

2022 Data Science Summer Institute

# Discovery of Small-Molecule Inhibitors of SARS-CoV-2 using Machine Learning

Hyojin Kim



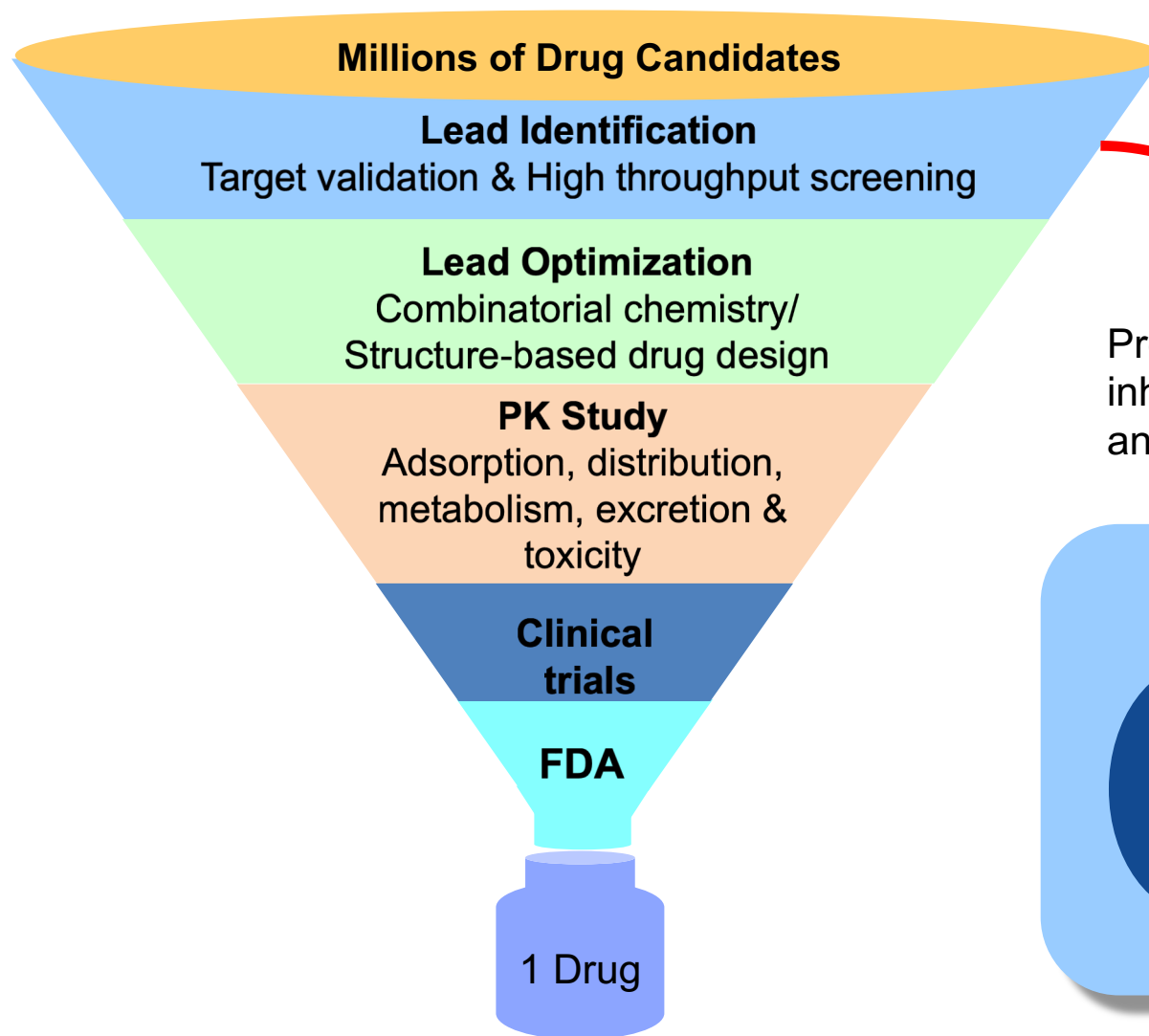
# Agenda

- **Background**
- **Task 1: Screen compounds using SMILES and molecular descriptor**
- **Task 2: Screen compounds using 3D atomic representation**
- **Optional Tasks: Model fusion and performance analysis**
- **Q/A**

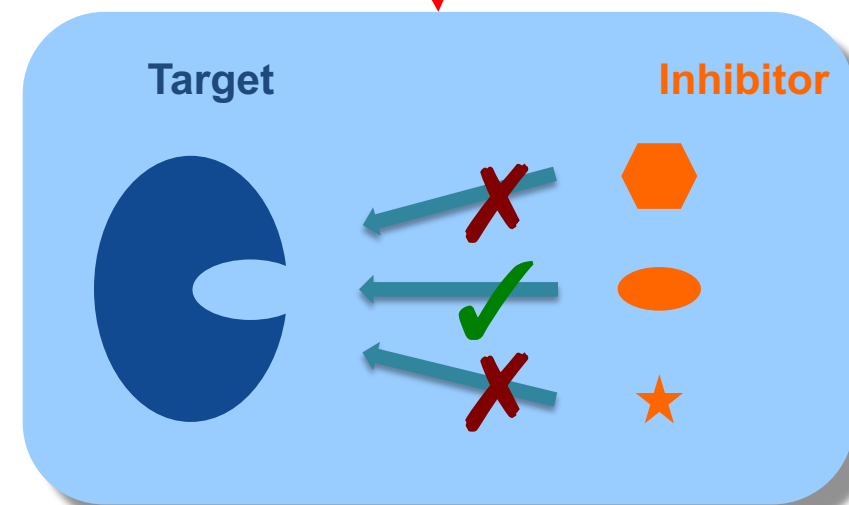
# Background



# Computational chemistry plays an important role in virtual high-throughput screening (vHTS) for drug design

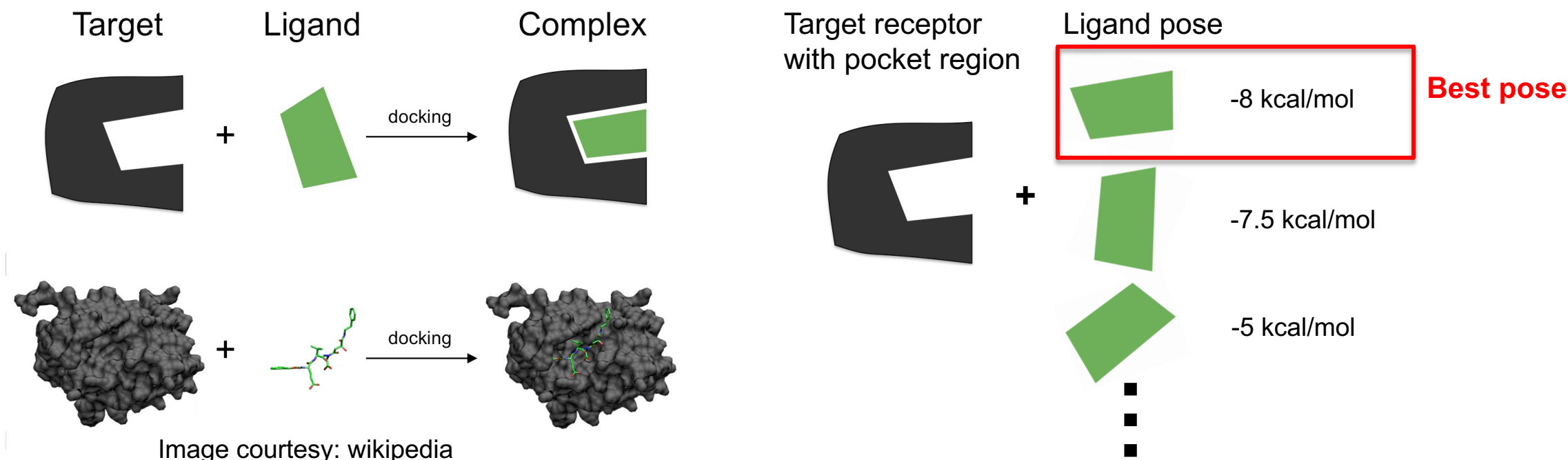


Predicting drug binding to help find inhibitors/activators for therapeutic targets and off-target binding (AutoDock Vina)



# Molecular docking is widely used to predict binding modes between ligand and protein in structural molecular biology

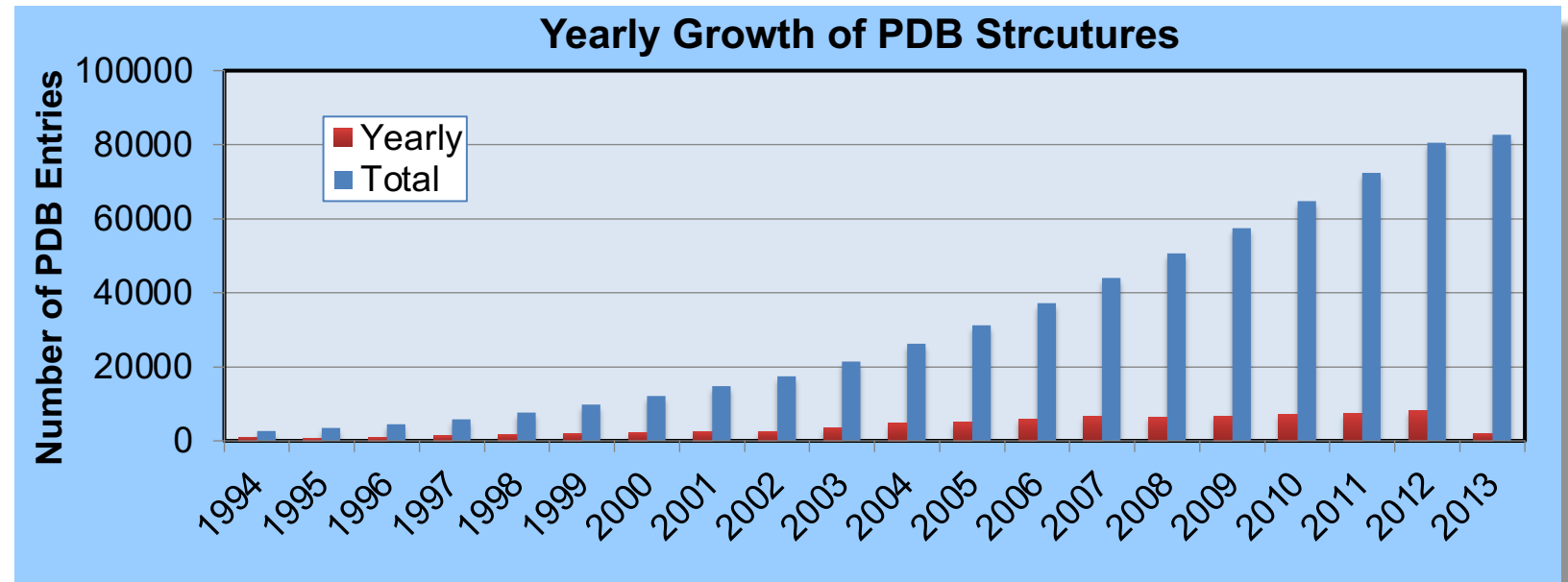
- Molecular docking is a method in computational chemistry to find molecule's pose (orientation) and conformation where the ligand and target receptor are bound to form a stable complex.
- The best pose is selected based on the scoring functions (binding free energy).
- Ligands are screened using their best score (lowest binding free energy).





# We need scalable molecular docking with high performance computing for rapid drug screening

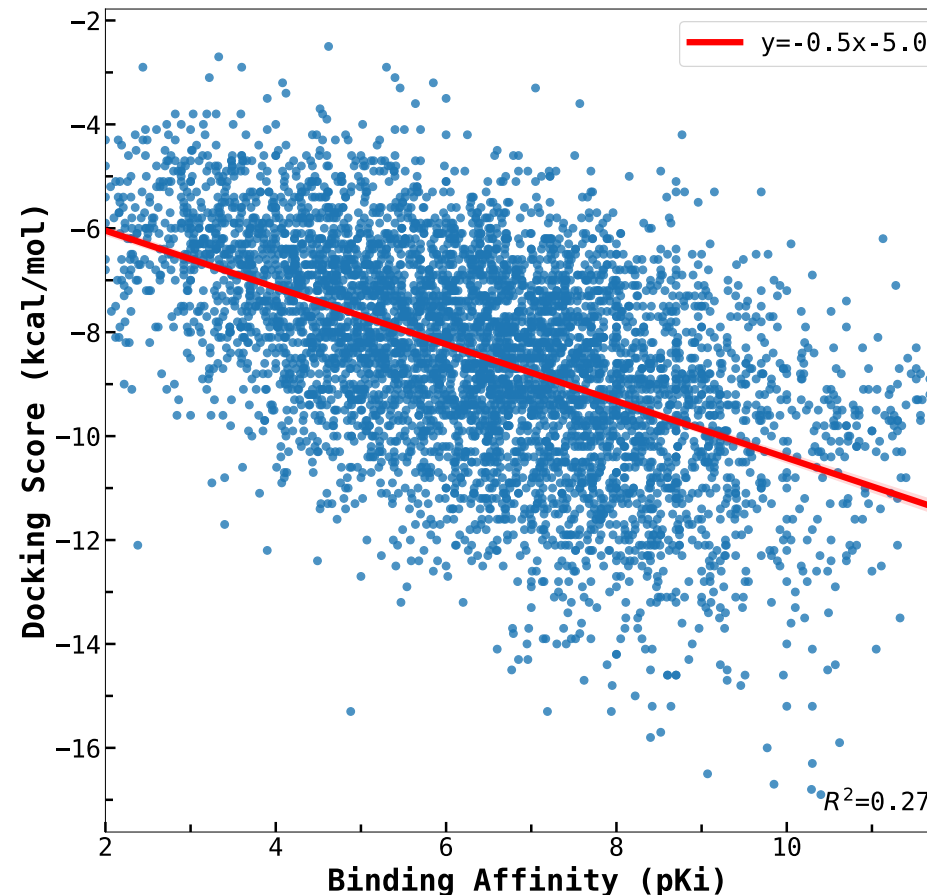
- 1 ligand docking into 1 target takes ~1 min. Rescoring 1 pose takes >10 min.
- Drug-like compounds  $\sim 10^{60}$  possibilities, more chemical compounds available every year
  - e.g., Enamine's REAL database comprises over 4.5 billion molecules



# More accurate screening solution is needed to complement molecular docking

- Low correlation between molecular docking's binding free energy and real binding affinity results in incorrect screening with many false positives.

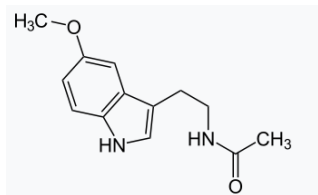
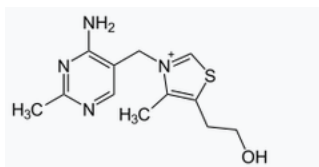
Autodock Vina  
for PDBbind 2020



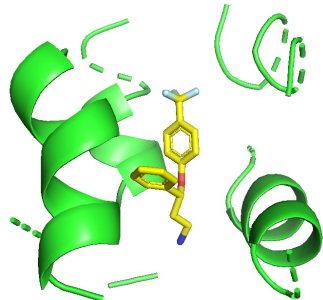
# Recent efforts using machine learning technology show great potential to improve vHTS

- Supervised machine learning approaches have been developed by training ML models on ligand-only or protein-ligand complex data to predict interaction (binding affinity).
- Another direction is use of unsupervised feature learning (VAE, attention, etc)
- ML models can complement molecular docking

Ligand only (for single protein target)



Protein-ligand complex



Machine Learning Models



Binding Affinity



# We have several choices to represent ligand and protein data

- Ligand data are usually derived from the simplified molecular-input line-entry system (SMILES) strings that encode molecular structures and specific instances in 1D string representation, from which we can extract various molecular descriptors or fingerprint
  - File format: .smi
- Ligand data can also be represented as 3D atomic data: list of atoms, each of which has 3D spatial coordinates and atomic features
  - File format: .pdb, .pdbqt, .sdf

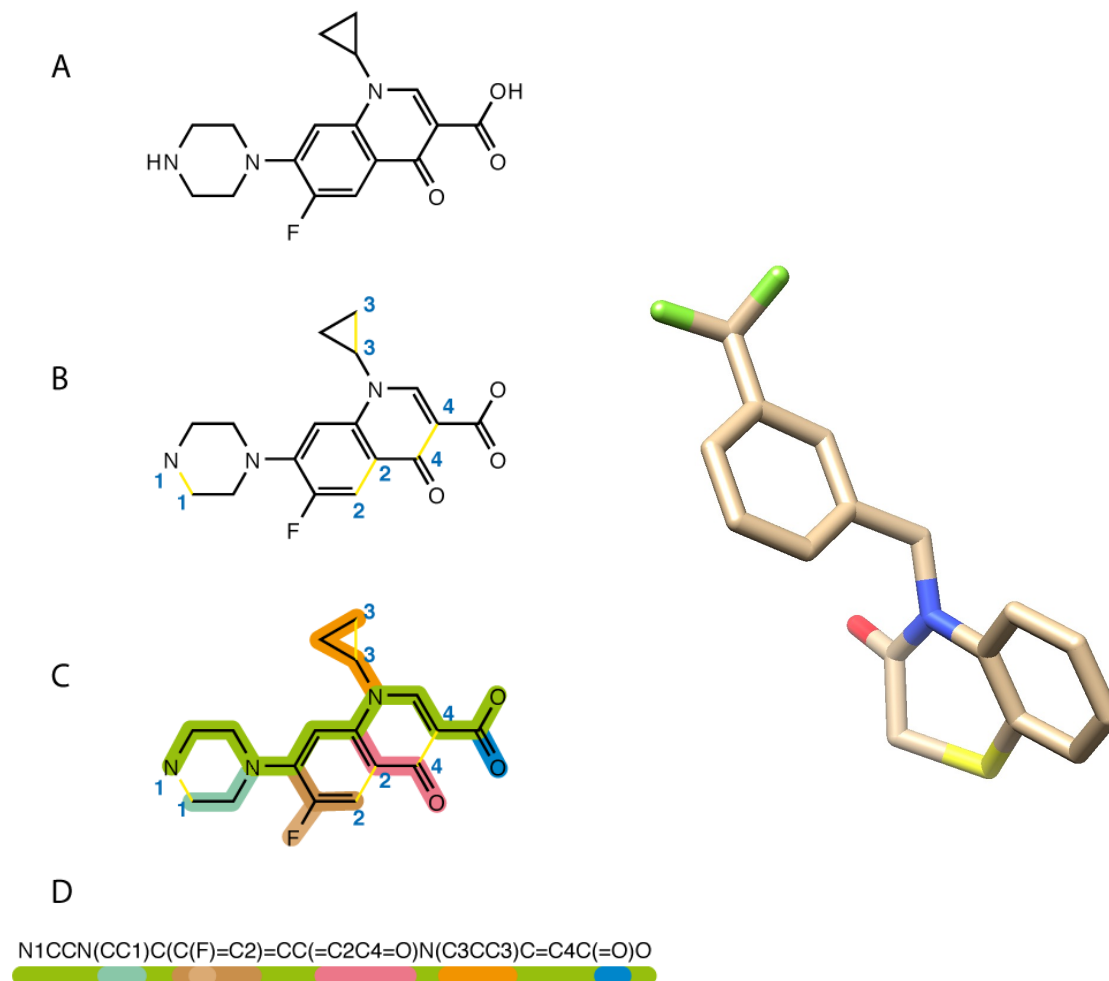


Image courtesy: wikipedia

# We have several choices to represent ligand and protein data

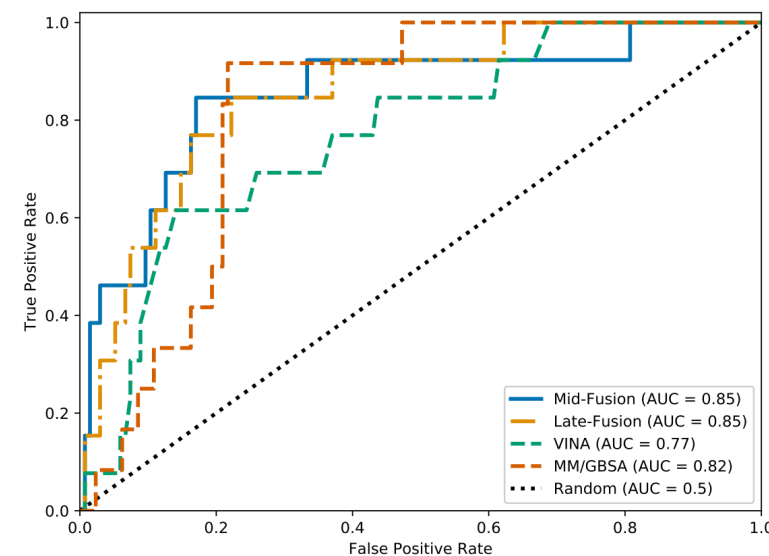
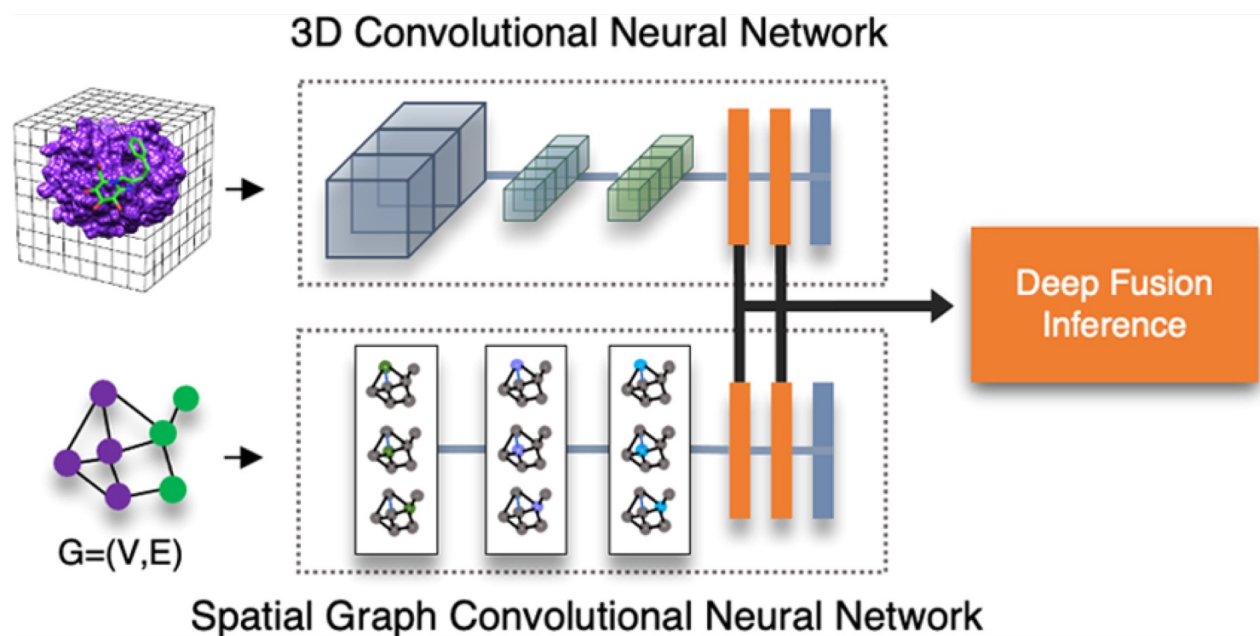
- Protein data can be derived from amino acid sequences and/or 3D atomic representations
  - File format for 3D structures: .pdb, .pdbqt

```
Mpro-N0029_OA_apo.pdb/1-103  1  I T V N V L A W L Y A A V I N G D R W F L N R F T T L N D F N L V A M K Y N Y E P L T Q D H 47
Mpro-N0029_OA_apo.pdb/1-103 48  V D I L G P L S A Q T G I A V L D M C A S L K E L L Q N G M N G R T I L G S A L L E D E F T P 94
Mpro-N0029_OA_apo.pdb/1-103 95  F D V V R Q C S G 103
```



# Our proposed ML approach using 3D structures of protein-ligand data outperforms ligand-only models

- We proposed to use two different ML architectures to complement each other (Fusion): 3D-CNN to capture 3D spatial information (shape), SG-CNN to capture relationships between ligand and protein atoms.
- Our ML approaches trained on protein-ligand data have proven to be fast, effective, compared to existing docking and molecular dynamics (MD) methods



Jones, Kim, et al., 2021

# Terminology

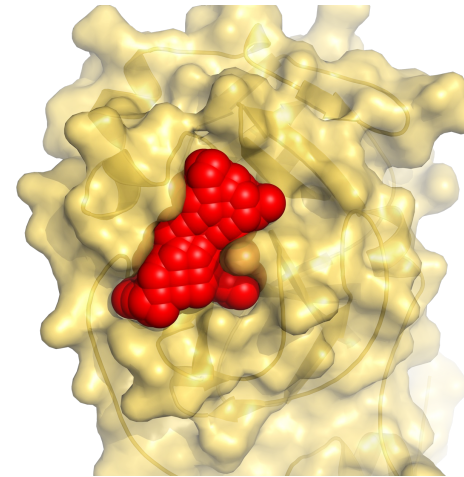
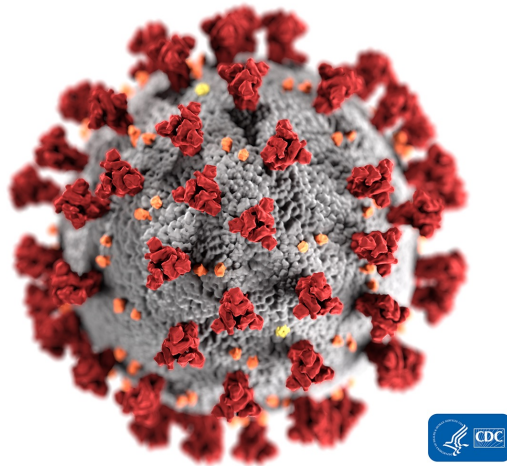
- Ligand: small molecule that forms a complex with a biomolecule (protein receptor)
- Binding affinity: strength of binding interaction between a ligand and a protein target
- Complex: protein-ligand structure where protein and ligand interact with each other
- Crystal structure: 3D arrangement of molecules or atoms throughout a crystal (considered a ground truth pose)
- Docking pose: complex data (predicted binding-conformation of ligand to target binding site) by molecular docking modeling methods (e.g., AutoDock VINA) used in structure-based drug development
- Protein Data Bank (PDB): a database for 3D structural data of biological molecules such as proteins and nucleic acids, each molecule data has its own PDB id
- SMILES: Simplified Molecular Input Line Entry System, 1D representation of a ligand structure

# Main Tasks



# DSSI 2022 is to develop machine learning methods to screen compounds targeting COVID-19

- The goal is to develop machine learning approaches to find potential inhibitors targeting main protease receptors of SARS-CoV-2 (“mpro”).
- We will use Postera Mpro dataset collecting experimentally measured binding affinity data, with LLNL’s curated compounds with binding affinity data.
- We will perform two challenge tasks to predict binding information using two different data representations.



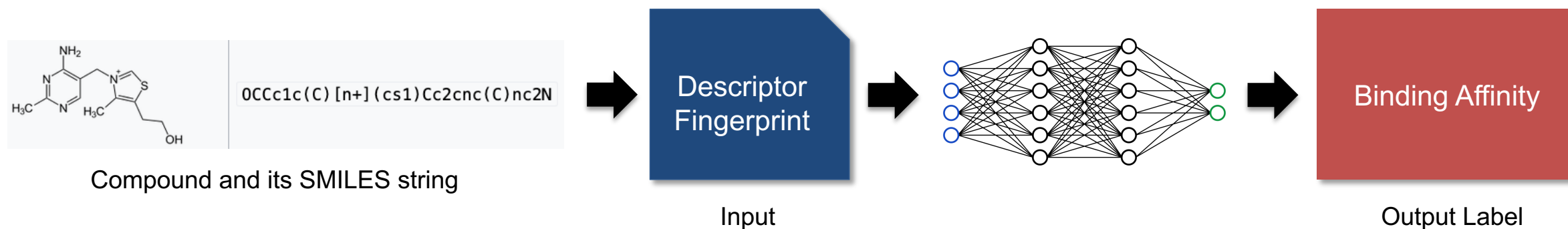
Main protease  
PDB ID: 6Y84  
Diamond, UK

Find inhibitors to stop viral replication



# Task 1 is to use molecular descriptors to predict binding affinity

- Molecular descriptors derived from SMILES string representation of ligand data will be used to predict binding affinity.
- We formulate binding affinity as a classification problem: bind vs. no-bind
- The goal is to develop ML models using random forest, neural network, etc.



# Task 1 dataset is based on rdkit feature descriptors

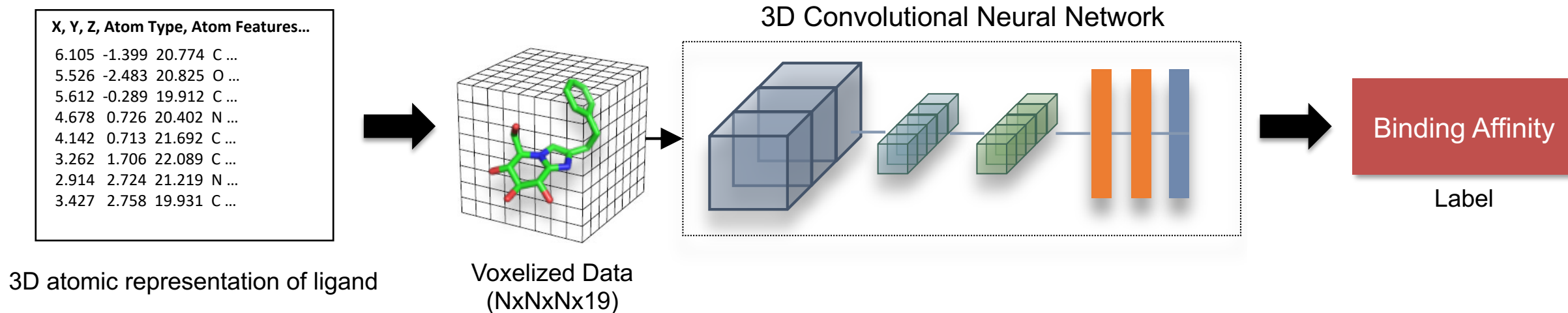
- The input data is ligand feature descriptors using rdkit
  - Rdkit is open-source cheminformatics python package
  - We used `rdkit.ML.Descriptors.MoleculeDescriptors.MolecularDescriptorCalculator`
- The descriptor contains various compound and chemical attributes such as
  - Molecular weights, average molecular weight of the molecule ignoring hydrogens, number of radical electrons, number of valence electrons, etc
  - 208 feature vectors
- The input file is **mpro\_exp\_data2\_rdkit\_feat.csv**
  - `cmpd_id`: compound id
  - `smiles`: SMILES string
  - `label`: 0 (no-bind) or 1 (bind)
  - `subset`: train/validation/test split (1955 training, 113 validation, 240 test samples)
  - `feat_1` to `feat_208`: rdkit descriptor features
  - Input csv contains **NaN**
- Instead of the provided descriptors, SMILES or molecular descriptors or fingerprints can be used.
  - [https://rdkit.org/UGM/2012/Landrum\\_RDKit\\_UGM.Fingerprints.Final.pptx.pdf](https://rdkit.org/UGM/2012/Landrum_RDKit_UGM.Fingerprints.Final.pptx.pdf)

# In Task 1, you will develop a ML method to read input molecular descriptors to classify bind or no-bind of the compound

- Any ML method can be used to predict interaction between ligands and the main protease.
  - Random forest
  - Neural networks with fully connected layers
  - RNNs (LSTM, BERT)
- Model performance can be measured using classification metrics: accuracy, precision, recall and F1 score
- Receiver operating characteristic (ROC) curve and Area Under the Curve (AUC) can also be reported

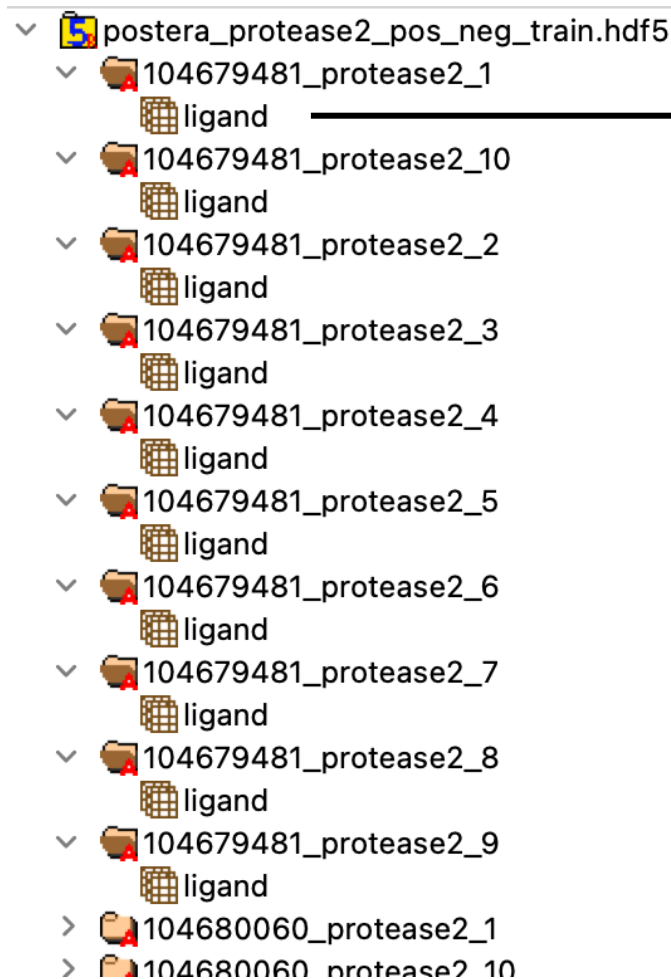
# Task 2 is to develop ML models to predict binding affinity using 3D structures of ligand data

- Now let's use 3D structure information of ligand data
- The goal is to develop a method to predict interaction as a classification problem using any CNN method to interpret 3D atomic representation



# Task 2 dataset contains 3D atomic representations for ligands

- There are three hierarchical data format (HDF5) files for training, validation and test:
  - postera\_protease2\_pos\_neg\_train.hdf5
  - postera\_protease2\_pos\_neg\_val.hdf5
  - postera\_protease2\_pos\_neg\_test.hdf5
- Each ligand contains 10 different poses (but without the target protein receptor) with a label (0 or 1), each of which contains 3D information in “ligand”
  - Keys: compound id + protease2 + pose id
- Each 3D ligand data is a multi-dimensional array whose size is 100x21 where 100 is the total number of atoms and 21 is the feature size: 3 for xyz coordinates and 19 features
  - In the case of a ligand with less than 100 atoms, the rest rows will be zero padded.



- Atom coordinates (xyz)
- Atomic numbers (10 one-hot encodings)
- Heavy valence
- Hetero valence
- Partial charge
- Mol code
- Hydrophobic
- Aromatic
- Acceptor
- Donor
- Ring

HDF viewer: <https://www.hdfgroup.org/downloads/hdfview/>

# You will use a convolutional neural network architecture to read 3D ligand data

- We recommend 3D-CNNs with 3D convolutional filters/kernels, OR graph-CNNs. But any CNN method can be used to interpret 3D ligand data.
- For 3D-CNNs, the ligand array data needs to be represented into a 3D voxel grid
- For Graph-CNNs, the ligand data needs to be represented as a graph
- Choice of voxel grid dimension depends on memory, resolution, collision, sparsity, filter size, etc
- Model performance report using classification metrics: accuracy, precision, recall and F1 score
- ROC curve and AUC can also be reported

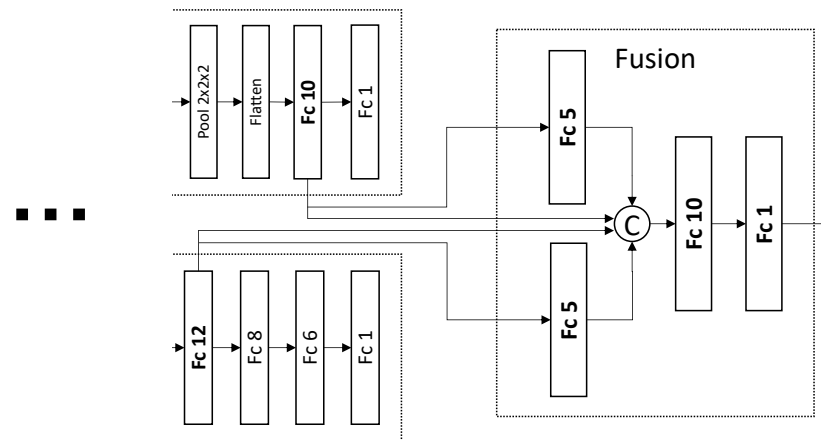


# Optional Tasks



# Optional Task 1: Two ML models using different feature representations could be complementary to improve prediction

- The benefit to fuse different ML models lies in combining their possibly complementary feature representations.
  - e.g., video activity recognition by bridging feature difference between multiple input types (visual and temporal data)
- Our "Fusion" framework shows an improved predictive power on the protein-ligand binding affinity.
- The simplest way to fuse two models is to average the final layer's activations.
- Another way would be to fuse intermediate layer's activations in a separate neural network.



Jones, Kim et al., 2021

# Optional Task 2: Comprehensive performance and behavior analysis is crucial for model interpretability

- Comprehensive model analysis includes 1) input feature importance, 2) spatial region relevance, 3) performance by compound family/similarity
  - Cluster compounds by SMILES, fingerprint, descriptors
  - Analyze performance per compound group



#### **Disclaimer**

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.