

Docking Deeply: Molecular Docking with Deep Learning Potentials.

David Koes

 @david_koes

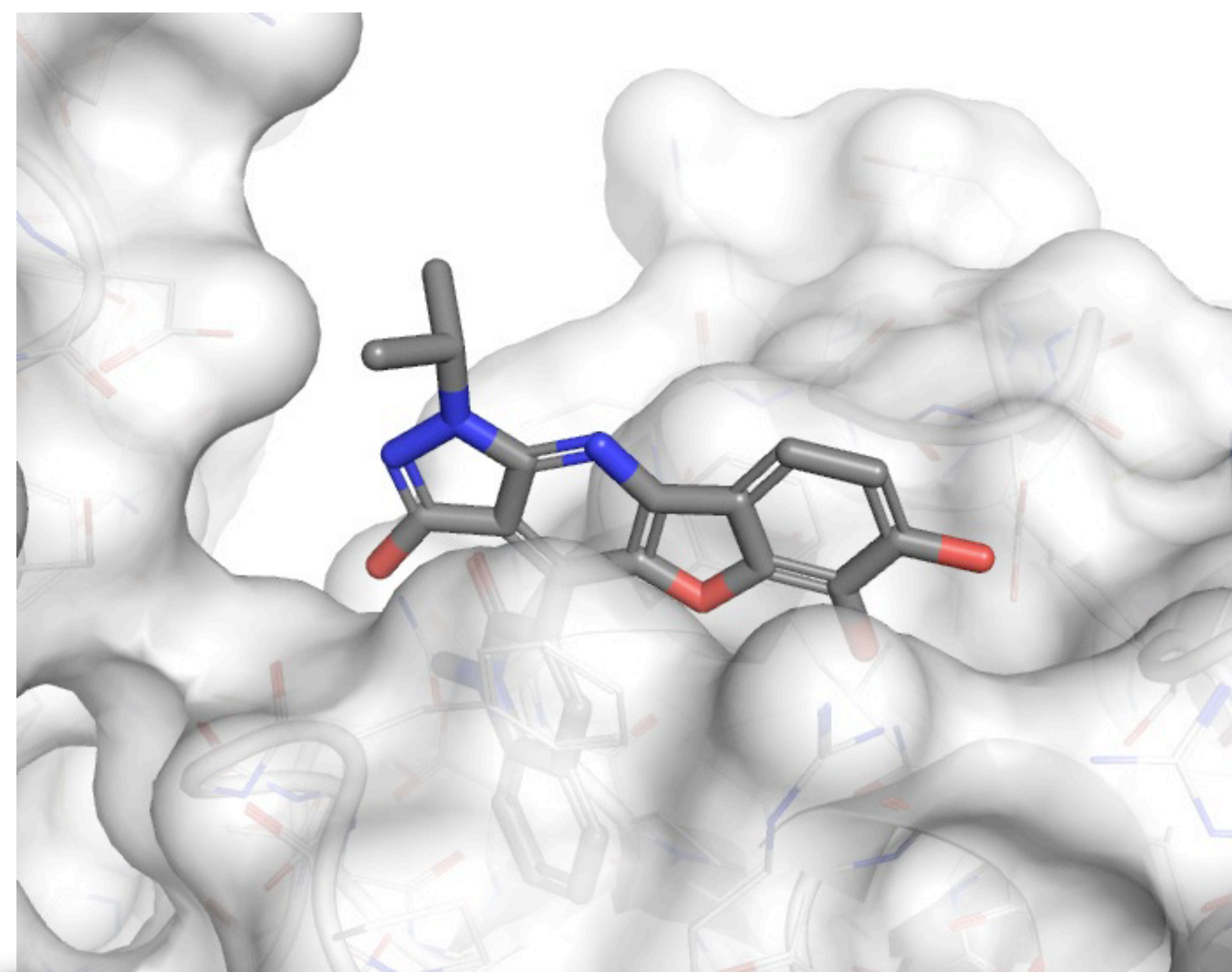


2018 Workshop on Free Energy Methods, Kinetics and
Markov State Models in Drug Design
Cambridge, MA
May 16, 2018

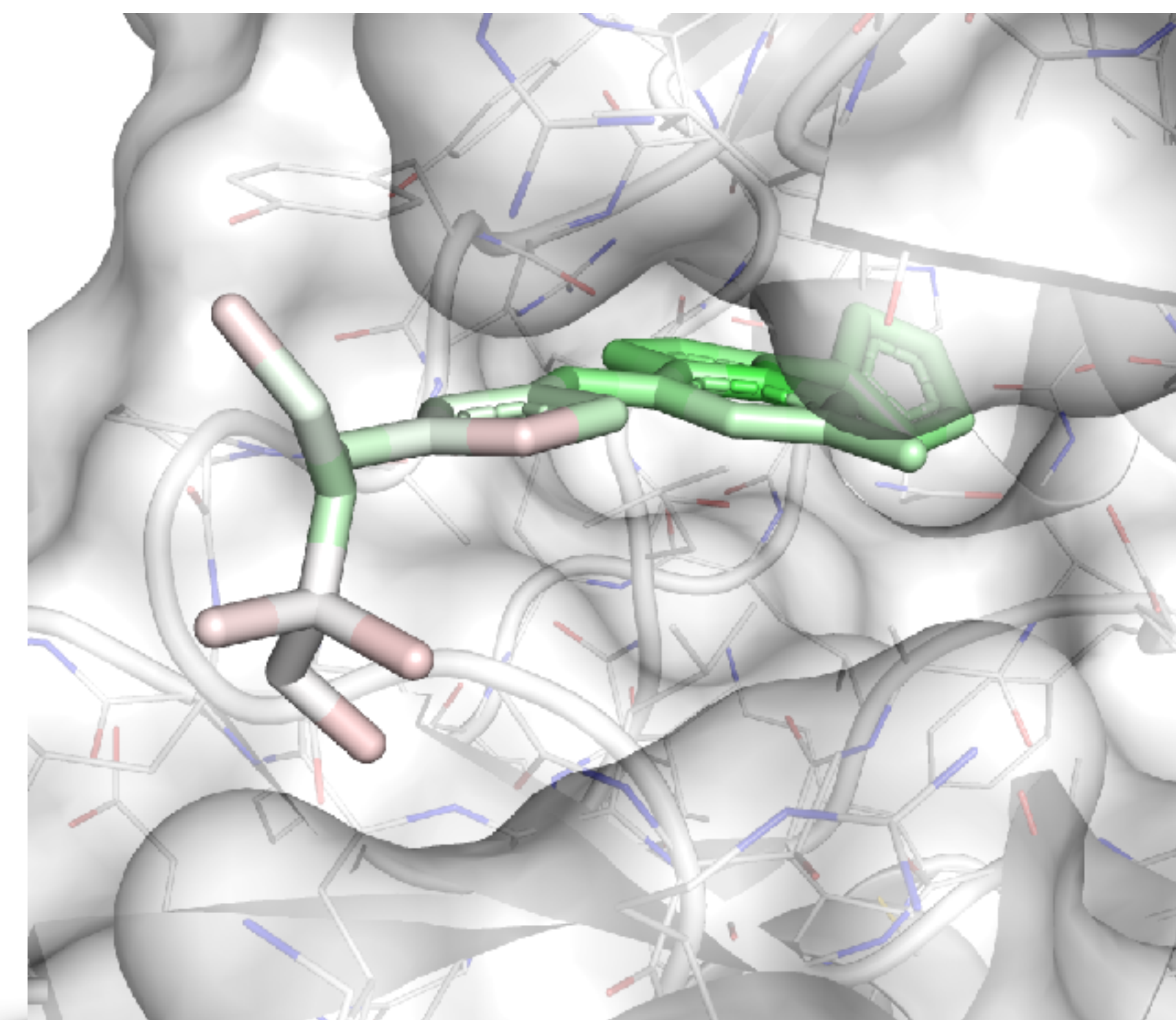


Structure Based Drug Design

Virtual Screening



Lead Optimization



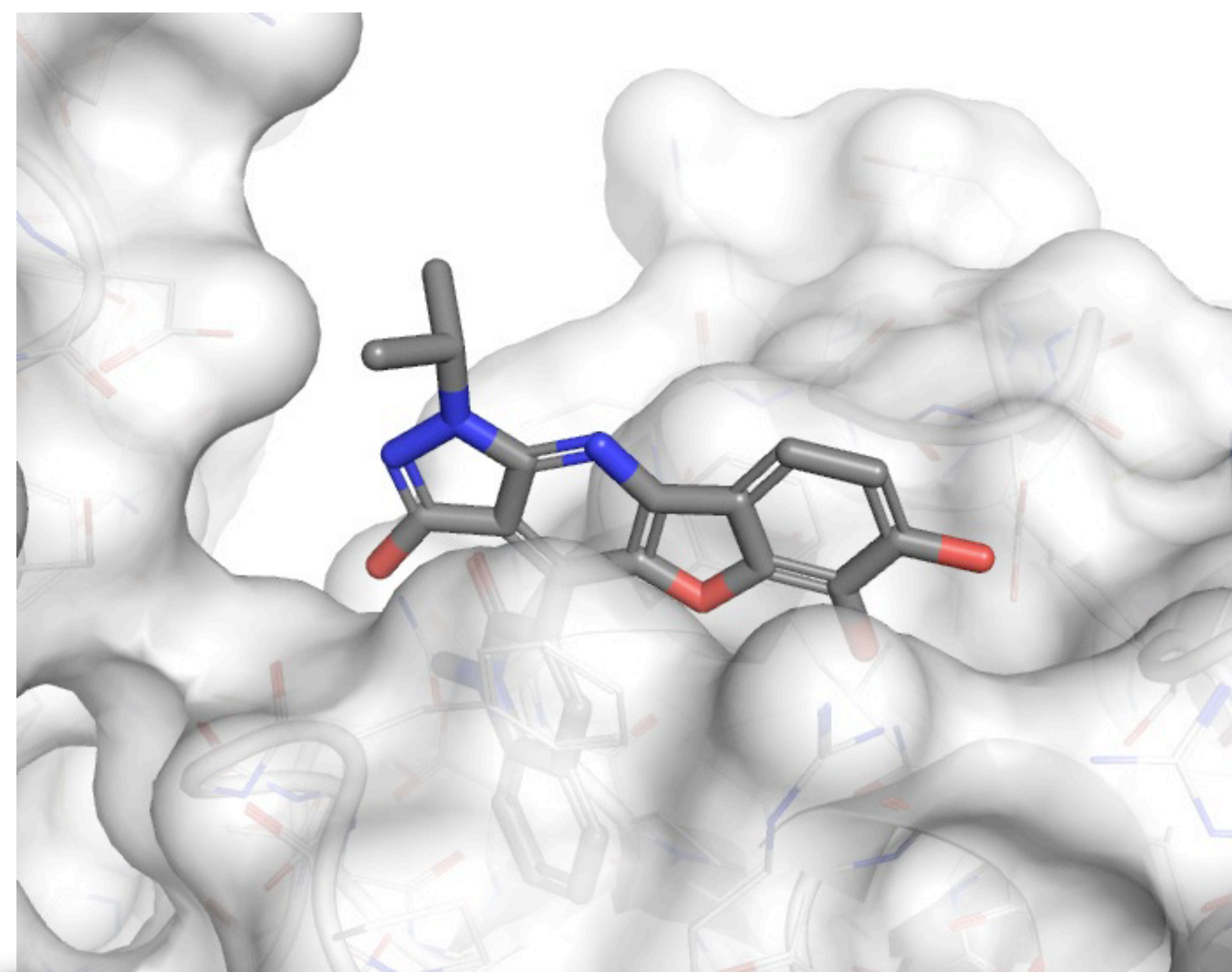
Pose Prediction

Binding Discrimination

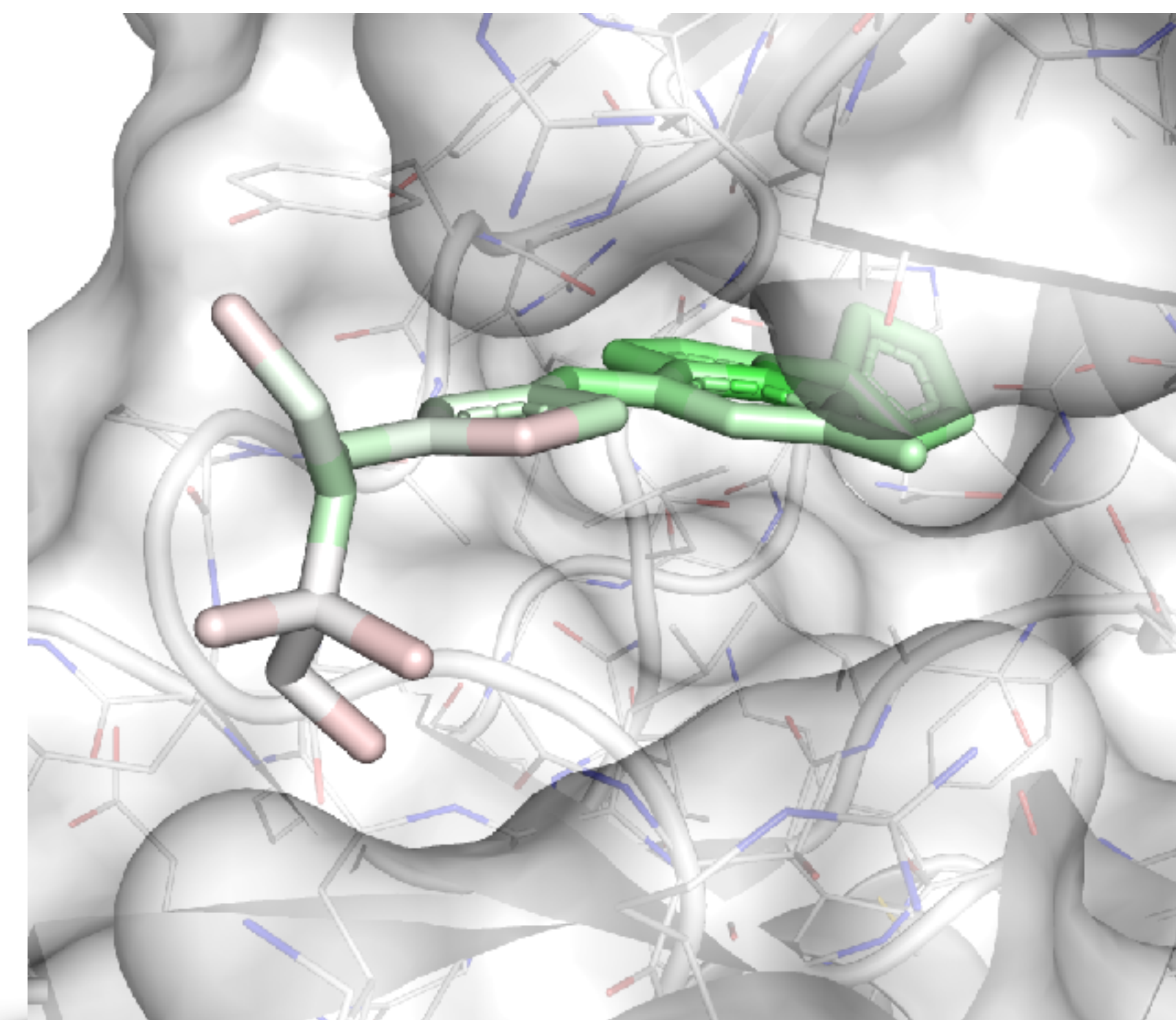
Affinity Prediction

Structure Based Drug Design

Virtual Screening



Lead Optimization

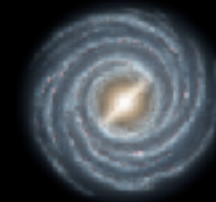


Pose Prediction

Binding Discrimination

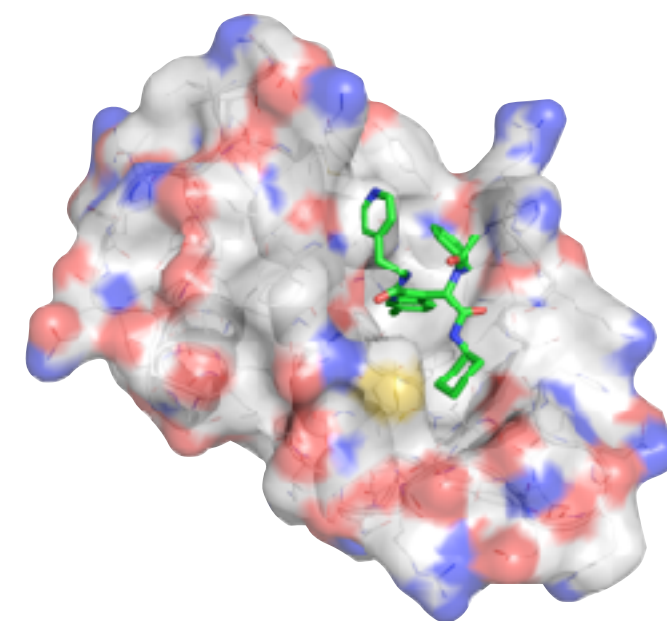
Affinity Prediction

Purchasable



Accessible

Drug Discovery Funnel

**Matching****Scoring****Dynamics**

Pharmit Search Engine

pharmit.csb.pitt.edu/search.html?SESSION=examples/4pps.json

Search MolPort
Pharmacophore Search -> Shape Filter
Load Receptor... Load Features...

Pharmacophore

- ☒ **HydrogenDonor**
(9.53,3.92,35.82) Radius 0.5
- ☒ **HydrogenAcceptor**
(9.53,3.92,35.82) Radius 0.5
- ☒ **HydrogenAcceptor**
(20.0,4.36,33.13) Radius 0.5
- ☒ **Hydrophobic**
(12.17,4.27,35.2) Radius 1.0
- ☒ **Hydrophobic**
(18.72,4.96,35.4) Radius 1.0
- ☒ **Hydrophobic**
(17.88,4.44,32.8) Radius 1.0
- ☒ **Hydrophobic**
(16.24,4.83,33.93) Radius 1.0

Add Sort

Shape

- ☐ Inclusive Shape
- ☐ Exclusive Shape

Filters

- ☐ Hit Reduction
- ☐ Hit Screening

Load Session... Save Session...

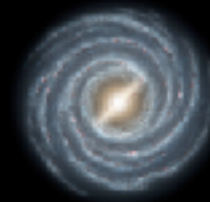
Pharmacophore Results

Name	RMSD	Mass	RBnds
MolPort-002-911-158	0.113	395	1
MolPort-000-705-595	0.124	330	0
MolPort-035-395-591	0.125	607	15
MolPort-002-509-936	0.132	314	0
MolPort-003-847-099	0.134	275	0
MolPort-002-741-818	0.147	351	0
MolPort-002-515-415	0.148	330	0
MolPort-009-018-993	0.150	300	1
MolPort-003-892-015	0.157	288	0
MolPort-003-941-332	0.164	272	0
MolPort-006-318-980	0.164	272	0
MolPort-000-720-575	0.164	272	0
MolPort-000-725-407	0.165	296	0
MolPort-039-348-092	0.169	378	1
MolPort-002-509-704	0.169	312	1
MolPort-002-520-588	0.170	375	3
MolPort-002-506-898	0.172	288	0
MolPort-006-069-030	0.173	607	15

Showing 1 to 18 of 1,335 hits
Previous 1 2 3 4 5 Next
Query took 2.235 seconds
Minimize Save...

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

Purchasable



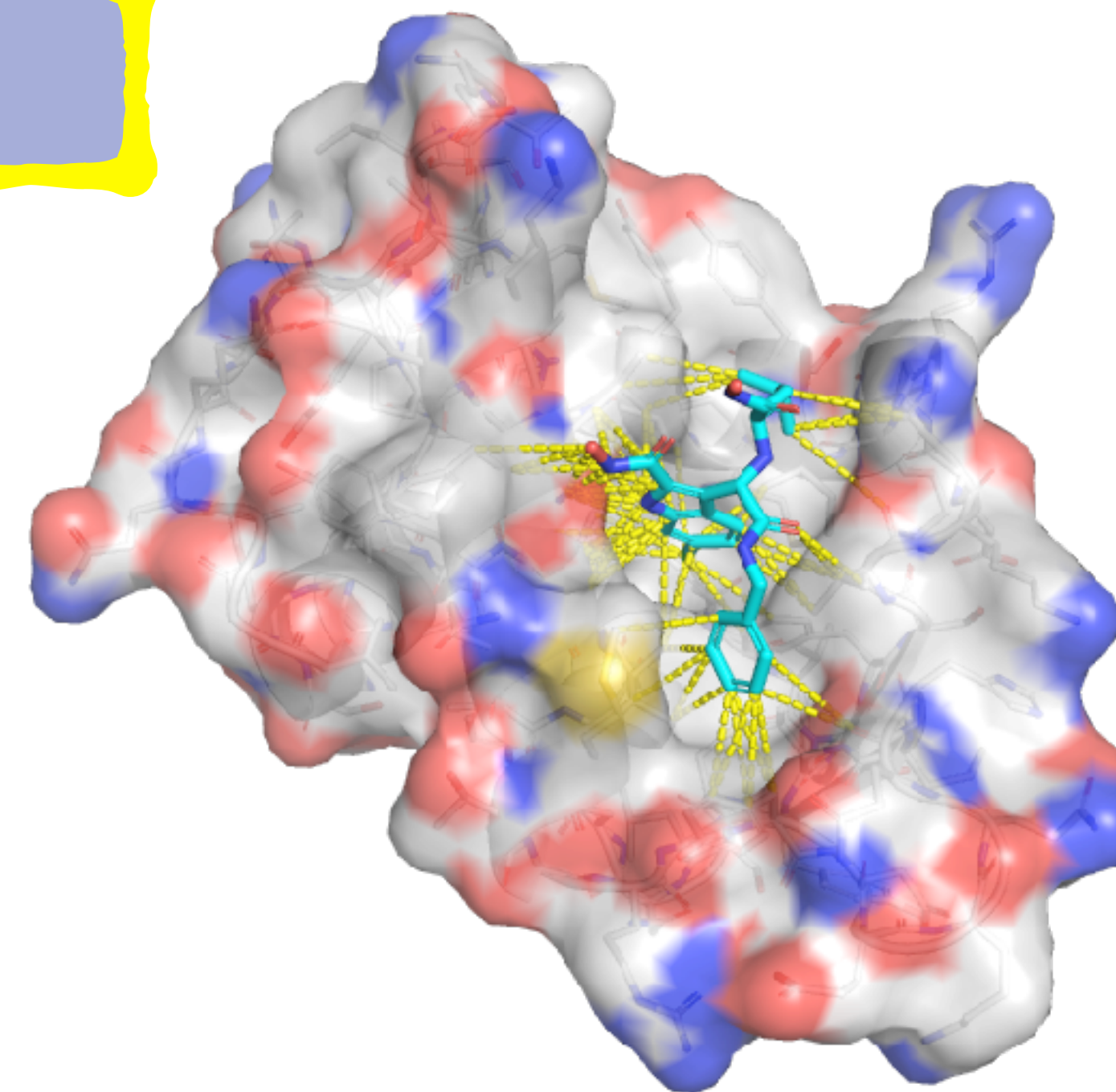
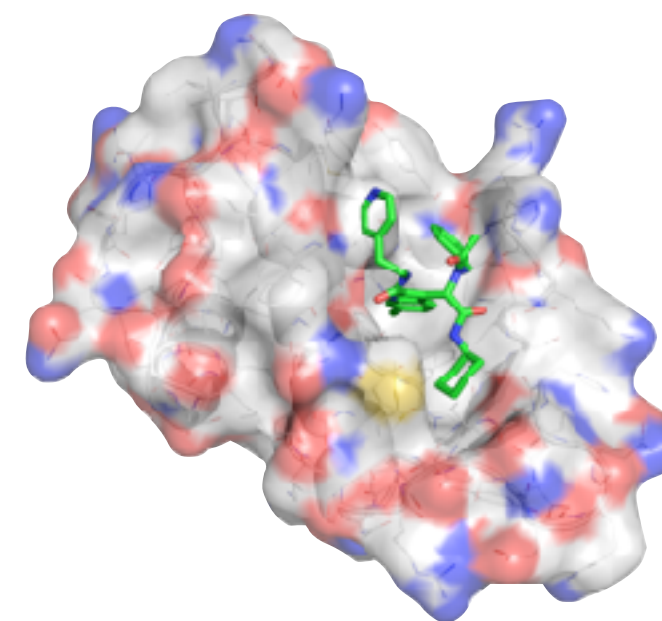
Accessible

Drug Discovery Funnel


Matching


Scoring

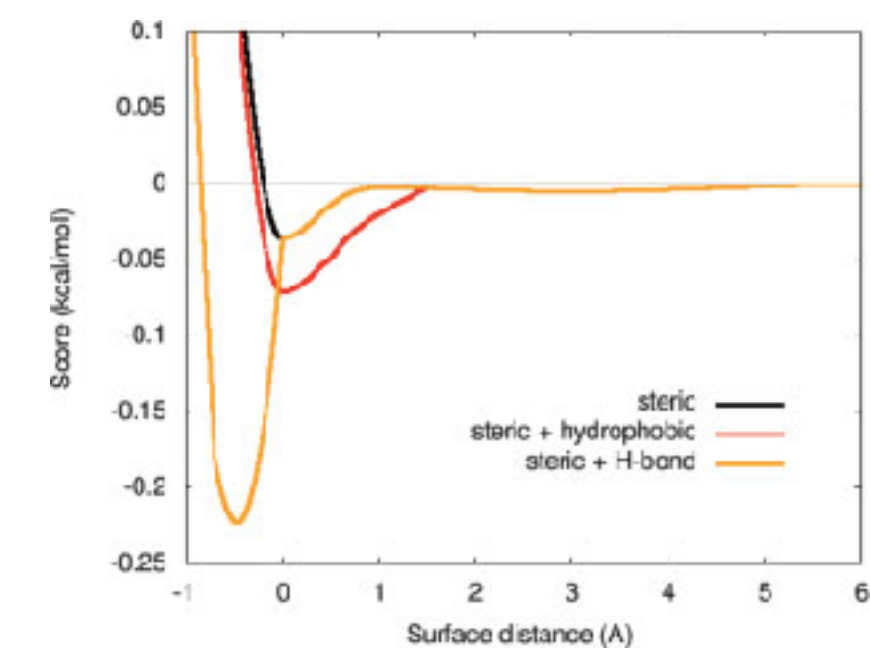

Dynamics



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

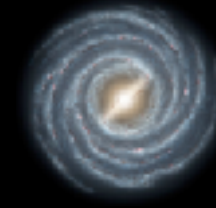
$$\text{hydrophobic}(d) = \begin{cases} w_{\text{hydrophobic}} & d < 0.5 \\ 0 & d > 1.5 \\ w_{\text{hydrophobic}}(1.5 - d) & \text{otherwise} \end{cases}$$

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



O. Trott, A. J. Olson, AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization and multithreading, *Journal of Computational Chemistry* 31 (2010) 455-461

Purchasable



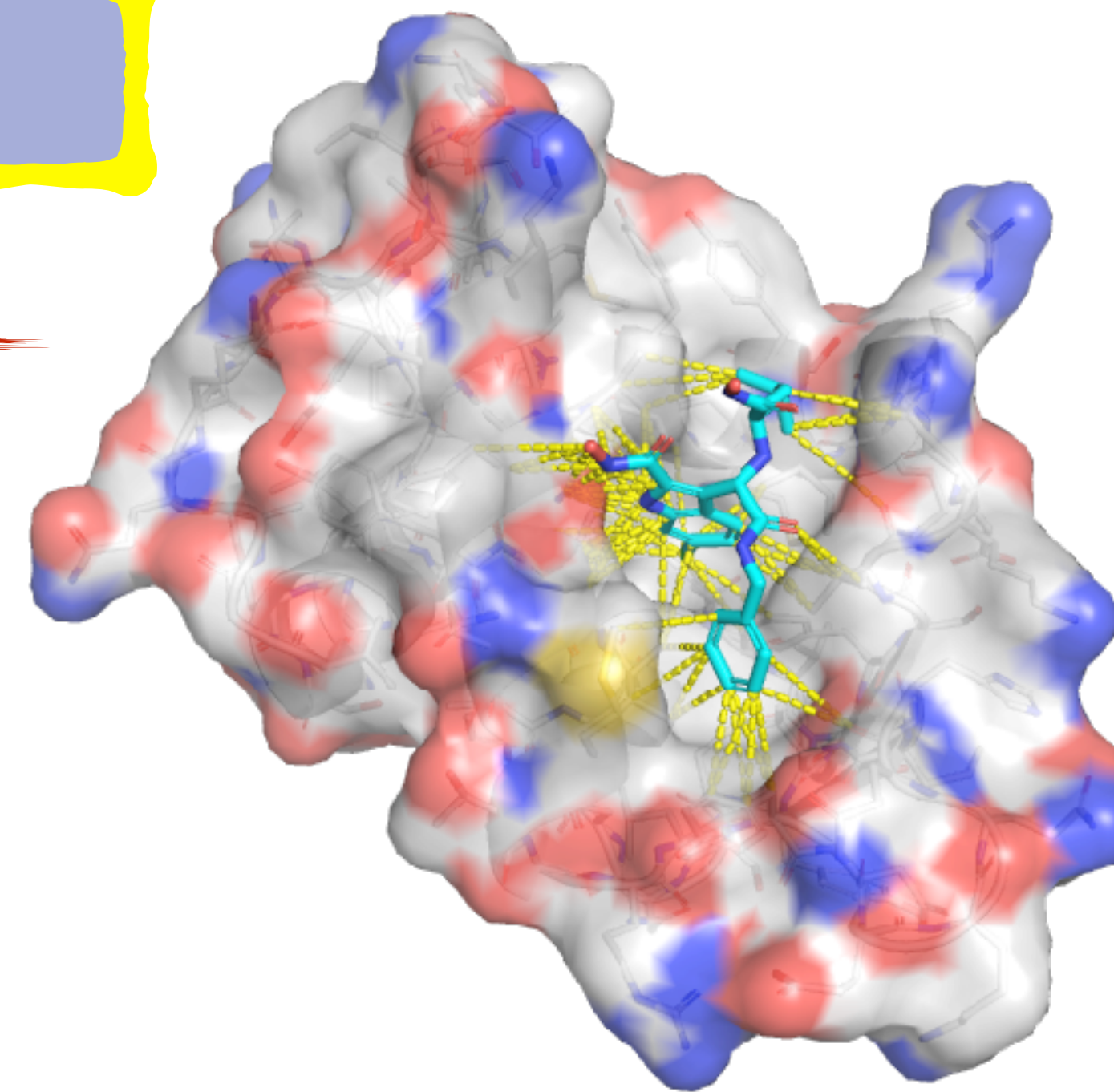
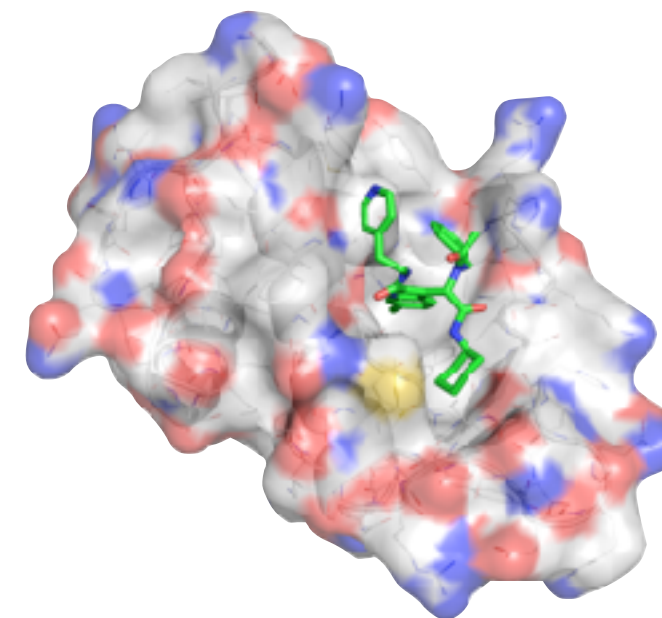
Accessible

Drug Discovery Funnel


Matching


Scoring

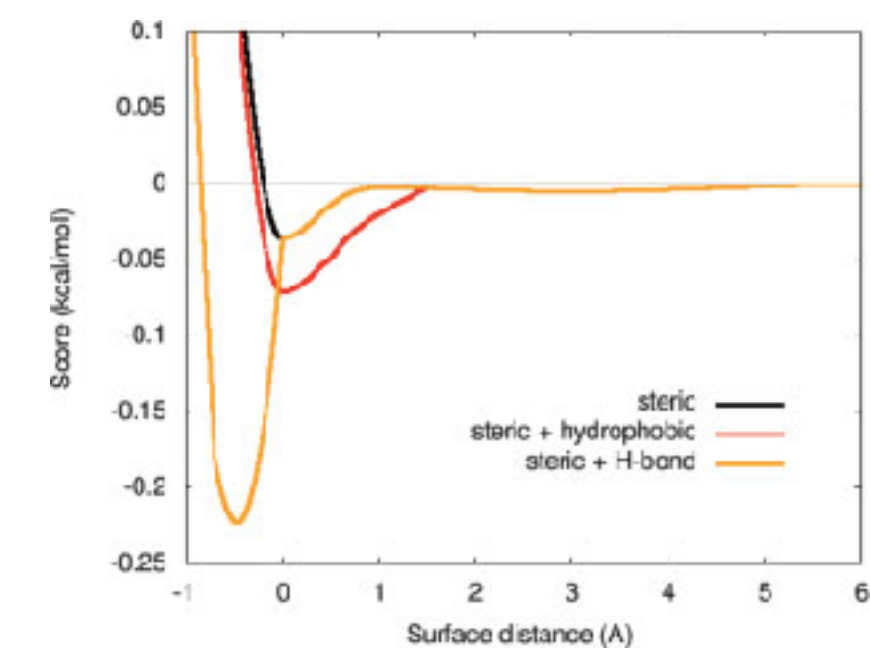

~~**Dynamics**~~



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

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

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



O. Trott, A. J. Olson, AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization and multithreading, *Journal of Computational Chemistry* 31 (2010) 455-461

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”

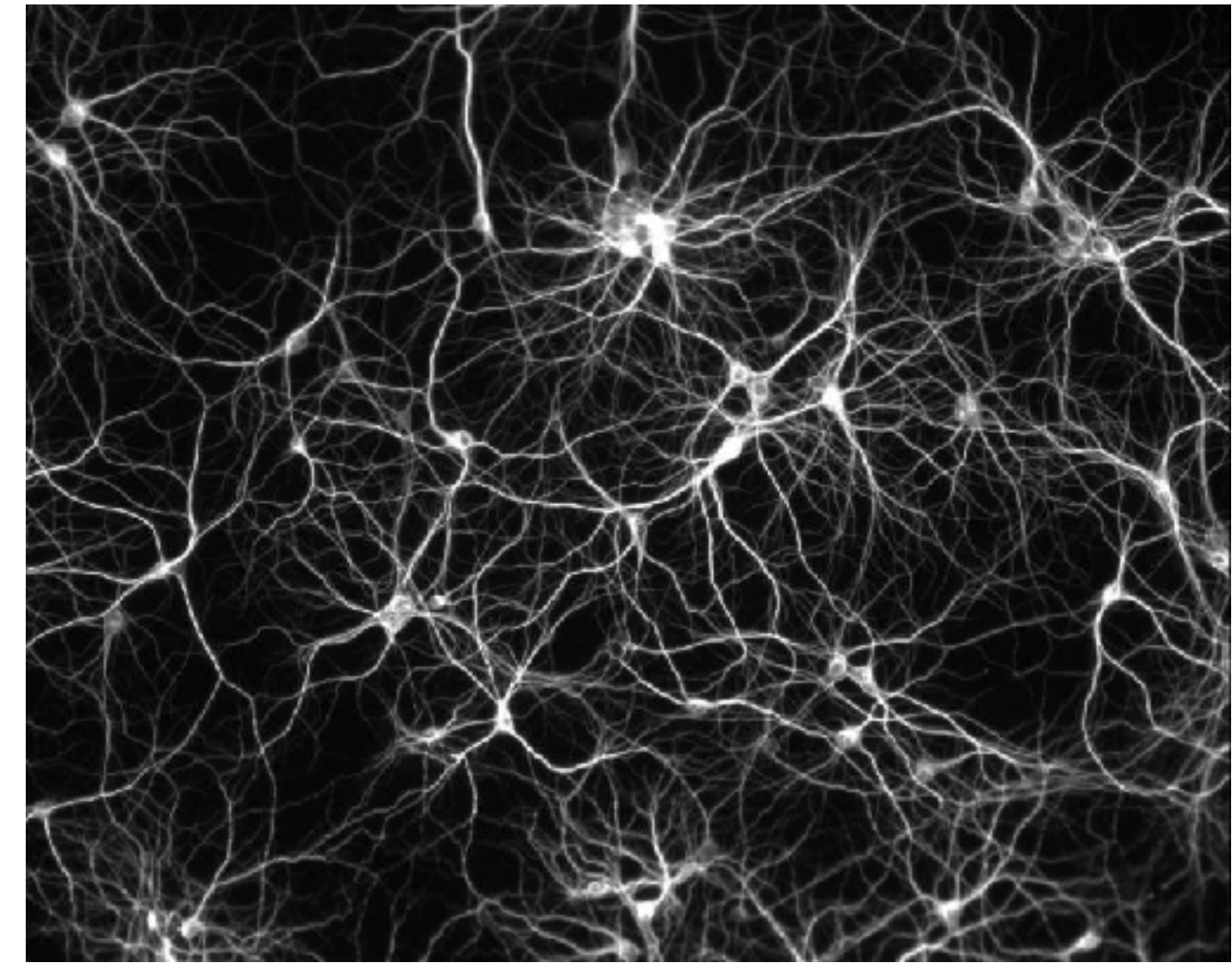
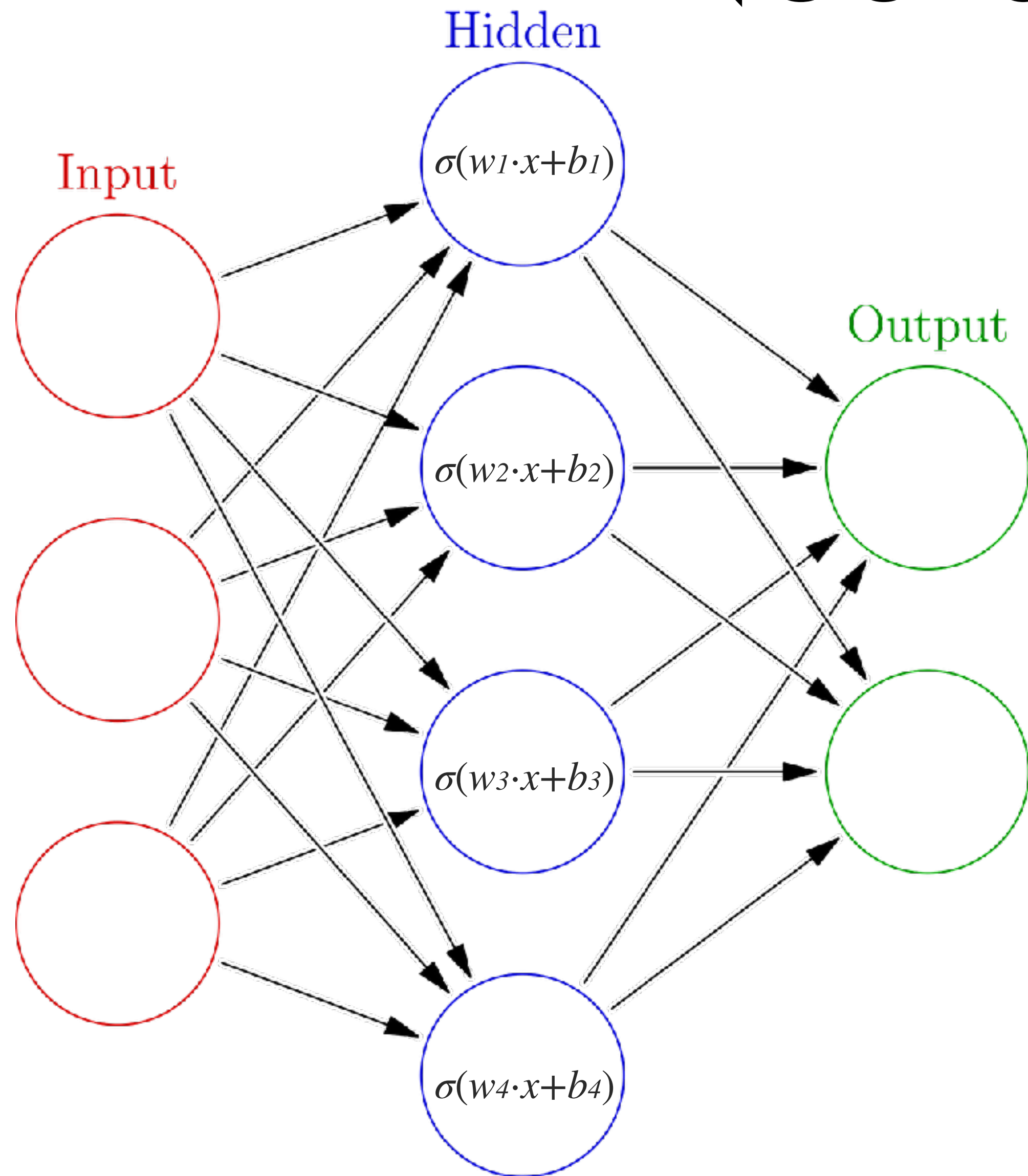
- 231,655,275 bioactivities in PubChem
- 125,526 structures in the PDB
- 16,179 annotated complexes in PDBbind



Machine Learning

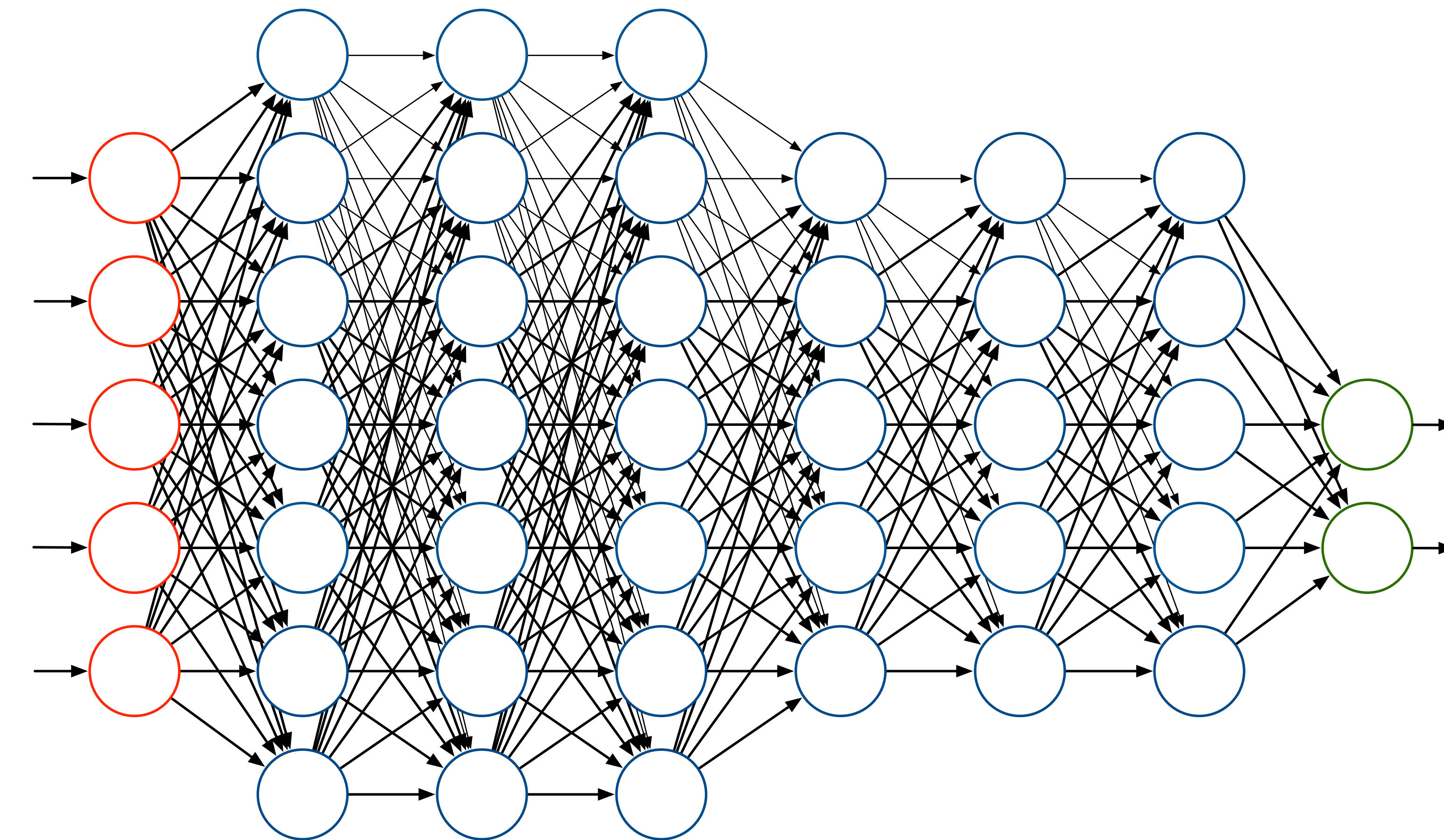


Neural Networks

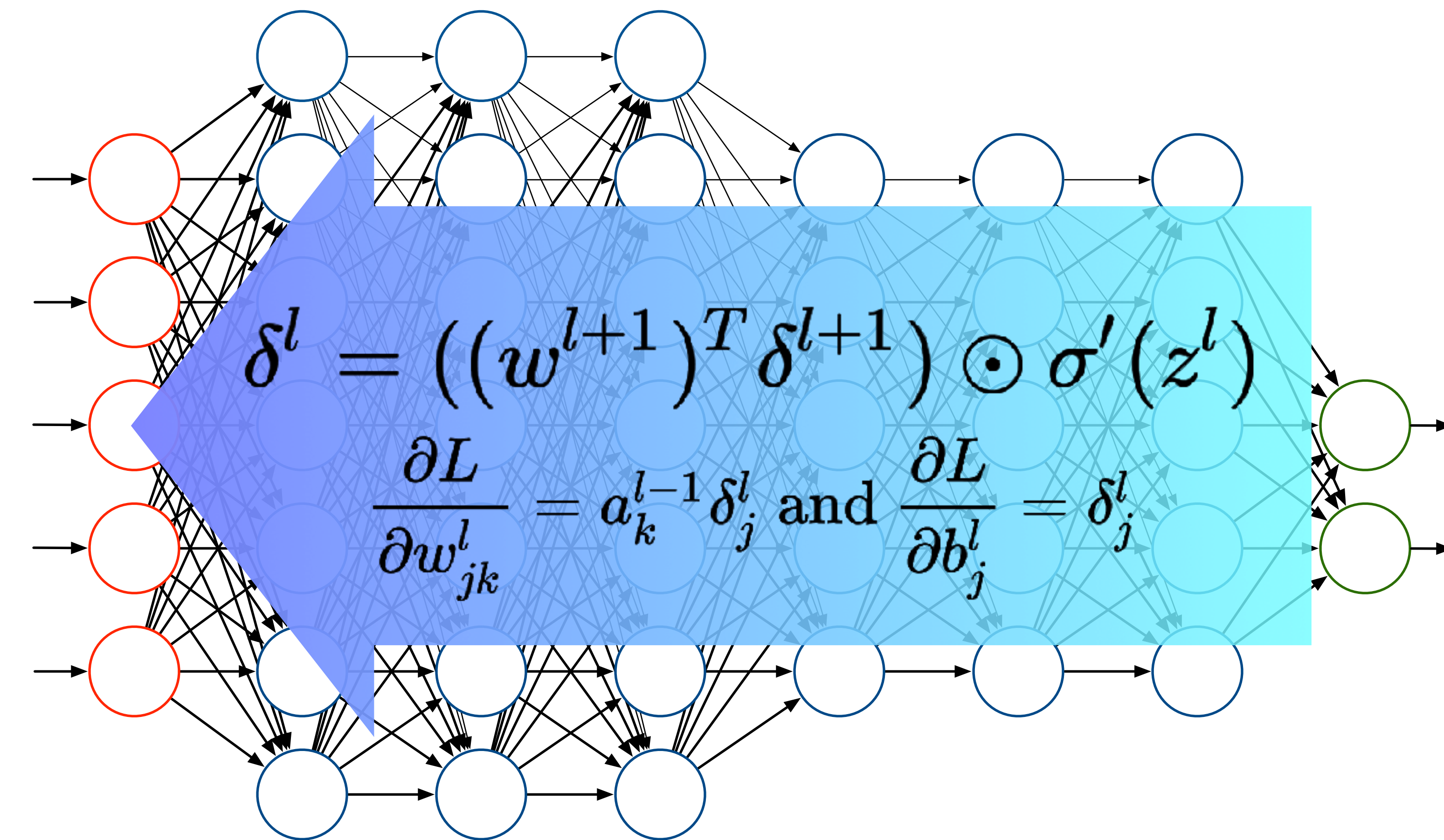


The **universal approximation theorem** states that, under reasonable assumptions, a feedforward **neural network** with a finite number of nodes **can approximate any continuous** function to within a given error over a bounded input domain.

Deep Learning



Deep Learning



Deep Learning

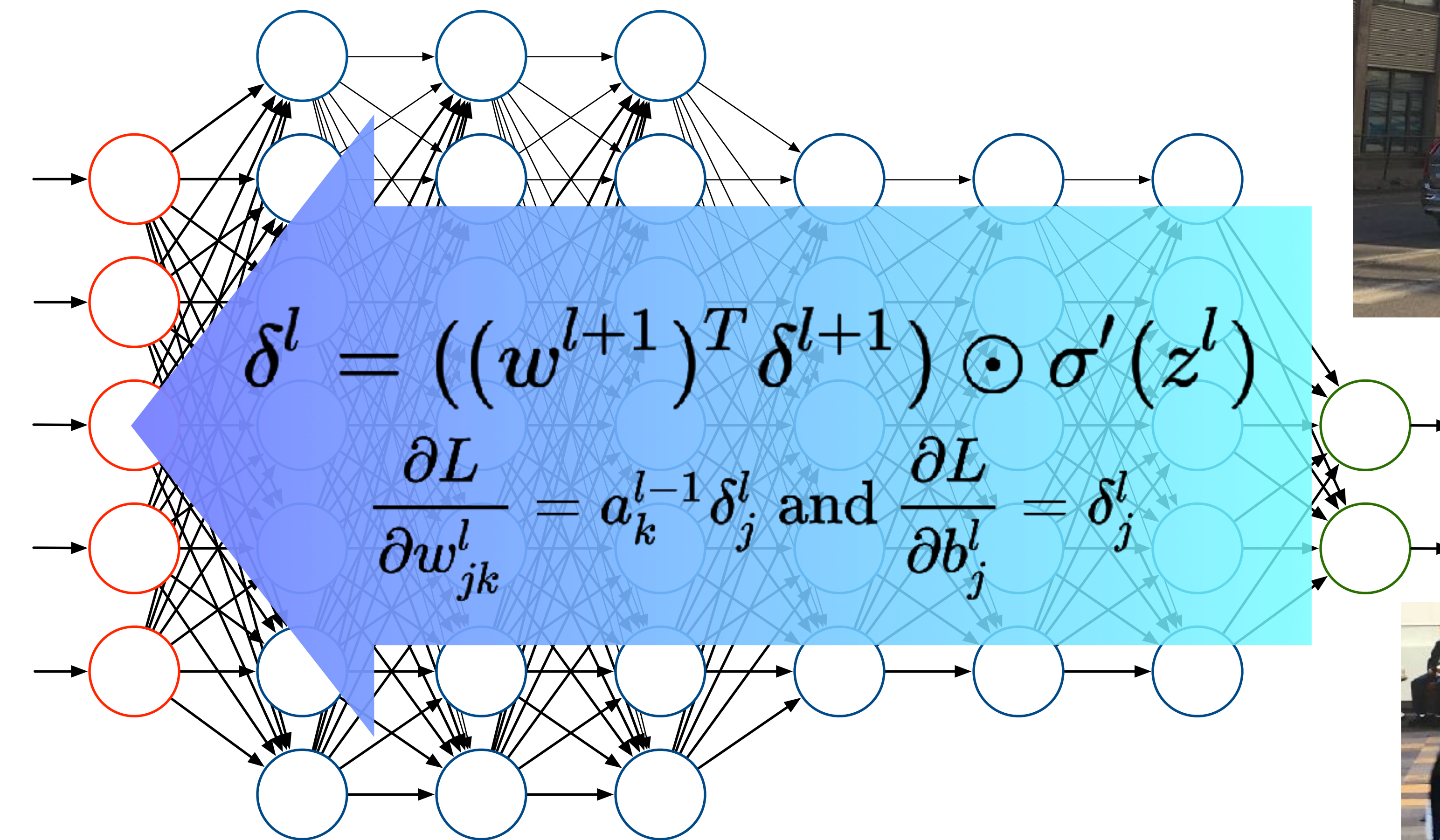


Image Recognition

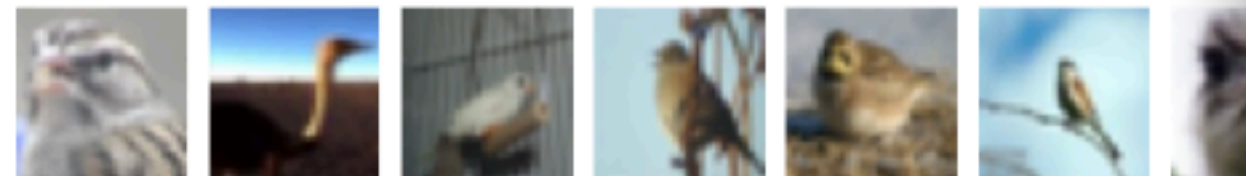
airplane



automobile



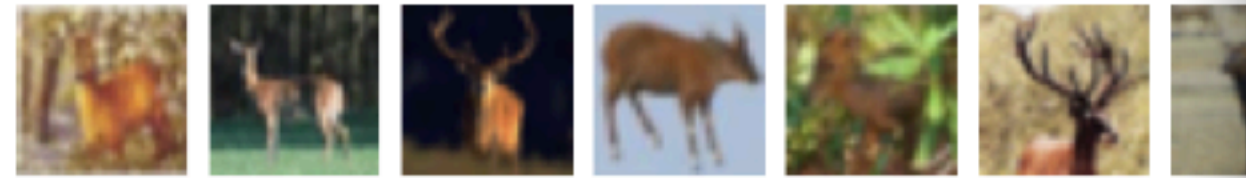
bird



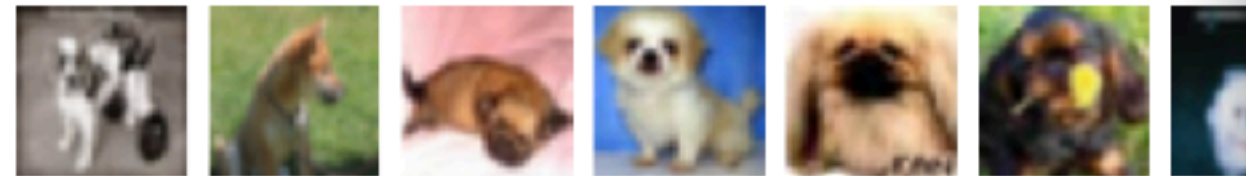
cat



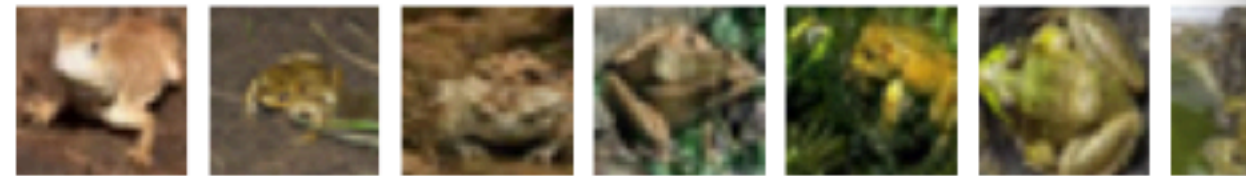
deer



dog



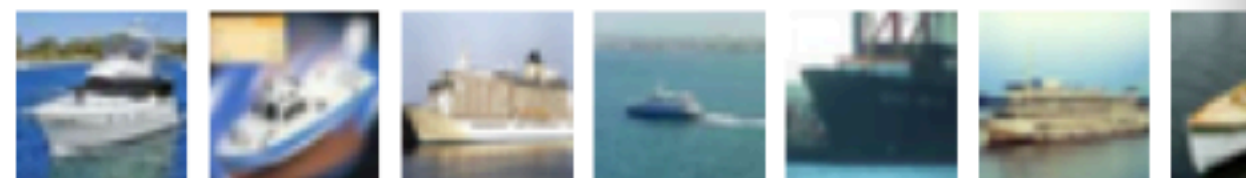
frog



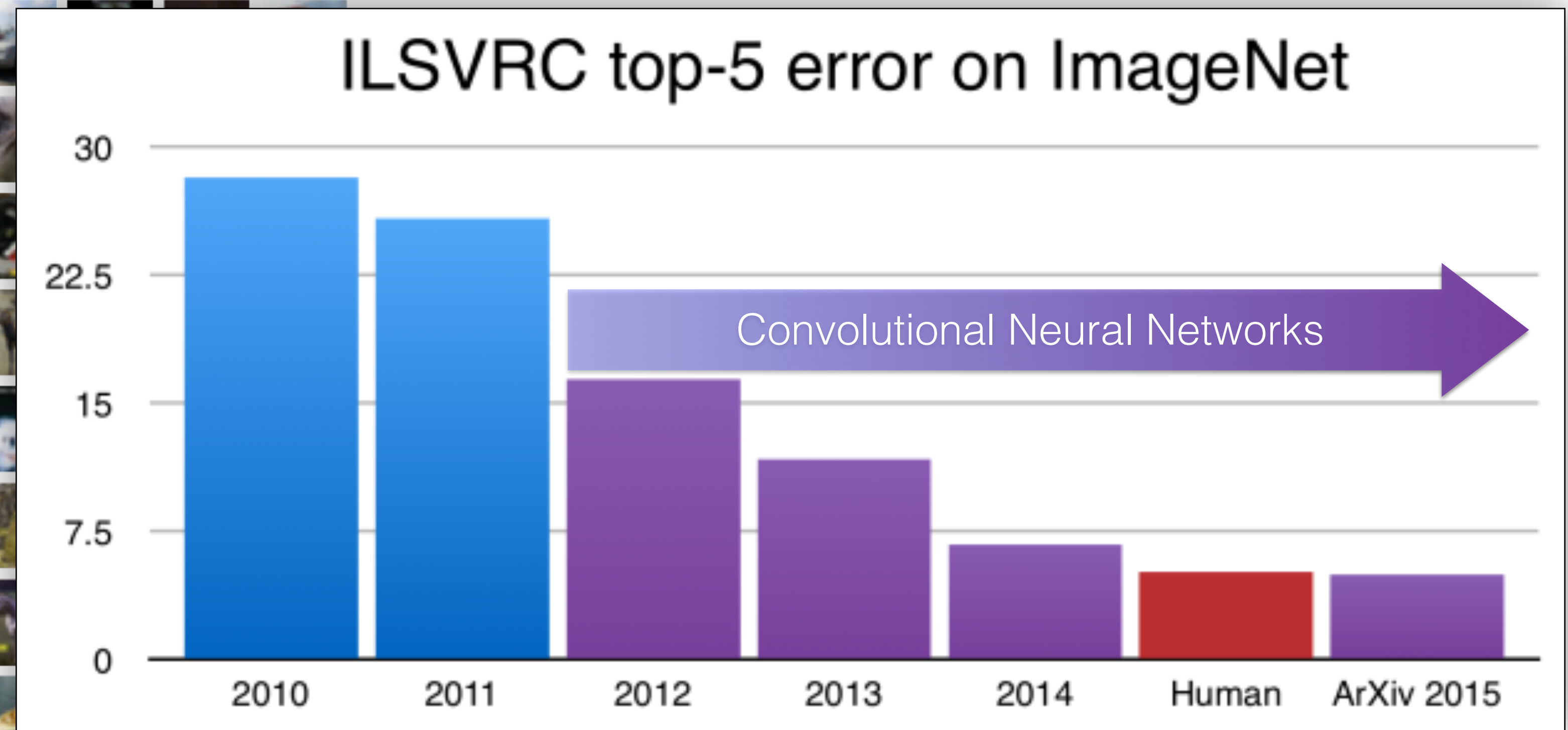
horse



ship

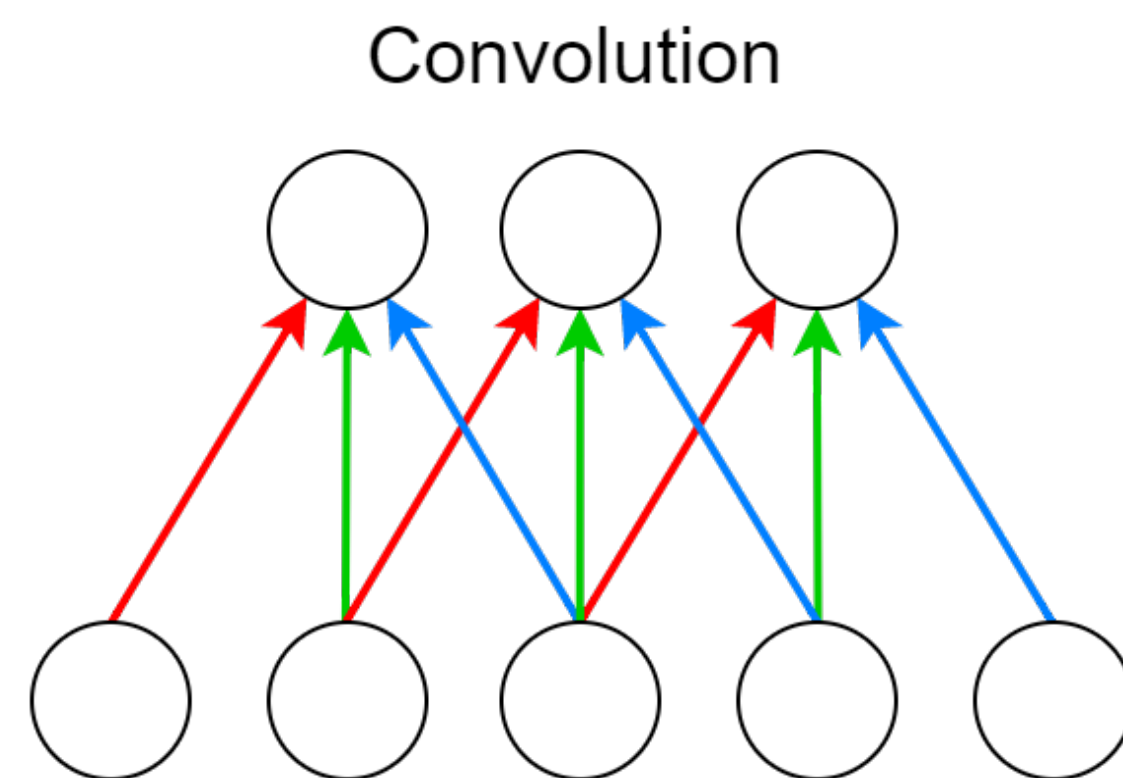
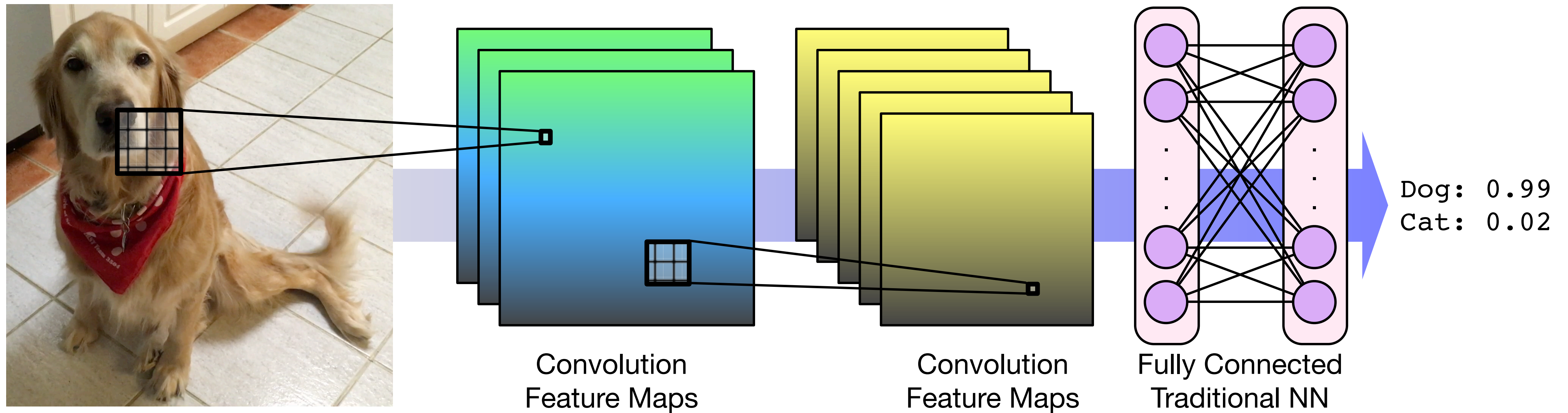


truck

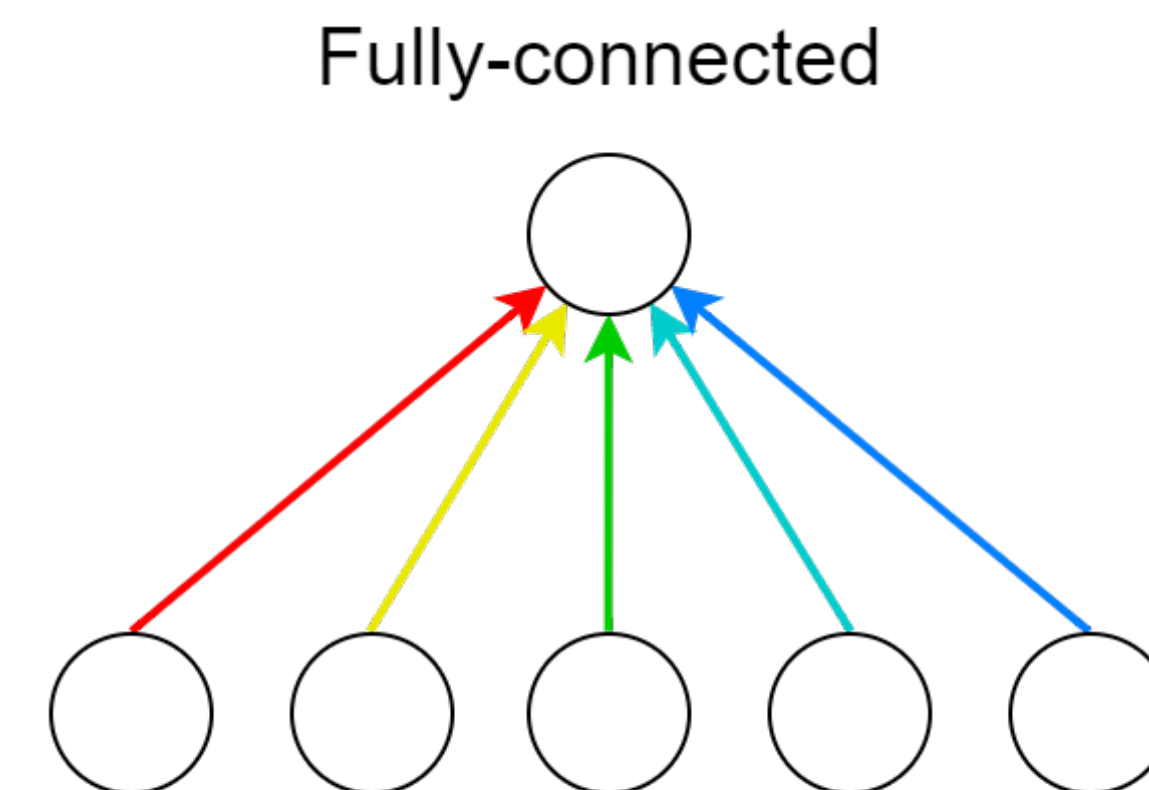


<https://devblogs.nvidia.com>

Convolutional Neural Networks

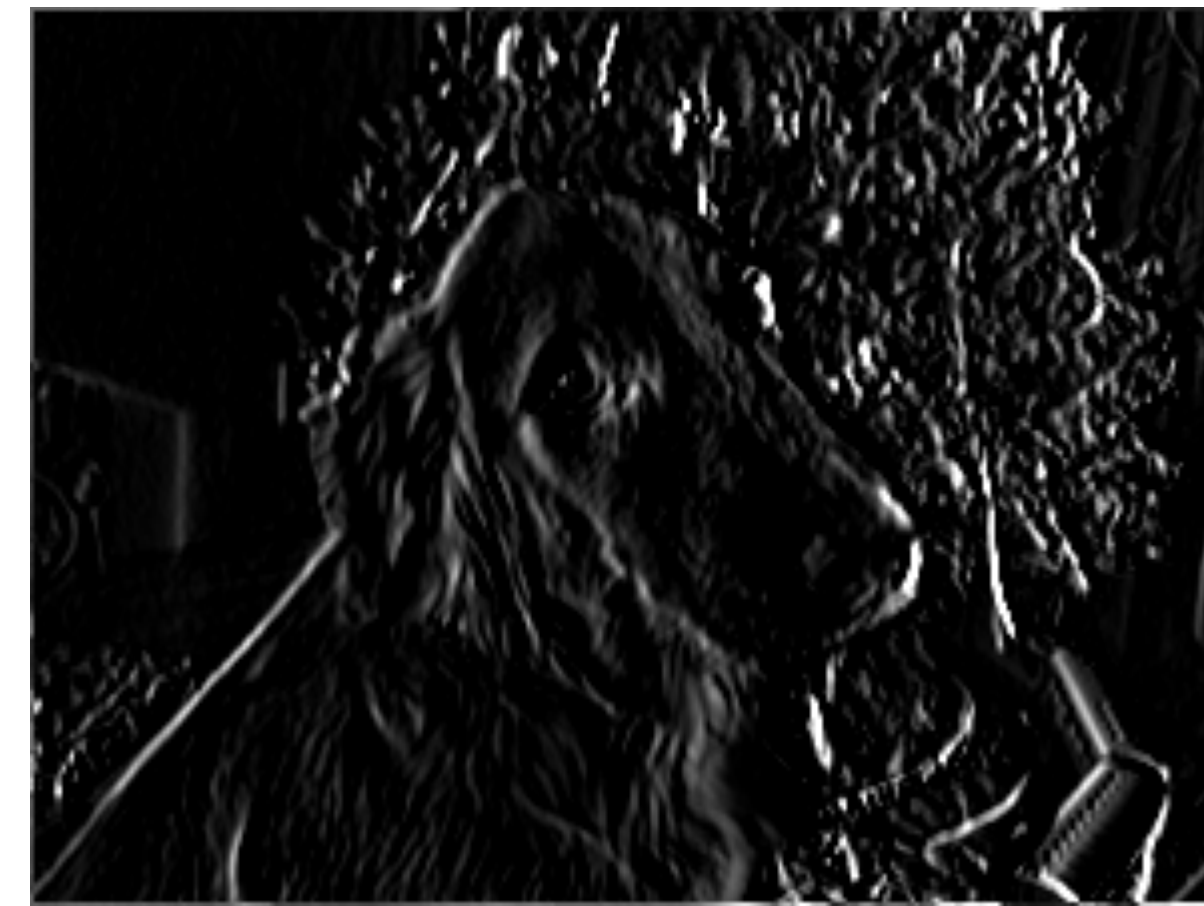
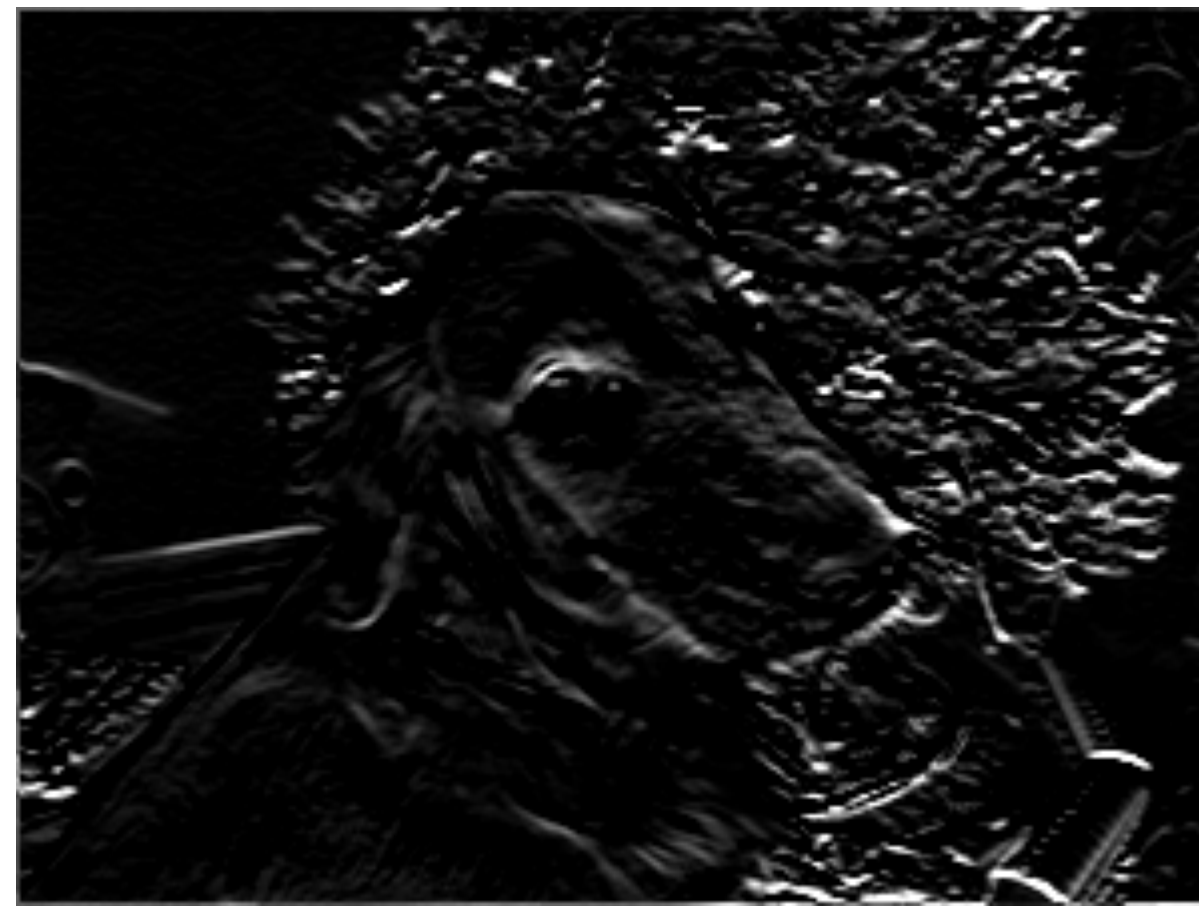


— weight 1
— weight 2
— weight 3



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

Convolutional Filters



-1	-1	-1
0	0	0
1	1	1

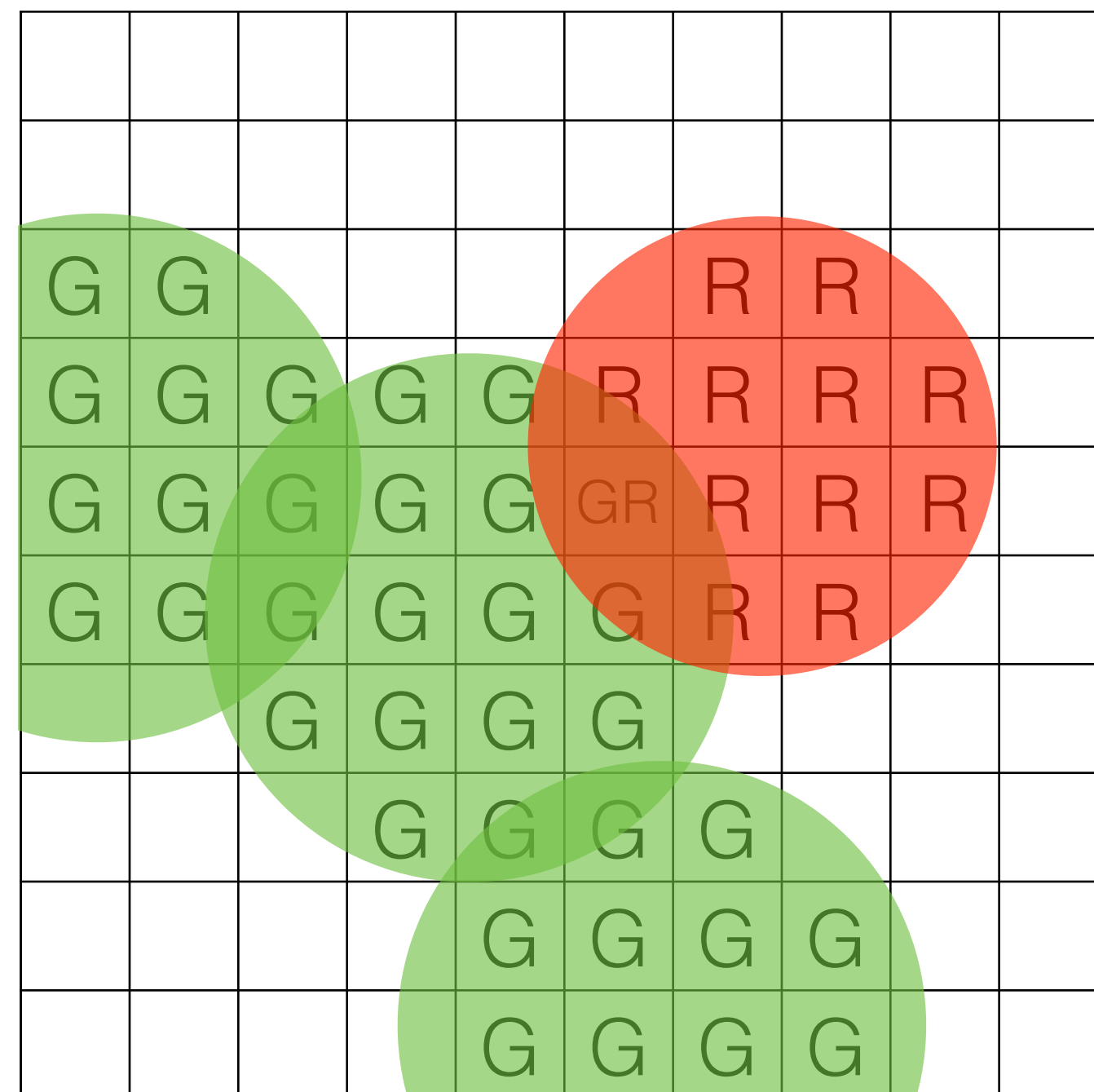
-1	0	1
-1	0	1
-1	0	1

-1	-1	-1
-1	8	-1
-1	-1	-1

CNNs for Protein-Ligand Scoring

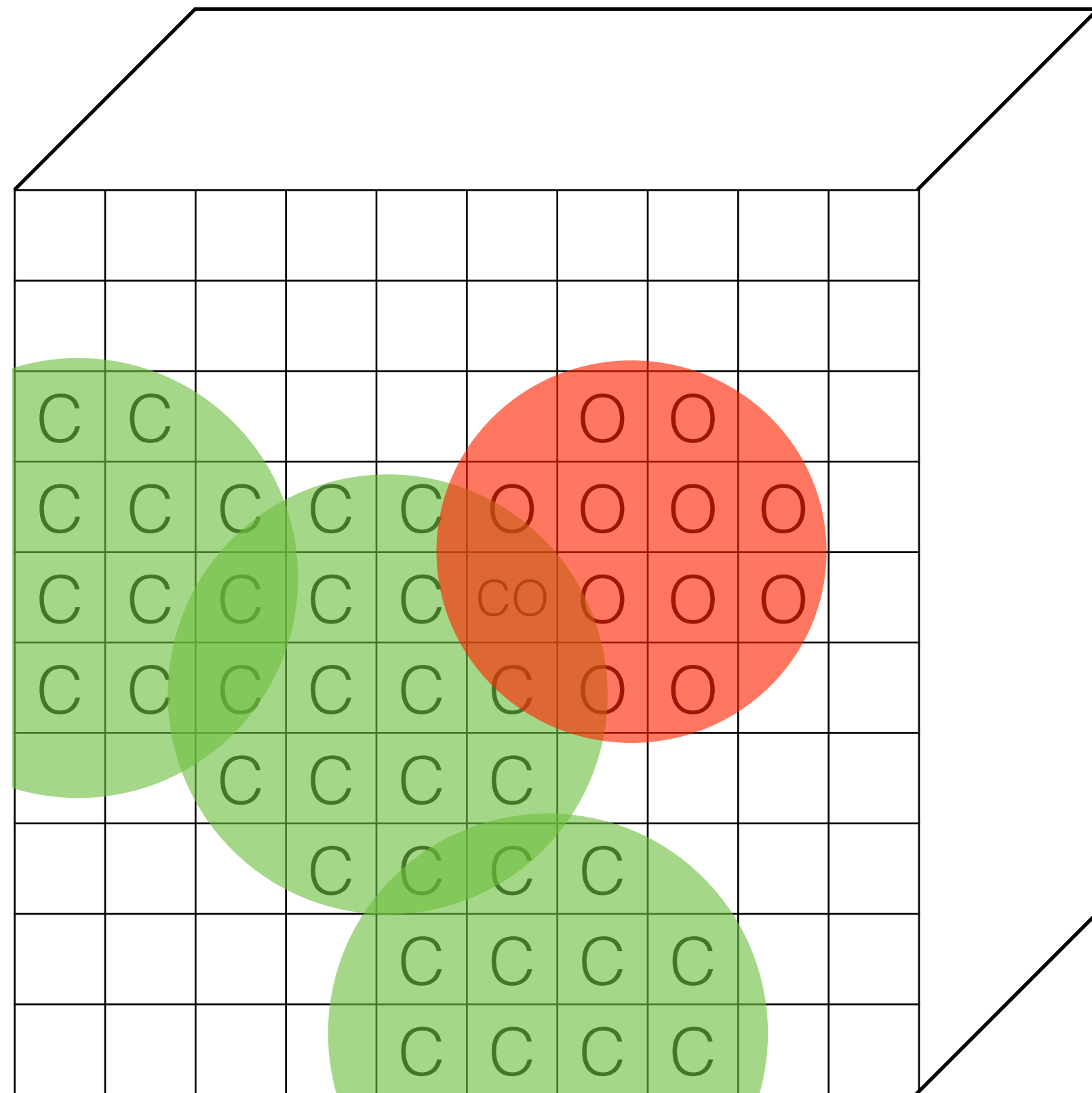


Protein-Ligand Representation



(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**.

Training Data



Pose Prediction

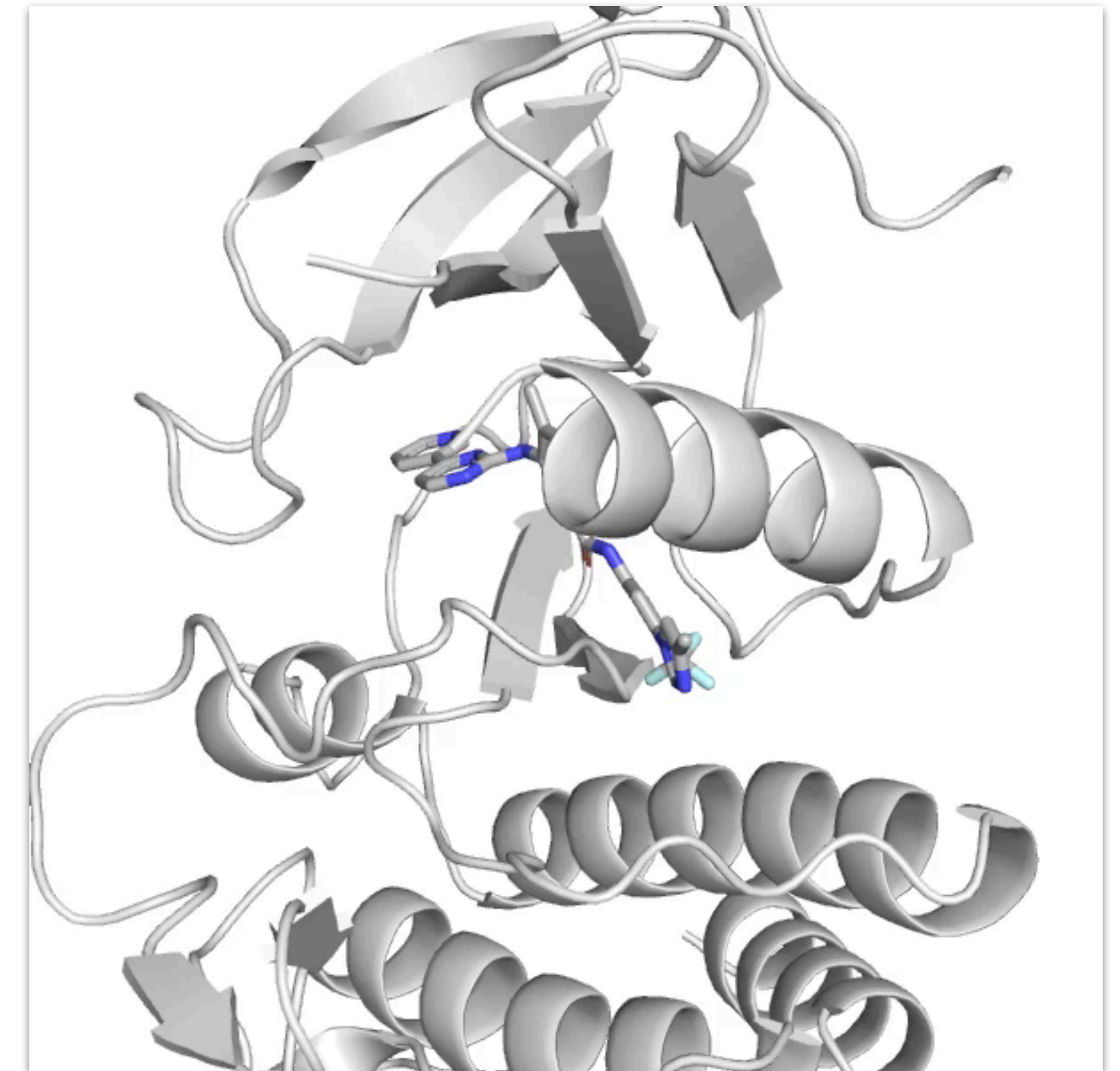
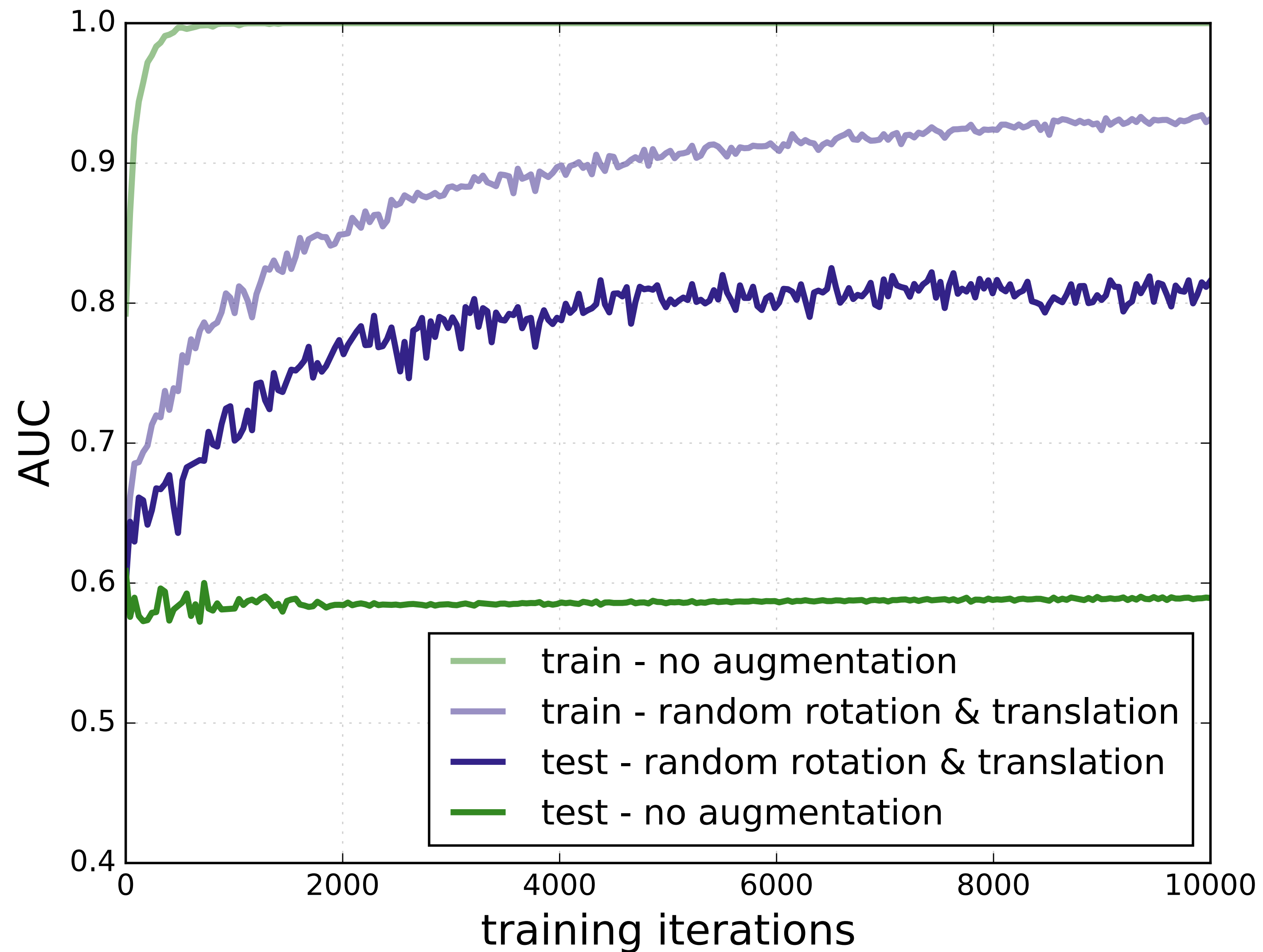
4056 protein-ligand complexes

- 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)

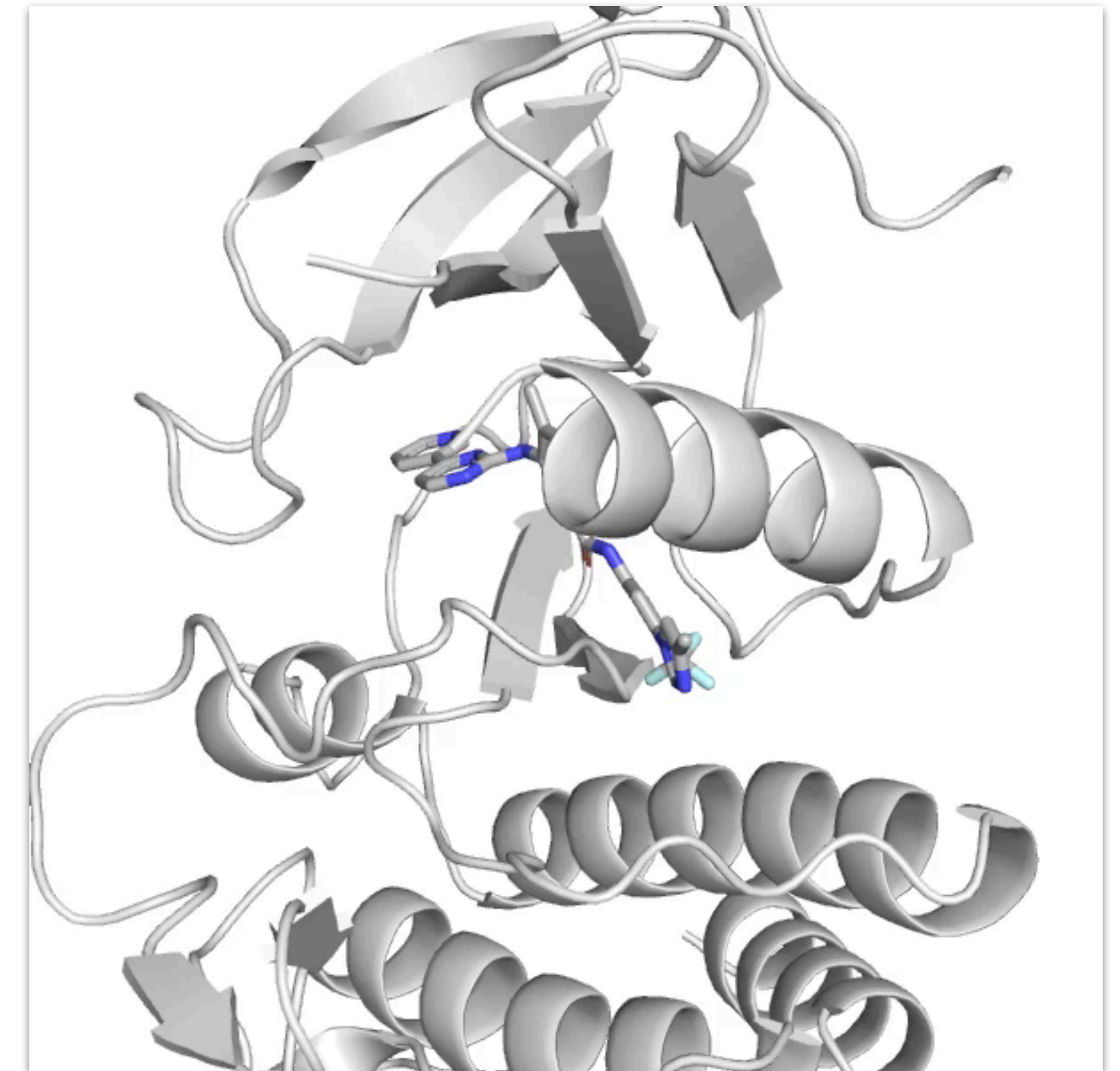
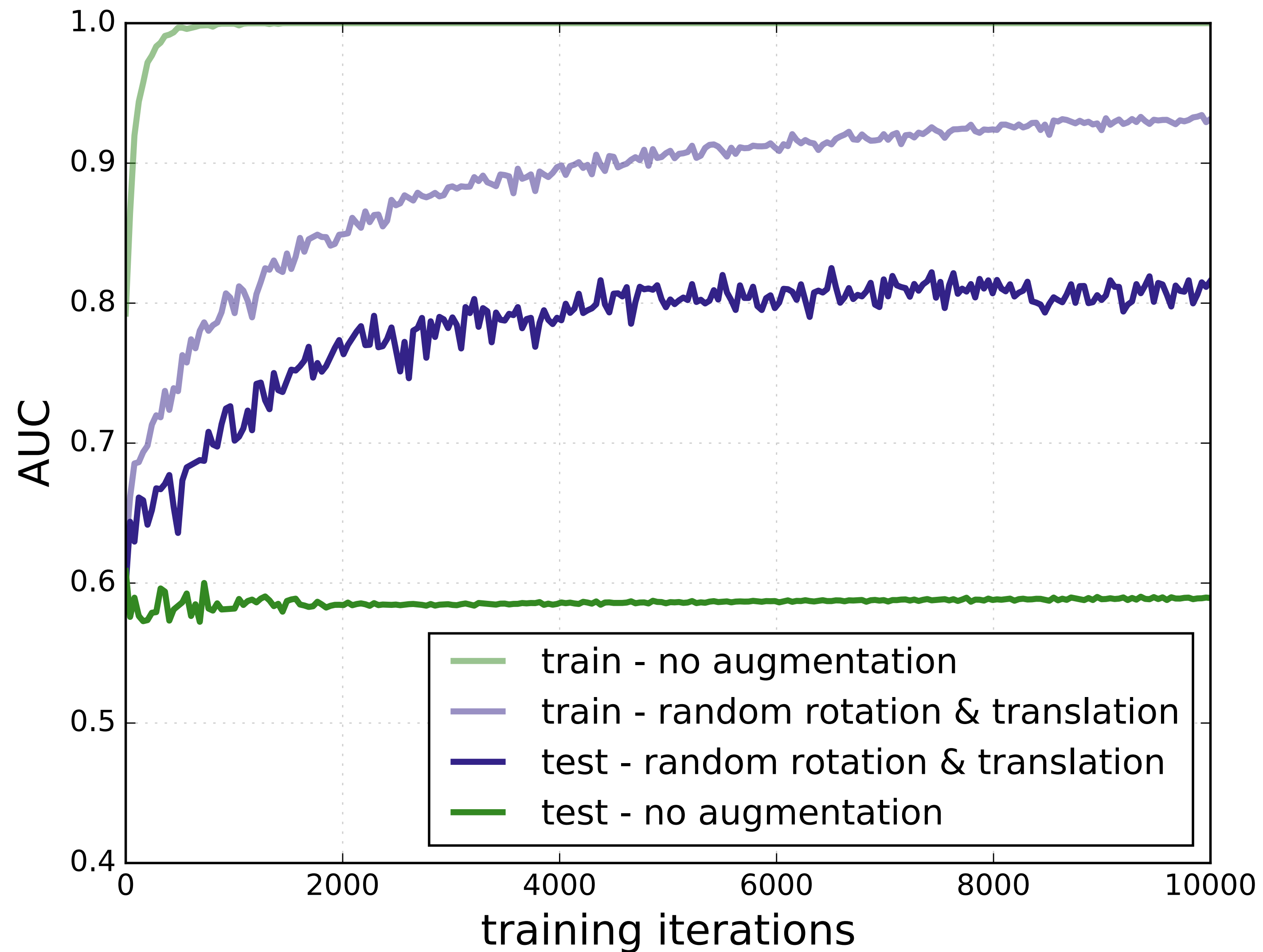
Affinity Prediction

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

Data Augmentation



Data Augmentation



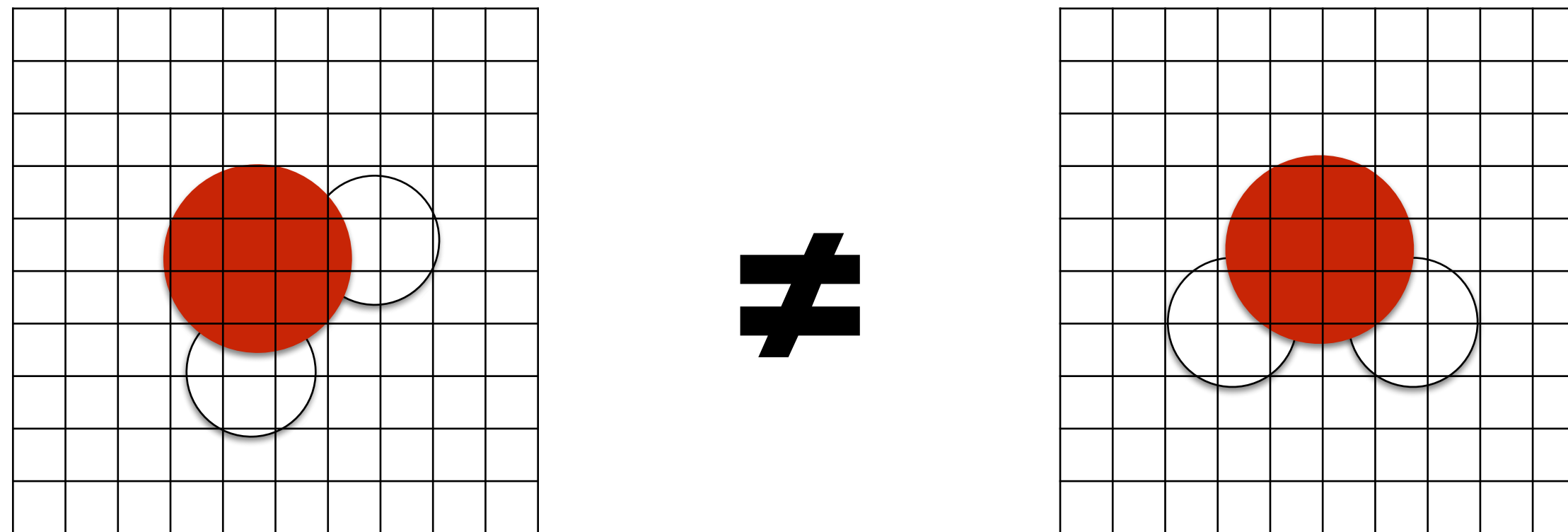
Why Grids?

Cons

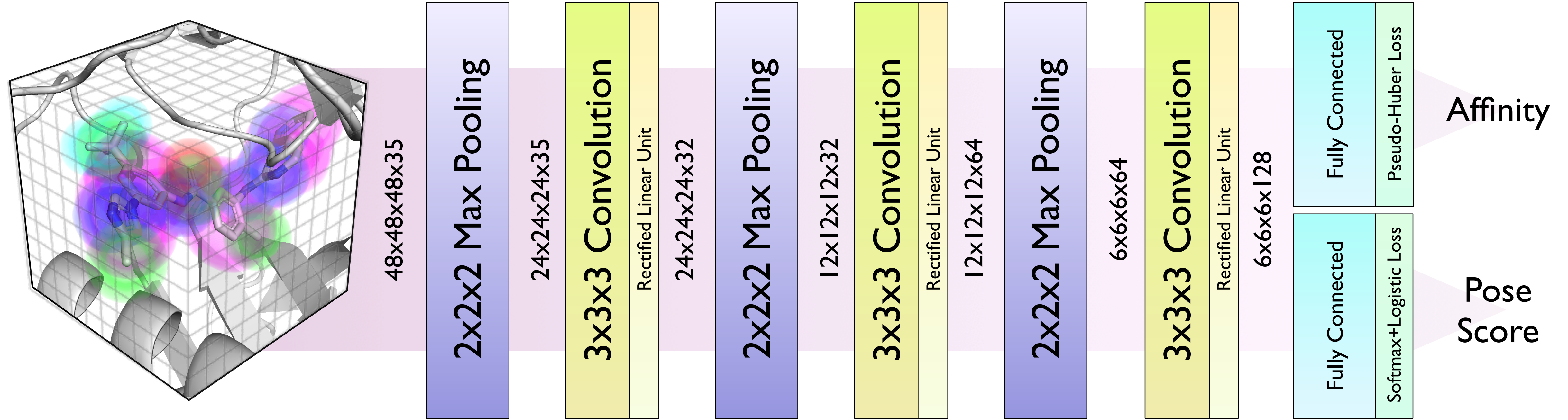
- *coordinate frame dependent*
- pairwise interactions not explicit

Pros

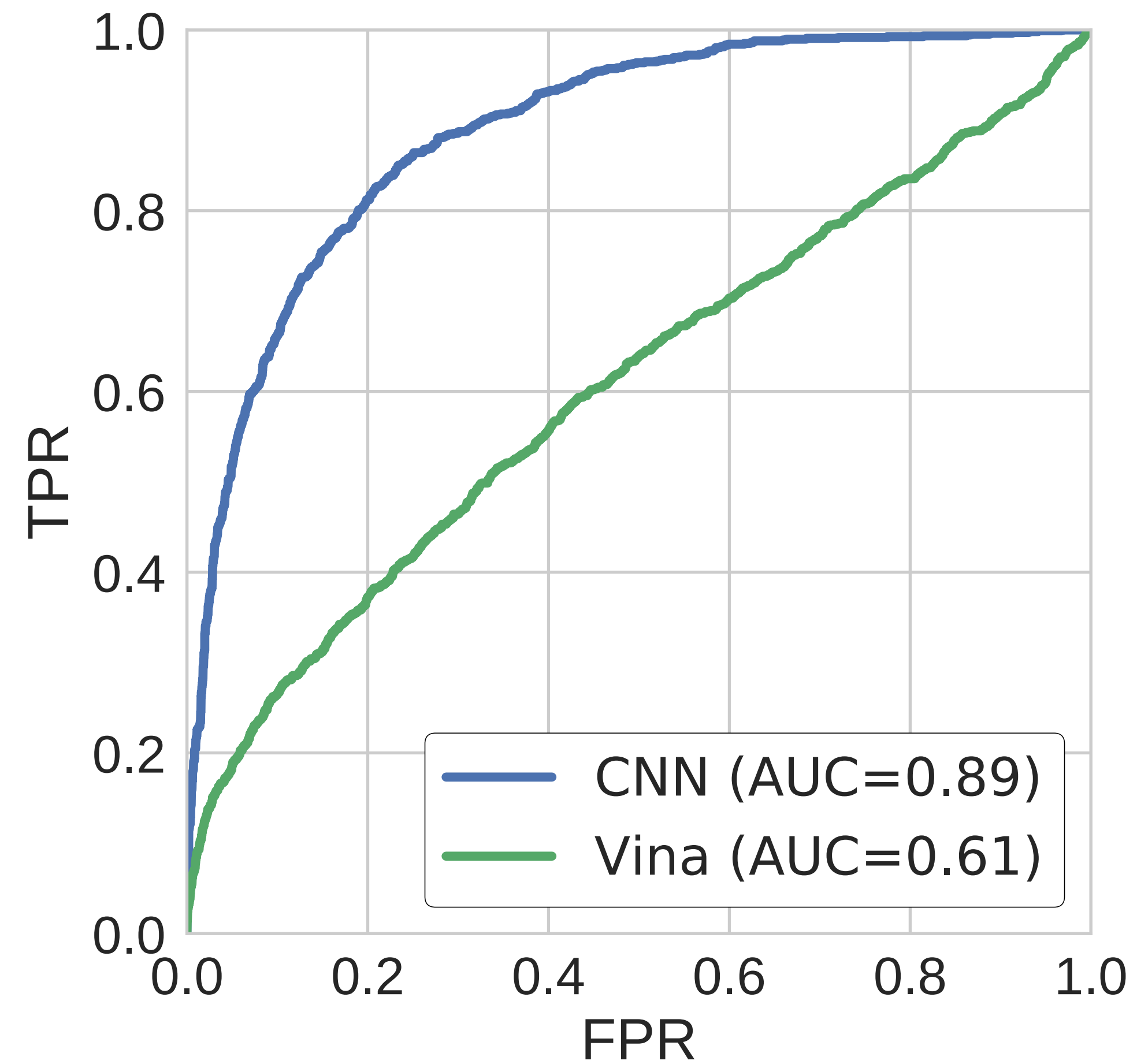
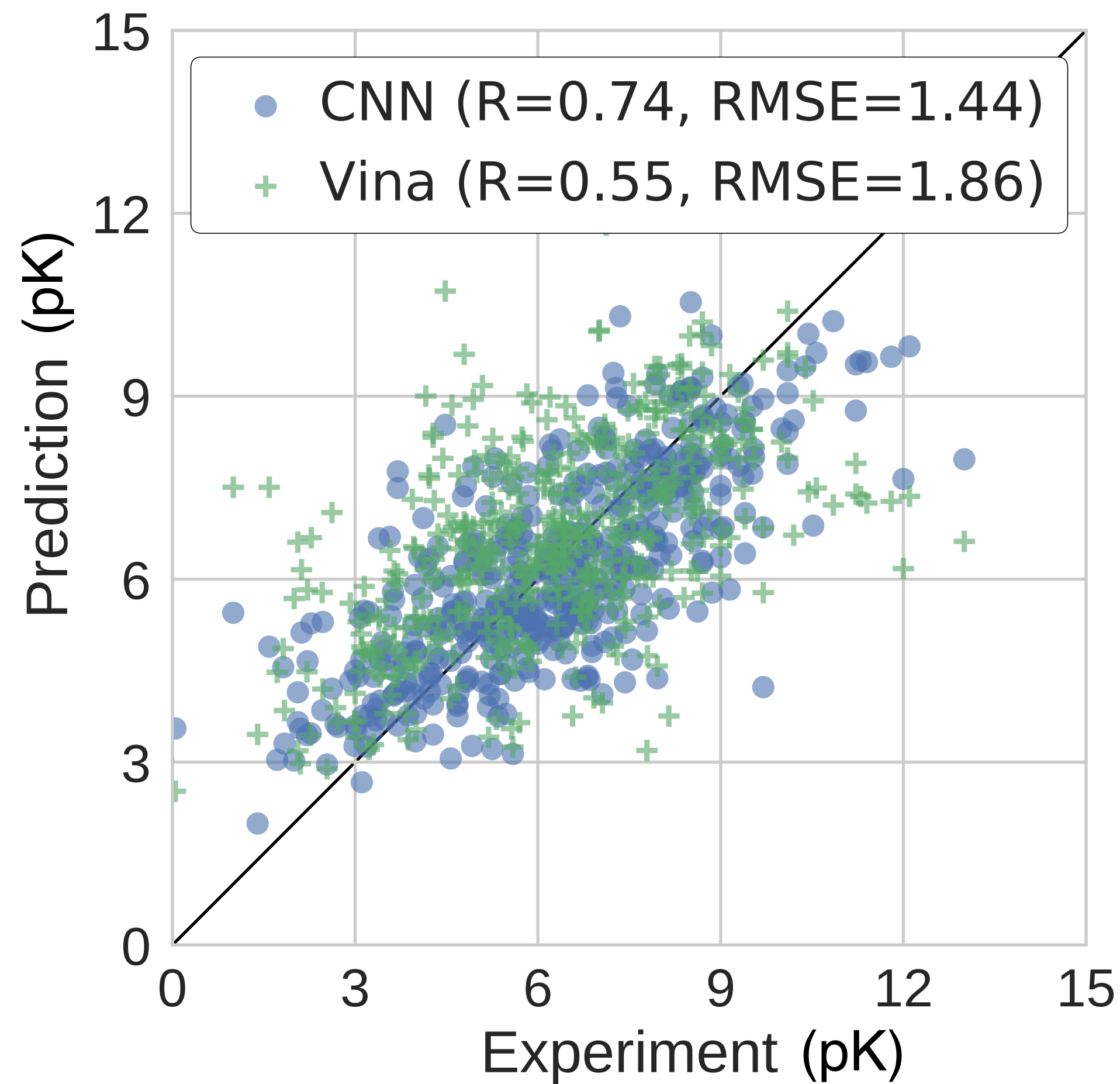
- clear spatial relationships
- amazingly parallel
- easy to interpret



Model

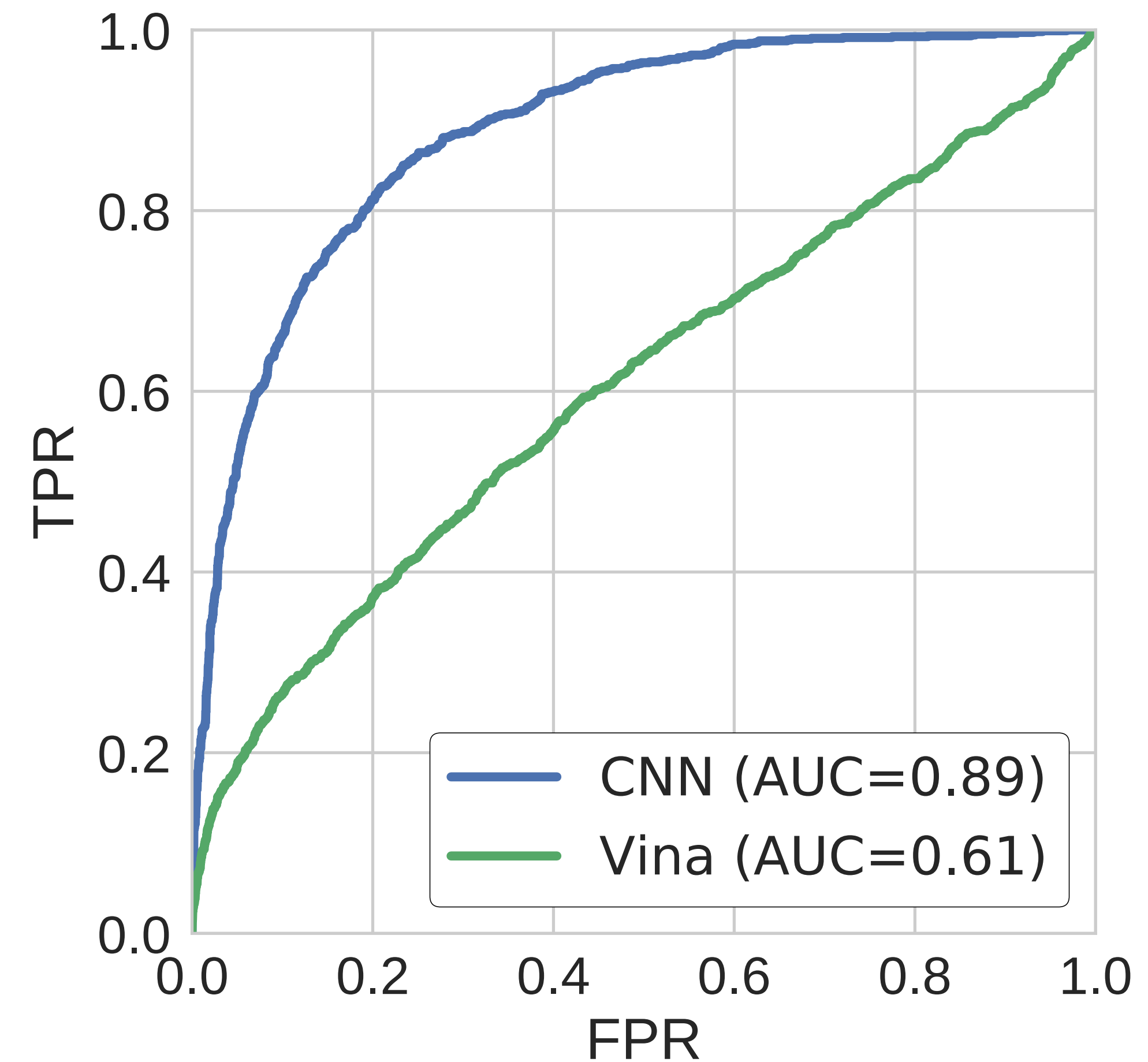
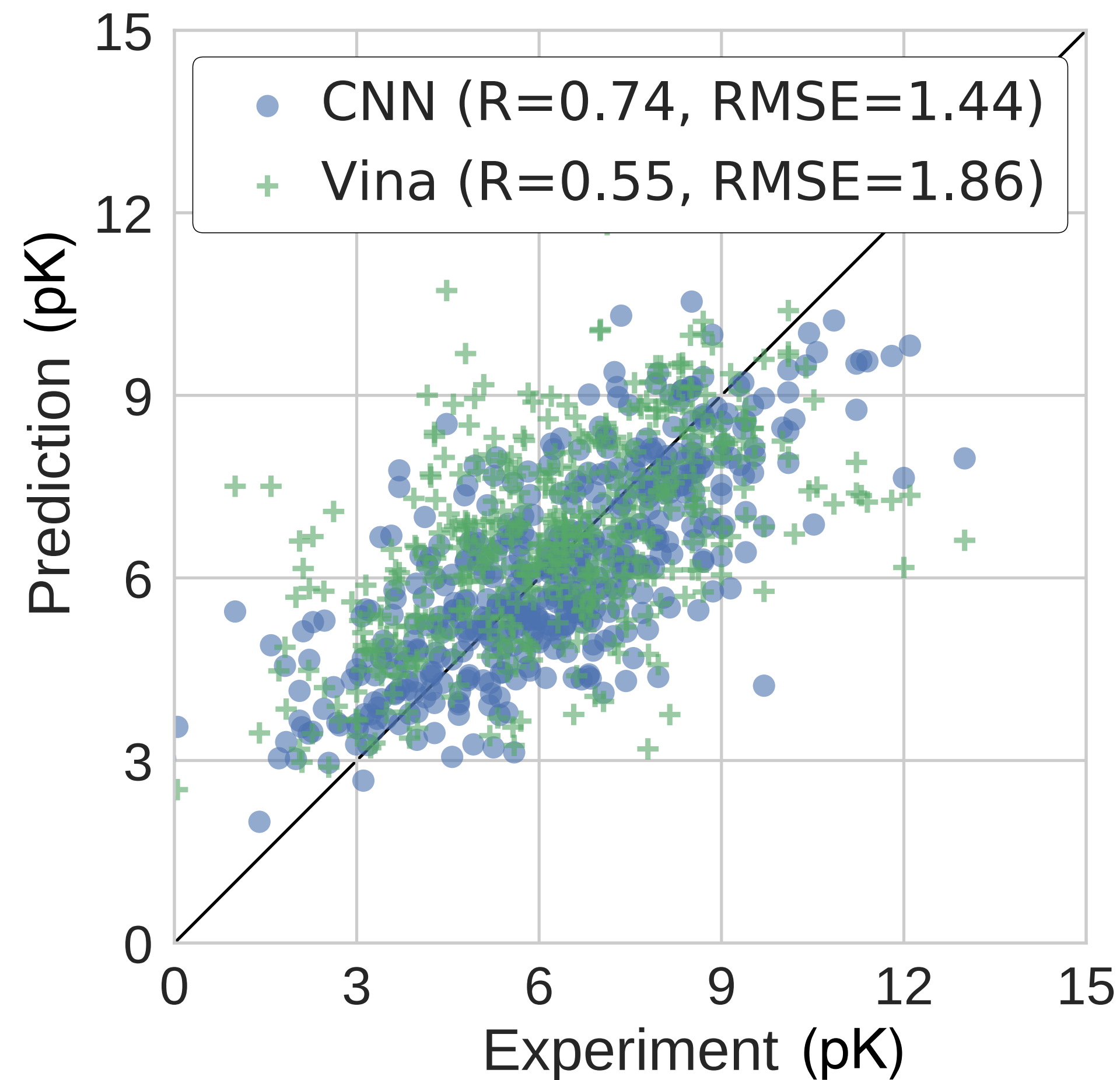


Results



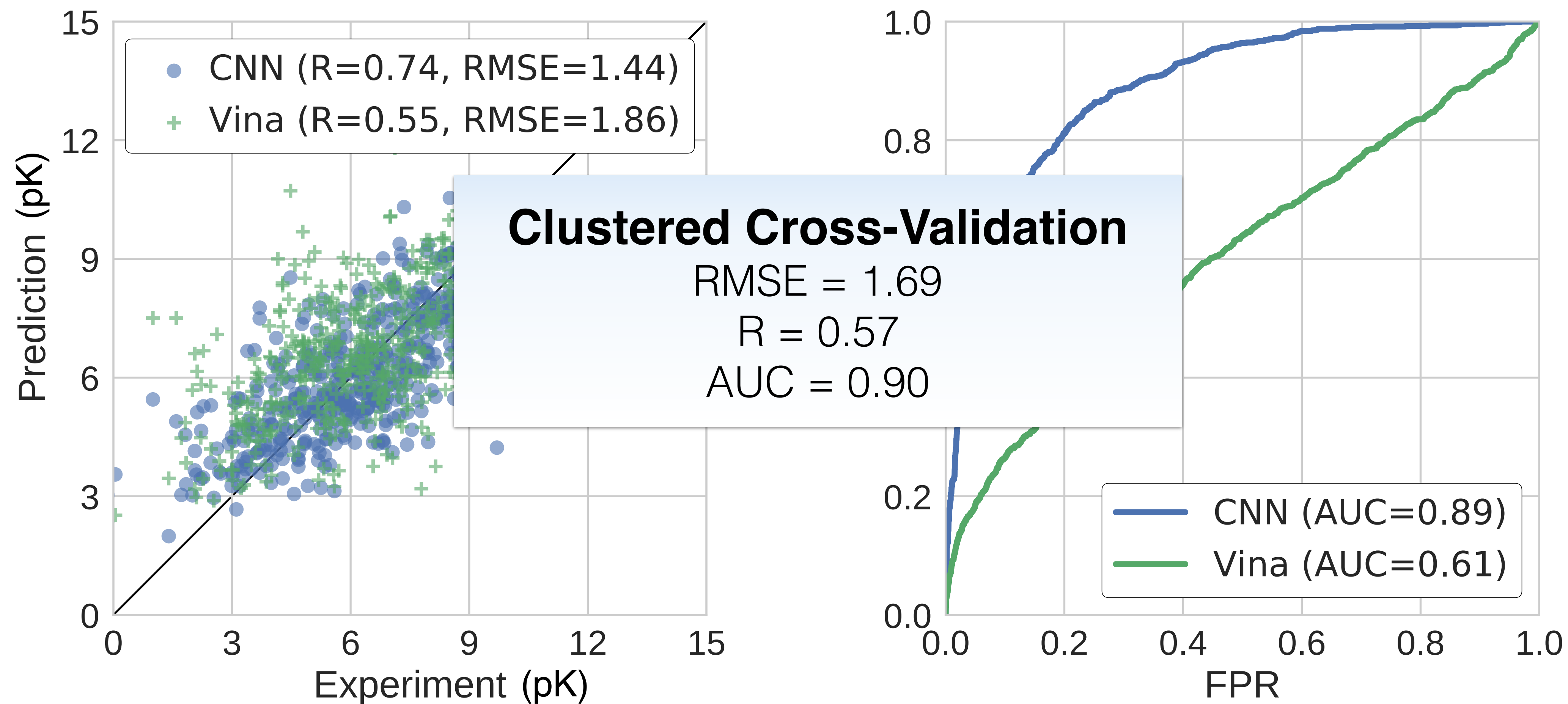
Trained on PDBbind refined; tested on CSAR

Results



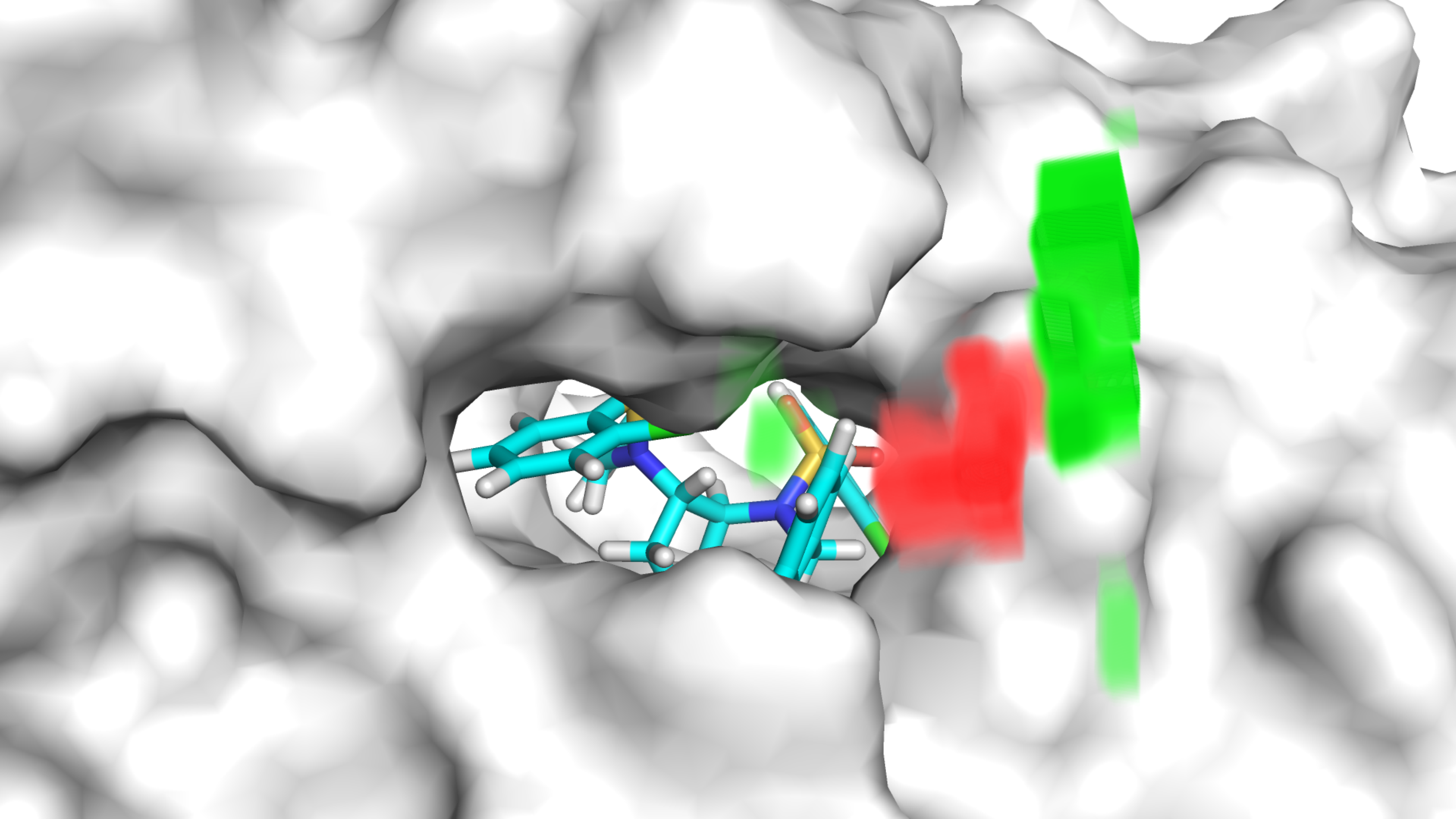
Trained on PDBbind refined; tested on CSAR

Results

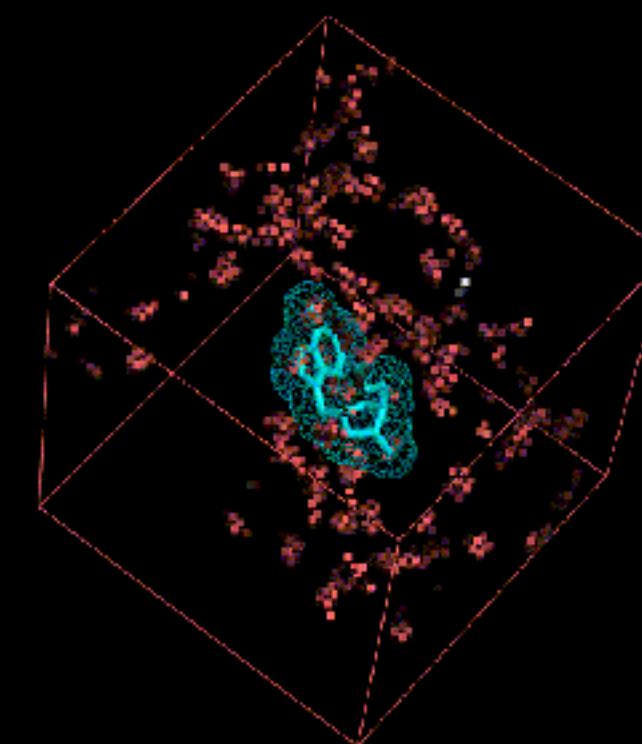
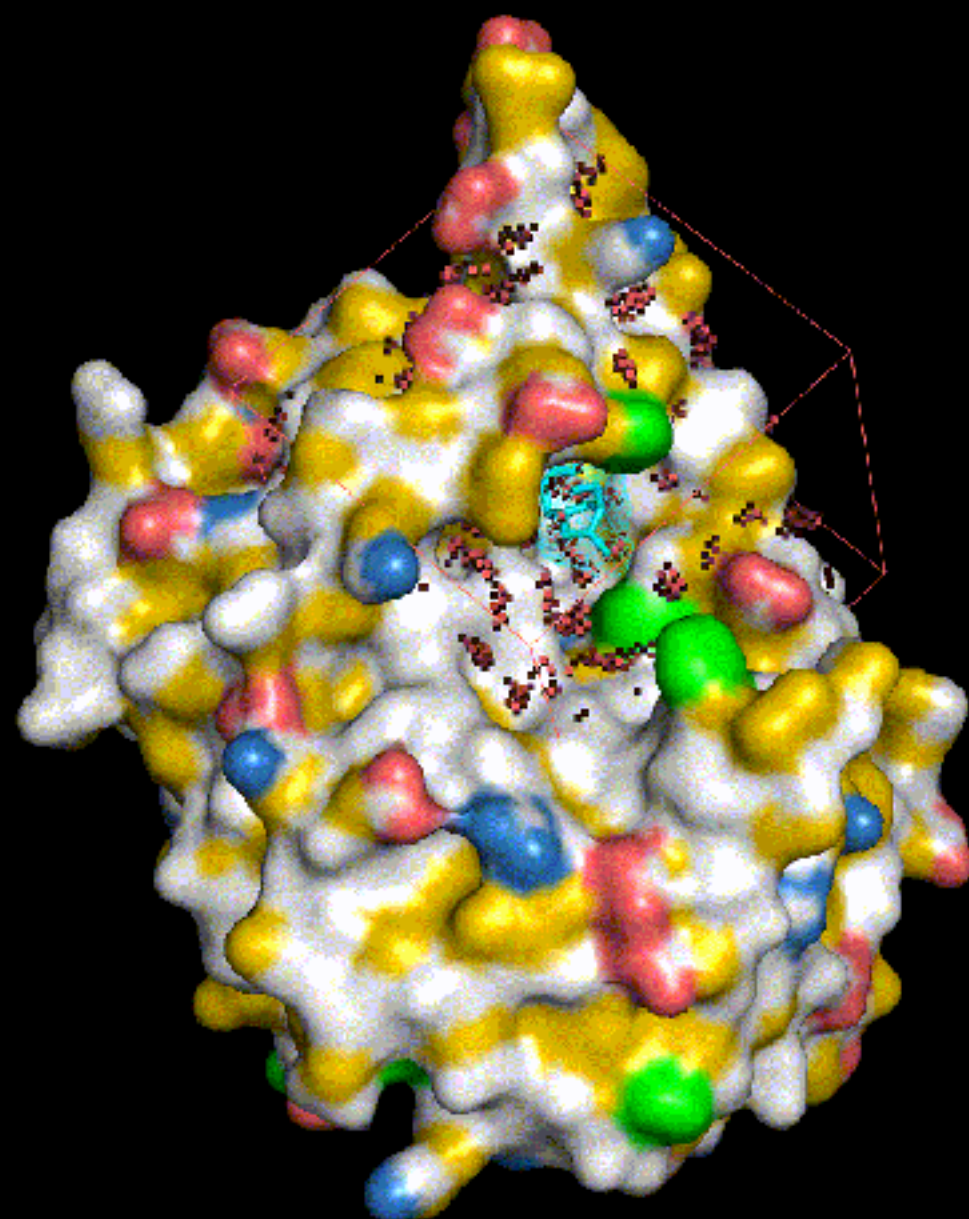


Trained on PDBbind refined; tested on CSAR

What about water?



Grid Inhomogeneous Solvation Theory



GIST analysis of 3BGS (purine nucleoside phosphorylase) active site



Tom Kurtzman



Eric Chen

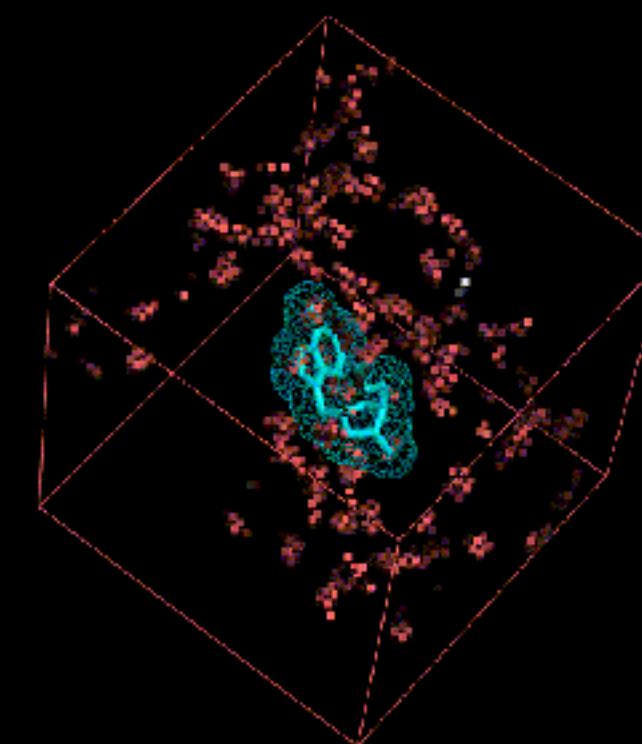
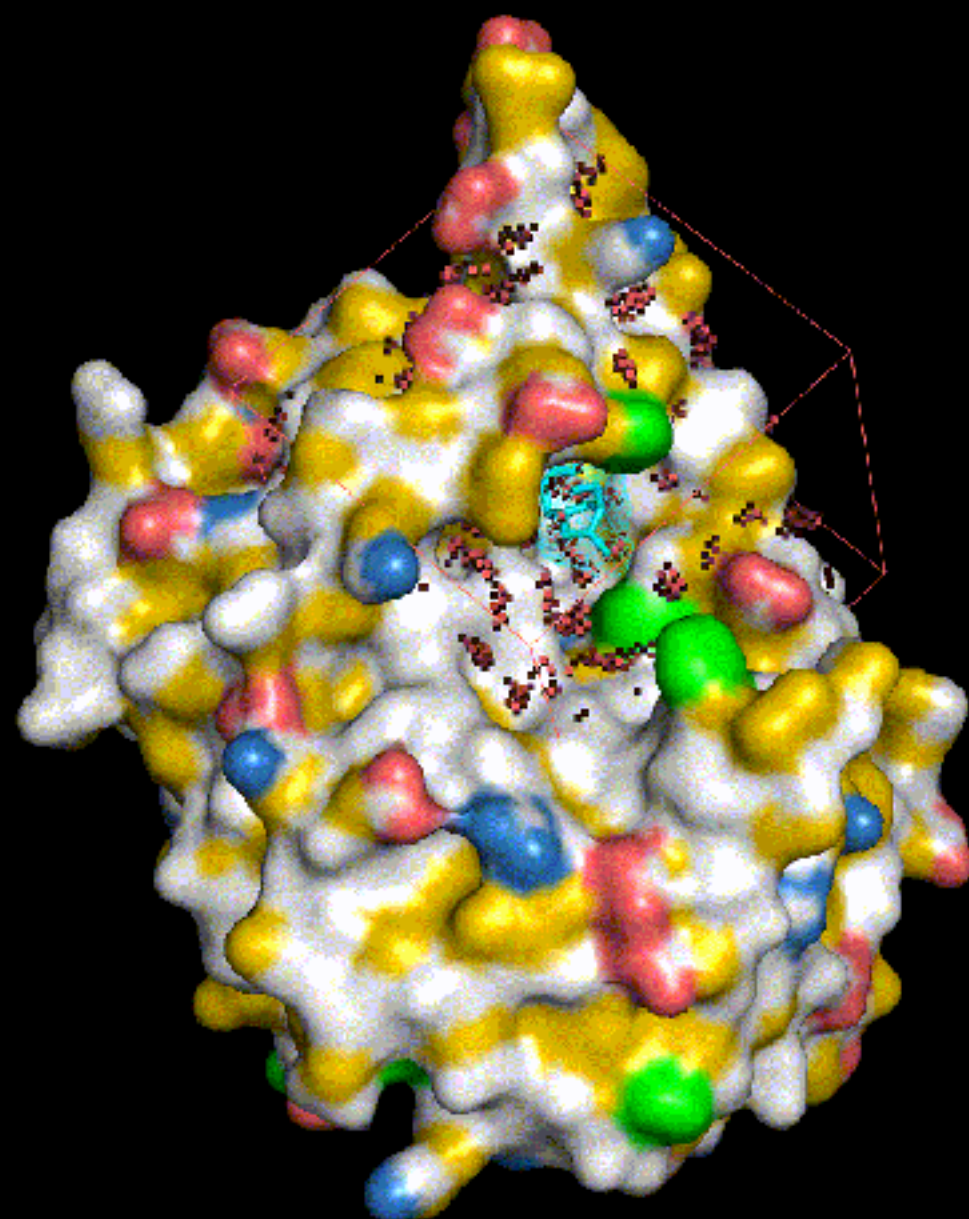


Steven Ramsey



Anthony Cruz-Balberdy

Grid Inhomogeneous Solvation Theory



GIST analysis of 3BGS (purine nucleoside phosphorylase) active site



Tom Kurtzman



Eric Chen

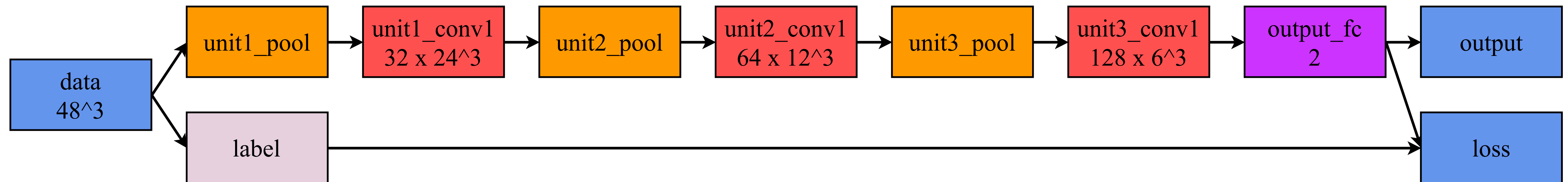


Steven Ramsey

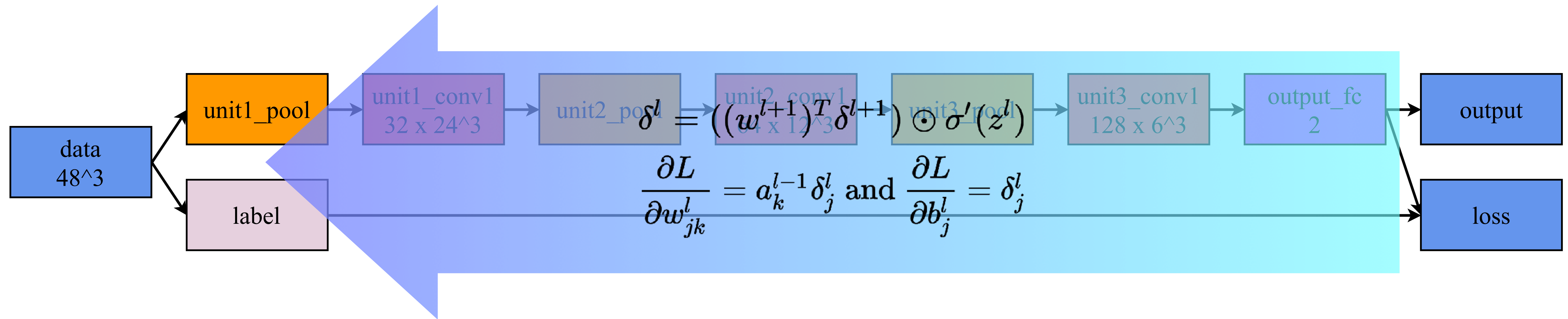


Anthony Cruz-Balberdy

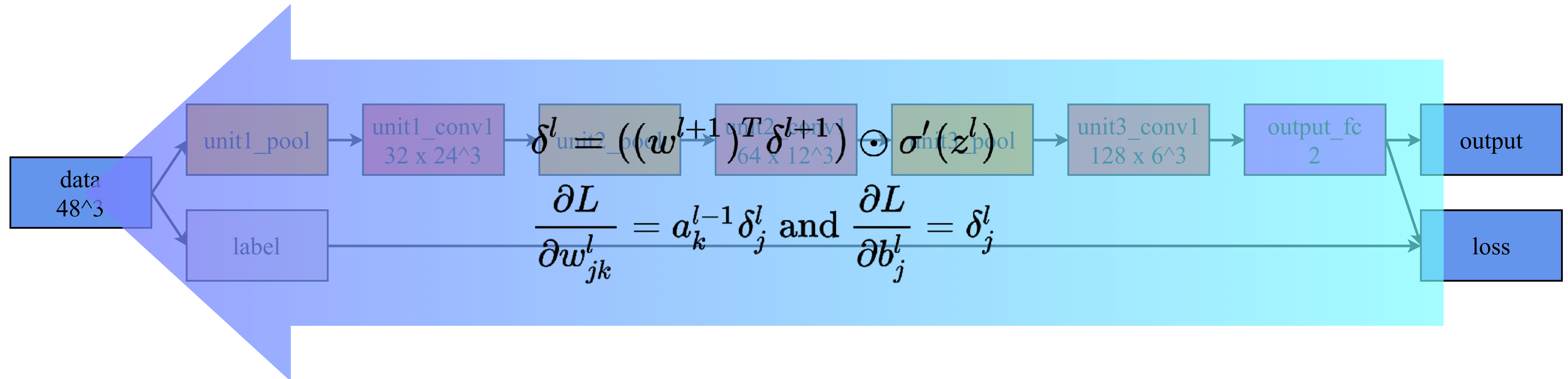
Beyond Scoring



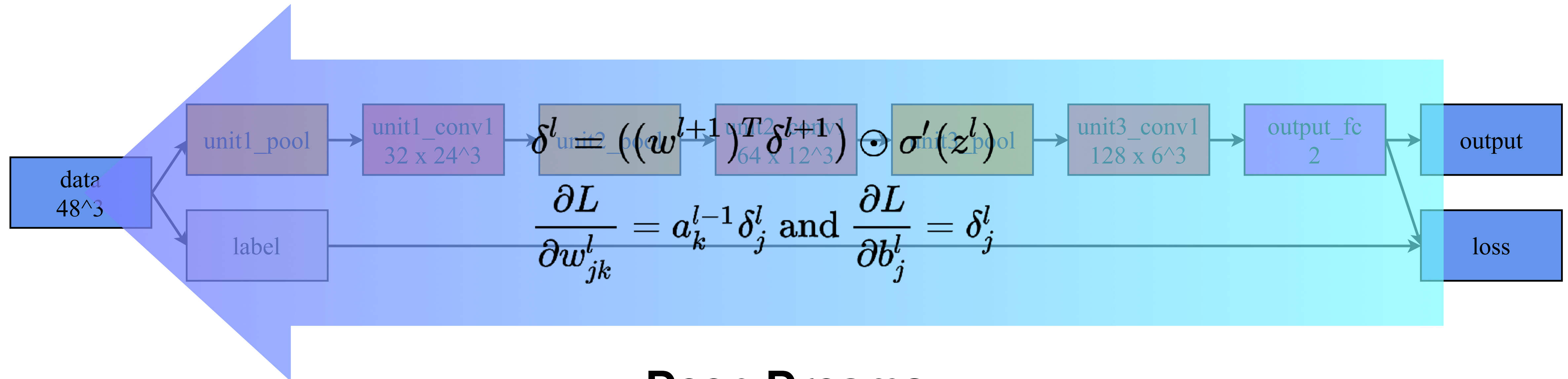
Beyond Scoring



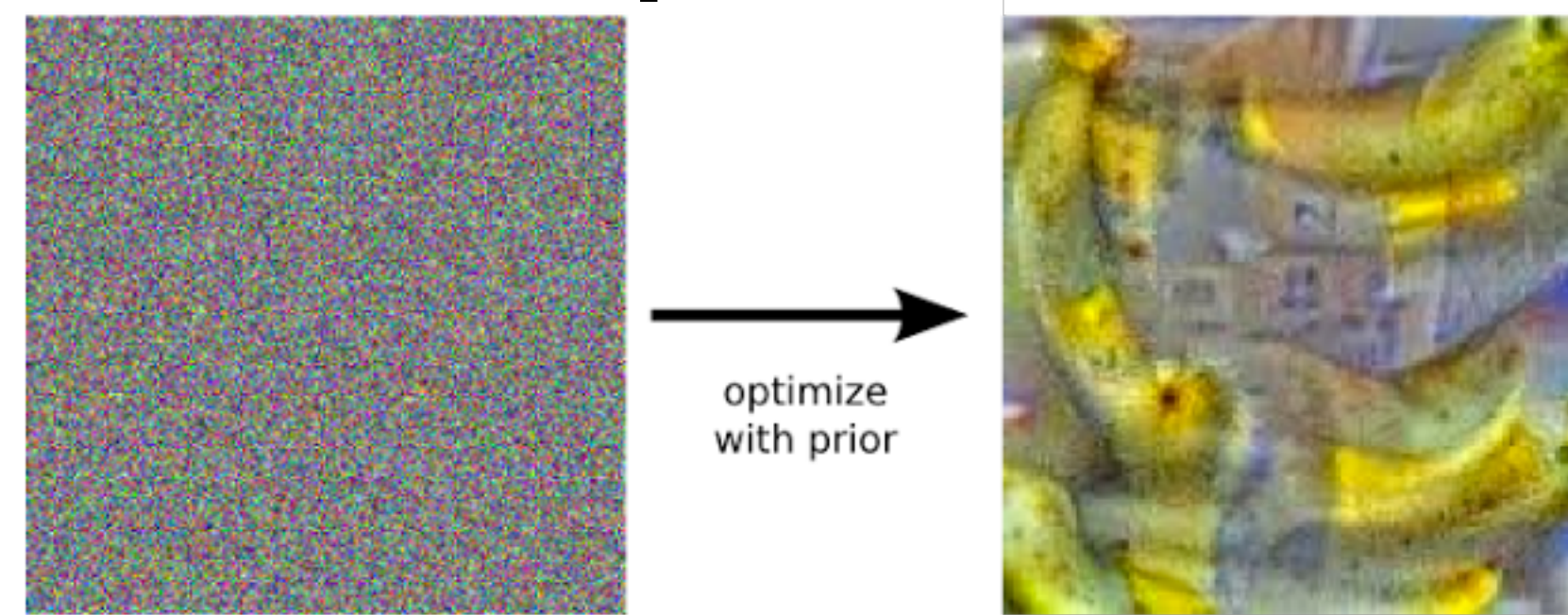
Beyond Scoring



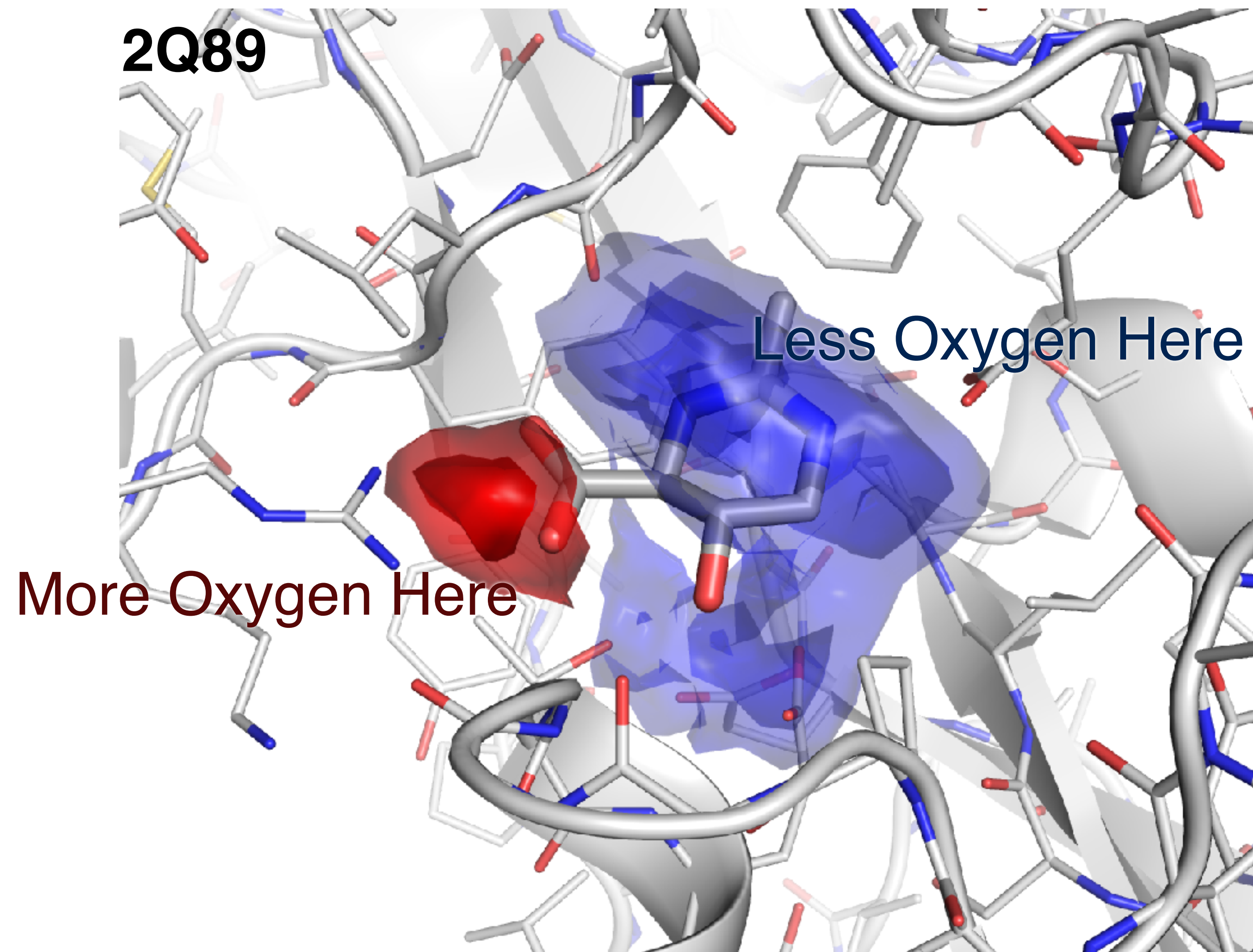
Beyond Scoring



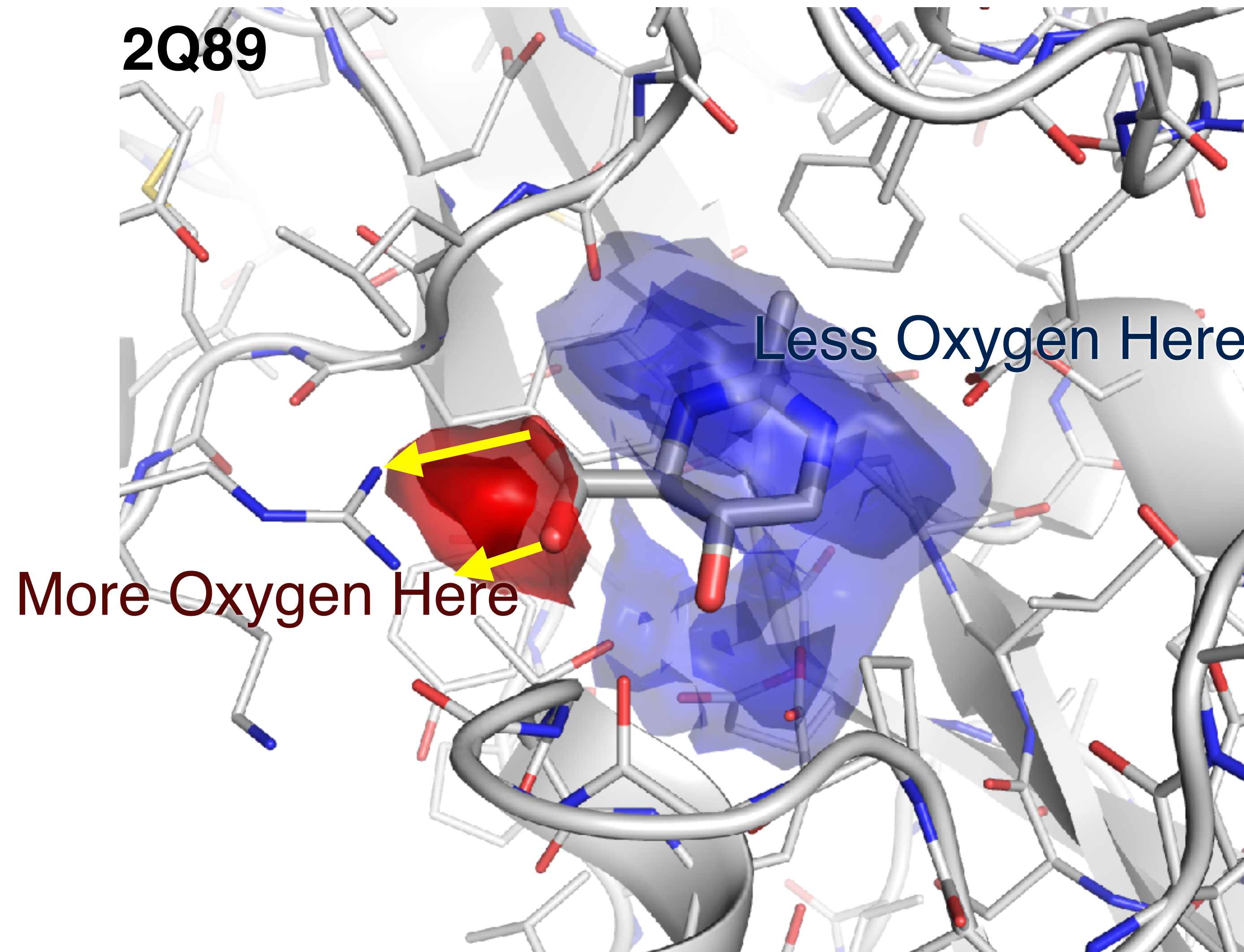
Deep Dreams



Beyond Scoring



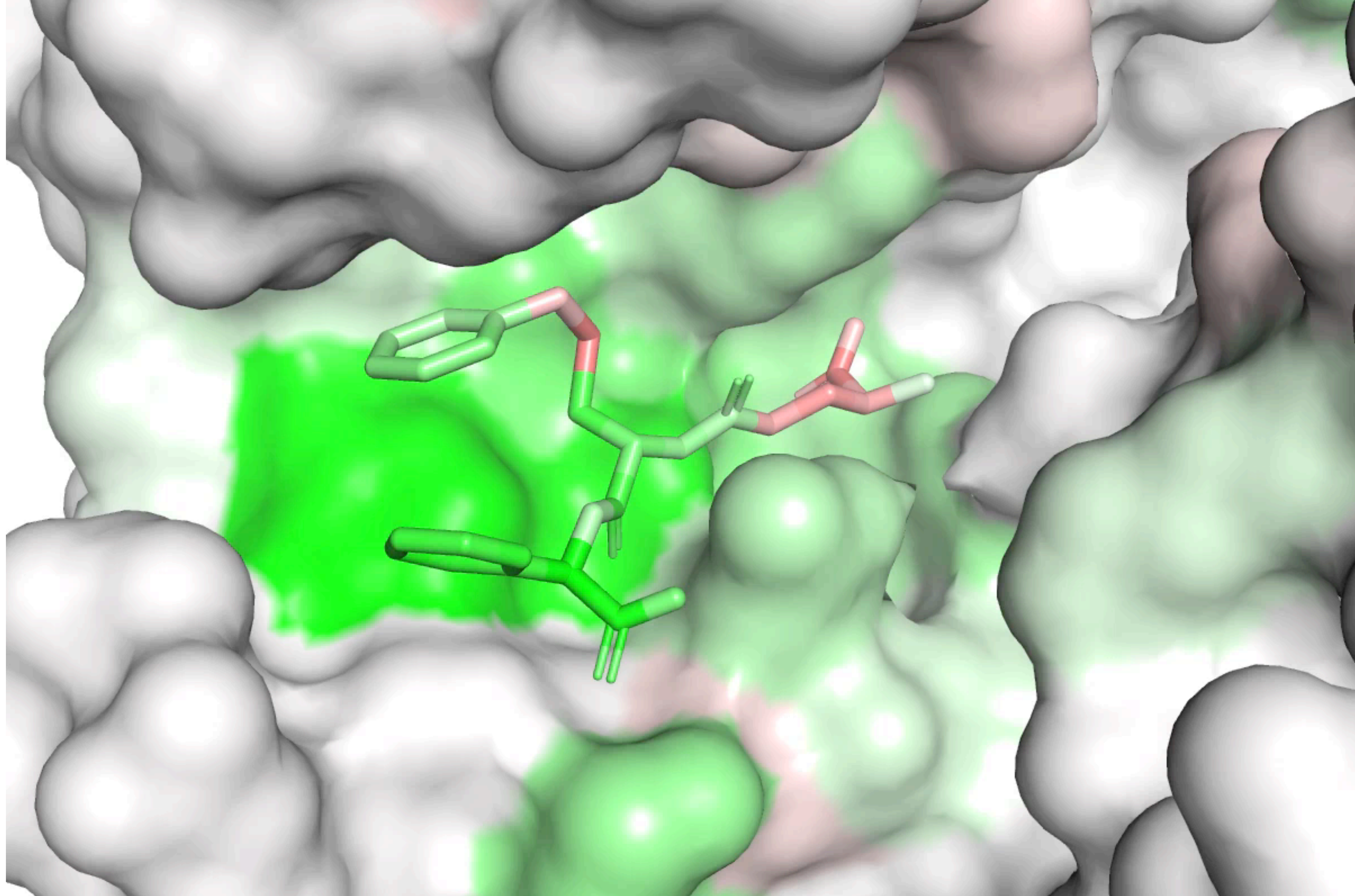
Beyond Scoring

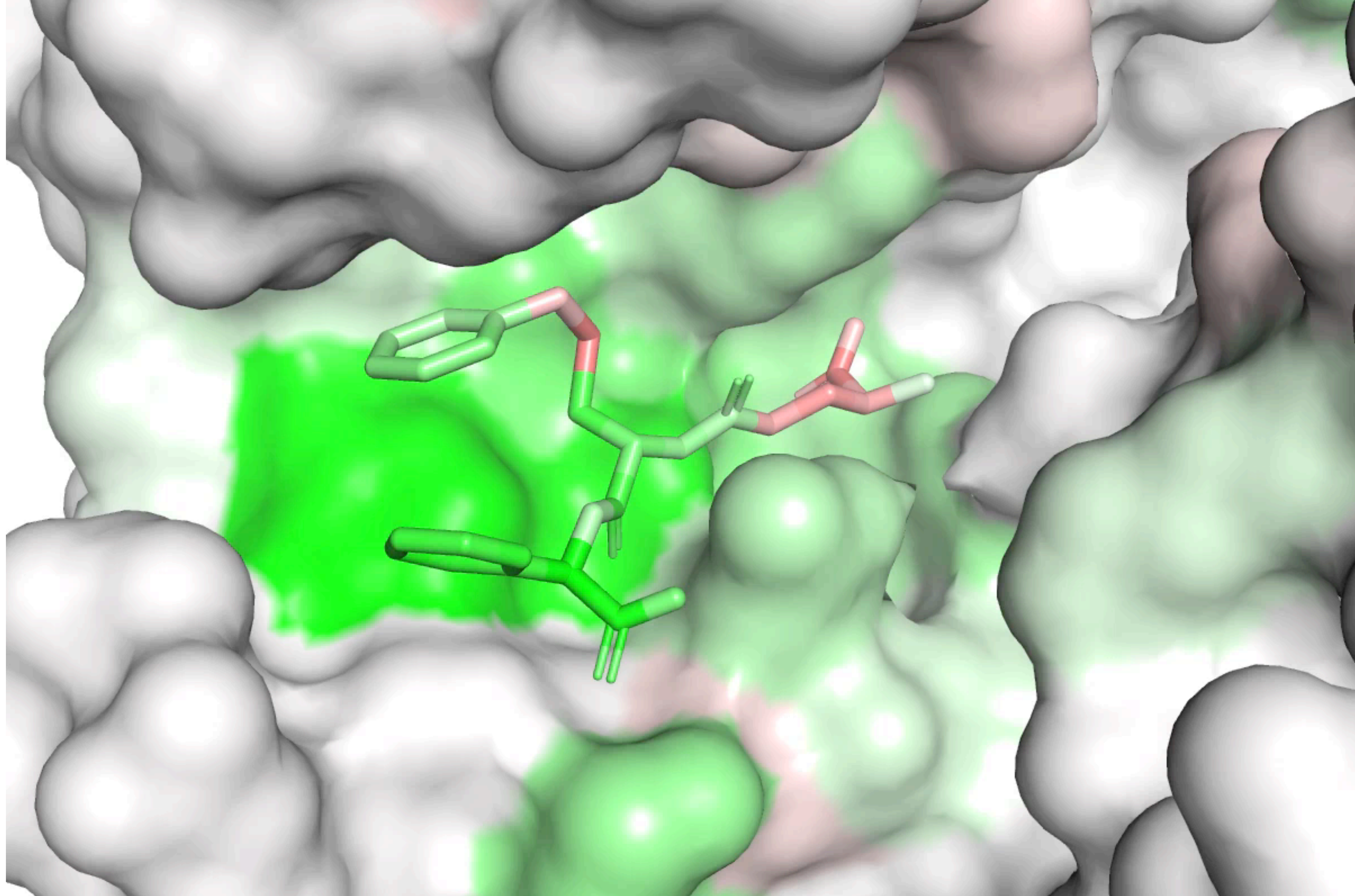


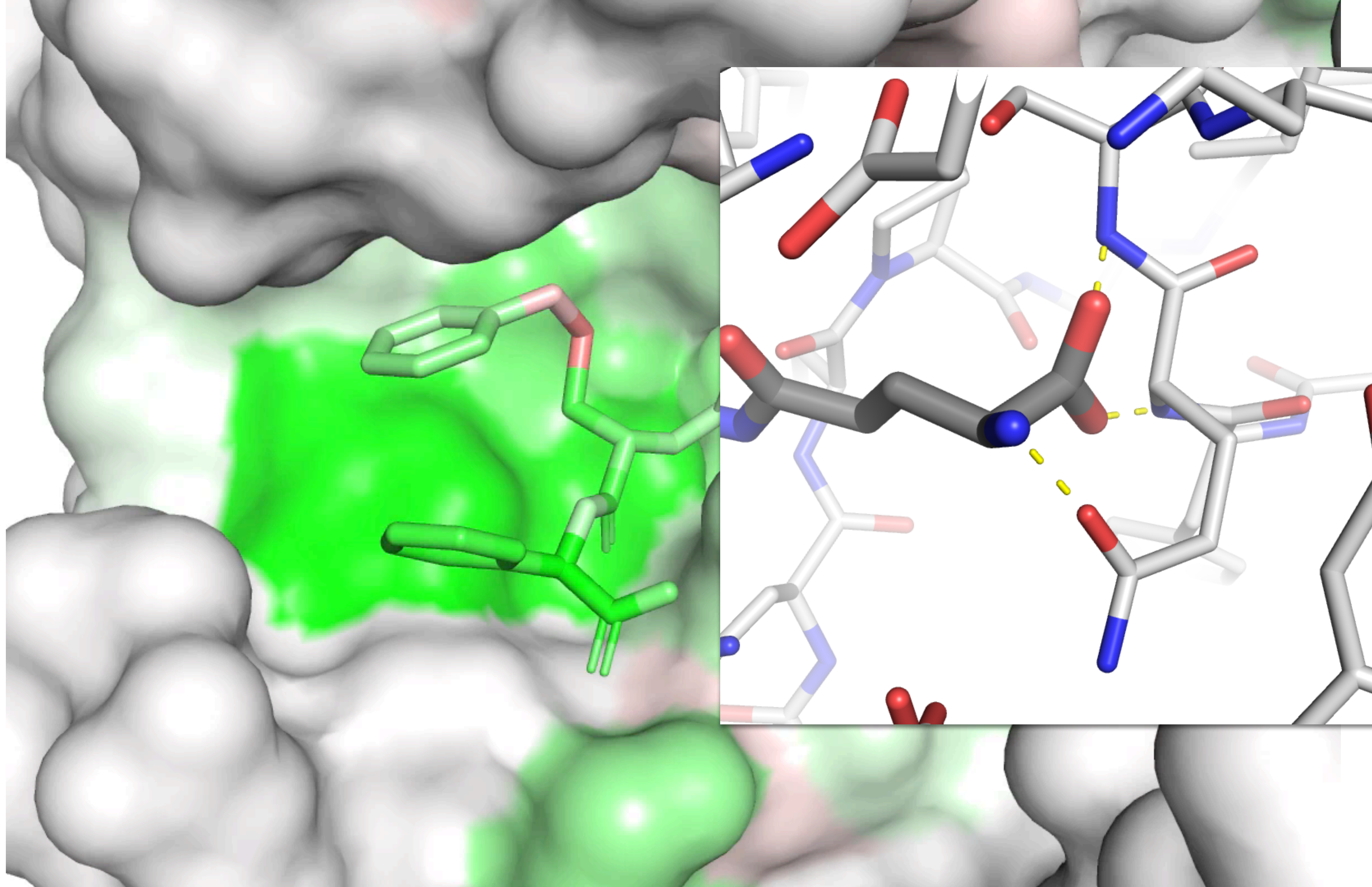
$$\frac{\partial L}{\partial A} = \sum_{i \in G_A} \frac{\partial L}{\partial G_i} \frac{\partial G_i}{\partial D} \frac{\partial D}{\partial A}$$

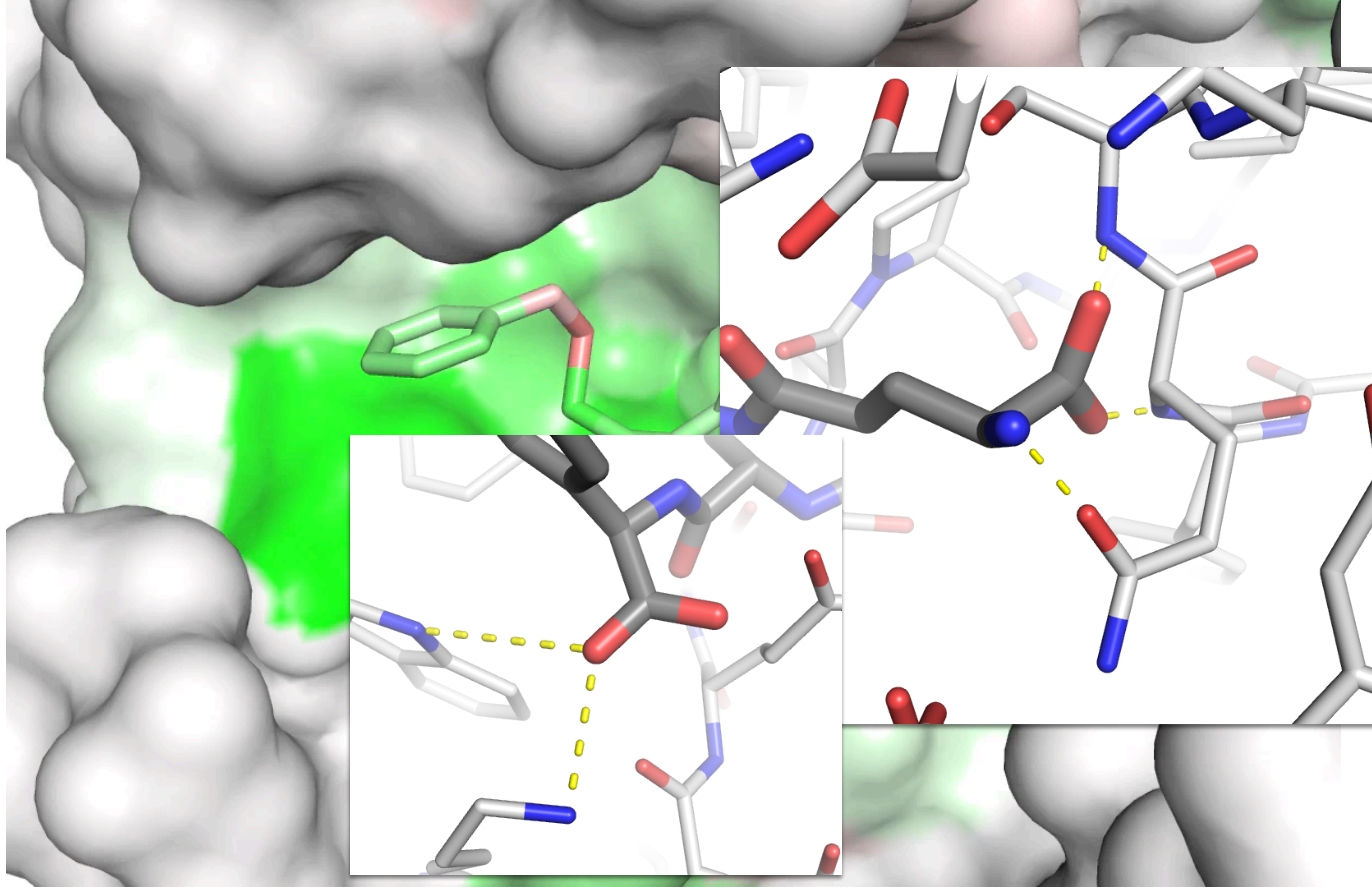
unit1_pool

label

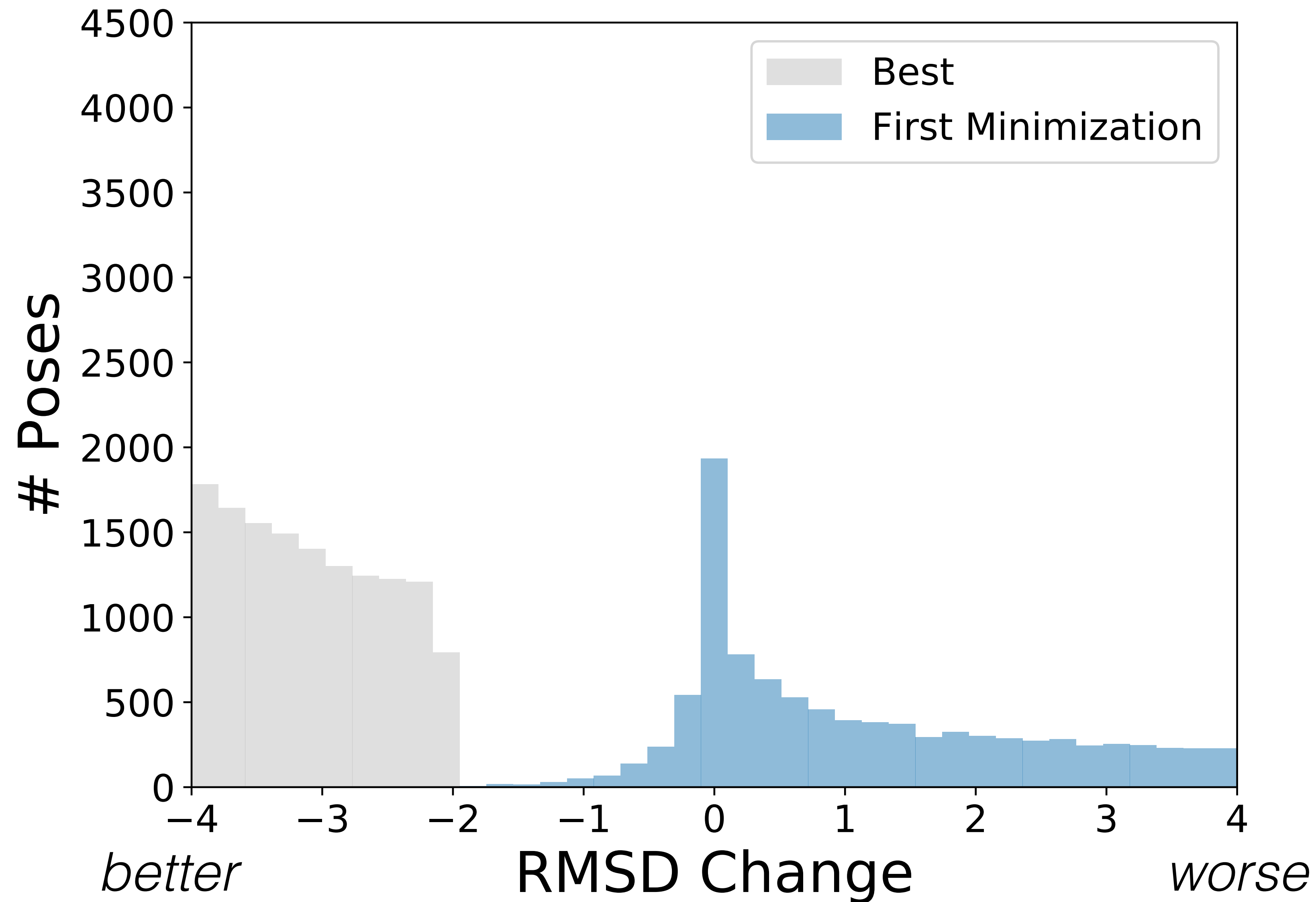




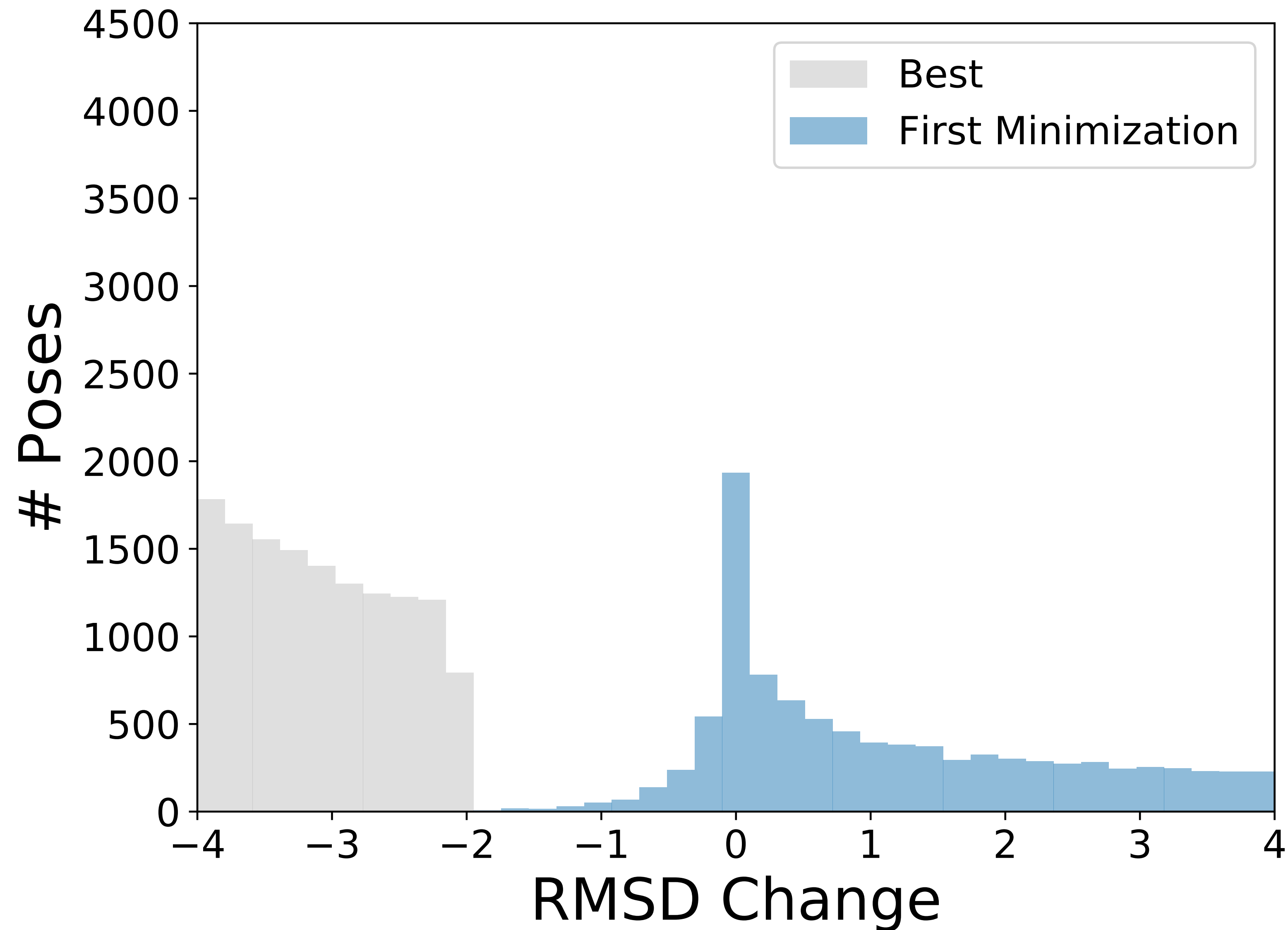




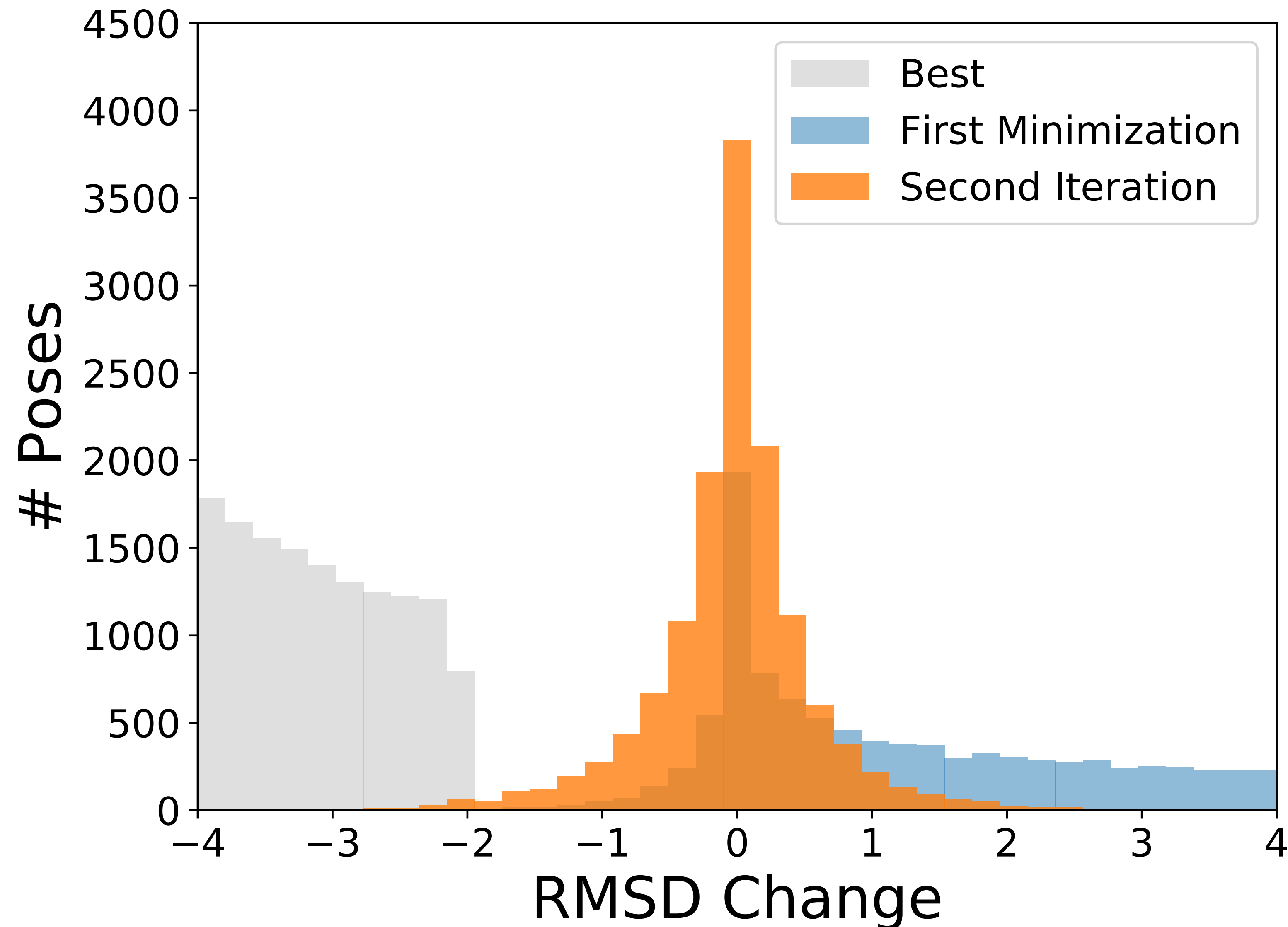
Minimizing Low RMSD Poses



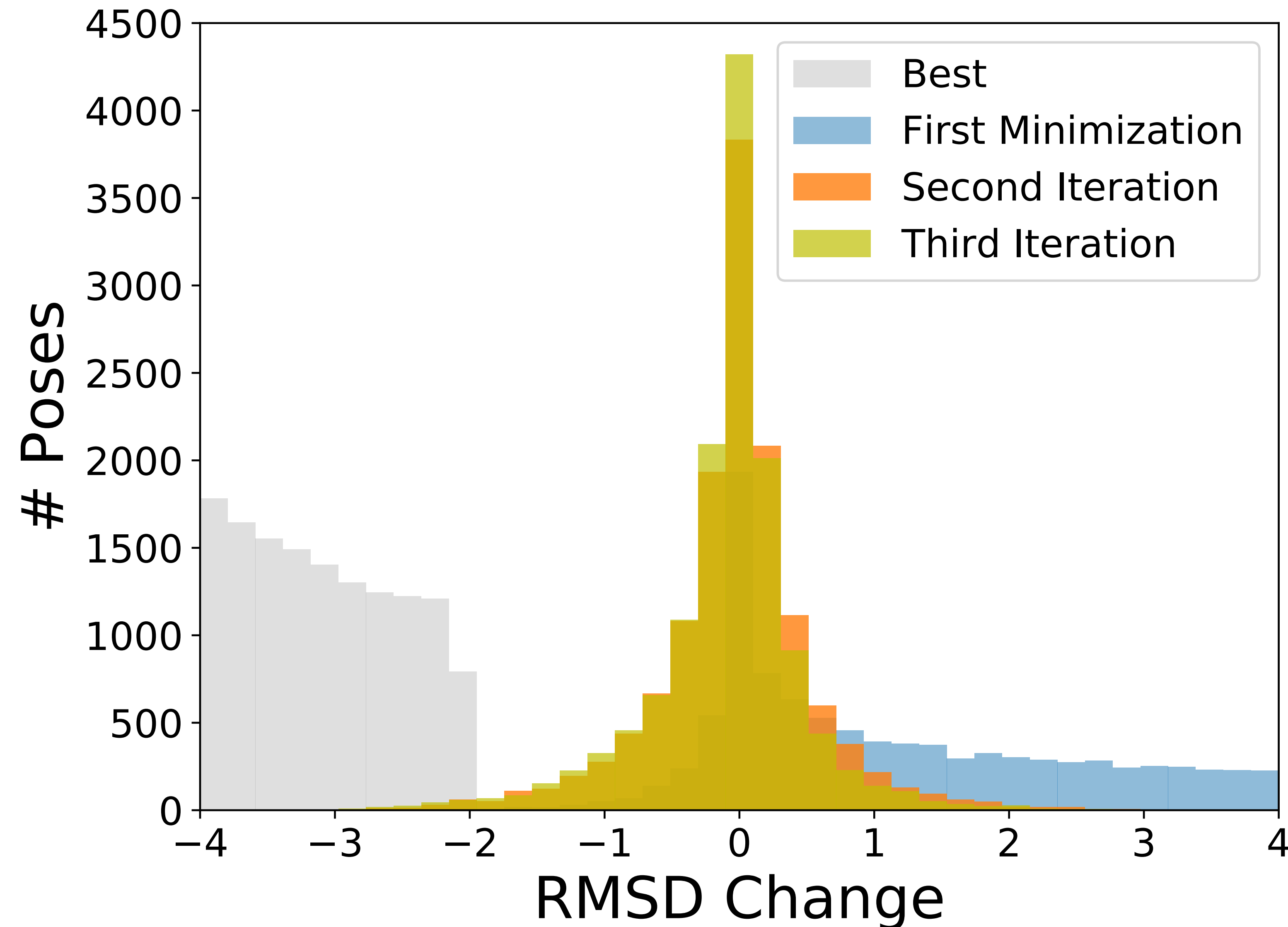
Iterative Refinement



Iterative Refinement



Iterative Refinement



Docking

vina/smina/gnina

Sampling

MCMC

MCMC

MCMC

MCMC

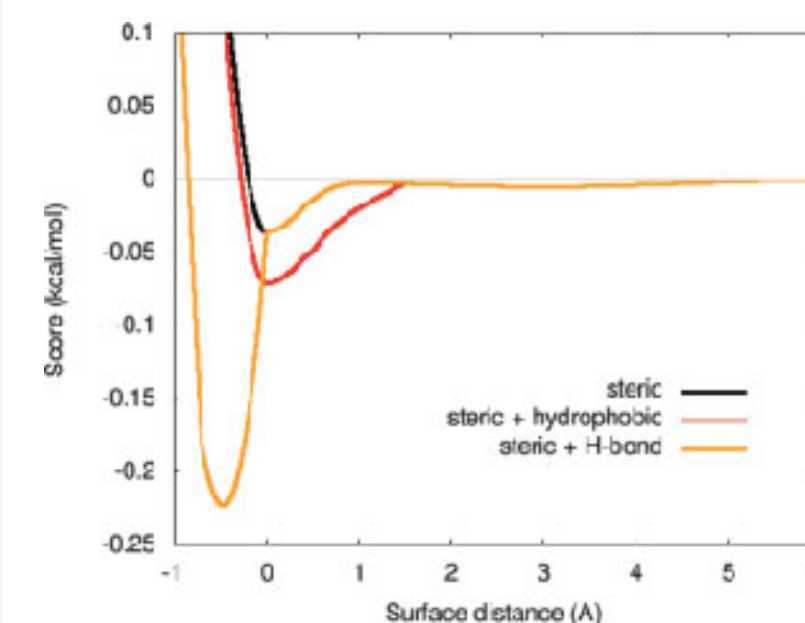
MCMC

⋮

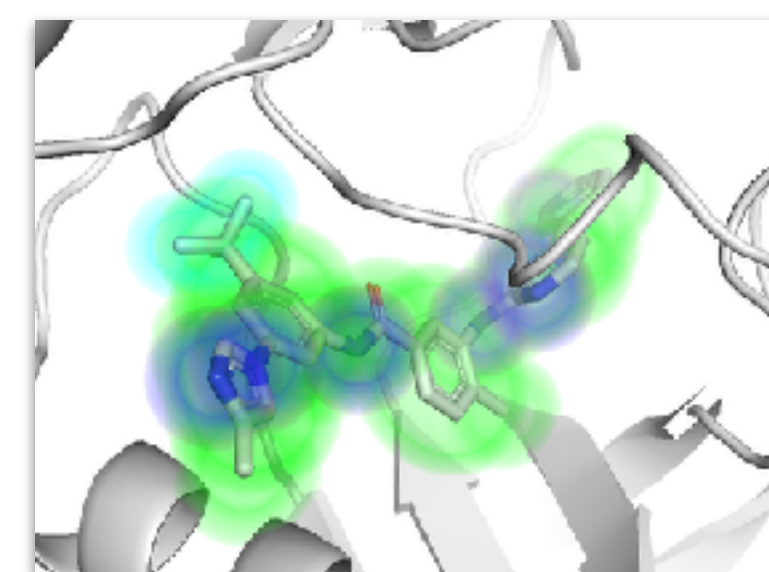
*N (50) independent Monte Carlo chains
Scored with grid-accelerated Vina
Best identified pose retained*

best
poses

Refinement



Vina



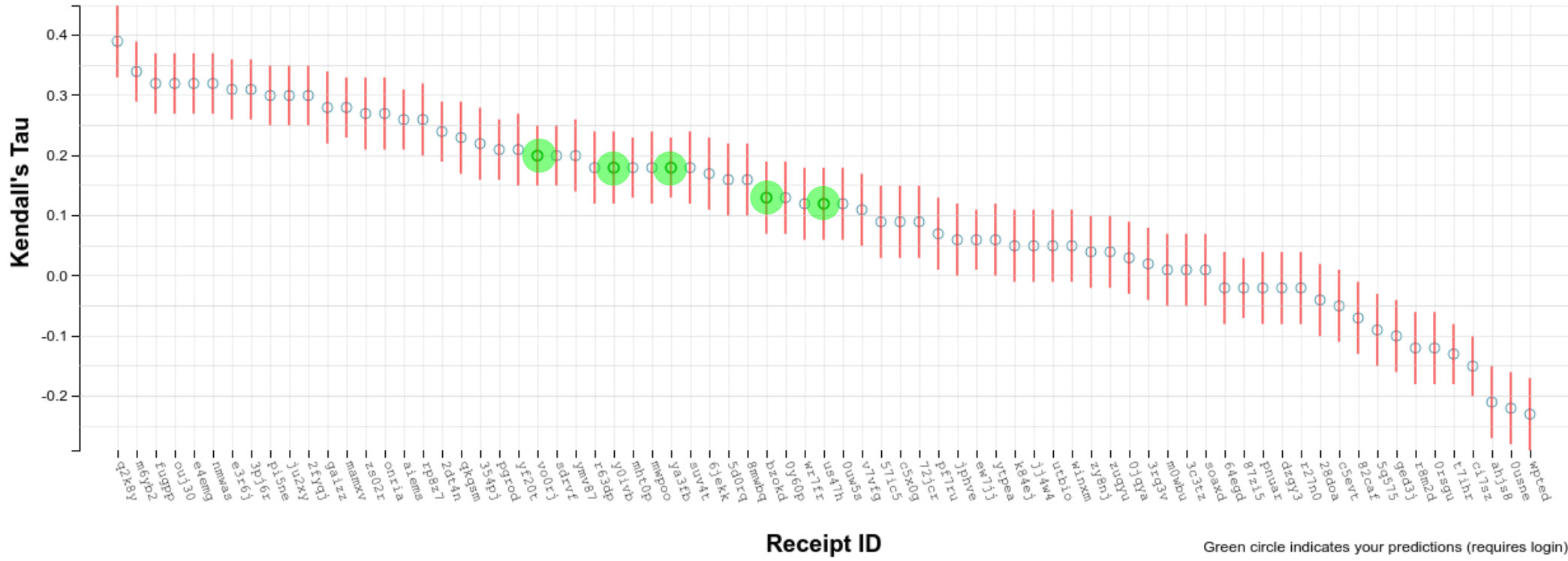
CNN

Rescoring
CNN
pose
affinity

D3R Grand Challenge 3

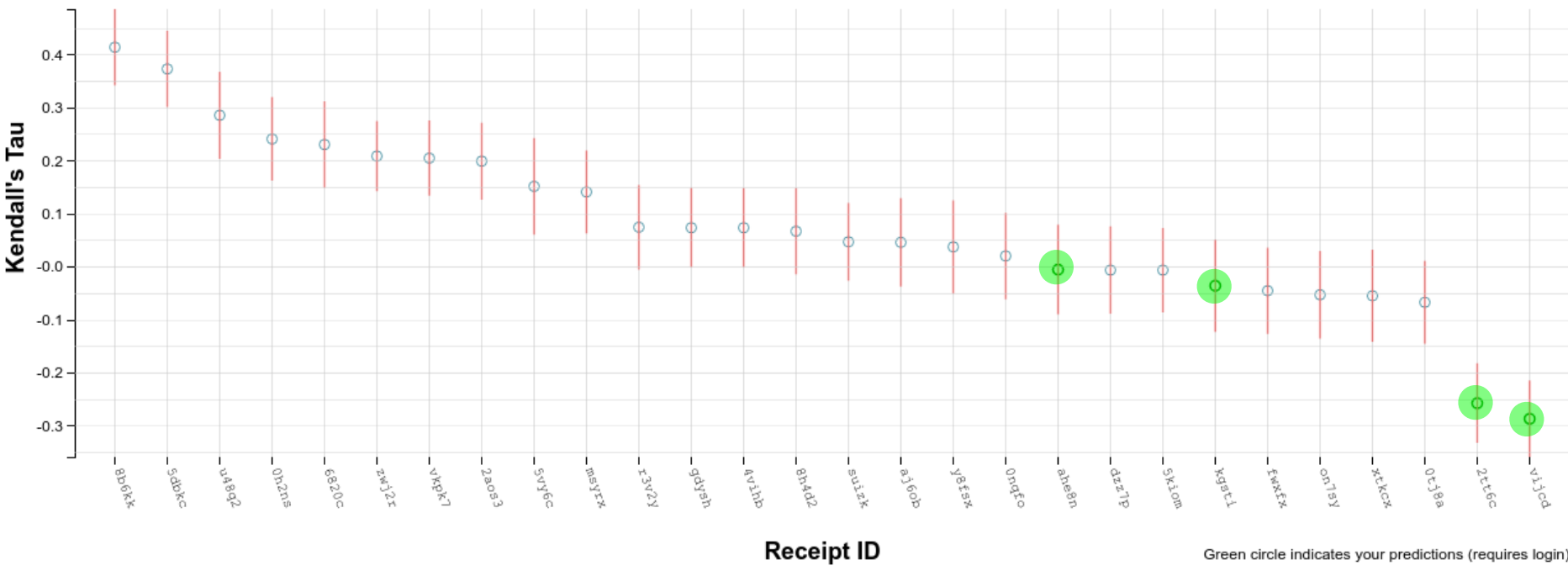
Grand Challenge 3 - CatS_stage2

Affinity Ranking - Kendall's Tau



Grand Challenge 3 - p38a

Affinity Ranking - Kendall's Tau



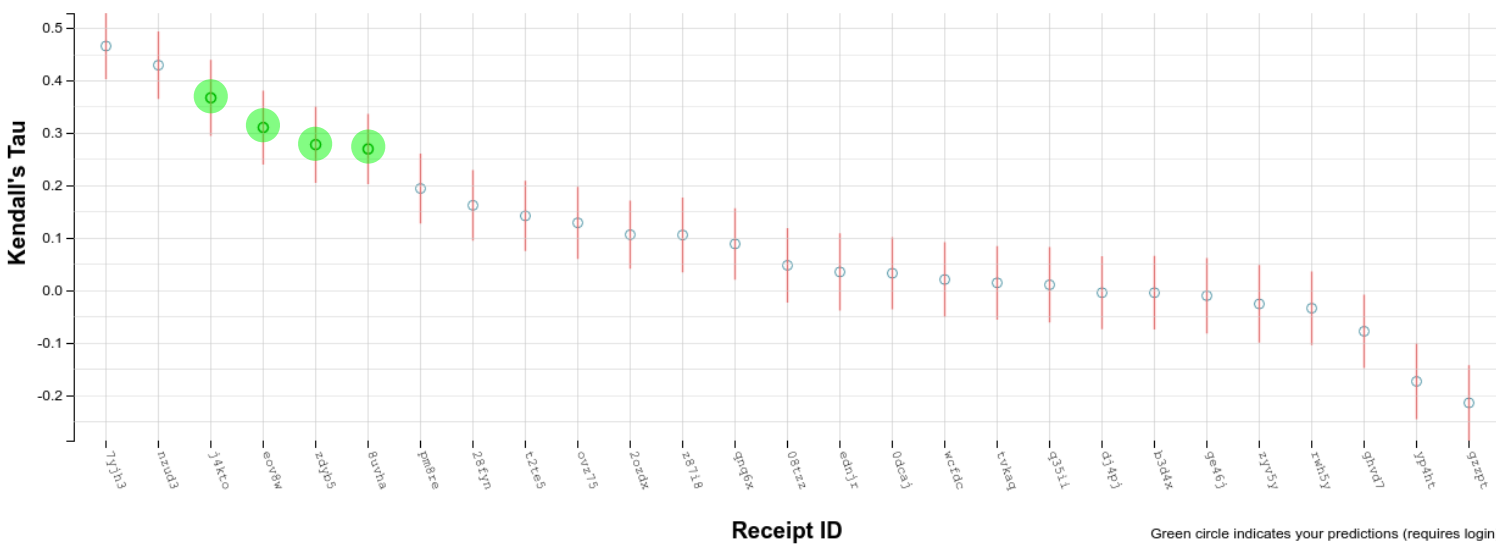
Grand Challenge 3 - VEGFR2

Affinity Ranking - Kendall's Tau



Grand Challenge 3 - JAK2_SC2

Affinity Ranking - Kendall's Tau



Grand Challenge 3 - TIE2

Affinity Ranking - Kendall's Tau



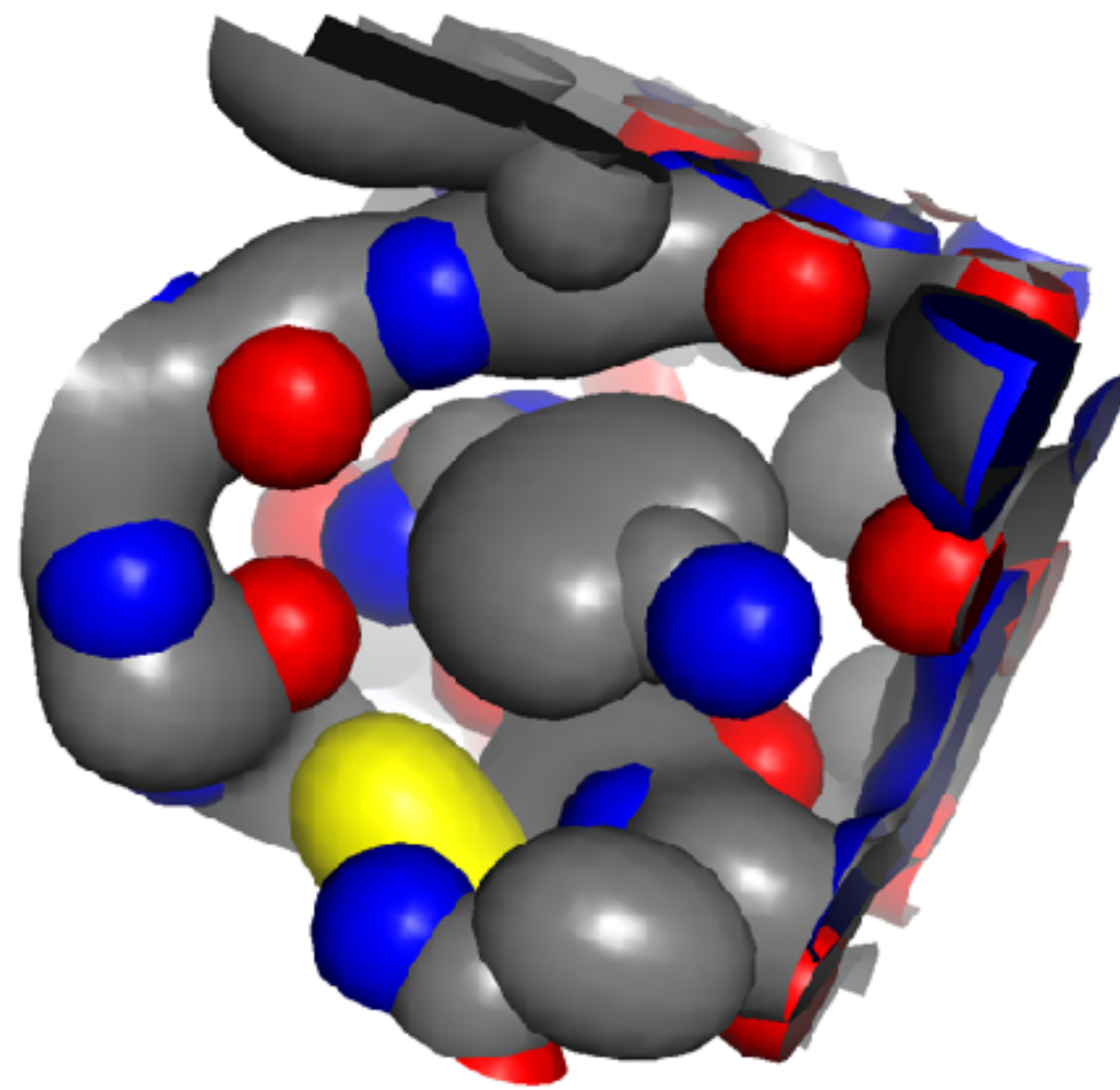
Grand Challenge 3

Spearman Correlation

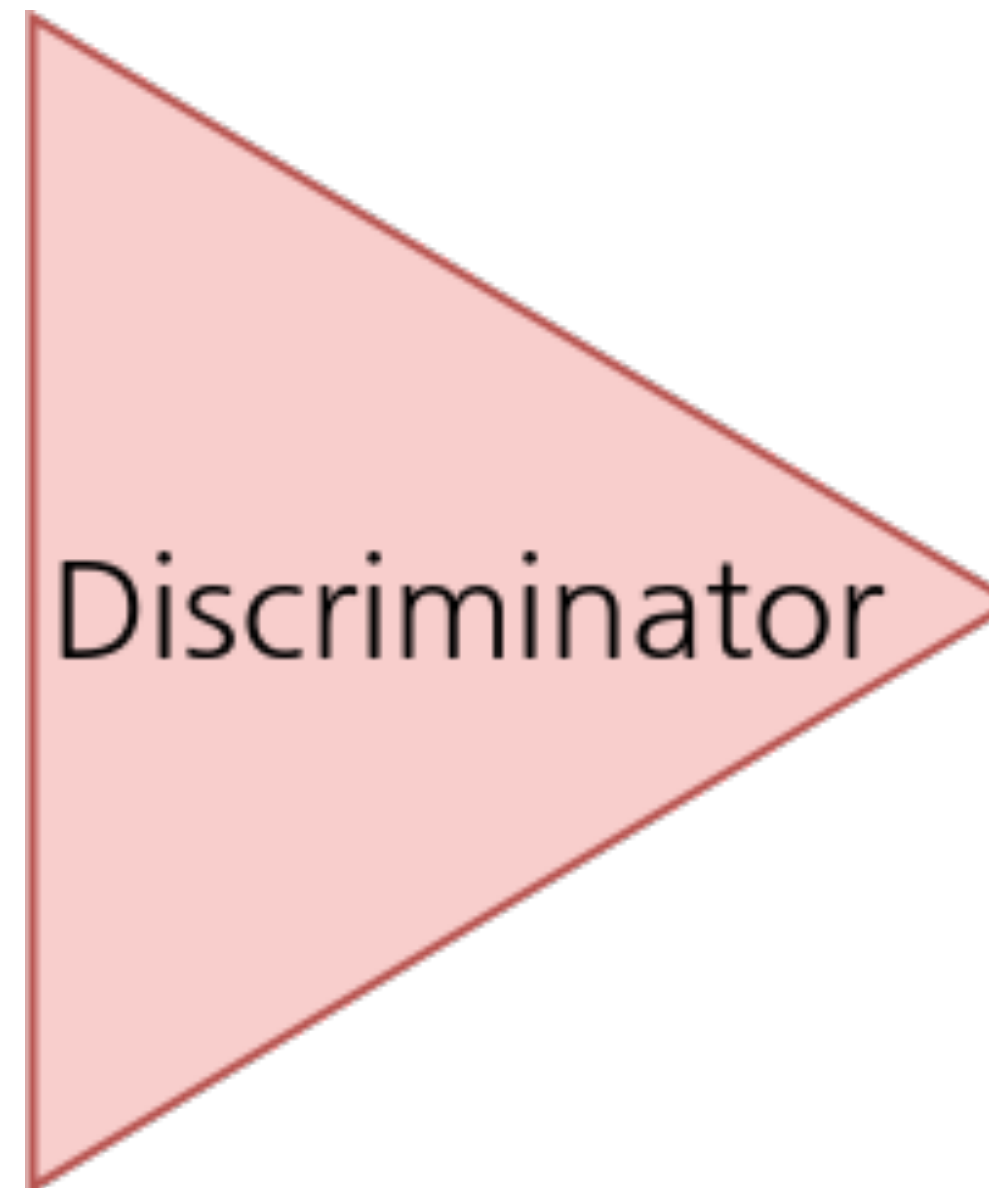
	cnn_docked_affinity	cnn_rescore_affinity	cnn_docked_scoring	cnn_rescore_scoring	vina
cat	0.0701	0.154	-0.0351	0.178	0.179
p38a	-0.0784	-0.116	-0.329	-0.305	-0.0631
vegfr2	0.366	0.484	0.434	0.448	0.414
jak2	0.428	0.338	0.39	0.27	0.106
jak2_sub3	0.68	0.369	-0.372	0.159	-0.633
tie2	0.648	0.835	0.136	-0.078	0.561
abl1	0.634	0.745	0.005	0.182	0.713

**and now for something
completely different...**

Discriminative Models



receptor & ligand grid

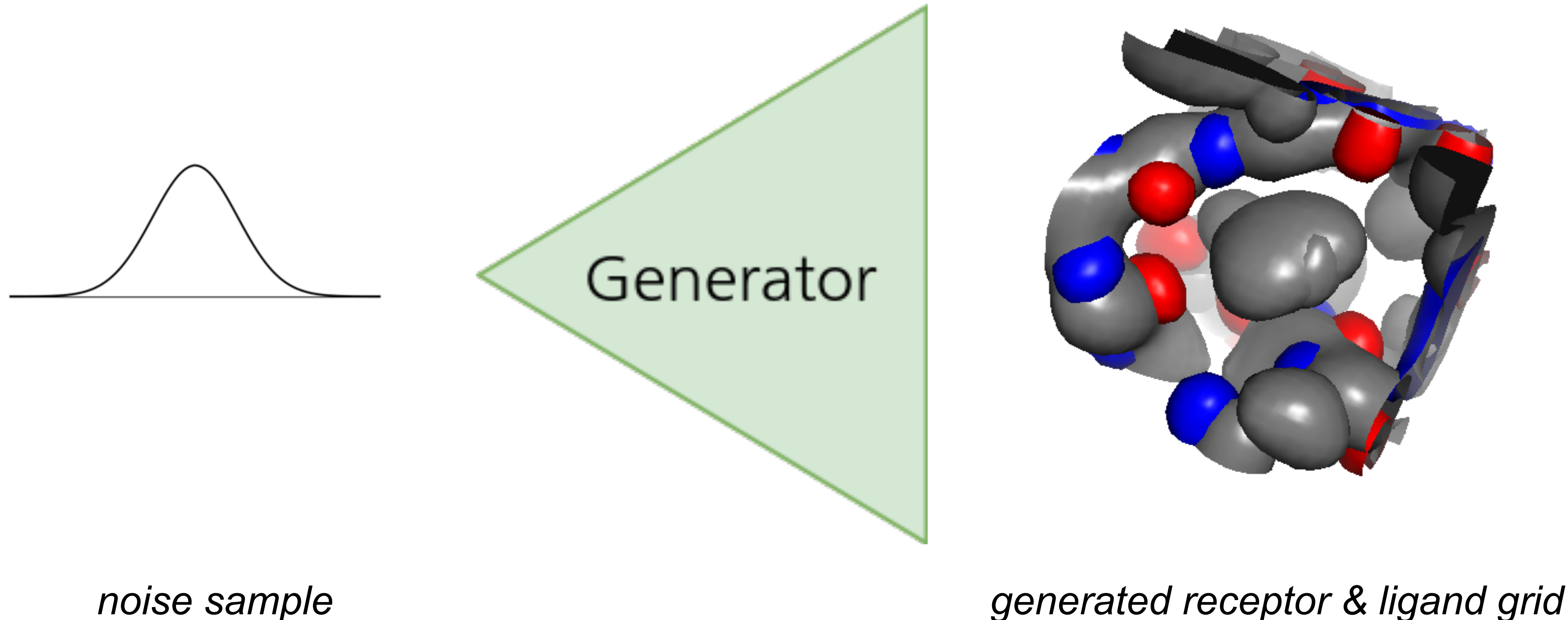


active/decoy

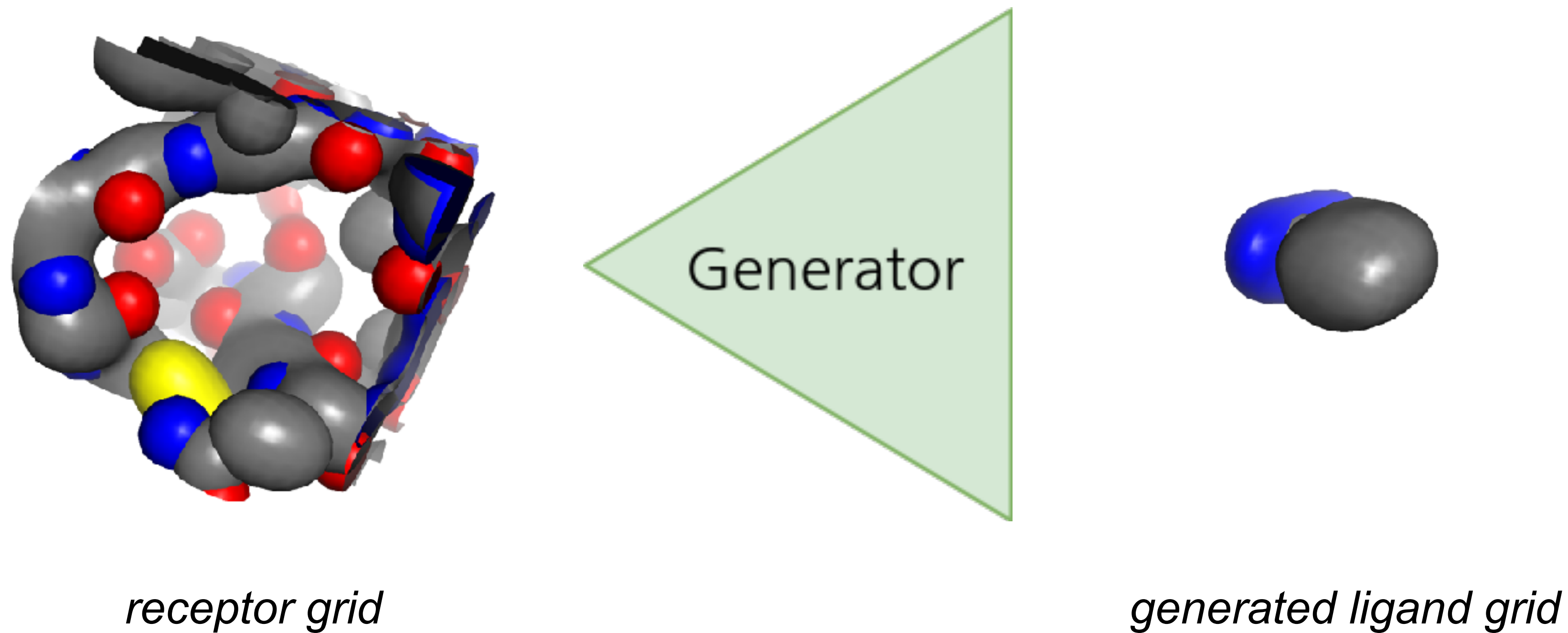
predicted class

Generative Models

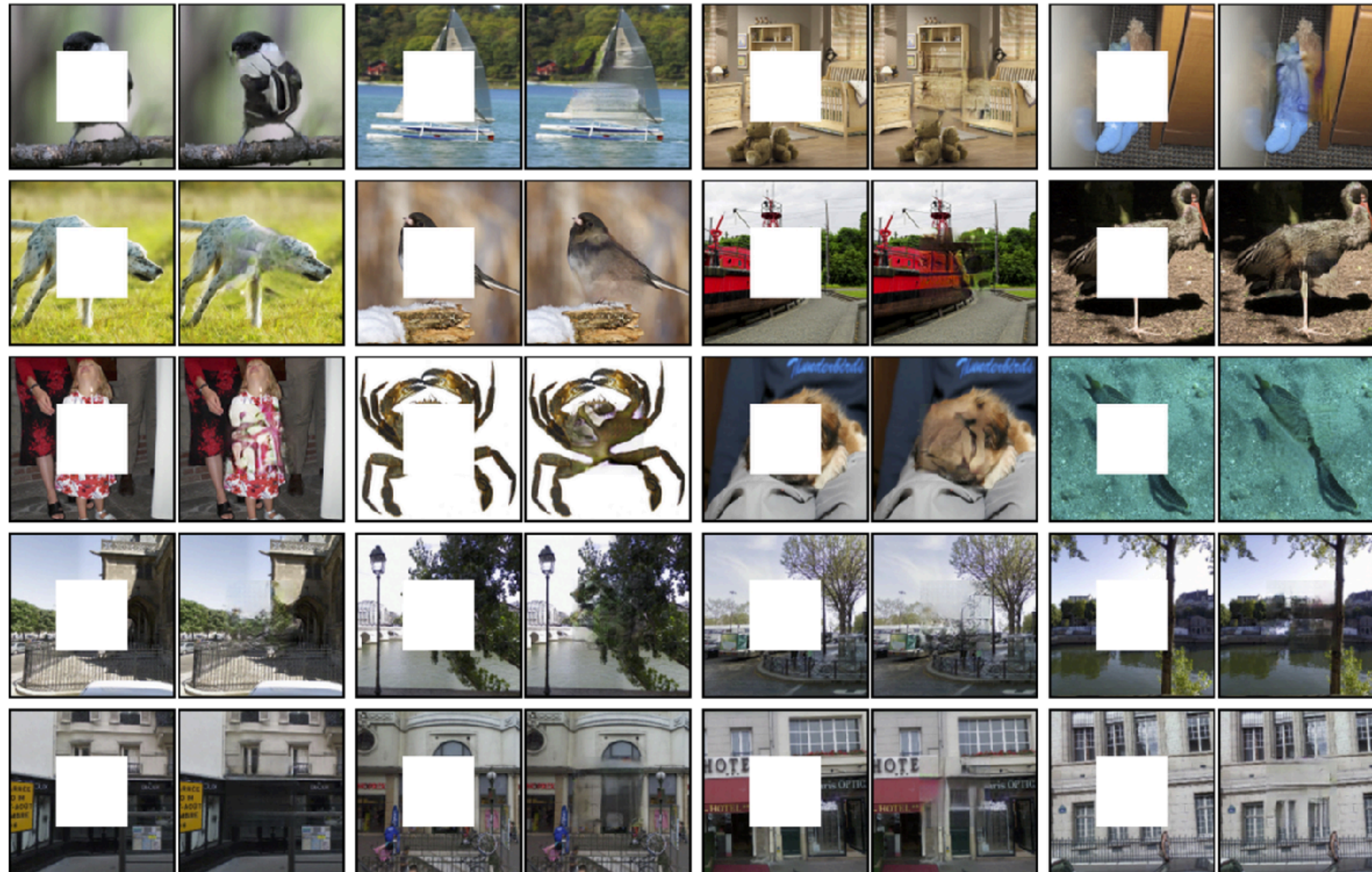
Generative models approximate a data distribution directly. They can map samples from one distribution (noise or input data) to realistic samples from an output distribution of interest.



Context Encoding



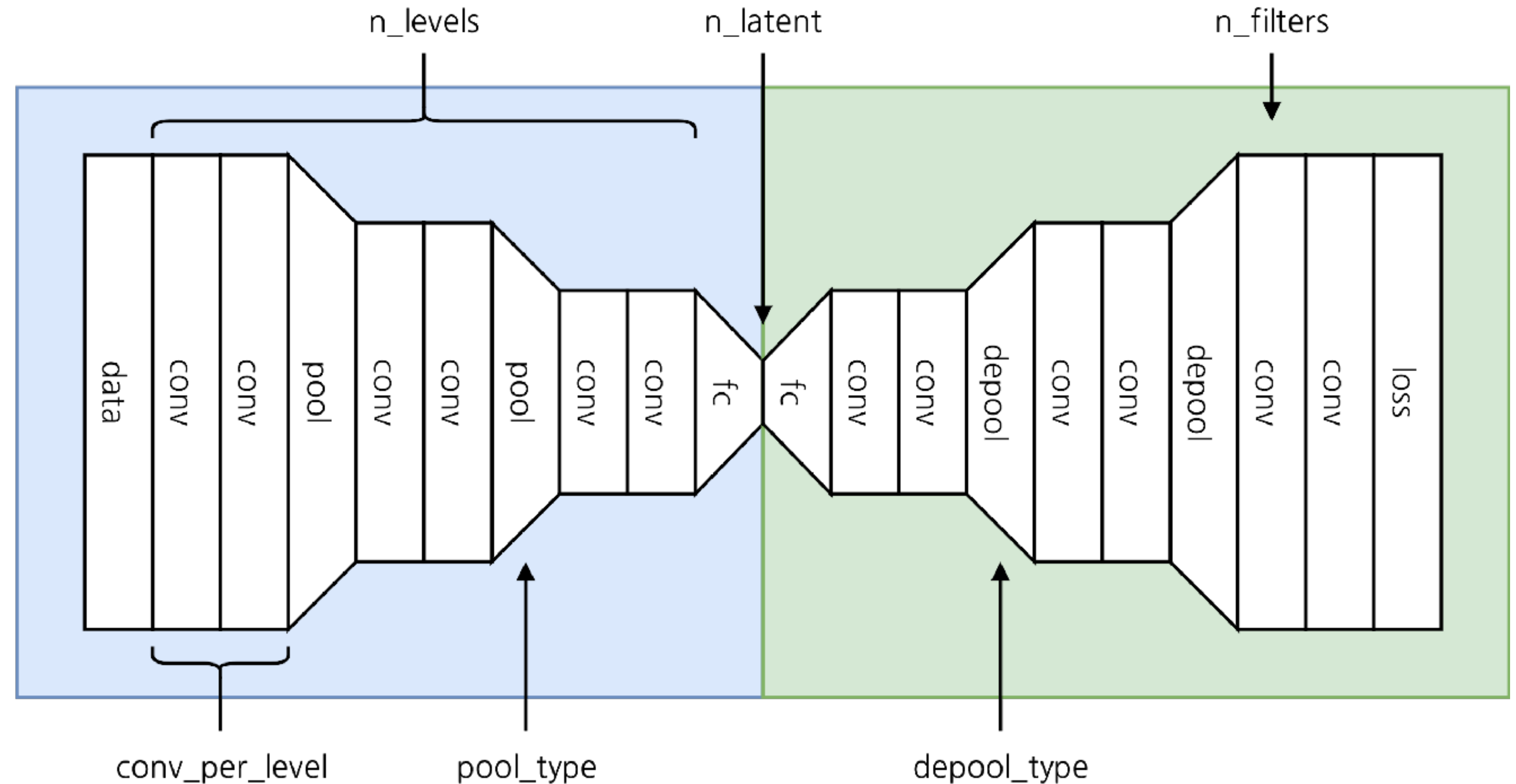
Context Encoding



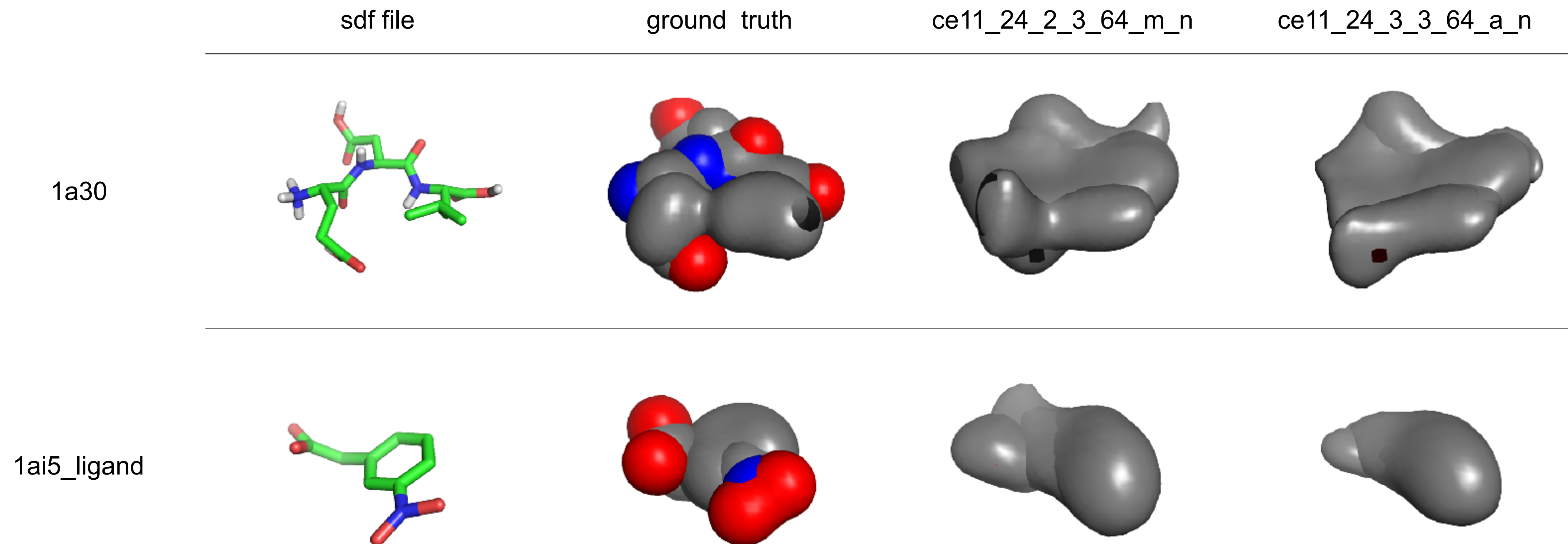
http://people.eecs.berkeley.edu/~pathak/context_encoder/

Model Architecture

- data_dim (24)
- resolution (0.5, 1.0)
- n_levels (3, 4, 5)
- conv_per_level (1, 2, 3)
- n_filters (16, 32, 64)
- width_factor (1, 2)
- n_latent (512, 1024)
- pool_type
 - max pooling
 - average pooling
 - strided convolution
- depool_type
 - nearest-neighbor
 - strided deconvolution
- loss_types
 - L2 loss



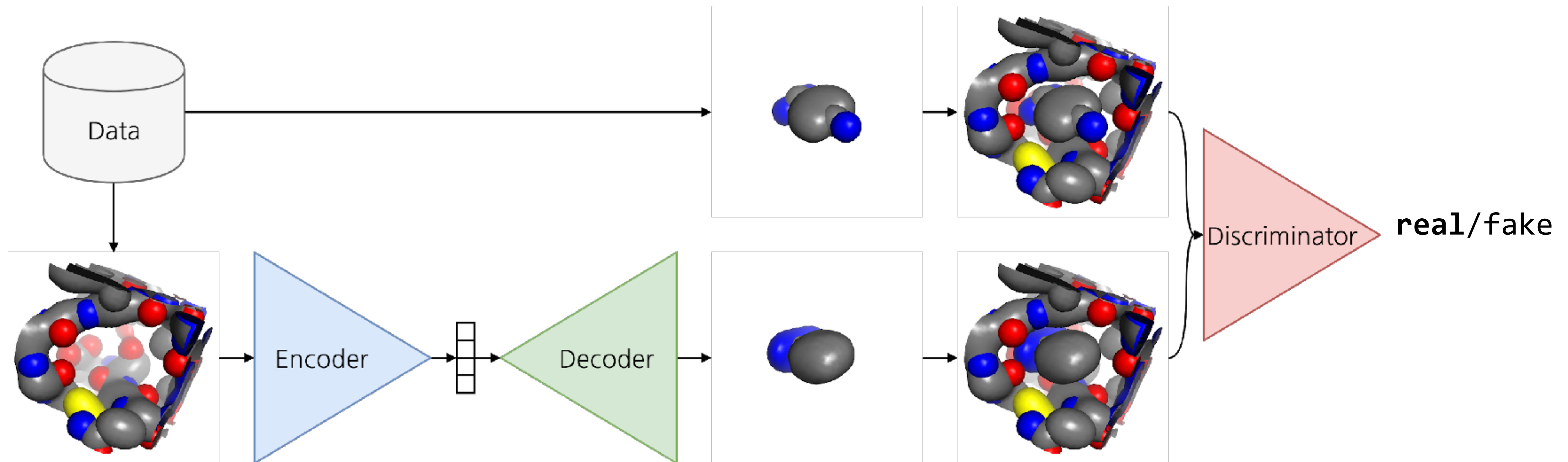
Context Encoder Examples



Generative Adversarial Networks

A discriminator network is trained to distinguish real vs. fake receptor-ligand grids

A generator network (context encoder) is trained to produce output that fools the discriminator



PROGRESSIVE GROWING OF GANs FOR IMPROVED QUALITY, STABILITY, AND VARIATION

Tero Karras
NVIDIA

Timo Aila
NVIDIA

Samuli Laine
NVIDIA

Jaakko Lehtinen
NVIDIA
Aalto University



<https://youtu.be/G06dEcZ-QTg>

PROGRESSIVE GROWING OF GANs FOR IMPROVED QUALITY, STABILITY, AND VARIATION

Tero Karras
NVIDIA

Timo Aila
NVIDIA

Samuli Laine
NVIDIA

Jaakko Lehtinen
NVIDIA
Aalto University



<https://youtu.be/G06dEcZ-QTg>

Preliminary GAN Examples

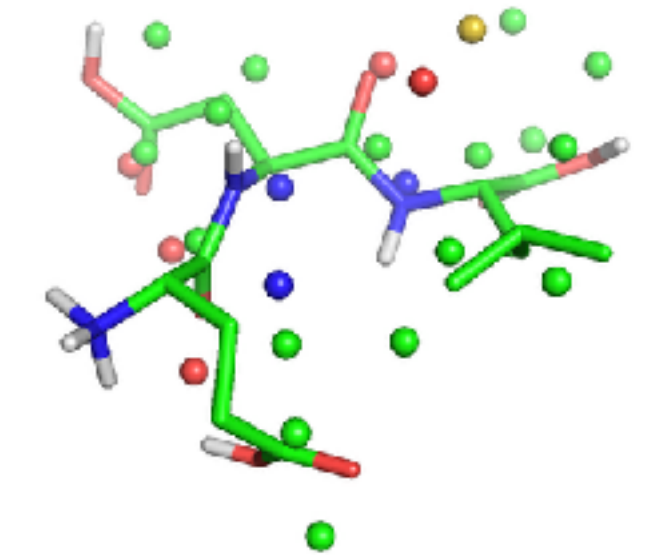
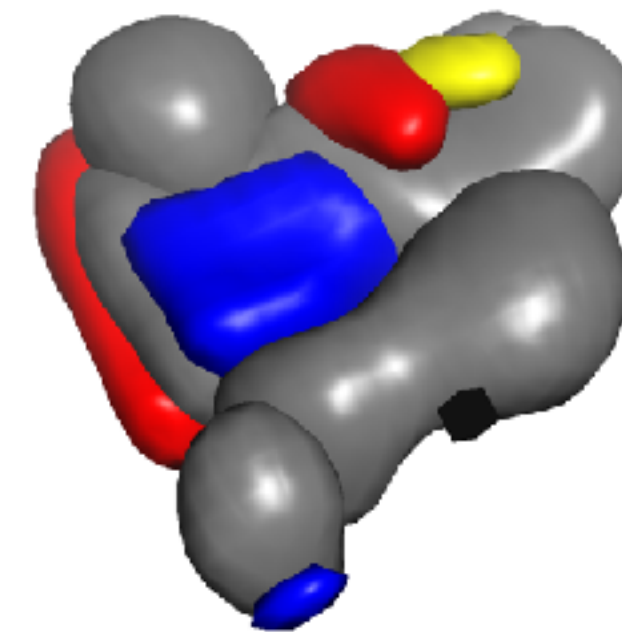
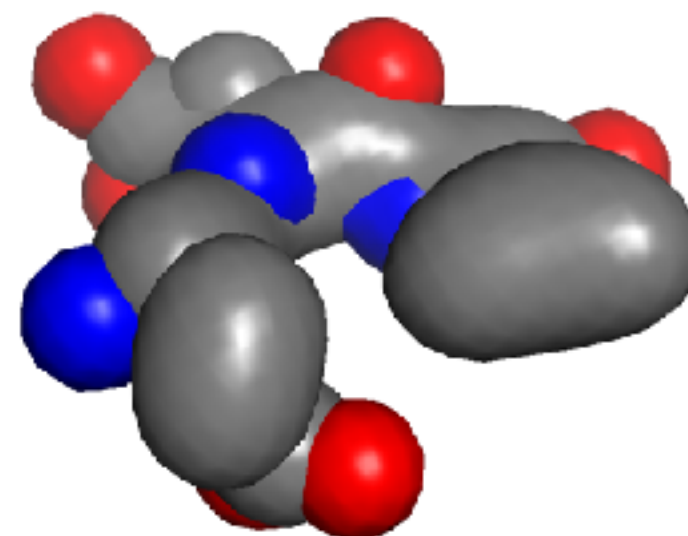
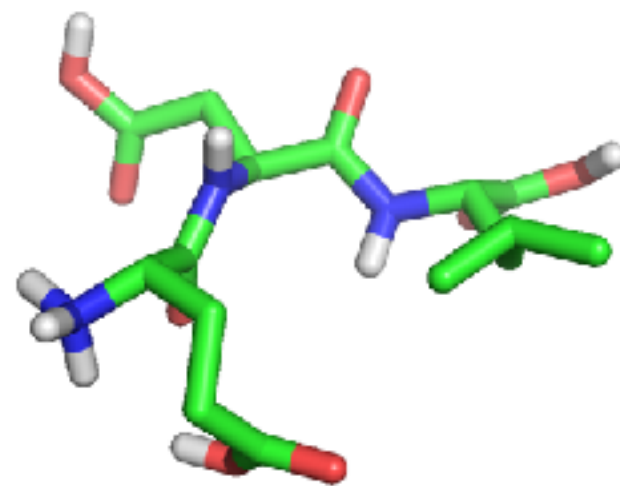
ligand

ground truth

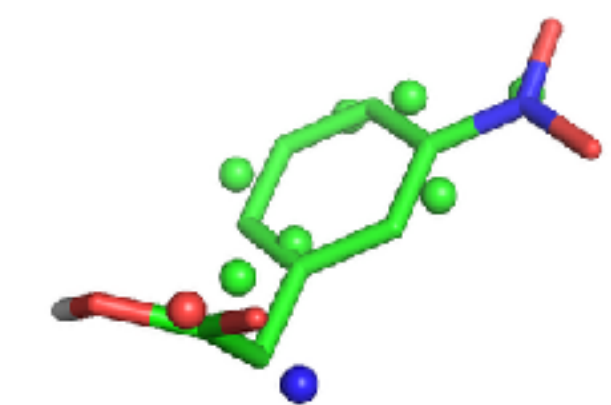
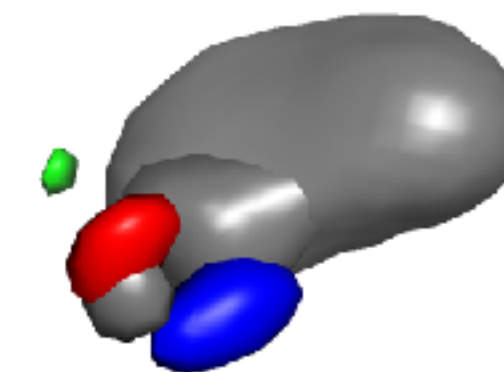
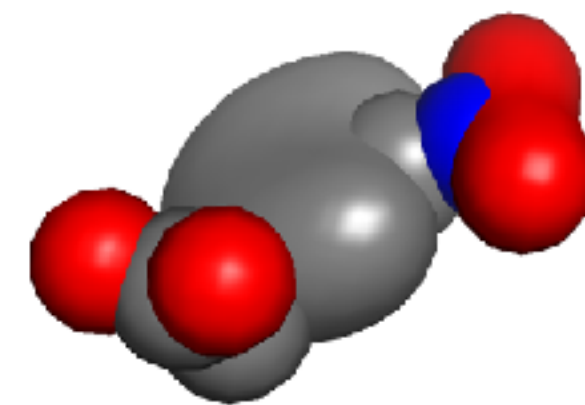
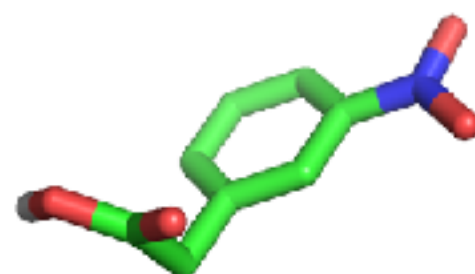
GAN

fit atoms

1a30



1ai5

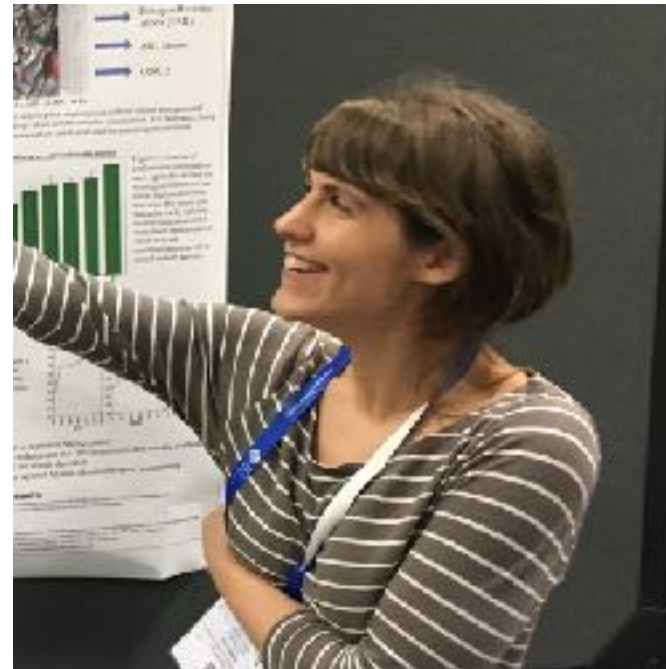


Preliminary GAN Examples

	ligand	ground truth	GAN	fit atoms
1a30				
1ai5				

<http://torch.ch/blog/2015/11/13/gan.html>

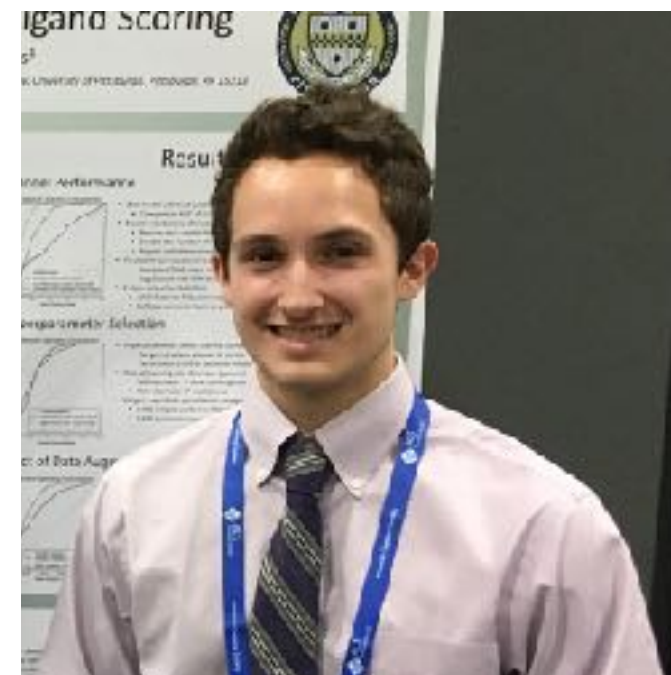
Acknowledgements



Jocelyn Sunseri



Josh Hochuli



Matt Ragoza

Group Members

Jocelyn Sunseri

Jonathan King

Paul Francoeur

Matt Ragoza

Josh Hochuli

Pulkit Mittal

Alec Helbling

Gibran Biswas

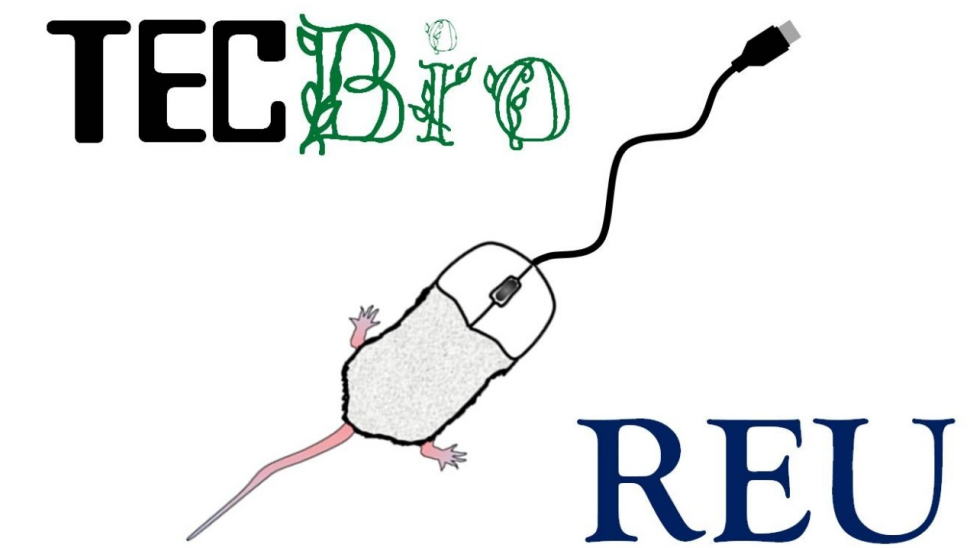
Sharanya Bandla

Faiha Khan

Lily Turner



Department of
Computational and
Systems Biology



AI GRANT



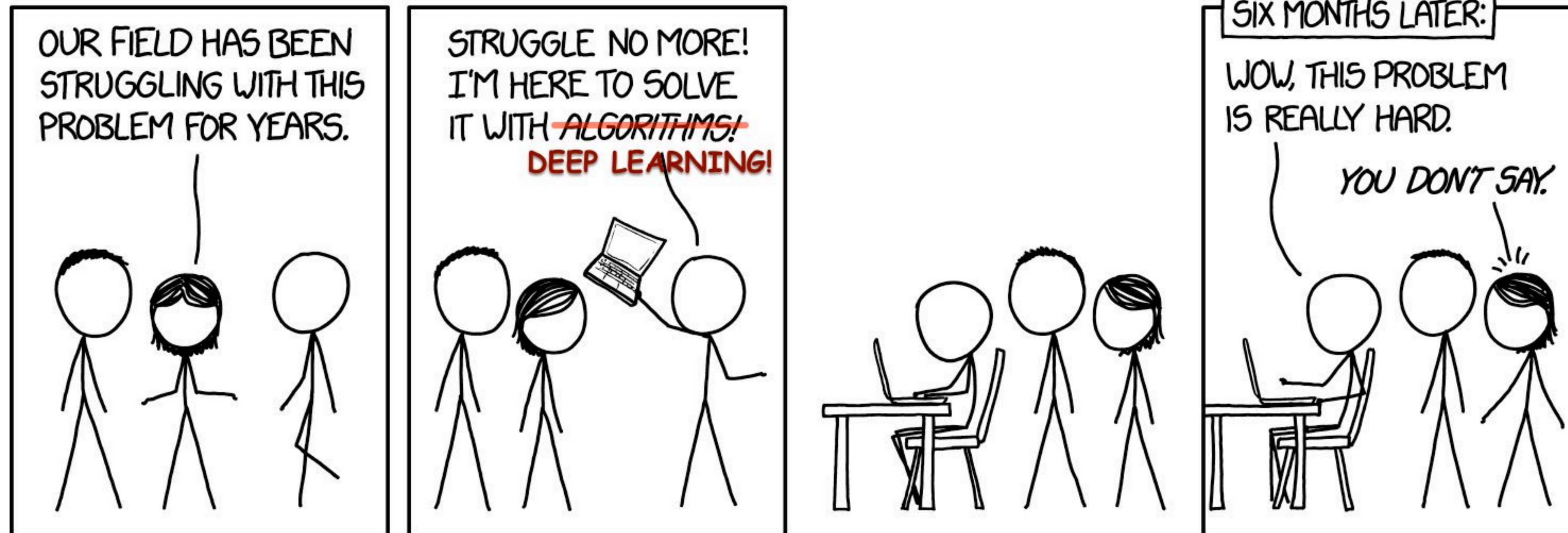
National Institute of
General Medical Sciences
[R01GM108340](#)



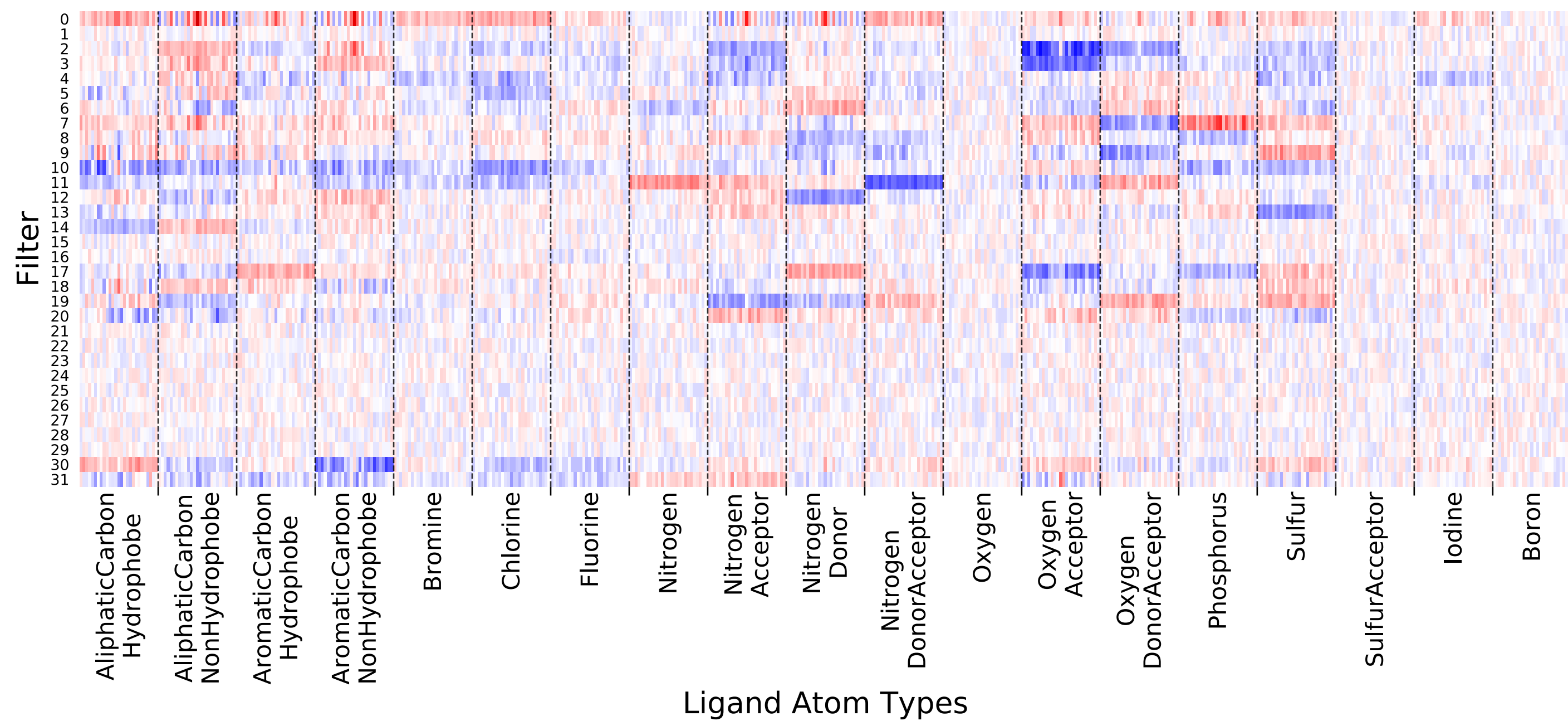
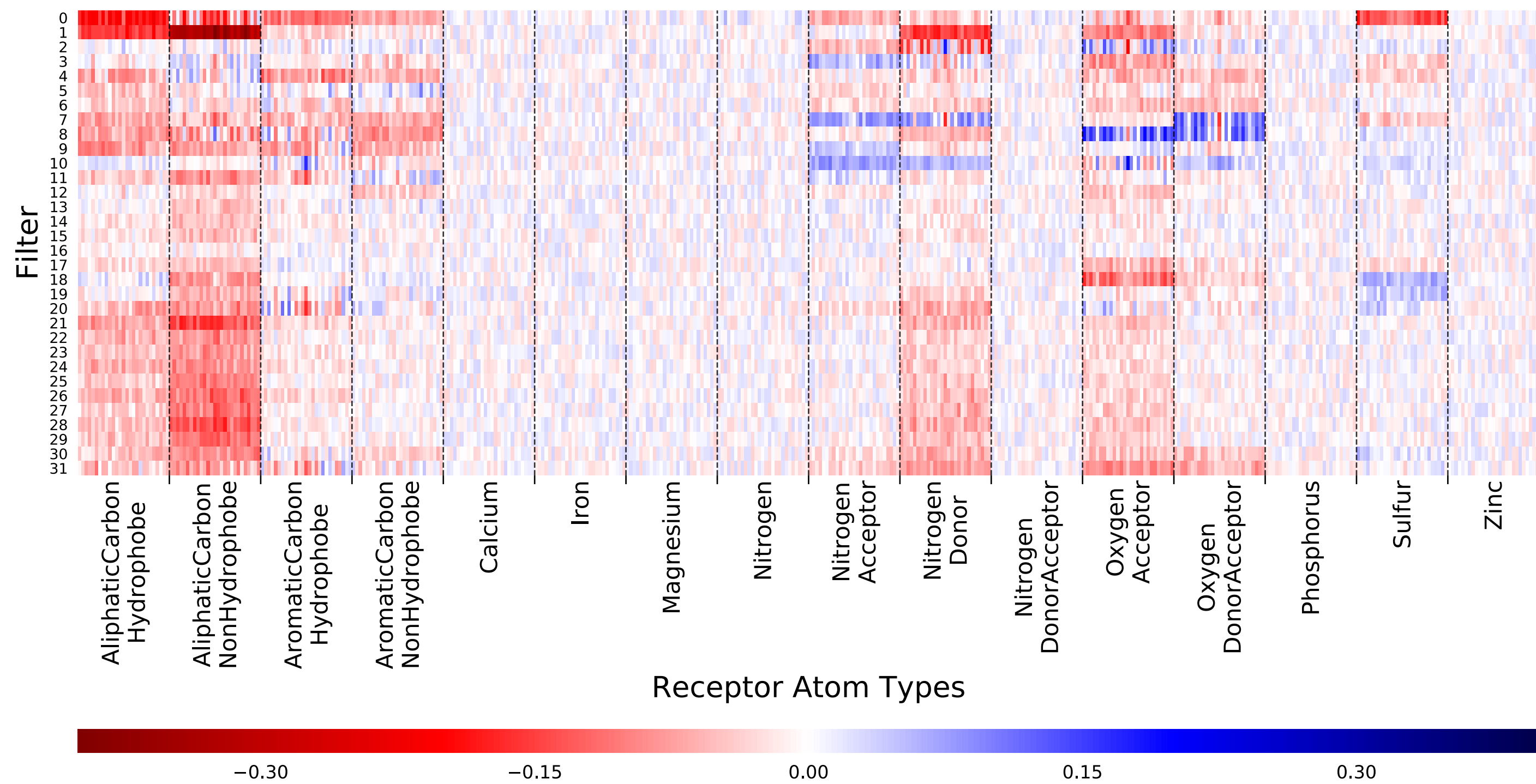
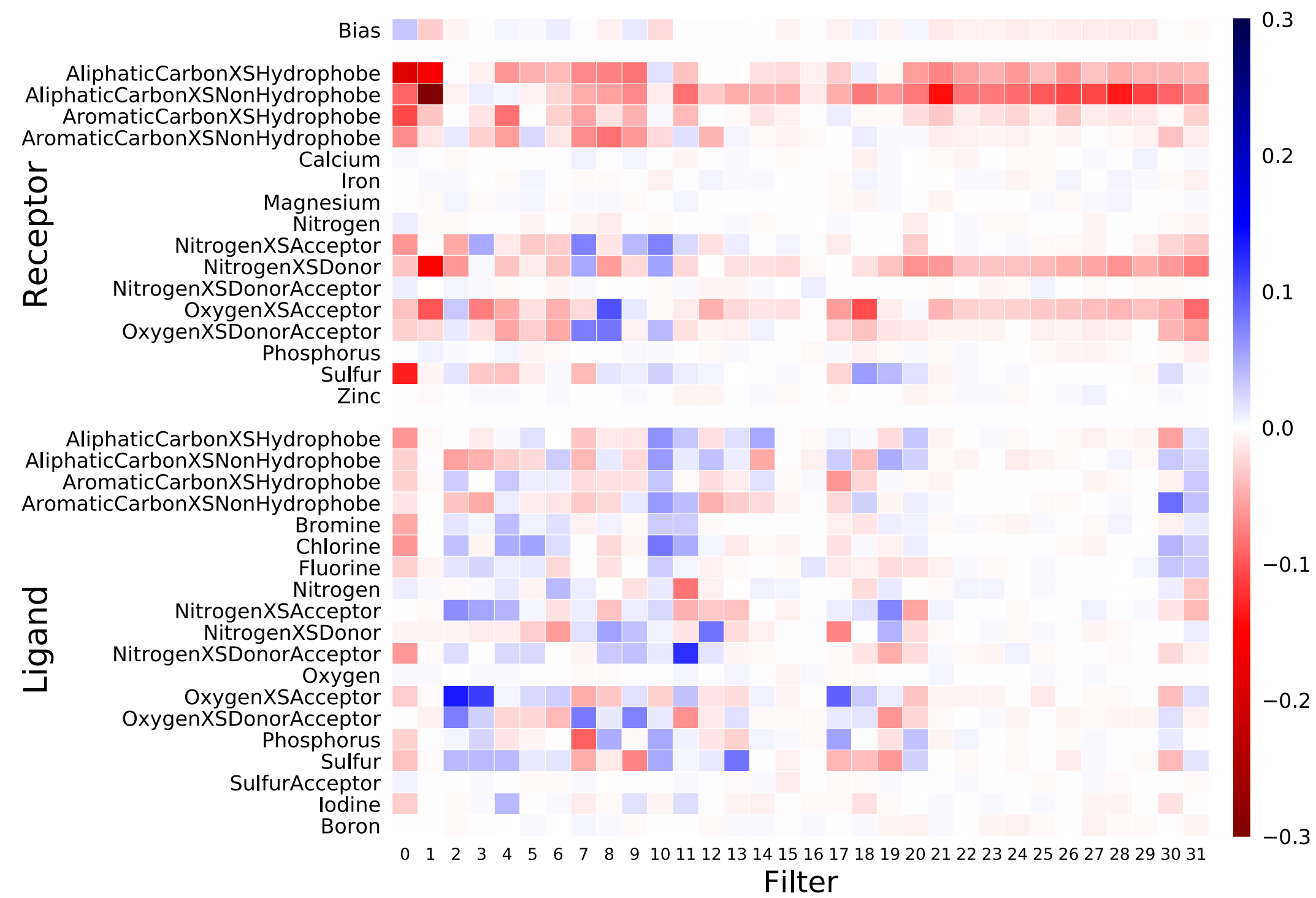
 github.com/gnina

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

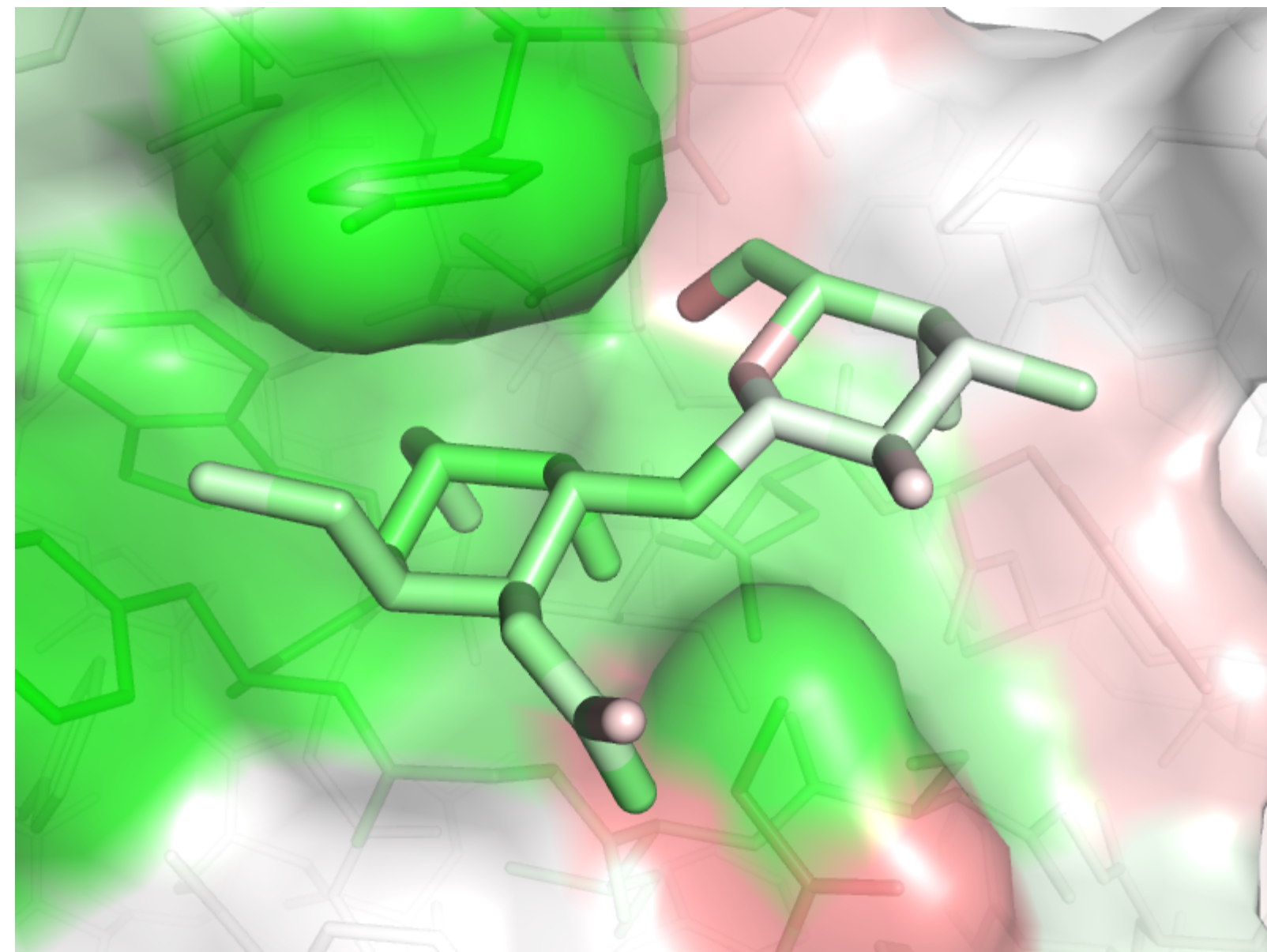
 @david_koes



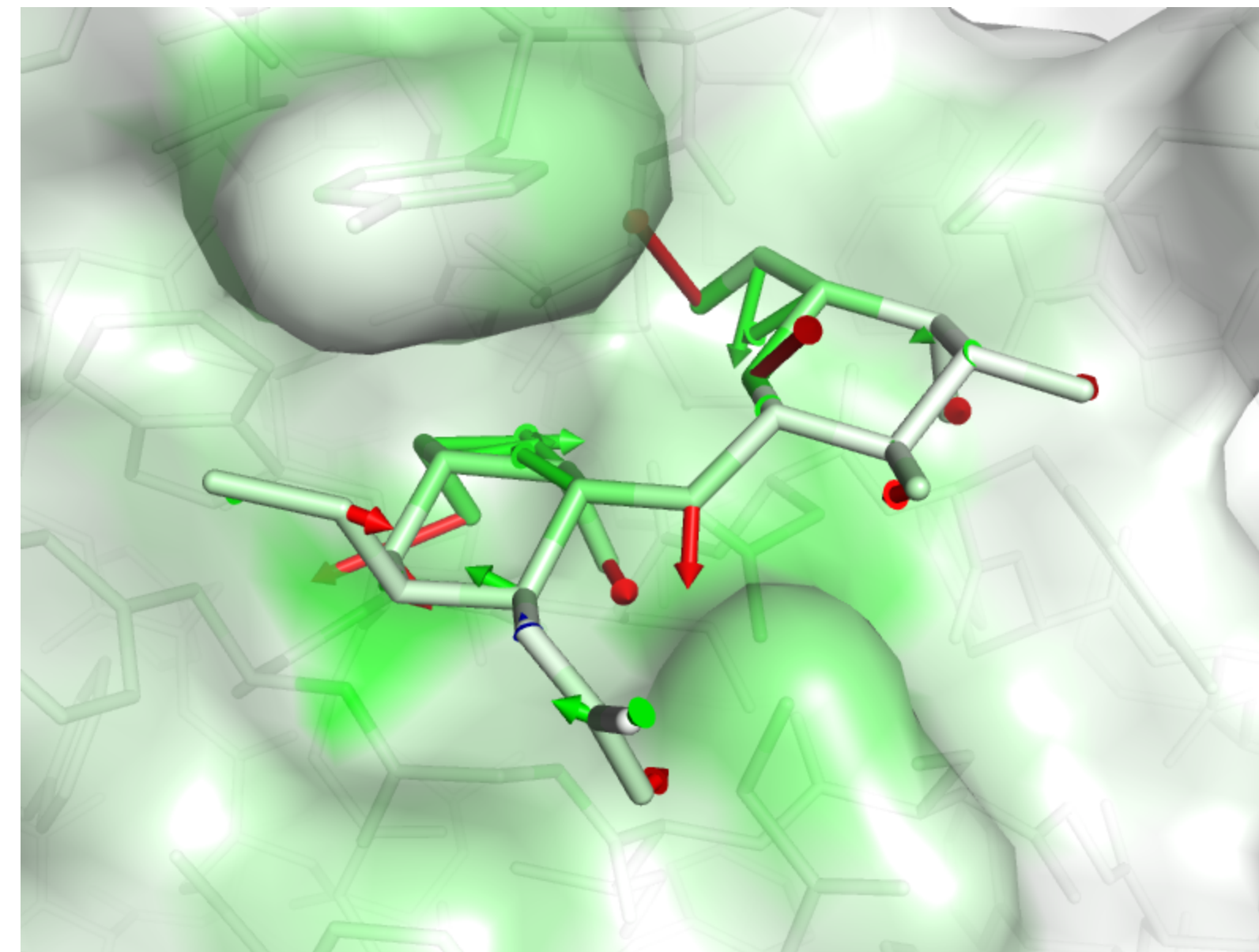
Filter Visualization



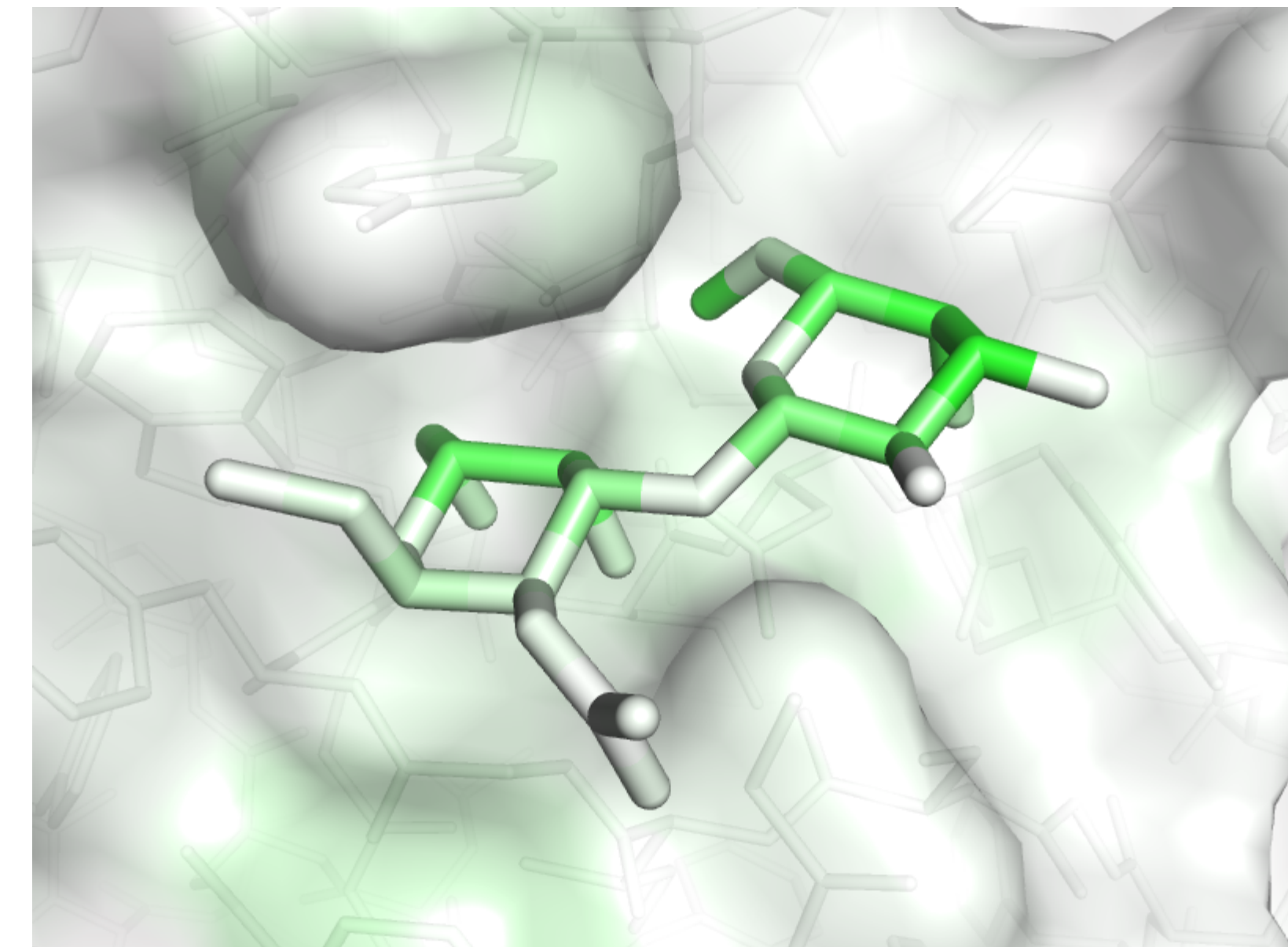
Visualization



masking



gradients



layer-wise relevance

