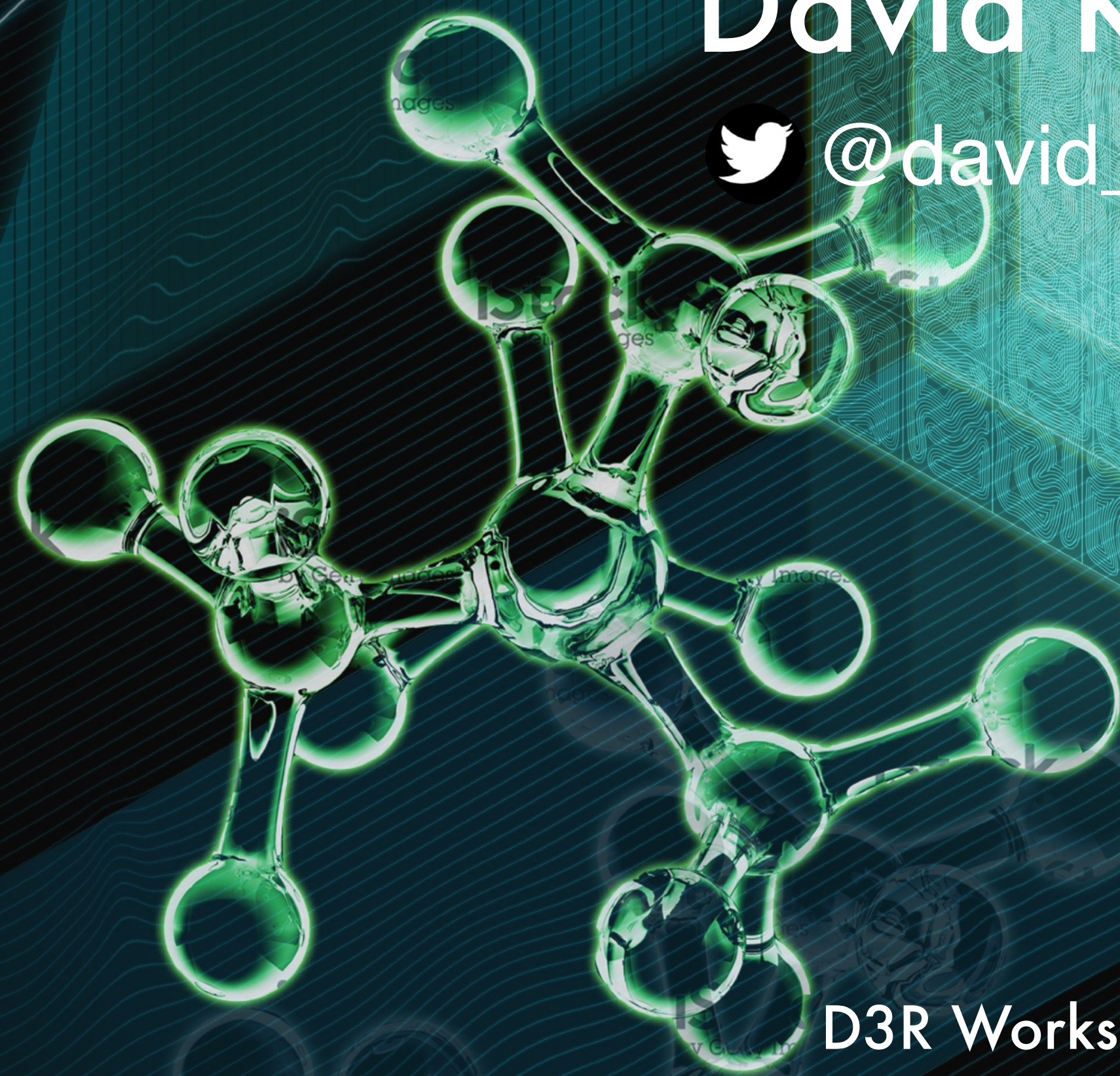# Protein-Ligand Scoring with Convolutional Neural Networks

## David Koes

@david_koes
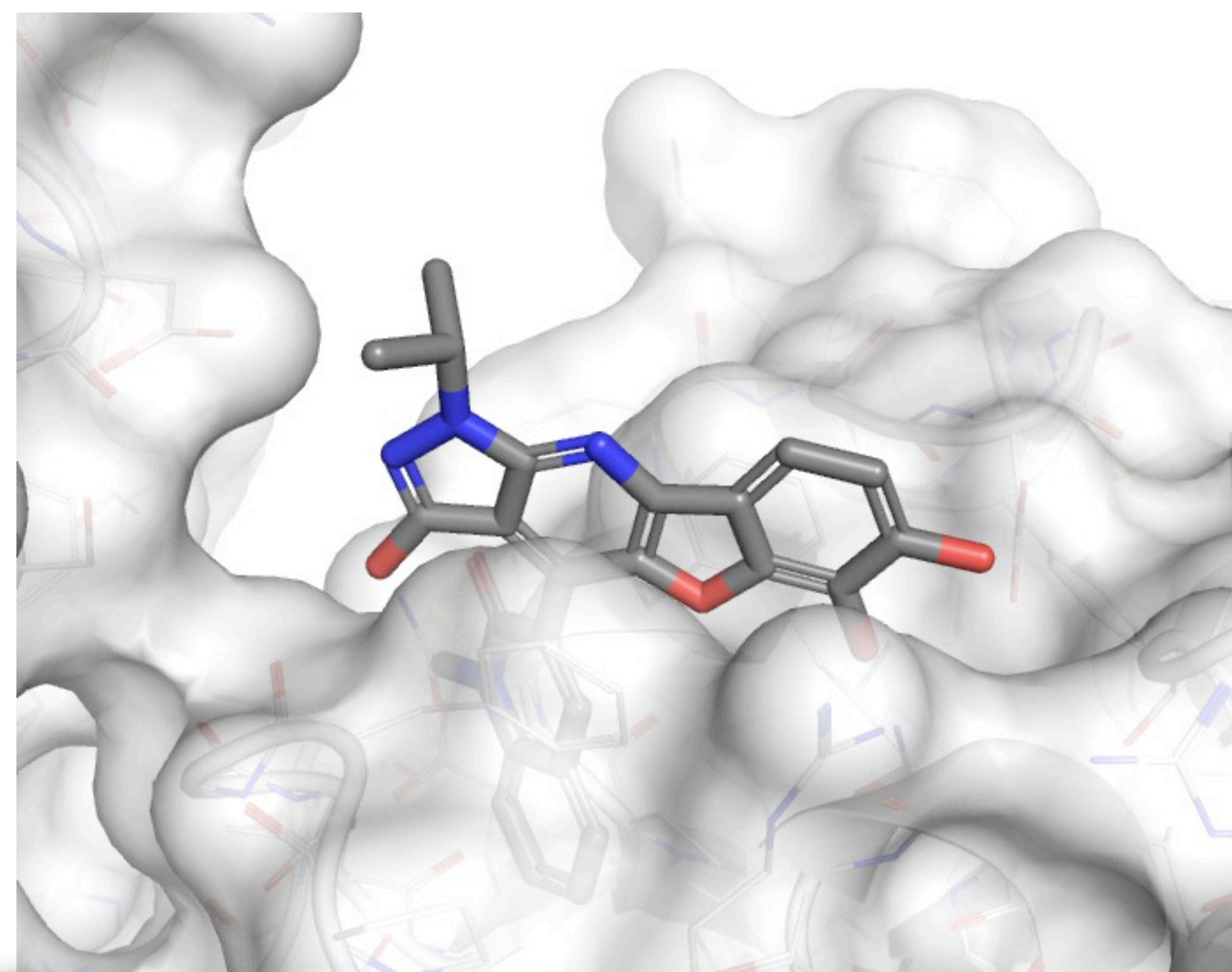
D3R Workshop
San Diego
February 22, 2018
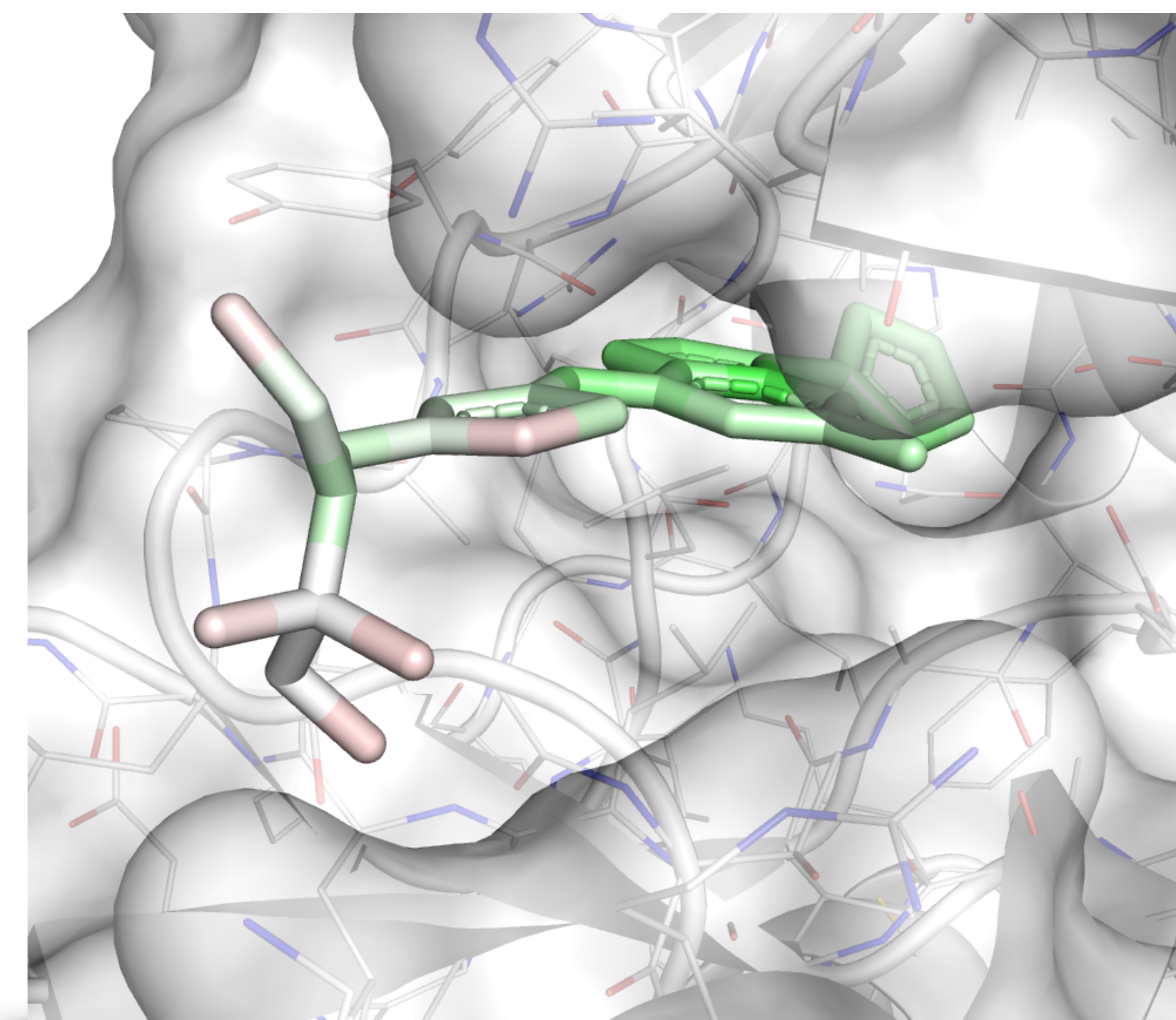
# Structure Based Drug Design
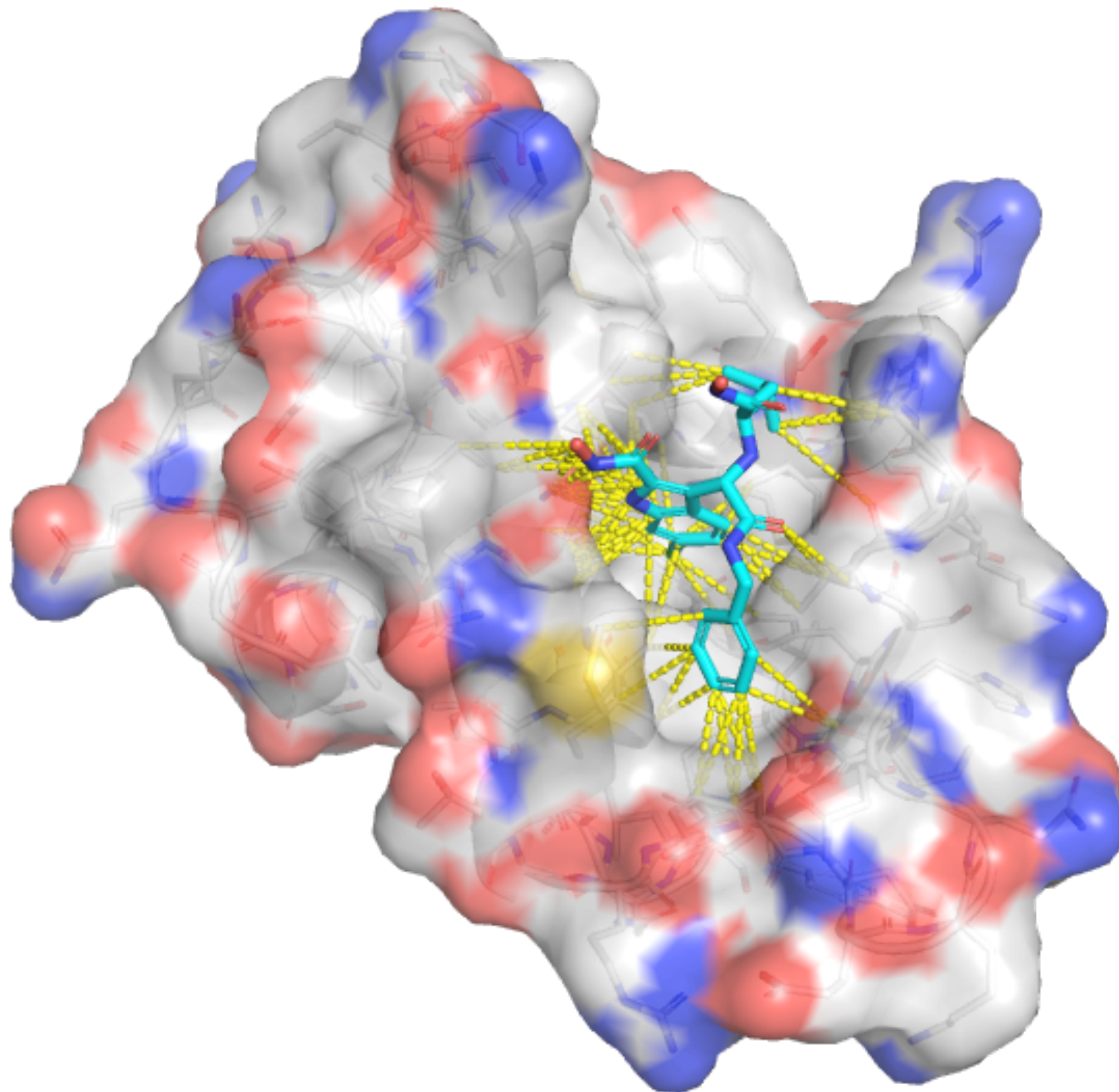
**Virtual Screening**

**Lead Optimization**





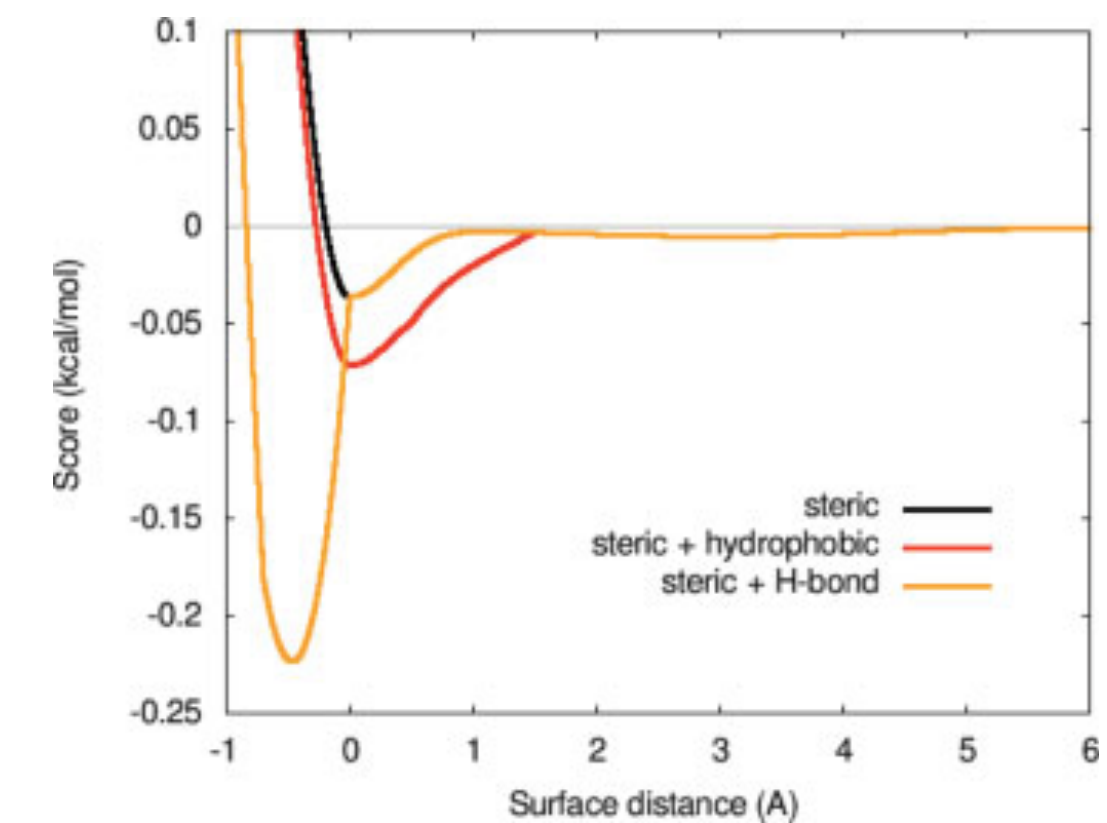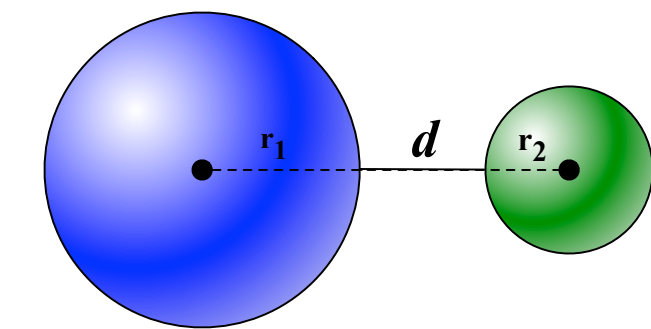**Pose Prediction** | Binding Discrimination | **Affinity Prediction**

# Protein-Ligand Scoring

## AutoDock Vina



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

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

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

3

# Can we do better?

Accurate pose prediction, binding discrimination, **and** affinity prediction without sacrificing performance?
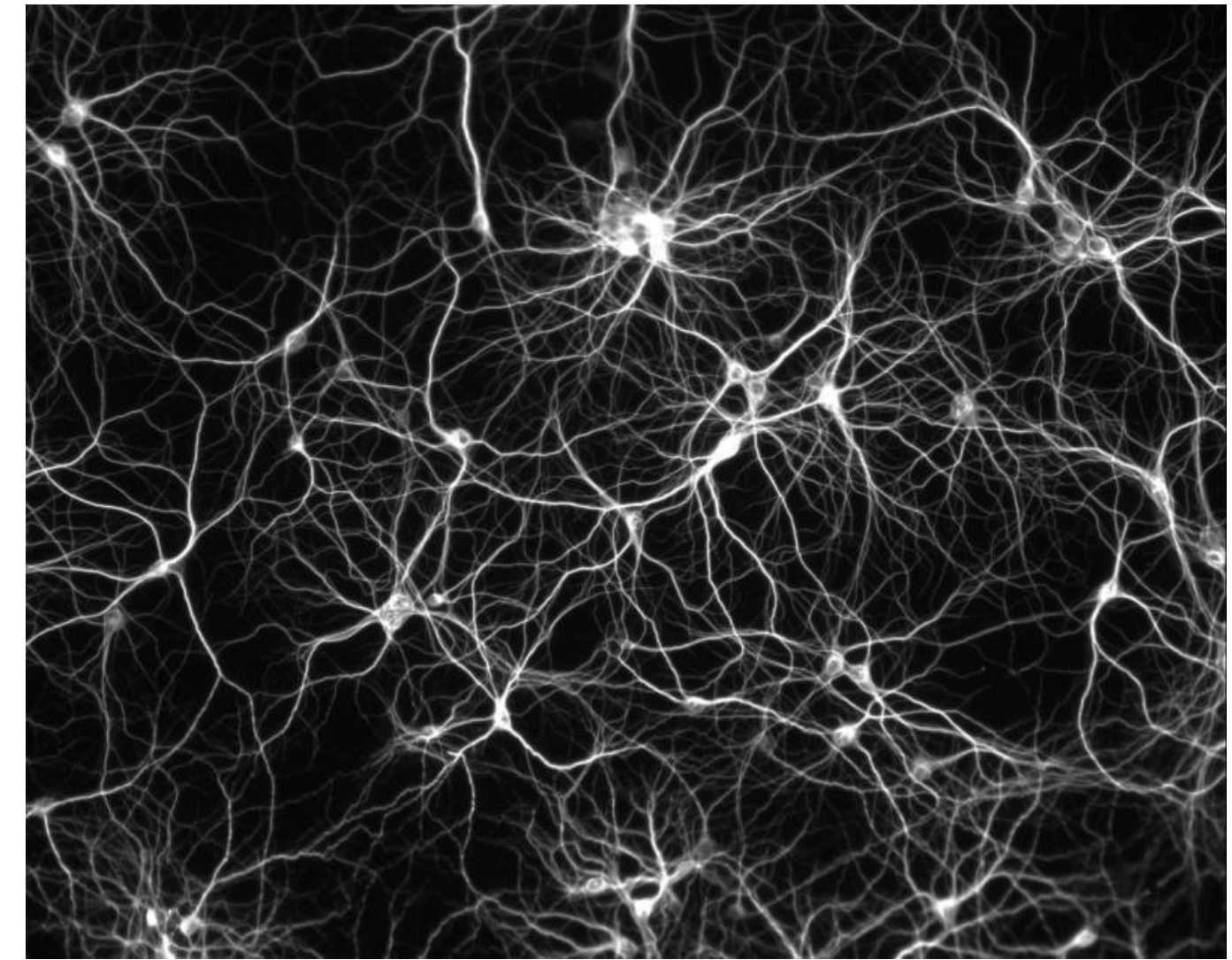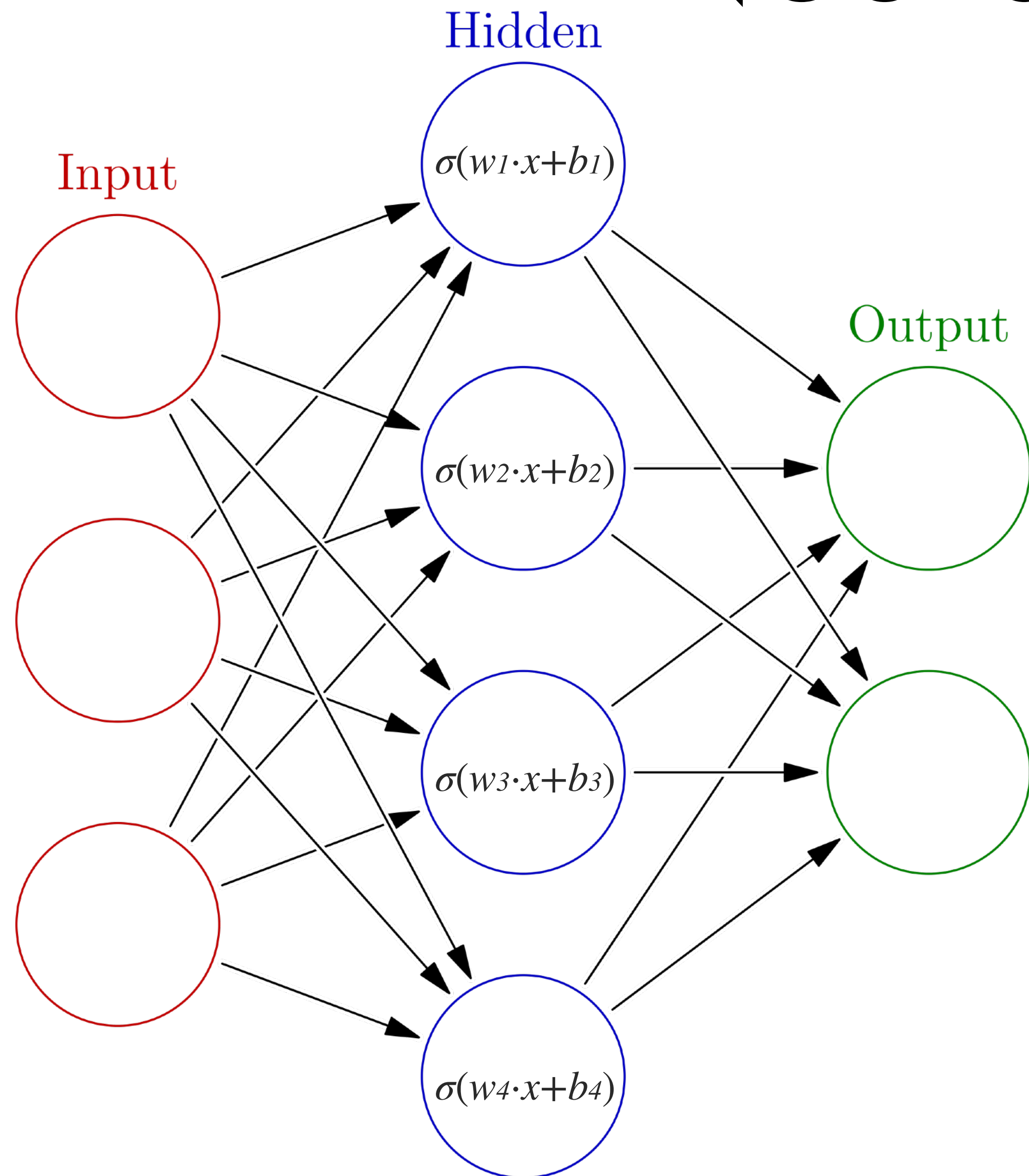
**Key Idea:** Leverage "big data"
- 231,655,275 bioactivities in PubChem
- 125,526 structures in the PDB
- 16,179 annotated complexes in PDBbind
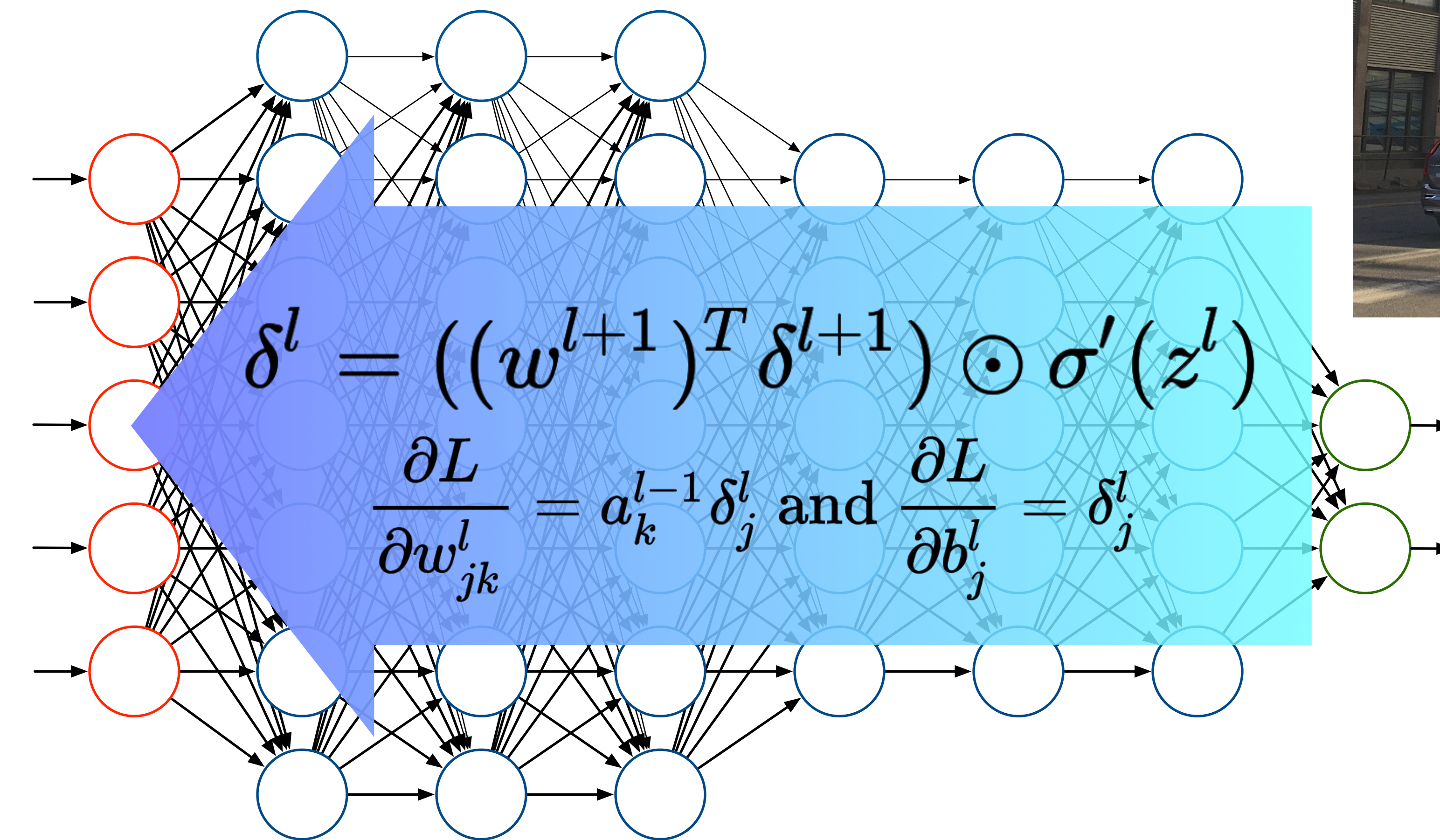
# Machine Learning

Features $X \rightarrow$ **Model** $\rightarrow y$ Prediction

# Neural Networks

Hidden

Input

$\sigma(w_1 \cdot x + b_1)$

Output

$\sigma(w_2 \cdot x + b_2)$

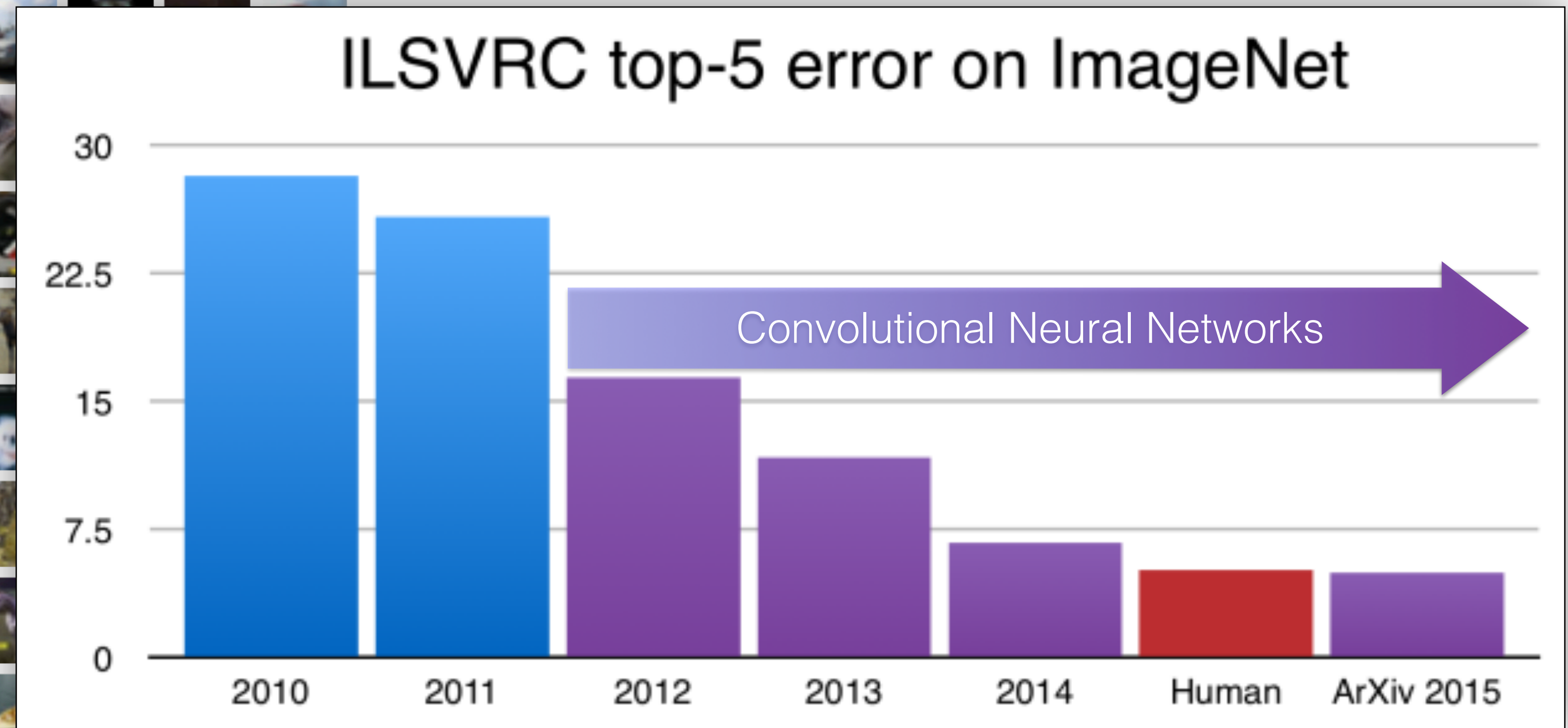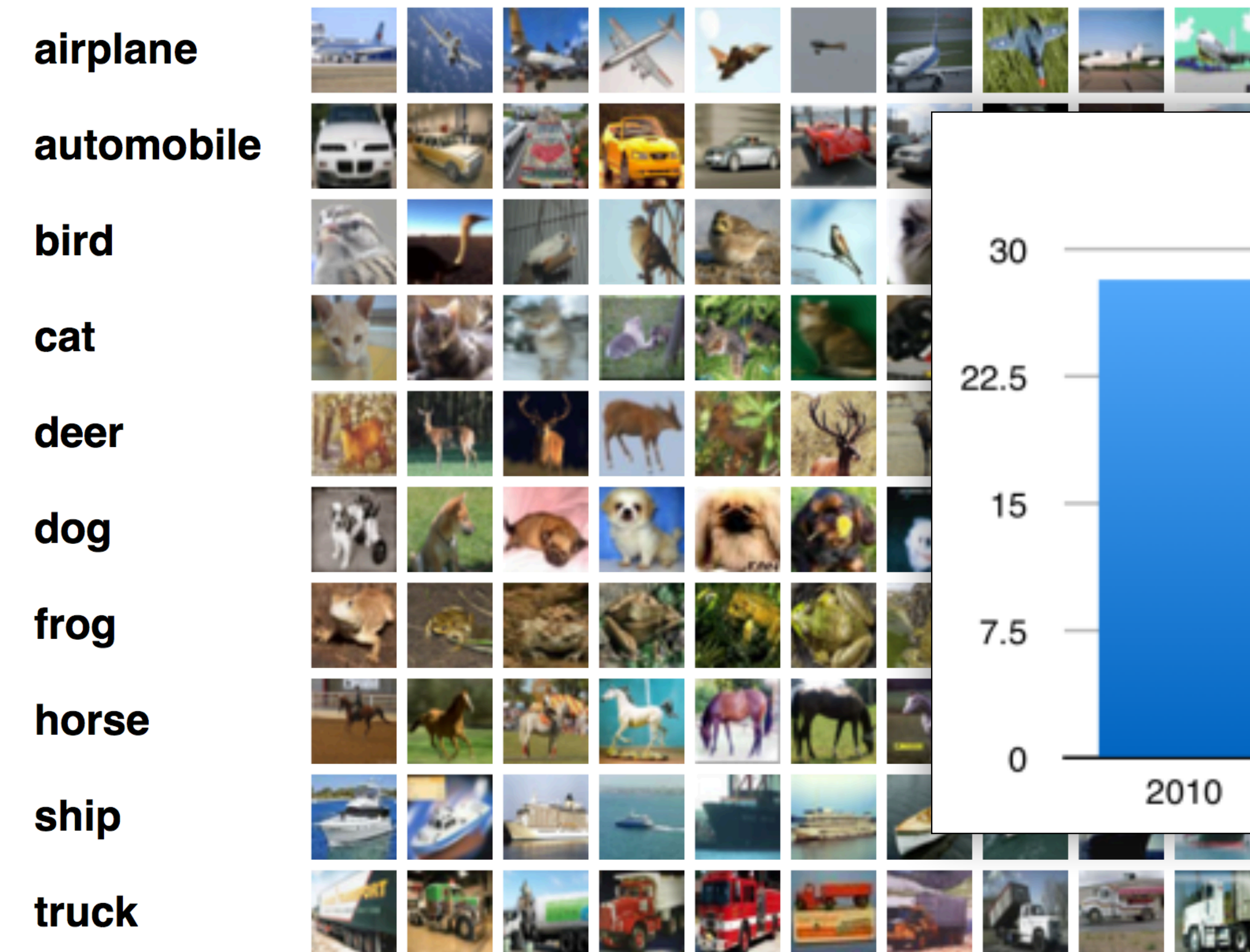$\sigma(w_3 \cdot x + b_3)$

$\sigma(w_4 \cdot x + b_4)$

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.
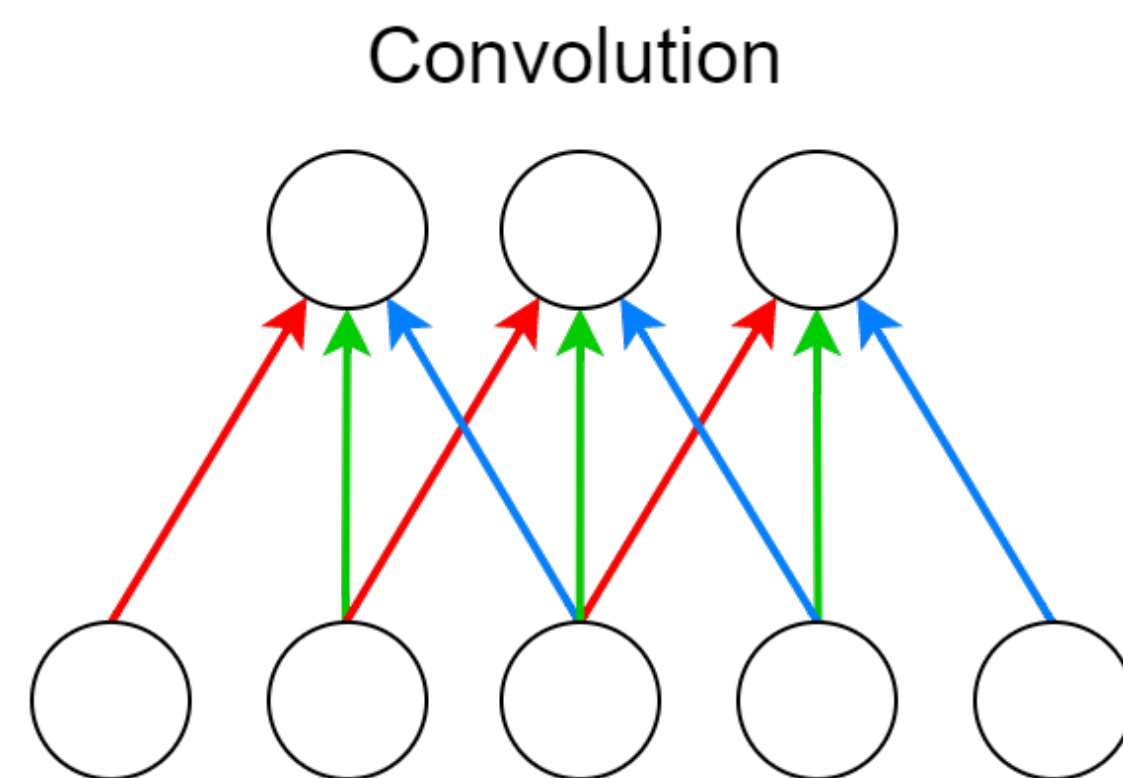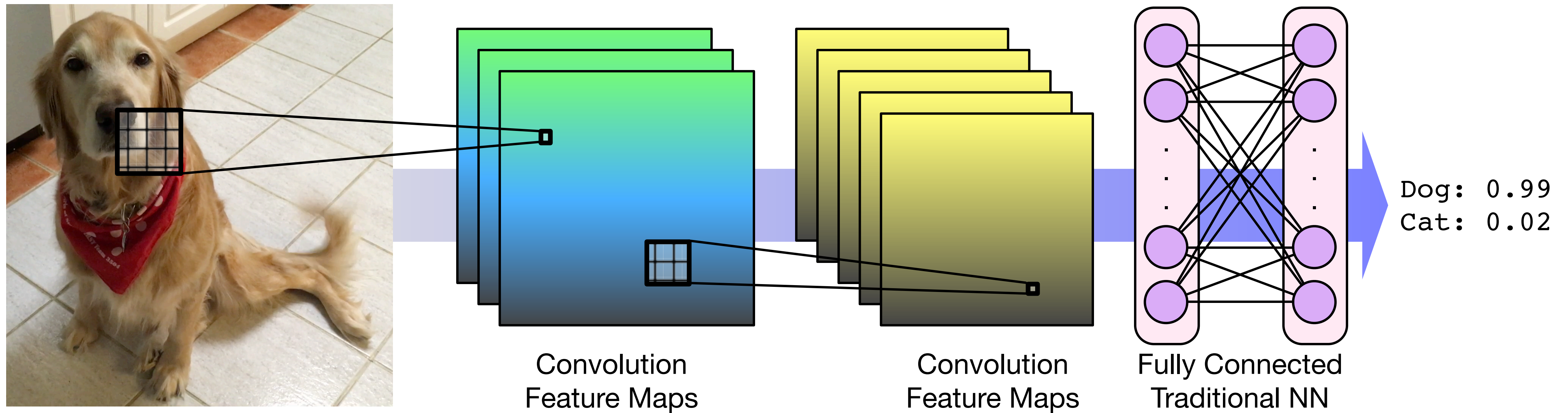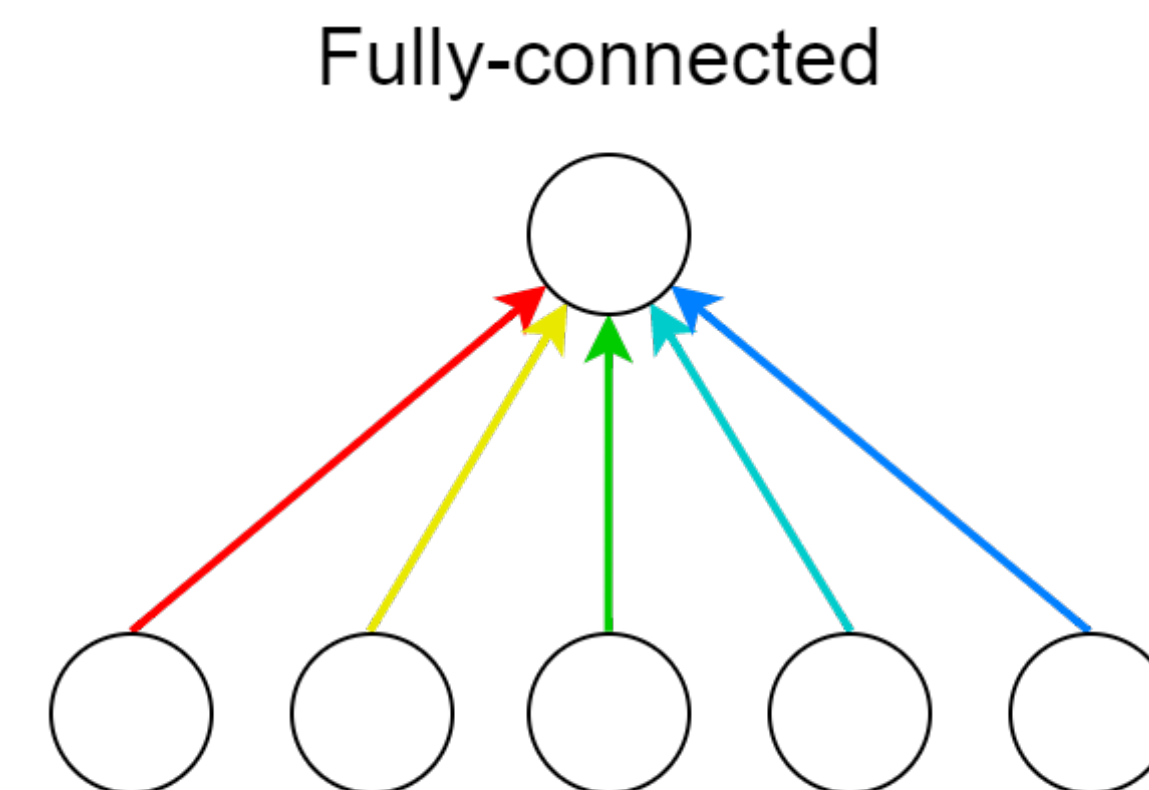
6

# Deep Learning



$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

$$\frac{\partial L}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial L}{\partial b_j^l} = \delta_j^l$$

# Image Recognition



airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

## ILSVRC top-5 error on ImageNet

Convolutional Neural Networks

2010  2011  2012  2013  2014  Human  ArXiv 2015

https://devblogs.nvidia.com

# Convolutional Neural Networks



Convolution
Feature Maps

Convolution
Feature Maps

Fully Connected
Traditional NN

Dog: 0.99
Cat: 0.02

Convolution

weight 1
weight 2
weight 3

Fully-connected

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



**CNN**

Pose Prediction

Binding
Discrimination

Affinity Prediction

# 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**.

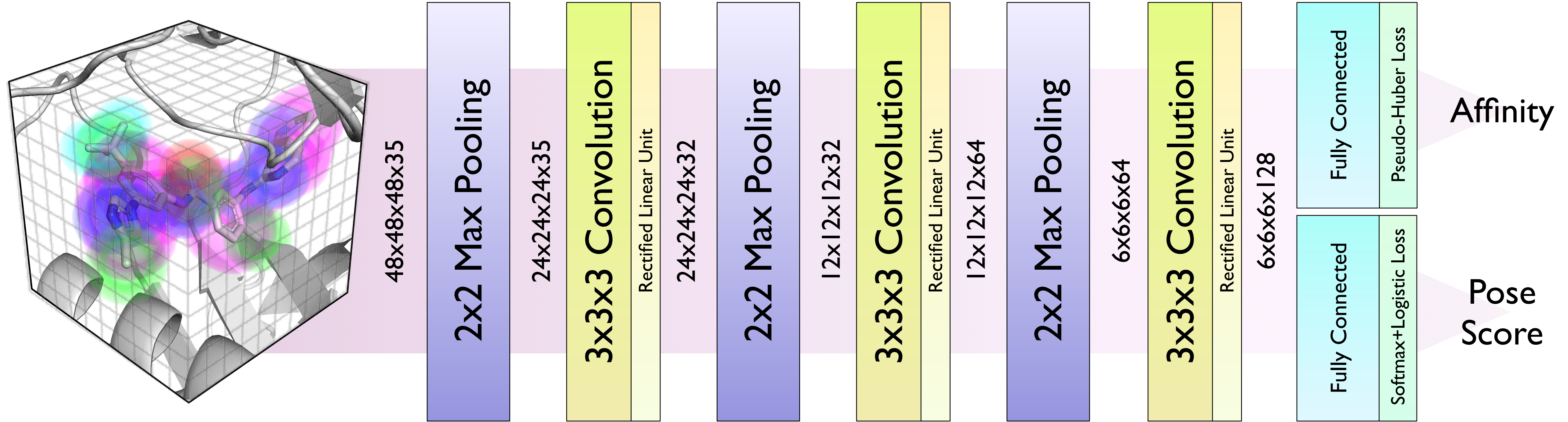# Training Data



## Pose Prediction

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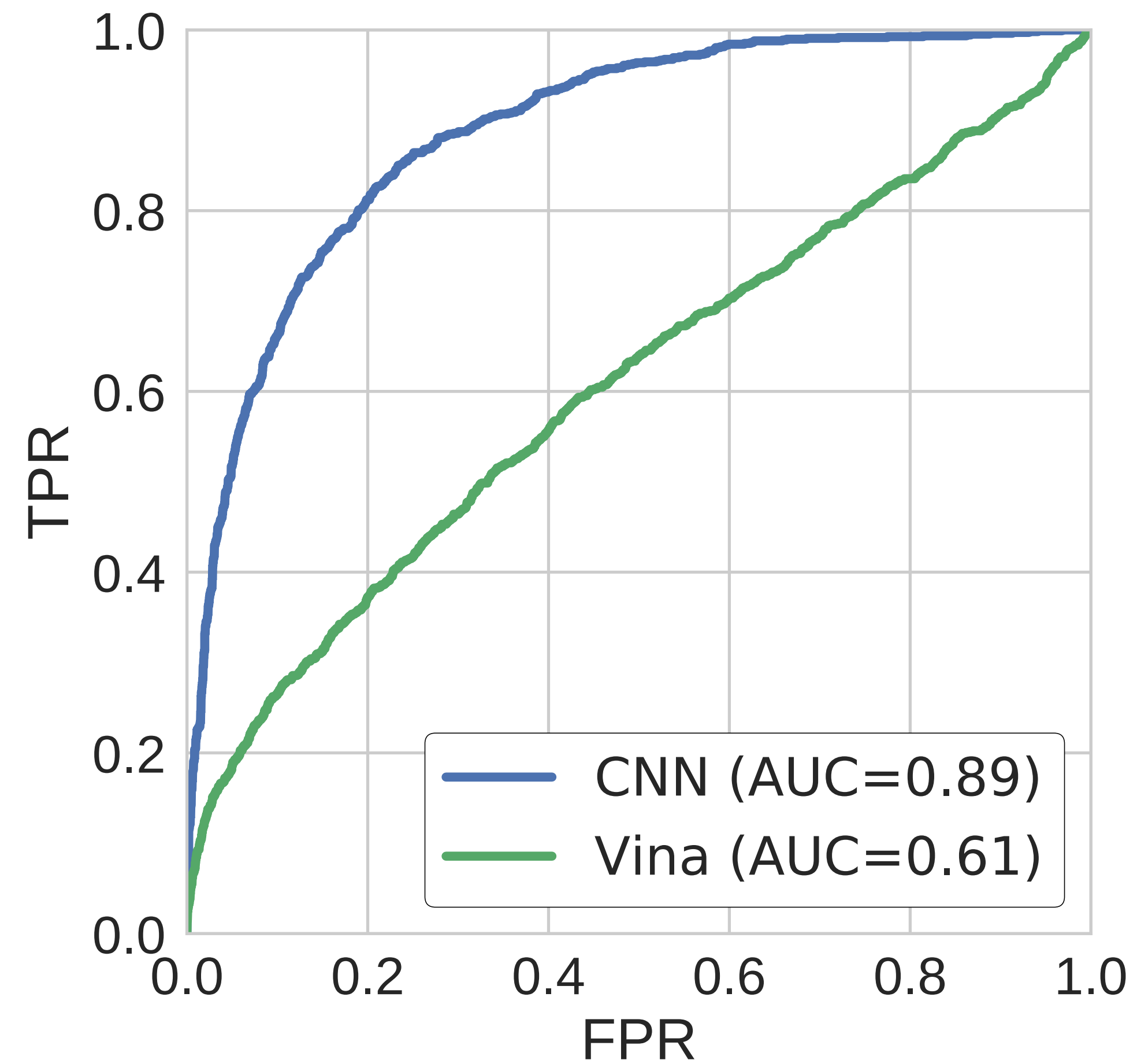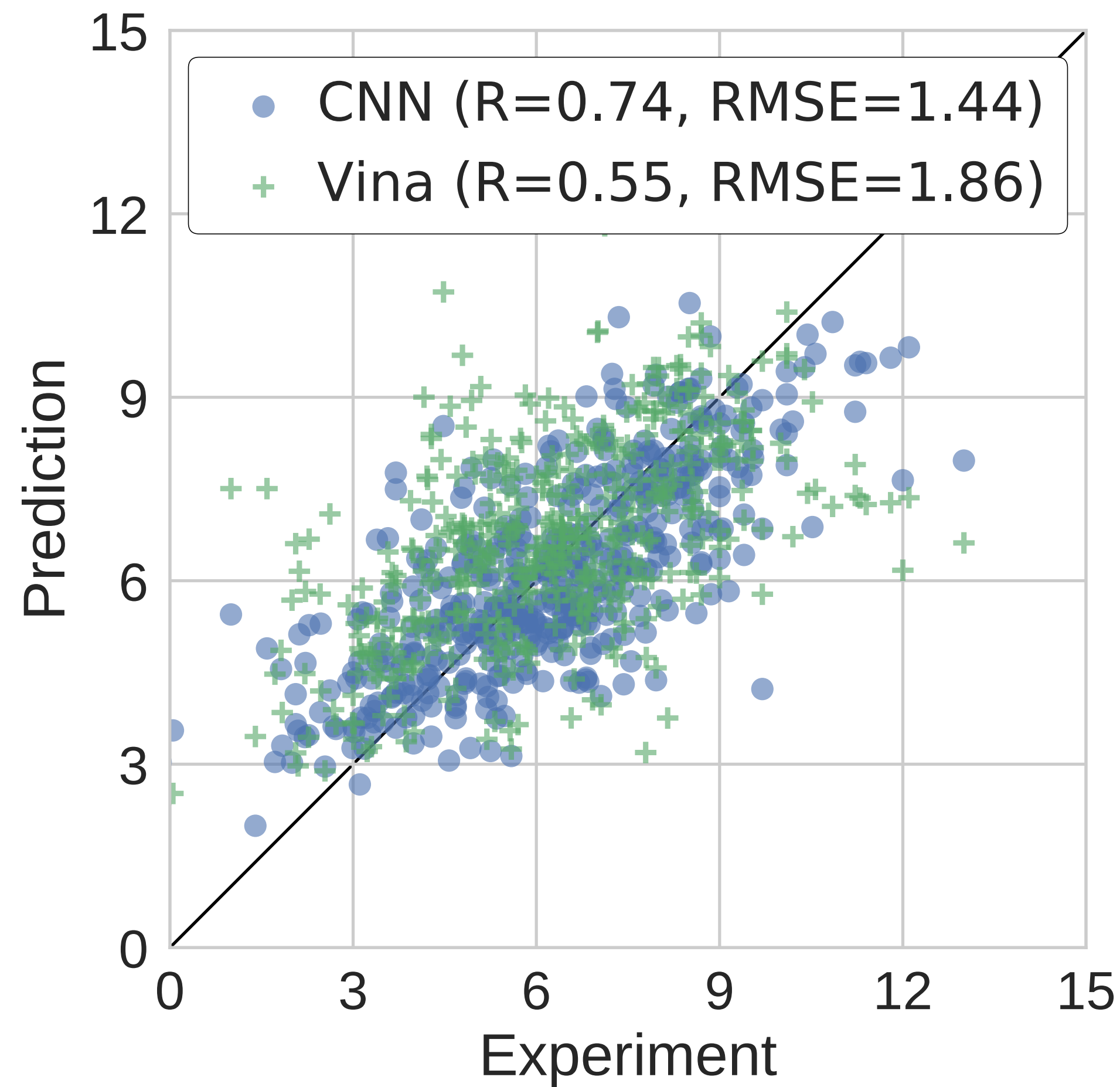  - 8,688  <2Å RMSD (actives)
  - 76,743 >4Å RMSD (decoys)

## Affinity Prediction

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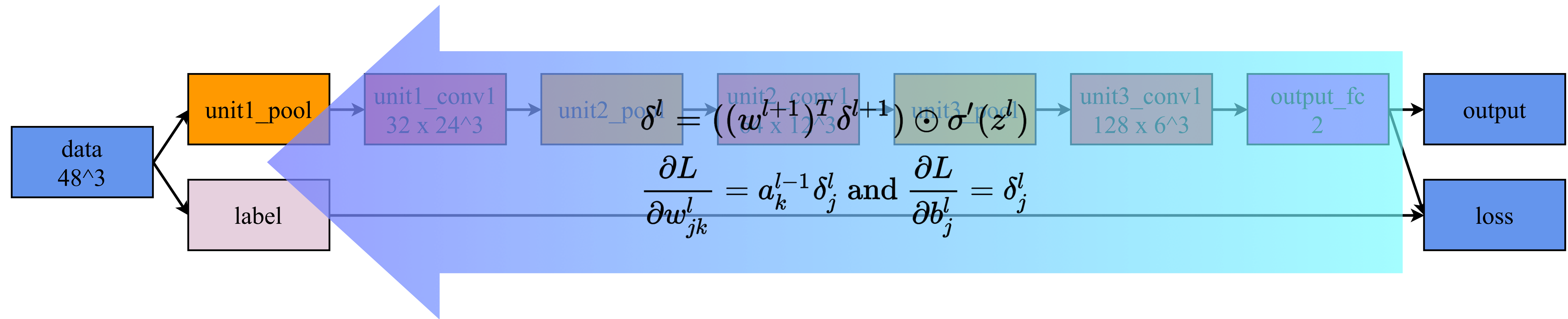- assign known affinity
- **regression problem**

13

# Model

# Results



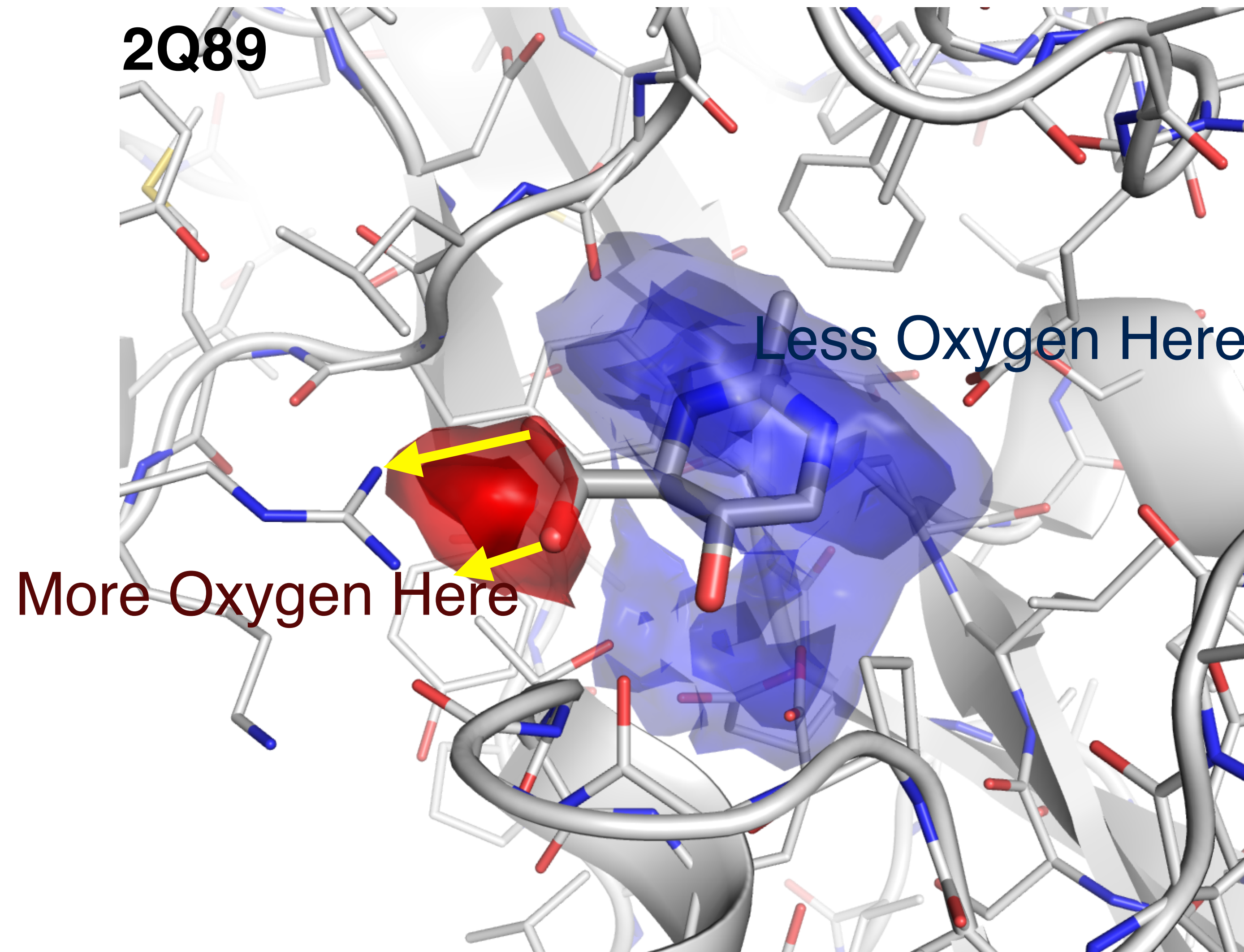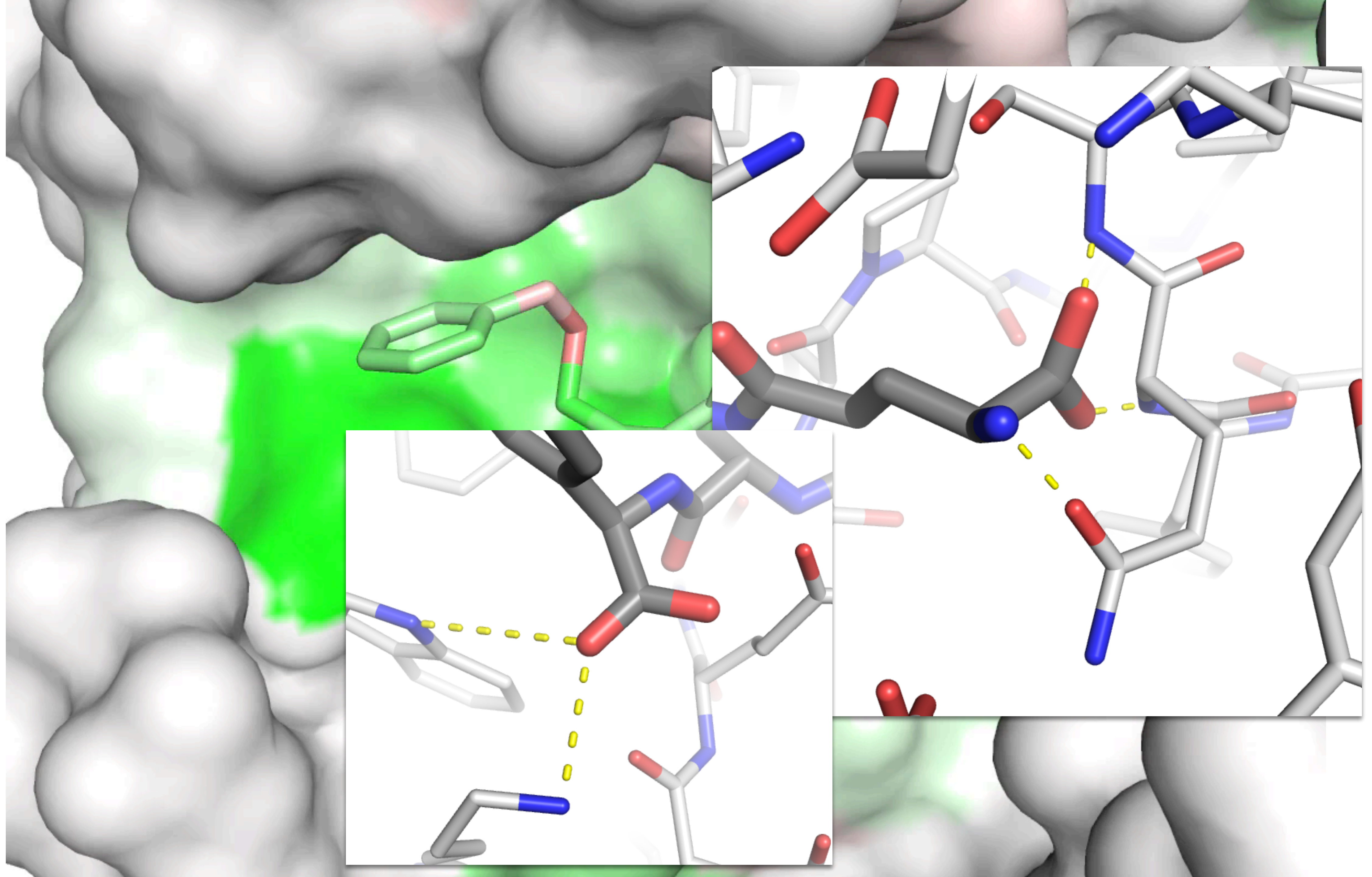Trained on PDBbind refined; tested on CSAR

# Beyond Scoring



$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

$$\frac{\partial L}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial L}{\partial b_j^l} = \delta_j^l$$

## Deep Dreams



optimize
with prior

https://research.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html
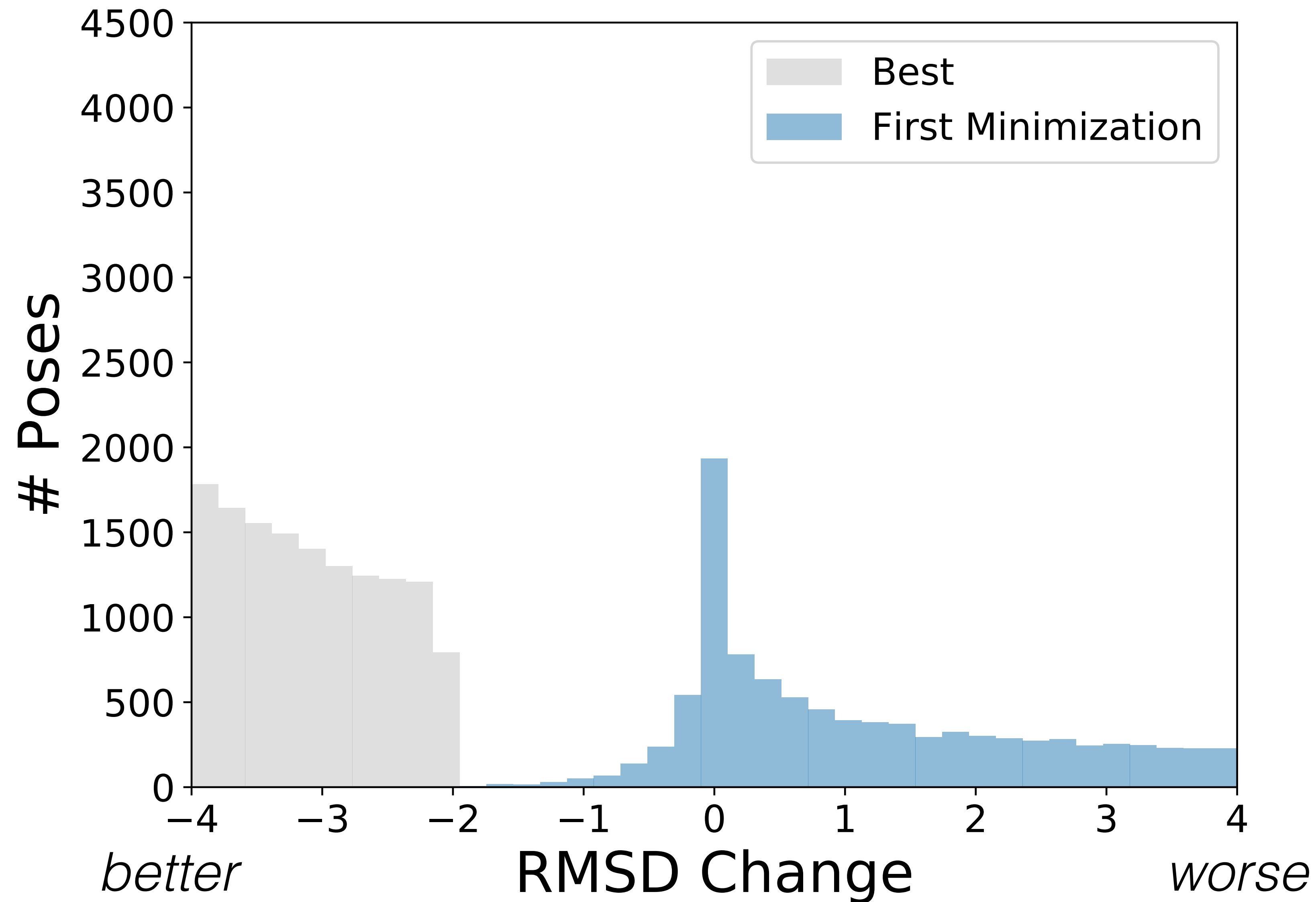
# Beyond Scoring

**2Q89**

Less Oxygen Here

More Oxygen Here

$$\frac{\partial L}{\partial A} = \sum_{i \in G_A} \frac{\partial L}{\partial G_i} \frac{\partial G_i}{\partial D} \frac{\partial D}{\partial A}$$
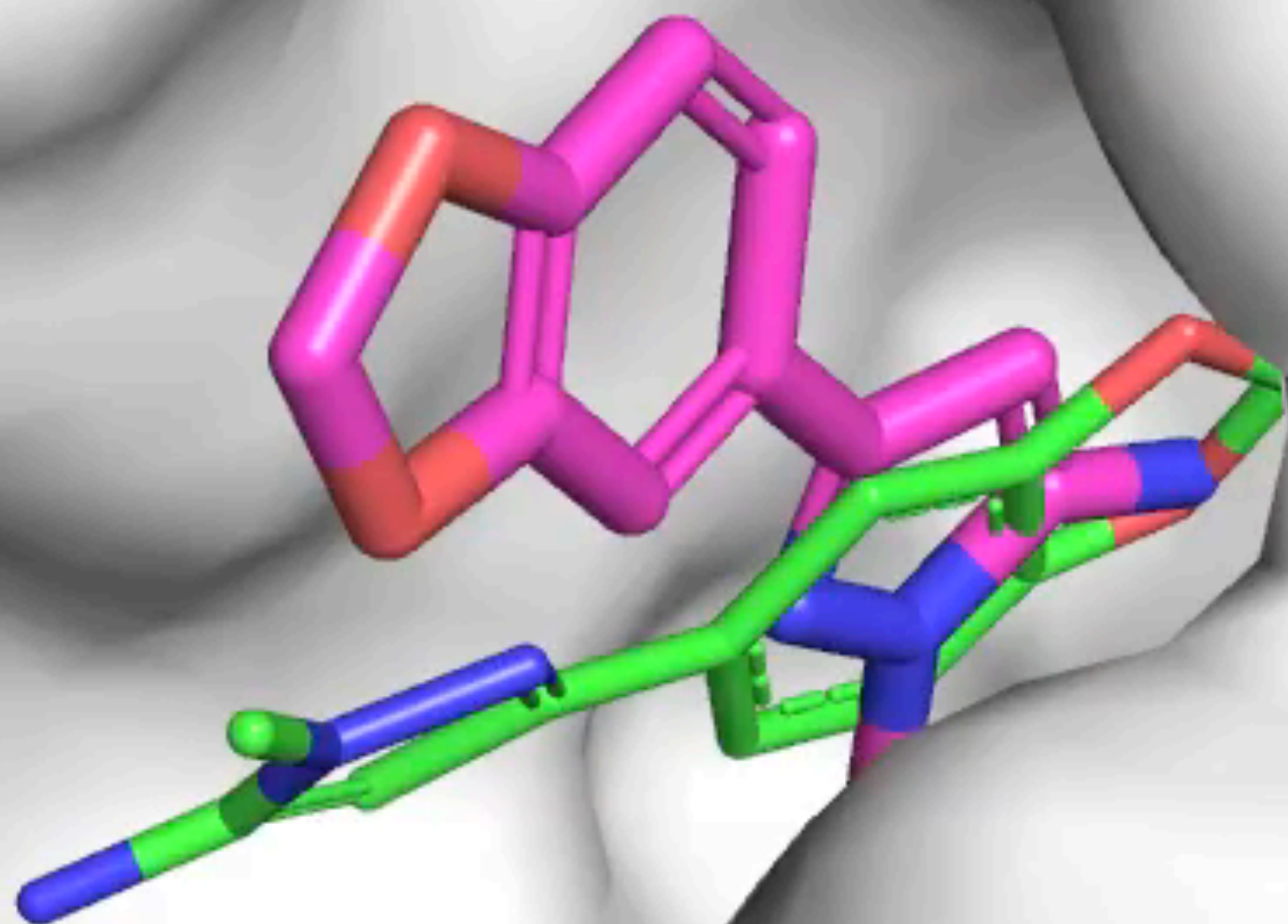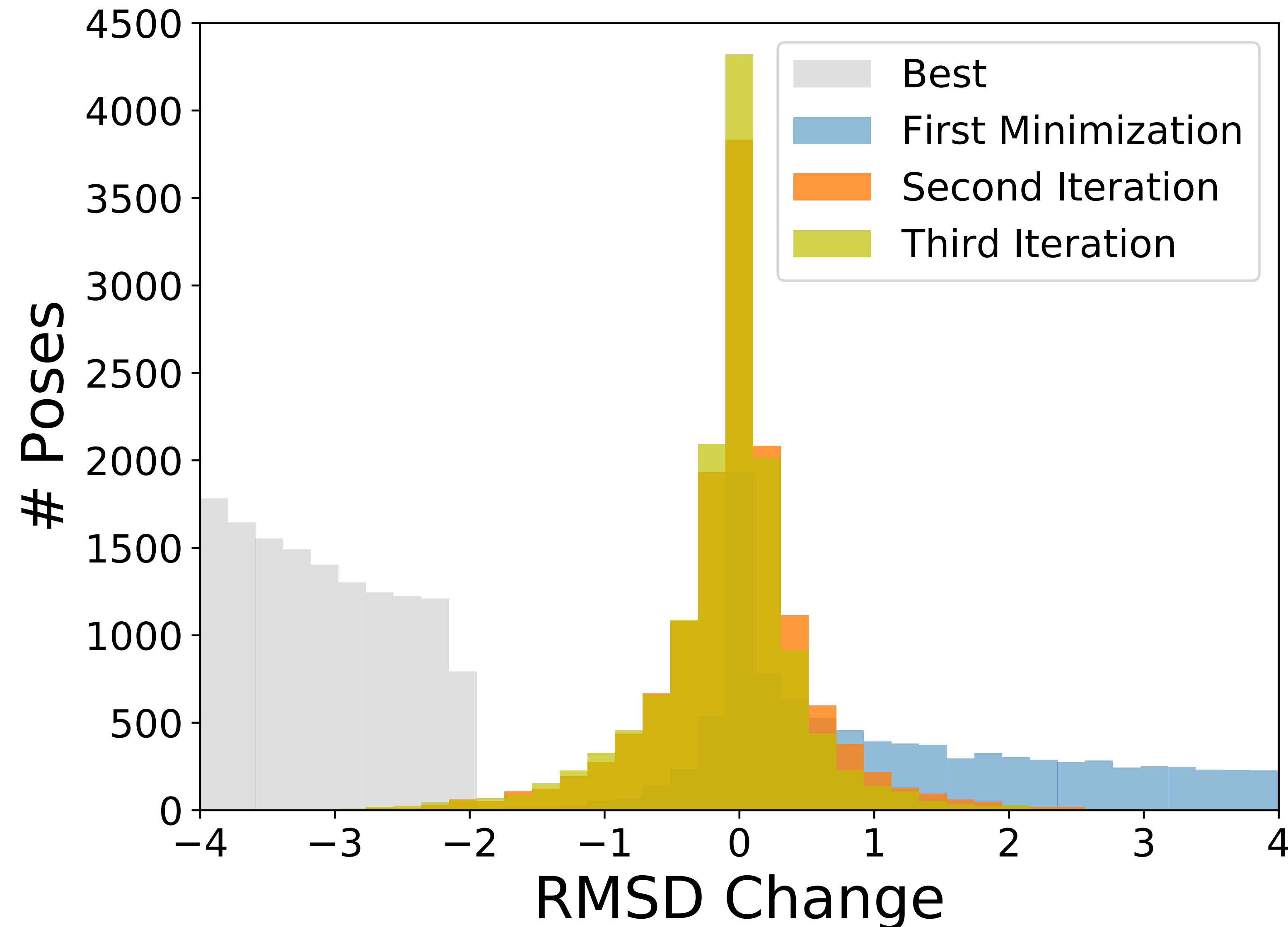
# Minimizing Low RMSD Poses

# Iterative Refinement
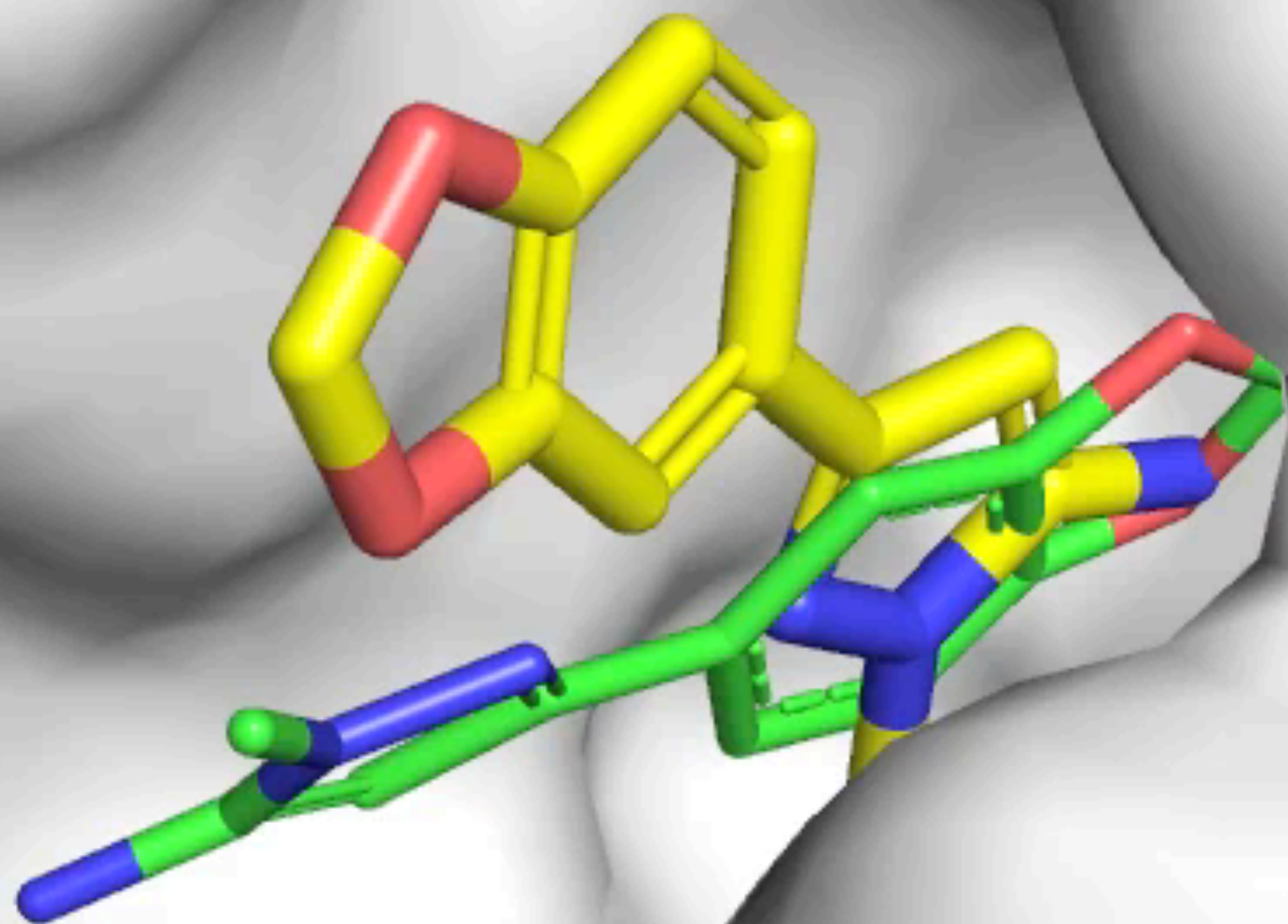
# Docking

## vina/smina/gnina

### Sampling

### Refinement



MCMC

MCMC

MCMC

MCMC
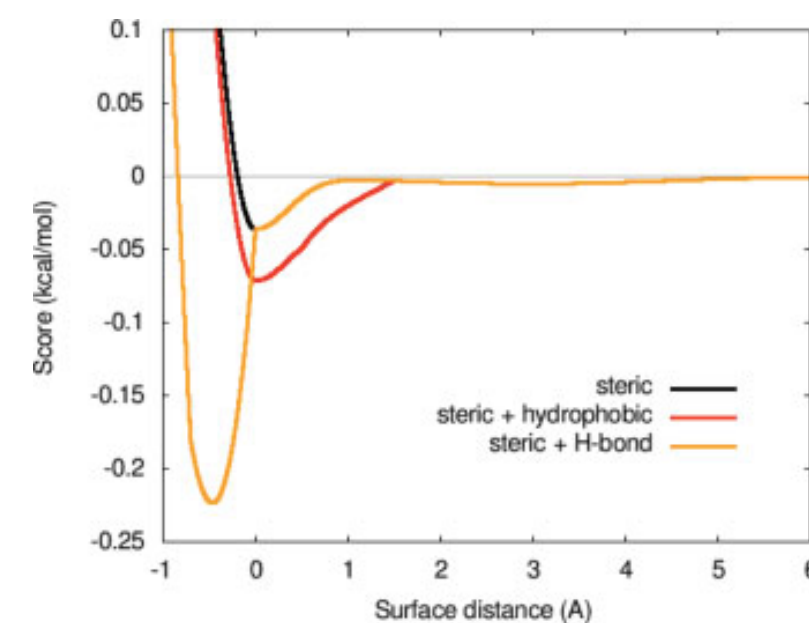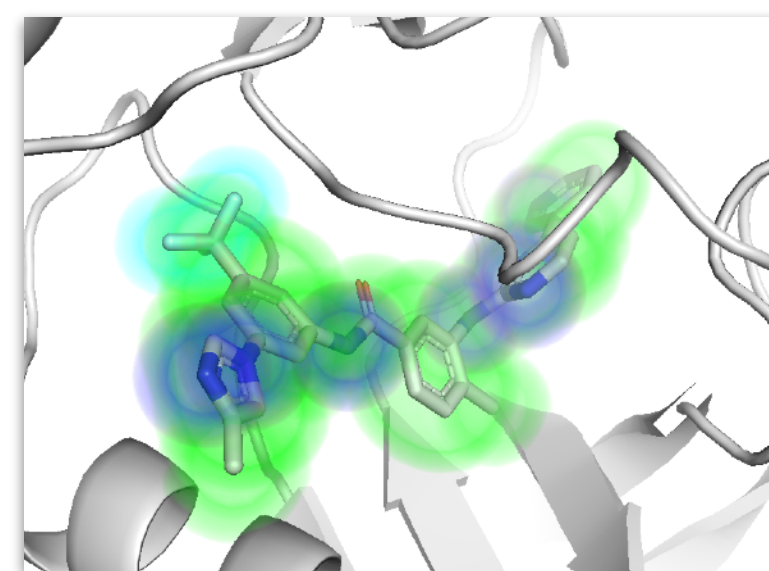
MCMC

⋮

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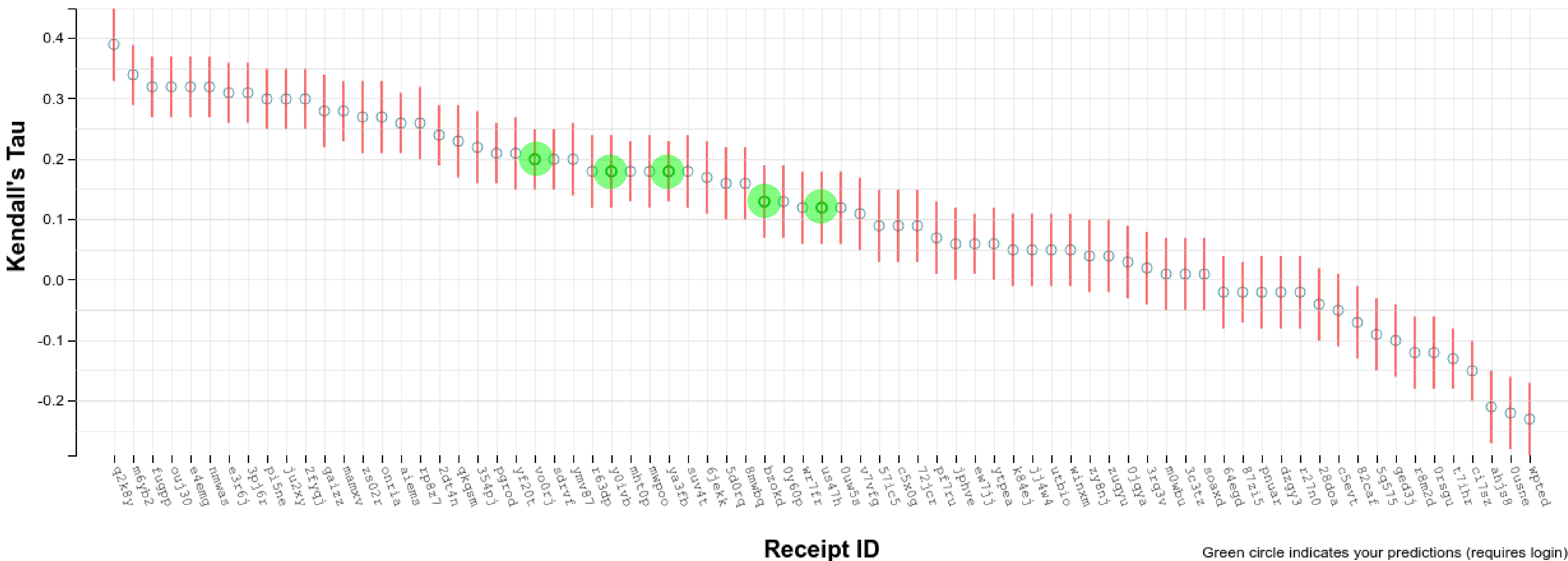
best
poses

**Vina**

**CNN**

Rescoring

**CNN**
pose
affinity

# D3R Results

# Grand Challenge 3



**Grand Challenge 3 - CatS_stage2**

**Affinity Ranking - Kendall's Tau**

Green circle indicates your predictions (requires login)

25

# Grand Challenge 3



**Grand Challenge 3 - JAK2_SC2**

**Affinity Ranking - Kendall's Tau**

Green circle indicates your predictions (requires login)

# Grand Challenge 3



**Grand Challenge 3 - p38a**

**Affinity Ranking - Kendall's Tau**

Green circle indicates your predictions (requires login)

# Grand Challenge 3



**Grand Challenge 3 - TIE2**
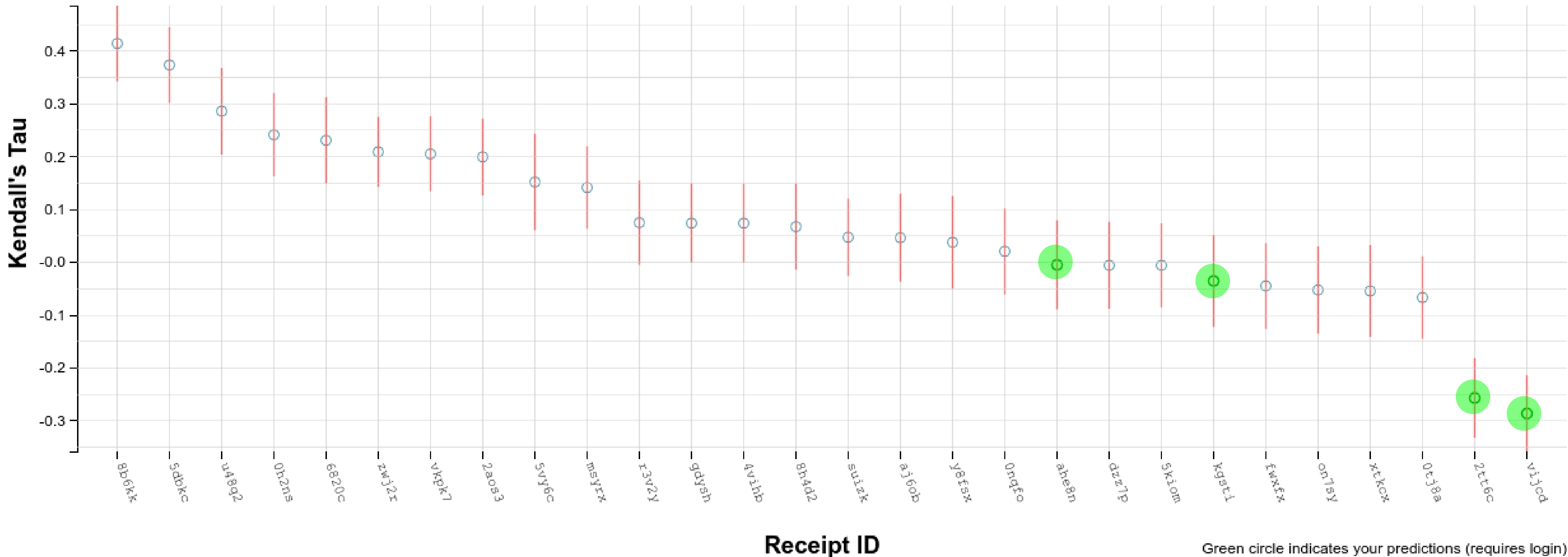
**Affinity Ranking - Kendall's Tau**

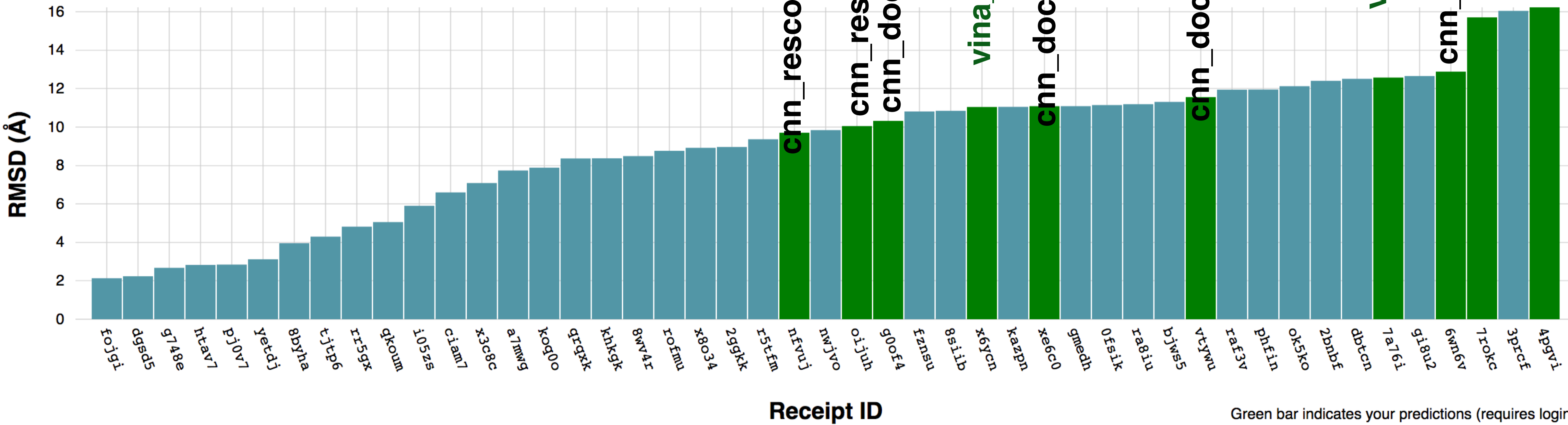Green circle indicates your predictions (requires login)

# Grand Challenge 3



**Grand Challenge 3 - VEGFR2**

**Affinity Ranking - Kendall's Tau**

# Grand Challenge 3

## Spearman Correlation

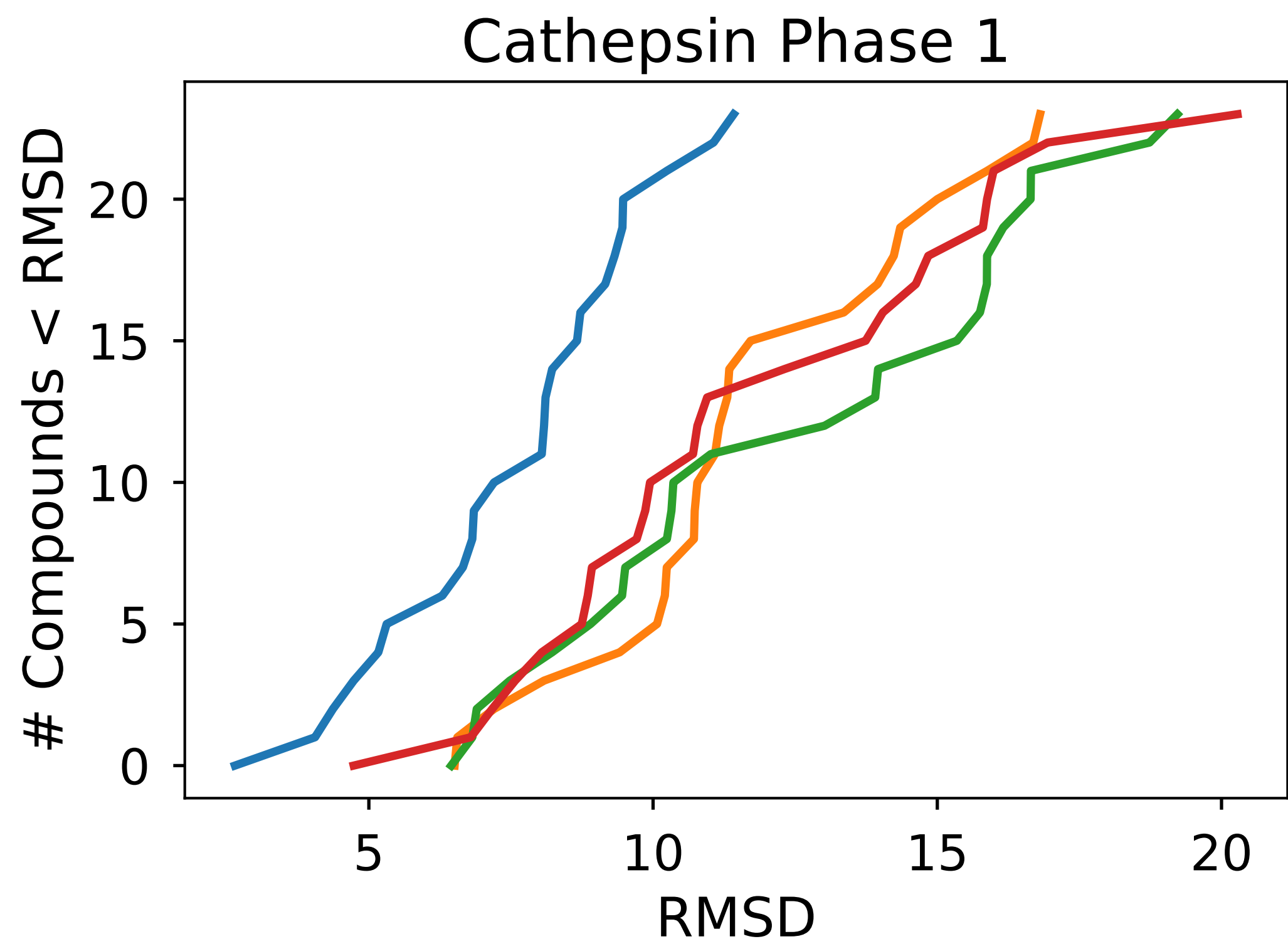|          | cnn_docked_affinity | cnn_rescore_affinity | cnn_docked_scoring | cnn_rescore_scoring | vina    |
|----------|---------------------|----------------------|--------------------|---------------------|---------|
| cat      | 0.0701              | 0.154                | -0.0351            | 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   |

# GC3: Pose Prediction

# GC3: Pose Prediction



Cathepsin Phase 1

cross-docking

Cathepsin Phase 1b

redocking

# Grand Challenge 1

# Grand Challenge 2



**Affinity Ranking (Stage 2) - Kendall's Tau**

*affinity rescore*

Kendall's Tau

Receipt ID

Green circle indicates your predictions (requires login)
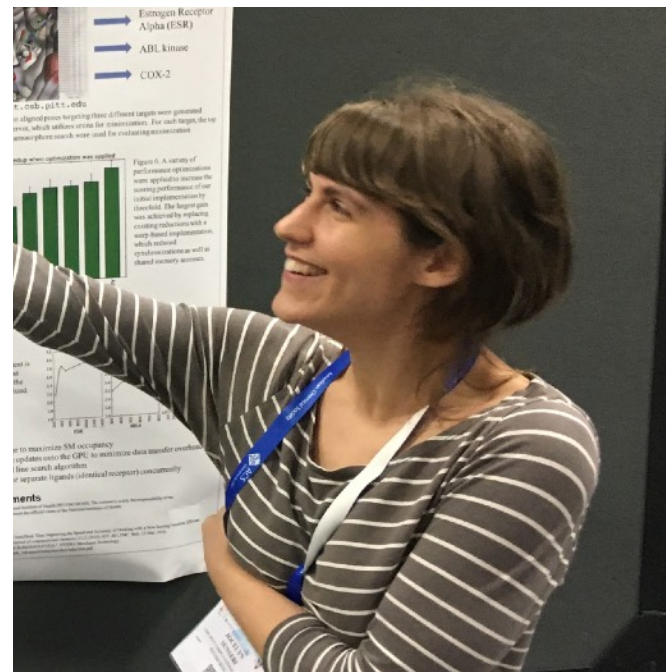
# Future Plans

Train CNN for docking

- iteratively train on docked poses

- train on cross-docked poses

- fully integrate CNN scoring into search
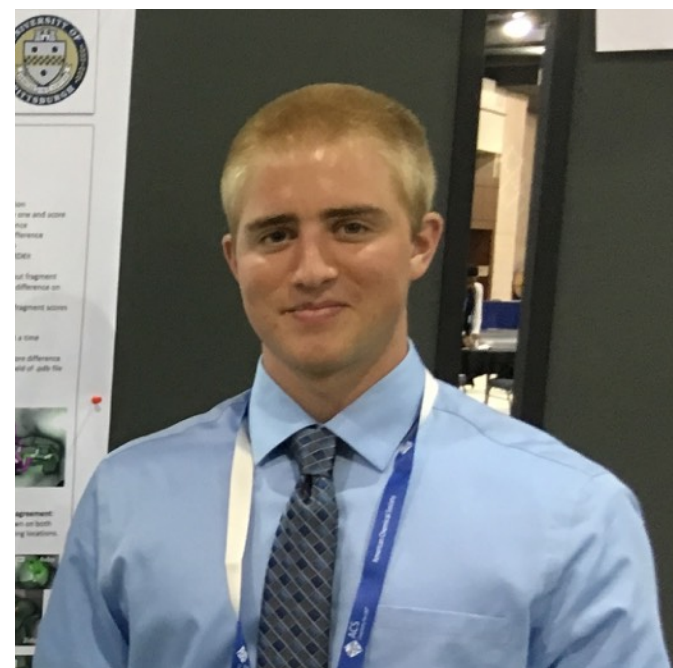
Continue to improve model/training parameters

Next Grand Challenge

- Finish fully automated predictions early
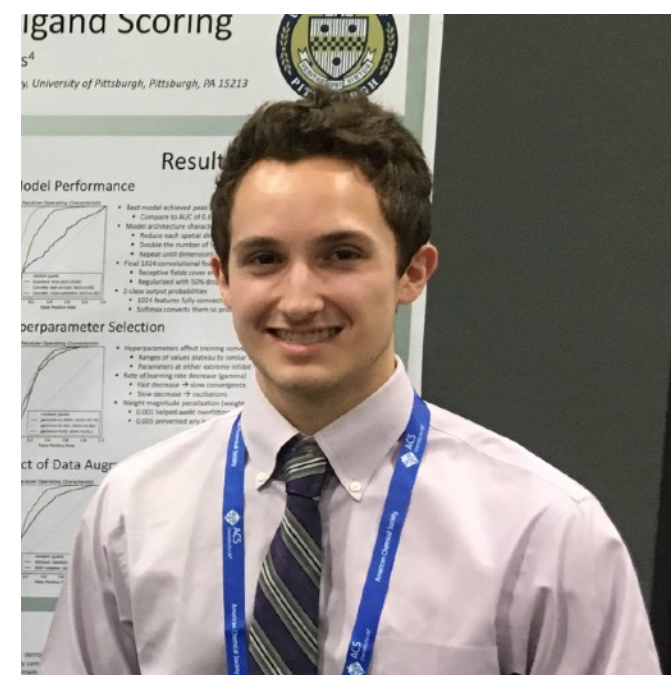
- Make automated+human insight submission

# Acknowledgements



Jocelyn Sunseri



Josh Hochuli



Matt Ragoza

**Group Members**
Jocelyn Sunseri
Jonathan King
Paul Francoeur
Matt Ragoza
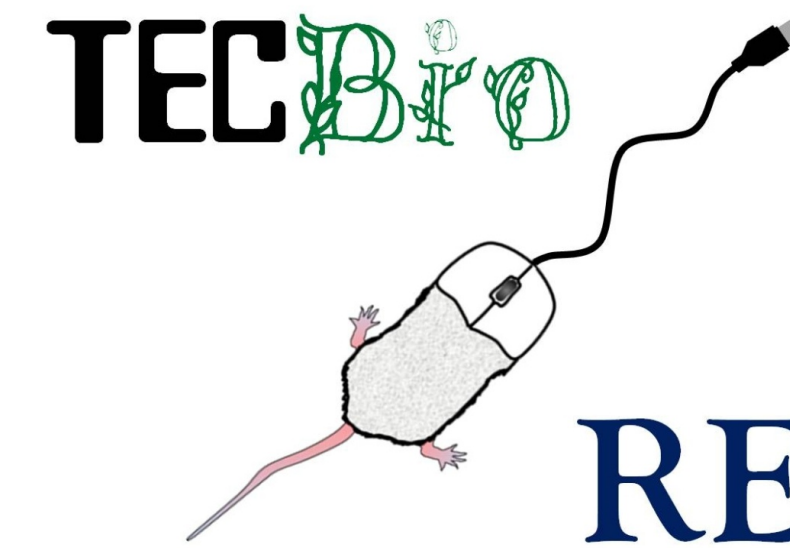Josh Hochuli
Pulkit Mittal
Alec Helbling
Gibran Biswas
Sharanya Bandla
Faiha Khan
Lily Turner
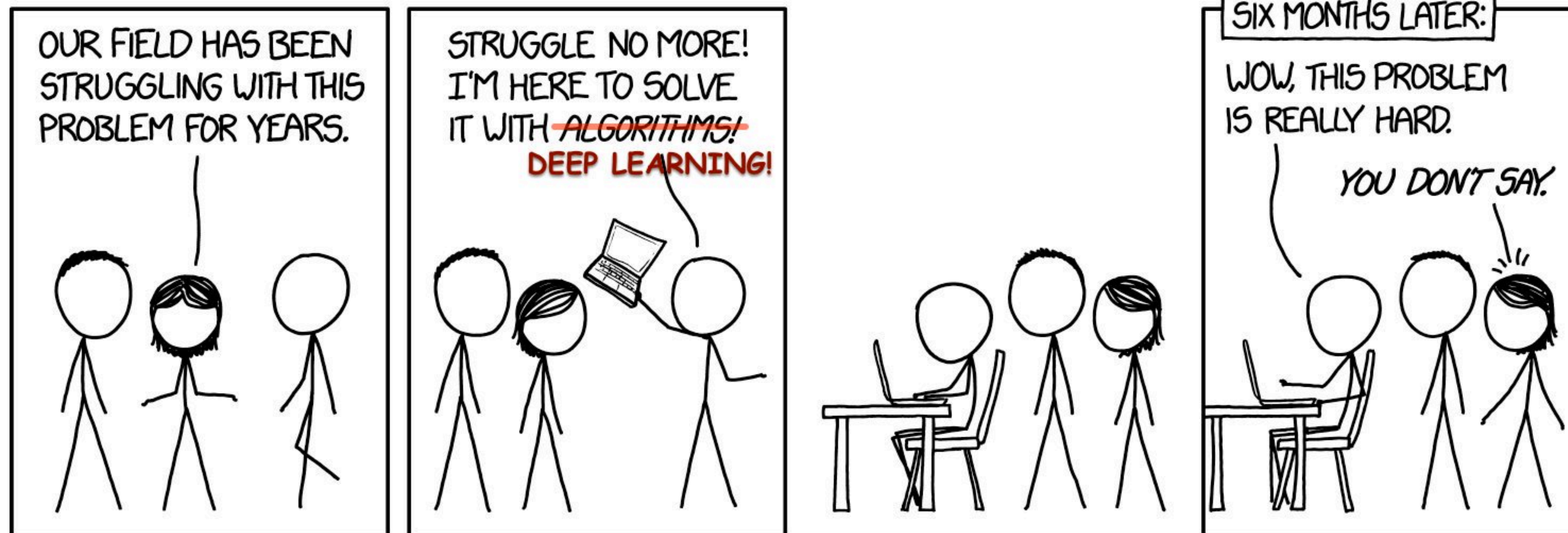


Department of
Computational and
Systems Biology

github.com/gnina

http://bits.csb.pitt.edu
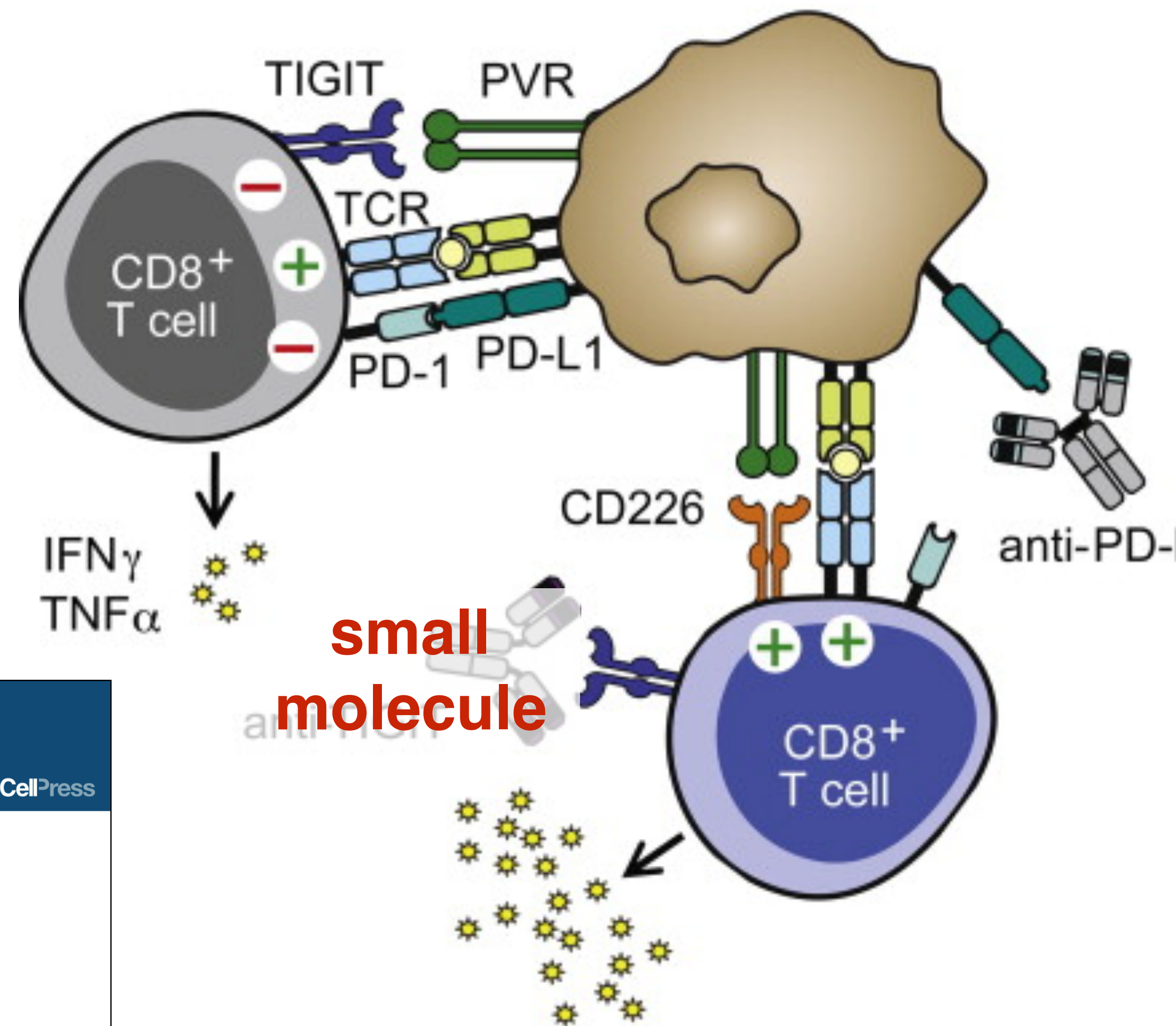
@david_koes

# Prospective Case Study: TIGIT

Can we block TIGIT/ PVR interaction with a small molecule?



**small molecule**

**The Immunoreceptor TIGIT Regulates Antitumor and Antiviral CD8+ T Cell Effector Function**

Robert J. Johnston,[1] Laetitia Comps-Agrar,[2] Jason Hackney,[3] Xin Yu,[1] Mahrukh Huseni,[4] Yagai Yang,[5] Summer Park,[6] Vincent Javinal,[5] Henry Chiu,[7] Bryan Irving,[1] Dan L. Eaton,[2] and Jane L. Grogan[1],*
[1]Department of Cancer Immunology
[2]Department of Protein Chemistry
[3]Department of Bioinformatics and Computational Biology
[4]Department of Oncology Biomarker Development
[5]Department of Translational Oncology
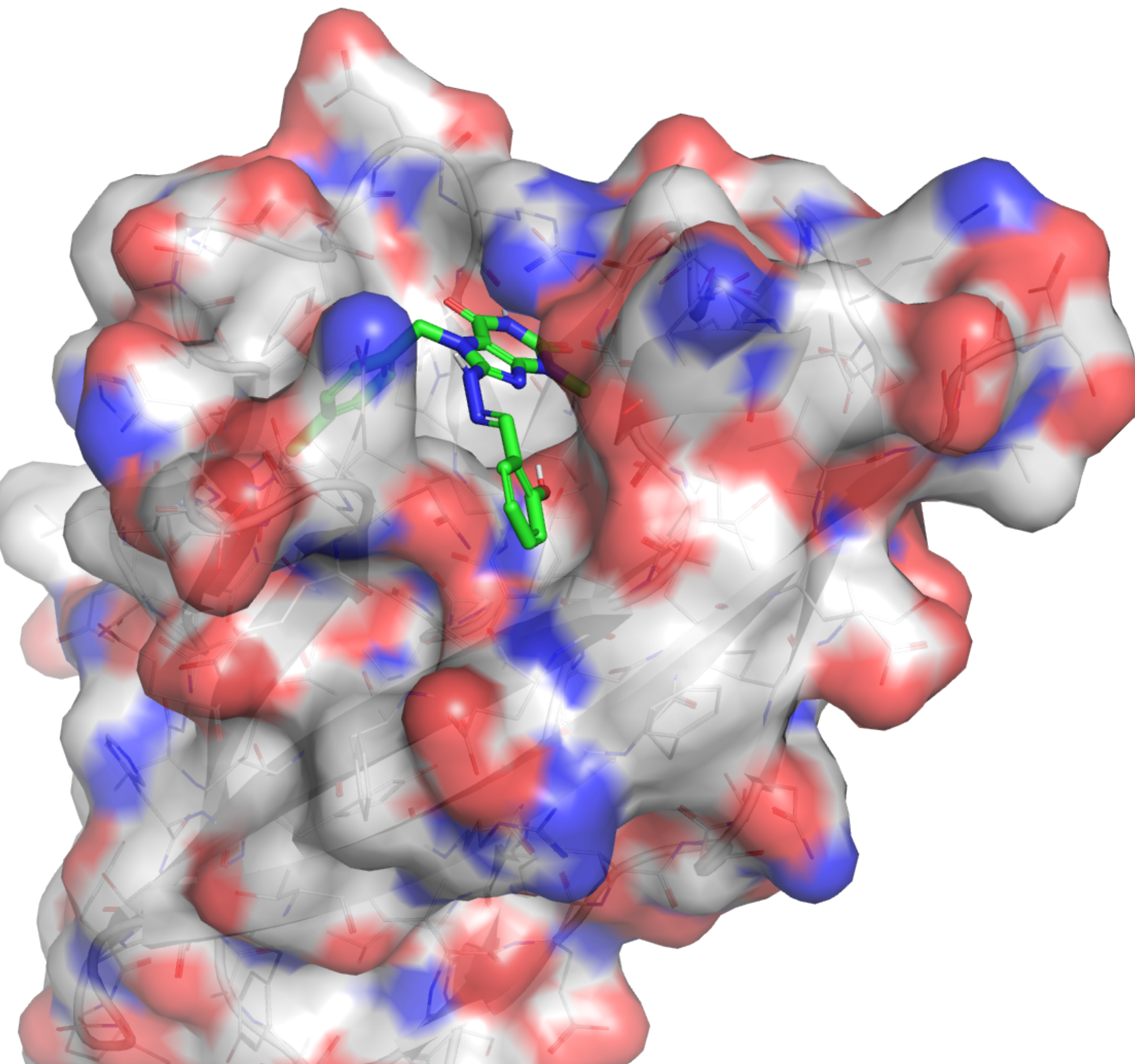[6]Department of Translational Immunology
[7]Department of Biochemical and Cellular Pharmacology
Genentech, 1 DNA Way, South San Francisco, CA 94080, USA
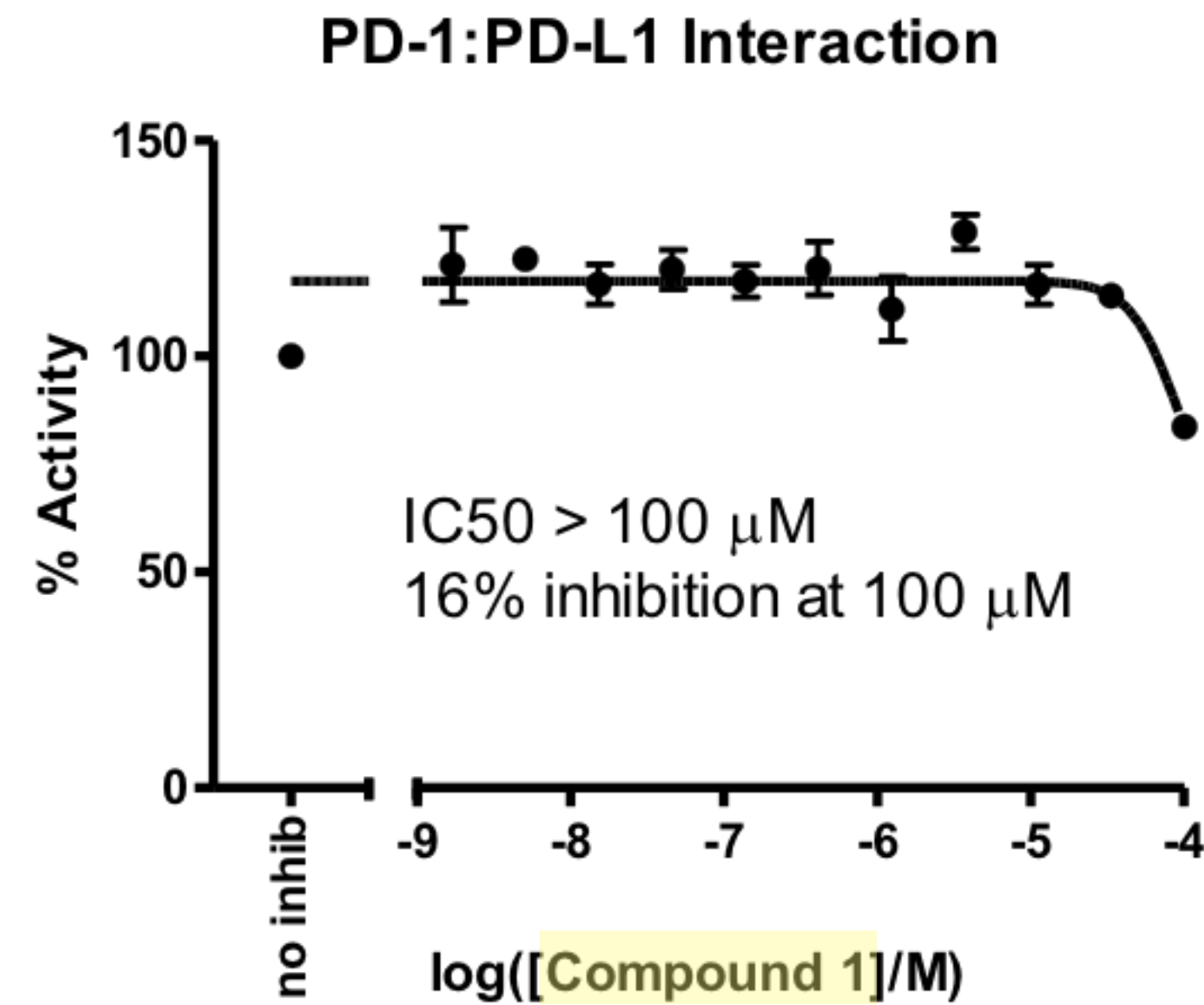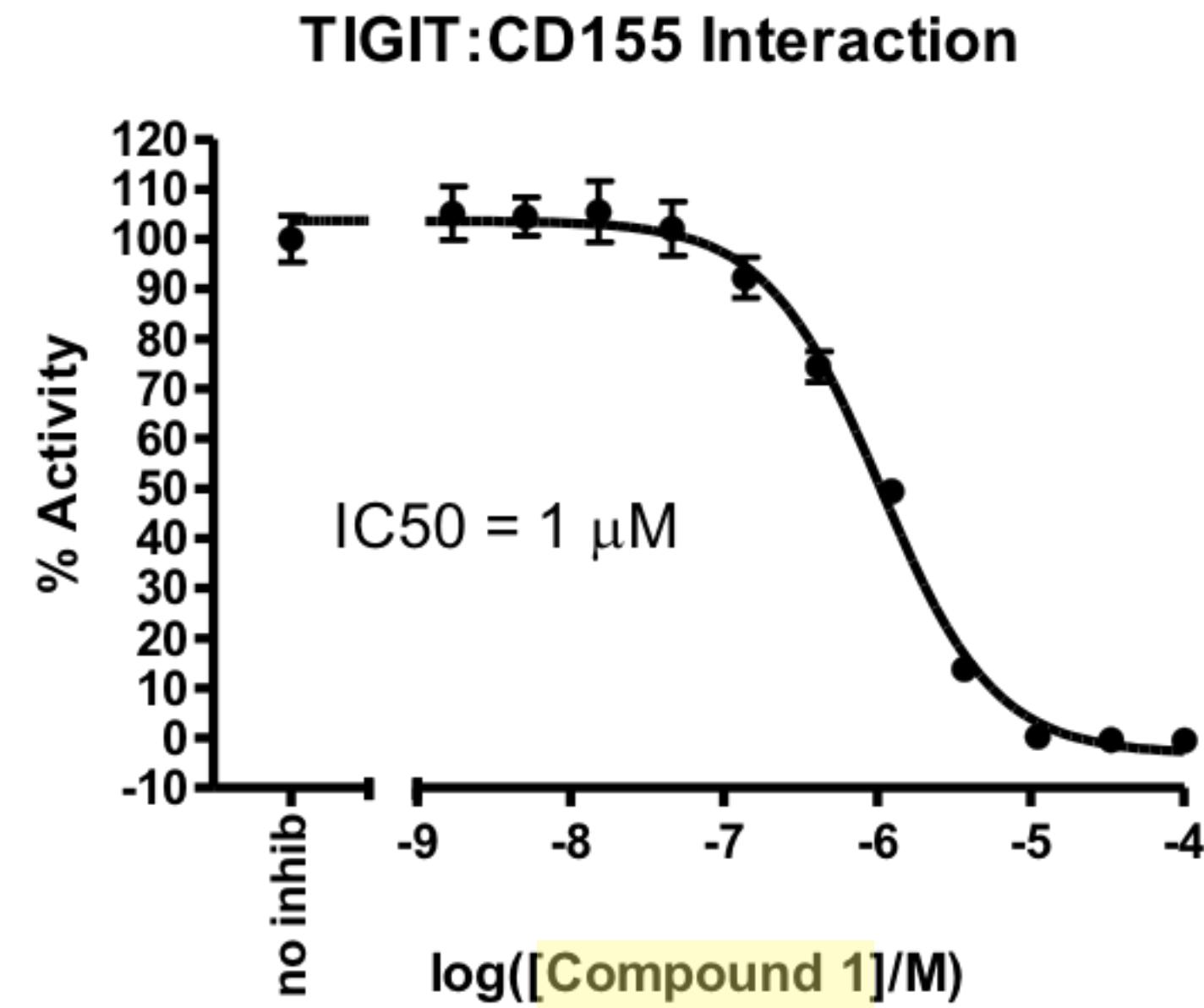*Correspondence: grogan.jane@gene.com
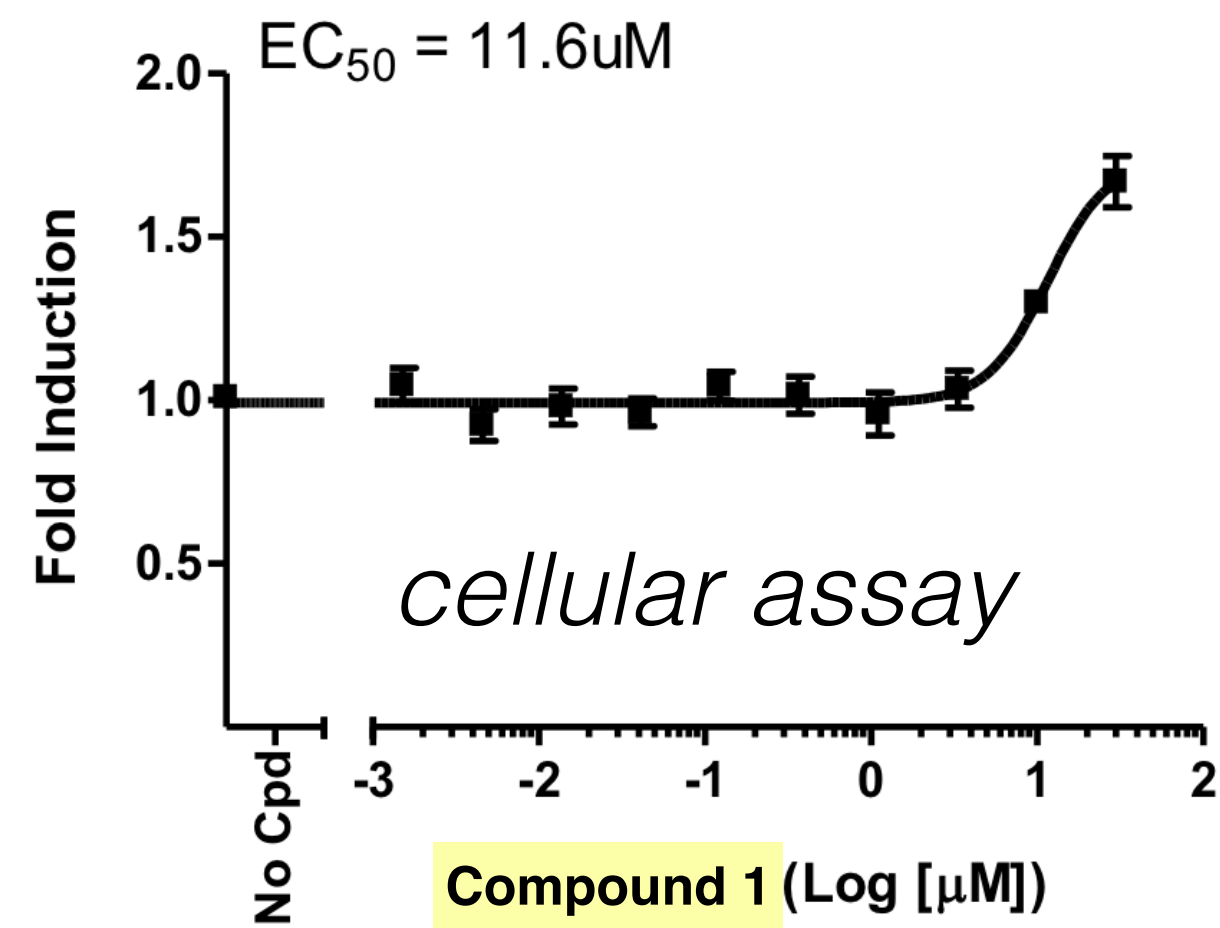http://dx.doi.org/10.1016/j.ccell.2014.10.018

# Screening



**10 diverse compounds selected for screening**
- **top ranked by Vina**
- **top ranked by CNN**

| Name | CNN Affinity | CNN Score | Vina |
|---|---|---|---|
| Compound 1 | 7.69807 | 0.994763 | 85.95 |
| Compound 2 | 5.57909 | 0.0180277 | -8.12632 |
| Compound 3 | 6.73692 | 0.0624742 | -9.81935 |
| Compound 4 | 6.87897 | 0.953488 | -3.81378 |
| Compound 5 | 6.32813 | 0.209807 | -8.60293 |
| Compound 6 | 5.689 | 0.0437 | -8.991 |
| Compound 7 | 4.368 | 0.022 | -9.34722 |
| Compound 8 | 4.81 | 0.072 | -6.81787 |
| Compound 9 | 5.22 | 0.032 | -6.264 |
| Compound 10 | 6.67 | 0.361 | 6.1053 |

# Results



TIGIT:CD155 Interaction

% Activity — Test Compound # (all at 100 uM except #2 and #3 at 50 uM)

Negative Control, Positive Control, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10

TIGIT:CD155 Interaction

IC50 = 1 $\mu$M

log([Compound 1]/M)

TIGIT:CD155 Interaction

IC50 ~ 14 $\mu$M
65% inhibition at 100 $\mu$M

log([Compound 4]/M)

EC$_{50}$ = 11.6uM

cellular assay

No Cpd — Compound 1 (Log [$\mu$M])

PD-1:PD-L1 Interaction

IC50 > 100 $\mu$M
16% inhibition at 100 $\mu$M

log([Compound 1]/M)

PD-1:PD-L1 Interaction

IC50 = 14 $\mu$M

log([Compound 4]/M)

But…

# Filter Visualization