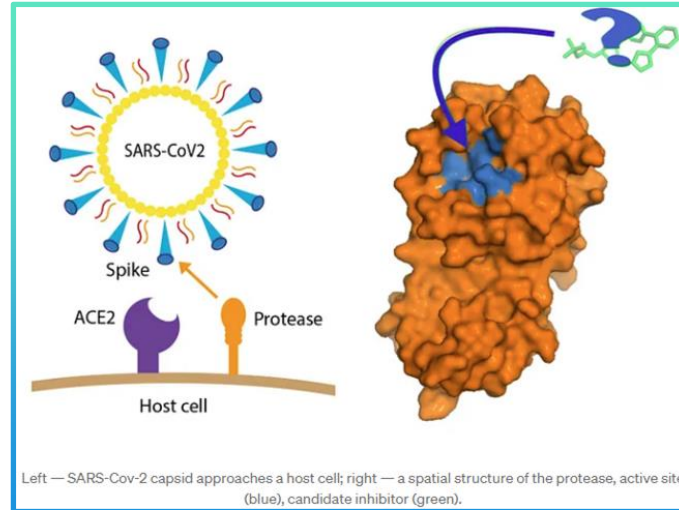# Potential Applications of AI in Drug Discovery Process



In-Jin, 2019, https://tcpharm.org/DOIx.php?id=10.12793/tcp.2019.27.3.87

# Importance of PLCs in Computer-aided Drug Design

# Protein-ligand Complexes (PLCs) & their role in Drug Design

- Proteins are essential biological molecules with diverse structures and functions.

- Ligands, or small molecules, can alter or assist protein structure and function by binding to proteins.

- Drug design often involves tuning druggable molecules to interact energetically with protein binding sites.

- Predicting binding affinity in PLCs is challenging but crucial for drug design.

- In-silico methods reduce production costs and enable the study of inaccessible molecular interactions



Left — SARS-Cov-2 capsid approaches a host cell; right — a spatial structure of the protease, active site (blue), candidate inhibitor (green).





A Guide to In Silico Drug Design, Chang, 2022

# Popular Existing Datasets & Limitations

PDBbind (2004 onwards, 23,496 PLC entries; http://www.pdbbind.org.cn/)

DUDE (2012 onwards, 22,886 active compounds; https://dude.docking.org/)

ONIONnet (2019 onwards; https://pubs.acs.org/doi/10.1021/acsomega.9b01997)

BindingDB (2007 onwards; https://www.bindingdb.org/rwd/bind/index.jsp)

AffinDB (2006 onwards; https://academic.oup.com/nar/article/34/suppl_1/D522/1132614)

Binding MOAD (2005 onwards, 41409 structures; http://www.bindingmoad.org/)
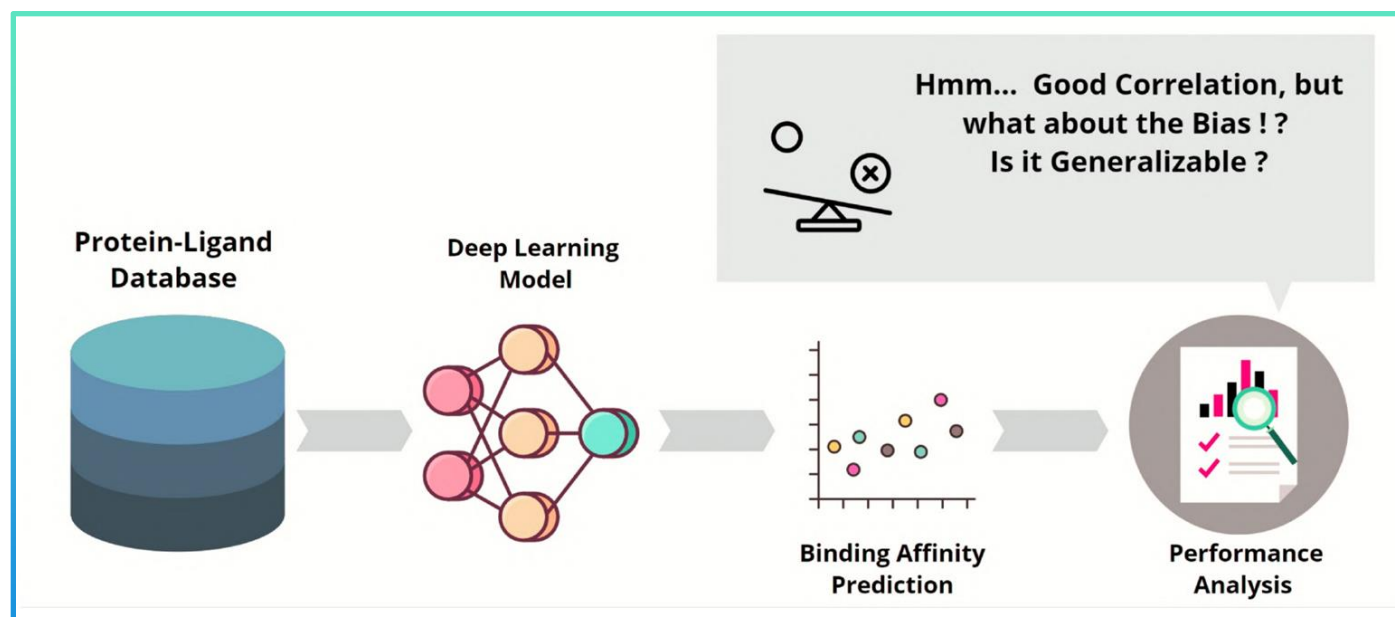
## Limitations

- Poor target and ligand diversity
- Low transferability to broad drug targets
- Lack of high-energy data (both structural & thermodynamic)
- Low volume
- High variability in validation quality – experimental errors from different labs & time

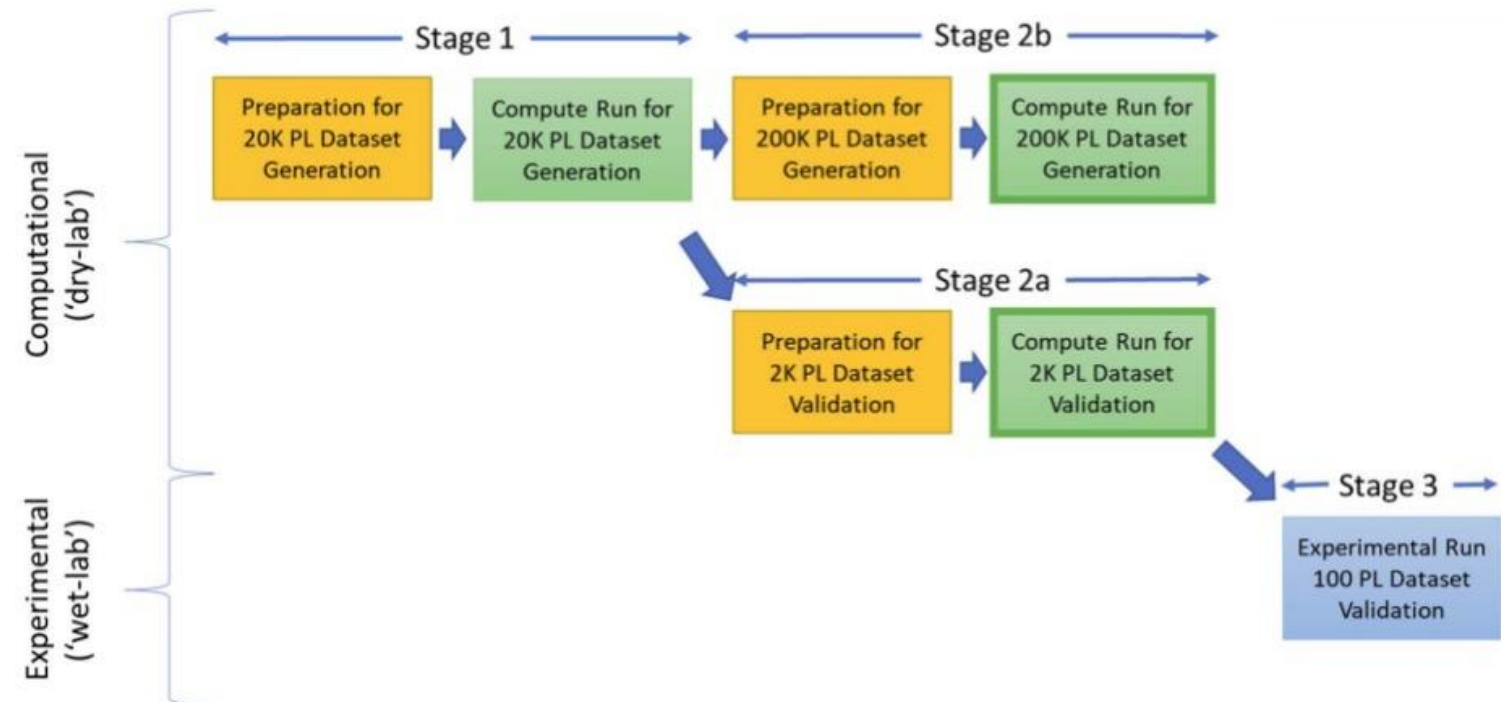# Motivation behind calculating PLC Binding Affinity

- Machine learning based scoring functions, for predicting binding affinity, have acceptable evaluation scores.

- Yet, they fail to perform similarly in virtual screenings.

- Hidden biases plausibly originate from data obtained from different experimental protocols.

- Inspiration to create a homogeneously computed Binding Free Energy of PLCs.



Latent Biases in Machine Learning Models for Predicting Binding Affinities Using Popular Data Sets, Kanakala, 2023

# AI³: Generate and Validate World's Largest PLC dataset

- The World's Largest Open PLC Dataset (AI³: OPLD) initiative aims to address corresponding dataset limitations.

- Collaboration between AWS, IIIT-H, INTEL, and Insilico Medicine.

- Phase 1 AI³ dataset, consists of ~20,000 PLCs and corresponding binding affinity.

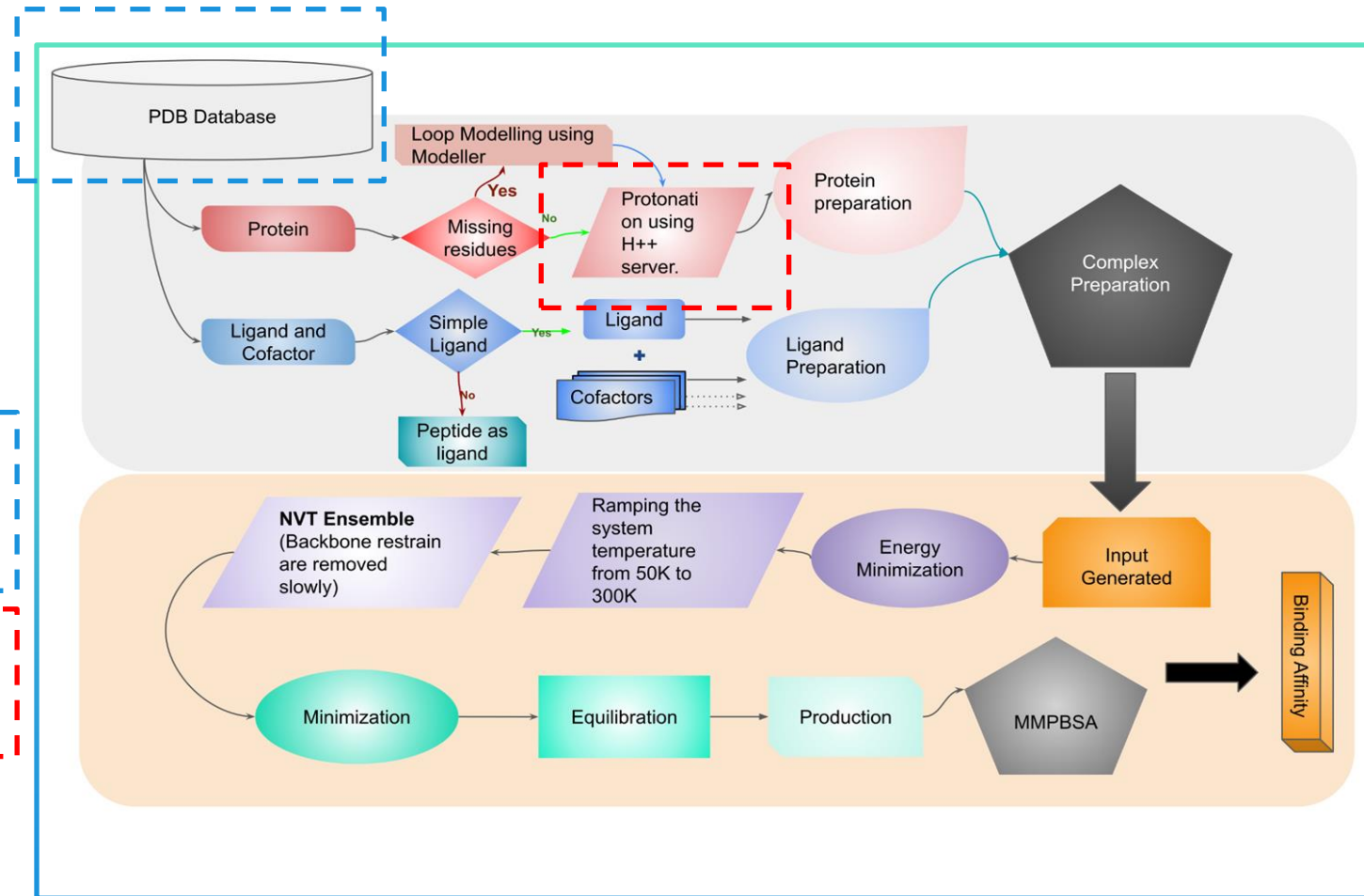- Goal: Create a dataset with ~220,000 entries, including the negative examples.

# Dataset Preparation

- The AI$^3$ (AWS-IIITH-Intel-Insilico) dataset is being prepared in two stages: Stage 1 (20 K bound PLCs) and Stage 2 (200 K unbound or partially bound PLCs).

- Protein structures are downloaded from RCSB PDB, and missing residues are modeled, modeler package UCSF Chimera.

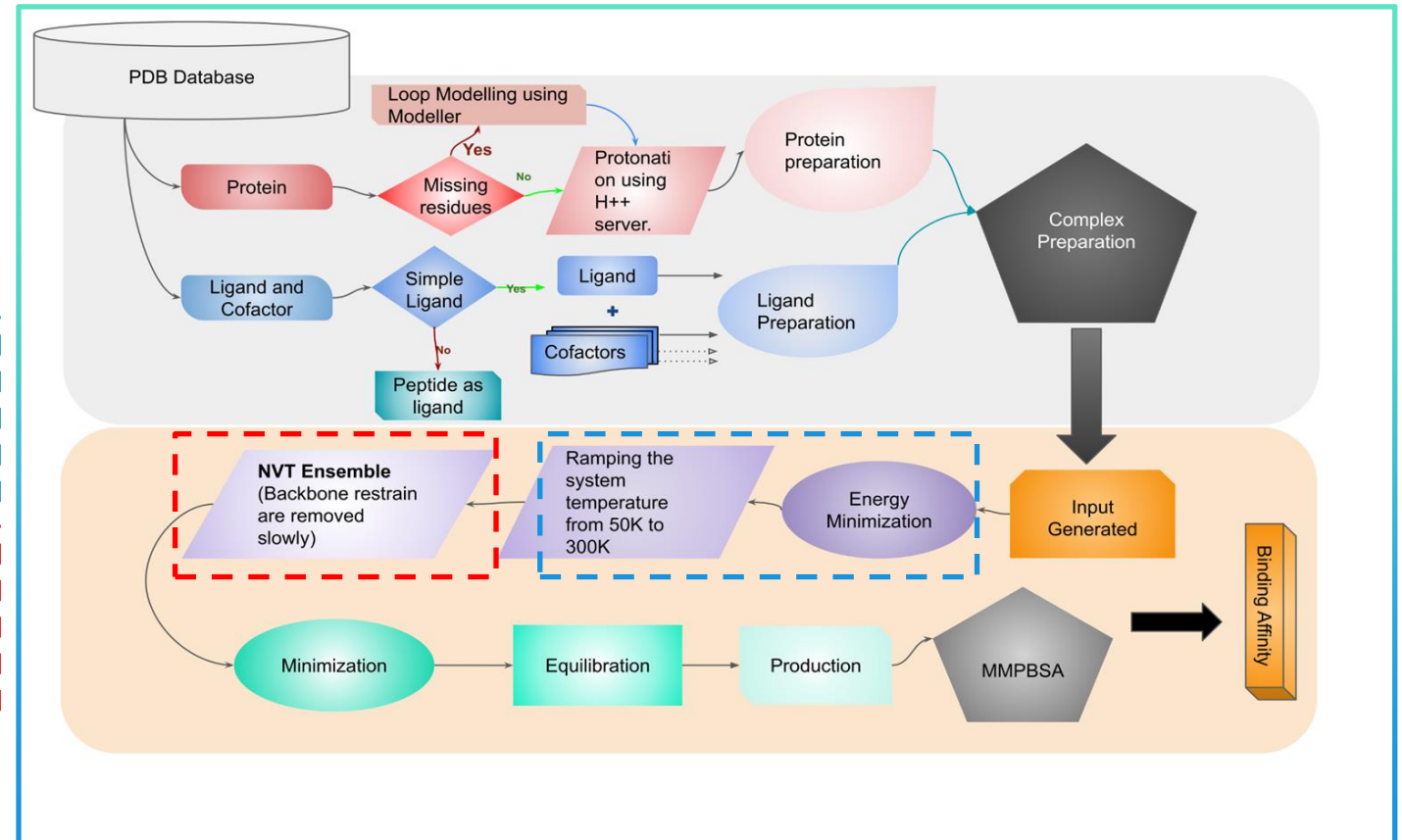- Protonation states are determined, from H++ server, at pH 7.4.

# MD simulations protocols

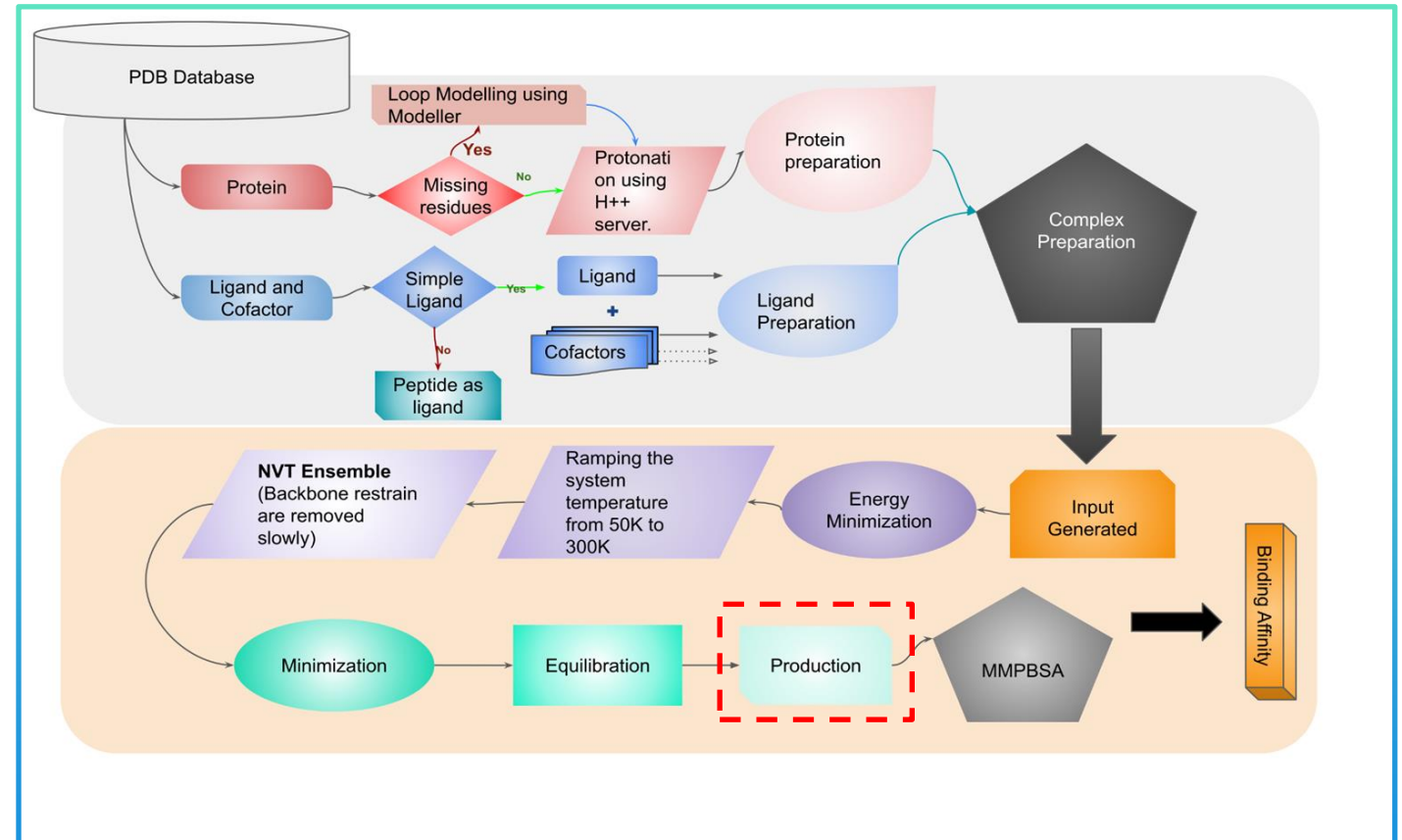- Amber ff1414 SB force field, TIP3P water model, GROMACS simulation package.

- Solvated PLC systems were subject to 2000 steps of steepest descent energy minimization and heating to 300 K.

- Backbone restrained were removed in a following NVT ensemble simulations, for 400 ps time length.
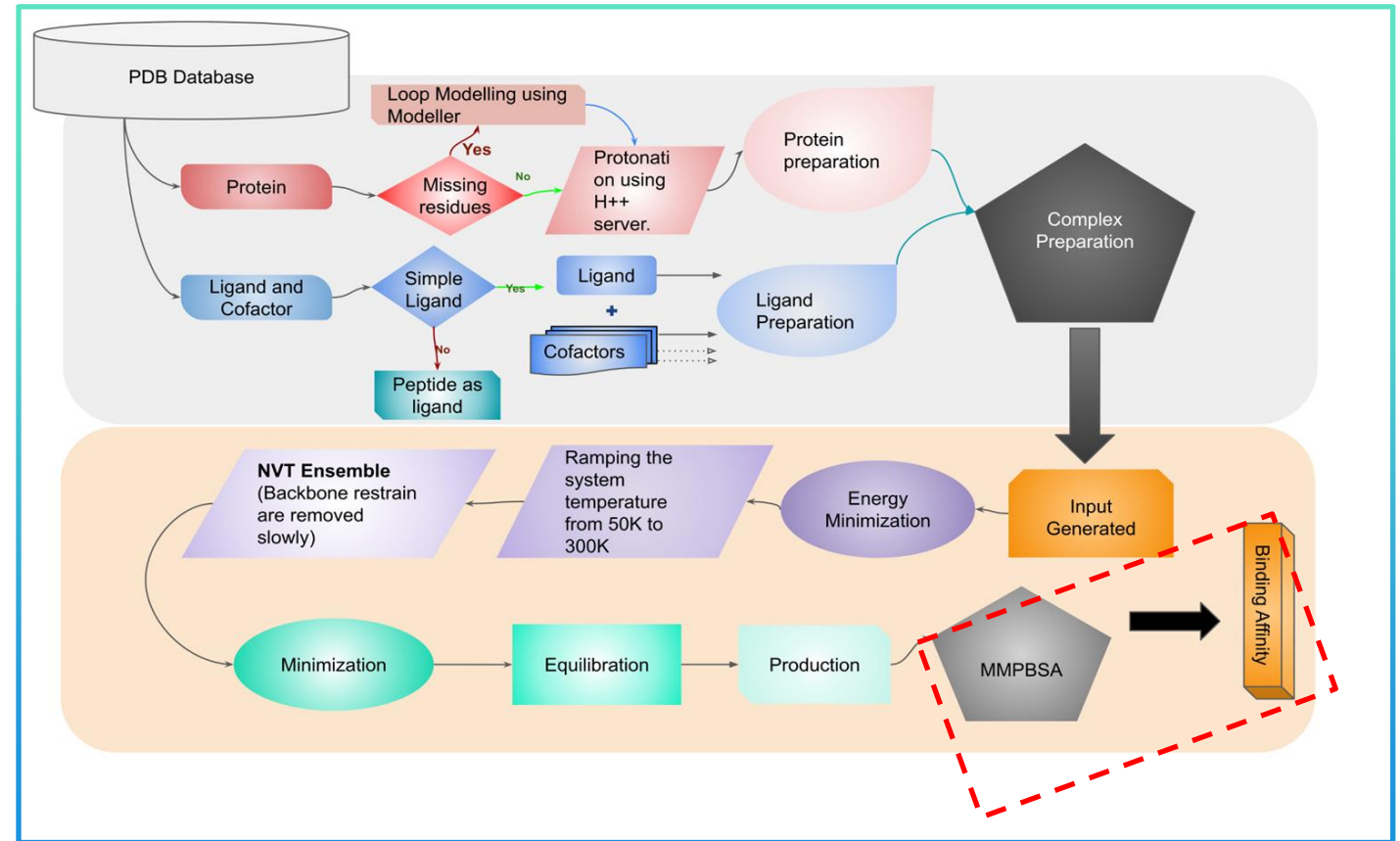
# MD simulations protocols

- Multiple short independent simulations are conducted to reduce uncertainty in predicting binding affinities.

- Each independent production MD runs are carried out under NPT conditions at 300K for 6 ns.

- Simulation frames are saved at regular intervals for analysis. The last 4 ns of data are used for binding free energy calculations.
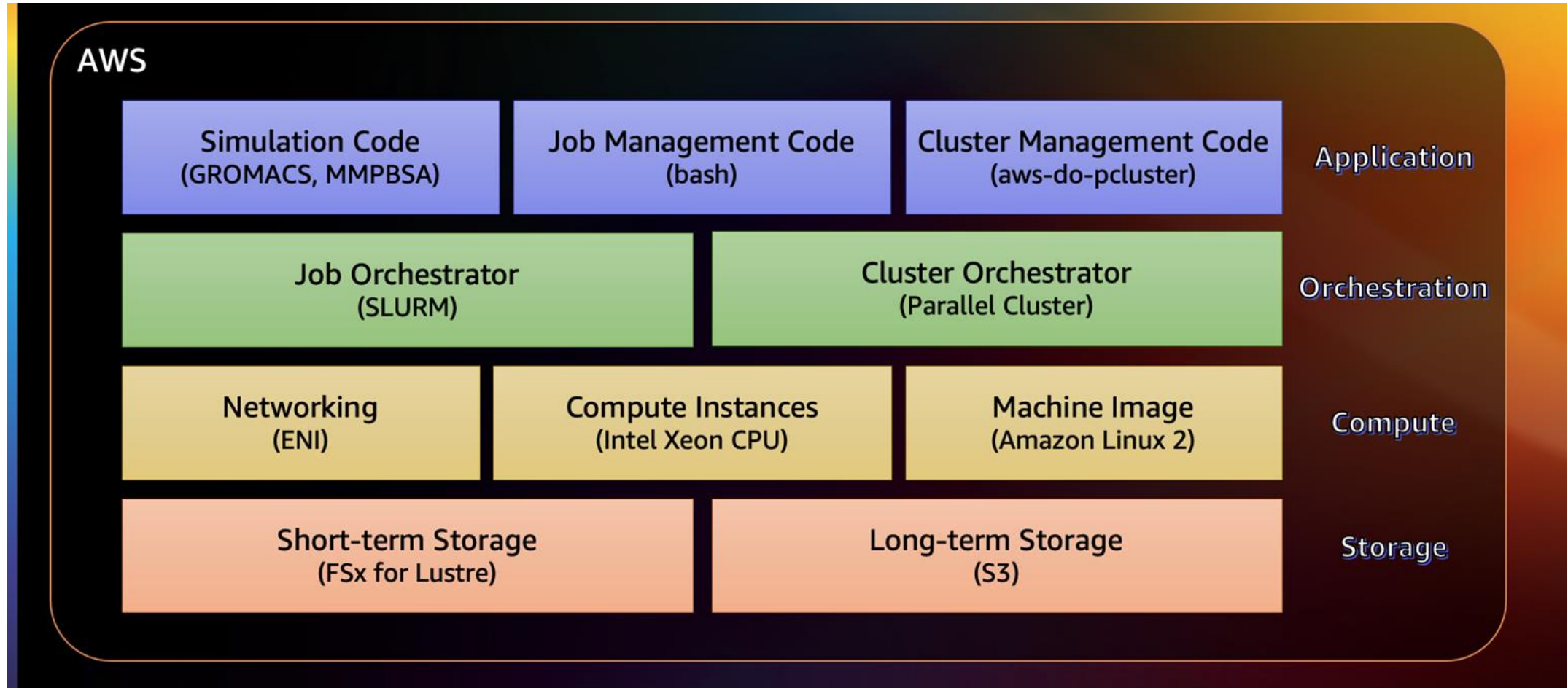
# Binding Free Energy Estimation

- The binding free energy is estimated using the molecular mechanics Poisson-Boltzmann surface area (MM-PBSA) approach.

- This method treats the solvent environment as a dielectric continuum. Polar and nonpolar solvation components are estimated.

- A single trajectory protocol is used to estimate binding affinities, ensuring robustness and accuracy.
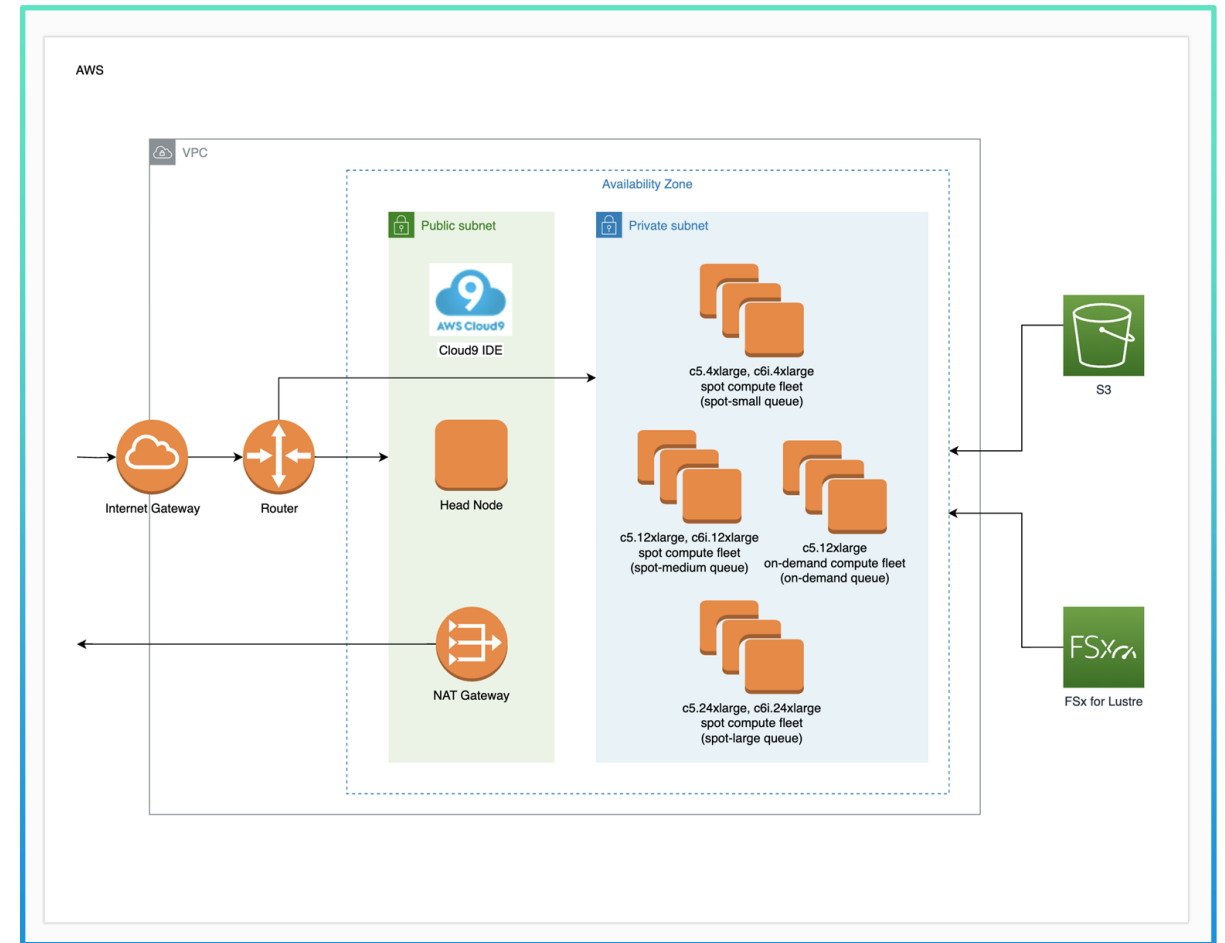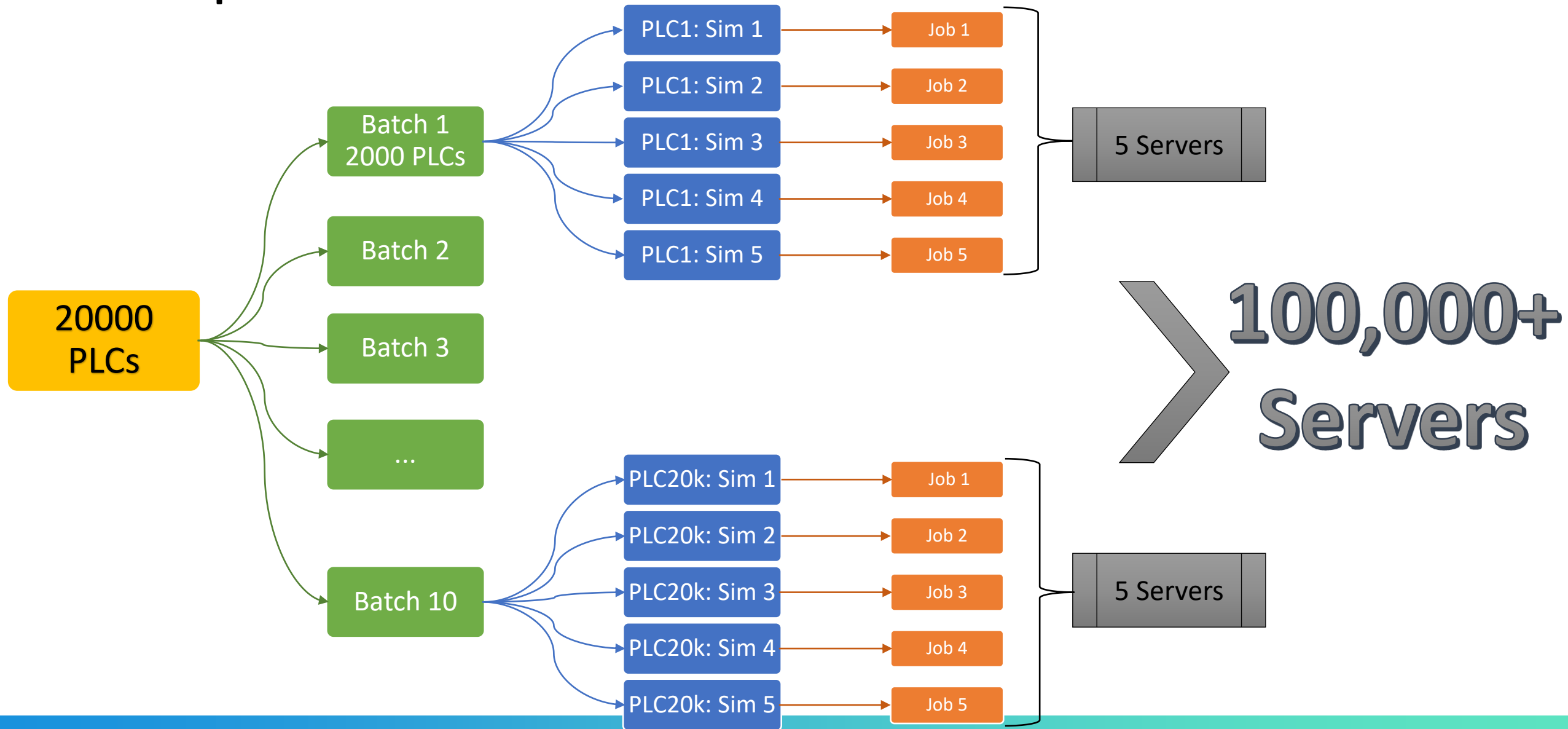
# Technology Stack and Architecture

# Execution and Deployment Architecture

| Lower Bound (Number of Atoms) | Higher Bound (Number of Atoms) | Number of CPU Cores | Time to Complete Job |
|---|---|---|---|
| 0 | 50,000 | 8 | 2-6 hrs |
| 50,001 | 100,000 | 24 | 6-8 hrs |
| 100,001 | 500,000 | 40 | 8-12 hrs |

# Leveraging AWS's "Planetary-scale" Computing Footprint

# Total Compute utilized … for Stage 1

**100,000+ Servers**

**2.3M+ CPU cores**

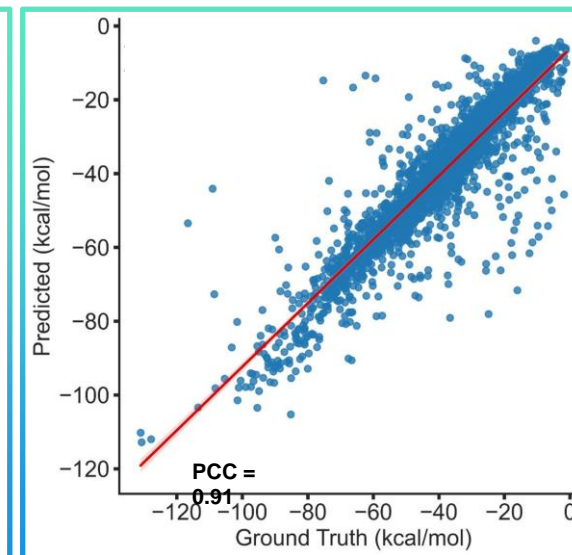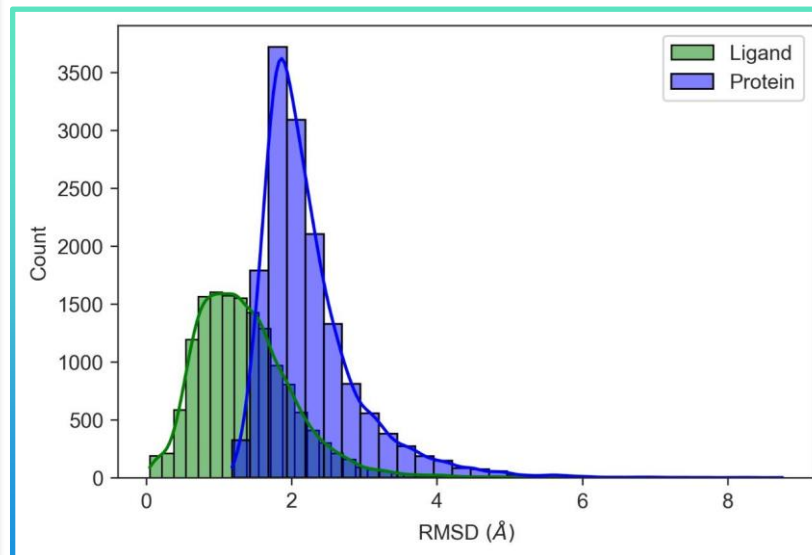**100M+ CPU core-hours**

20000 PLCs

3 Weeks

Calculated Observables

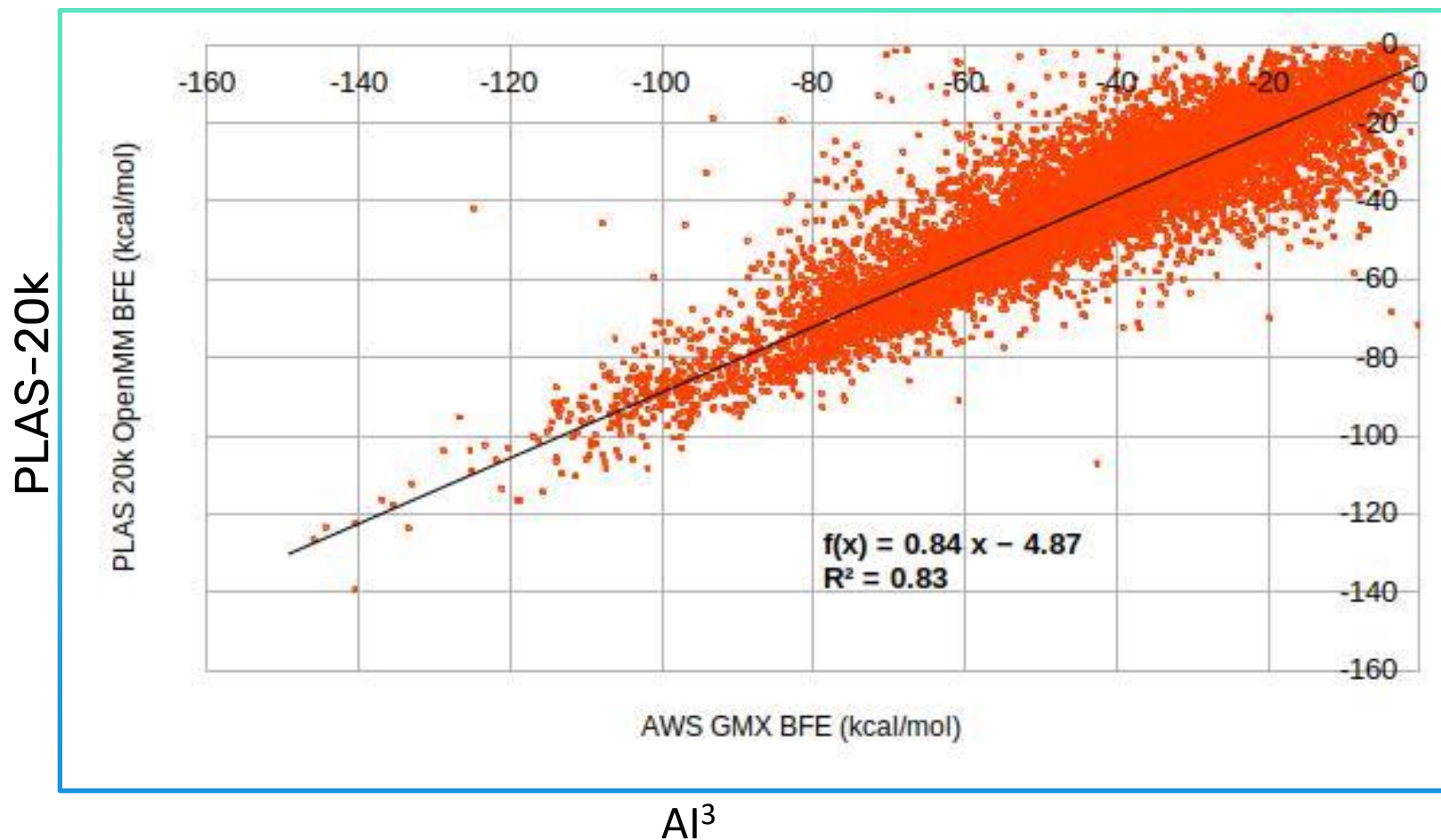# Stage 1: PLAS-20k as a Baseline Validated Dataset

**Training Deep-Learning Models with PLAS-20K**

- **Objective**: Accurate prediction of Protein-Ligand (PL) complex binding affinity.

- **Method**: Utilizing PLAS-20k dataset and deep learning model OnionNet.

-----

- **Results**: PLAS-20k achieves PCC of 0.91 (Strong correlation).

- RMSE of 8.15 kcal/mol (Accurate predictions).

-----

- **Significance**: PLAS-20k dataset is a powerful training resource.

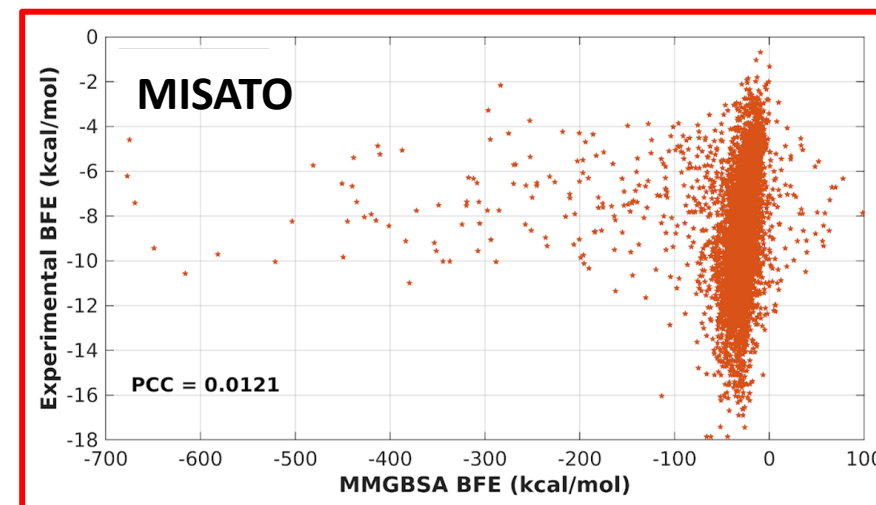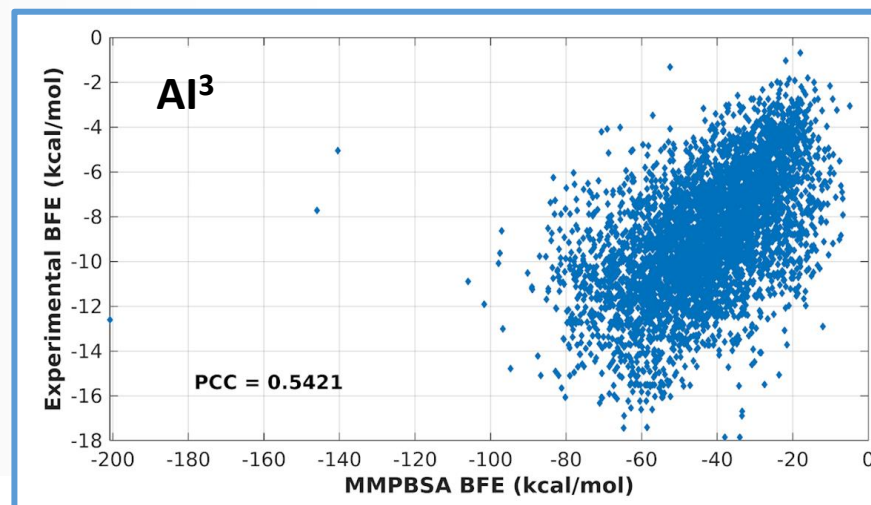- PLAS-20k demonstrates the potential of deep learning in binding affinity prediction.
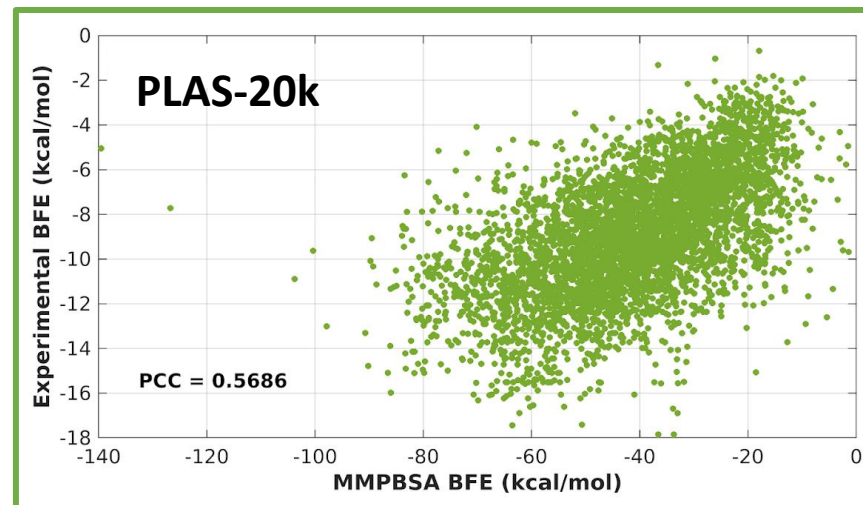


https://doi.org/10.26434/chemrxiv-2023-mg07d

# Stage 1: Results (PLAS-20k versus AI$^3$)



Chart axes: PLAS 20k OpenMM BFE (kcal/mol) versus AWS GMX BFE (kcal/mol)

$f(x) = 0.84\,x - 4.87$
$R^2 = 0.83$

PLAS–20k

AI$^3$

# Stage 1: Comparison with Other Datasets

| Dataset | PCC |
|---------|-----|
| PLAS-20k | 0.5686 |
| AI³ | 0.5421 |
| MISATO | 0.0121 |

# AI³: Stage 1

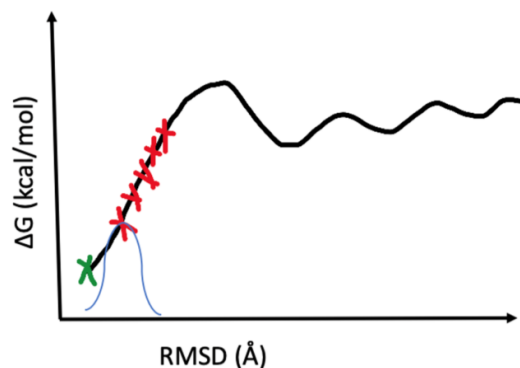| | | |
|---|---|---|
| Consortium to accelerate Protein-Ligand Datasets | Releasing Stage 1 of the World's largest Protein binding | Releasing Code for Augmenting the PLC Dataset on Cloud Infrastructure |
| AI³ | 4.6 TB | Github Repo (coming soon) |

# Future Plans for the AI$^3$ Dataset
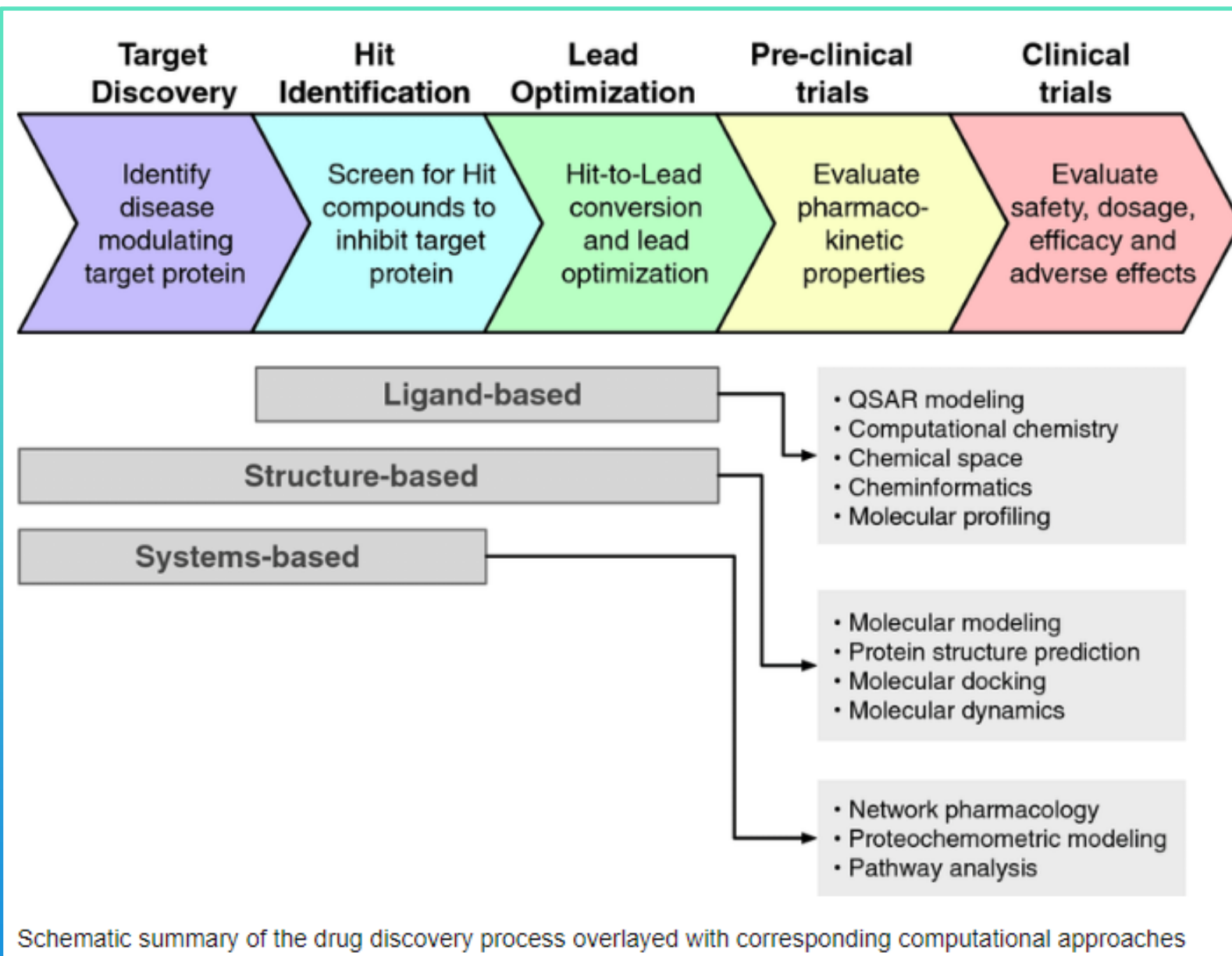


Stage 2b: Higher Energy PLCs

Generate ~220,000 PLC datasets.

Evaluate ~100 kinase inhibitors experimentally in the same lab conditions, and validate computational results.

Join us on the AI$^3$ journey to build the largest PLC Datasets & Accelerate Drug Discovery

Thank you

# Backup

Schematic summary of the drug discovery process overlayed with corresponding computational approaches

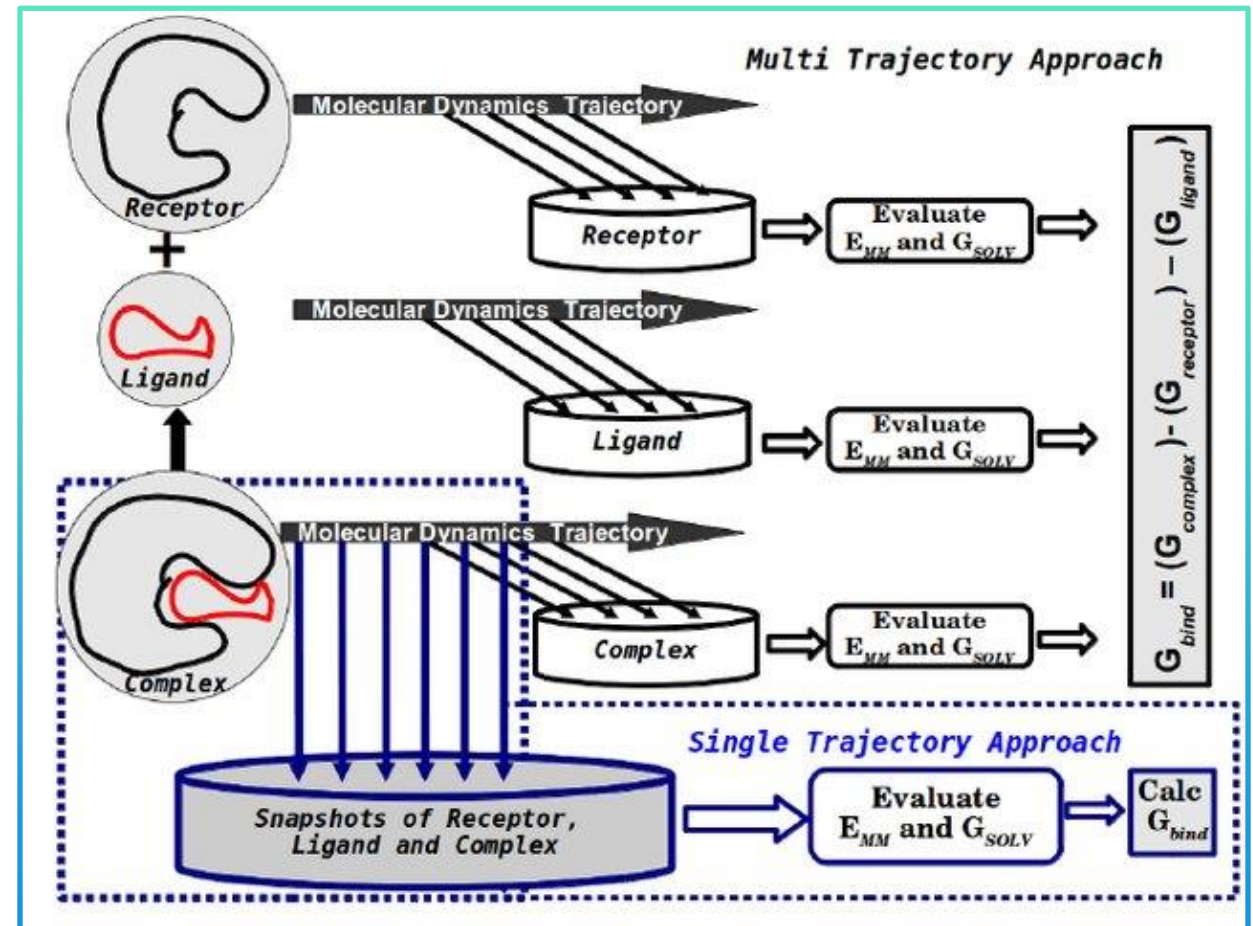Computer aided drug design methodology (1990s) (Nalini, 2020)

# Molecular Dynamics (MD) Simulations

- Molecular dynamics simulations consider protein conformational rearrangements during binding.

- Techniques like MM-PBSA and MM-GBSA calculate binding free energy.

- Post-processing methods, including thermodynamic integration and free-energy perturbation (FEP), contribute to binding free energy determination.

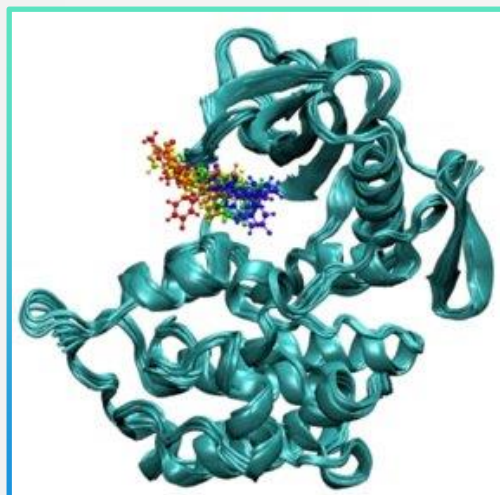$$\Delta G_{MM-PBSA} = \Delta E_{MM} + \Delta G_{Sol}$$

$$\Delta E_{MM} = \Delta E_{ele} + \Delta E_{vdw}$$
$$\Delta G_{Sol} = \Delta G_{pol} + \Delta G_{np}$$



In Silico Engineering of Proteins That Recognize Small Molecules, Mishra, 2012

# Stage 2b: Higher Energy PLCs (under development)



For each PLC, 10 partially bound (higher energy) structures are generated through steered MD simulation

Partially bound or Unbound Protein-Ligand (P-L) Binding Affinity Dataset

▶

Application: Negative examples in machine-learning model training for computational drug discovery.