

Protein-Ligand Scoring with Convolutional Neural Networks

David Koes

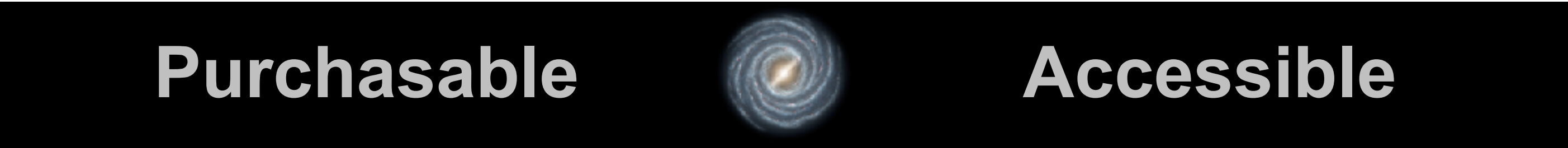


@david_koes

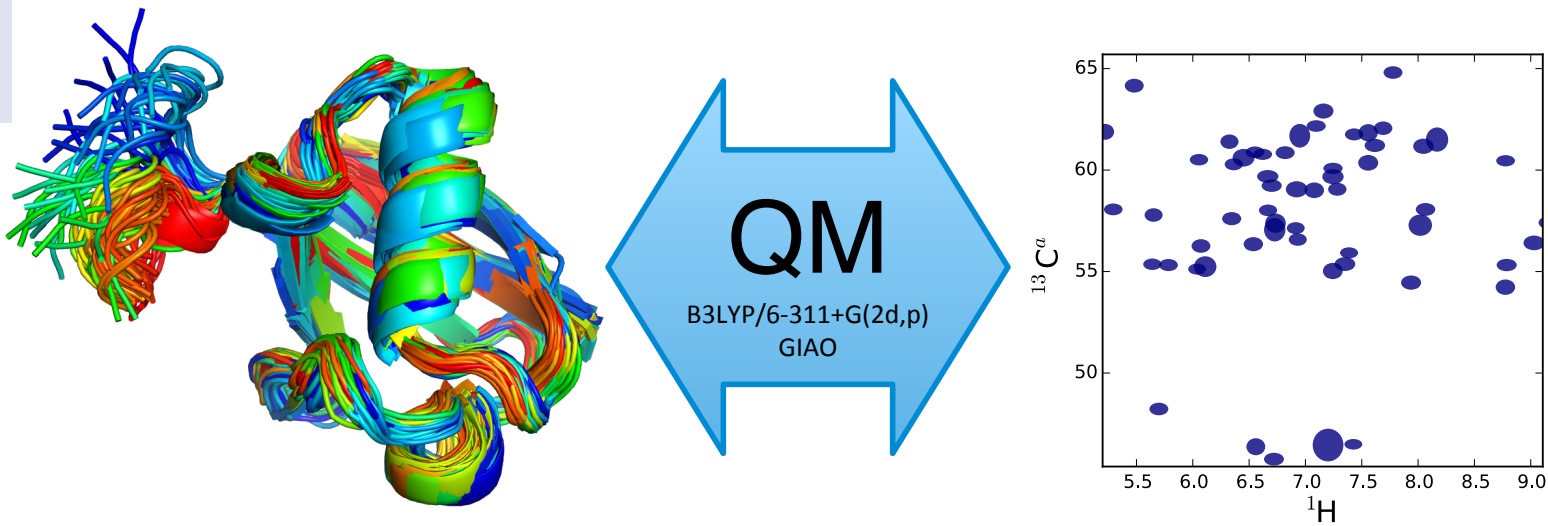
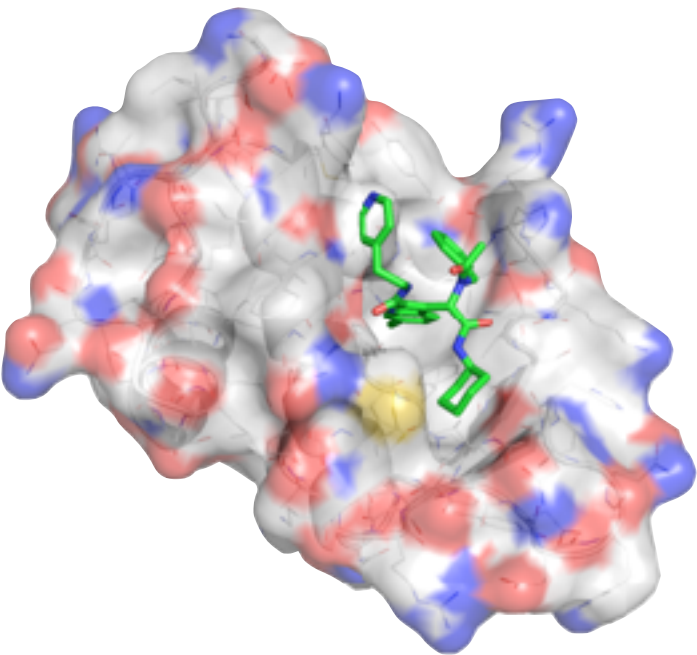
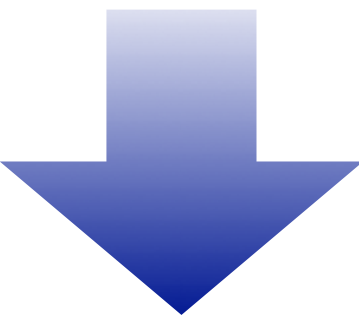


Novartis
June 23, 2017

Removing barriers to computational drug discovery one bit at a time



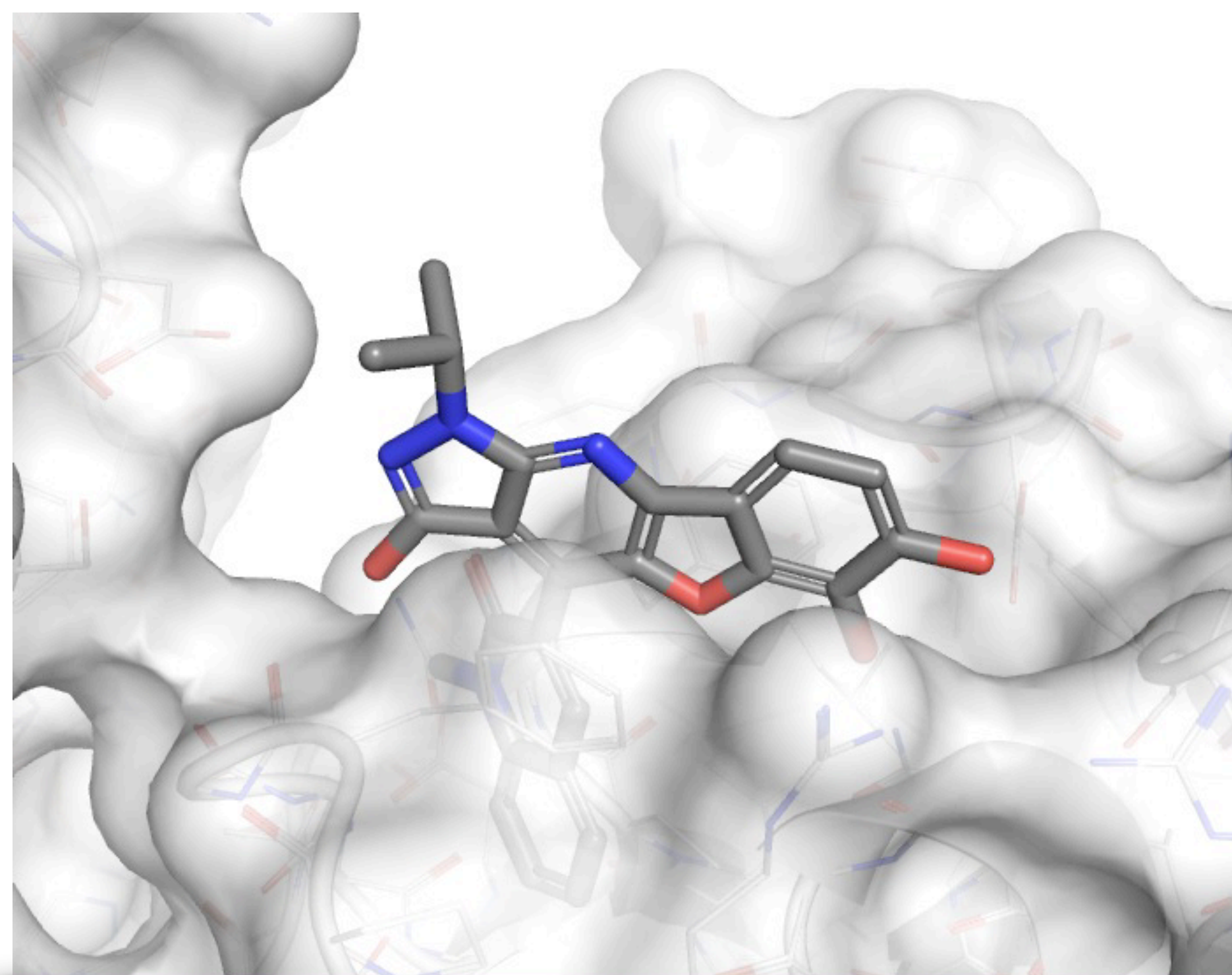
<http://pharmit.csb.pitt.edu>



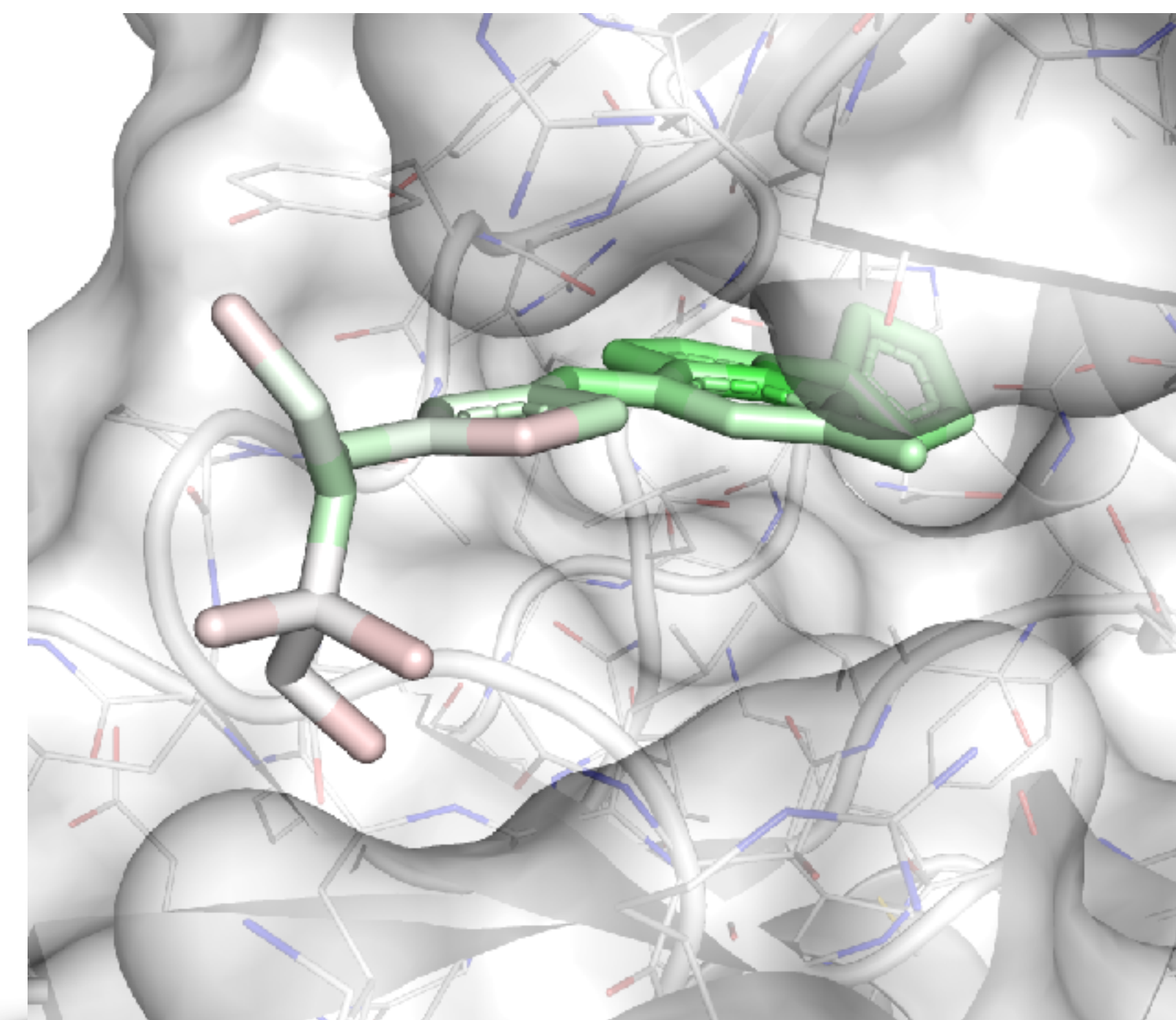
<https://github.com/dkoes/MD2NMR>

Structure Based Drug Design

Virtual Screening



Lead Optimization



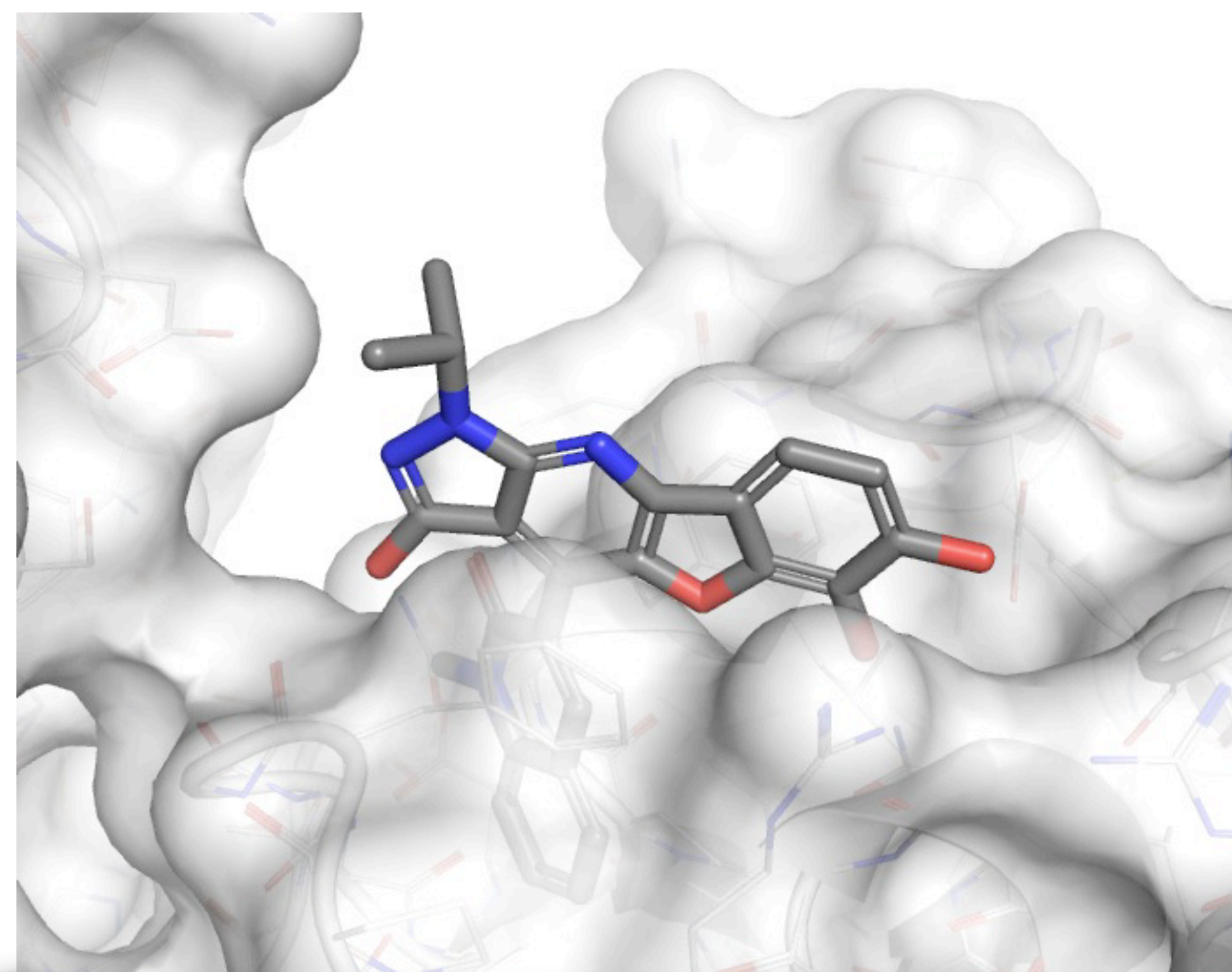
Pose Prediction

Binding Discrimination

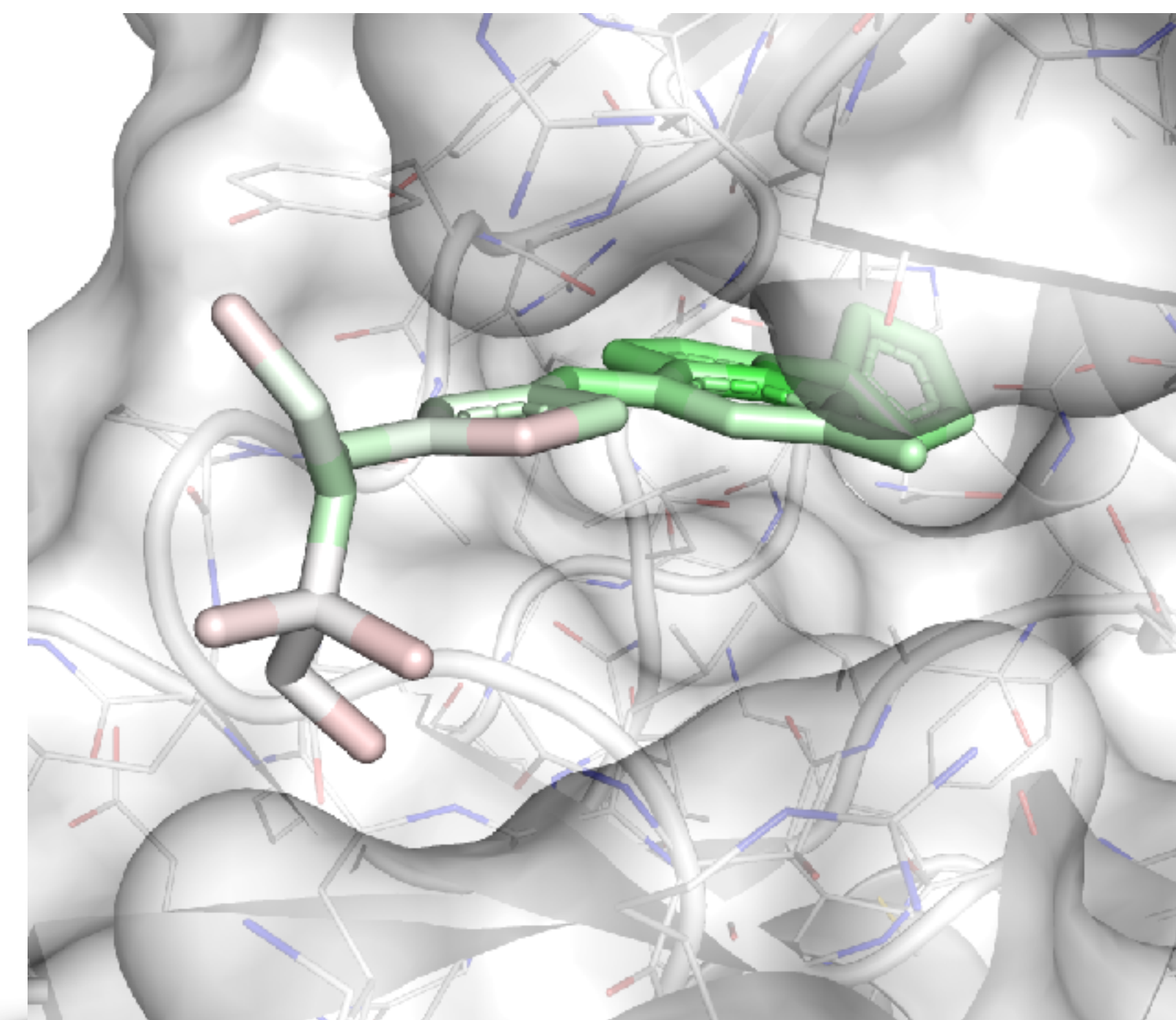
Affinity Prediction

Structure Based Drug Design

Virtual Screening



Lead Optimization

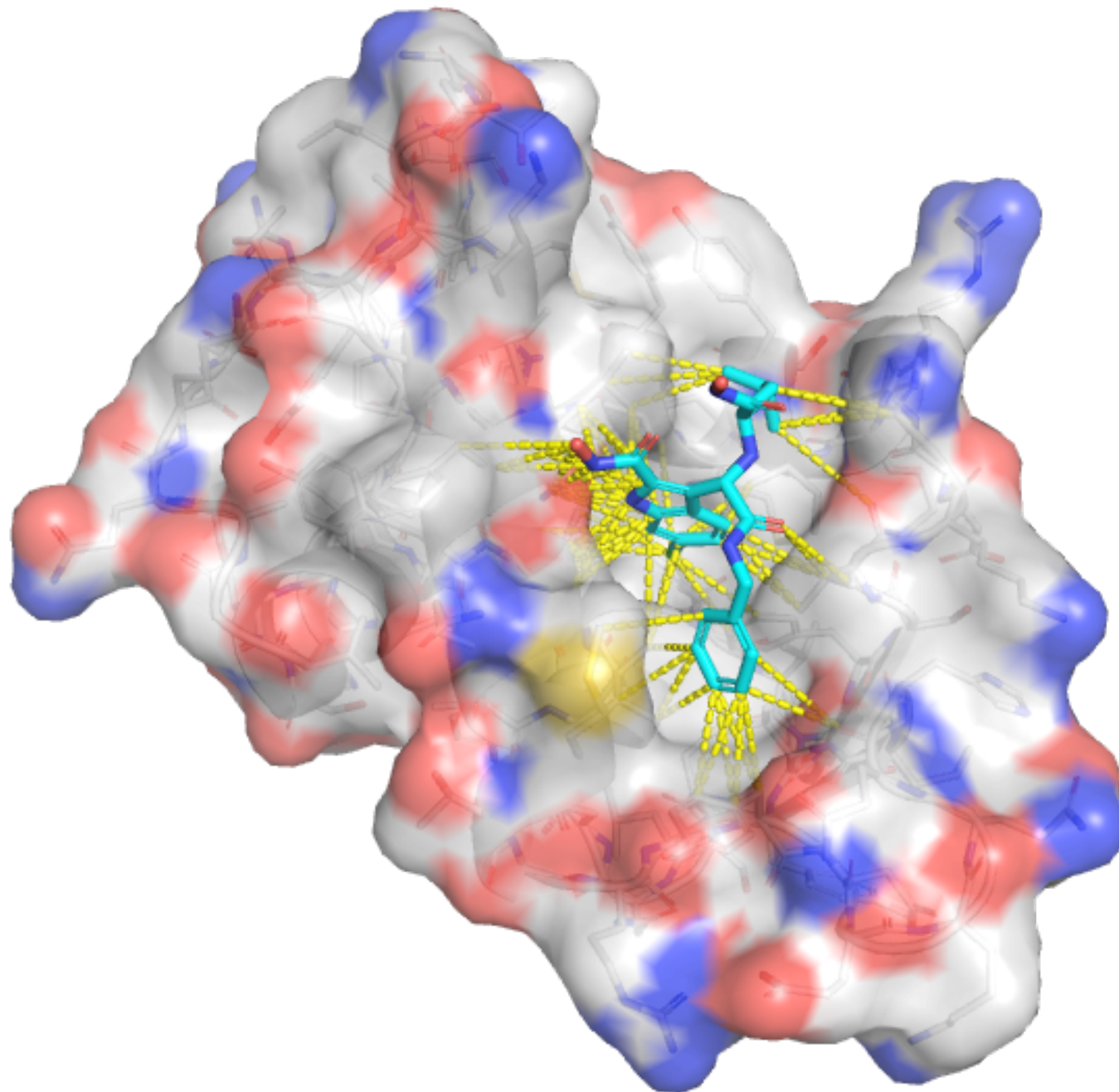


Pose Prediction

Binding Discrimination

Affinity Prediction

Protein-Ligand Scoring

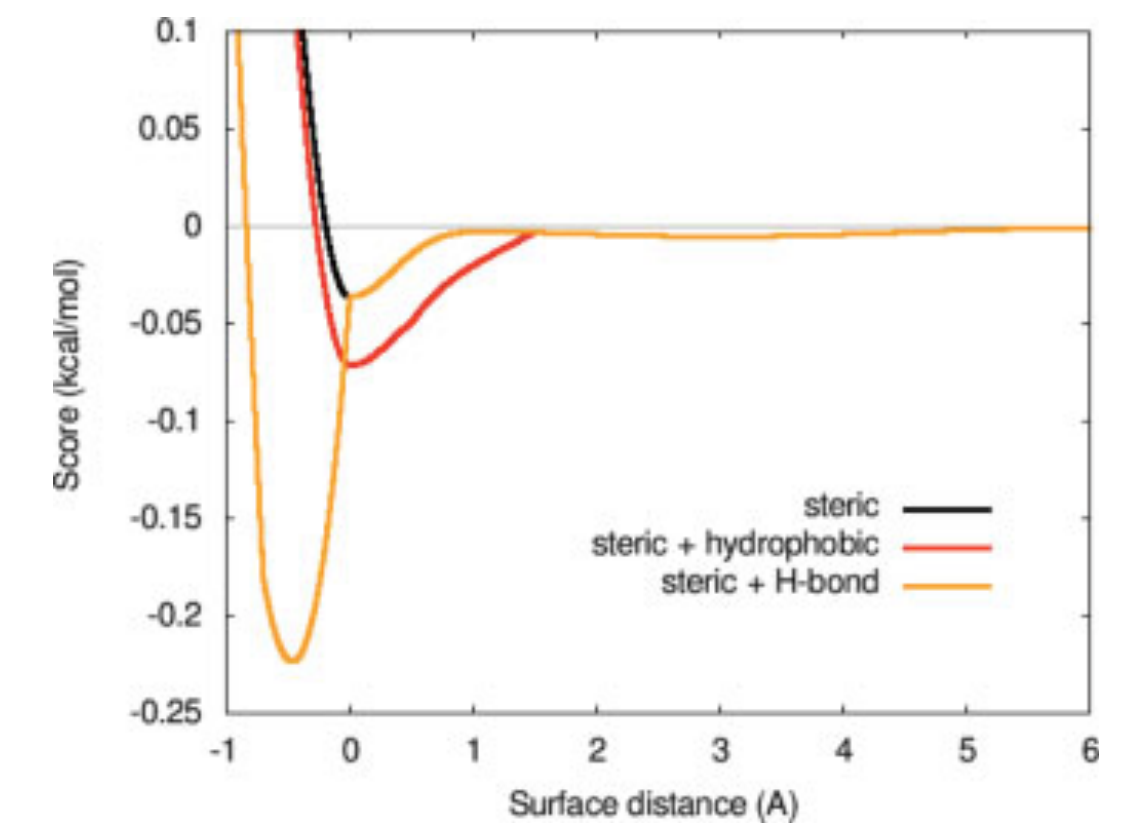
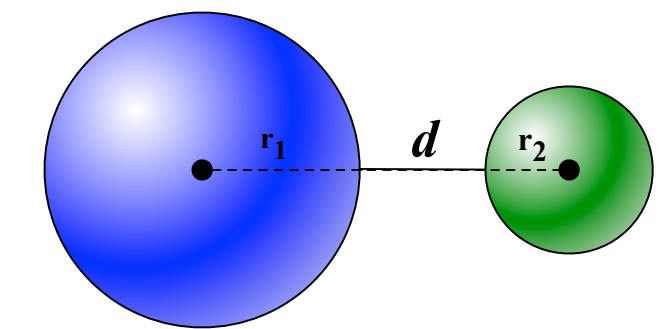


AutoDock Vina

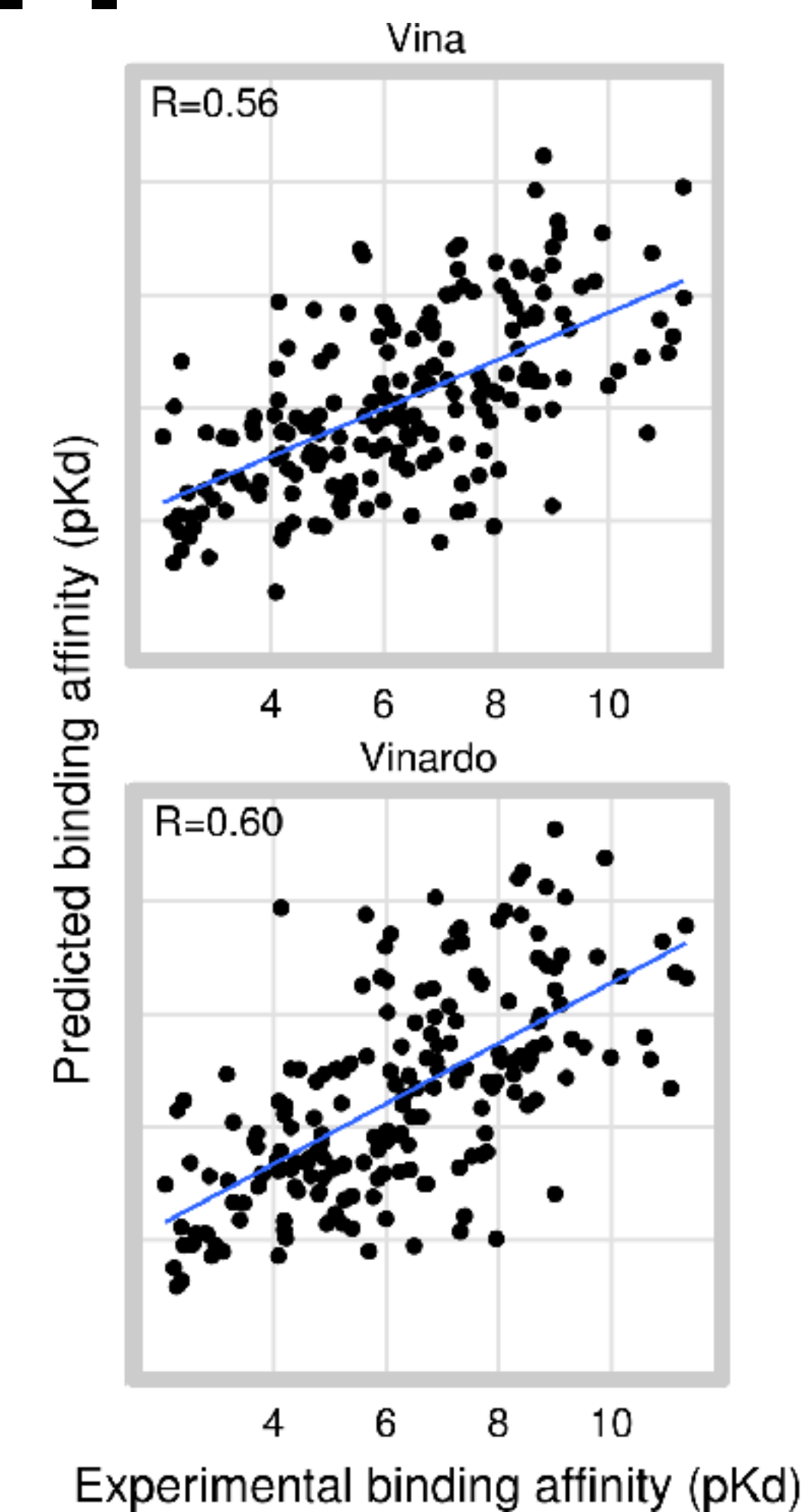
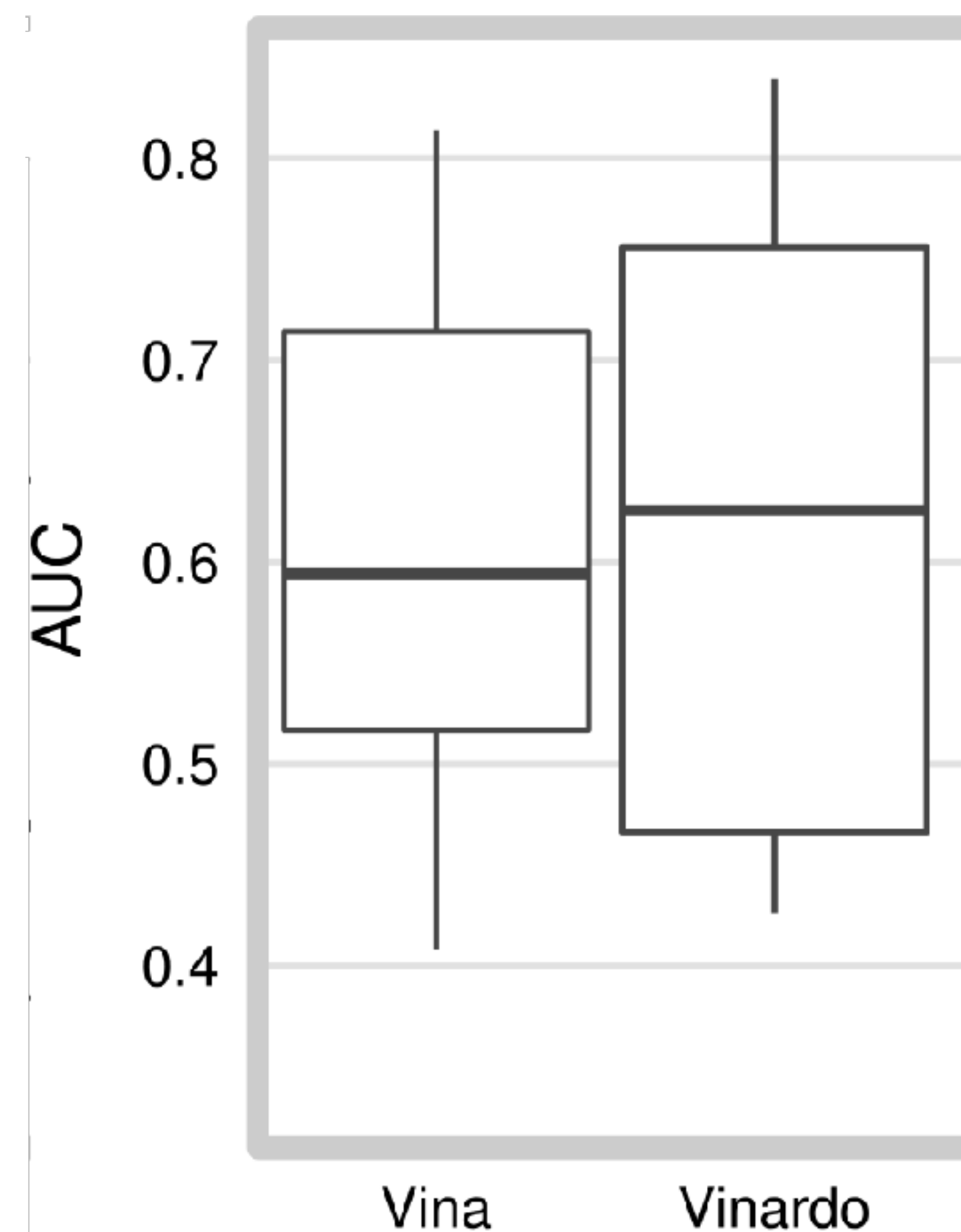
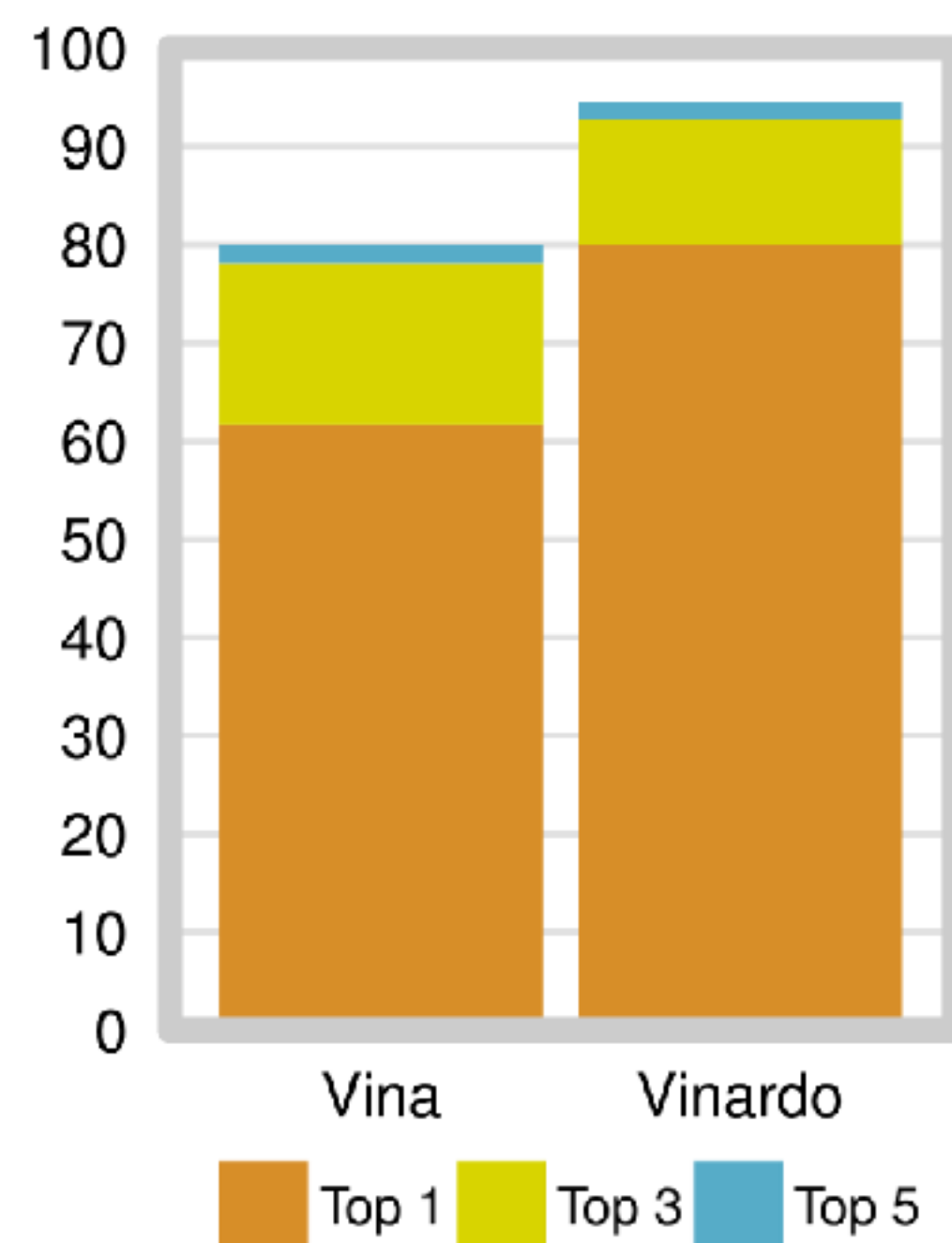
$$\begin{aligned} \text{gauss}_1(d) &= w_{\text{guass}_1} e^{-(d/0.5)^2} \\ \text{gauss}_2(d) &= w_{\text{guass}_2} e^{-((d-3)/2)^2} \\ \text{repulsion}(d) &= \begin{cases} w_{\text{repulsion}} d^2 & d < 0 \\ 0 & d \geq 0 \end{cases} \end{aligned}$$

$$\text{hydrophobic}(d) = \begin{cases} w_{\text{hydrophobic}} & d < 0.5 \\ 0 & d > 1.5 \\ w_{\text{hydrophobic}}(1.5 - d) & \text{otherwise} \end{cases}$$

$$\text{hbond}(d) = \begin{cases} w_{\text{hbond}} & d < -0.7 \\ 0 & d > 0 \\ w_{\text{hbond}}(-\frac{10}{7}d) & \text{otherwise} \end{cases}$$



State of the Art



Pose Prediction

Binding Discrimination

Affinity Prediction

Can we do better?

Accurate pose prediction, binding discrimination, **and** affinity prediction without sacrificing performance?



Can we do better?

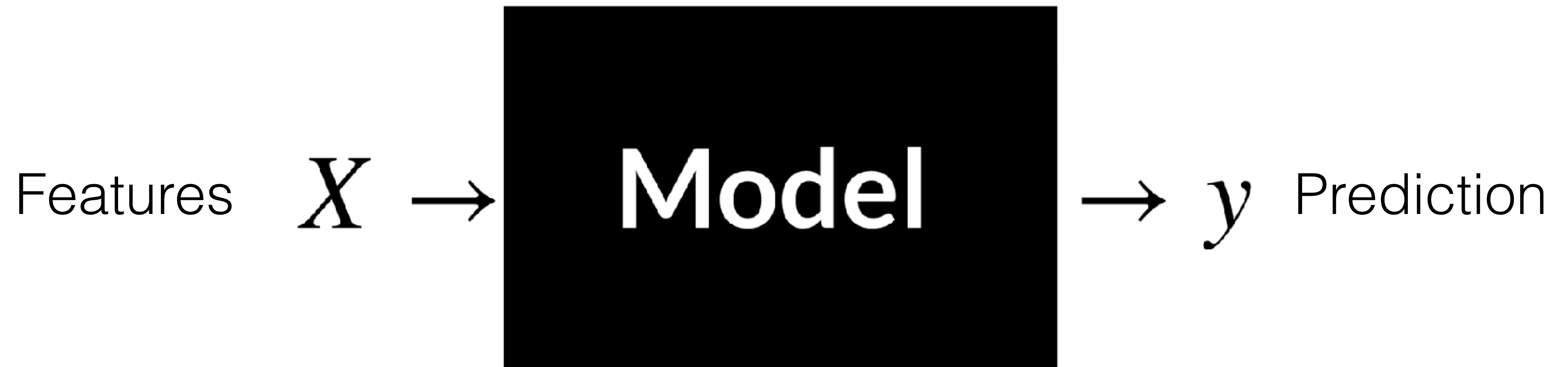
Accurate pose prediction, binding discrimination, **and** affinity prediction without sacrificing performance?

Key Idea: Leverage “big data”

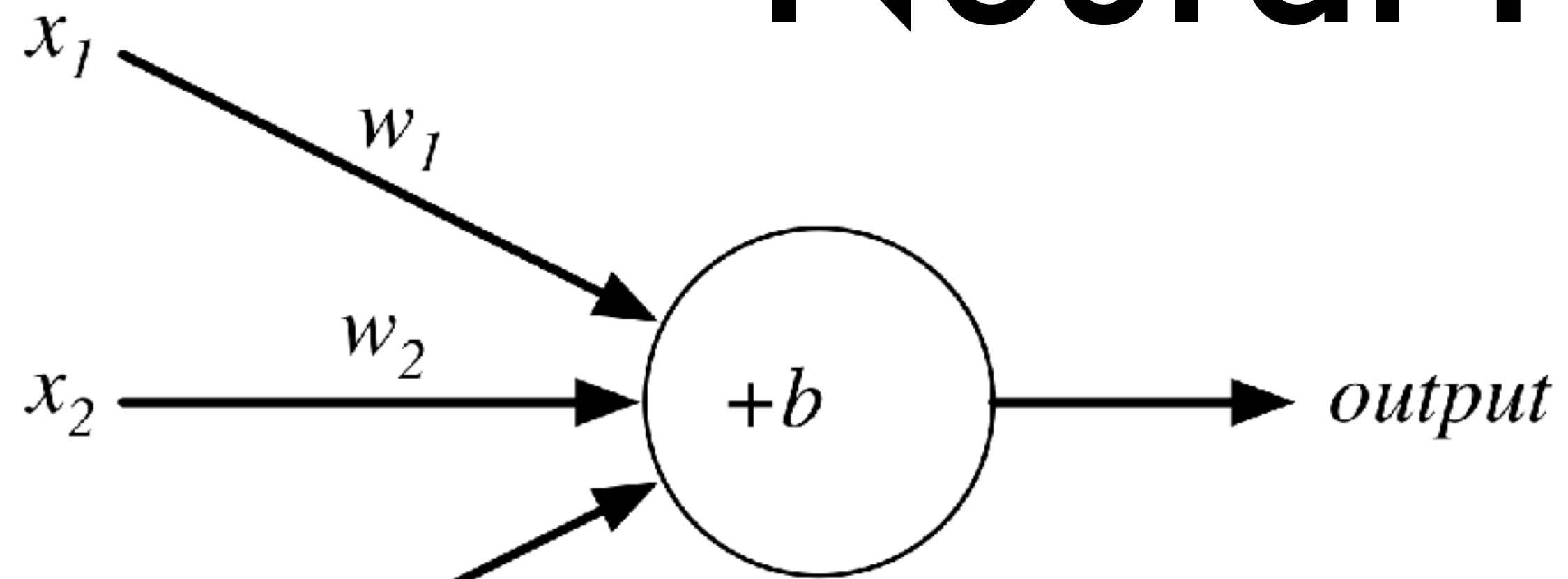
- 231,655,275 bioactivities in PubChem
- 125,526 structures in the PDB
- 16,179 annotated complexes in PDBbind



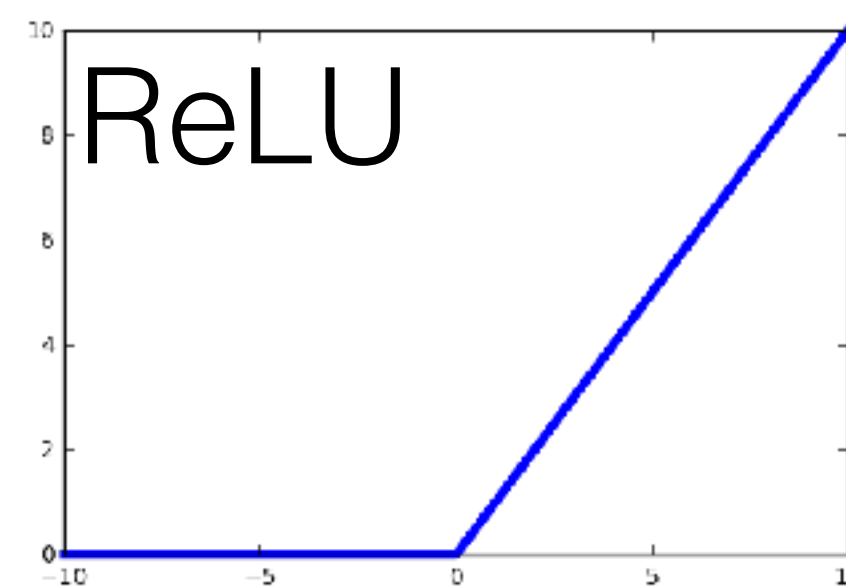
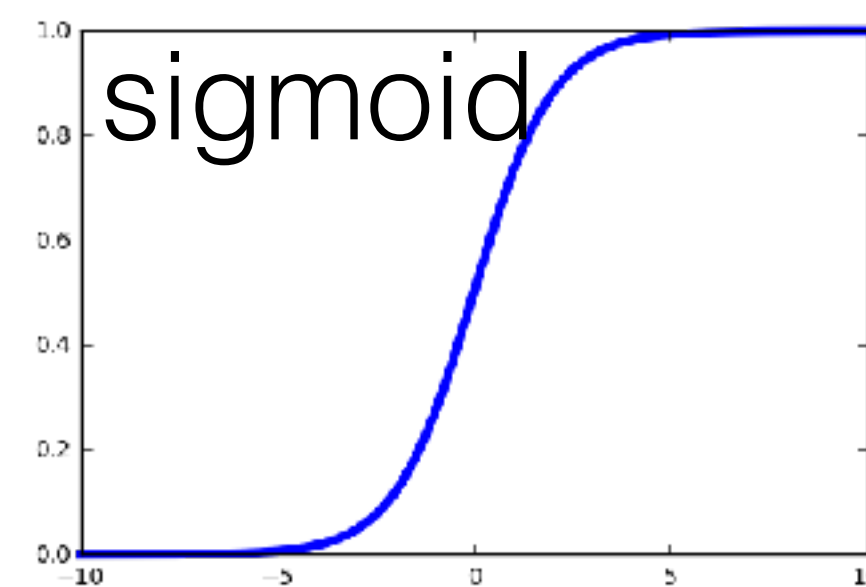
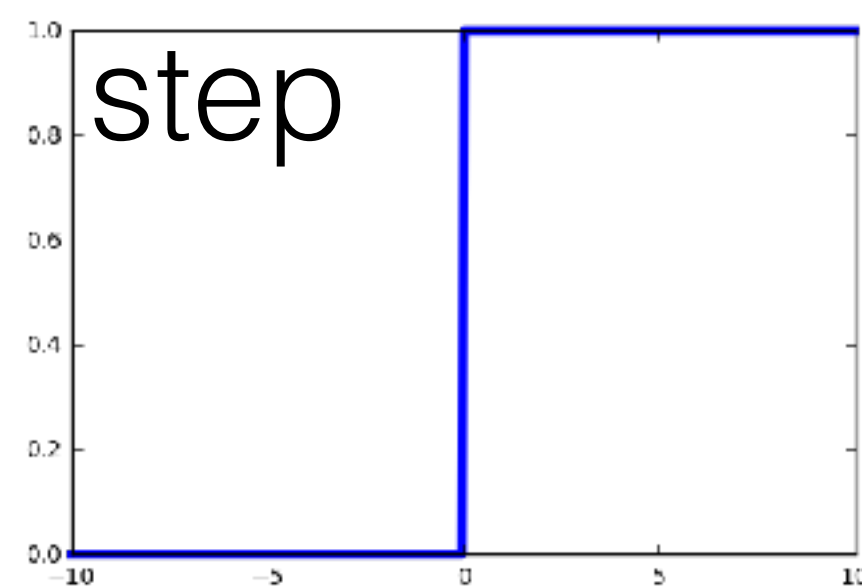
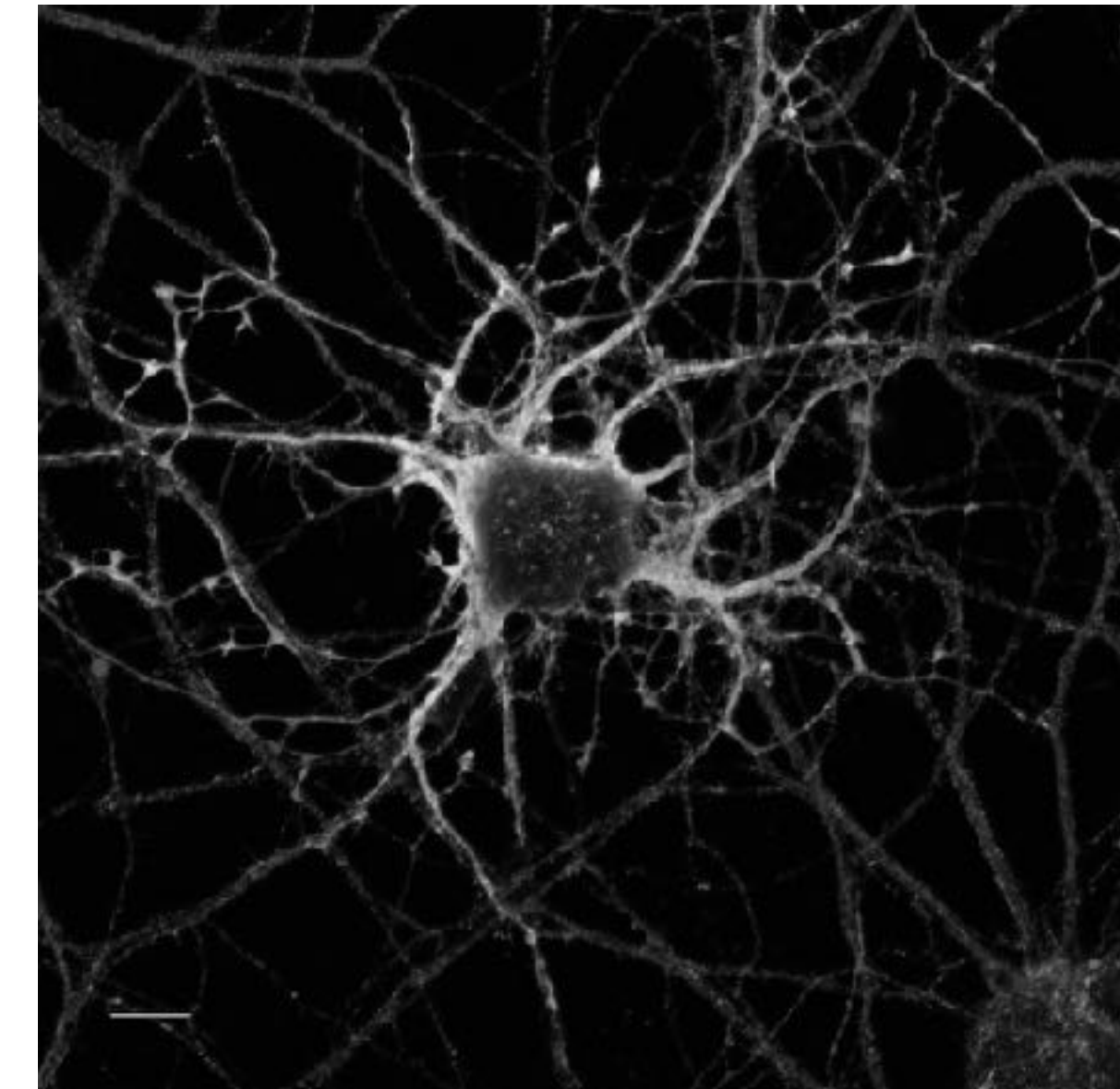
Machine Learning



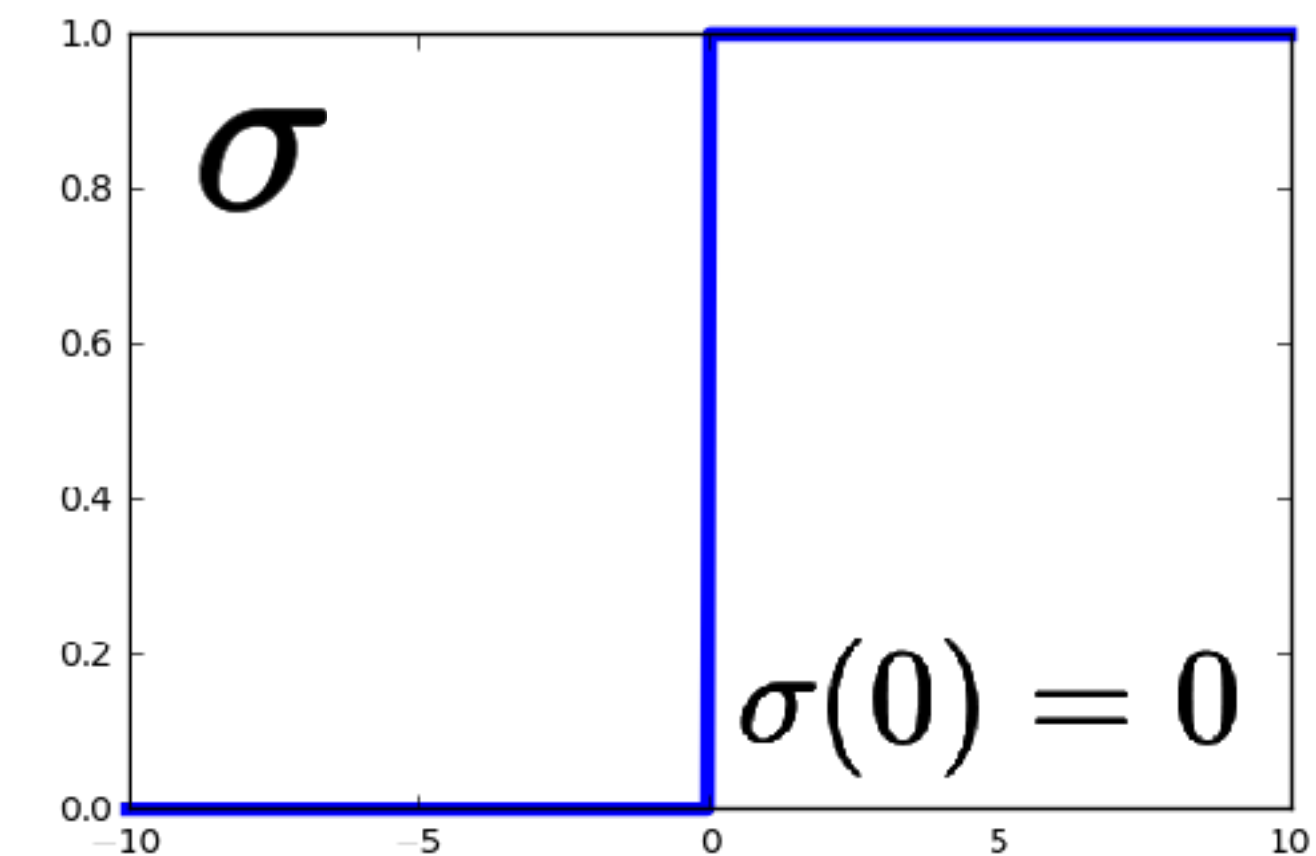
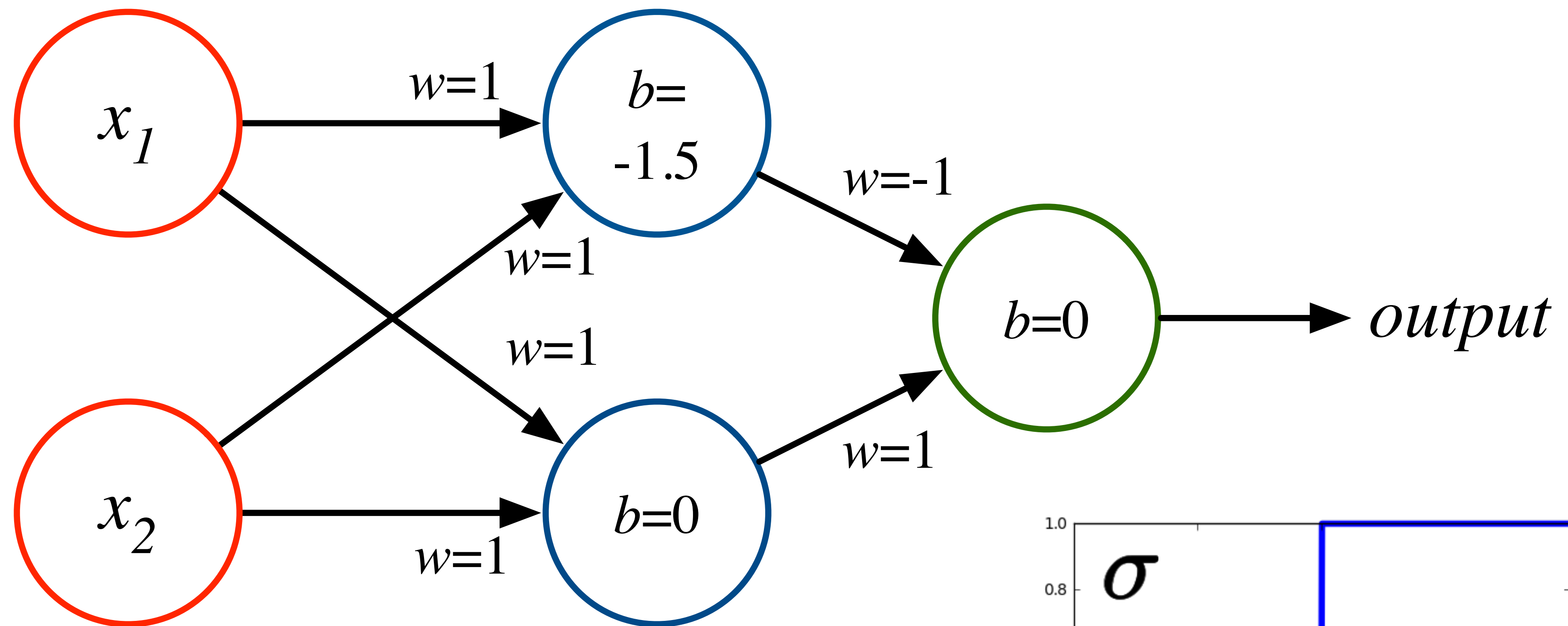
Neural Networks



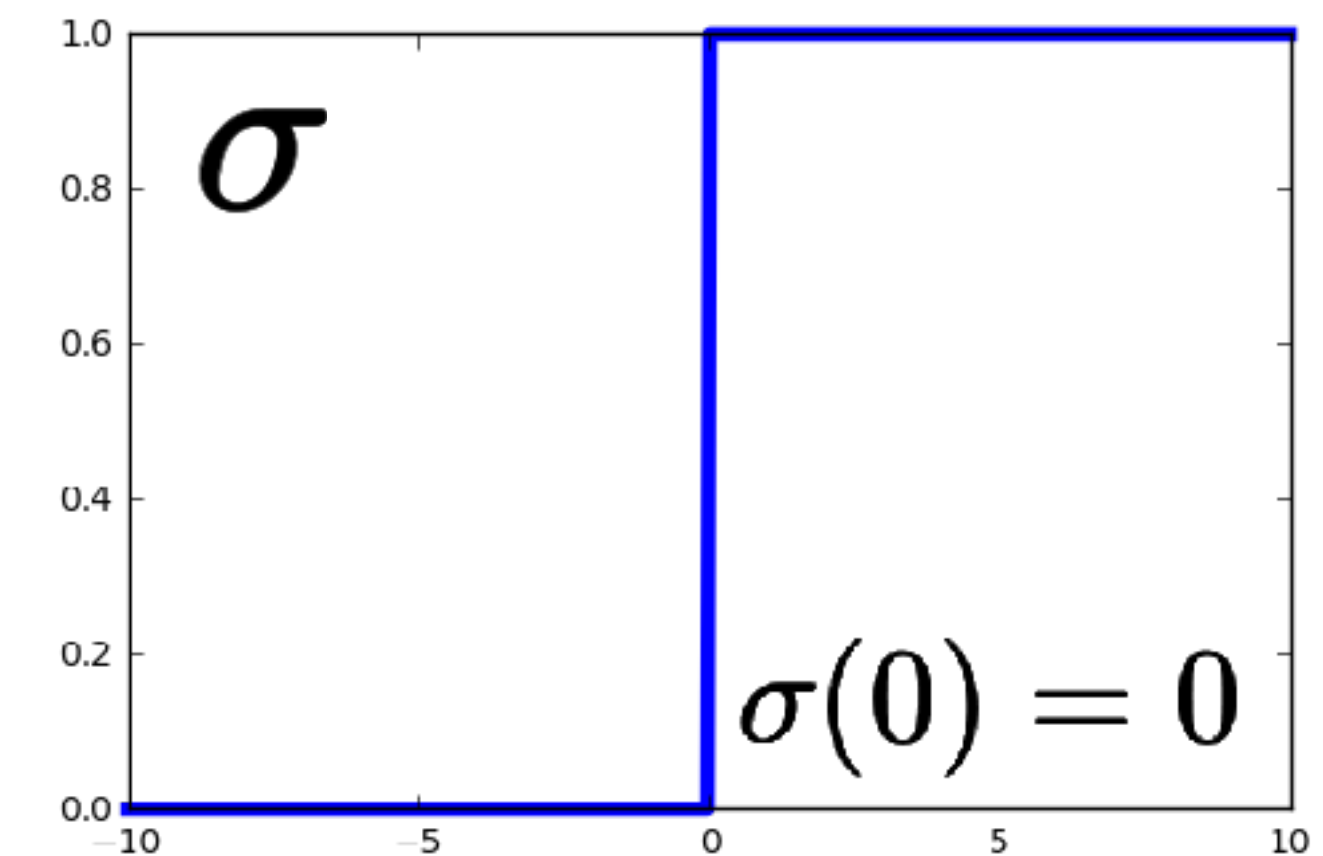
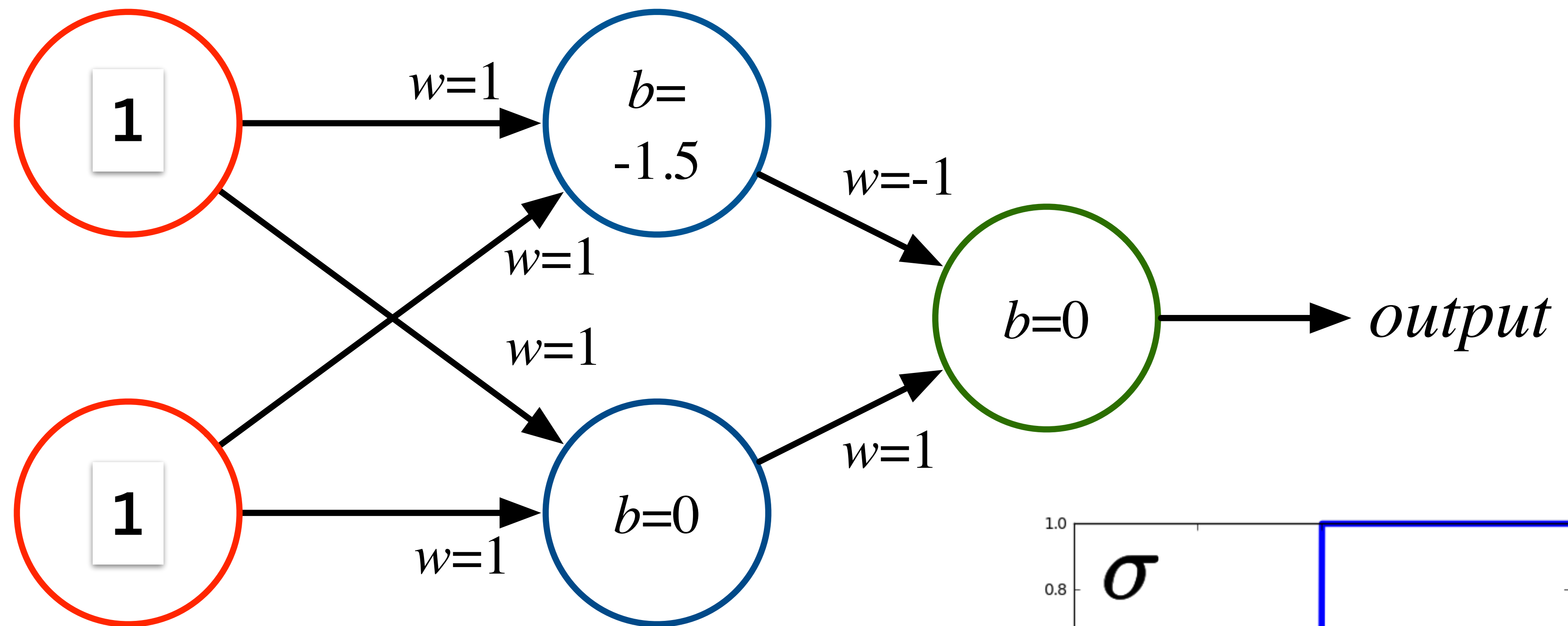
$$\text{output} = \sigma \left(\sum_i w_i x_i + b \right)$$



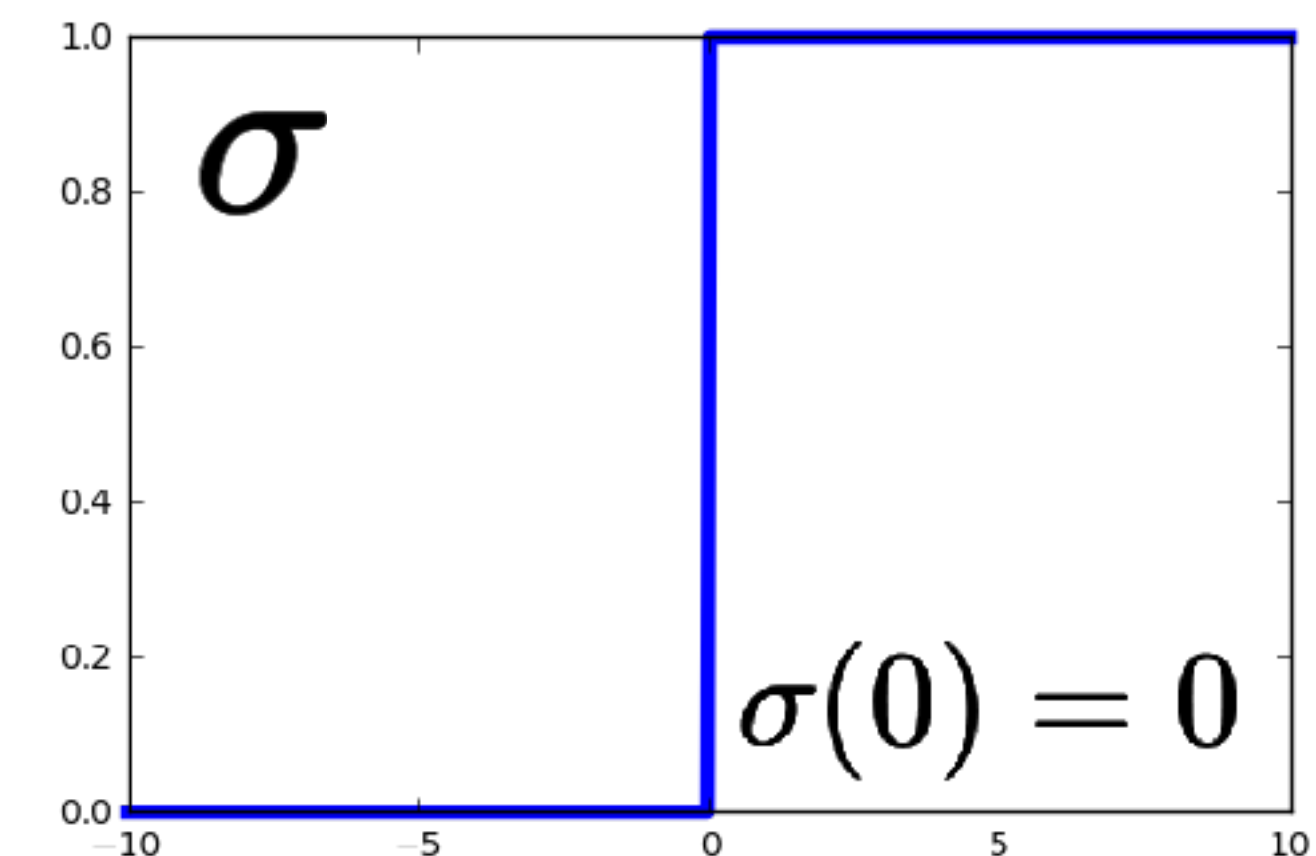
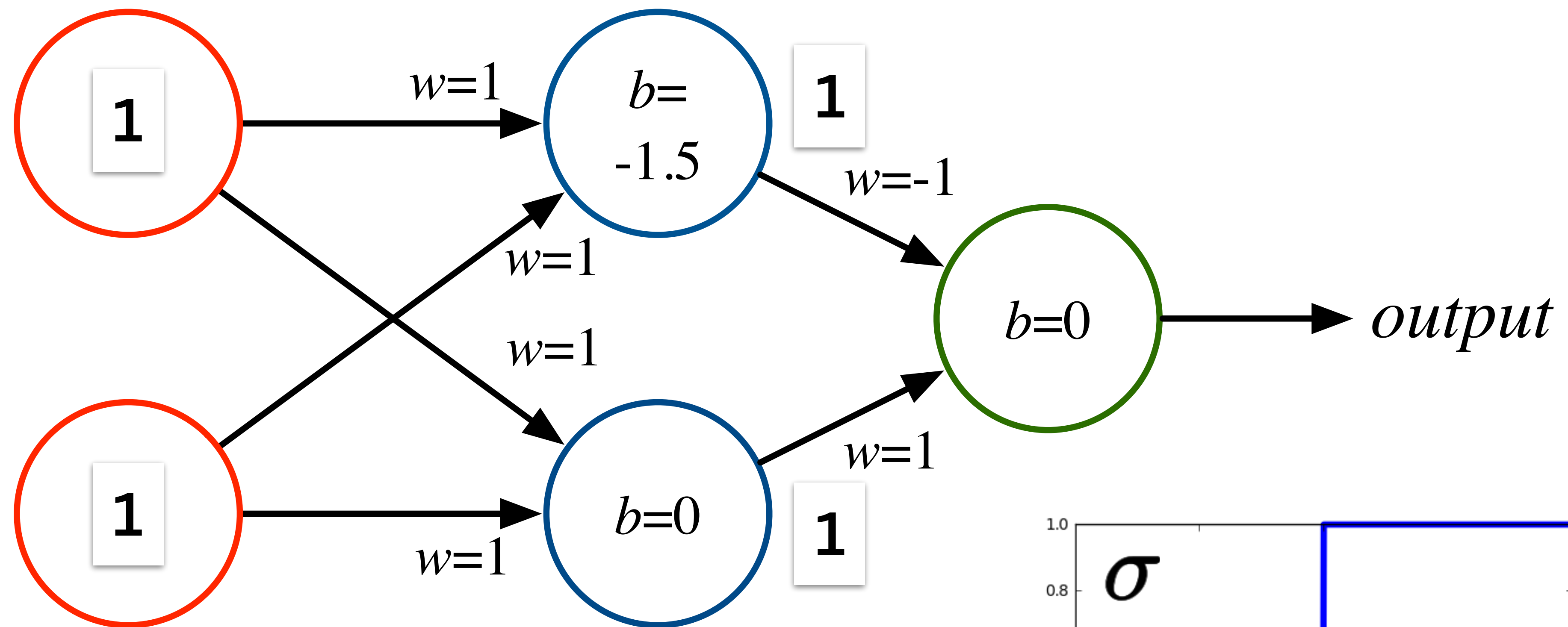
Neural Network Example



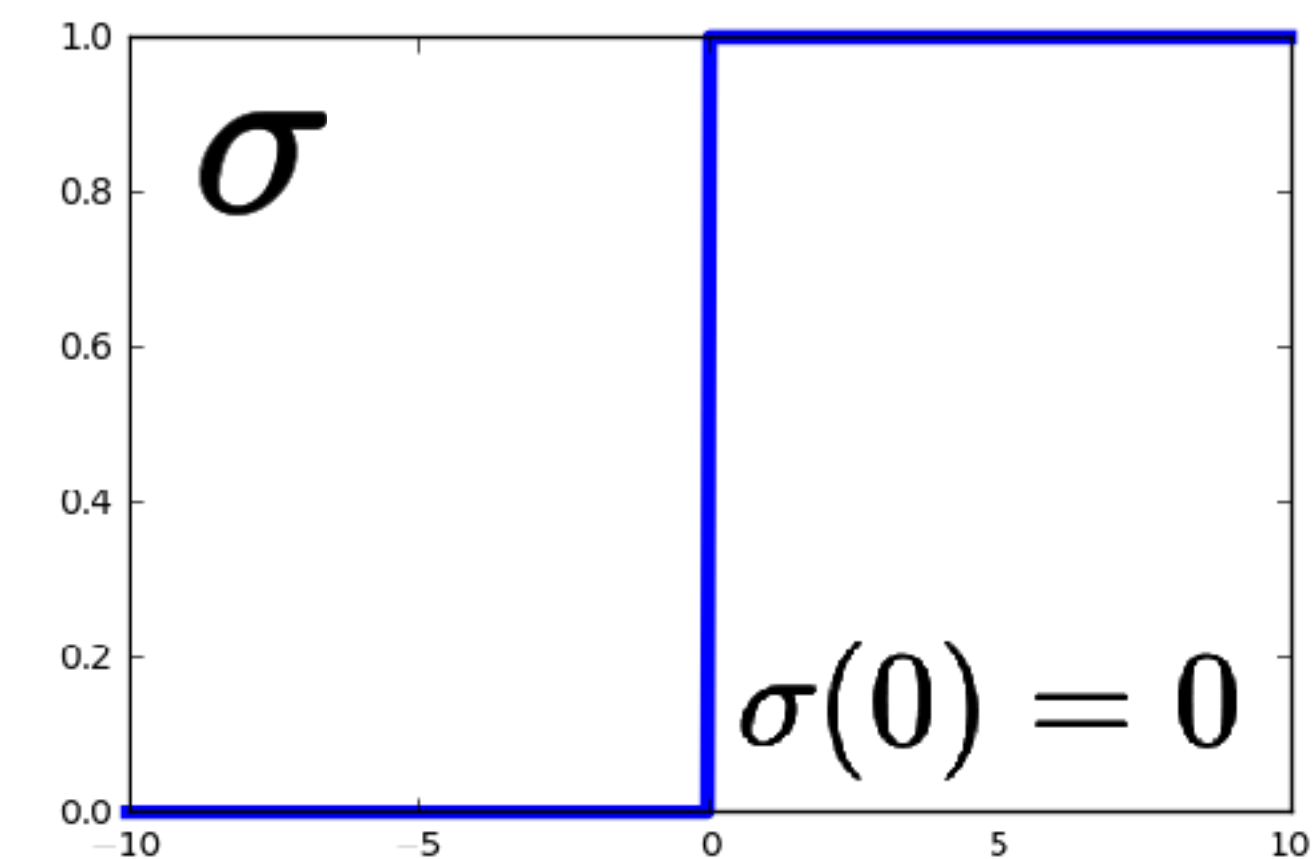
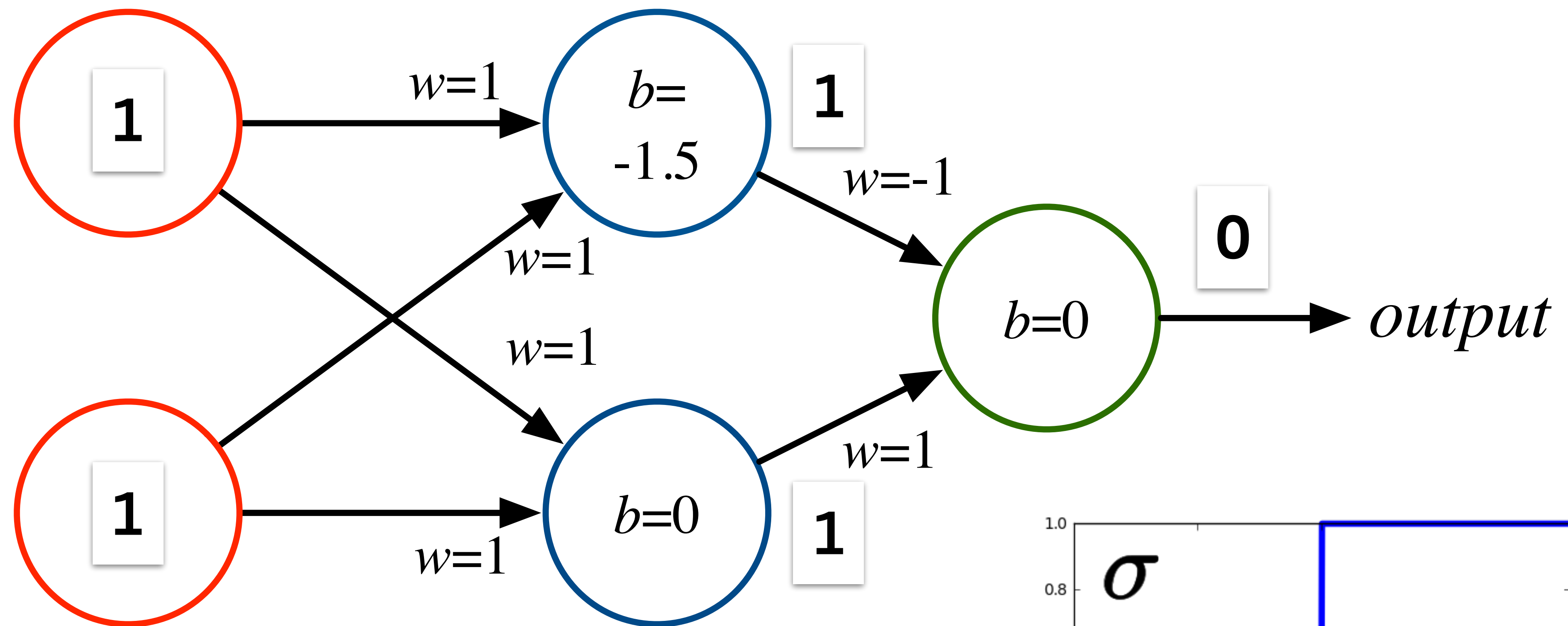
Neural Network Example



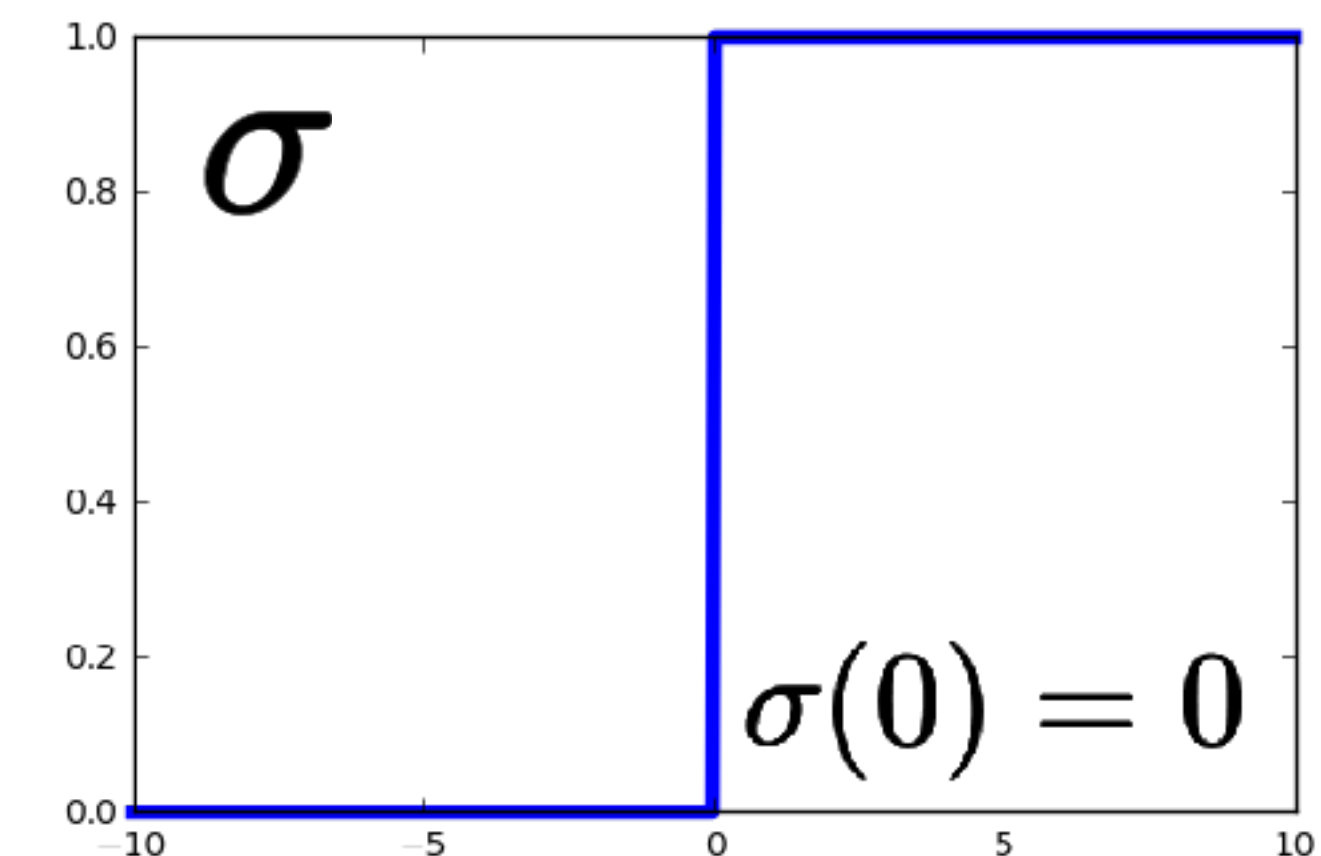
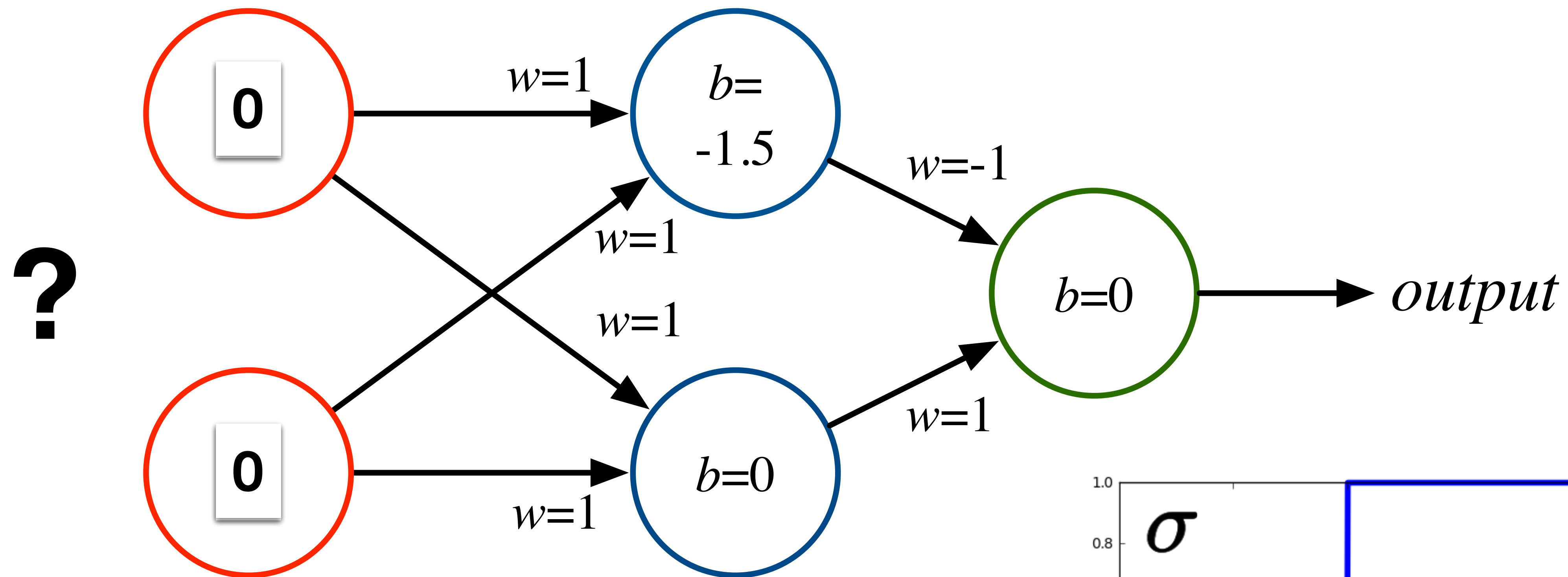
Neural Network Example



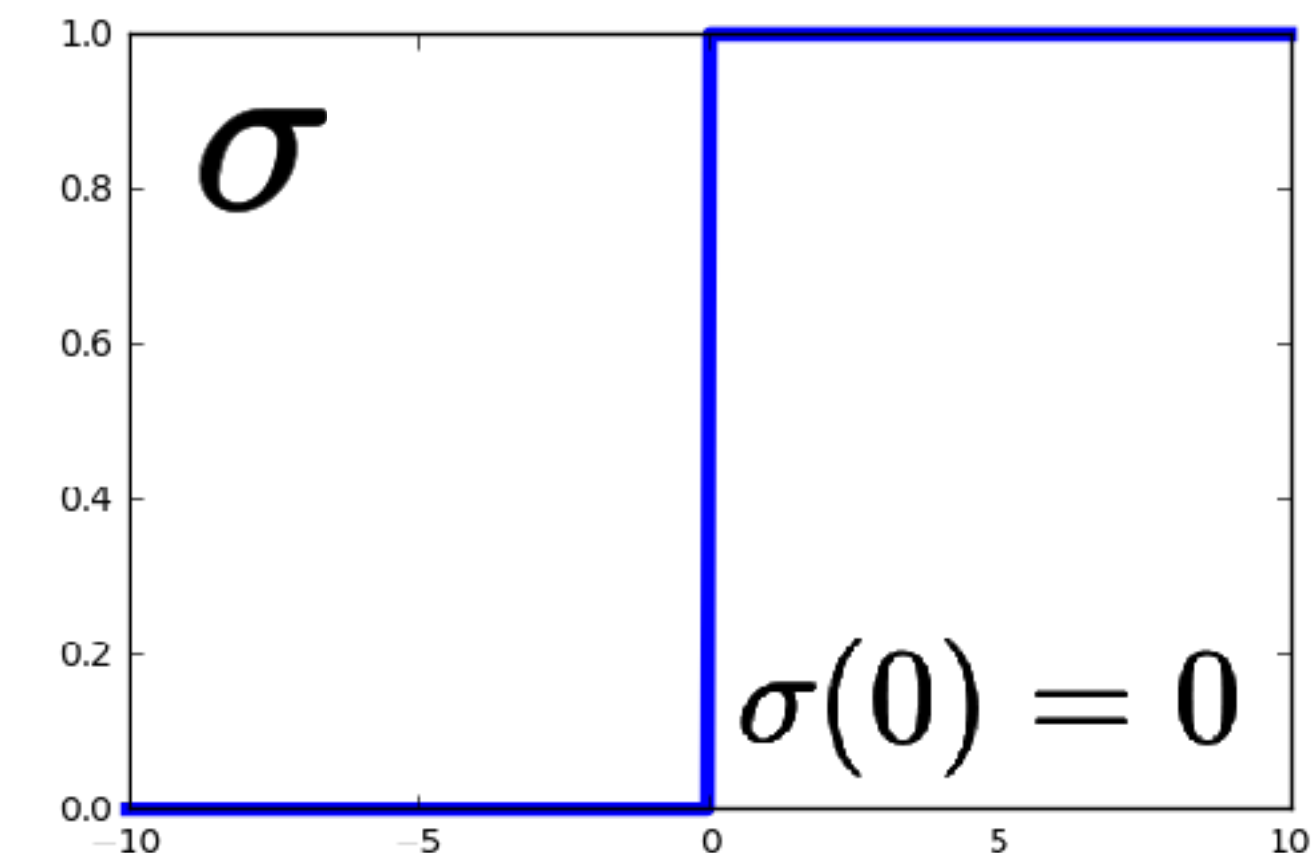
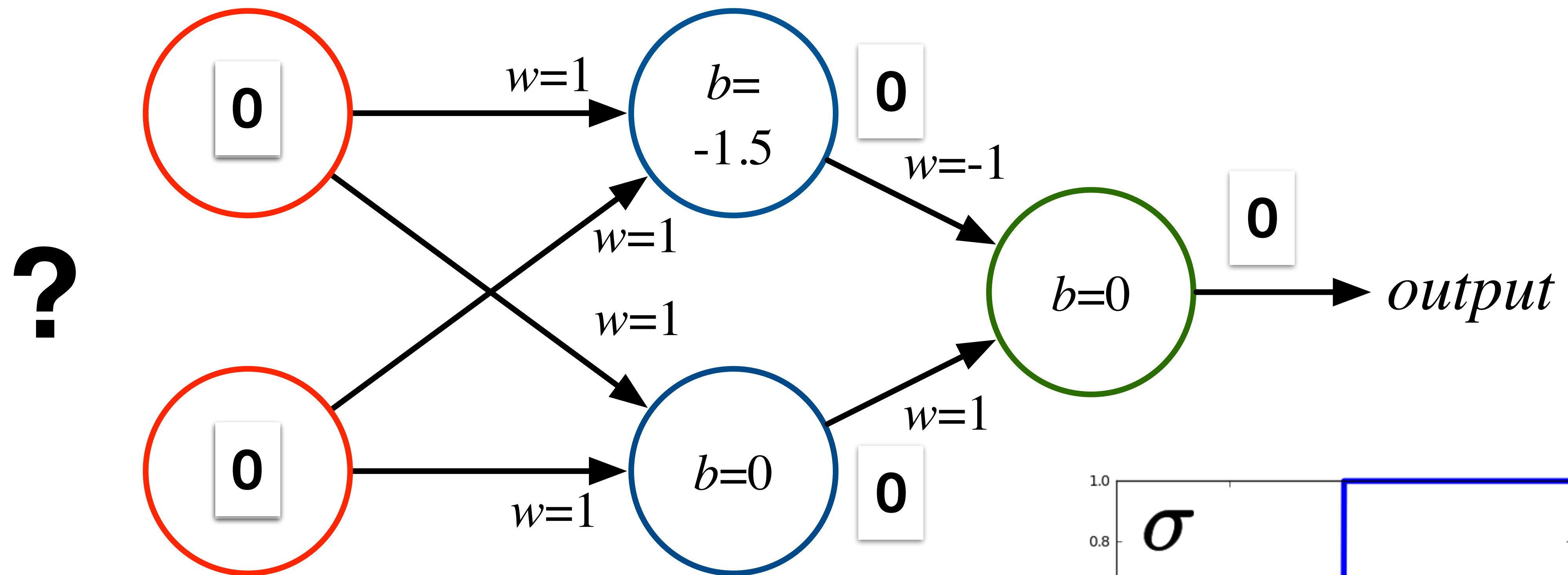
Neural Network Example



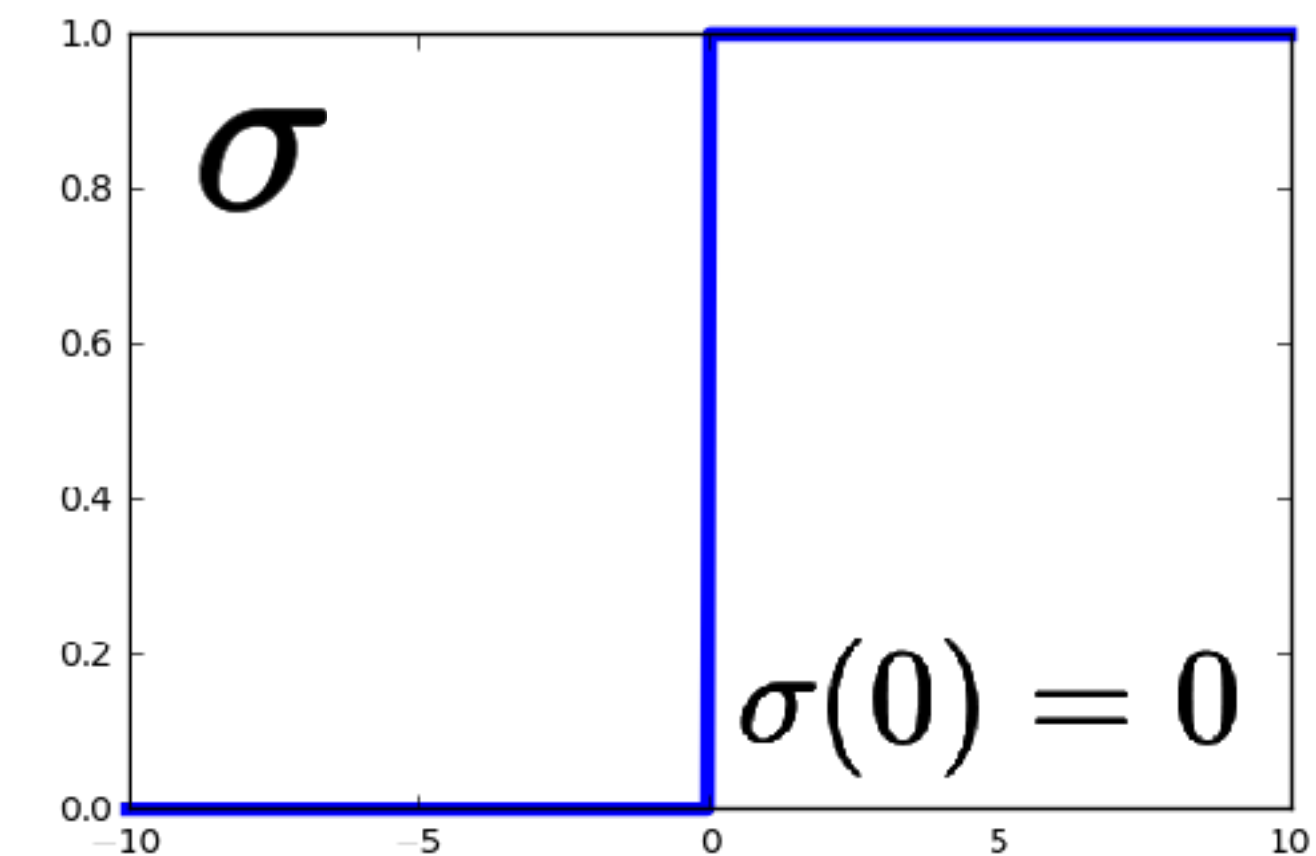
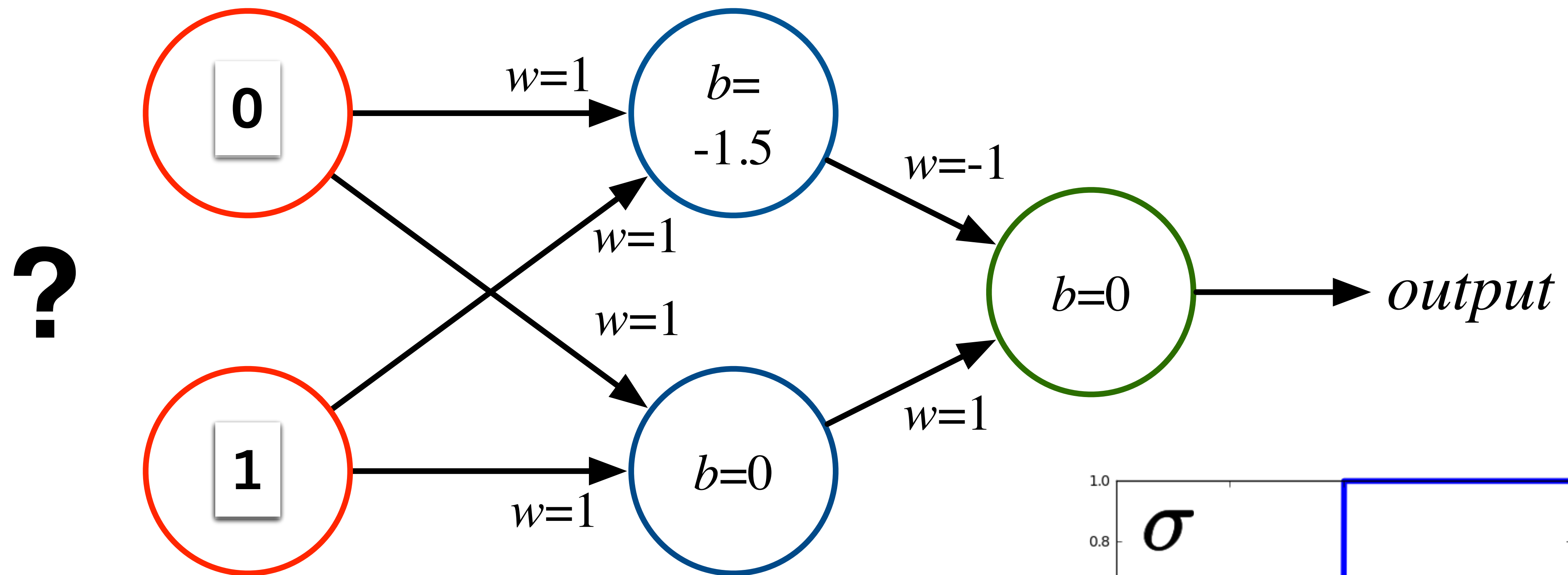
Neural Network Example



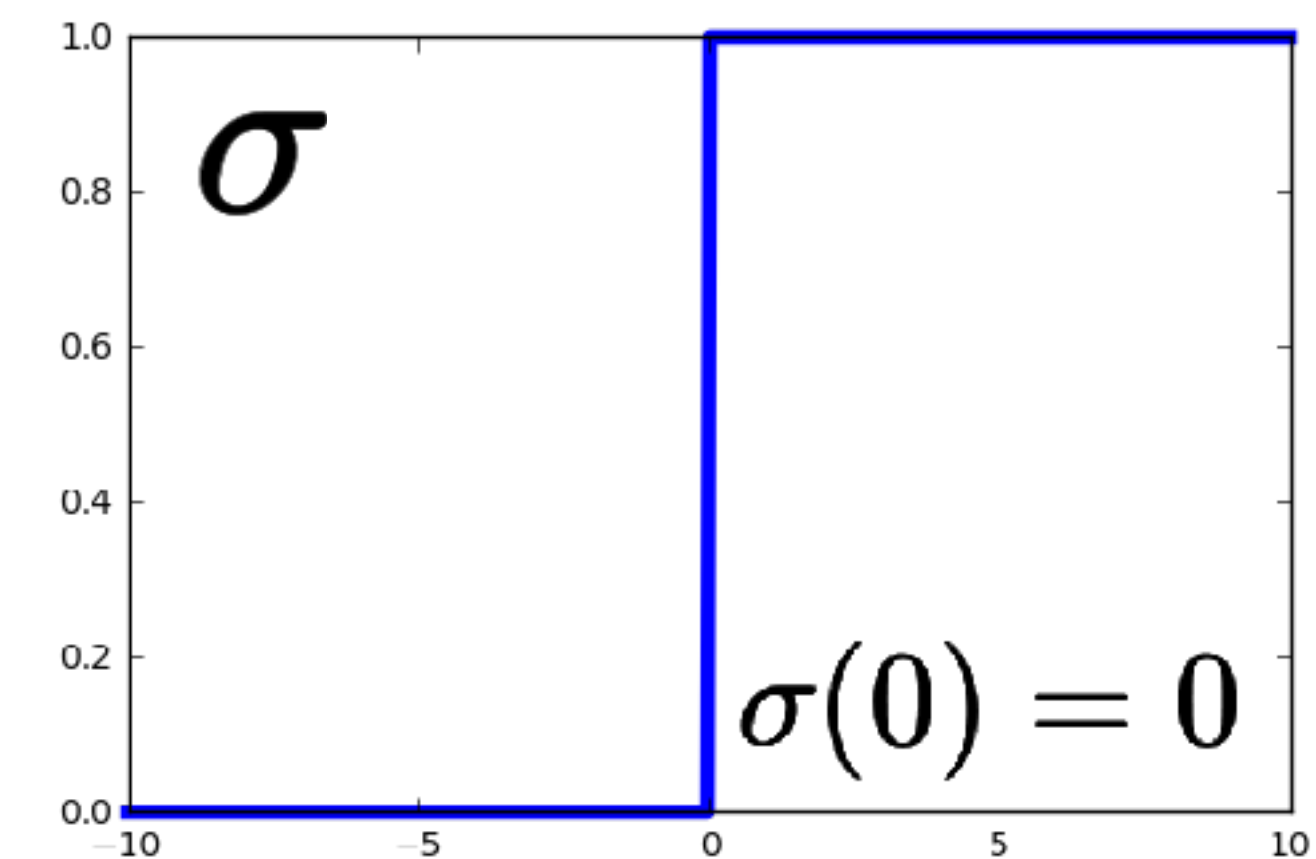
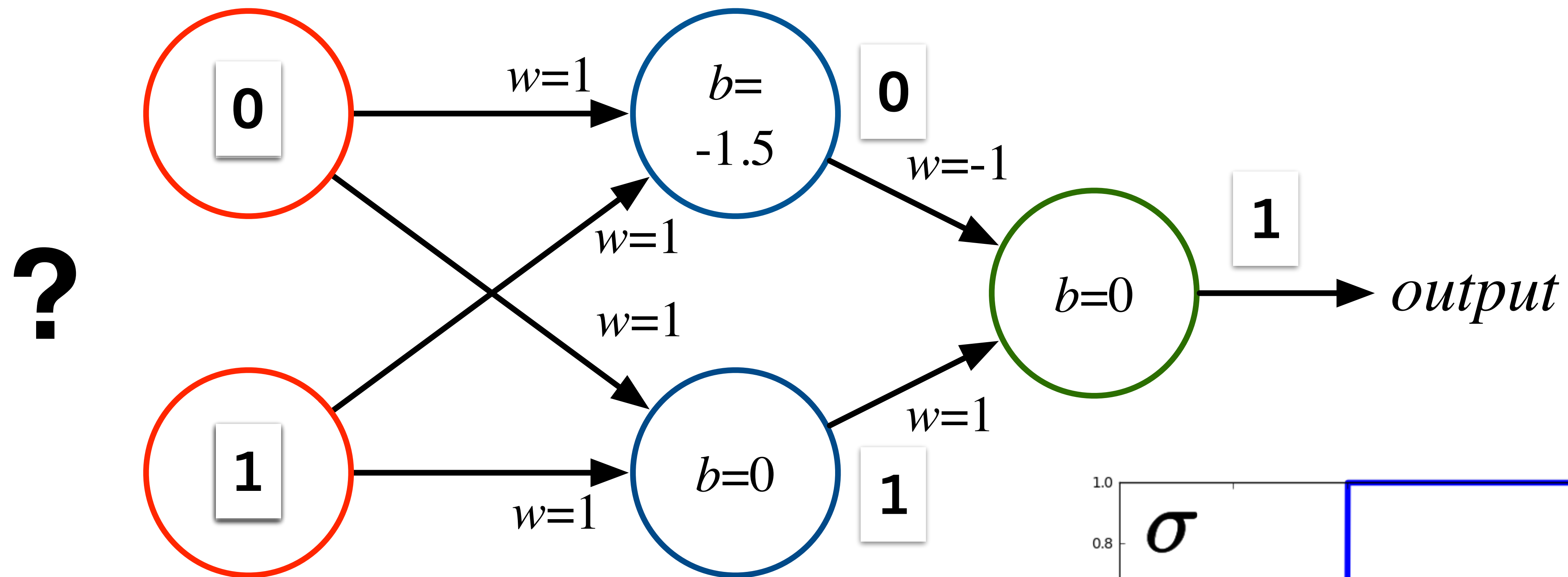
Neural Network Example



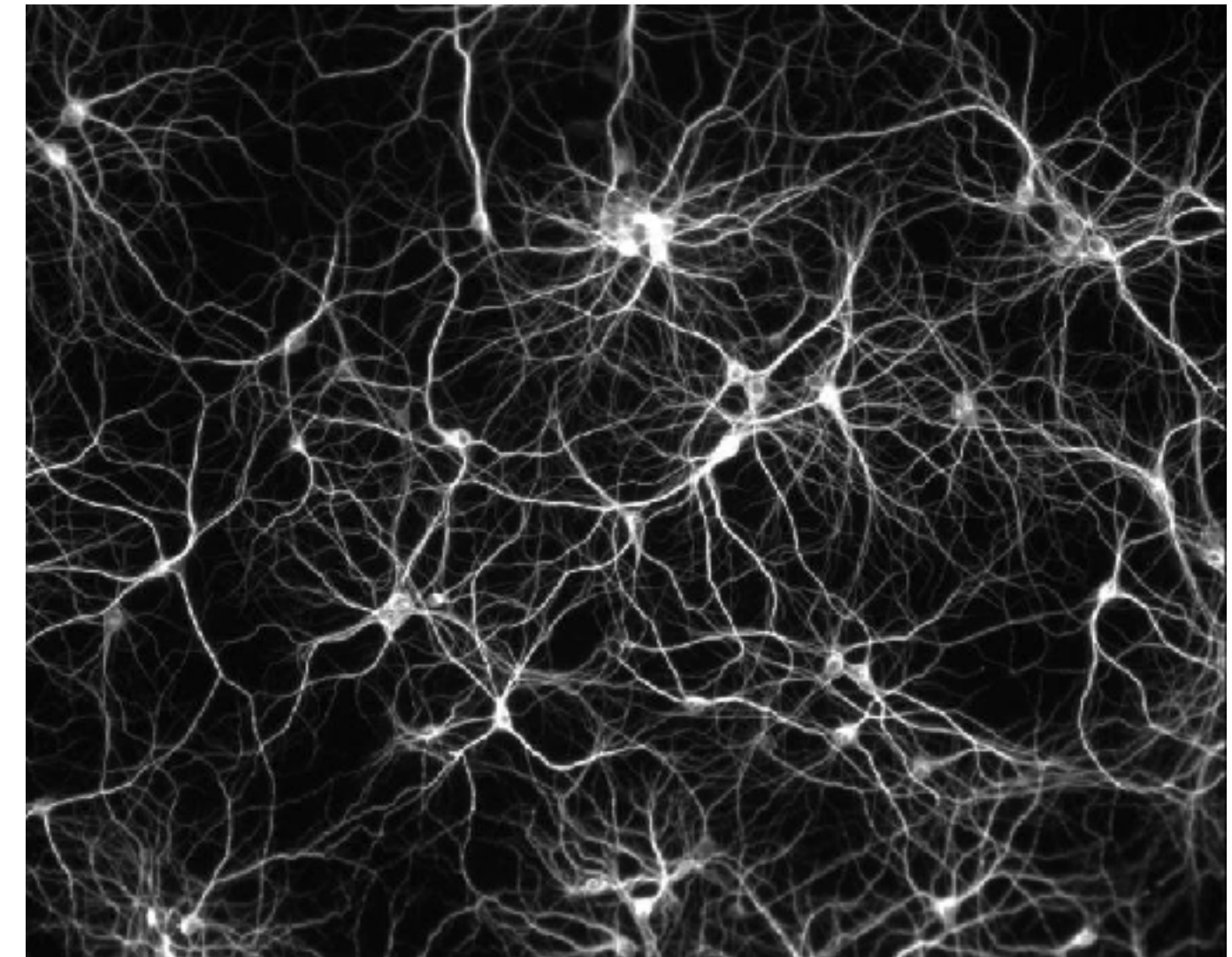
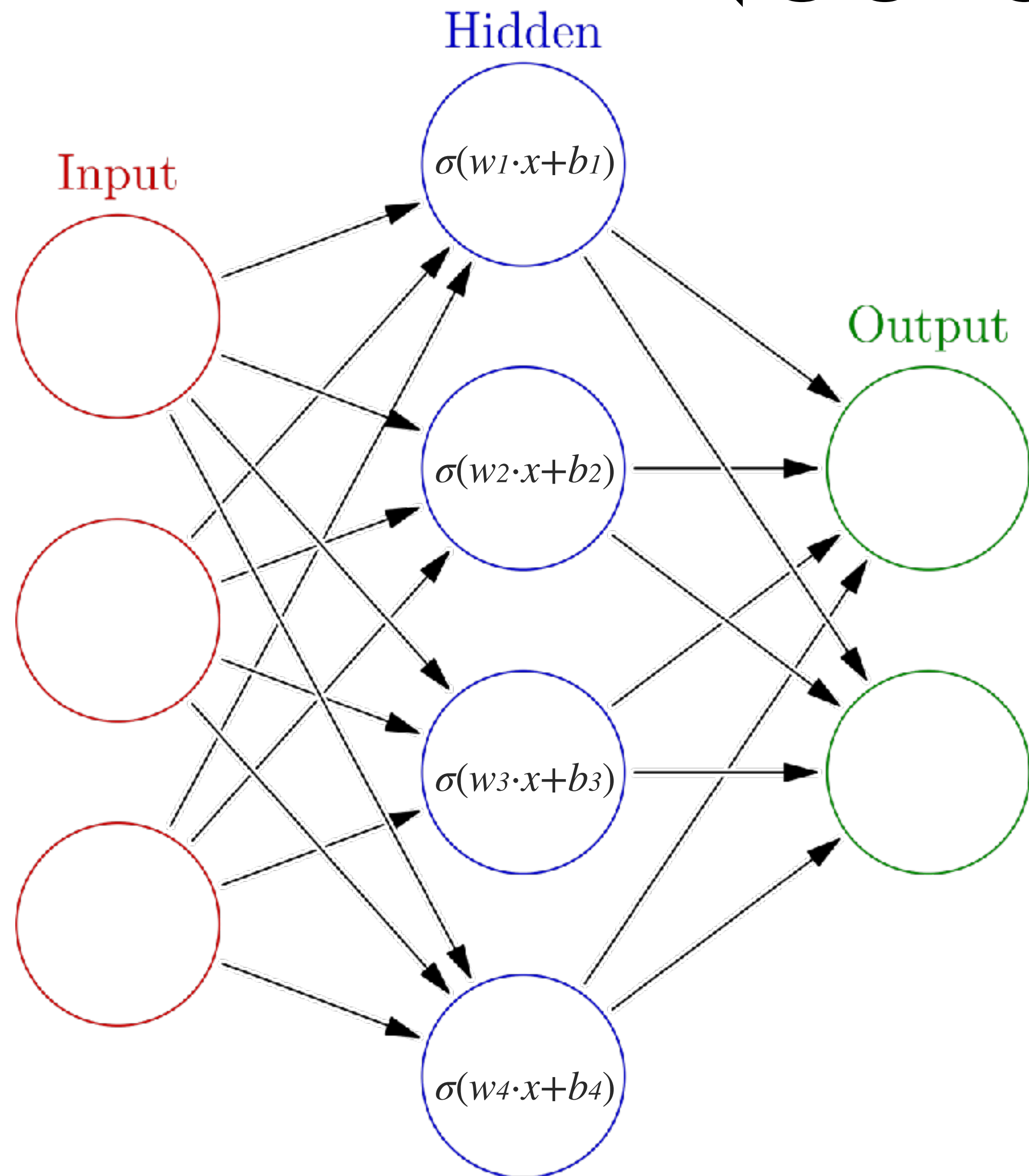
Neural Network Example



Neural Network Example

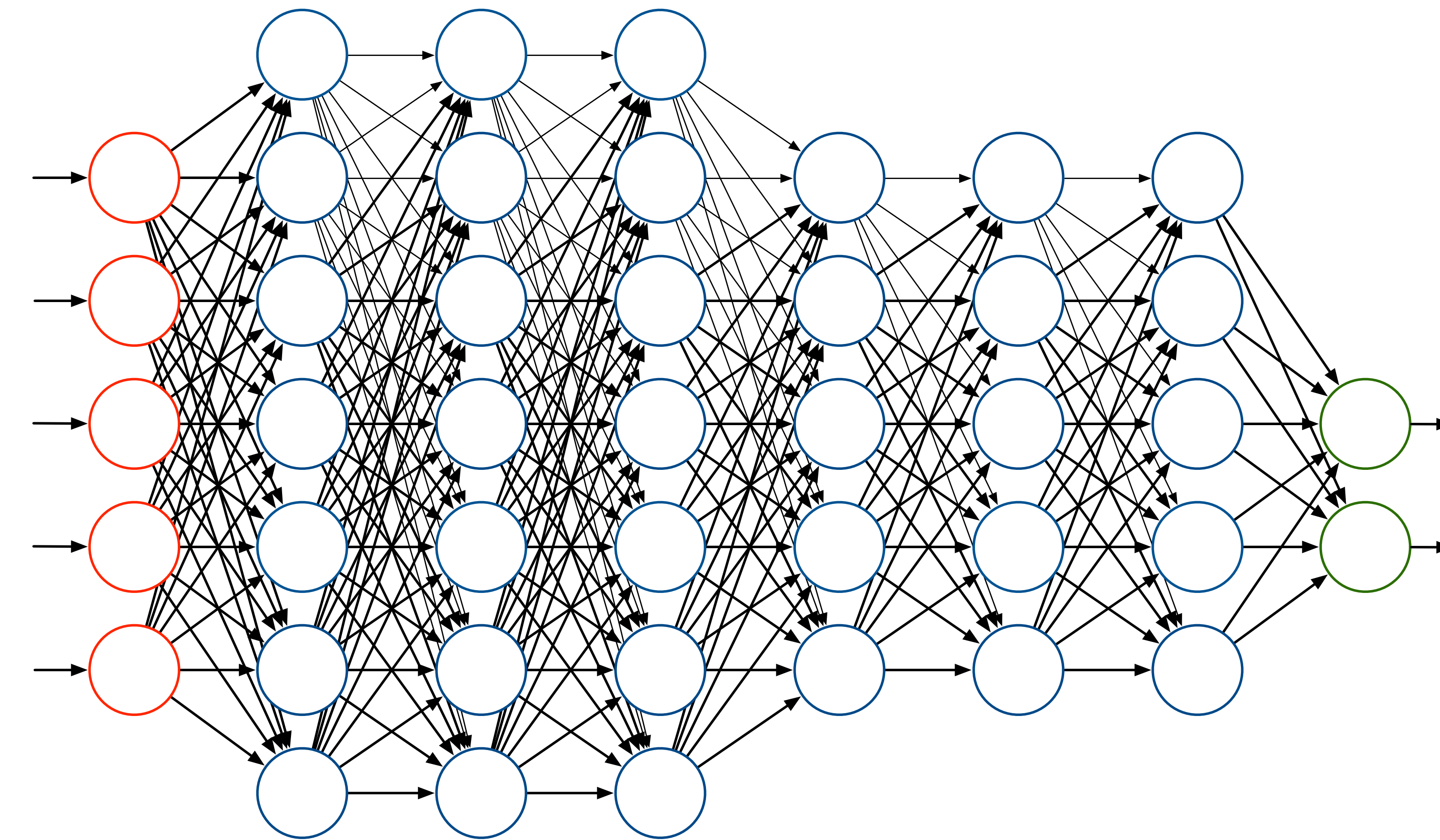


Neural Networks



The **universal approximation theorem** states that, under reasonable assumptions, a feedforward **neural network** with a finite number of nodes **can approximate any continuous** function to within a given error over a bounded input domain.

Deep Learning



Deep Learning

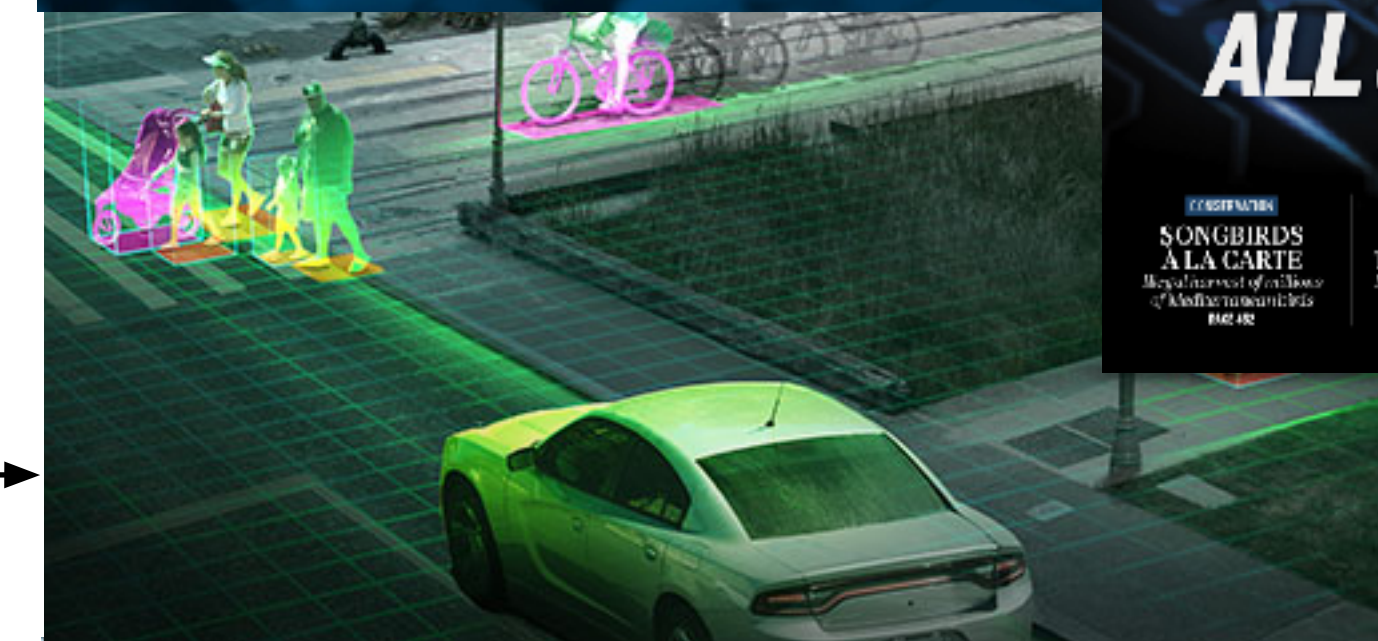
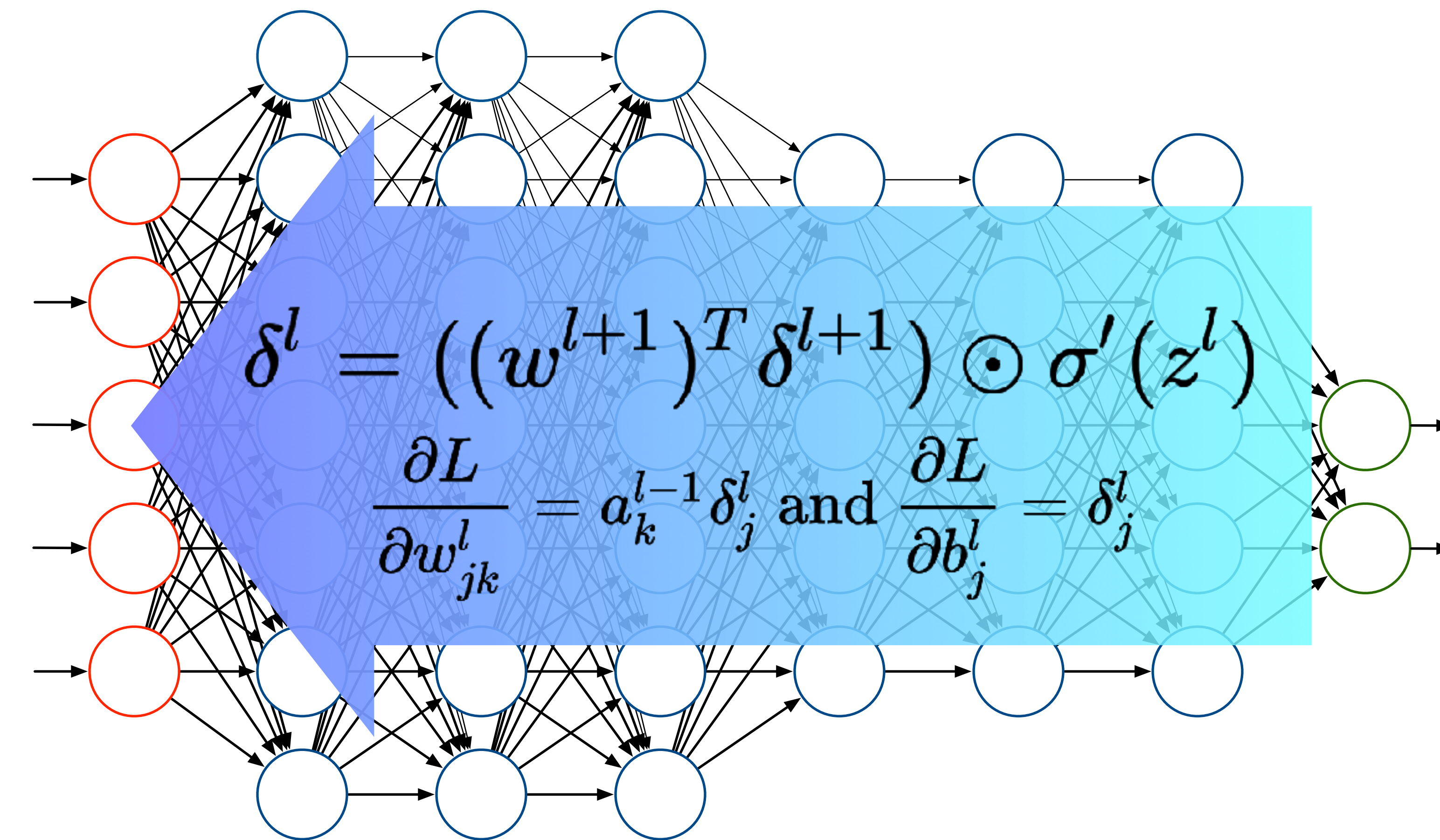
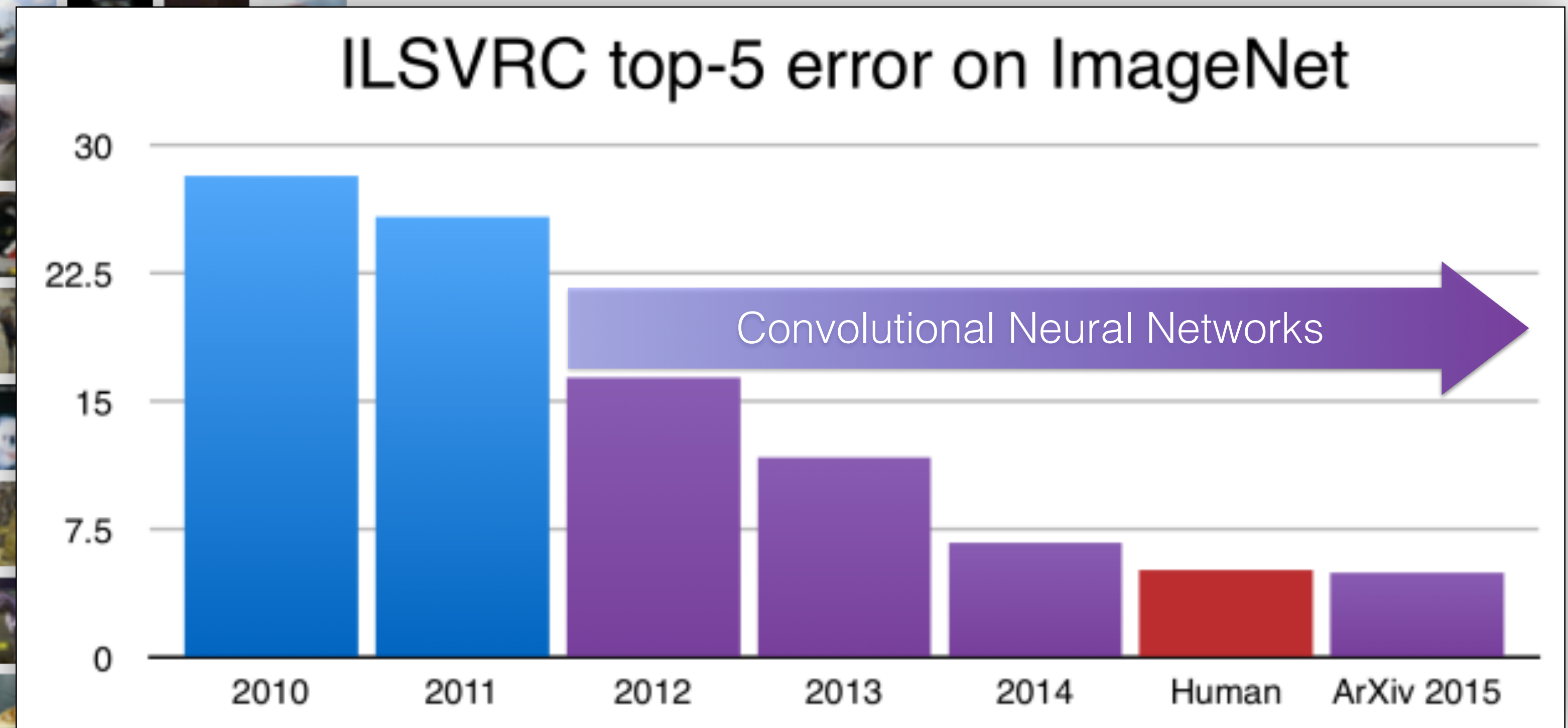
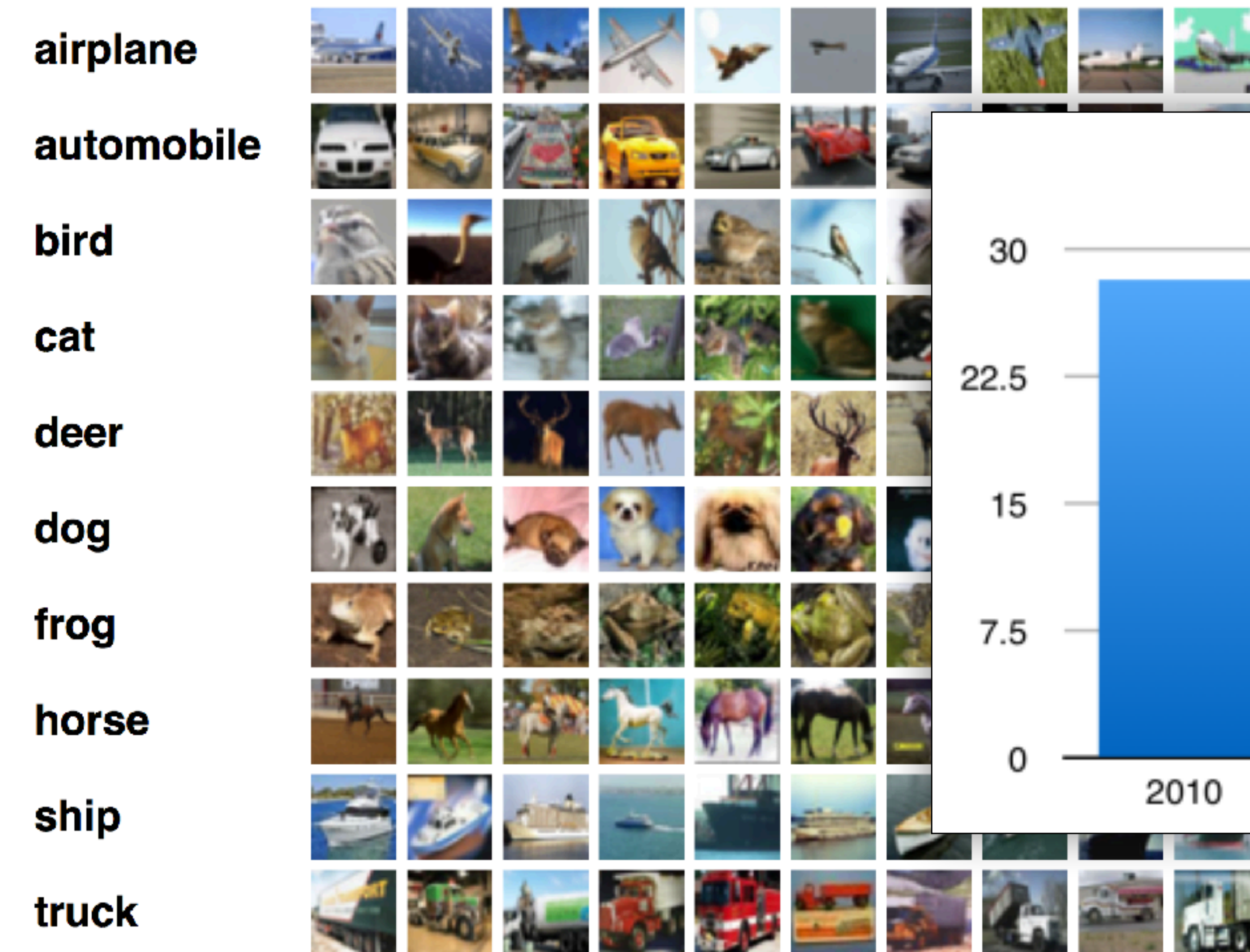
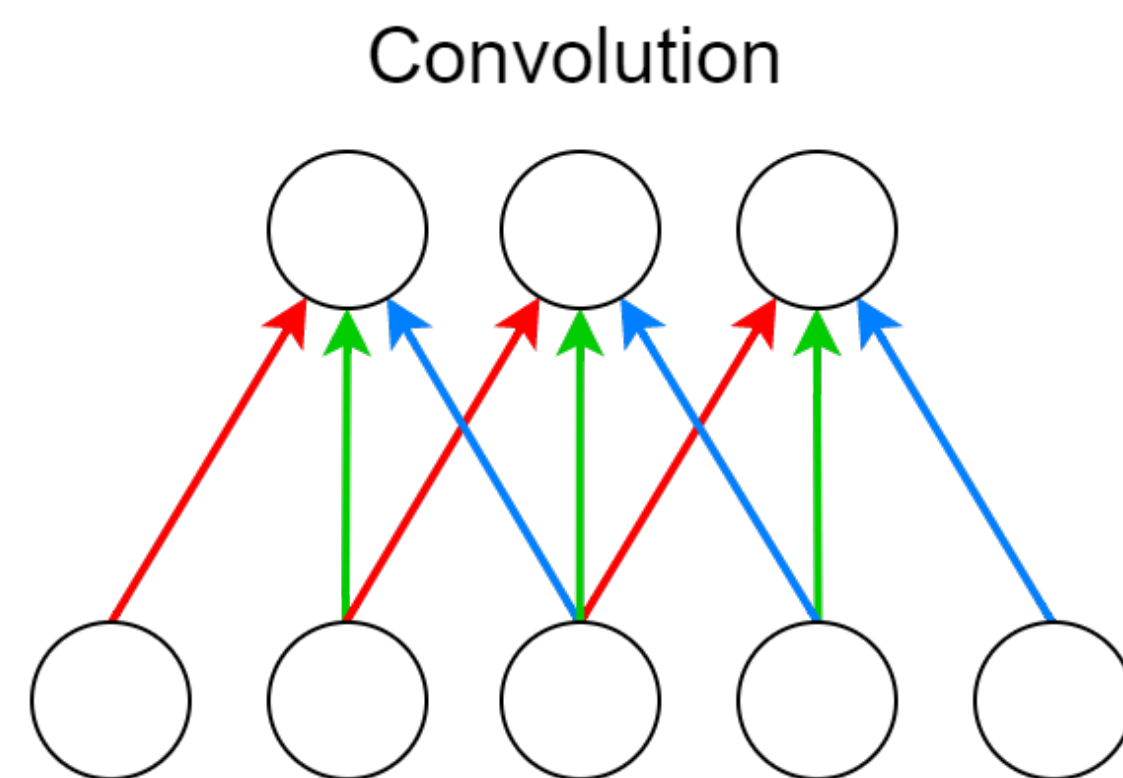
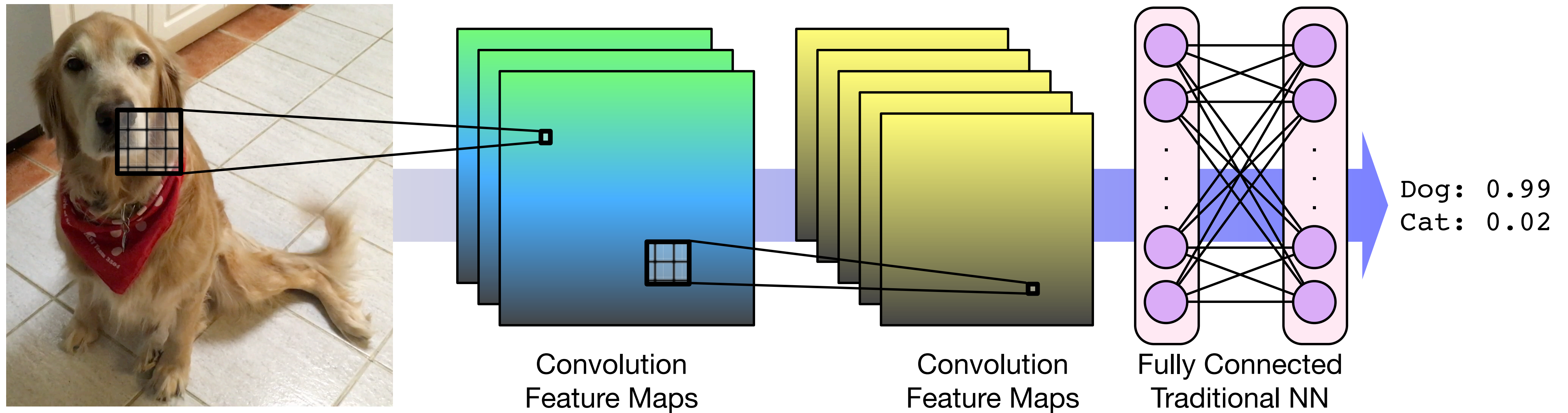


Image Recognition

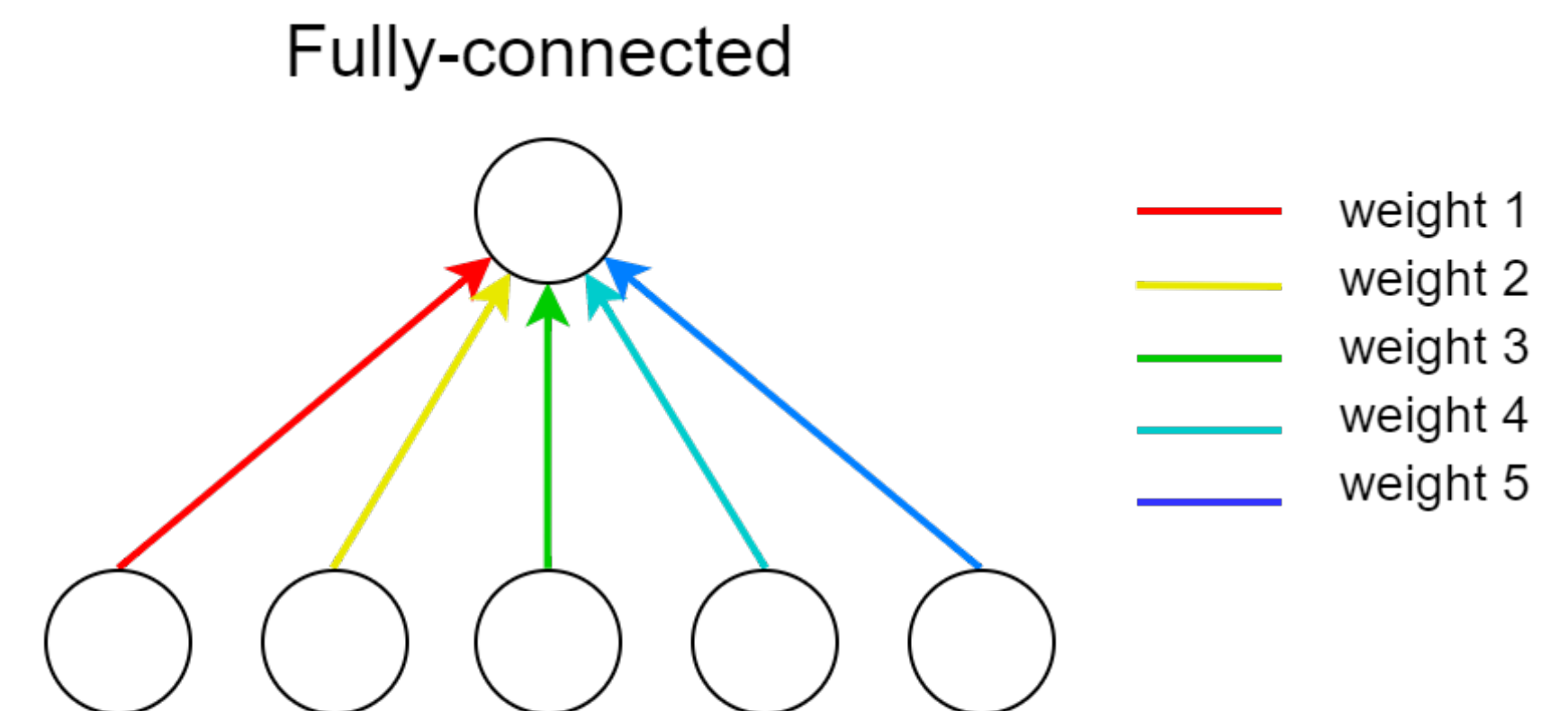


<https://devblogs.nvidia.com>

Convolutional Neural Networks



— weight 1
— weight 2
— weight 3

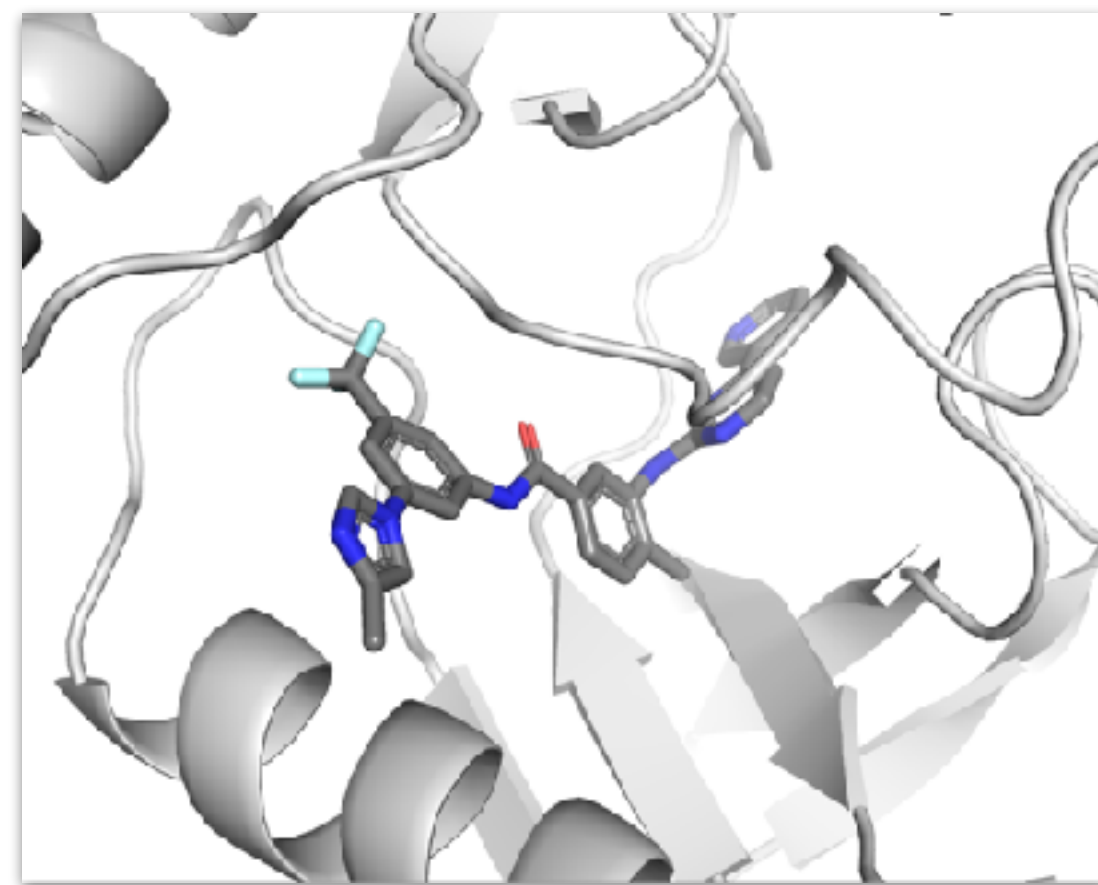


— weight 1
— weight 2
— weight 3
— weight 4
— weight 5

CNNs for Protein-Ligand Scoring



CNNs for Protein-Ligand Scoring



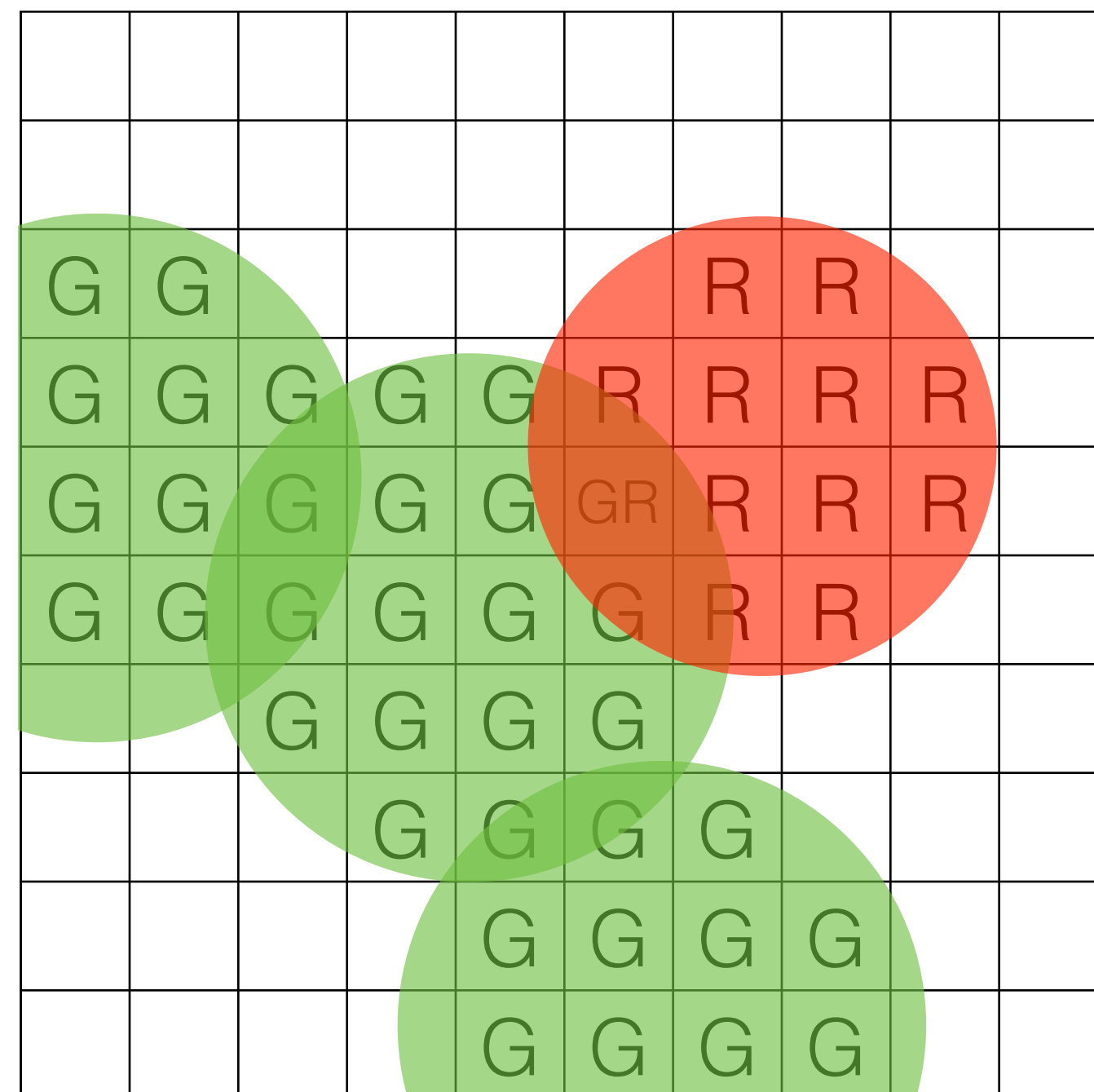
- Input representation
- Training
- Model optimization
- Visualize and Evaluation

Pose Prediction

Binding
Discrimination

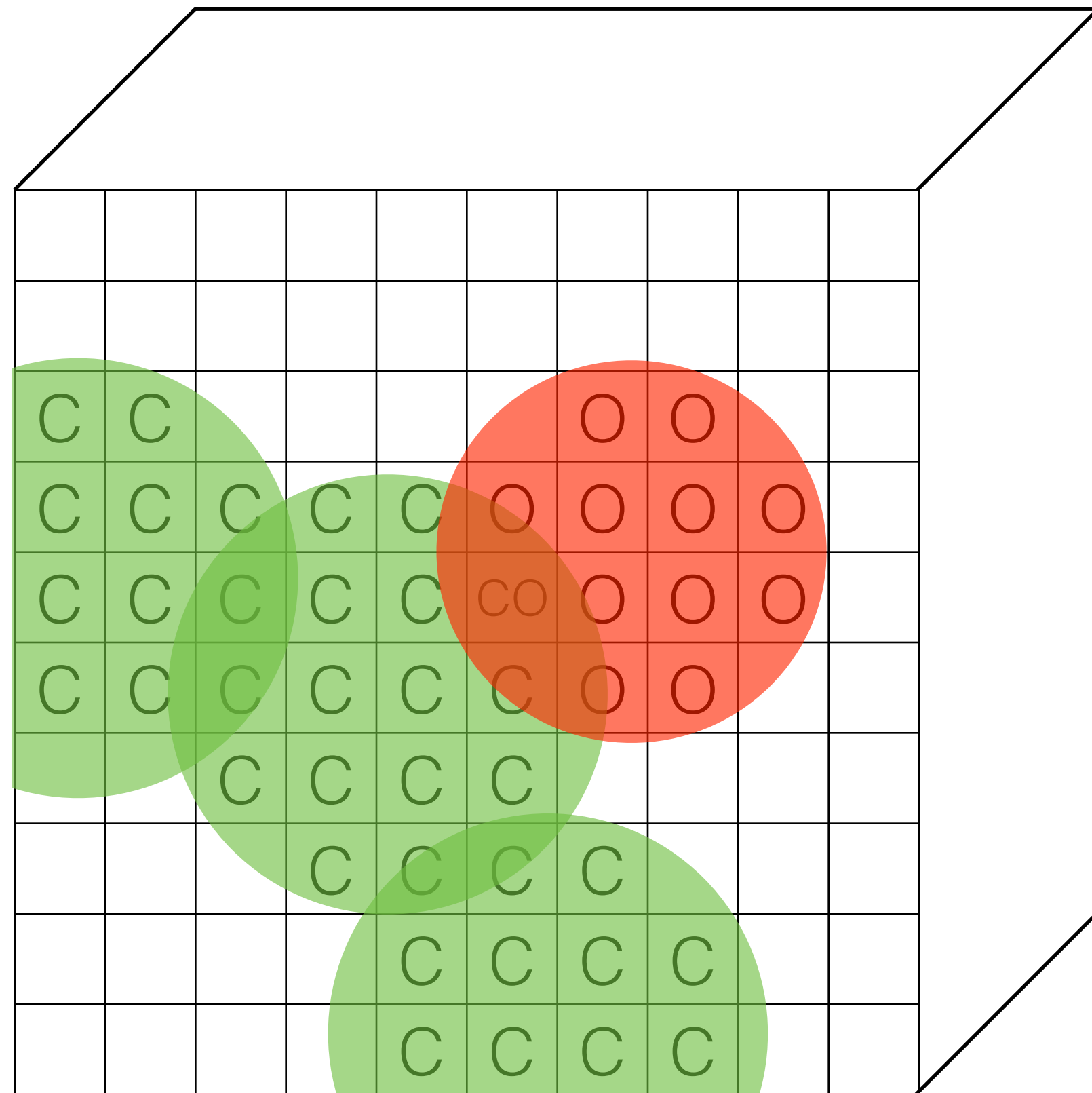
Affinity Prediction

Protein-Ligand Representation



(R,G,B) pixel

Protein-Ligand Representation



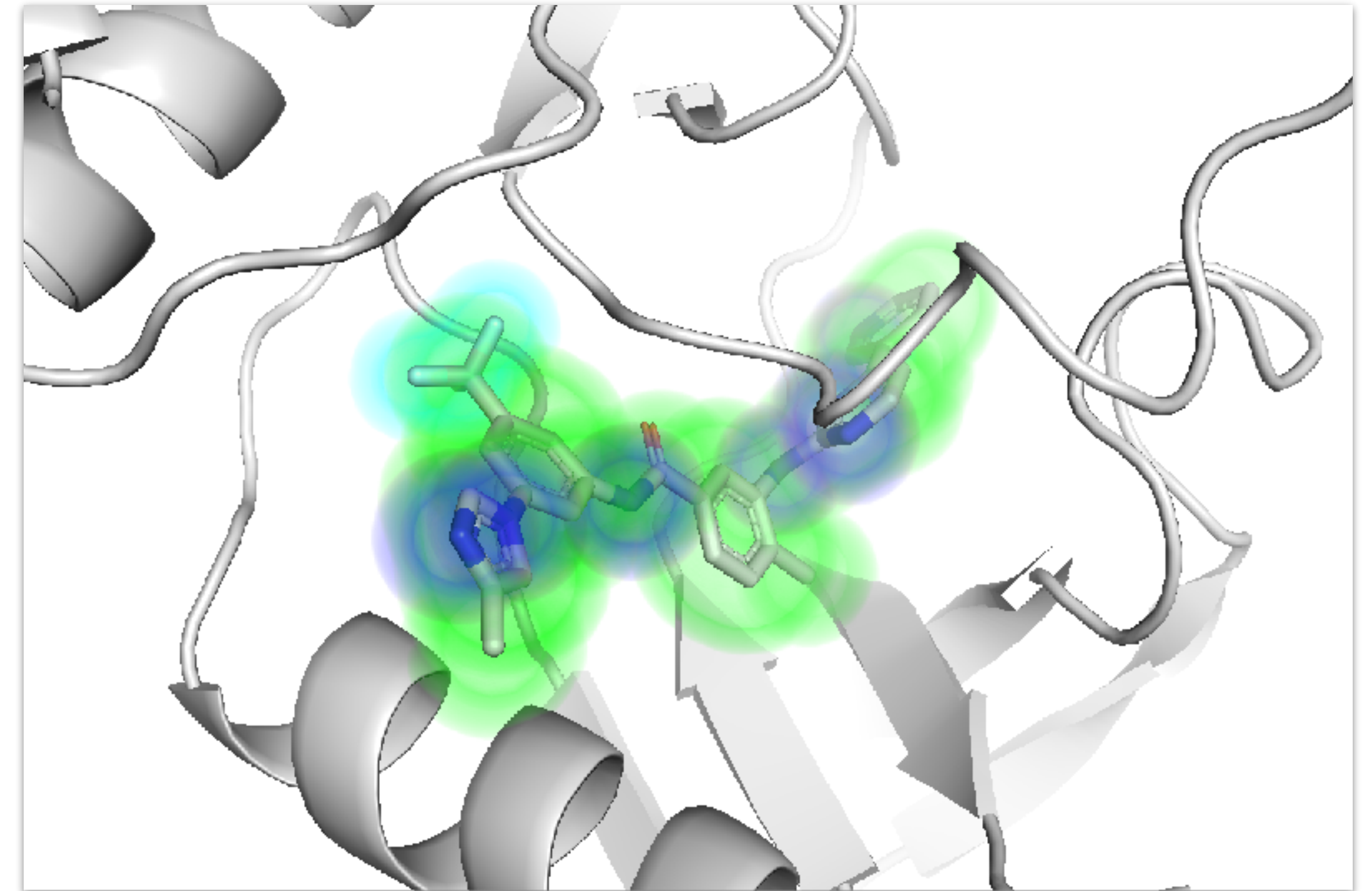
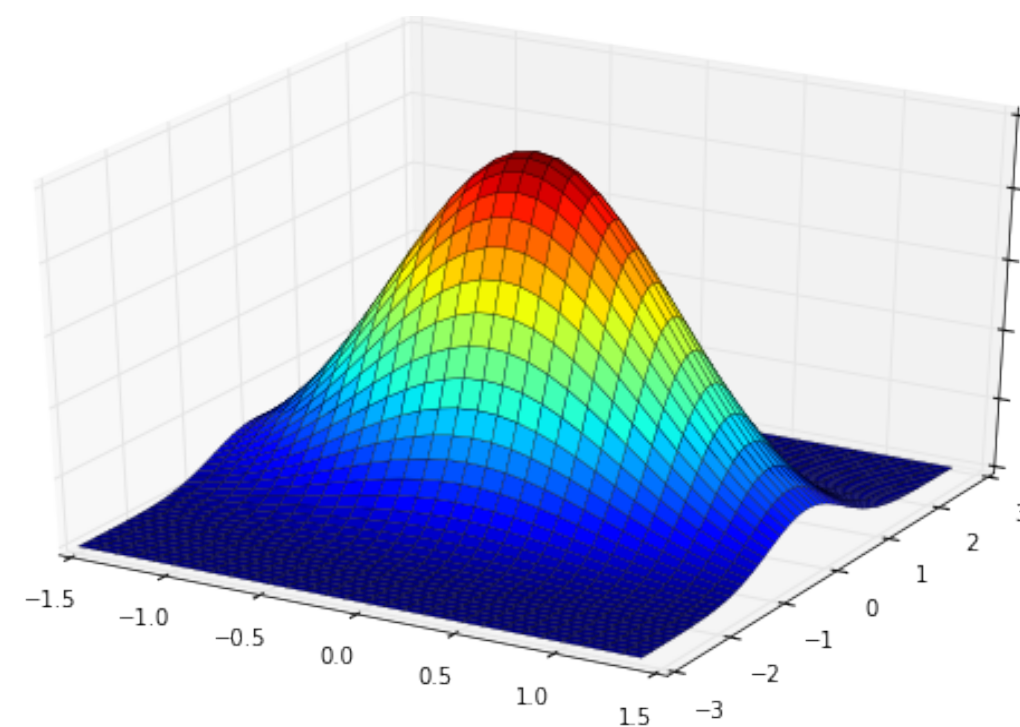
(R,G,B) pixel \rightarrow

(Carbon, Nitrogen, Oxygen,...) **voxel**

The only parameters for this representation are the choice of **grid resolution**, **atom density**, and **atom types**.

Atom Density

$$A(d, r) = \begin{cases} e^{-\frac{2d^2}{r^2}} & 0 \leq d < r \\ \frac{4}{e^2 r^2} d^2 - \frac{12}{e^2 r} d + \frac{9}{e^2} & r \leq d < 1.5r \\ 0 & d \geq 1.5r \end{cases}$$



Gaussian

Atom Types

Ligand

AliphaticCarbonXSHydrophobe
 AliphaticCarbonXSNonHydrophobe
 AromaticCarbonXSHydrophobe
 AromaticCarbonXSNonHydrophobe

Bromine

Chlorine

Fluorine

Iodine

Nitrogen

NitrogenXSAcceptor

NitrogenXSDonor

NitrogenXSDonorAcceptor

Oxygen

OxygenXSAcceptor

OxygenXSDonorAcceptor

Phosphorus

Sulfur

SulfurAcceptor

Receptor

AliphaticCarbonXSHydrophobe
 AliphaticCarbonXSNonHydrophobe
 AromaticCarbonXSHydrophobe
 AromaticCarbonXSNonHydrophobe

Calcium

Iron

Magnesium

Nitrogen

NitrogenXSAcceptor

NitrogenXSDonor

NitrogenXSDonorAcceptor

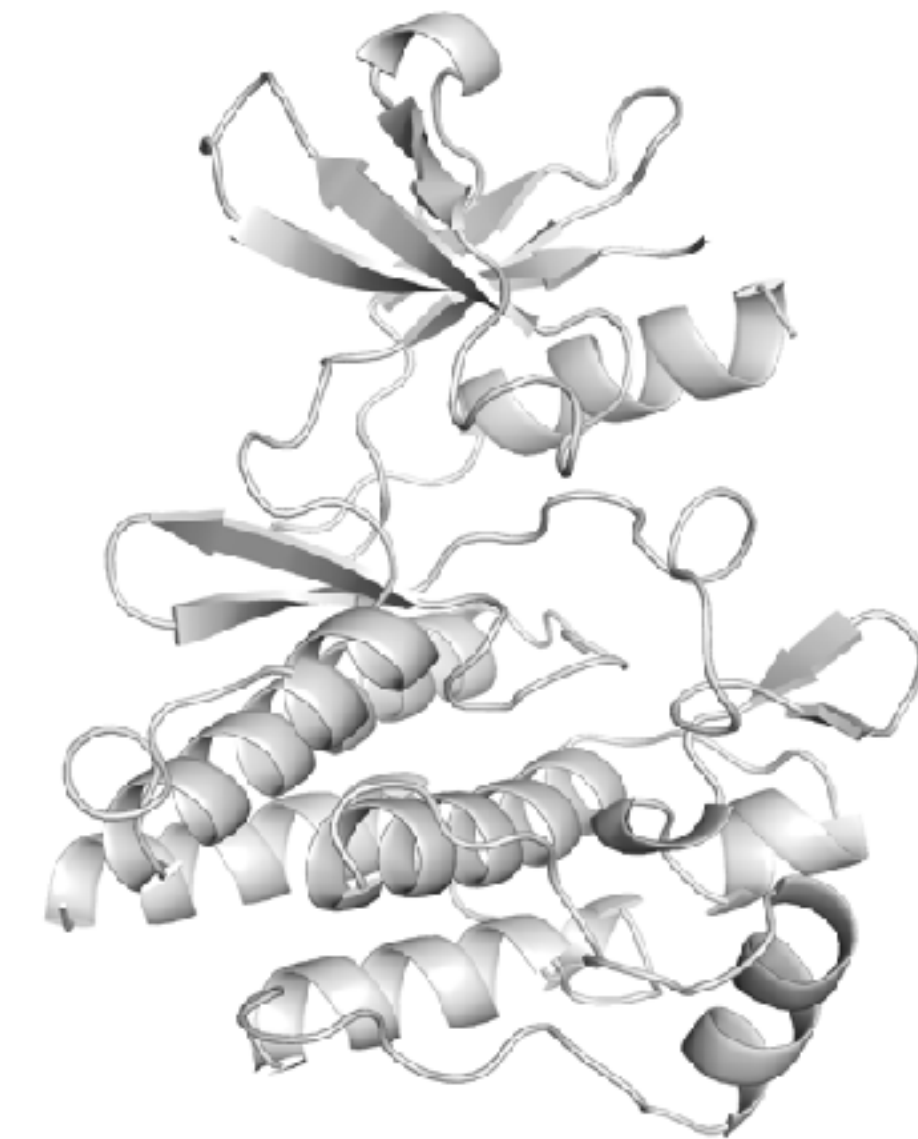
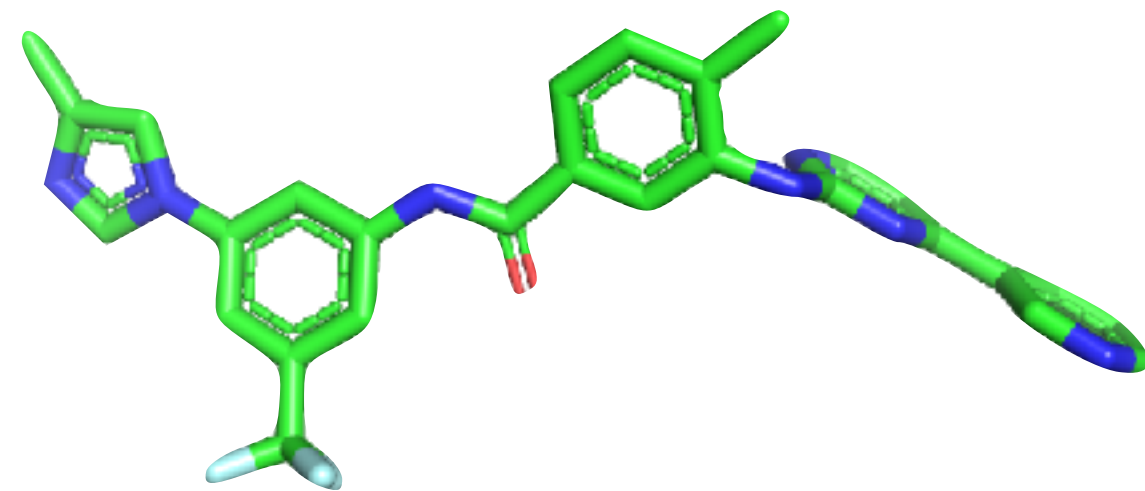
OxygenXSAcceptor

OxygenXSDonorAcceptor

Phosphorus

Sulfur

Zinc



Training Data

Pose Prediction



- 337 protein-ligand complexes
- curated for electron density
 - diverse targets
 - $<10\mu\text{M}$ affinity
 - **generate poses** with Vina
 - 745 $<2\text{\AA}$ RMSD (actives)
 - 3251 $>4\text{\AA}$ RMSD (decoys)



- 4056 protein-ligand complexes
- diverse targets
 - wide range of affinities
 - **generate poses** with AutoDock Vina
 - include minimized crystal pose
 - 8,688 $<2\text{\AA}$ RMSD (actives)
 - 76,743 $>4\text{\AA}$ RMSD (decoys)

Training Data

Binding Discrimination

D U D • E

102 targets

- 22,645 actives
- 1,407,145 decoys
- <10 μ M affinity
- **true poses unknown**
- **trust** docked poses

Affinity Prediction

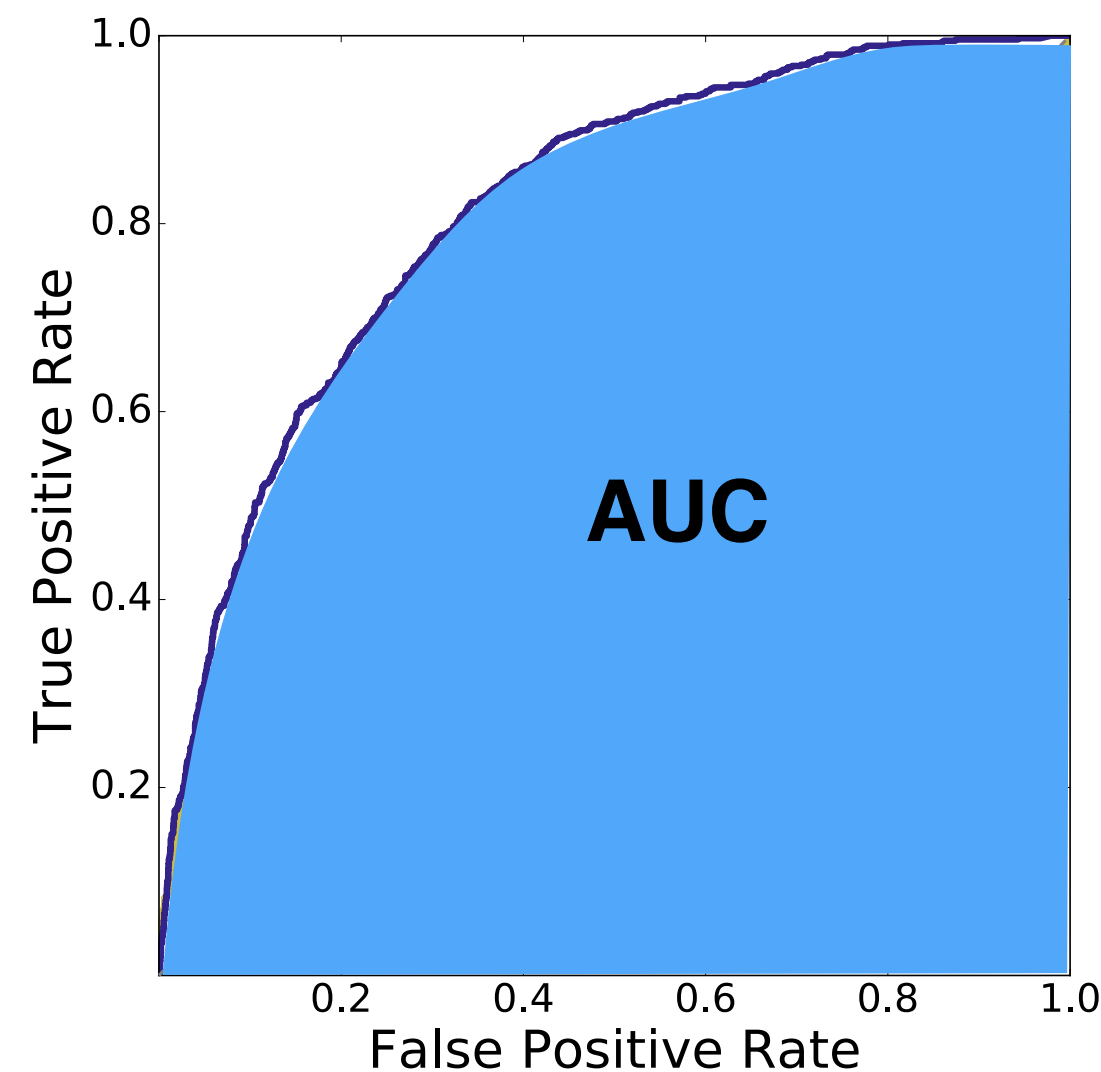


- 8,688 low RMSD poses
- assign known affinity
- **regression problem**

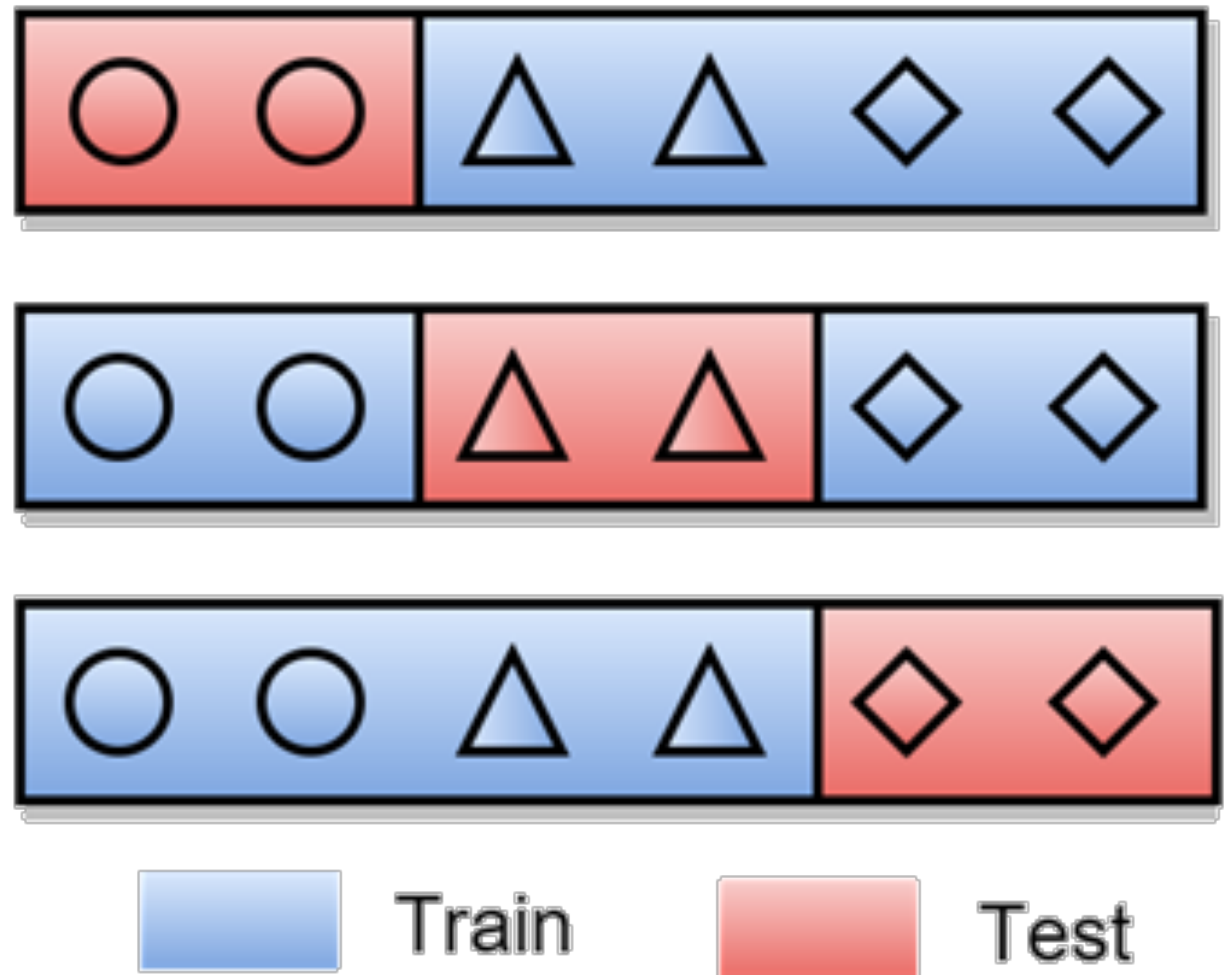
Model Evaluation

CSAR: >90% similar targets kept in same fold

DUD-E & PDBbind: >80% similar targets kept in same fold



Clustered Cross-validation



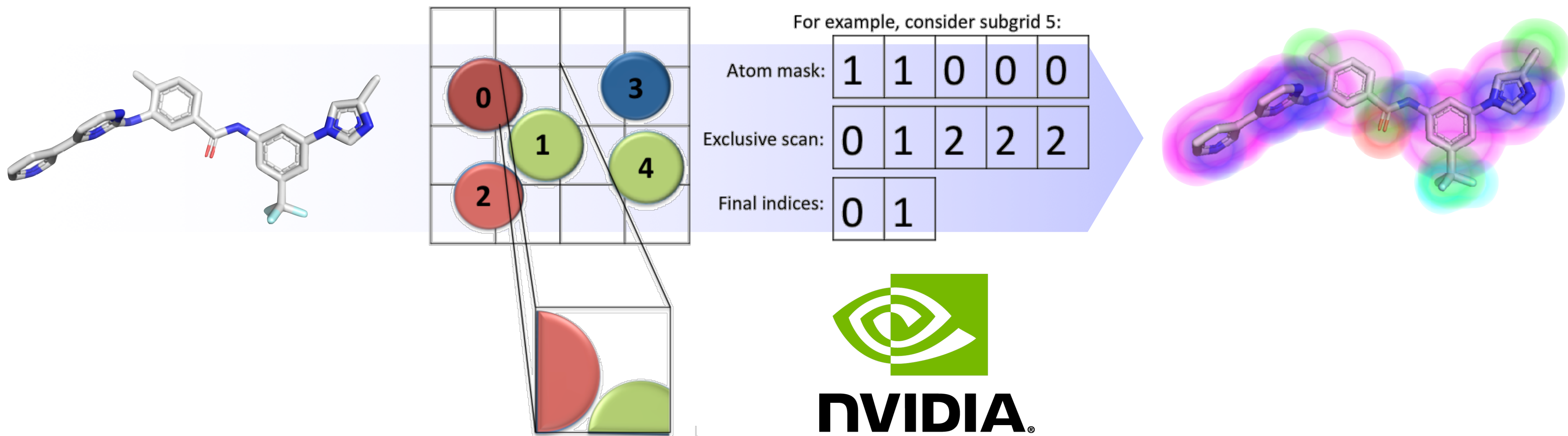
Model Training

Custom **MolGridDataLayer**

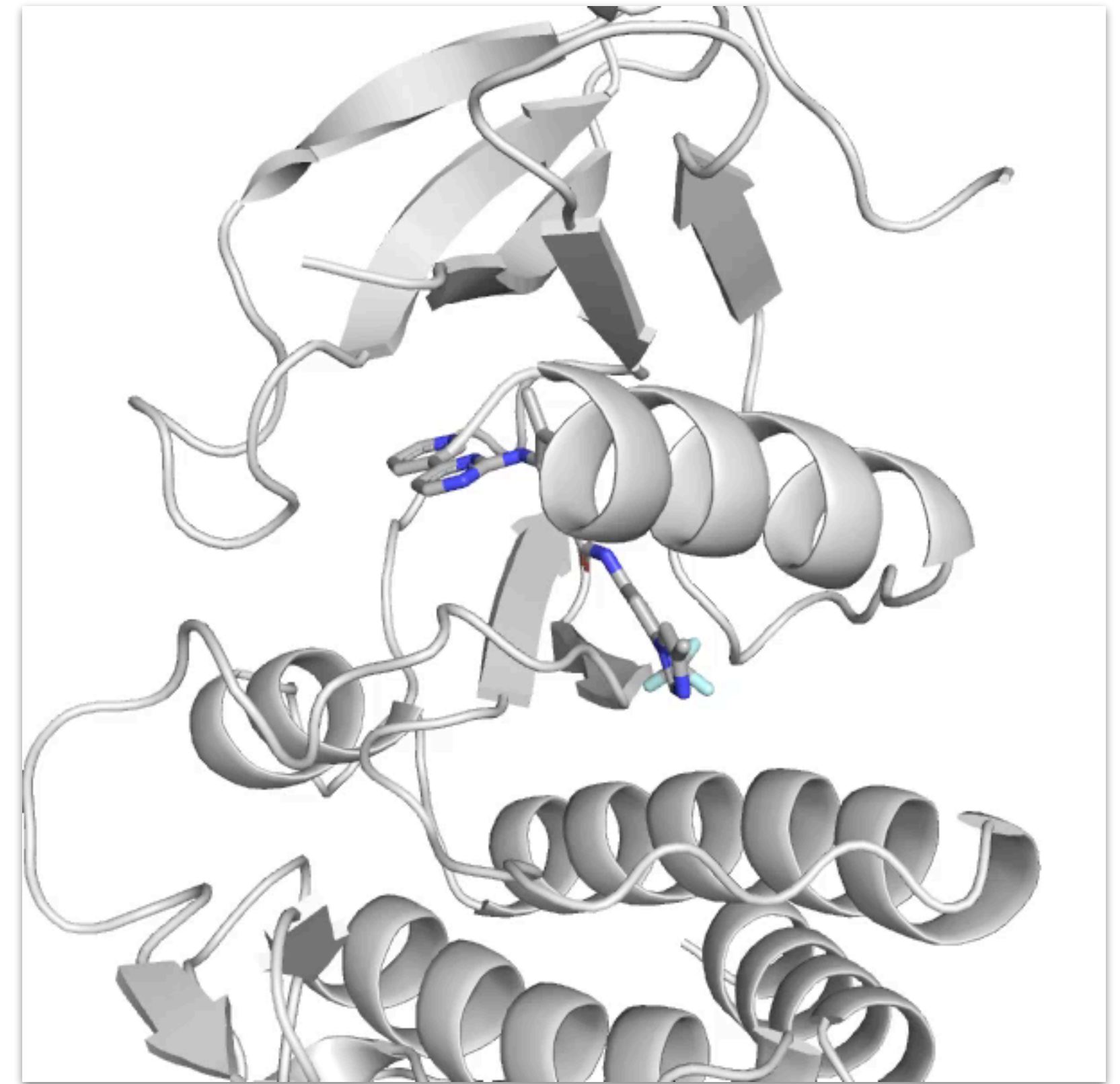
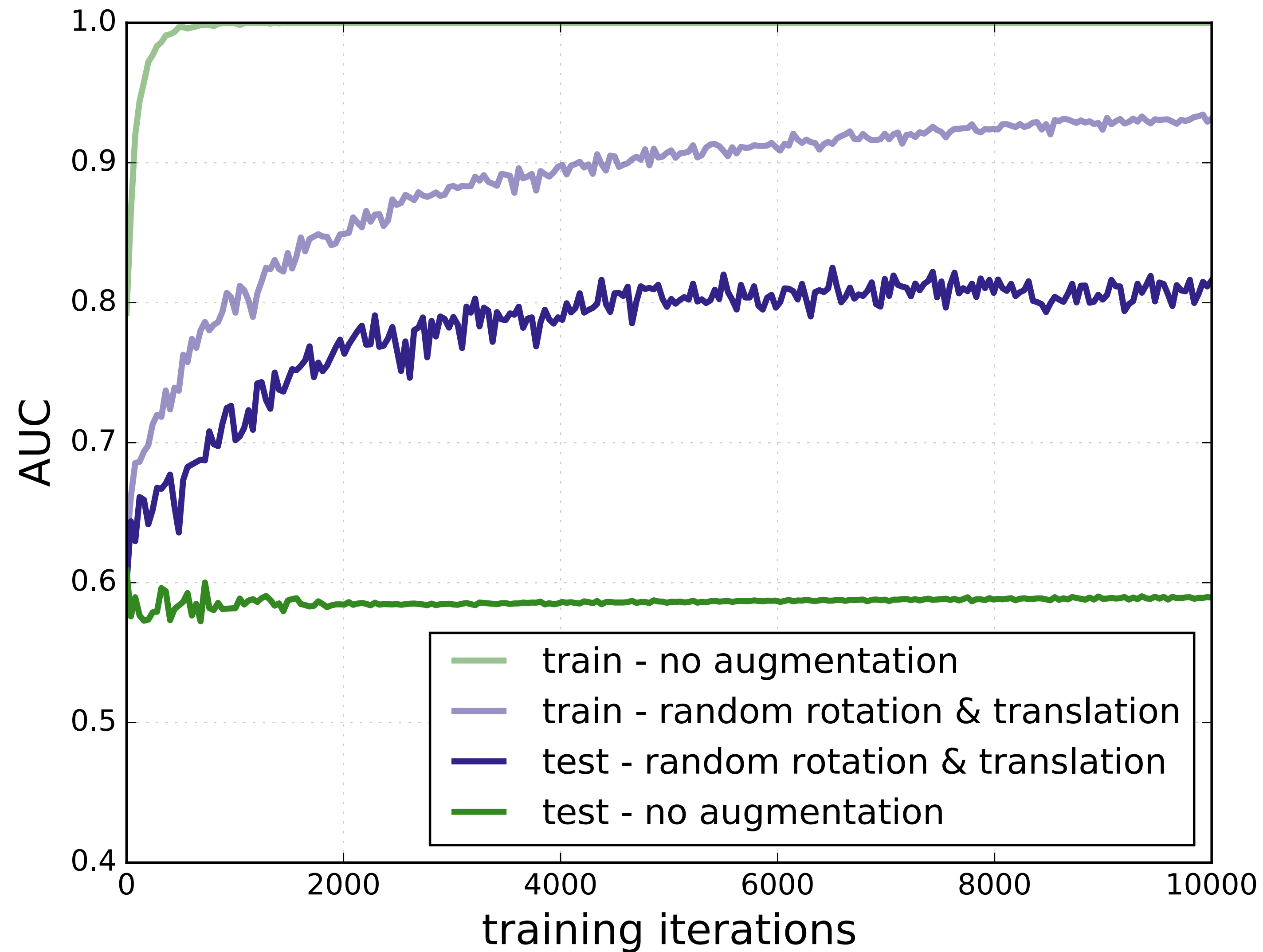
Parallelize over *atoms* to obtain a mask of atoms that overlap each grid region

Use exclusive scan to obtain a list of atom indices from the mask

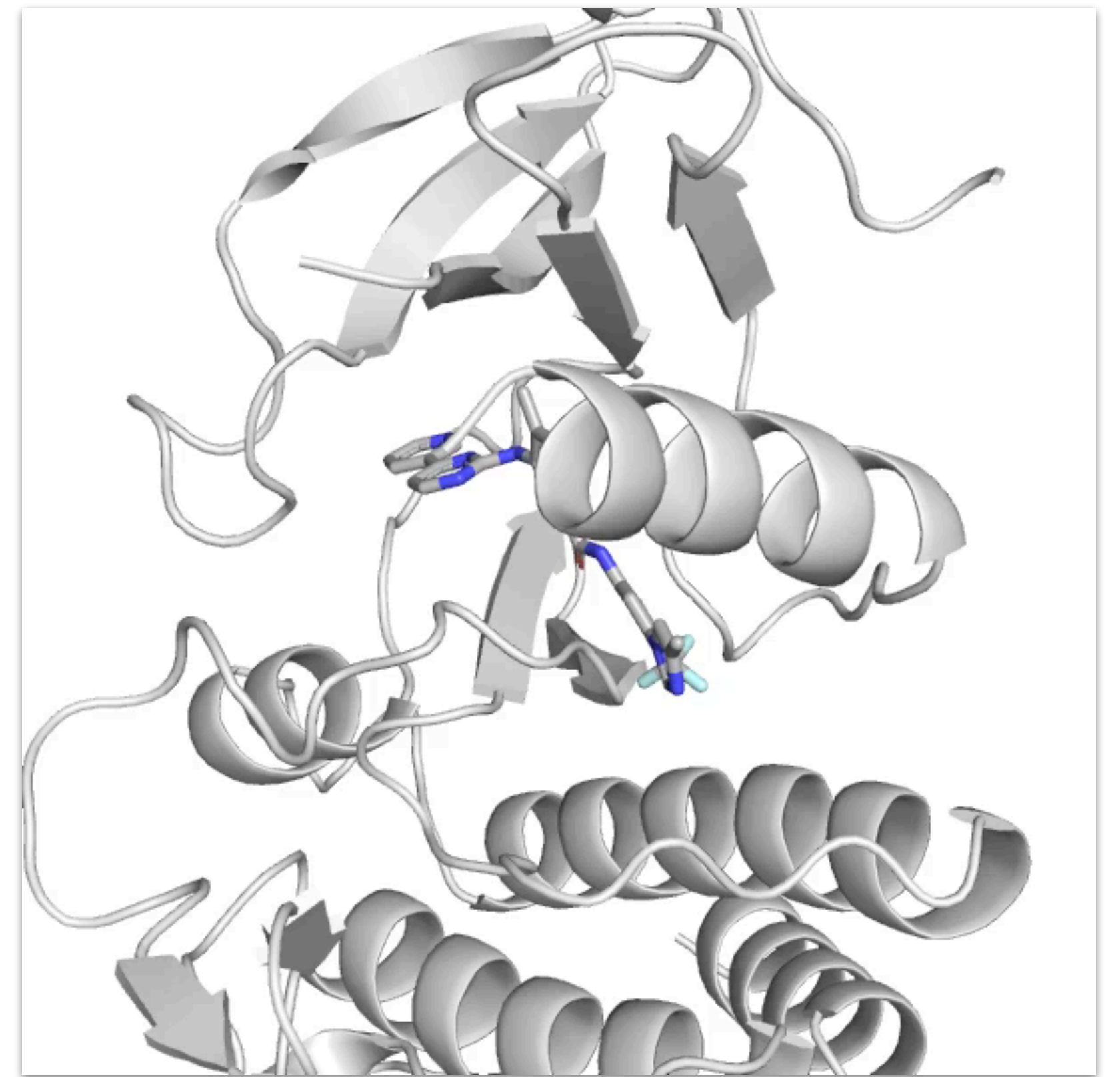
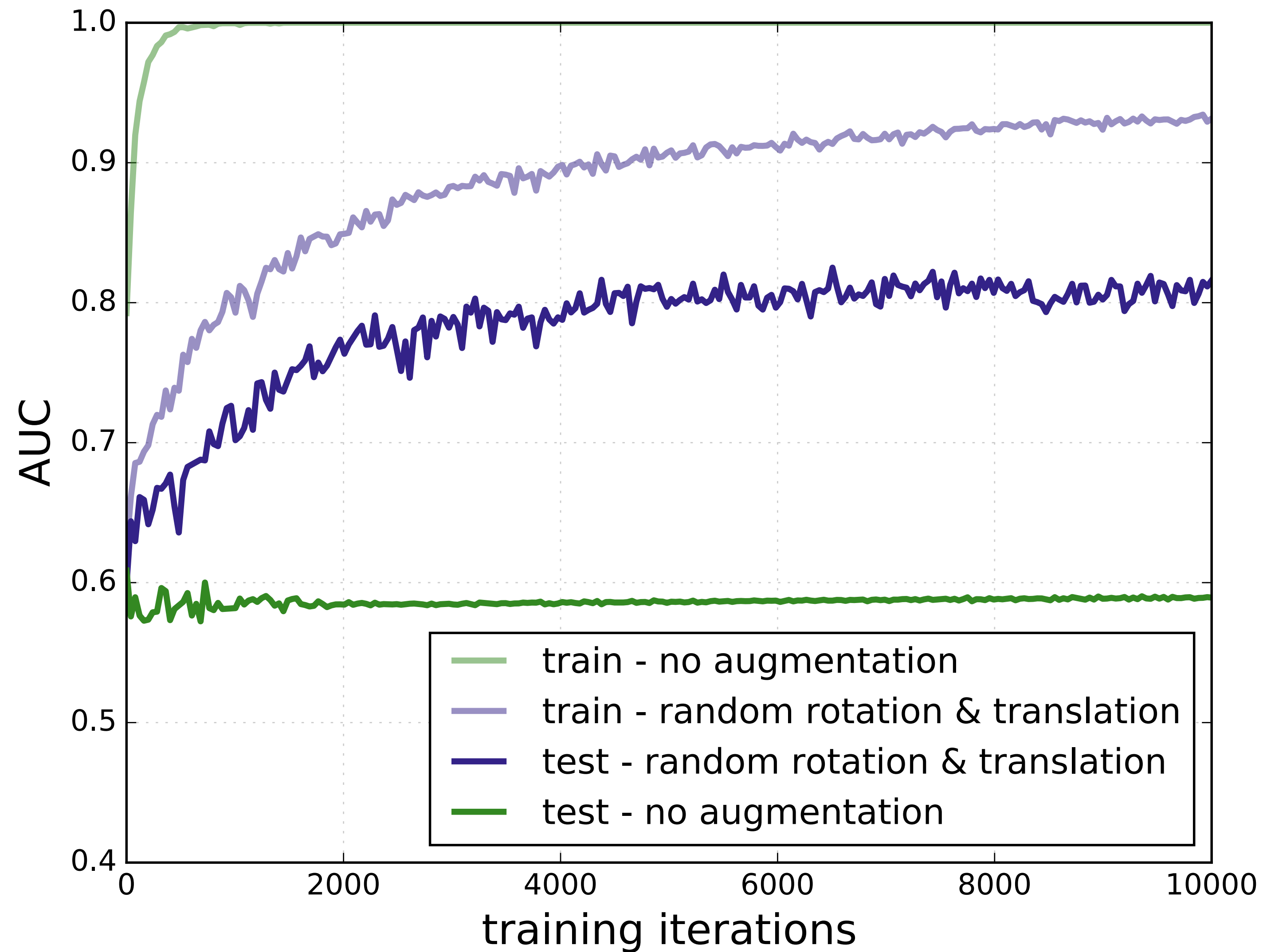
Parallelize over *grid points*, using reduced atom list to avoid $O(N_{\text{atoms}})$ check



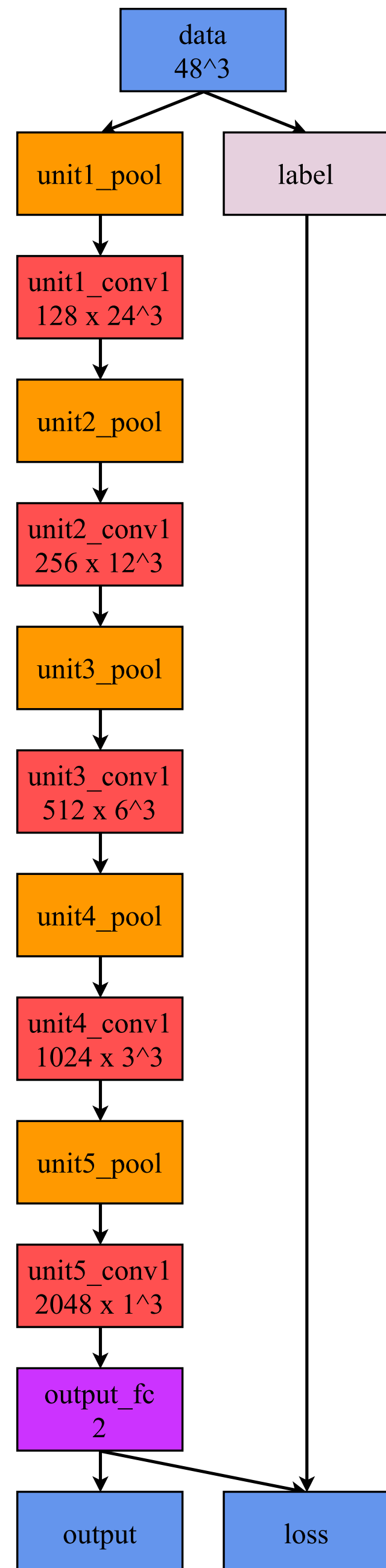
Data Augmentation



Data Augmentation



Model Optimization



Atom Types

- Vina (34)
- element-only (18)
- ligand-protein (2)

Atom Density Type

- Boolean
- Gaussian

Radius Multiple

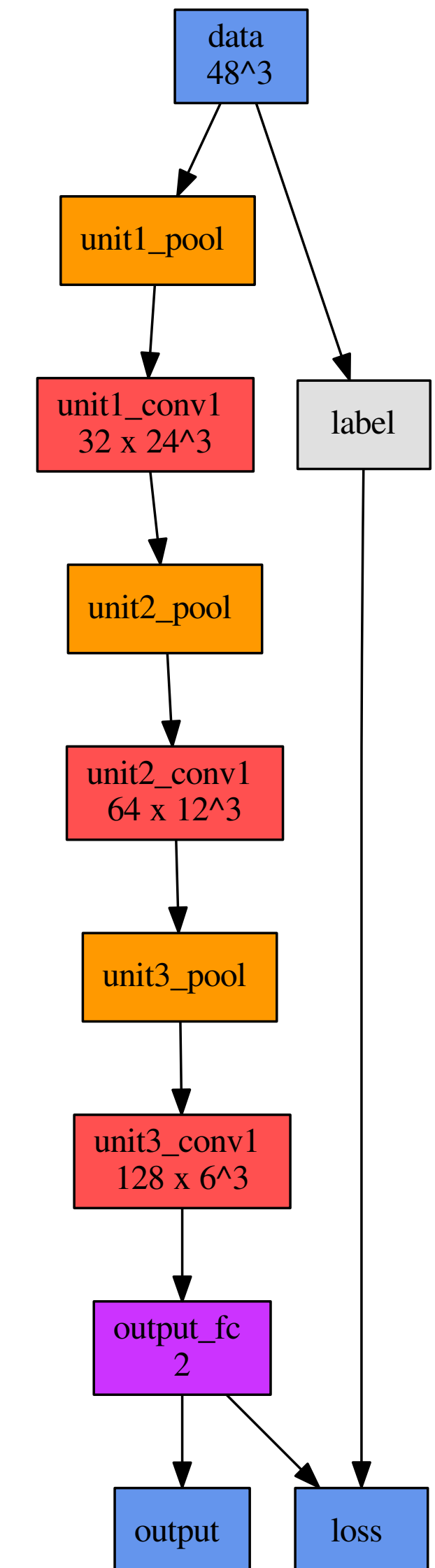
Resolution

Pooling

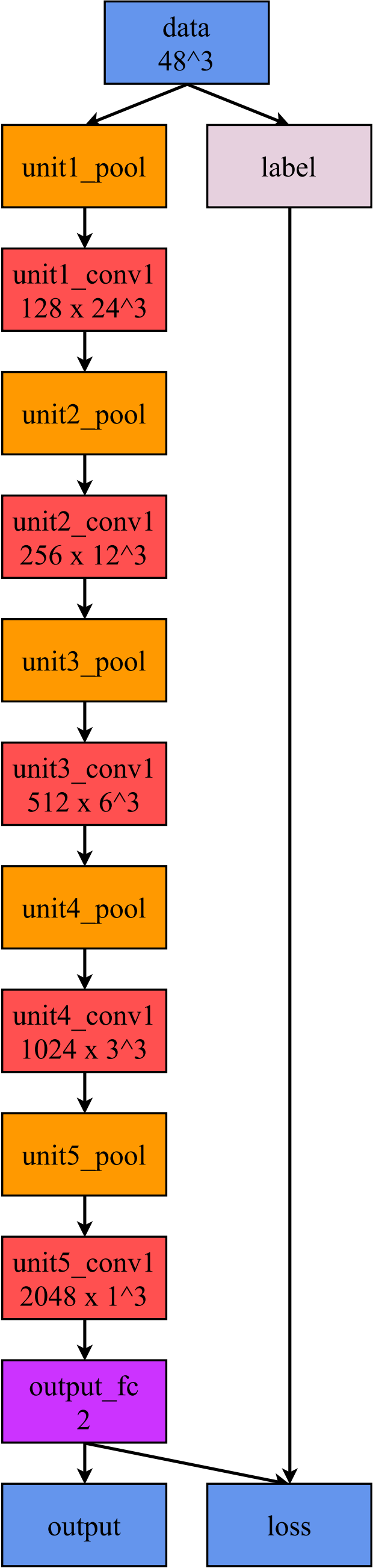
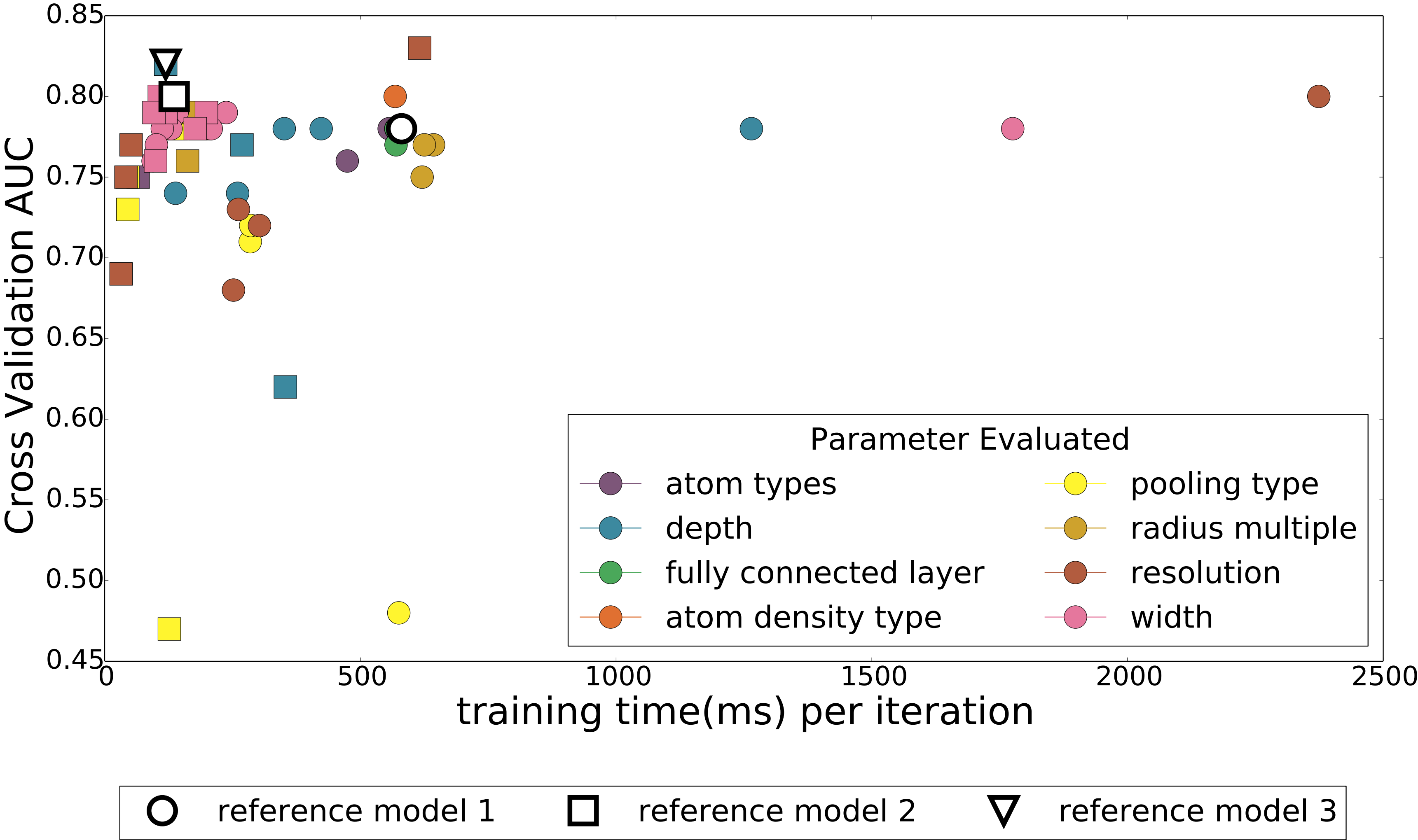
Depth

Width

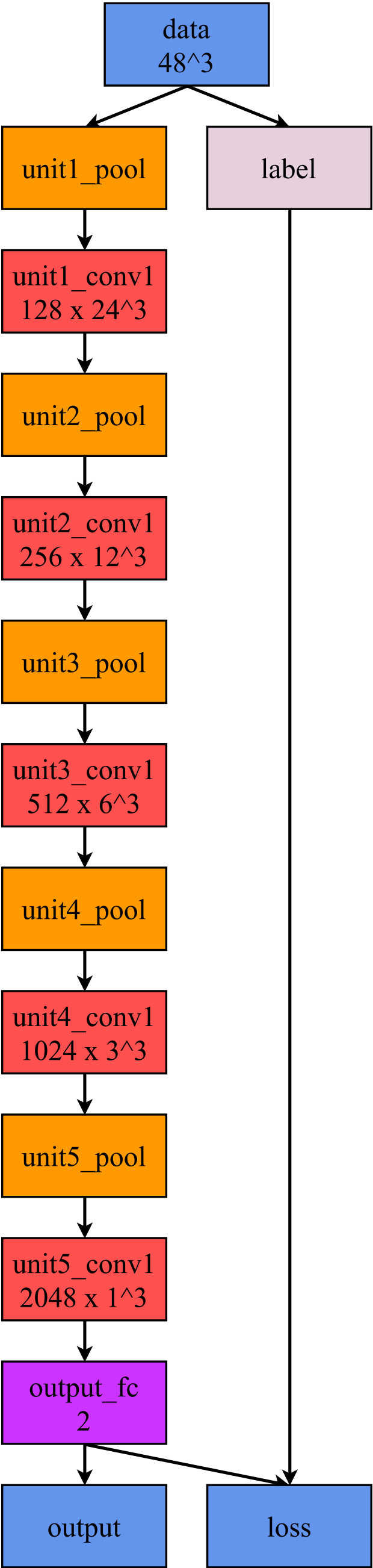
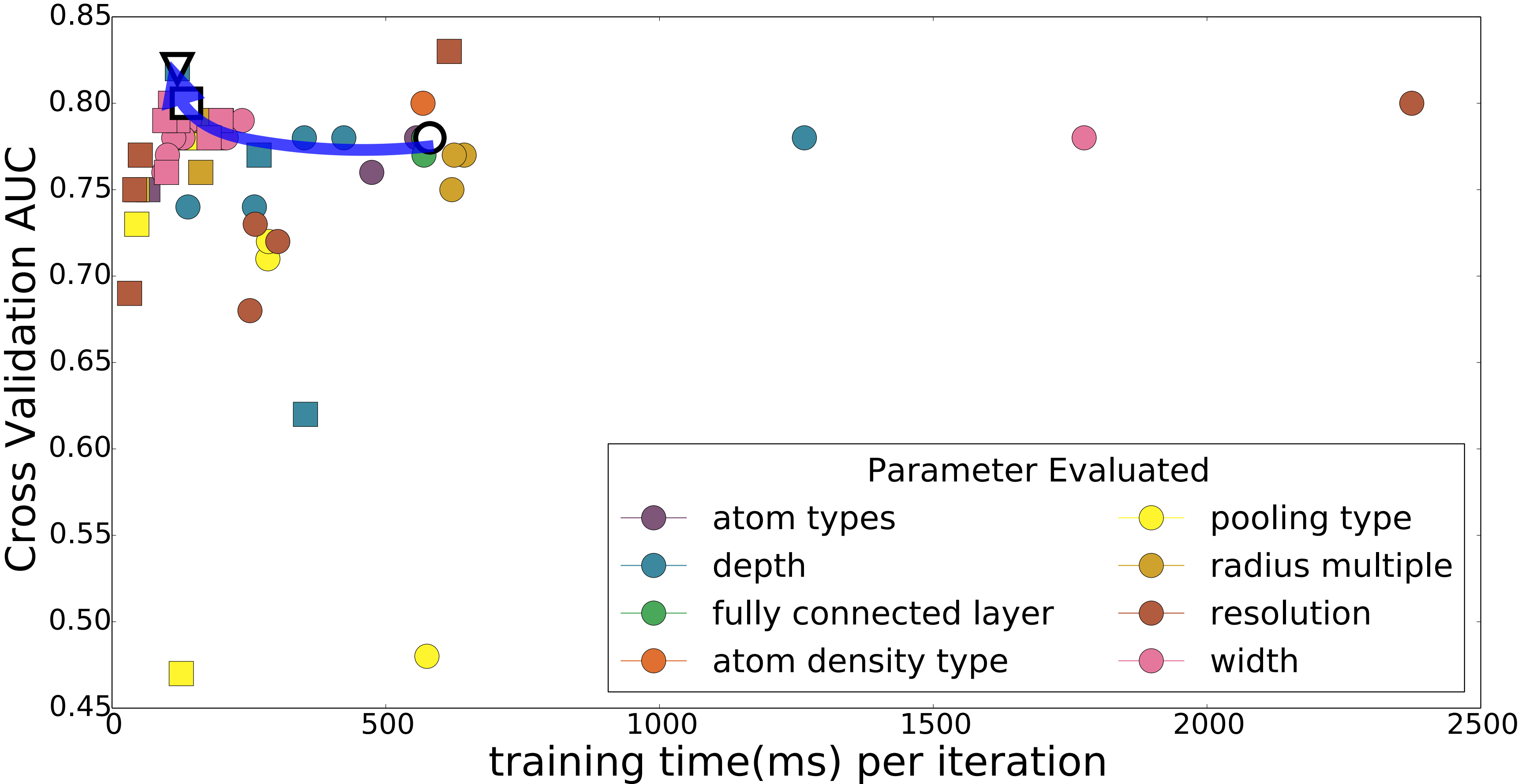
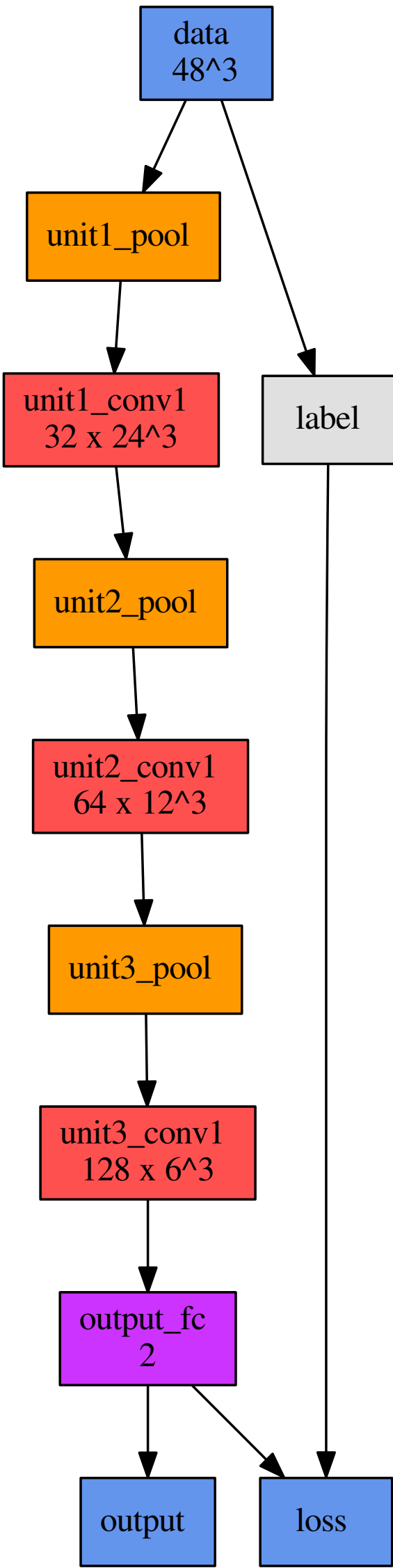
Fully Connected Layers



Model Optimization

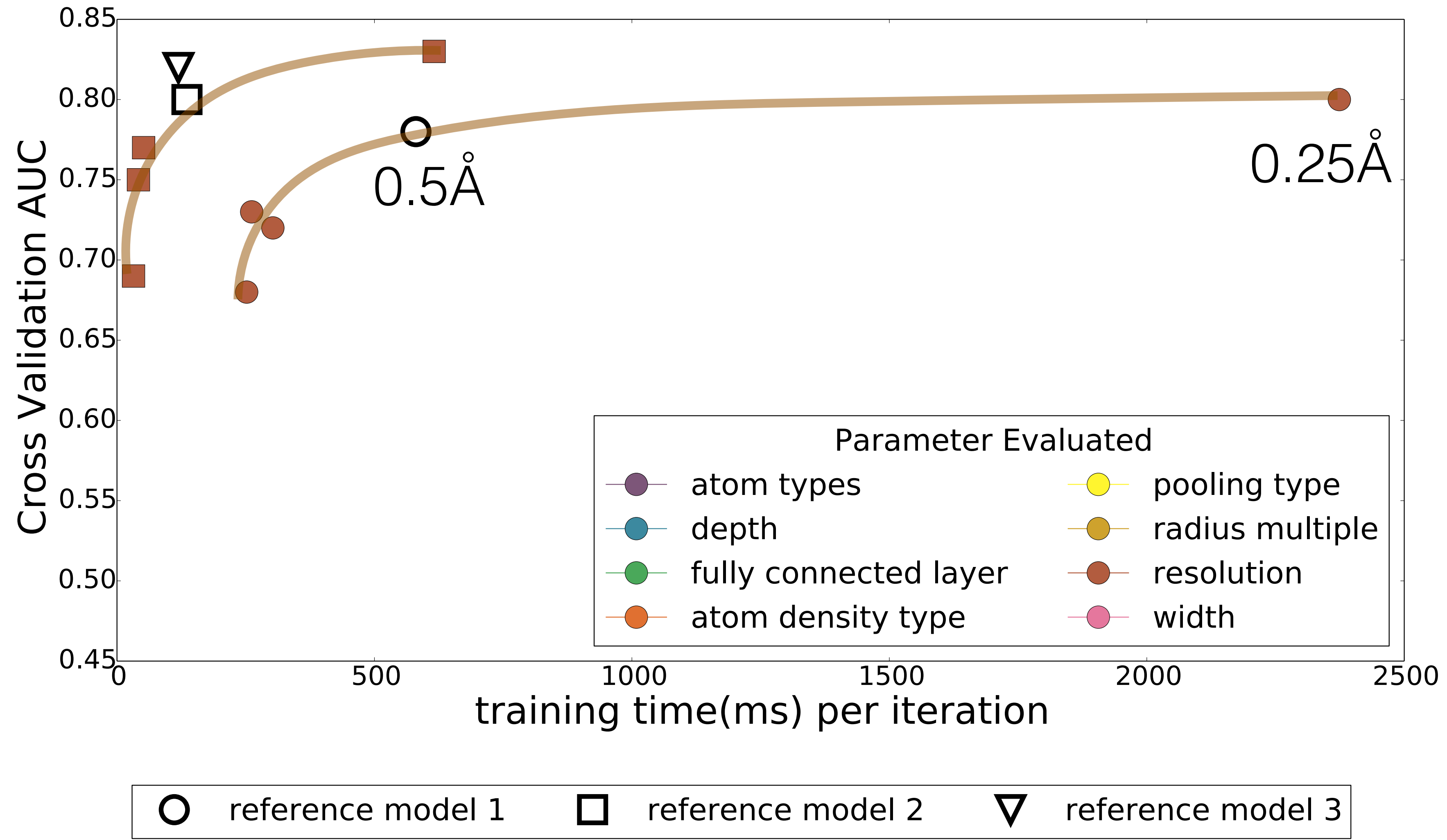


Model Optimization

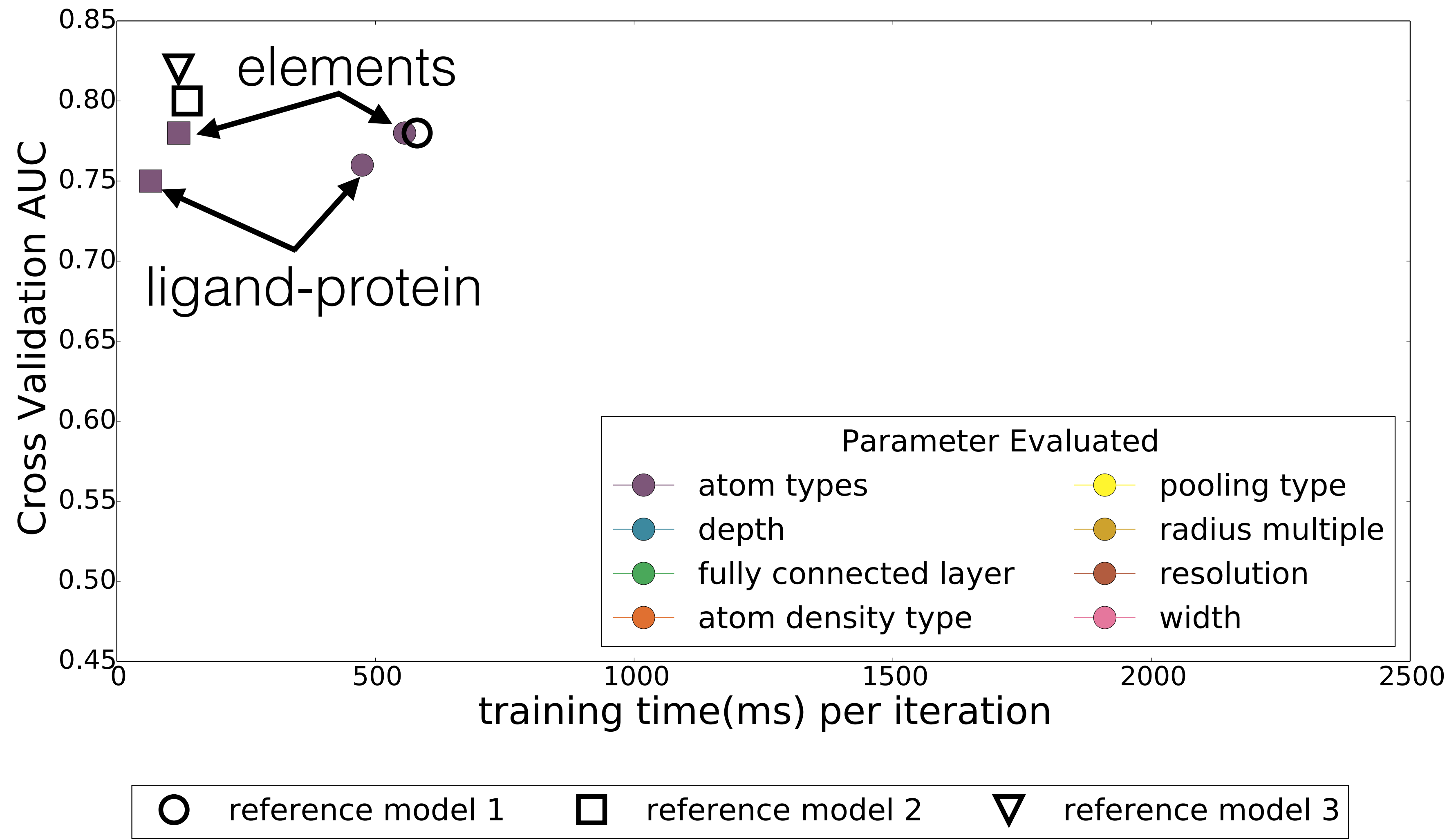


○ reference model 1 □ reference model 2 ▽ reference model 3

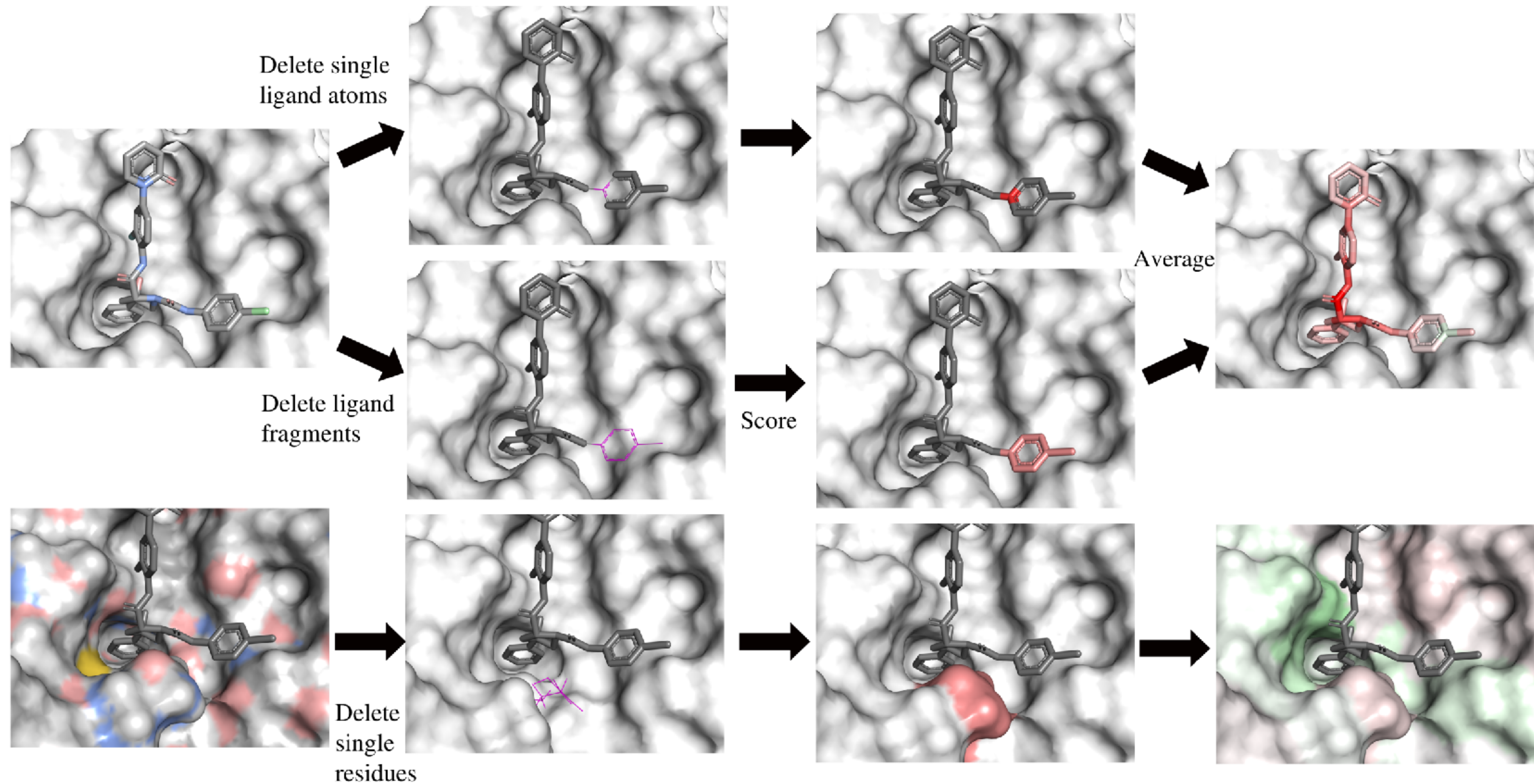
Grid Resolution



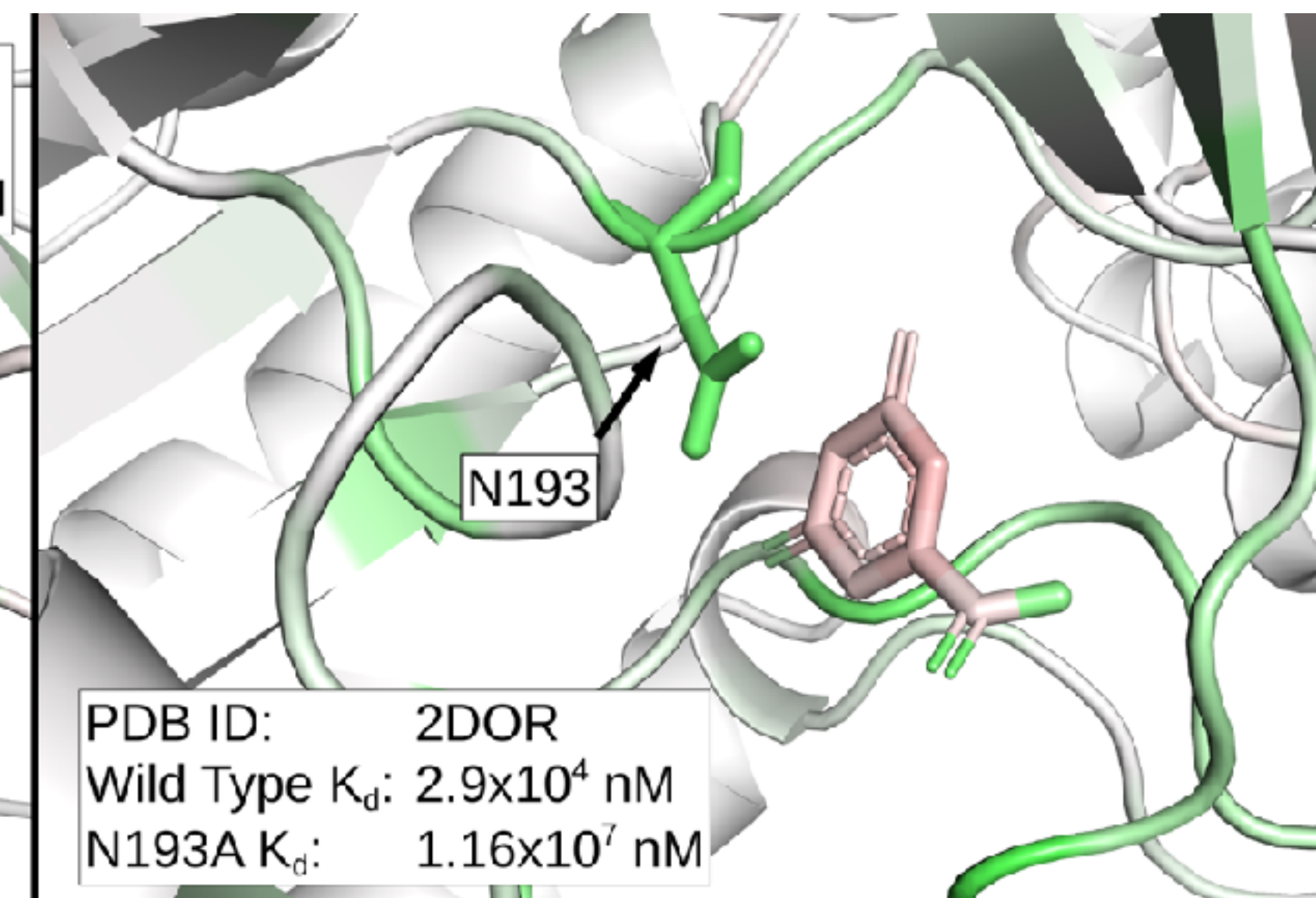
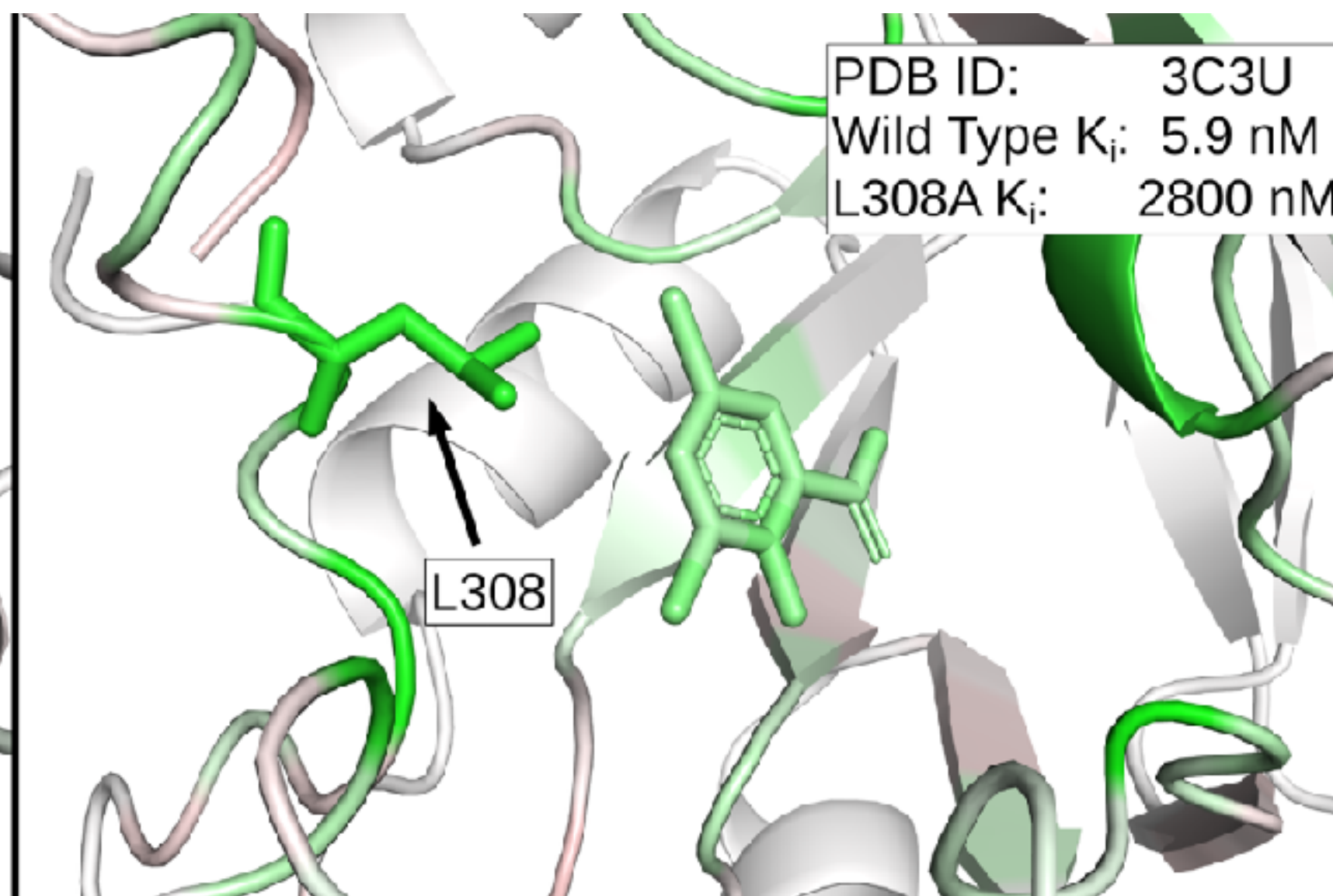
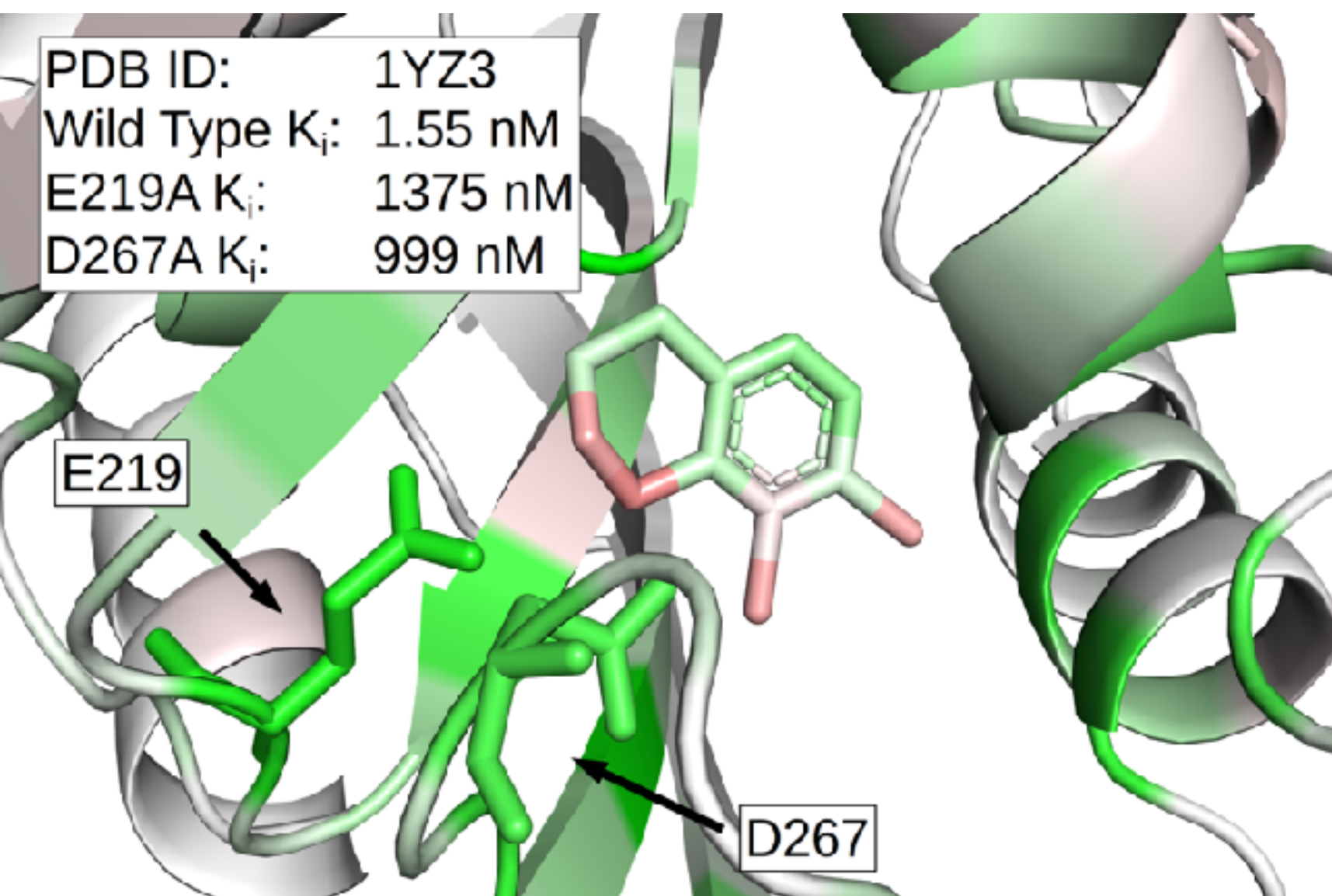
Atom Types



Visualization

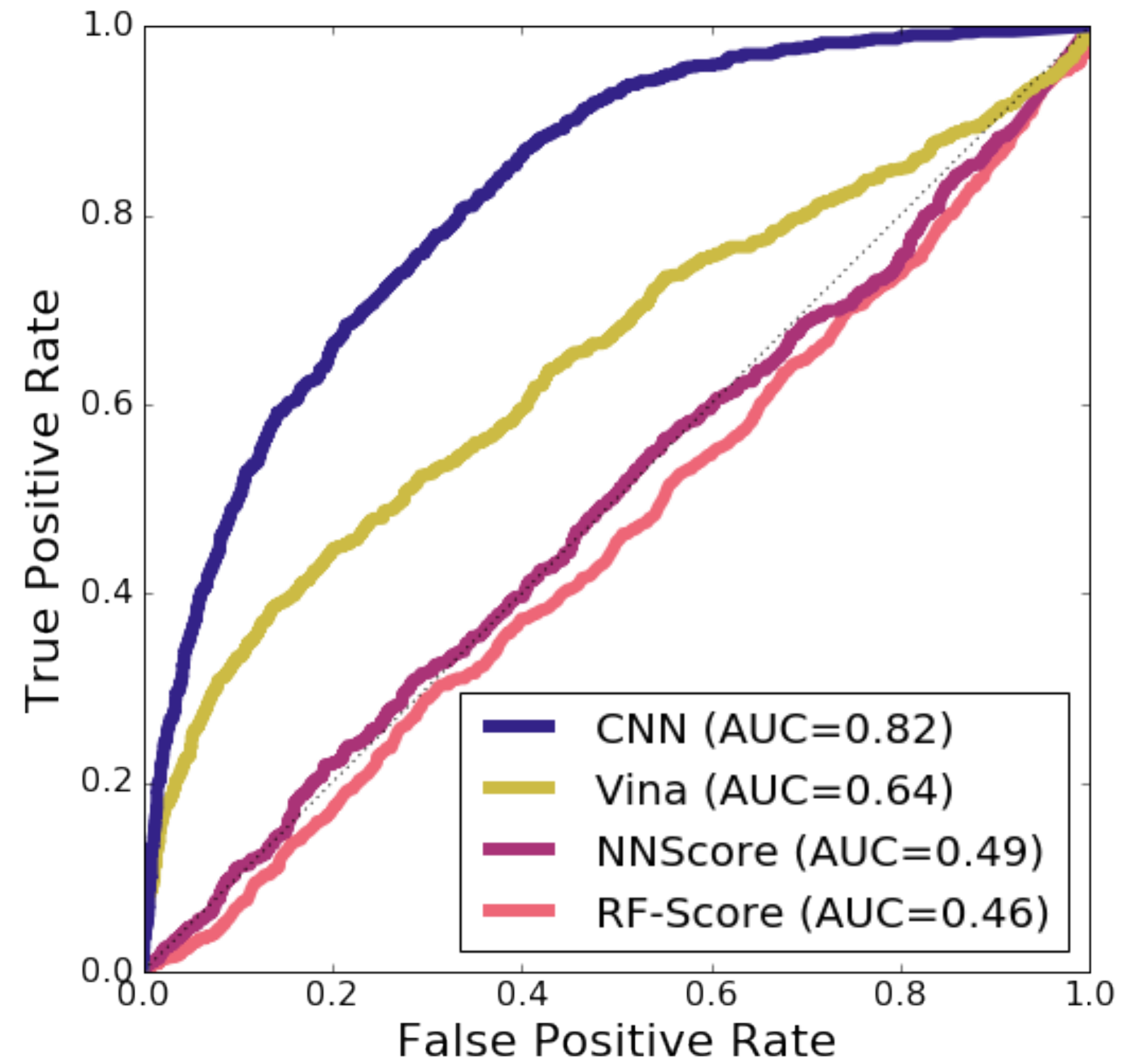


Visualizing Enzymes

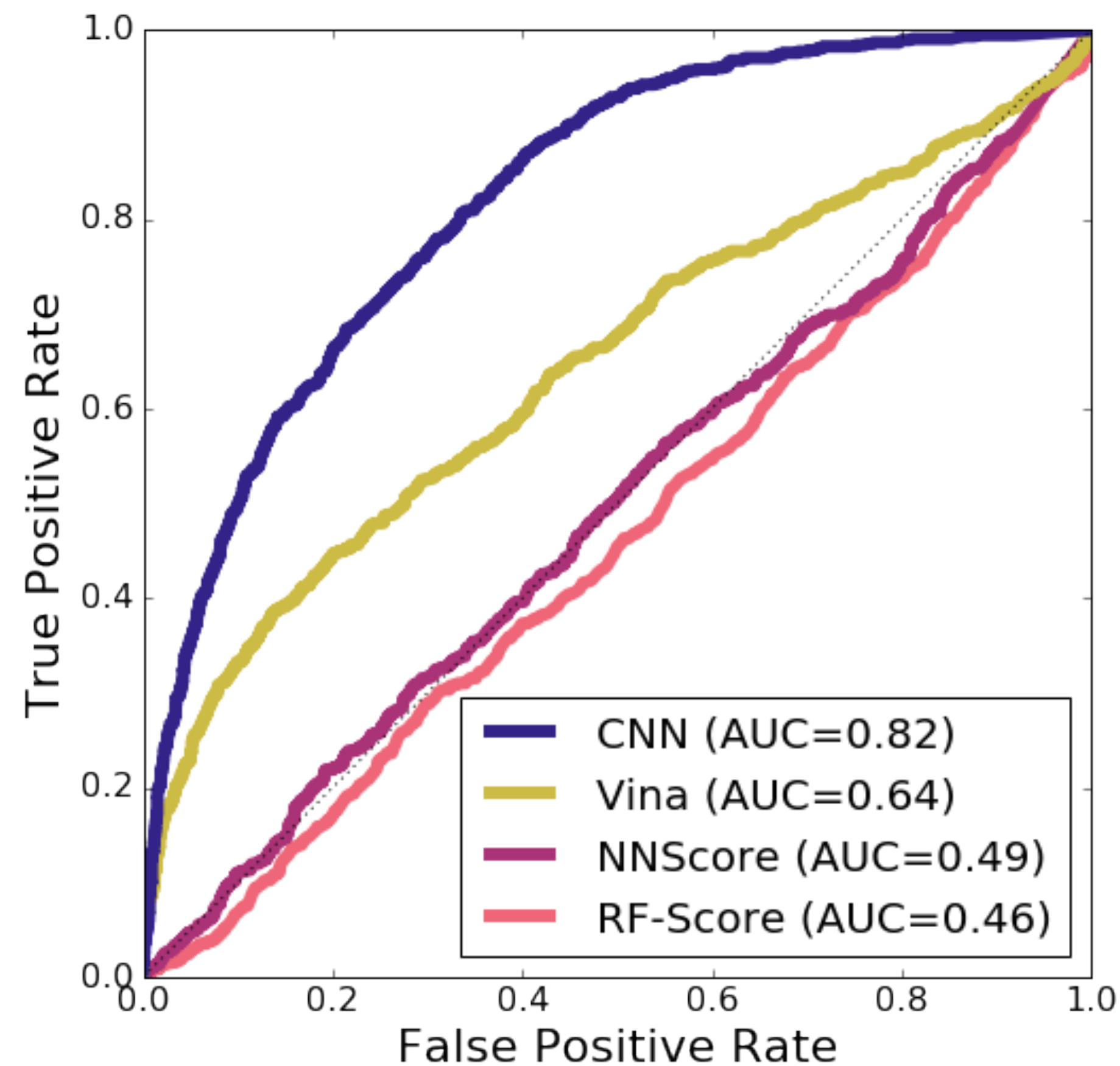


Cross-Validation Evaluation

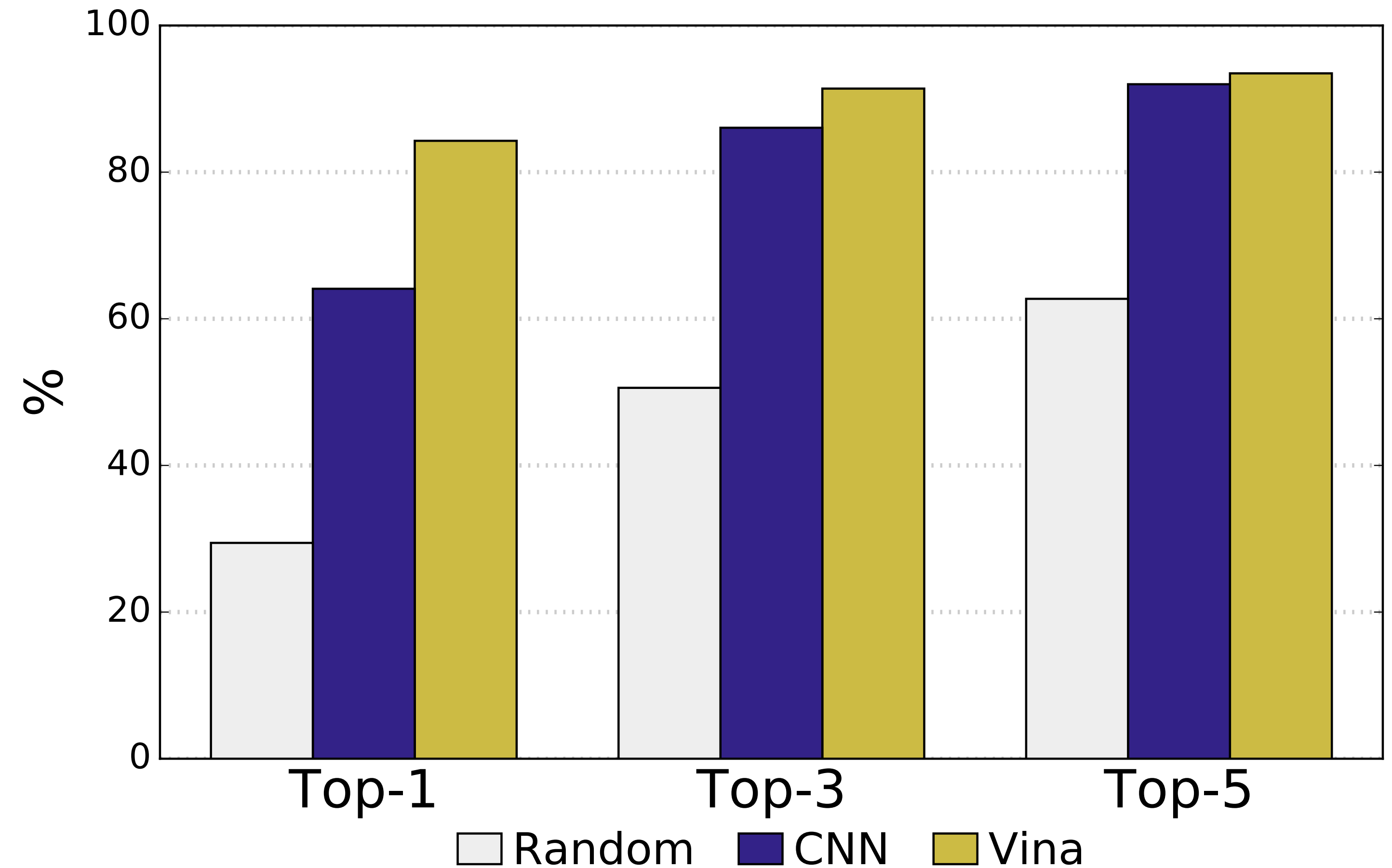
Pose Prediction (CSAR)



Pose Prediction (CSAR)

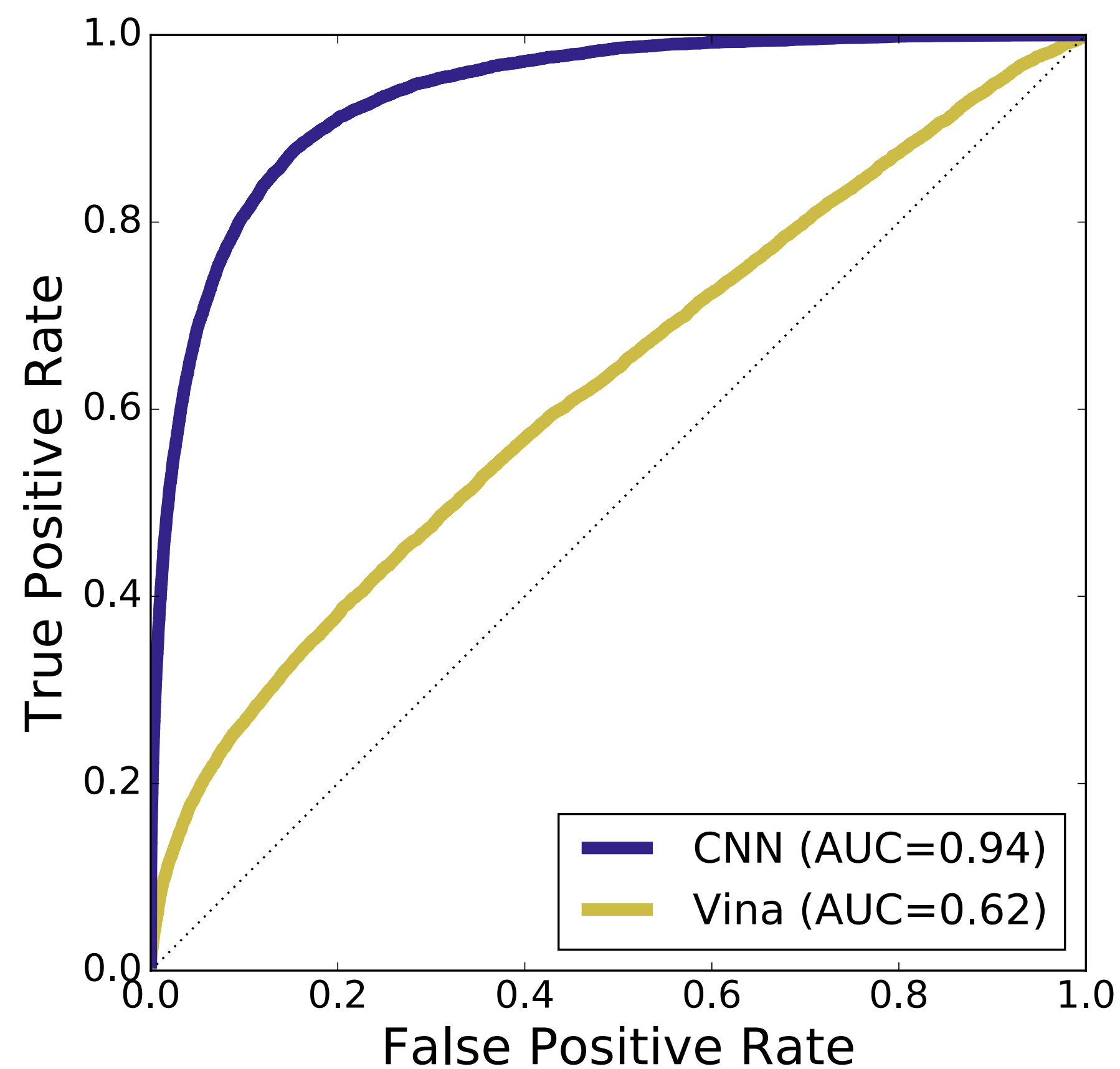


inter-target ranking

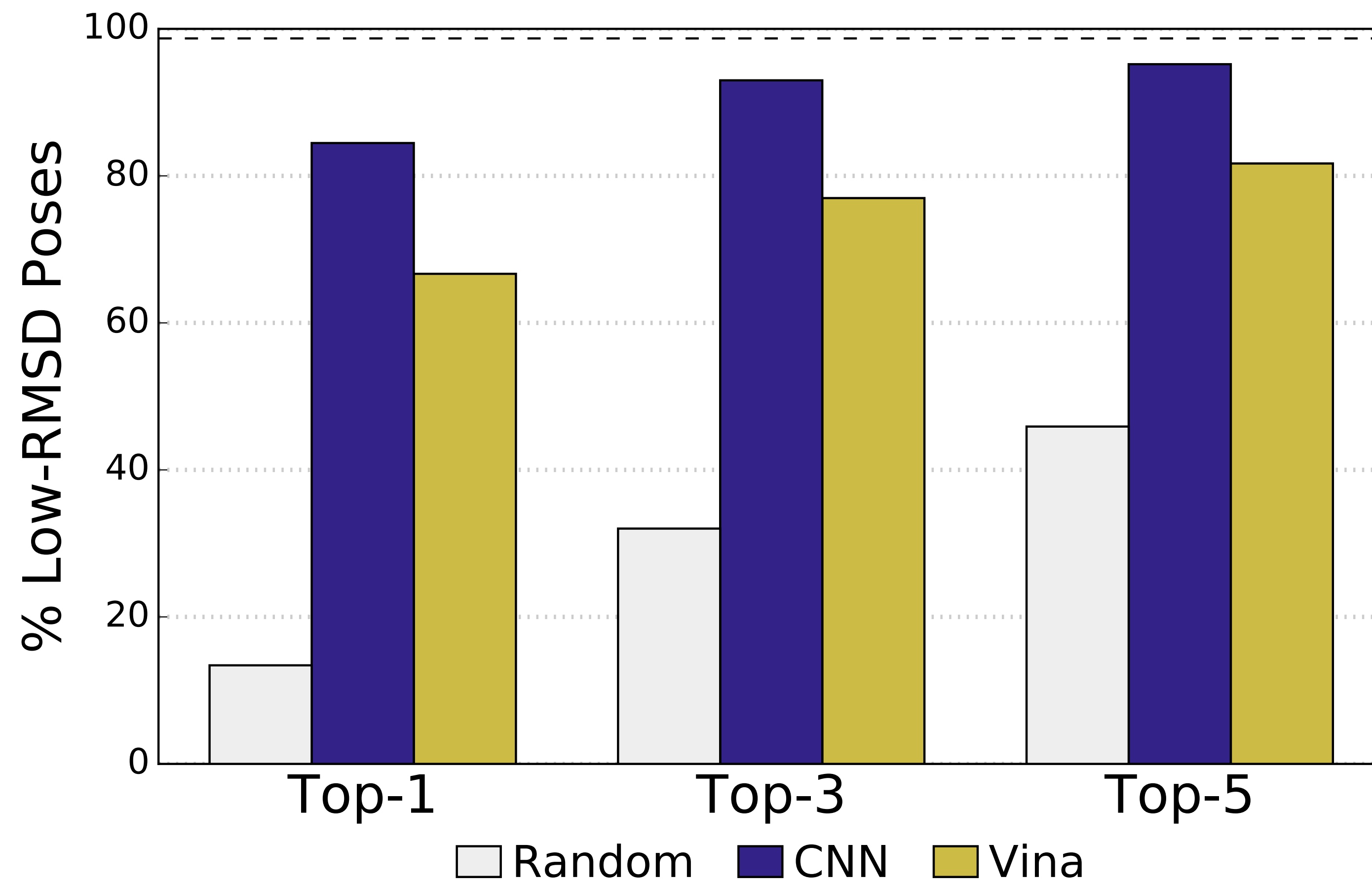


intra-target ranking

Pose Prediction (PDBbind)

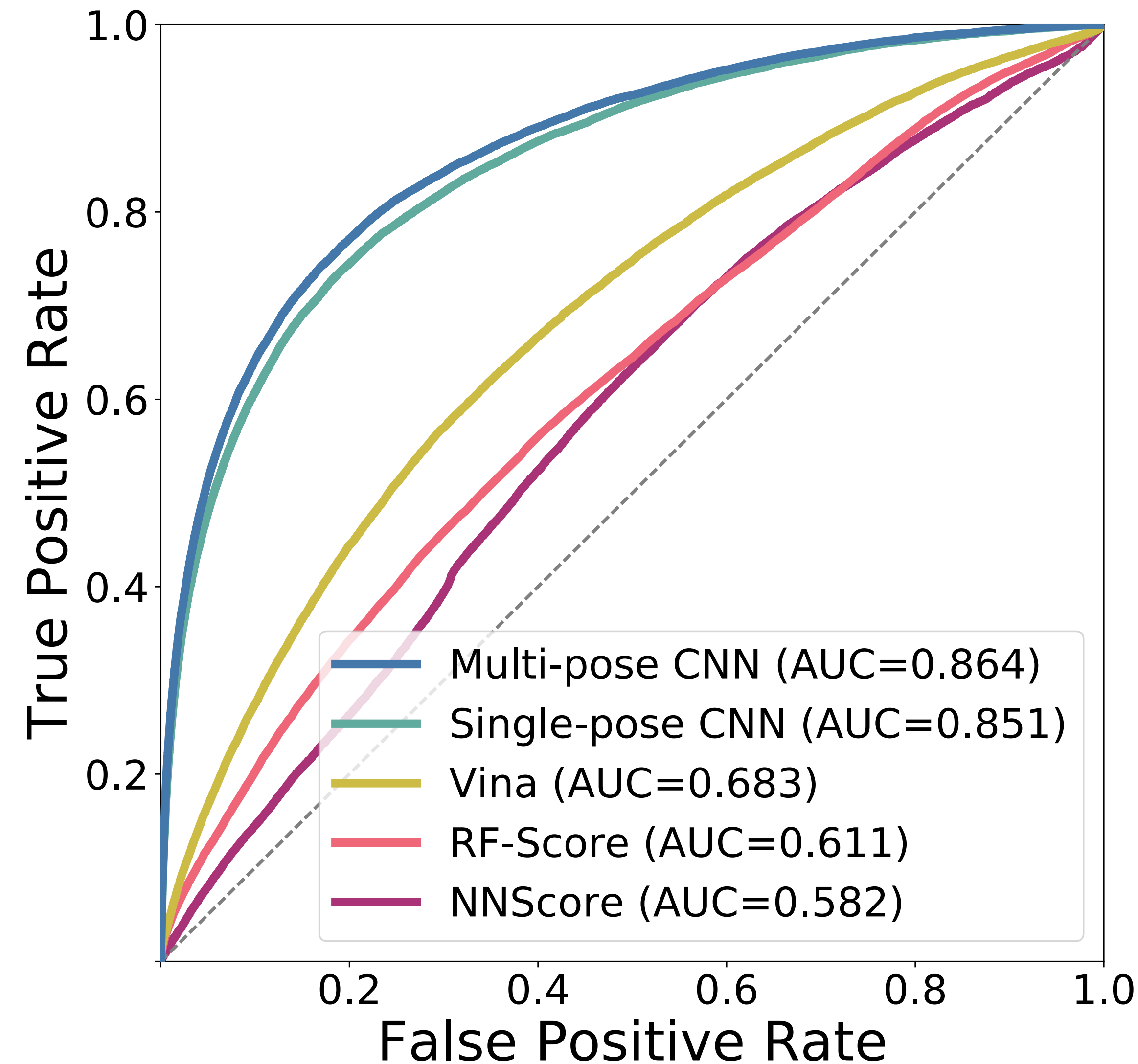


inter-target ranking



intra-target ranking

Binding Determination

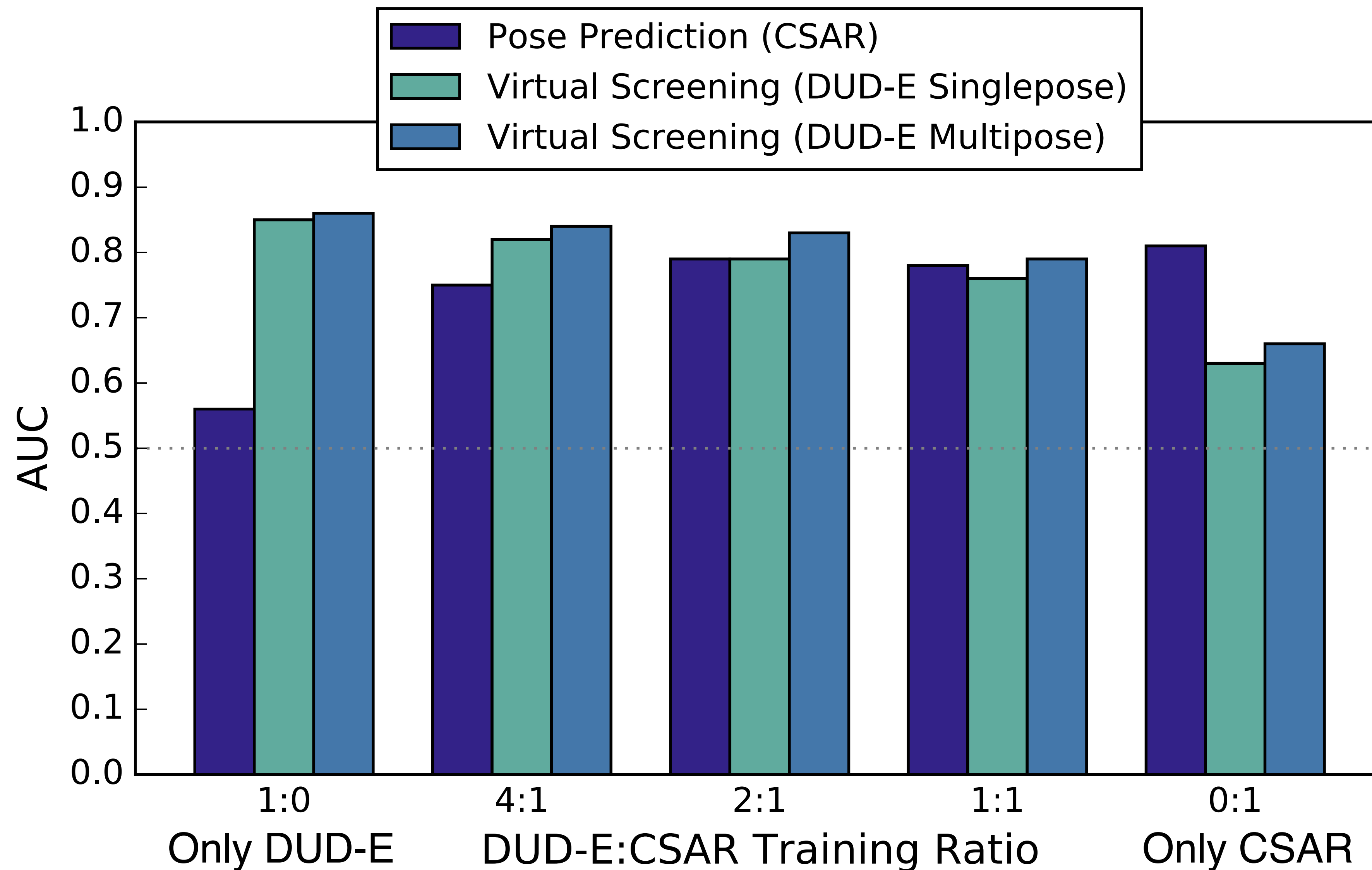


D U D • E

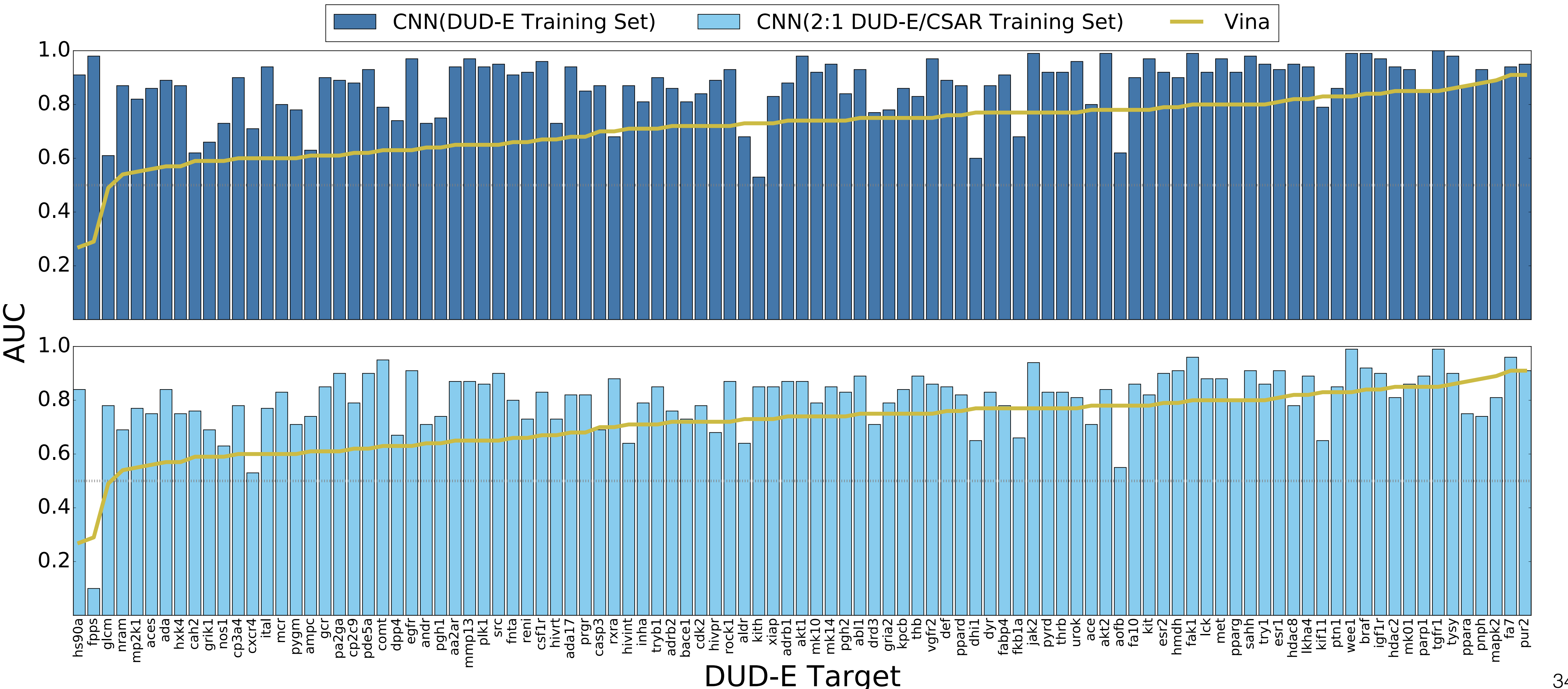
102 targets

- 22,645 actives
- 1,407,145 decoys
- <10 μ M affinity
- **true poses unknown**
- use top docked pose

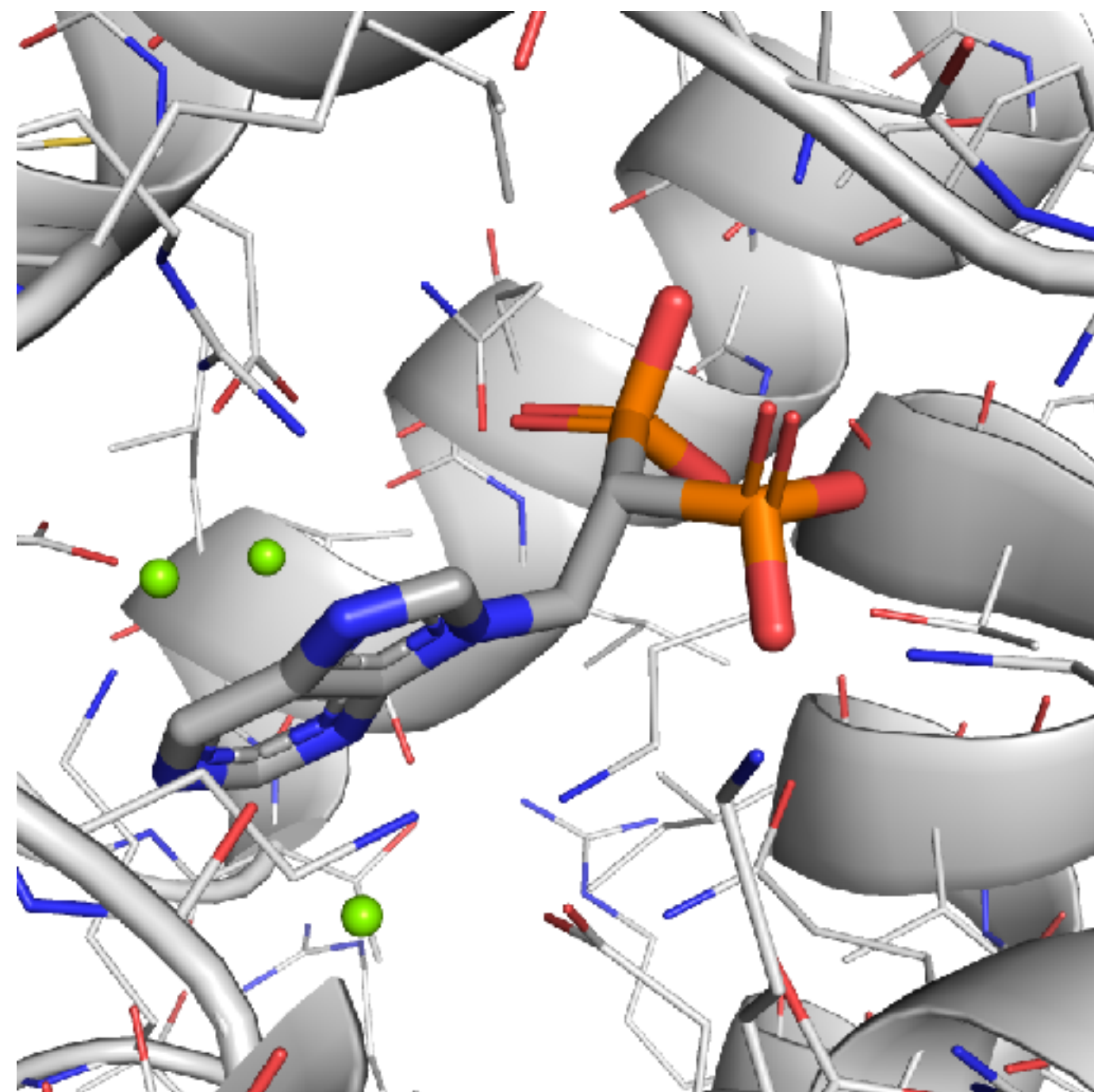
Combined Training



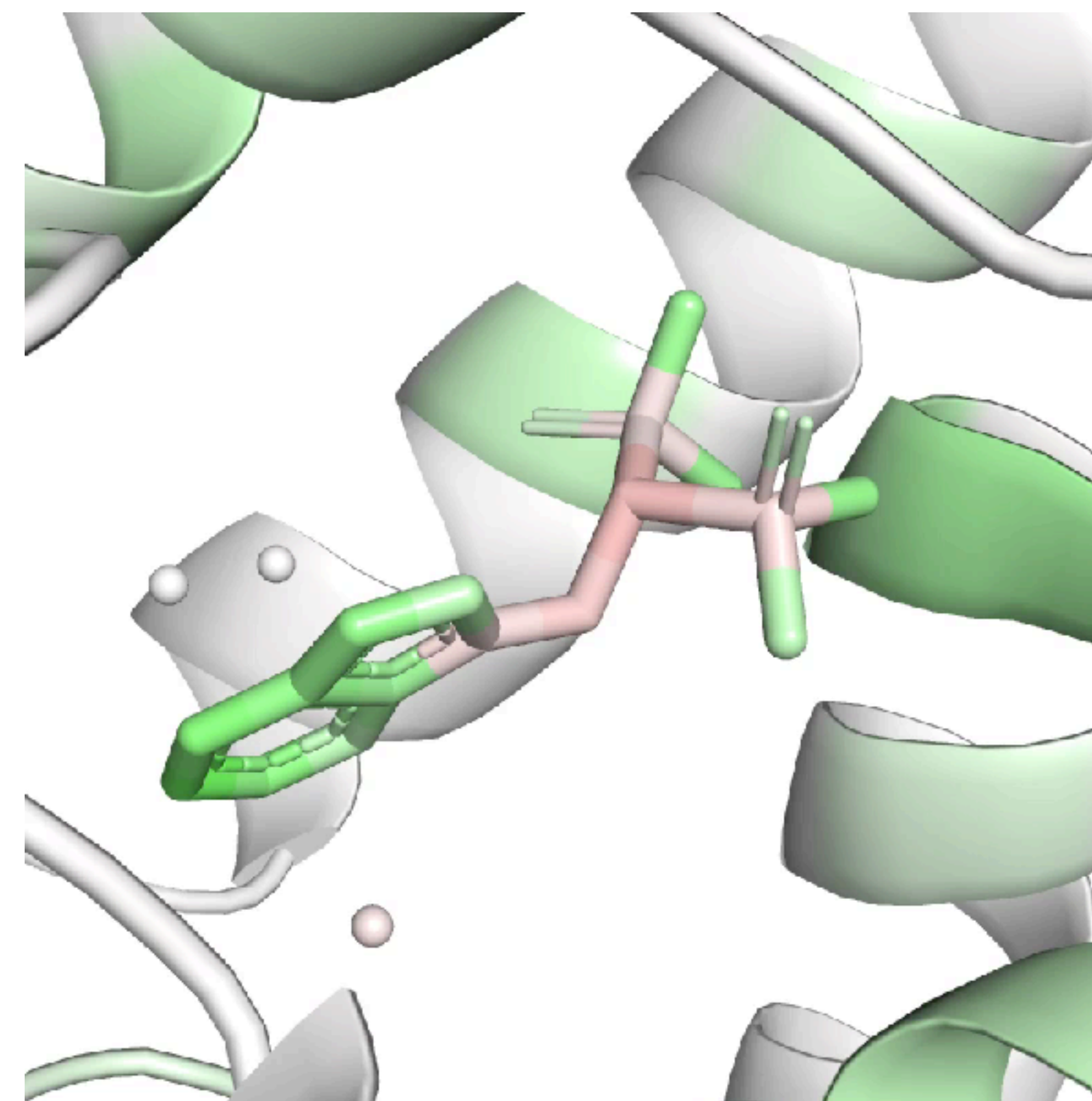
Binding Determination



fpps (farnesyl diphosphate synthase)

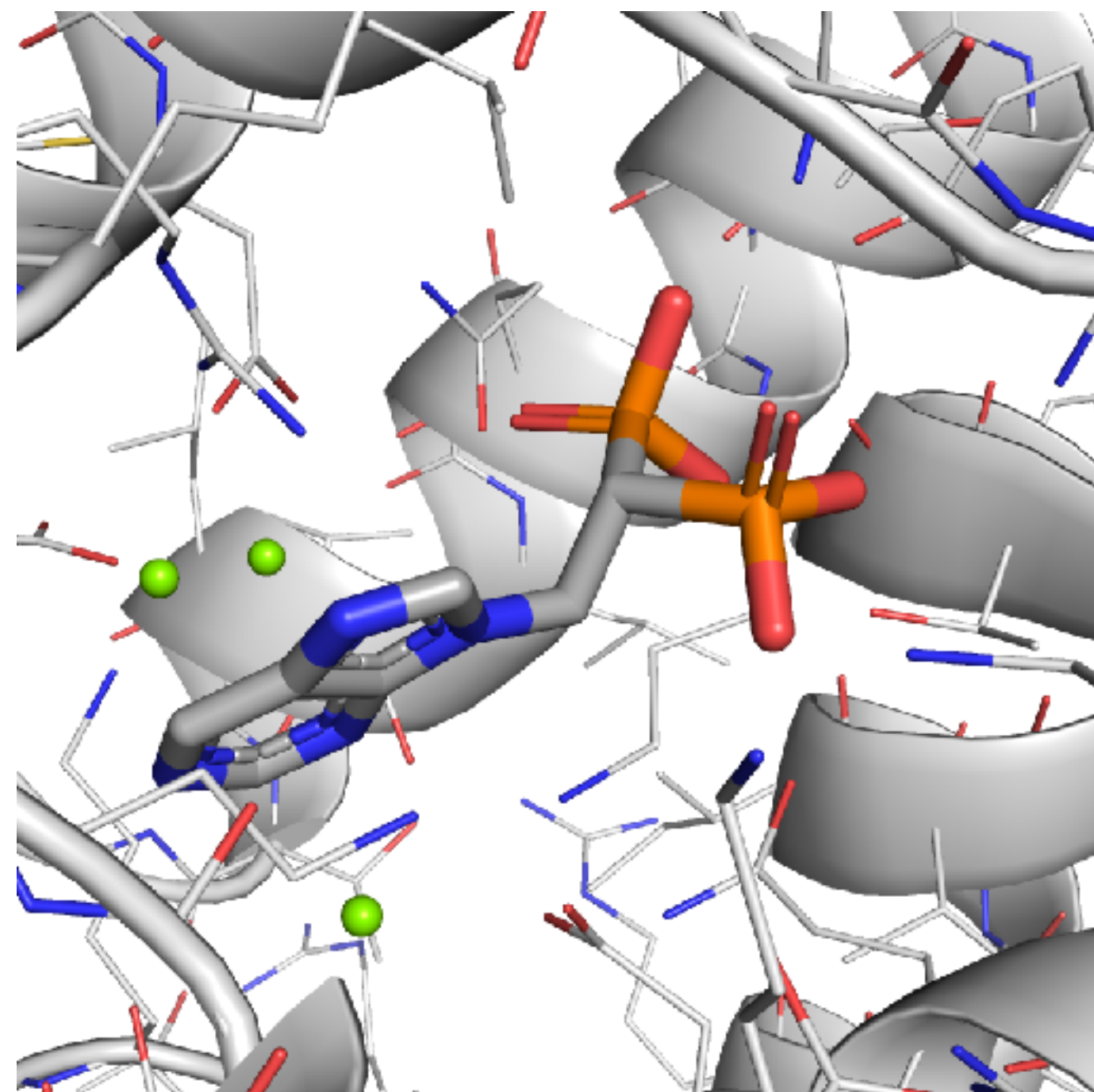


CHEMBL457424
Top Vina Pose (-8.2 kcal/mol)

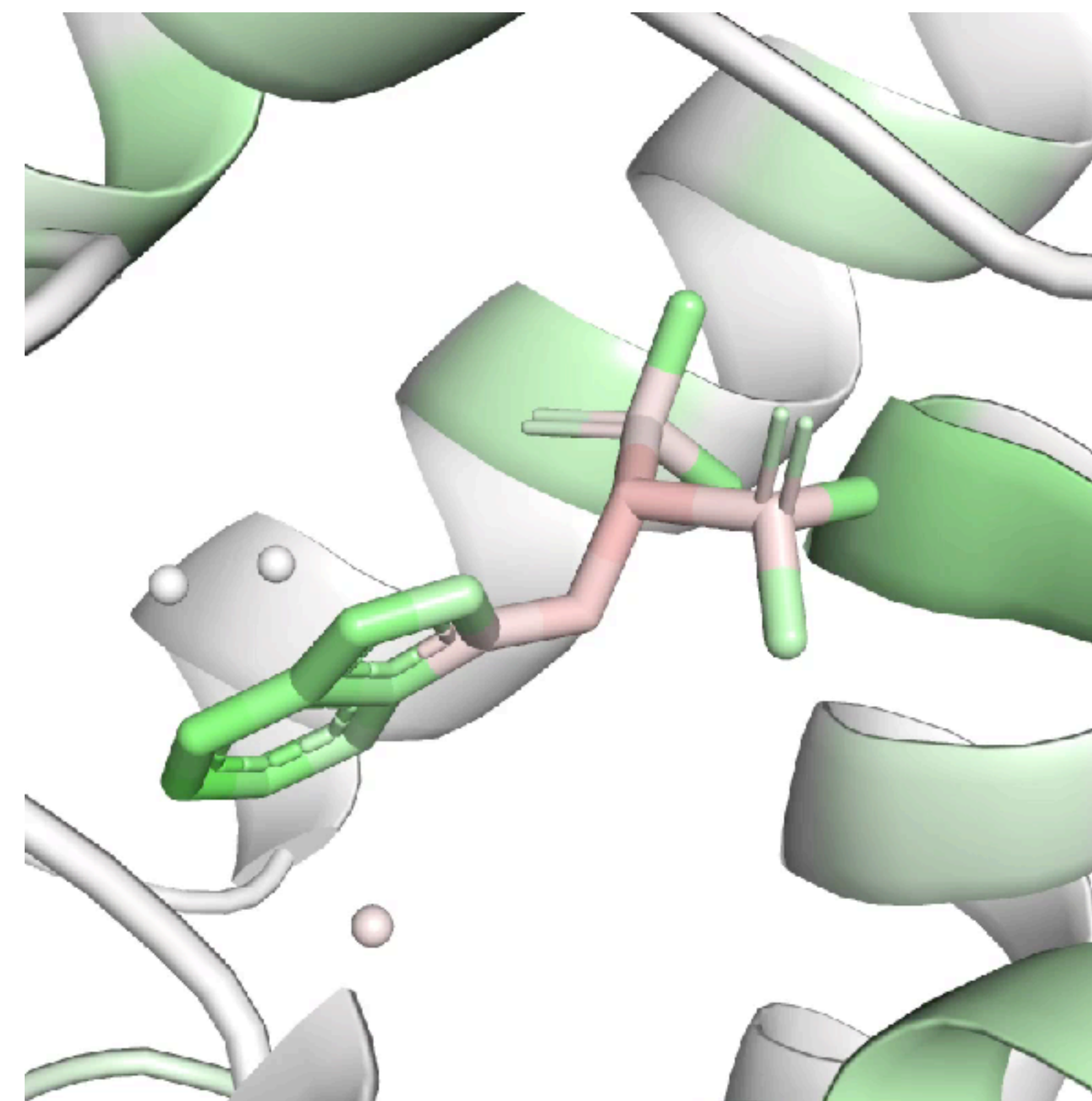


DUD-E Training Set
Score = 0.93 ± 0.03

fpps (farnesyl diphosphate synthase)

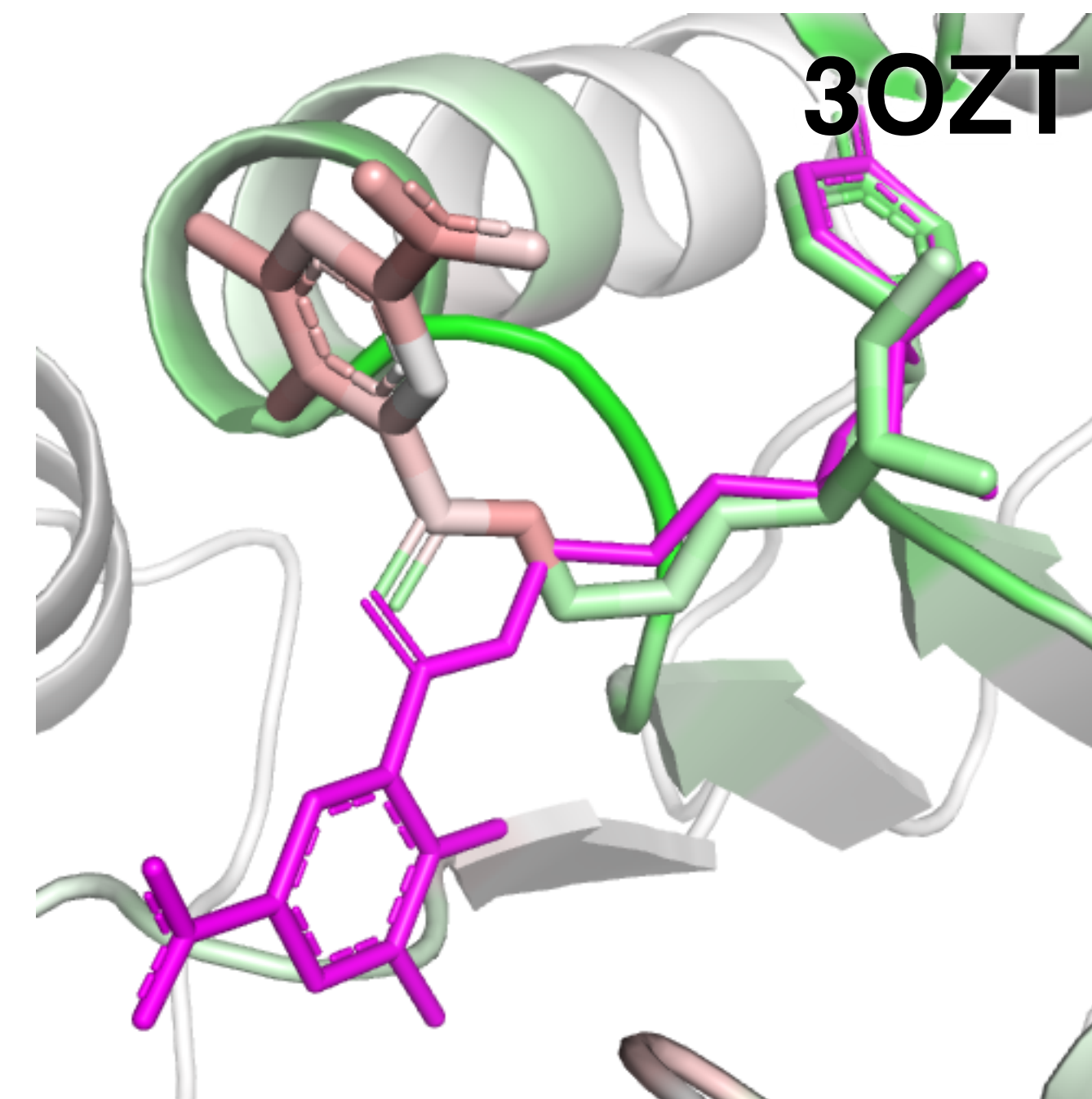
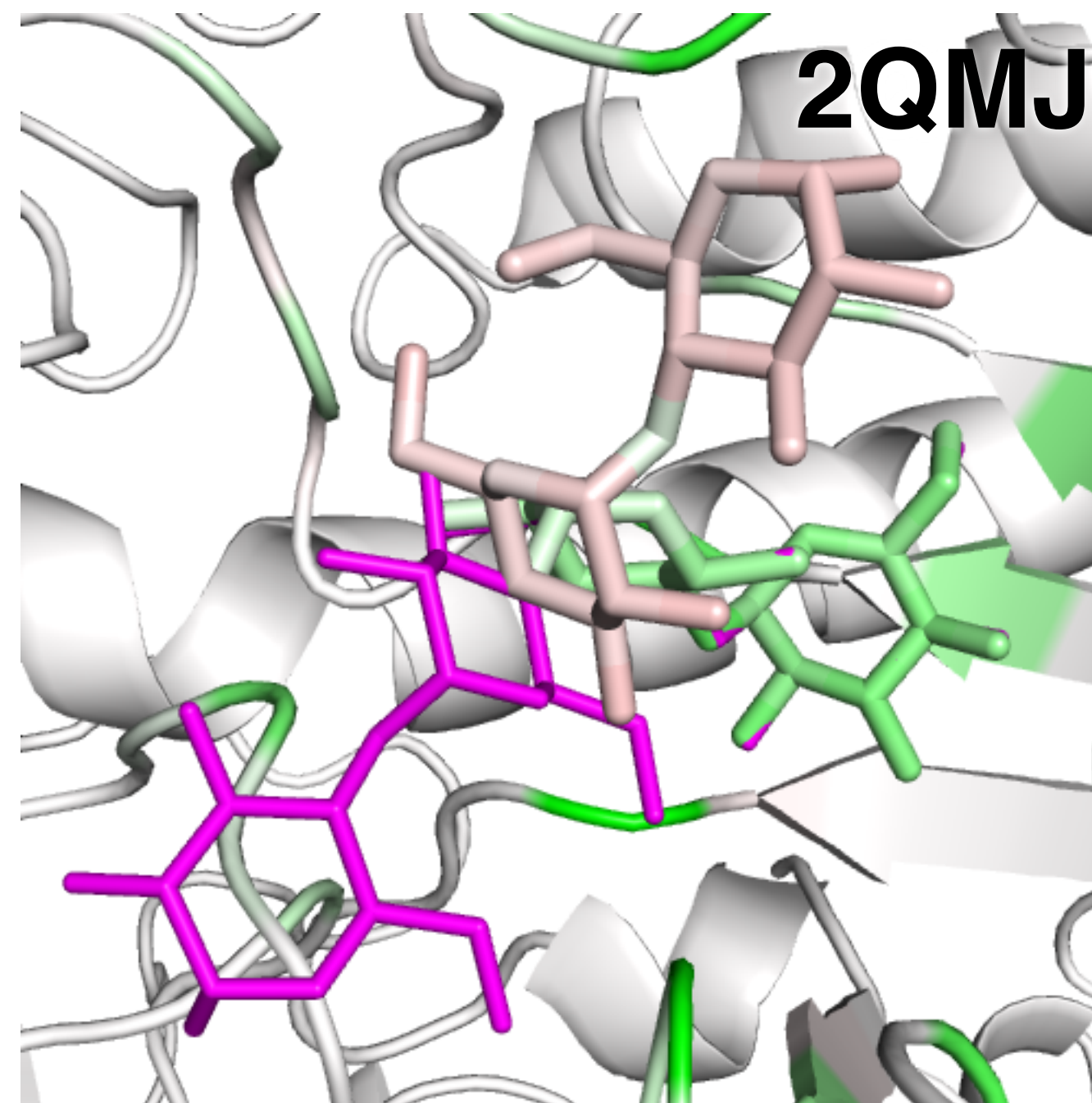
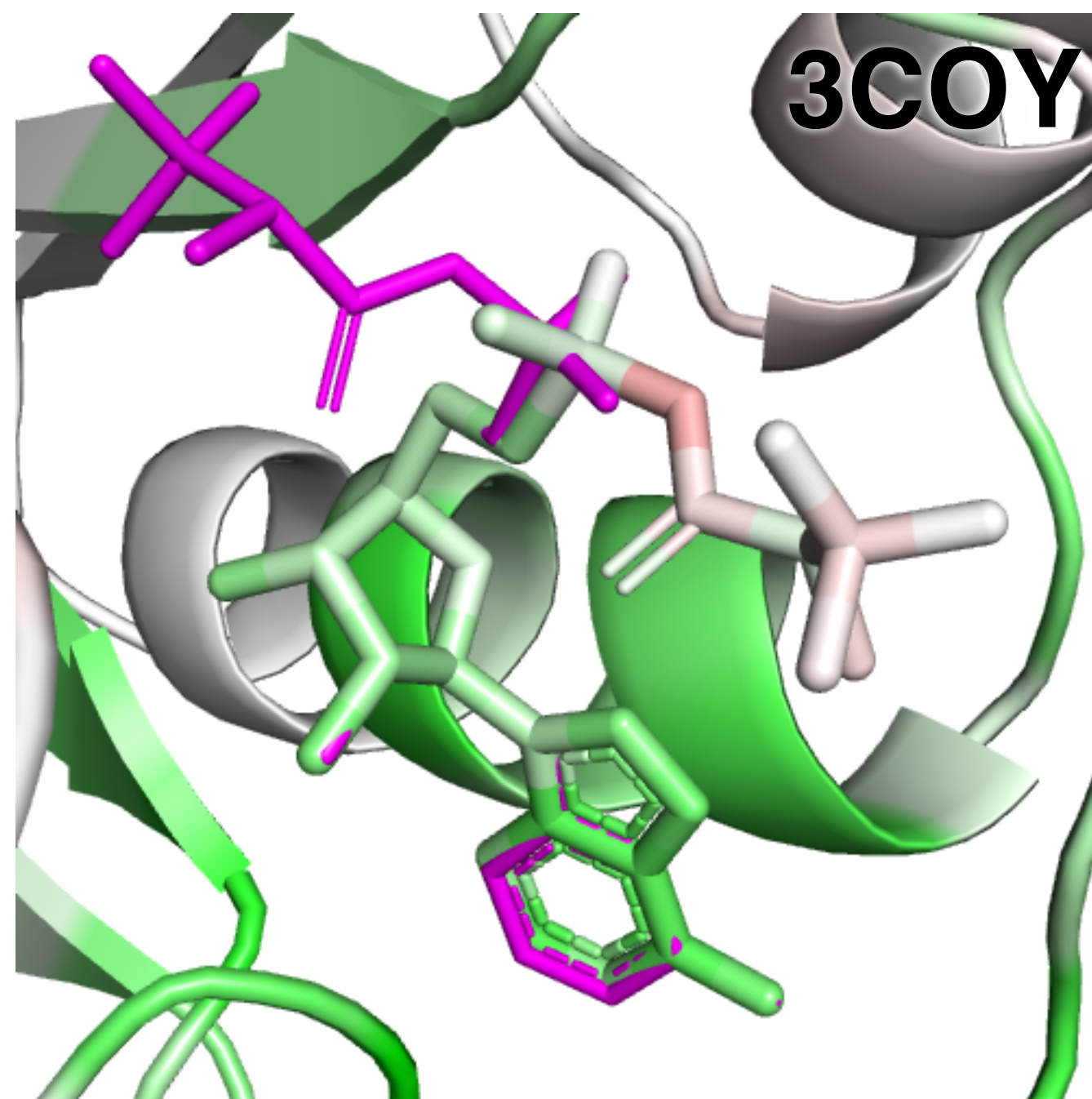


CHEMBL457424
Top Vina Pose (-8.2 kcal/mol)



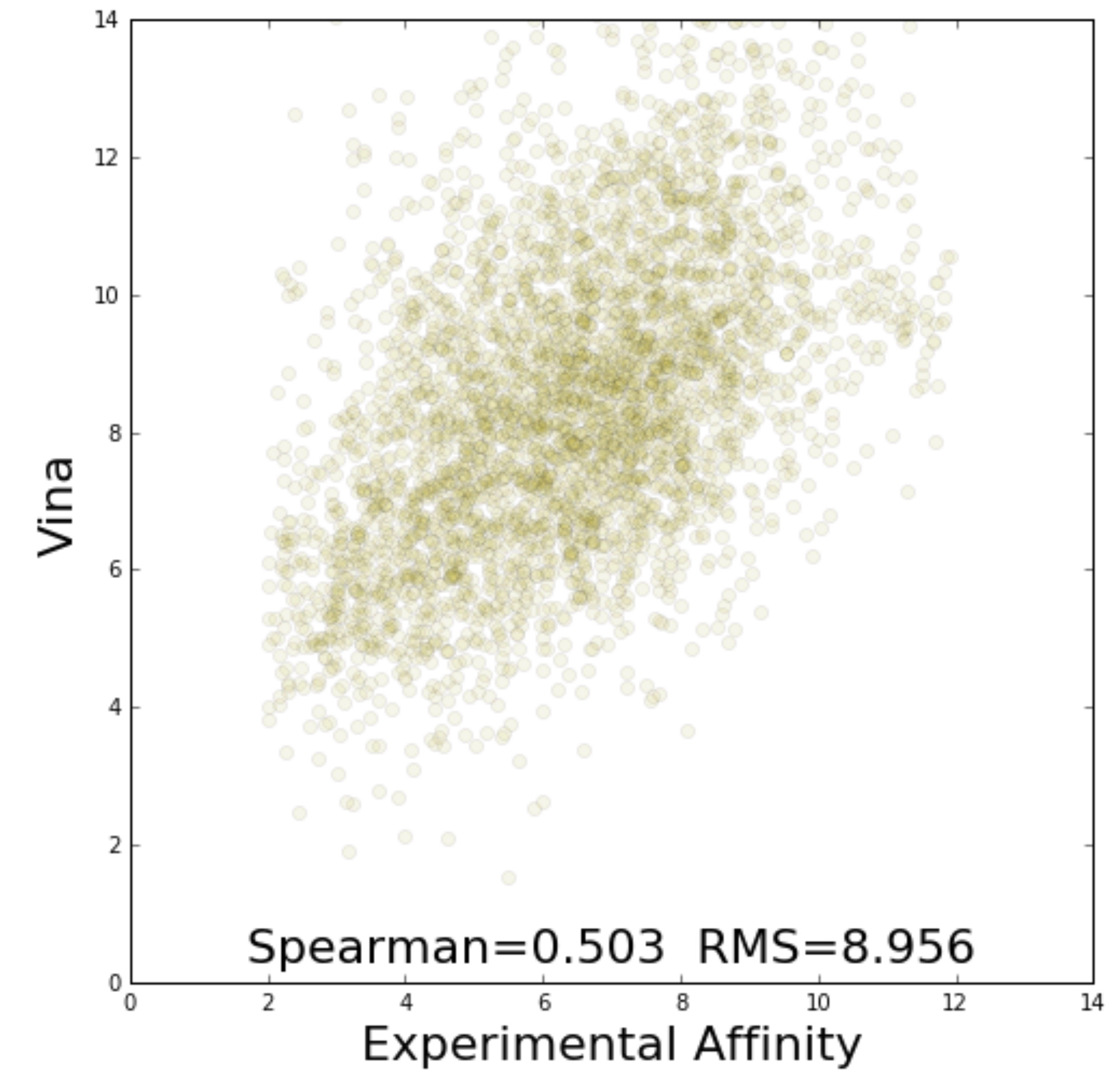
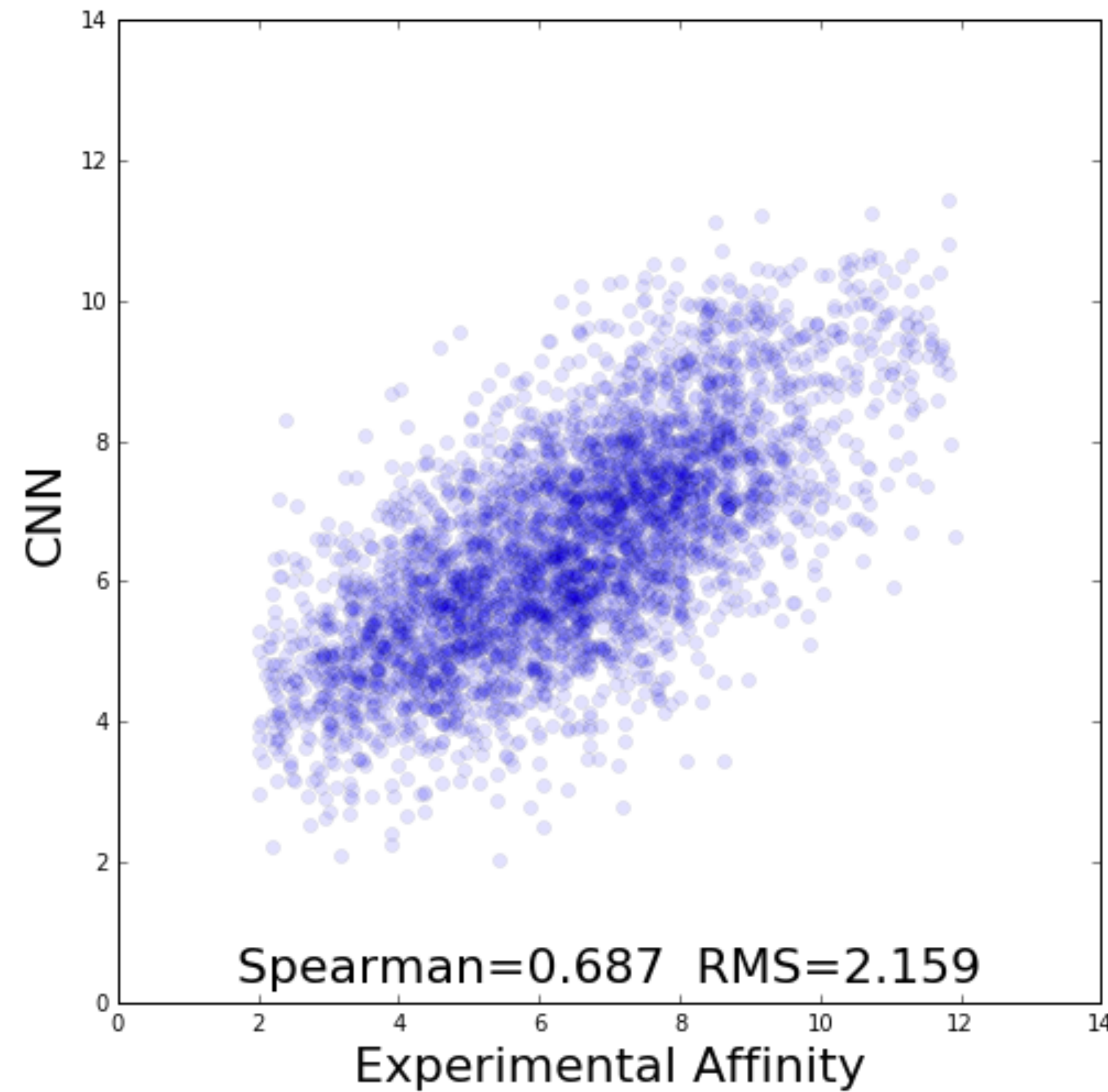
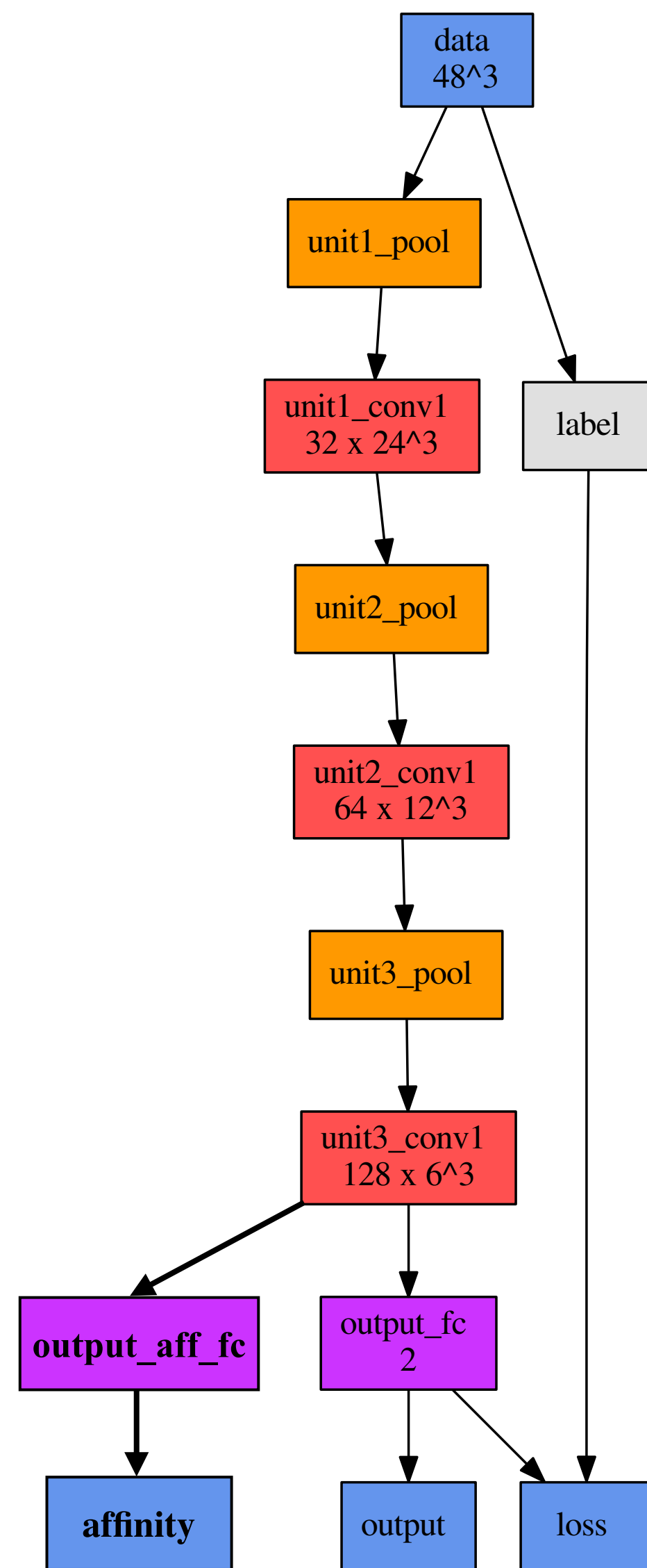
DUD-E Training Set
Score = 0.93 ± 0.03

Pose Sensitivity

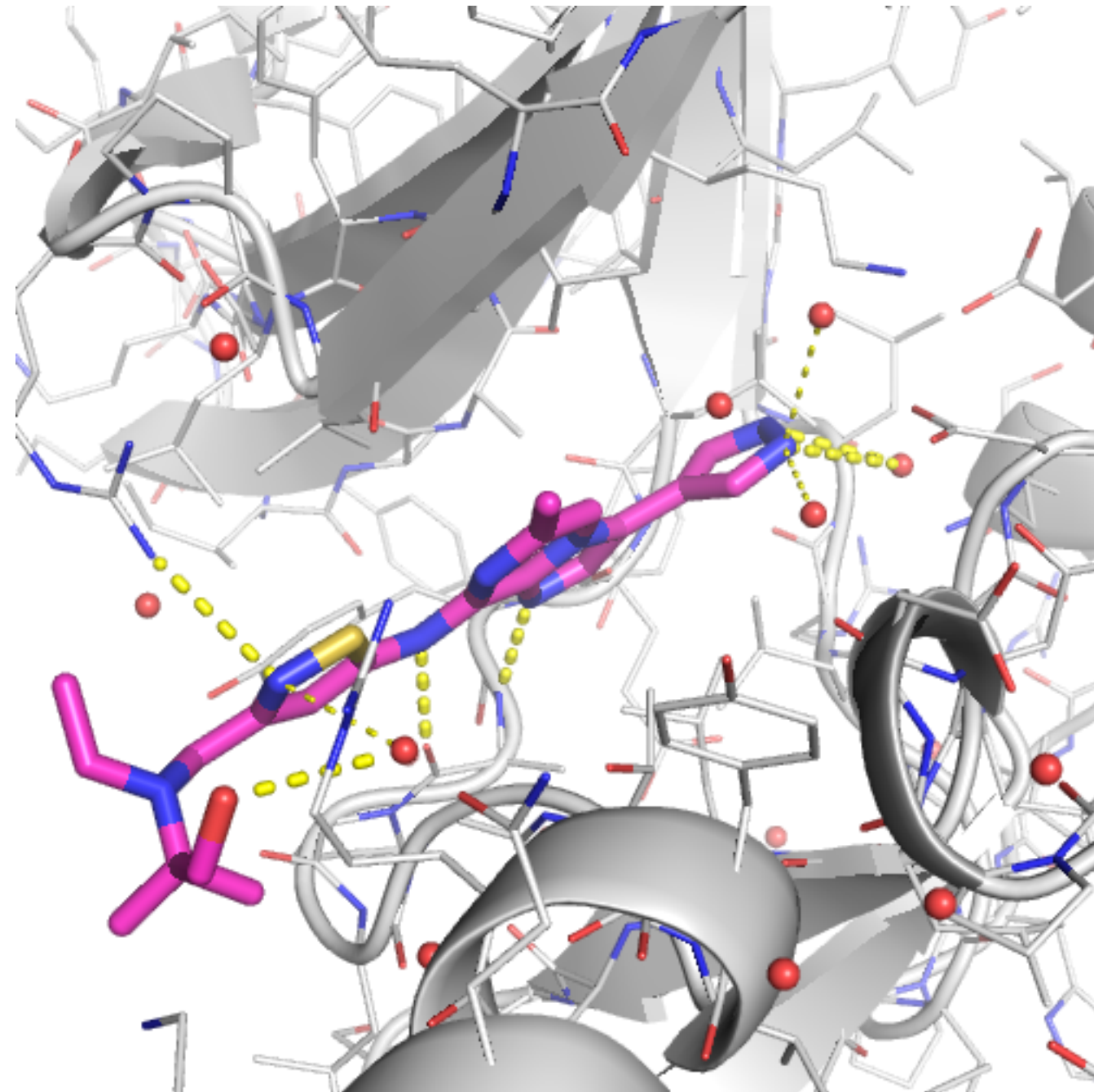


Partially Aligned Poses
Combined 2:1 Training Set

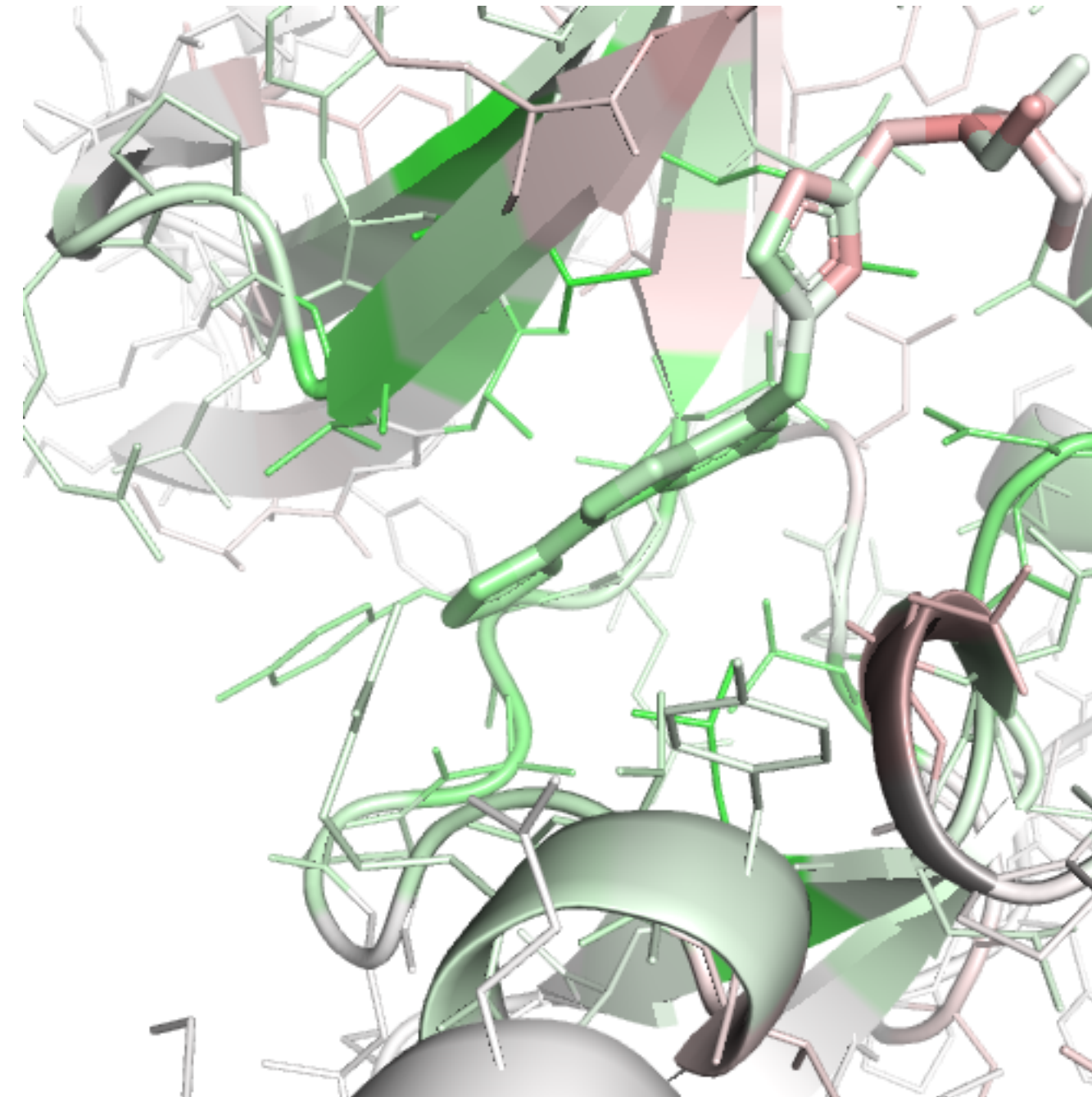
Affinity Prediction



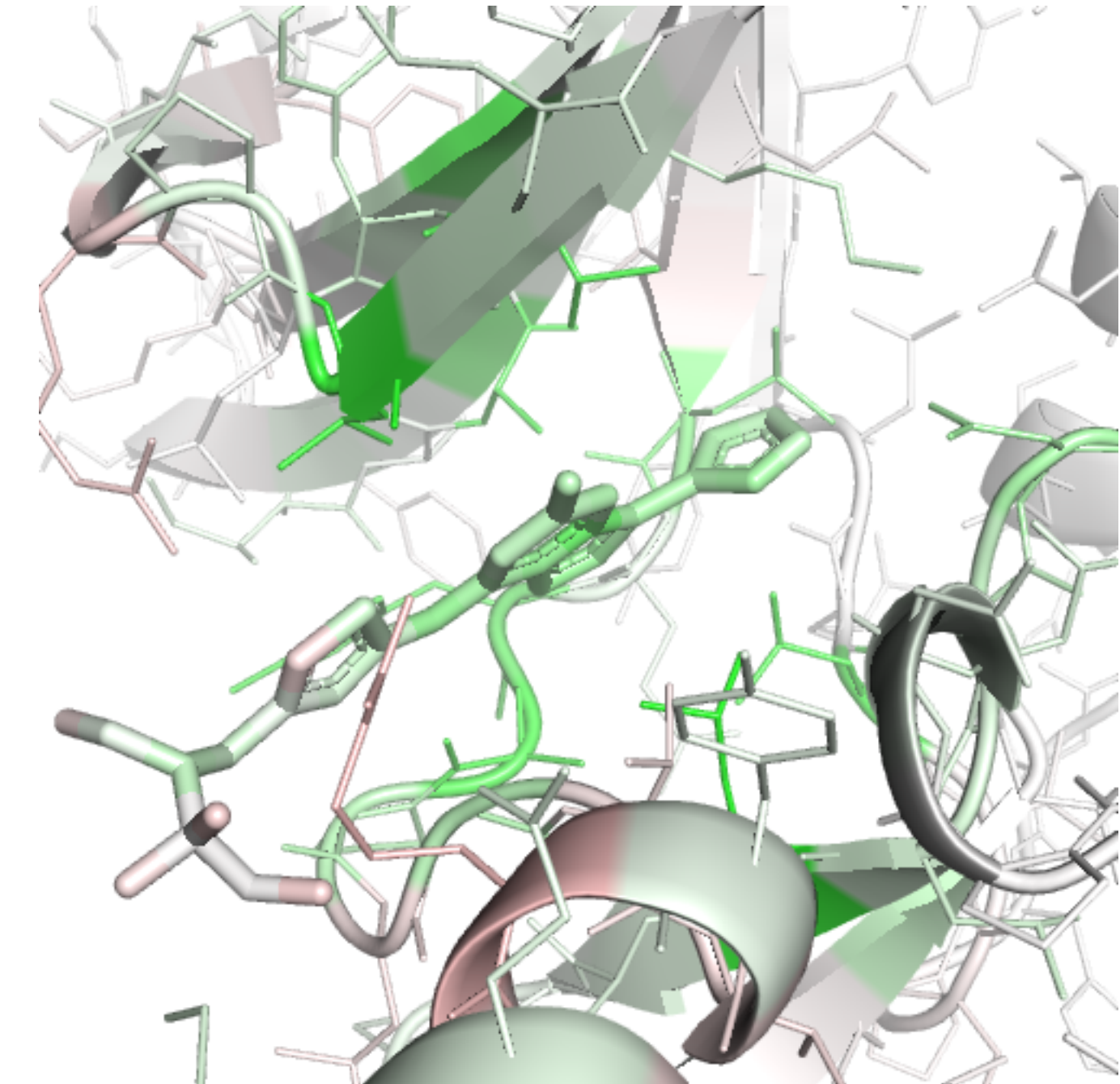
Examples



3MYG



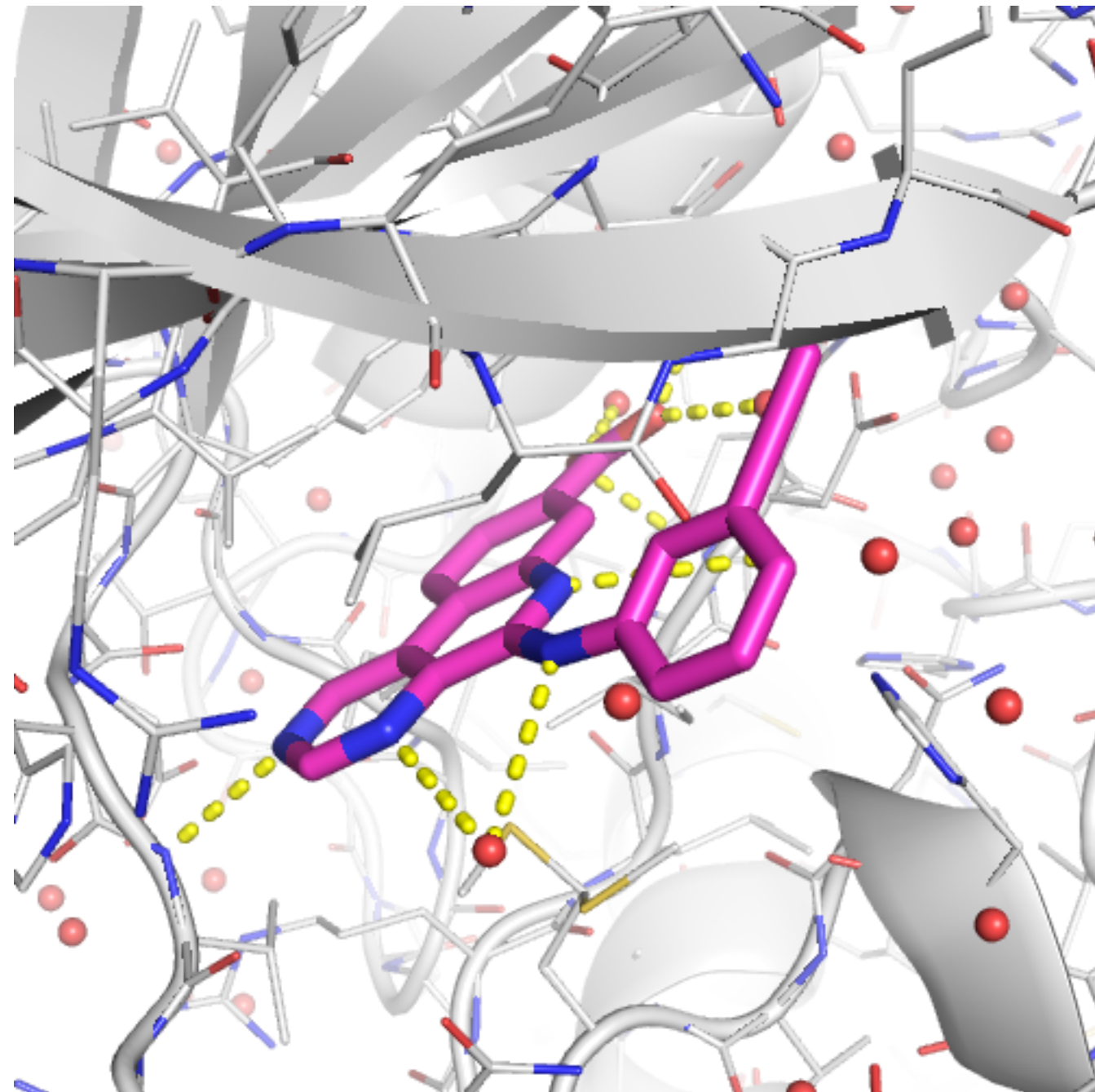
Vina (12.71Å)



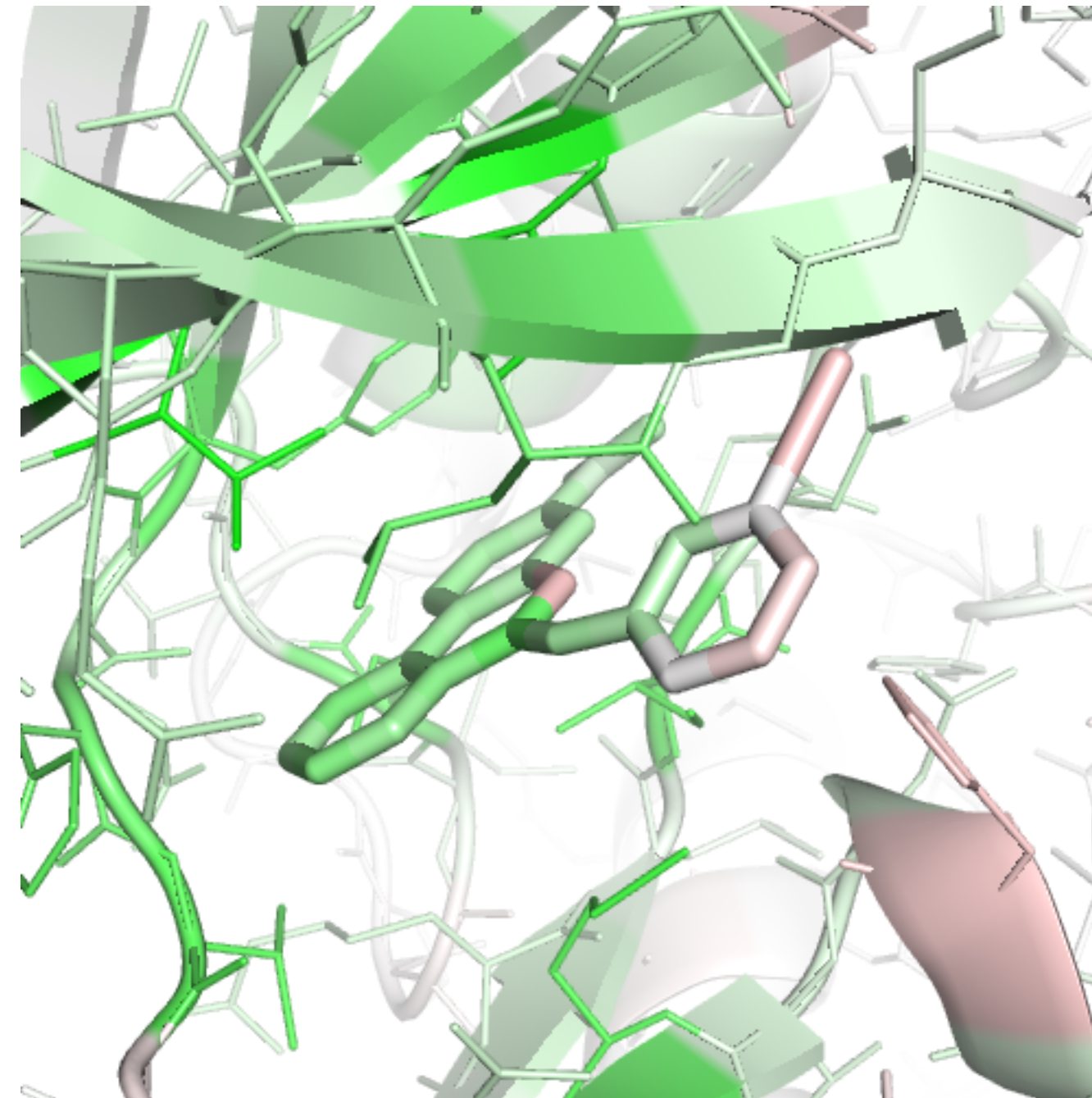
CNN (0.96Å)



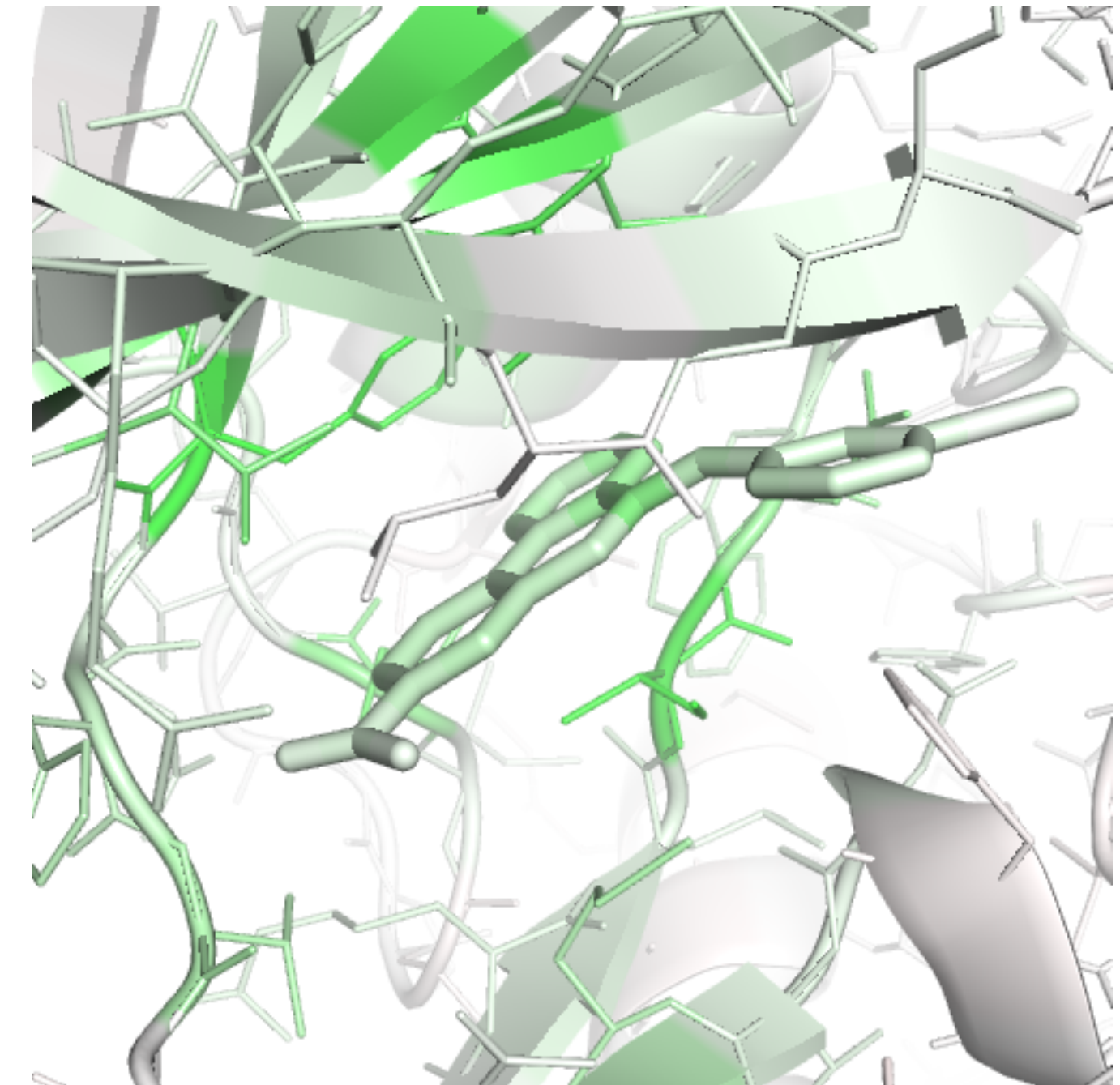
Examples



3PE2



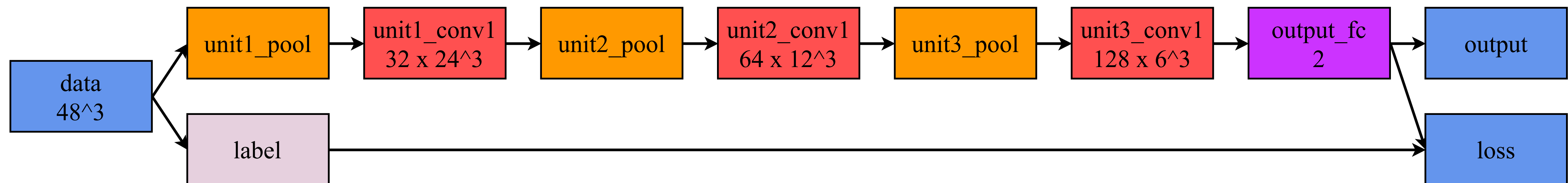
Vina (0.25Å)



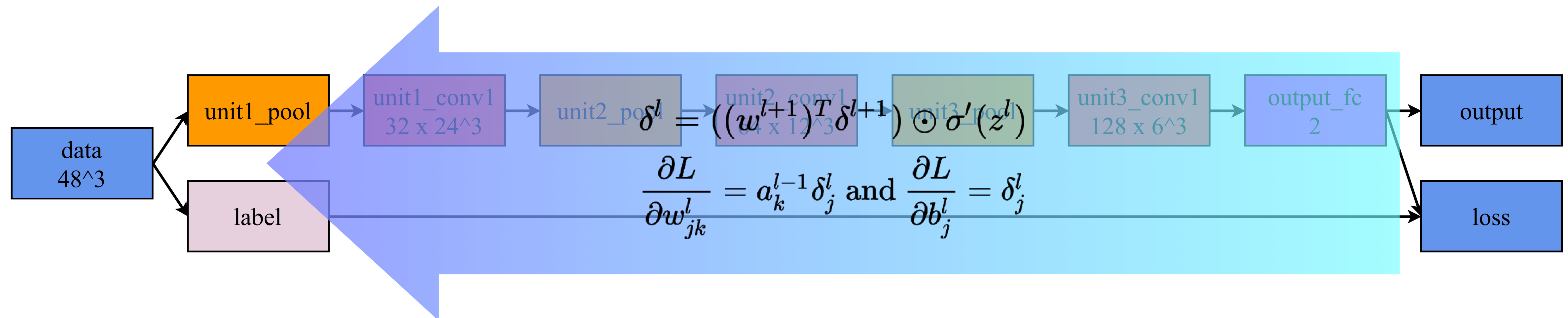
CNN (5.27Å)



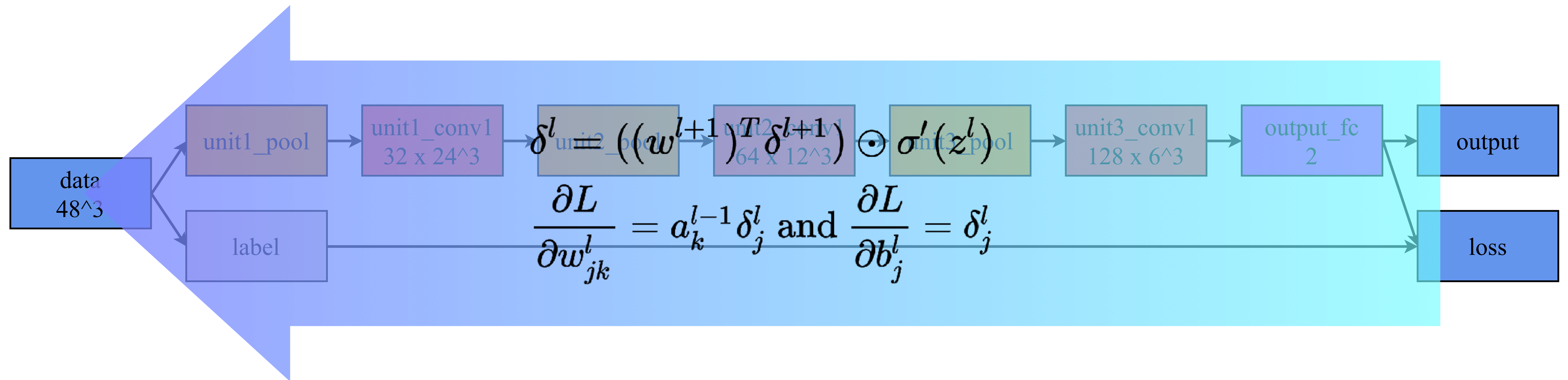
Beyond Scoring



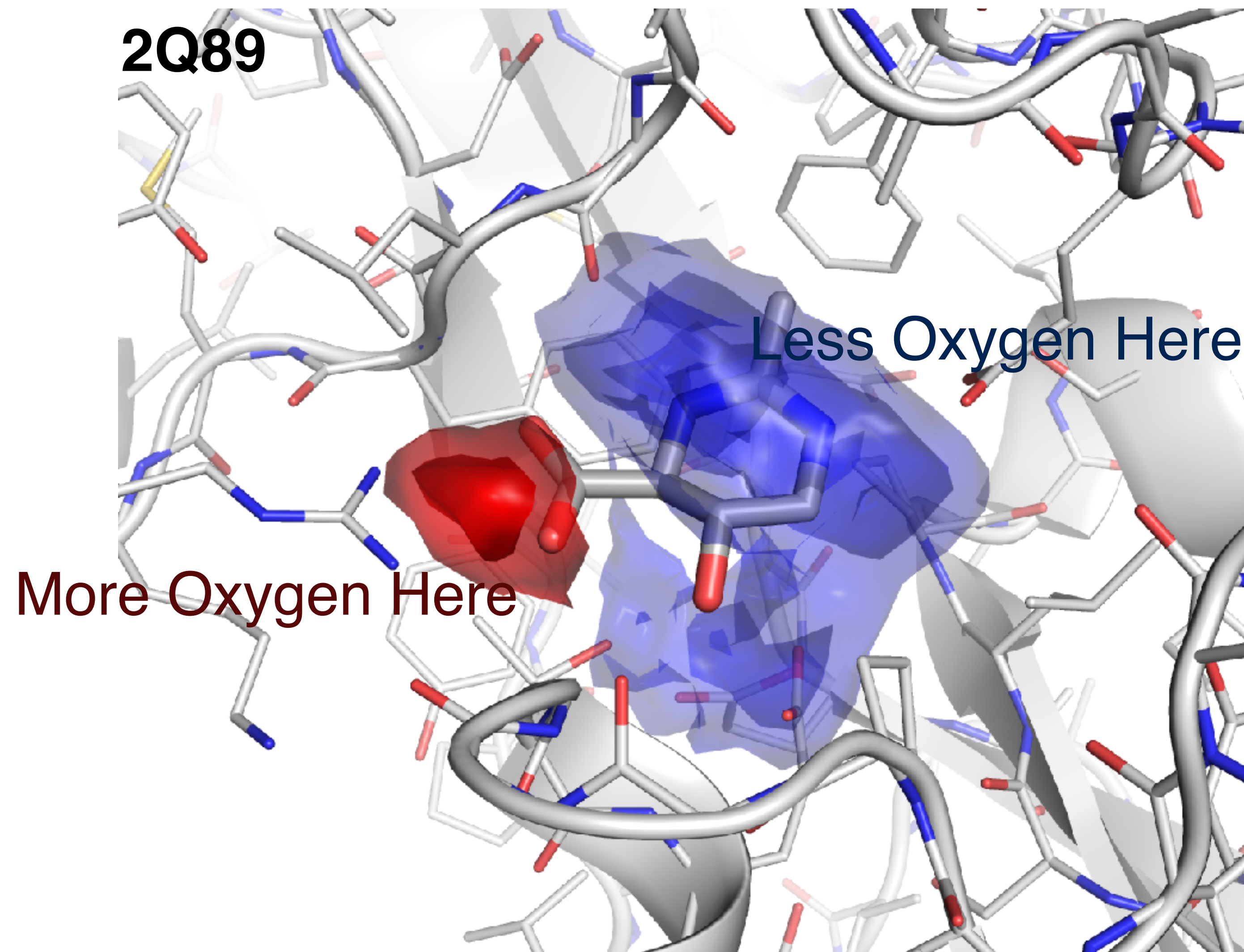
Beyond Scoring



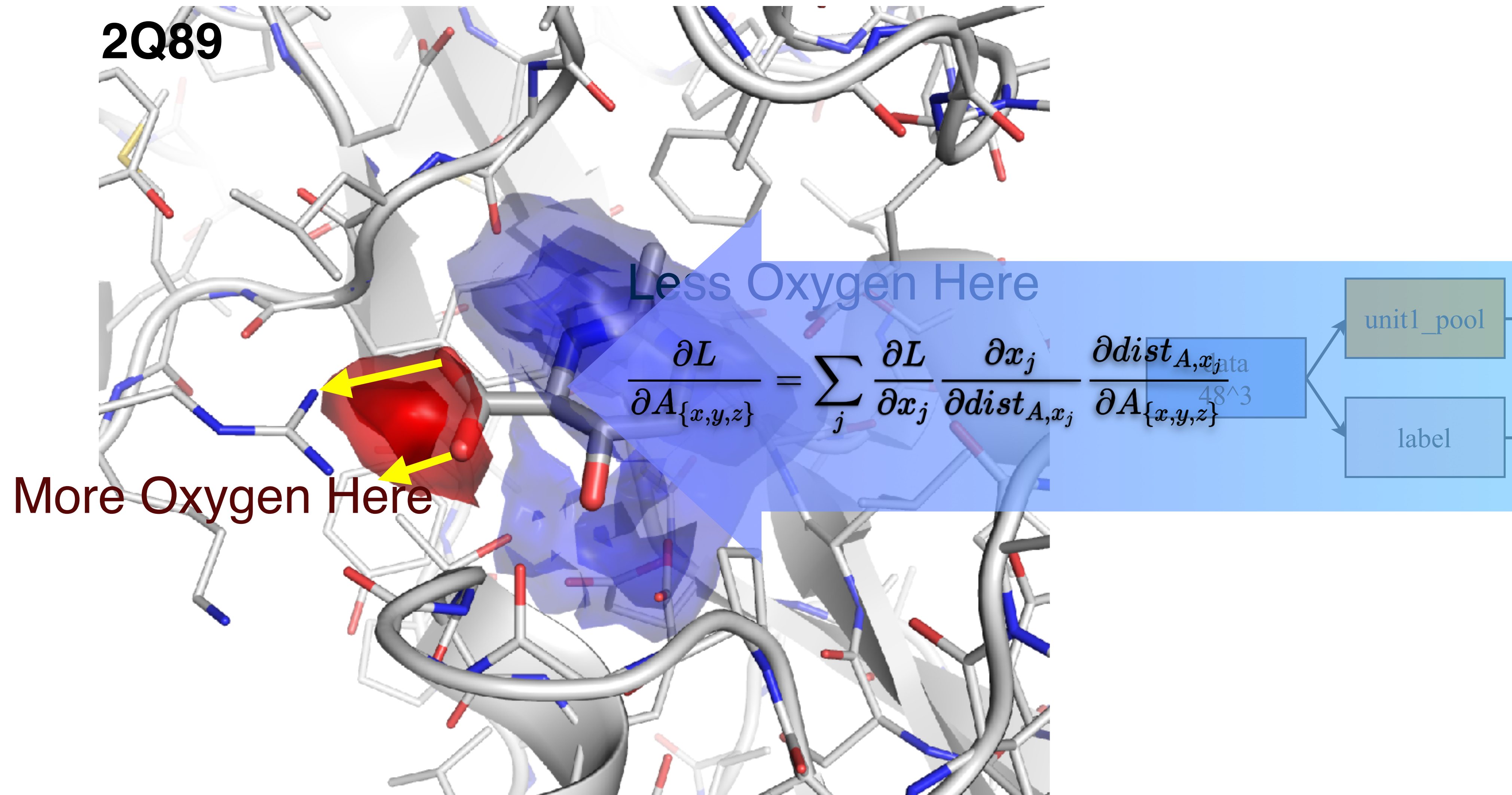
Beyond Scoring

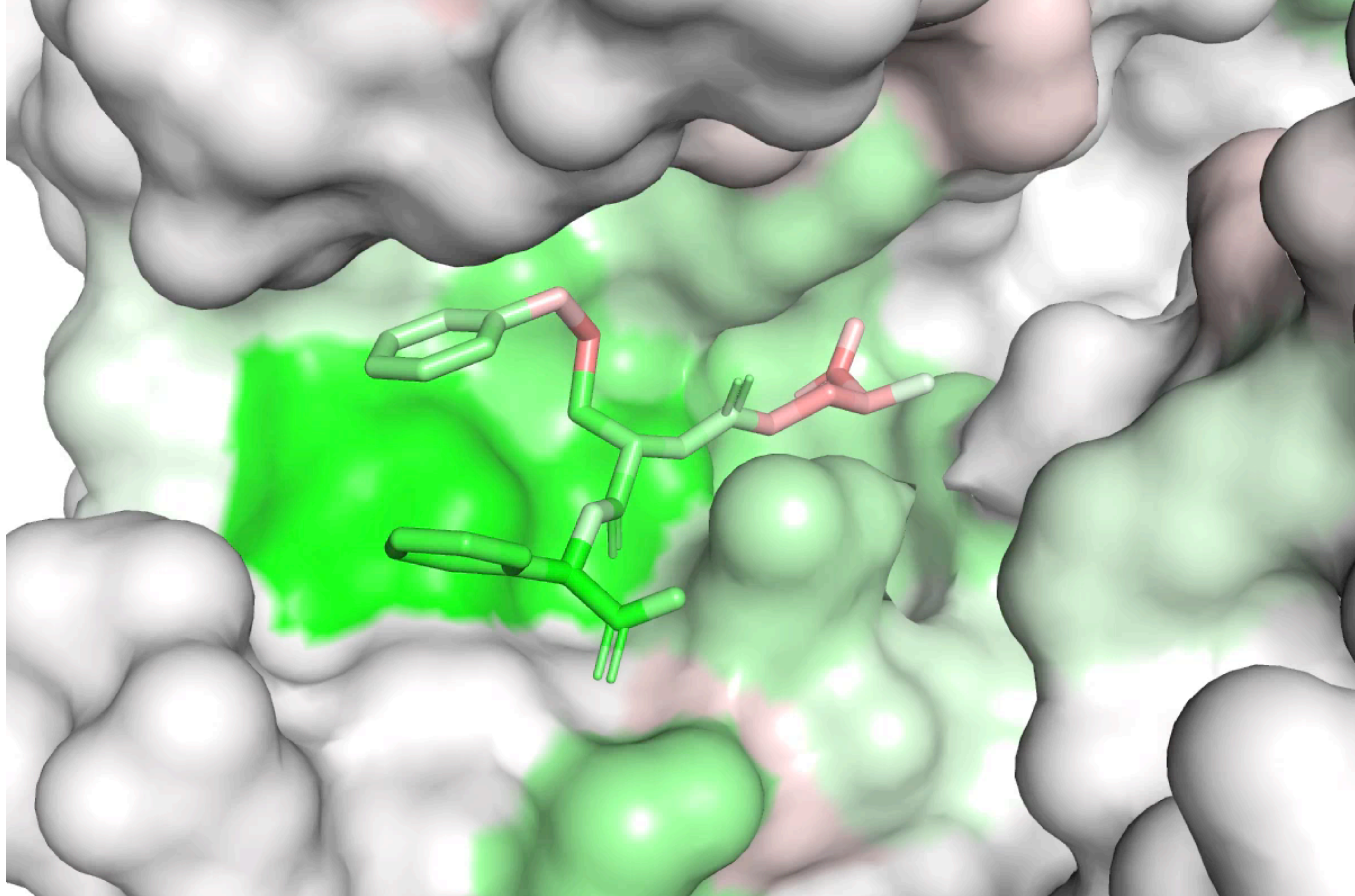


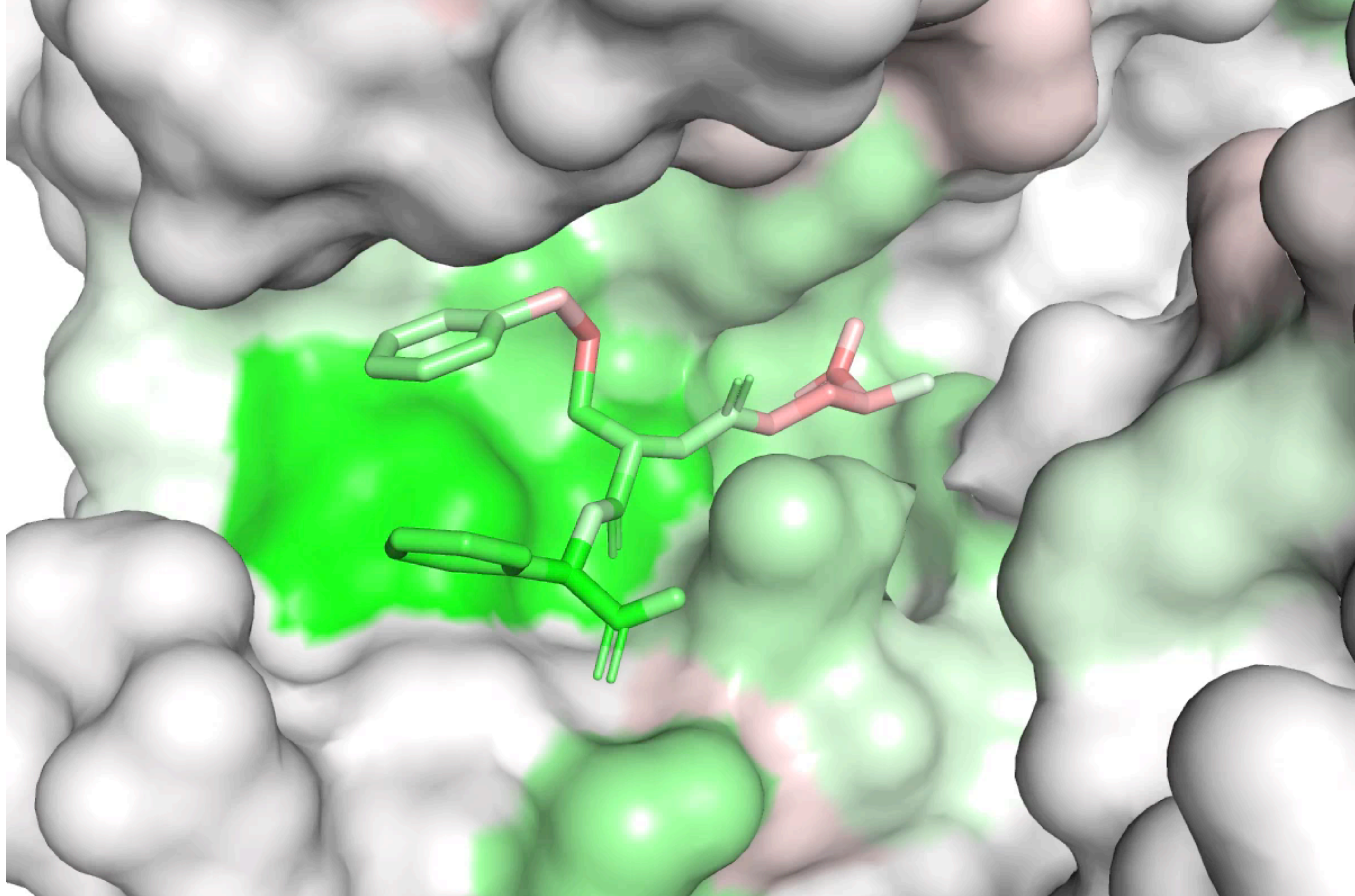
Beyond Scoring

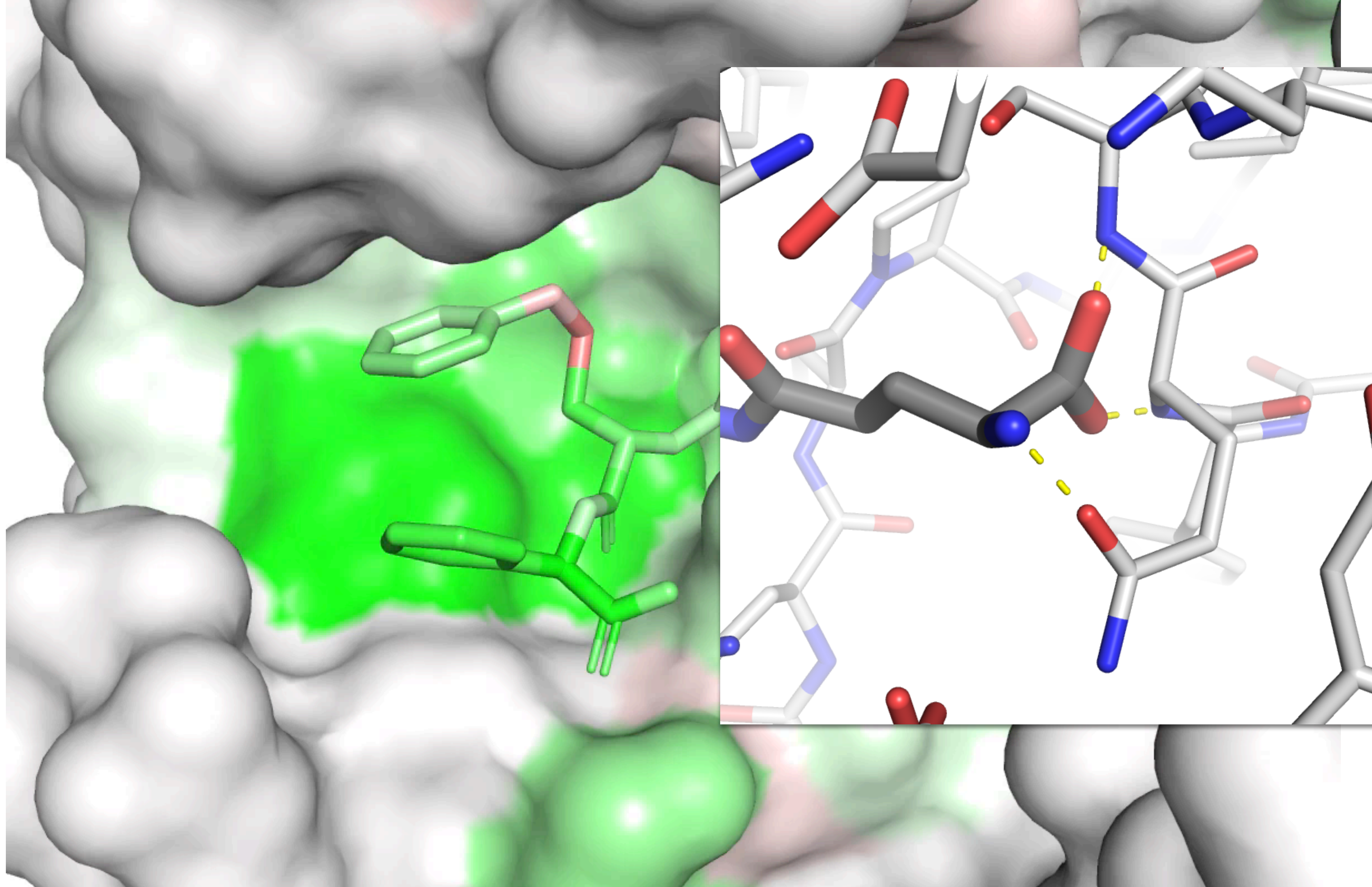


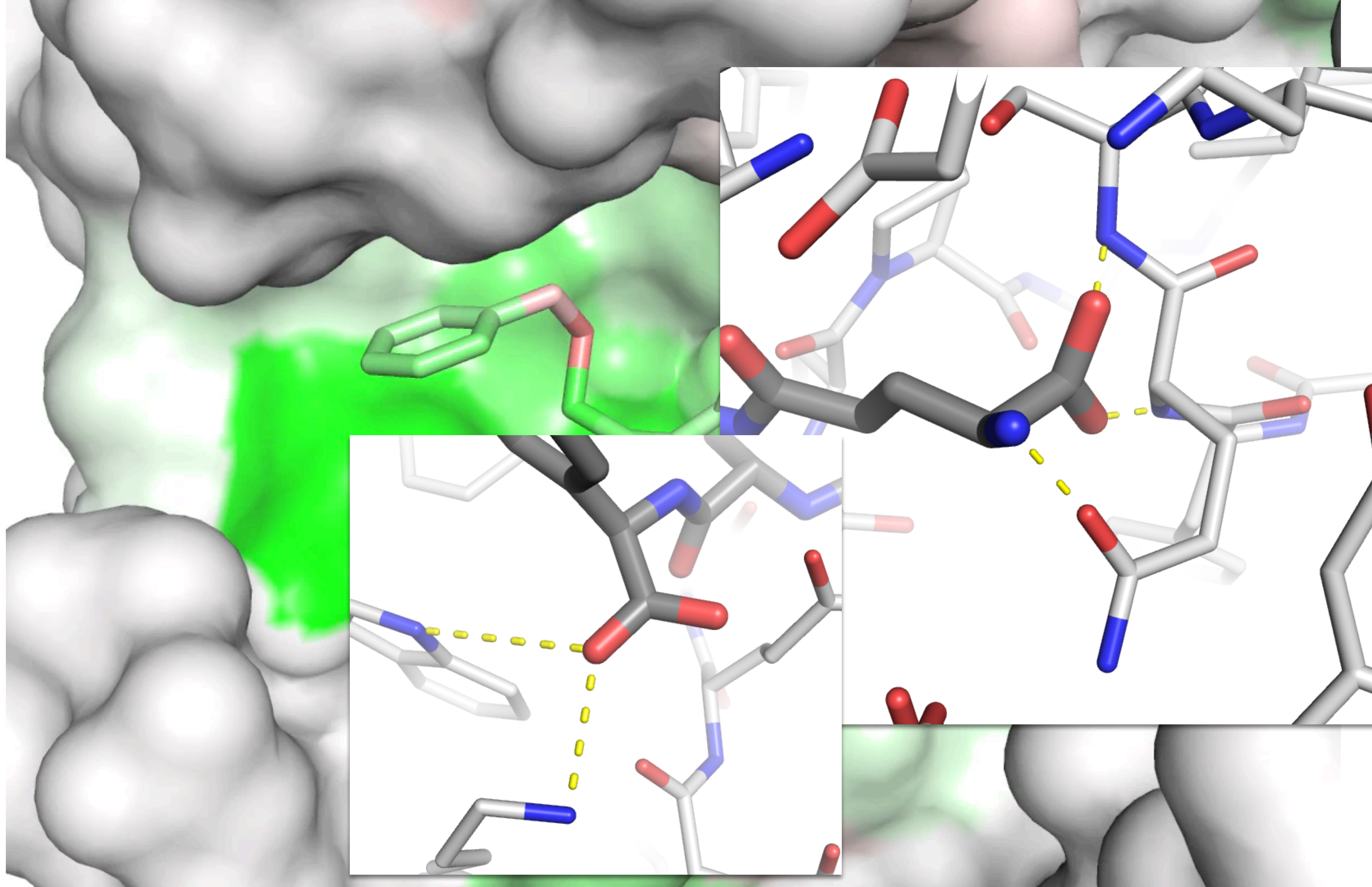
Beyond Scoring

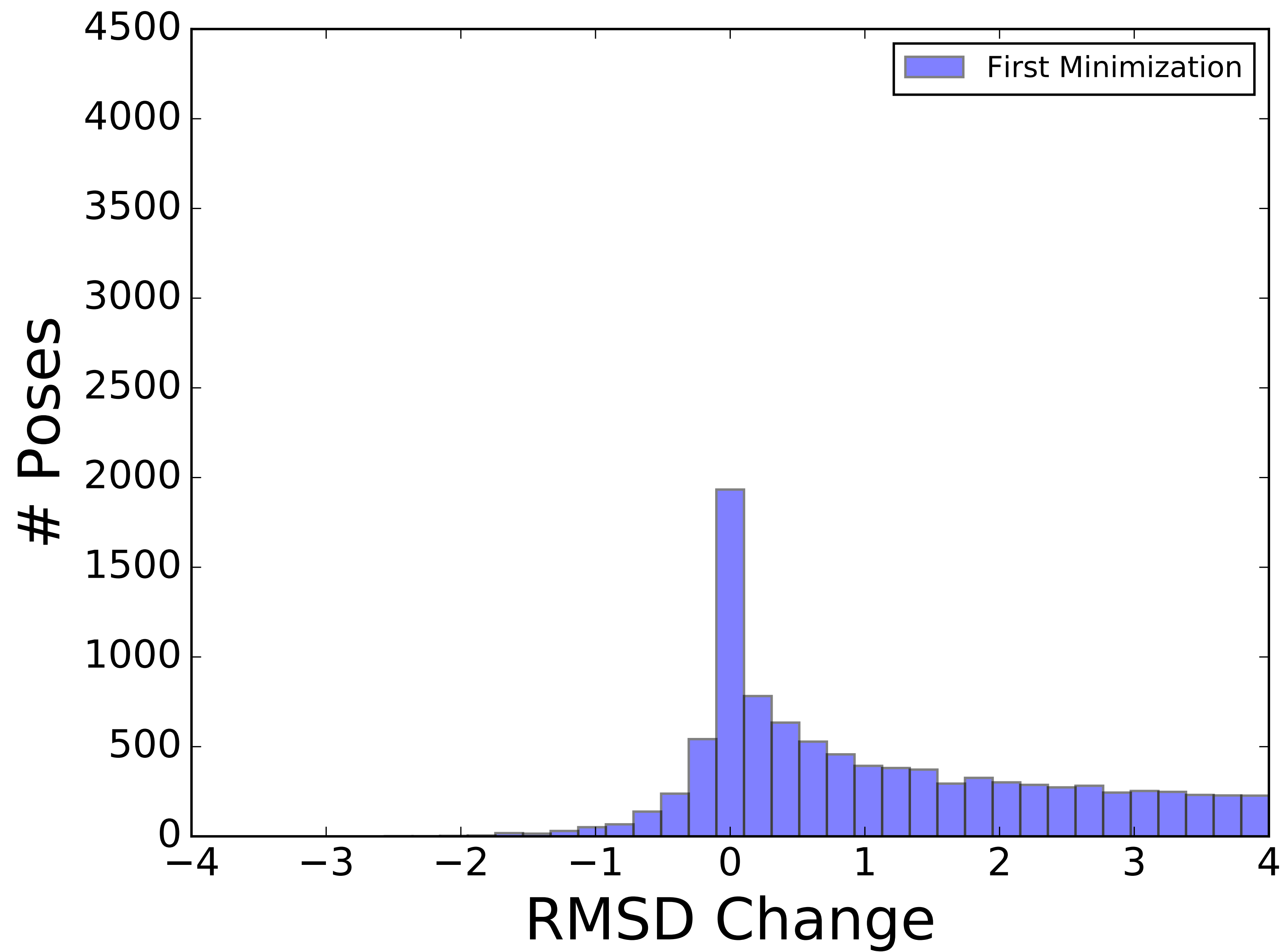


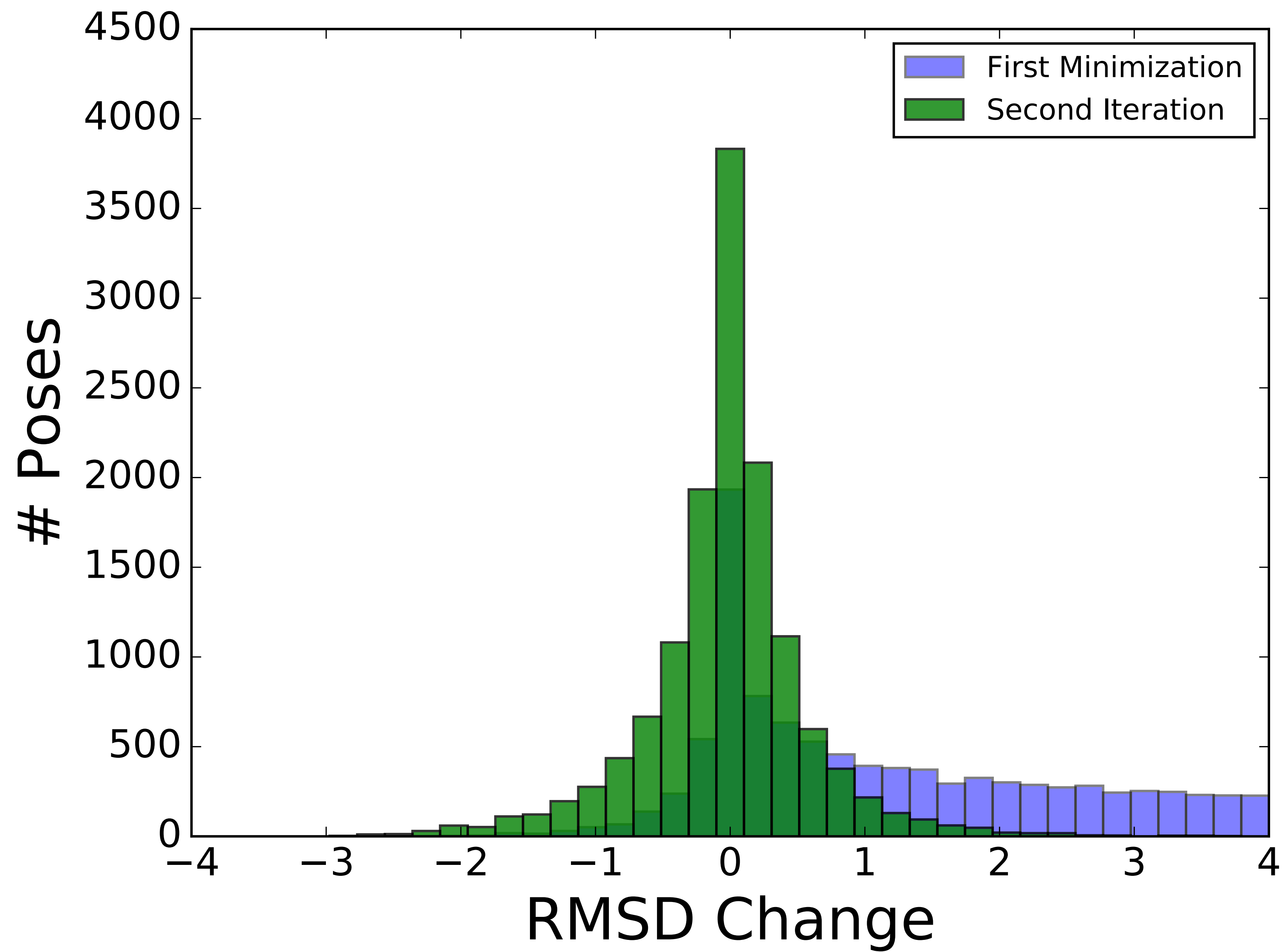


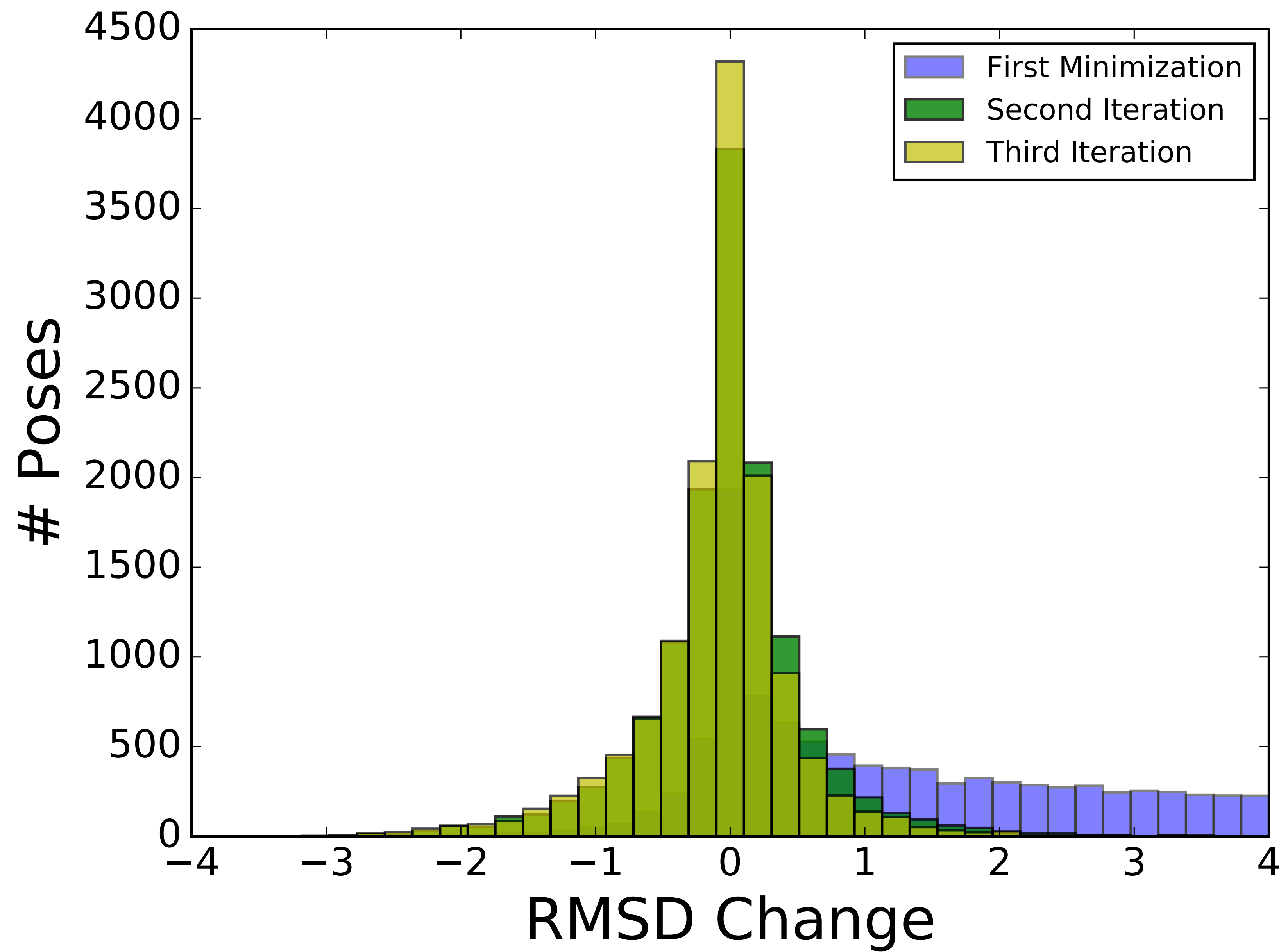












The Future

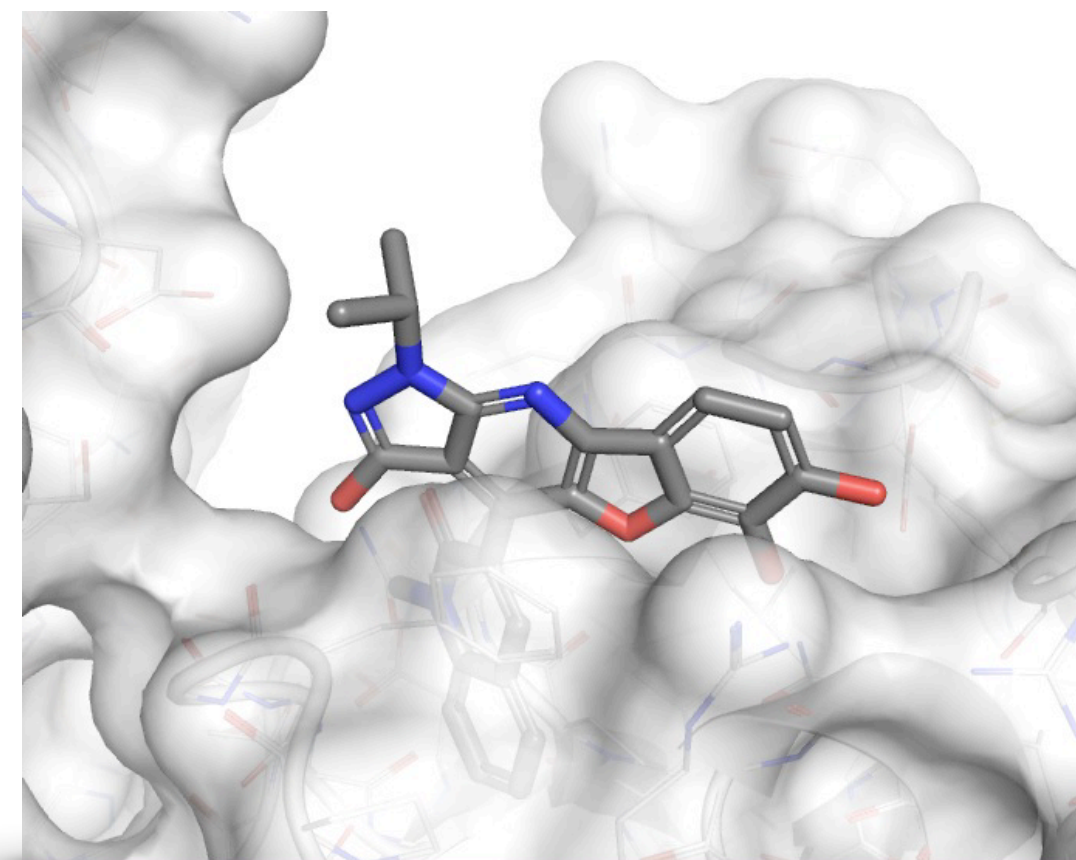
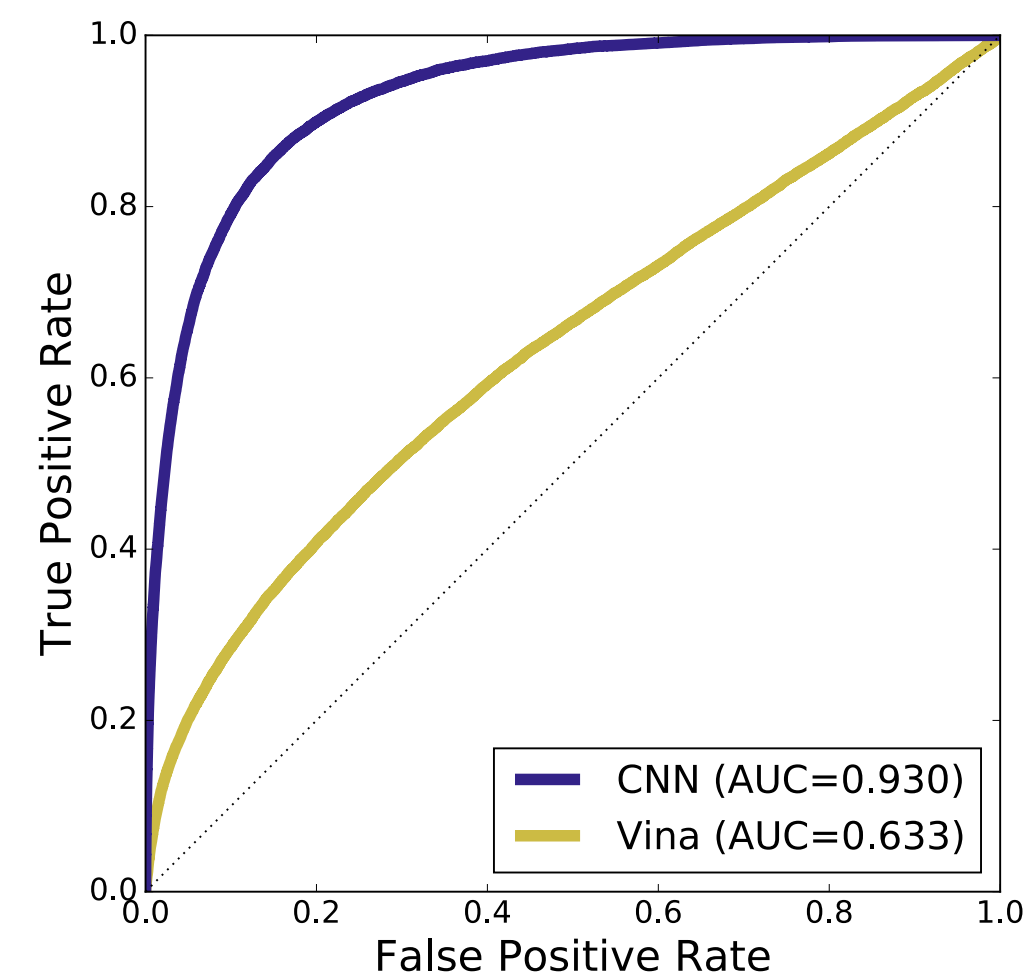
Pose
Selection



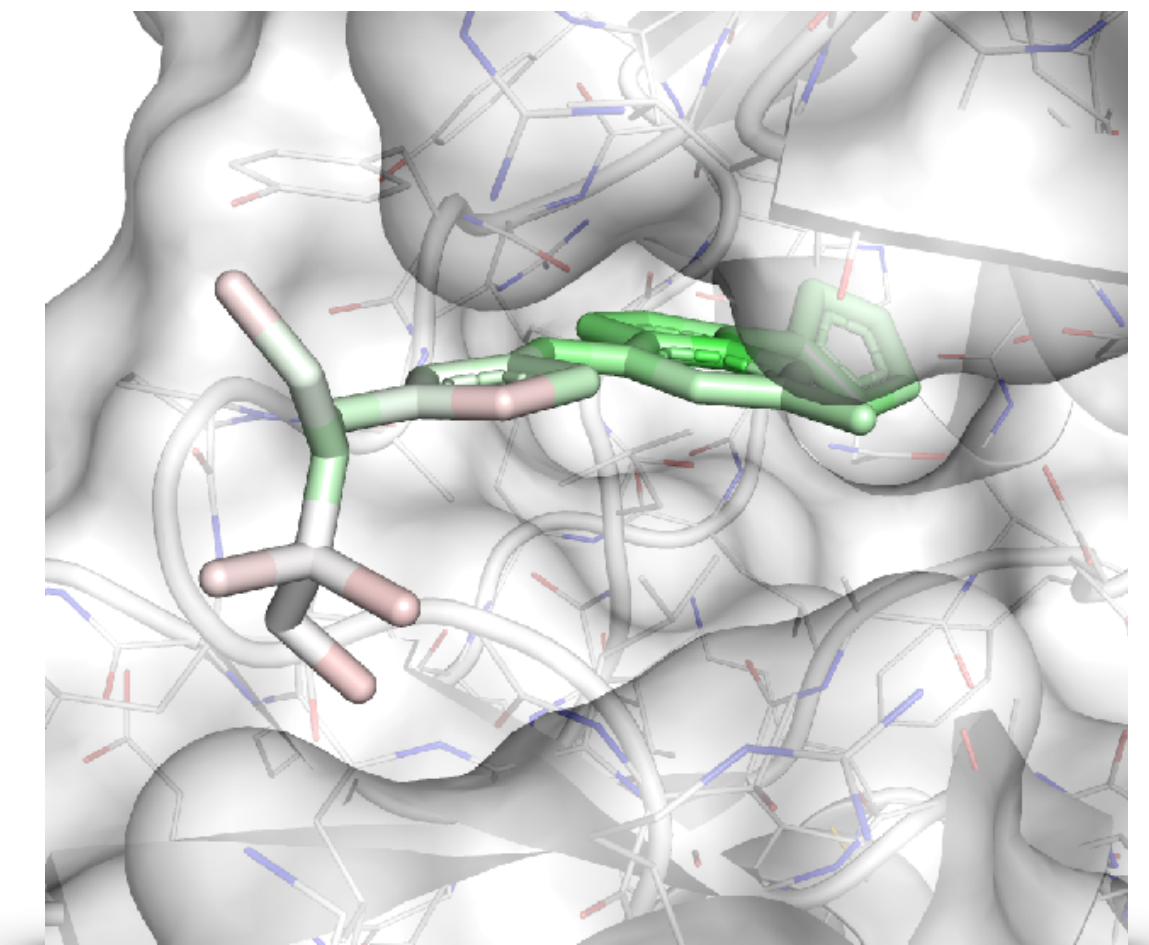
Pose
Generation



Compound
Generation



Virtual Screening



Lead Optimization

The Future

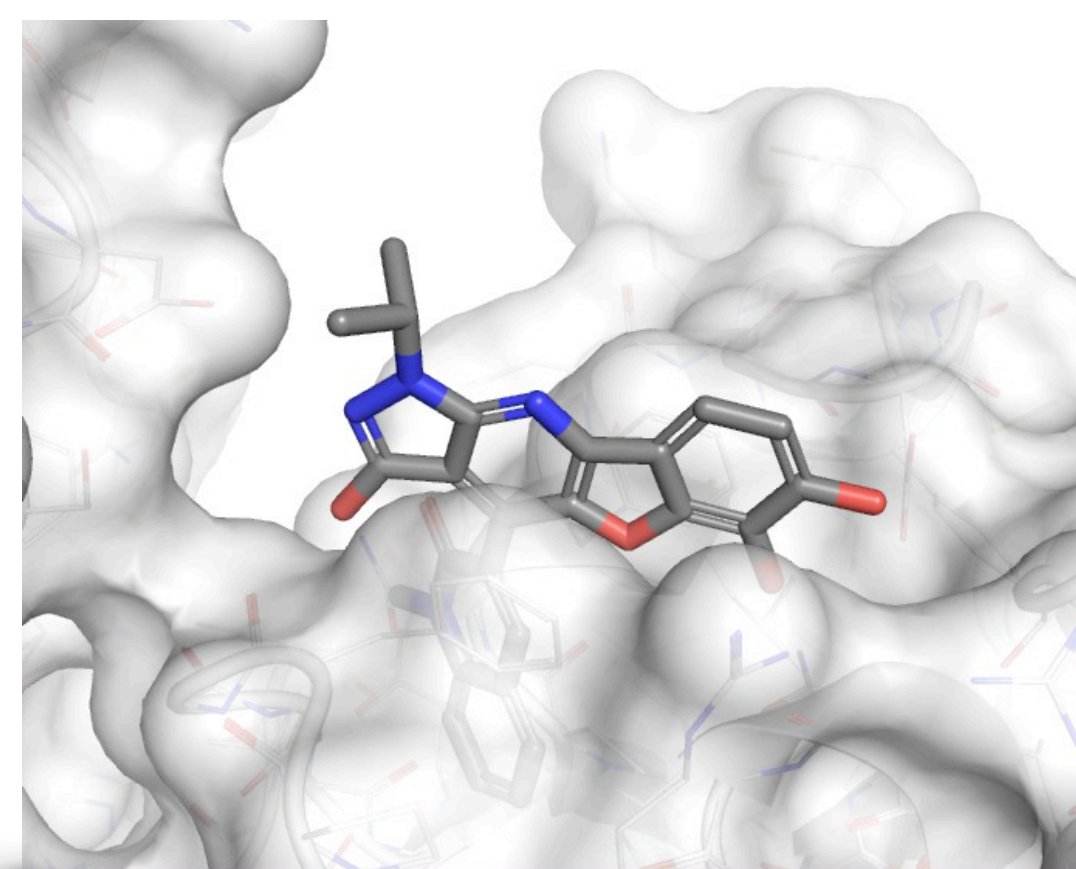
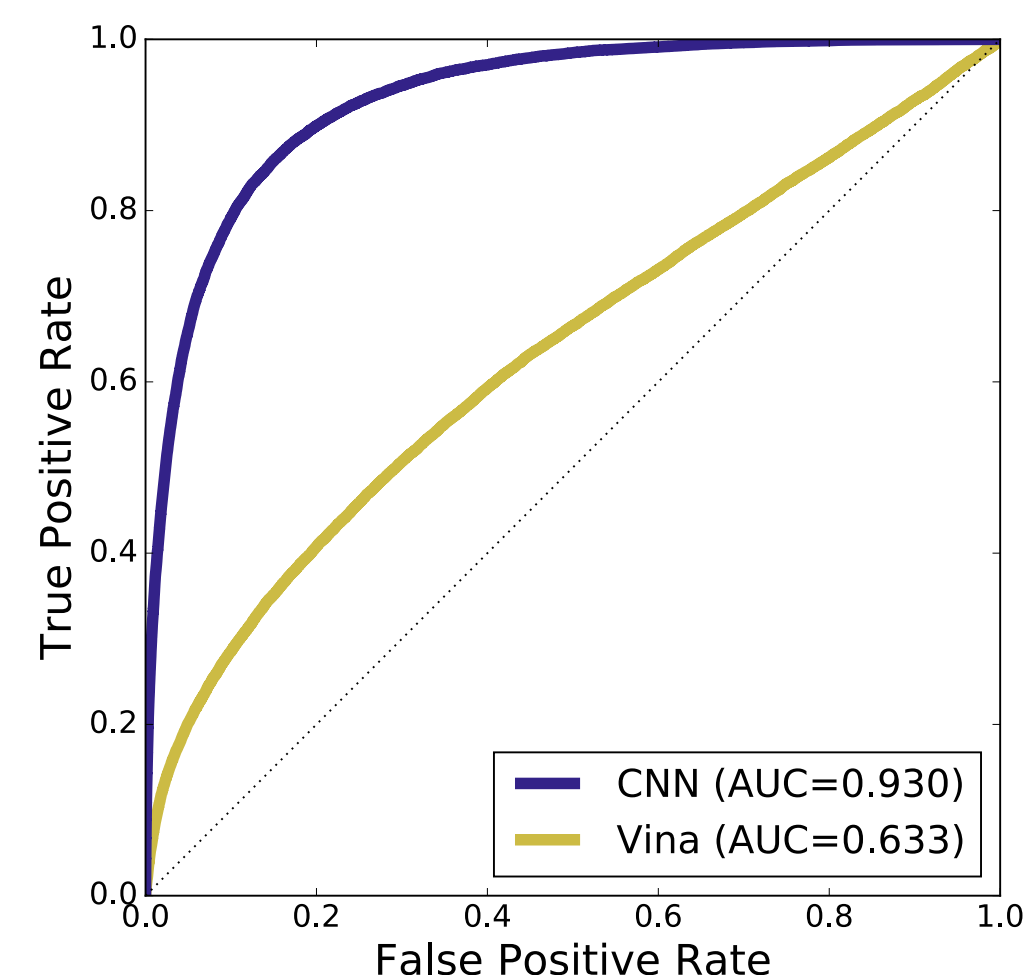
Pose
Selection



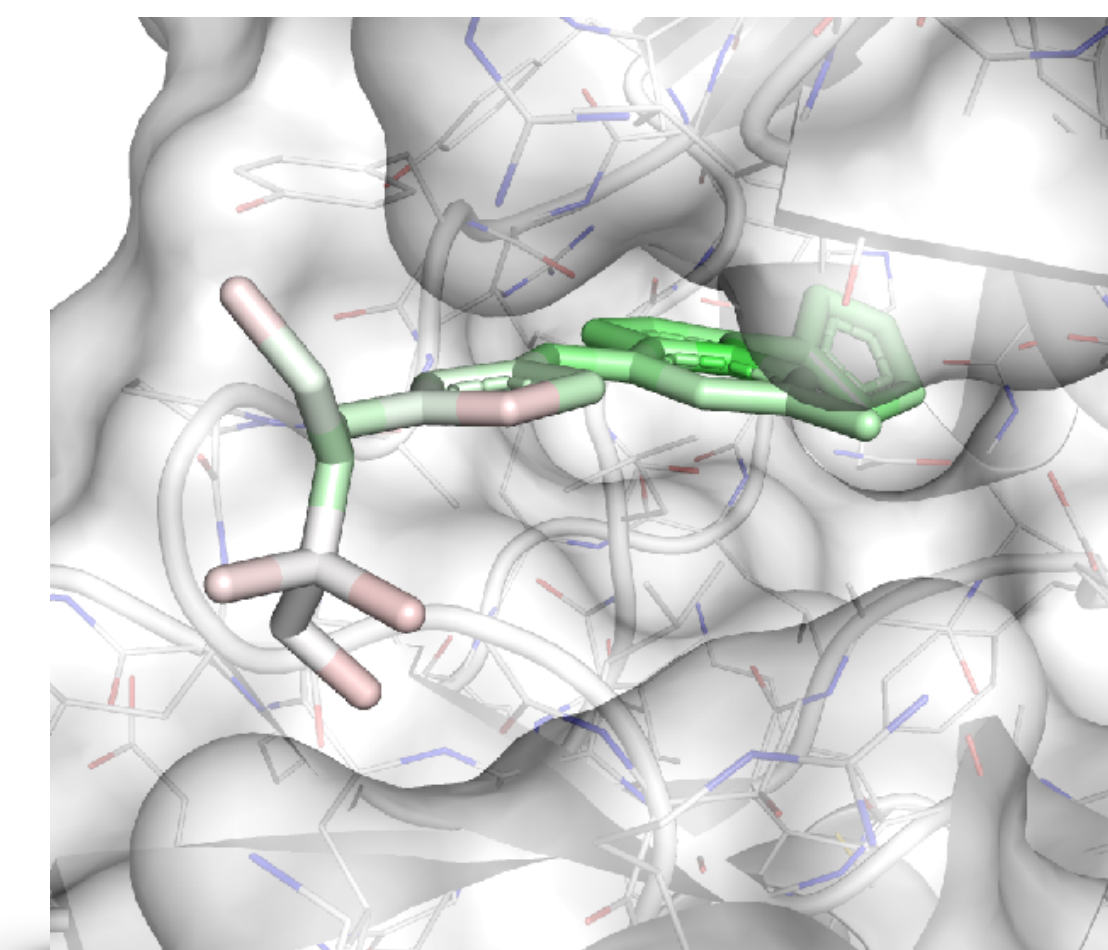
Pose
Generation



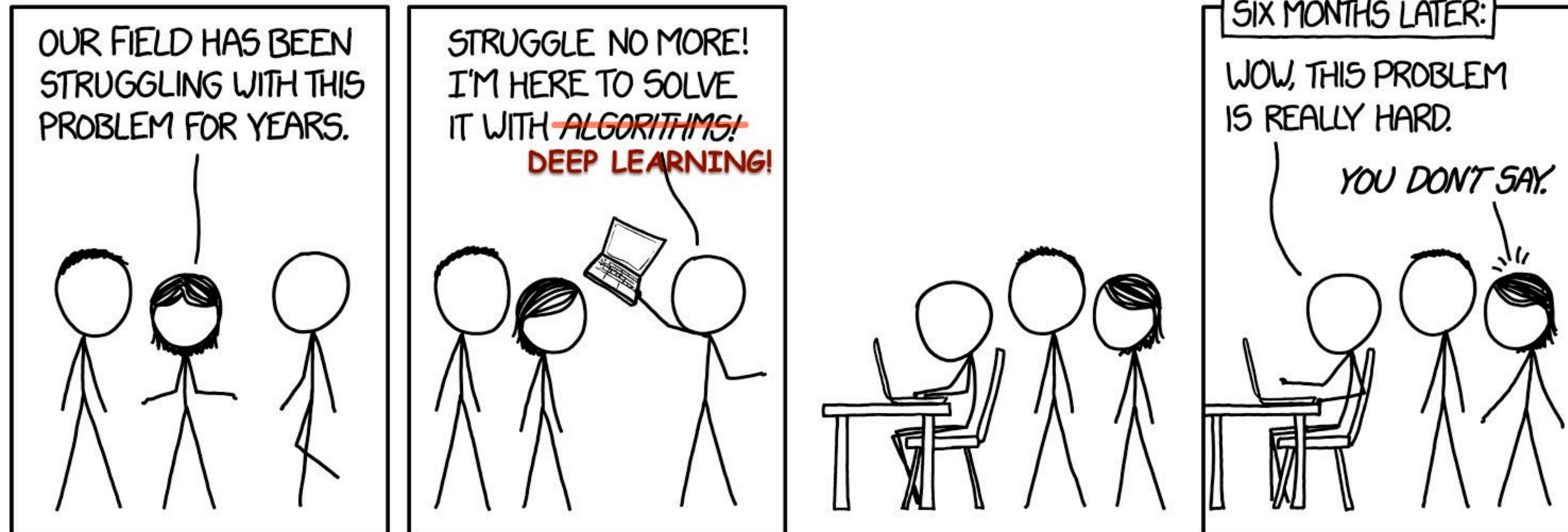
Compound
Generation



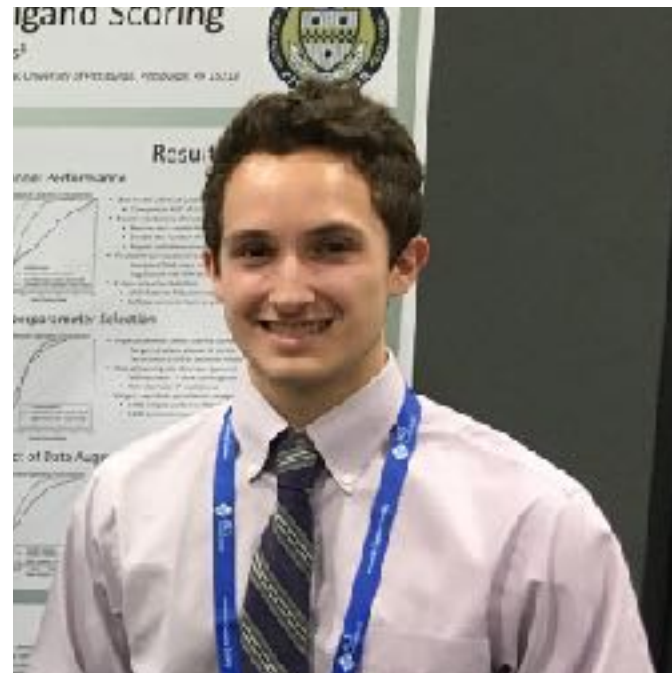
Virtual Screening



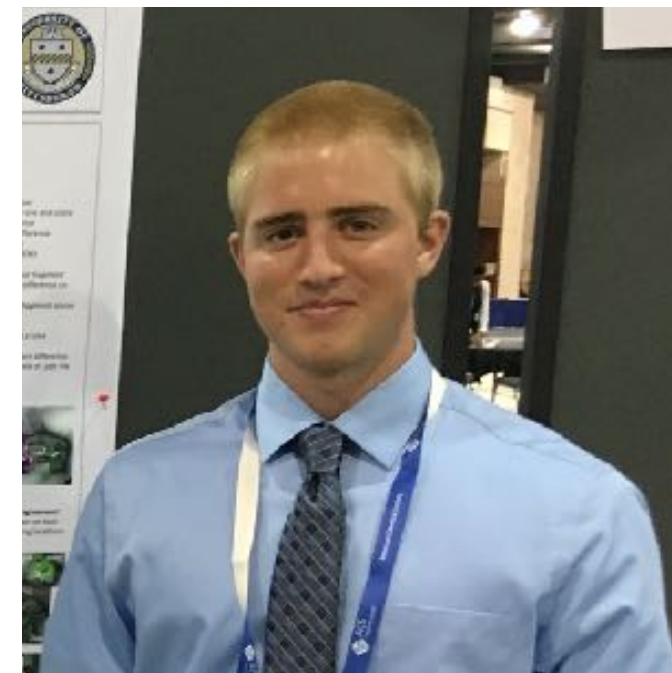
Lead Optimization



Acknowledgements



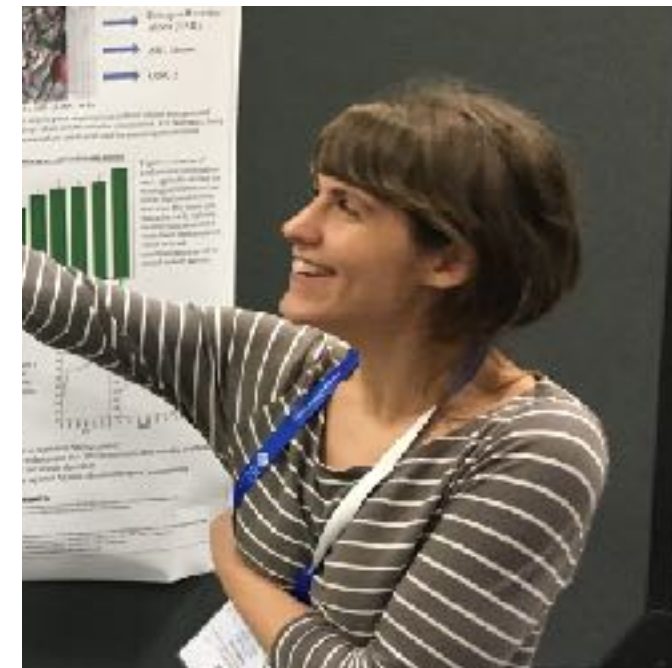
Matt Ragoza



Josh Hochuli



Elisa Idrobo



Jocelyn Sunseri

Group Members

Jocelyn Sunseri

Matt Ragoza

Josh Hochuli

Pulkit Mittal

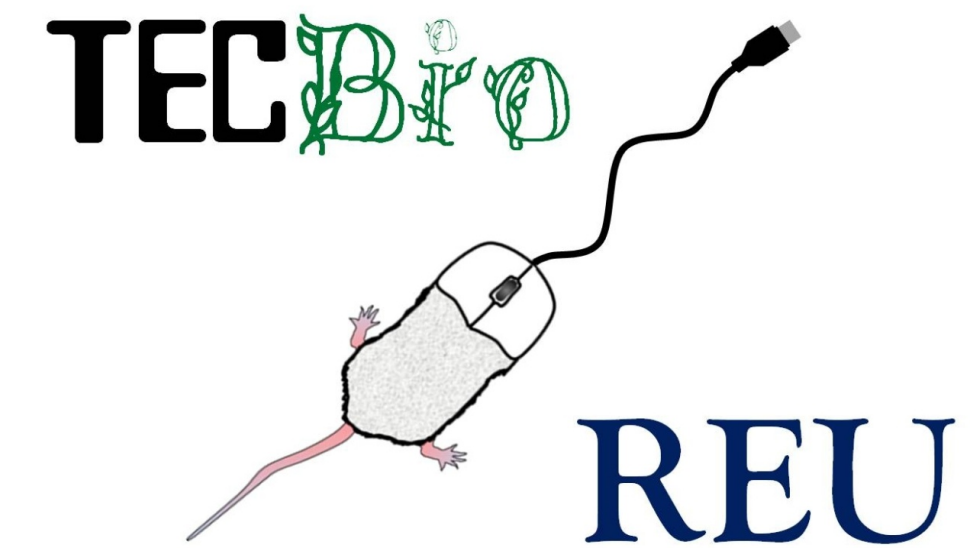
Alec Helbling

Tamar Skaist




Christopher Dunstan

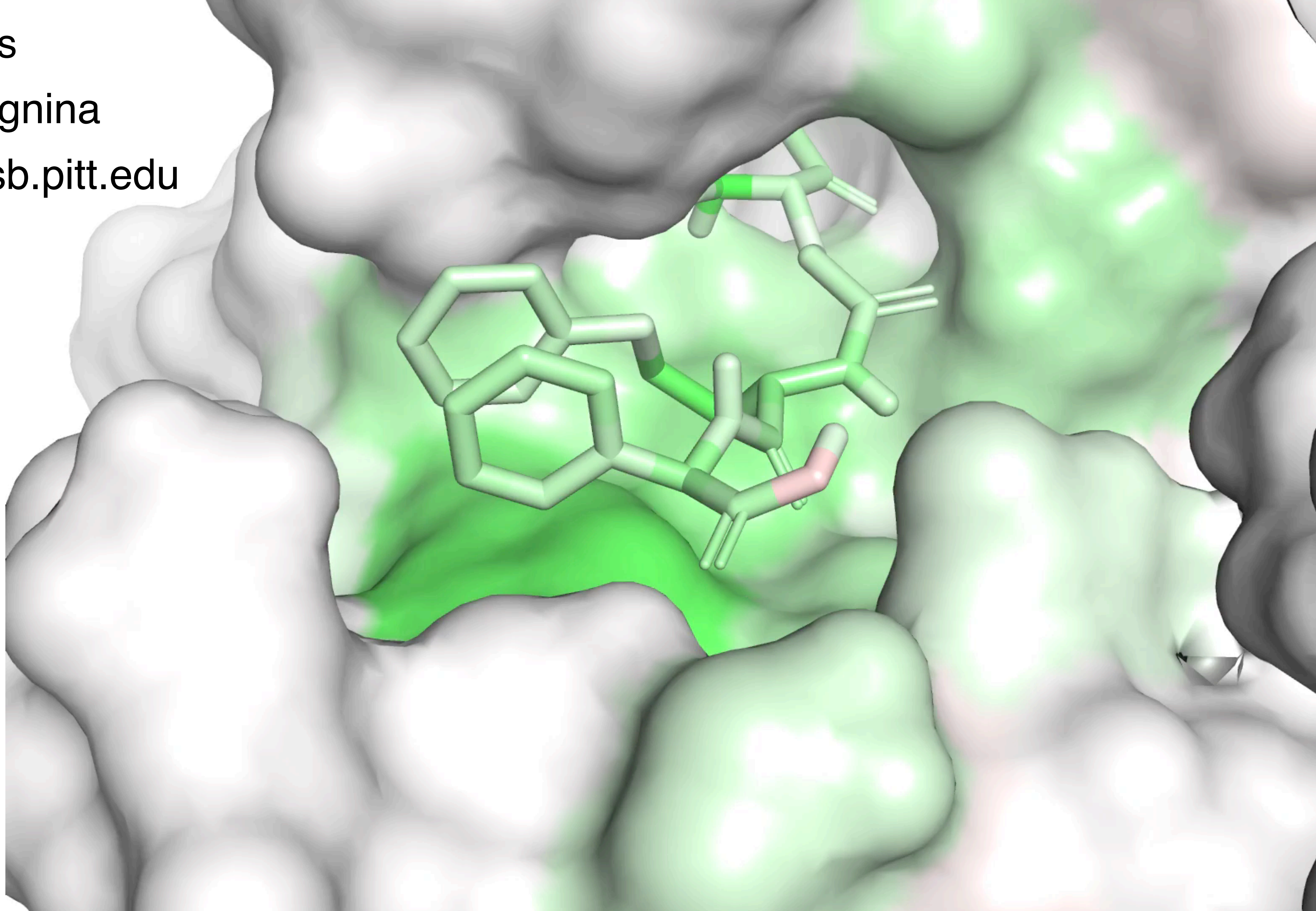





Department of
Computational and
Systems Biology



National Institute of
General Medical Sciences
R01GM108340

 @david_koes
 github.com/gnina
 <http://bits.csb.pitt.edu>



 @david_koes
 github.com/gnina
 <http://bits.csb.pitt.edu>

