# THE BIOPHARMACEUTICAL RESEARCH AND DEVELOPMENT PROCESS
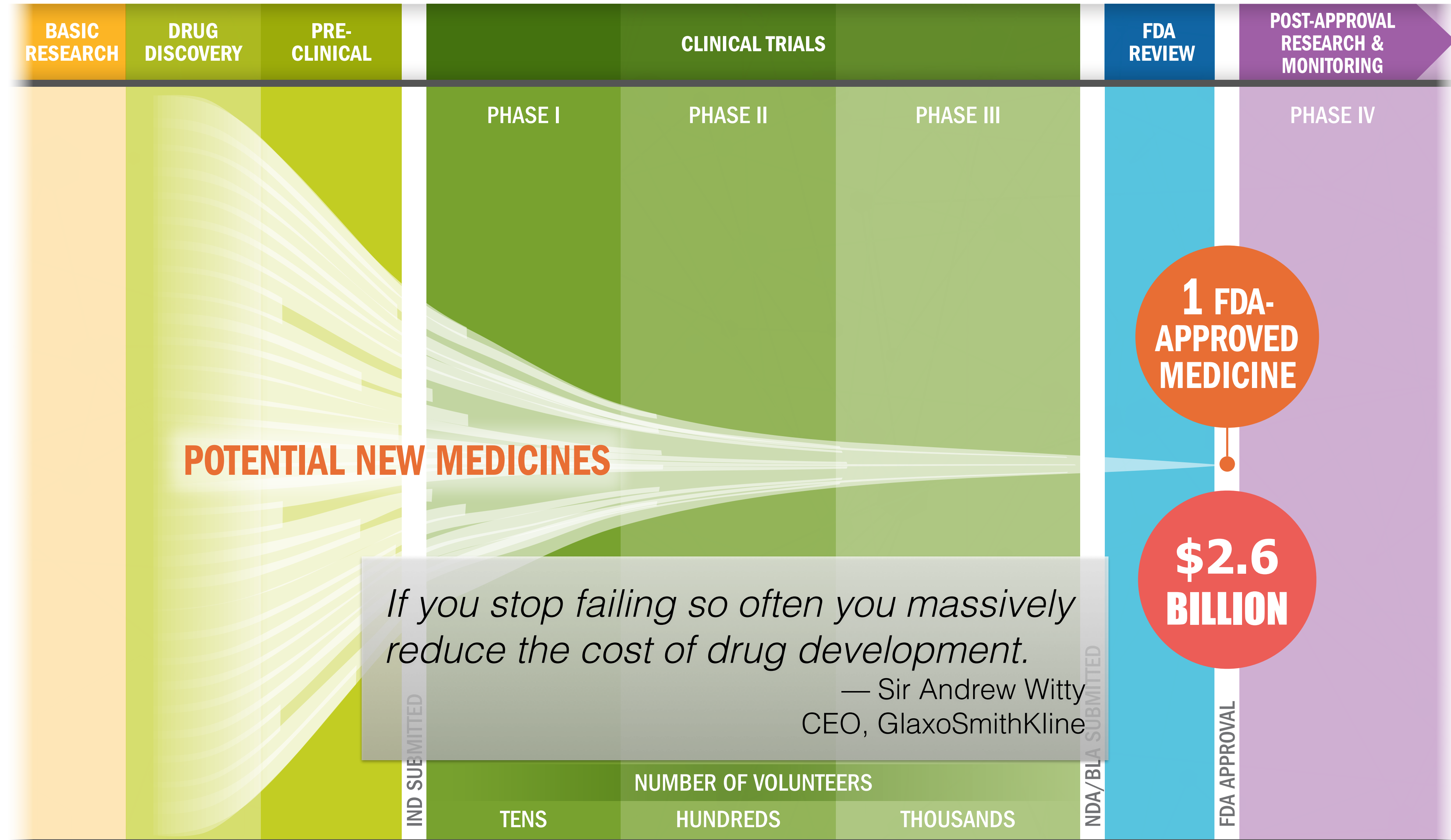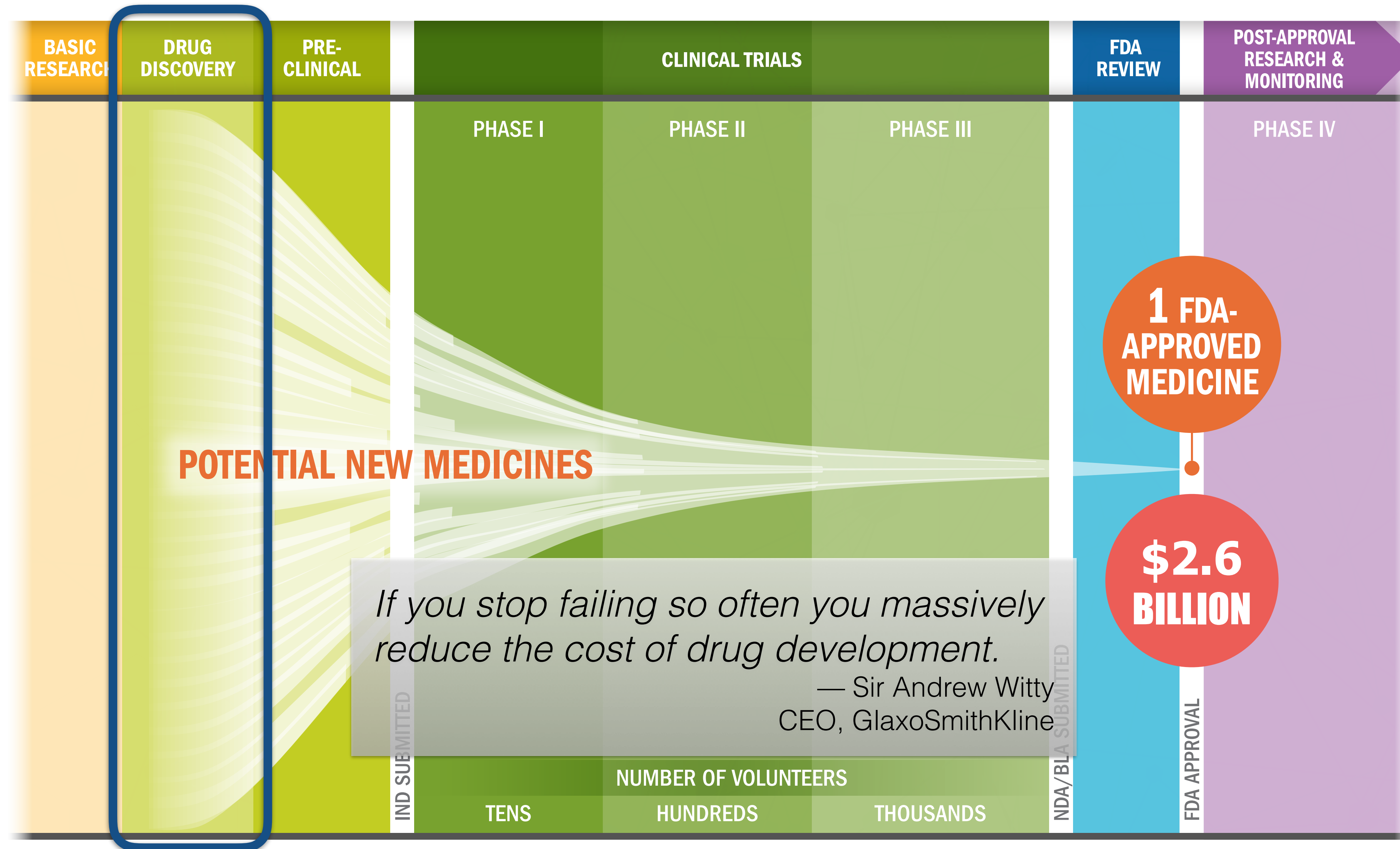


Source: Pharmaceutical Research and Manufacturers of America (http://phrma.org)

2

# THE BIOPHARMACEUTICAL RESEARCH AND DEVELOPMENT PROCESS



Source: Pharmaceutical Research and Manufacturers of America (http://phrma.org)

2

# THE BIOPHARMACEUTICAL RESEARCH AND DEVELOPMENT PROCESS



| BASIC RESEARCH | DRUG DISCOVERY | PRE-CLINICAL | CLINICAL TRIALS | | | FDA REVIEW | POST-APPROVAL RESEARCH & MONITORING |

PHASE I    PHASE II    PHASE III

PHASE IV

**POTENTIAL NEW MEDICINES**

*If you stop failing so often you massively reduce the cost of drug development.*
— Sir Andrew Witty
CEO, GlaxoSmithKline

**1** FDA-APPROVED MEDICINE

**$2.6 BILLION**

IND SUBMITTED

NDA/BLA SUBMITTED

FDA APPROVAL

NUMBER OF VOLUNTEERS
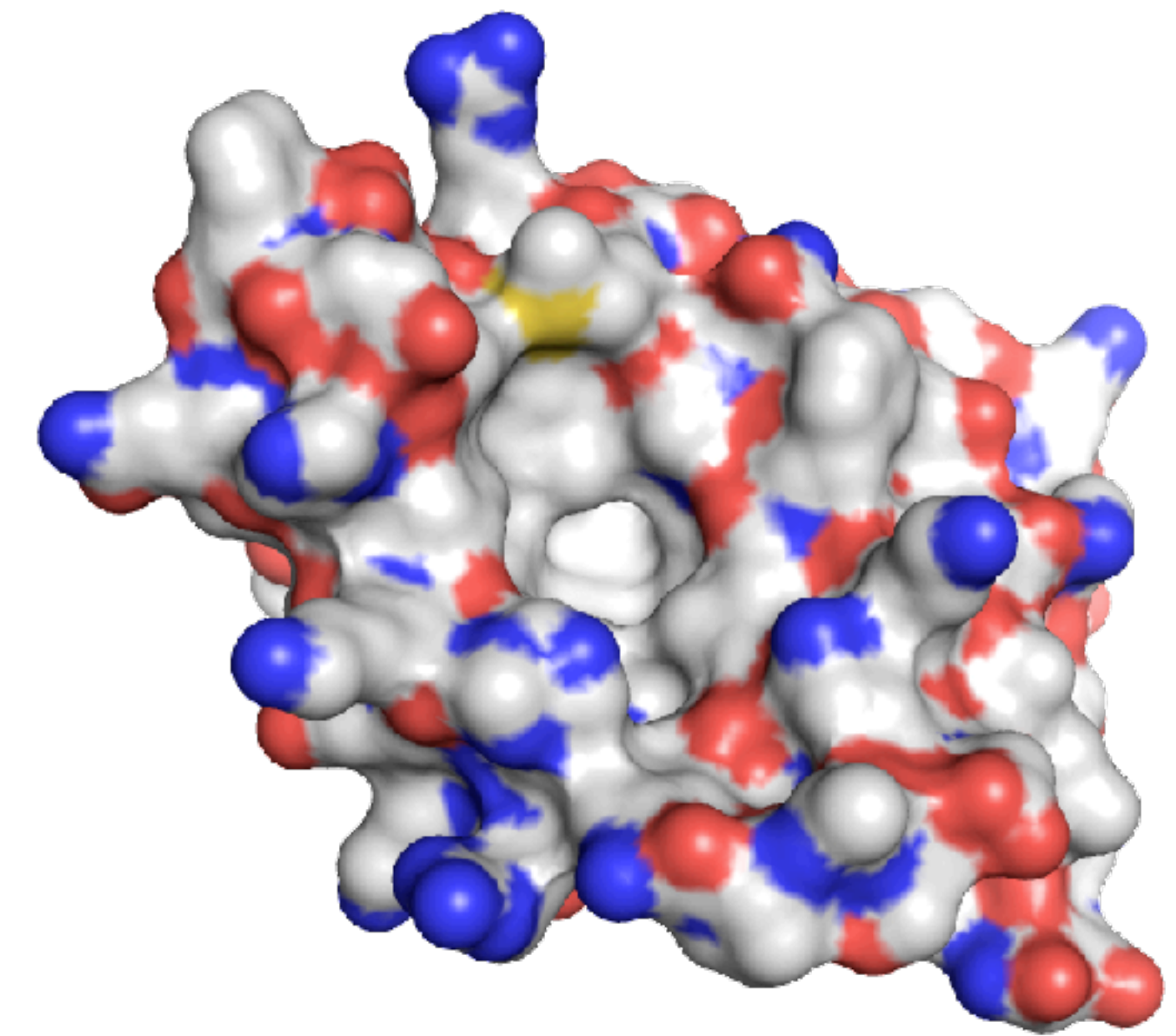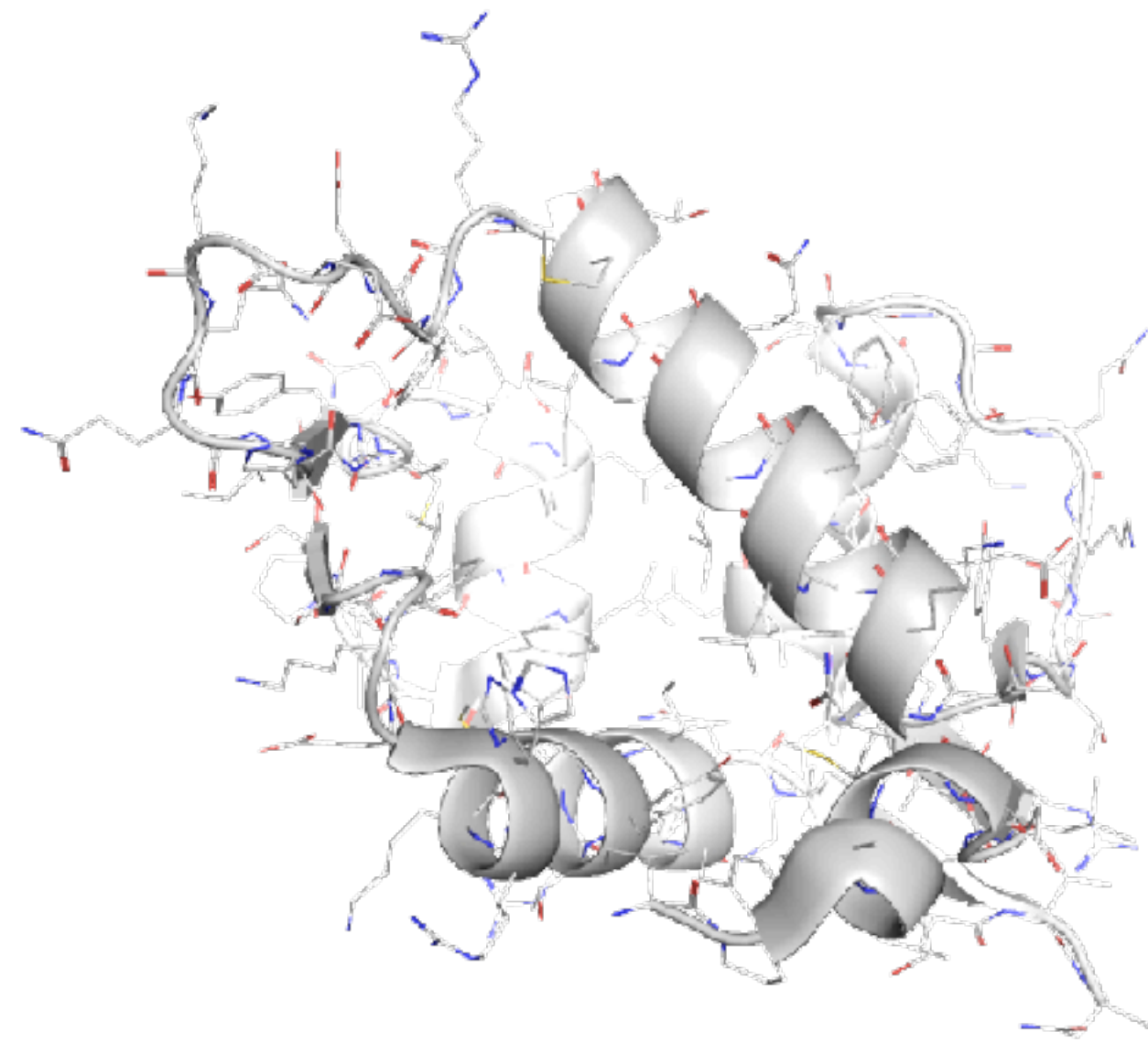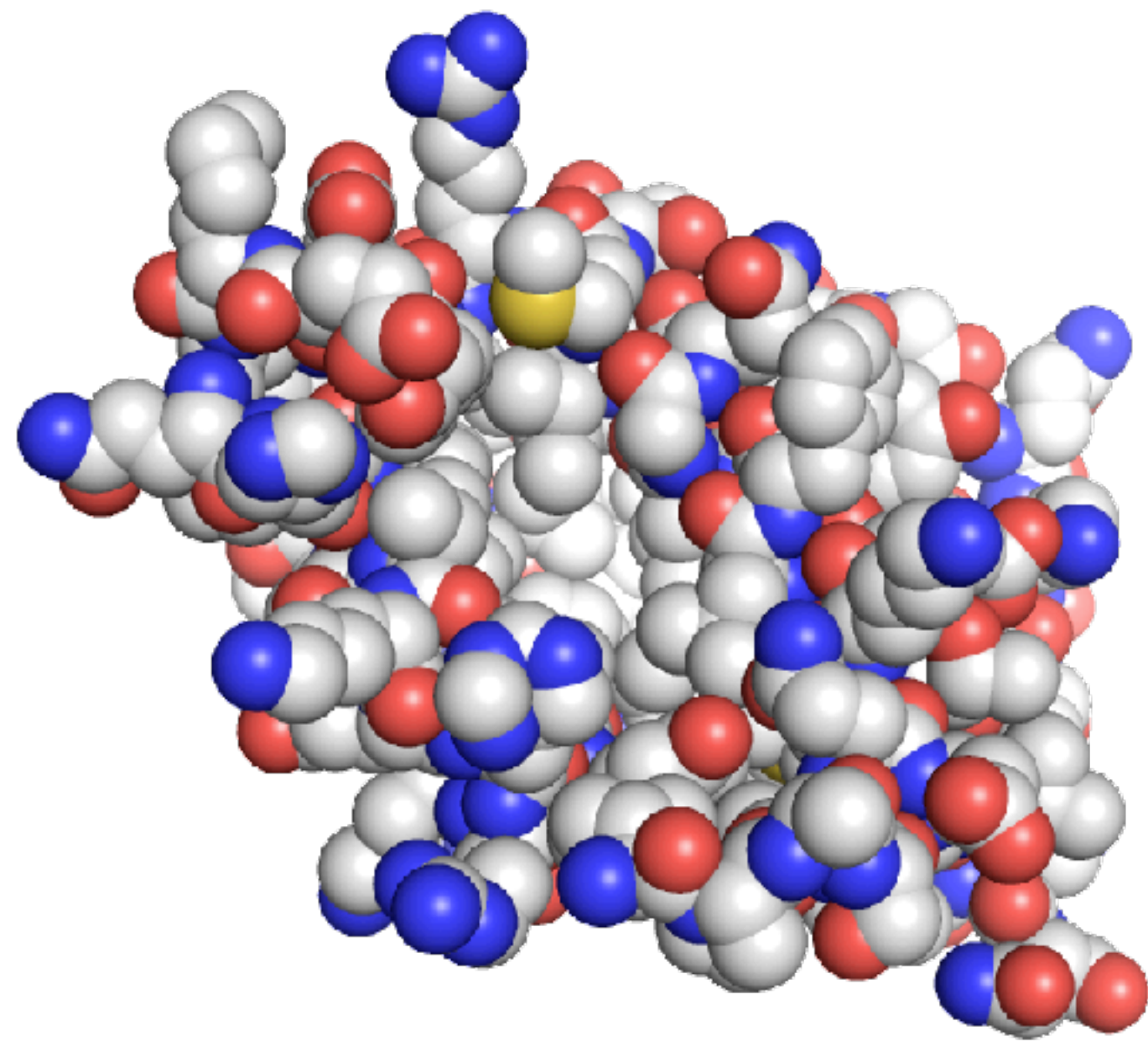
TENS    HUNDREDS    THOUSANDS

Source: Pharmaceutical Research and Manufacturers of America (http://phrma.org)

1. Does the compound do what you want it to?

2. Does the compound **not** do what you **don't** want it to?

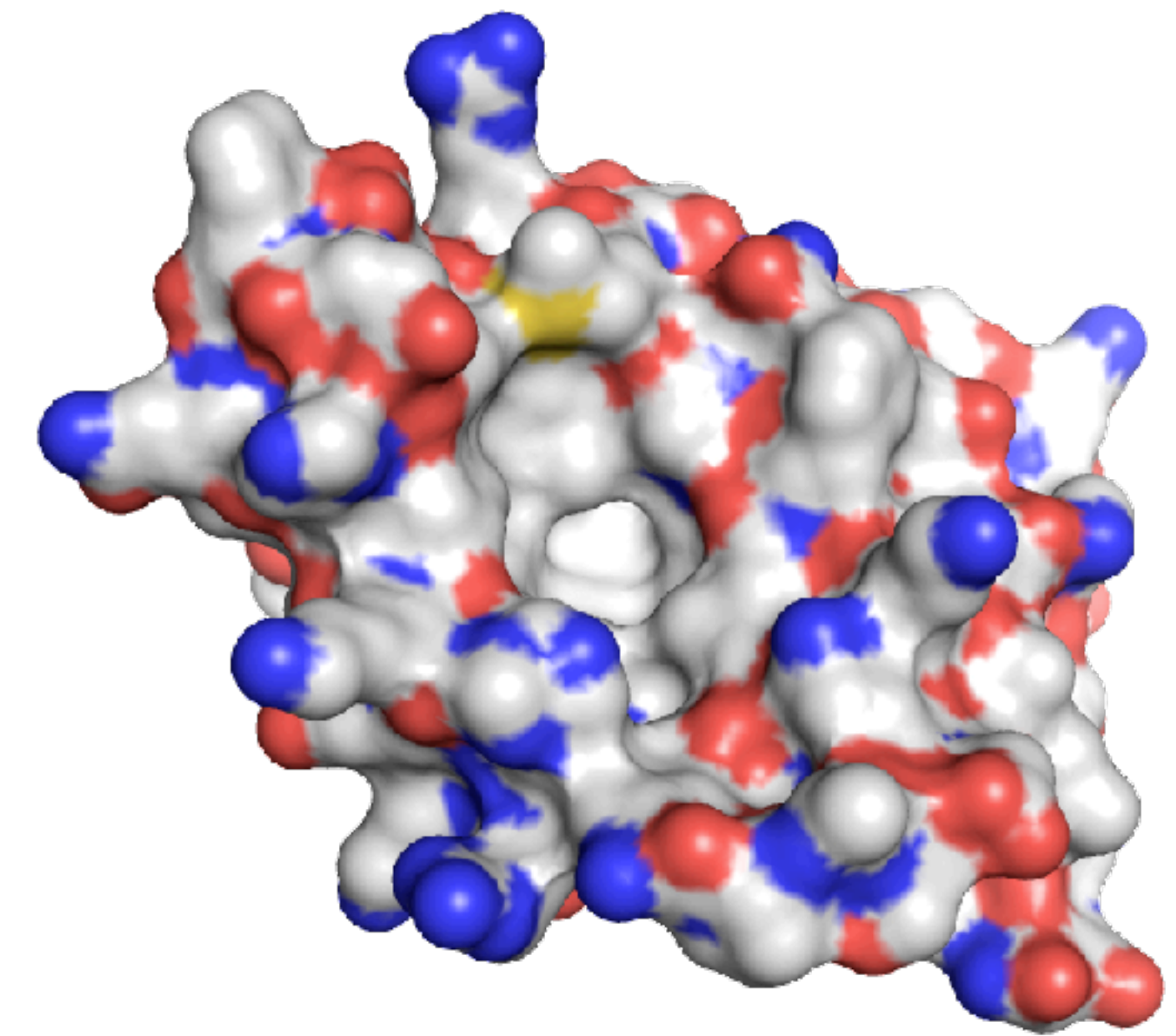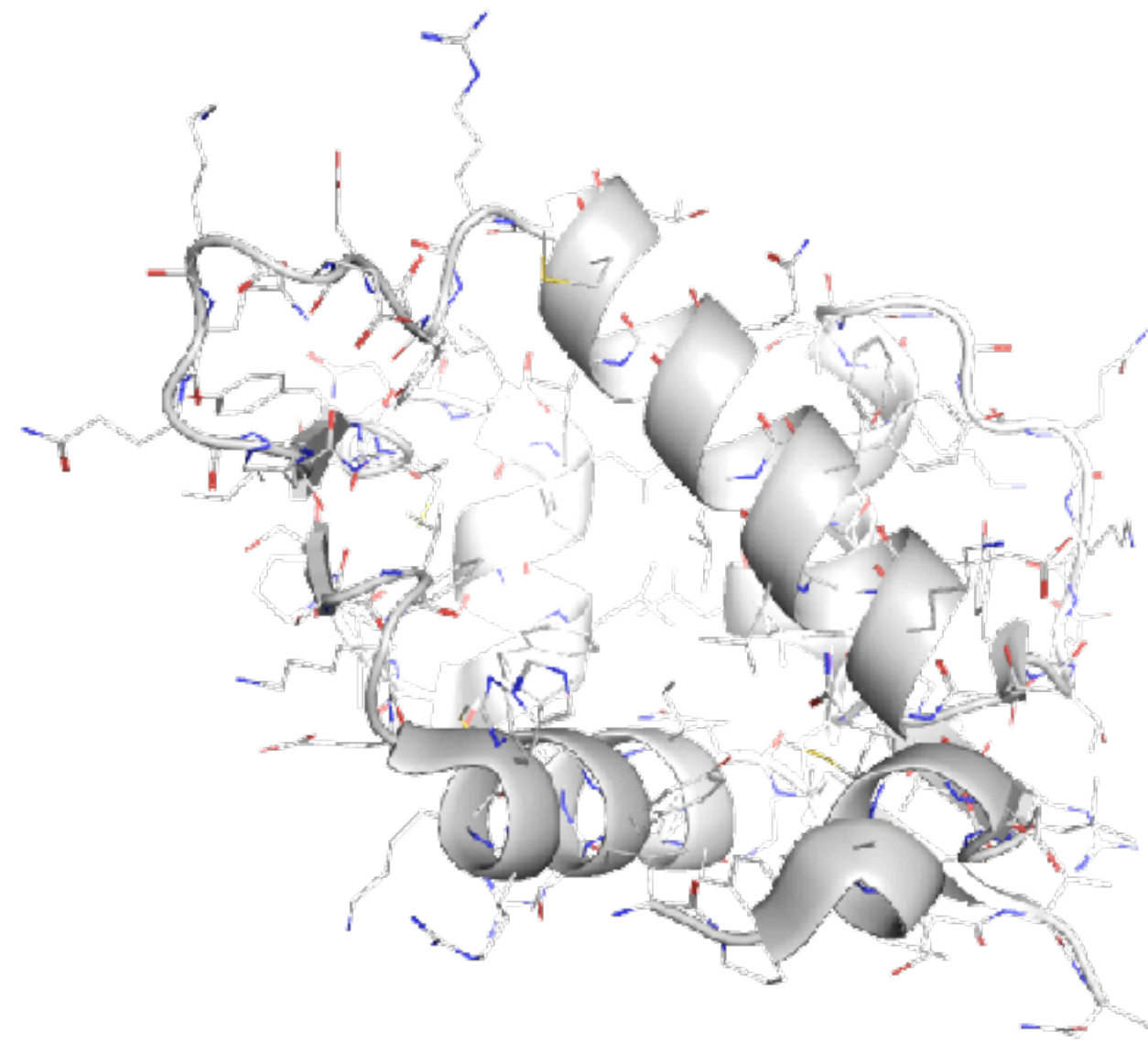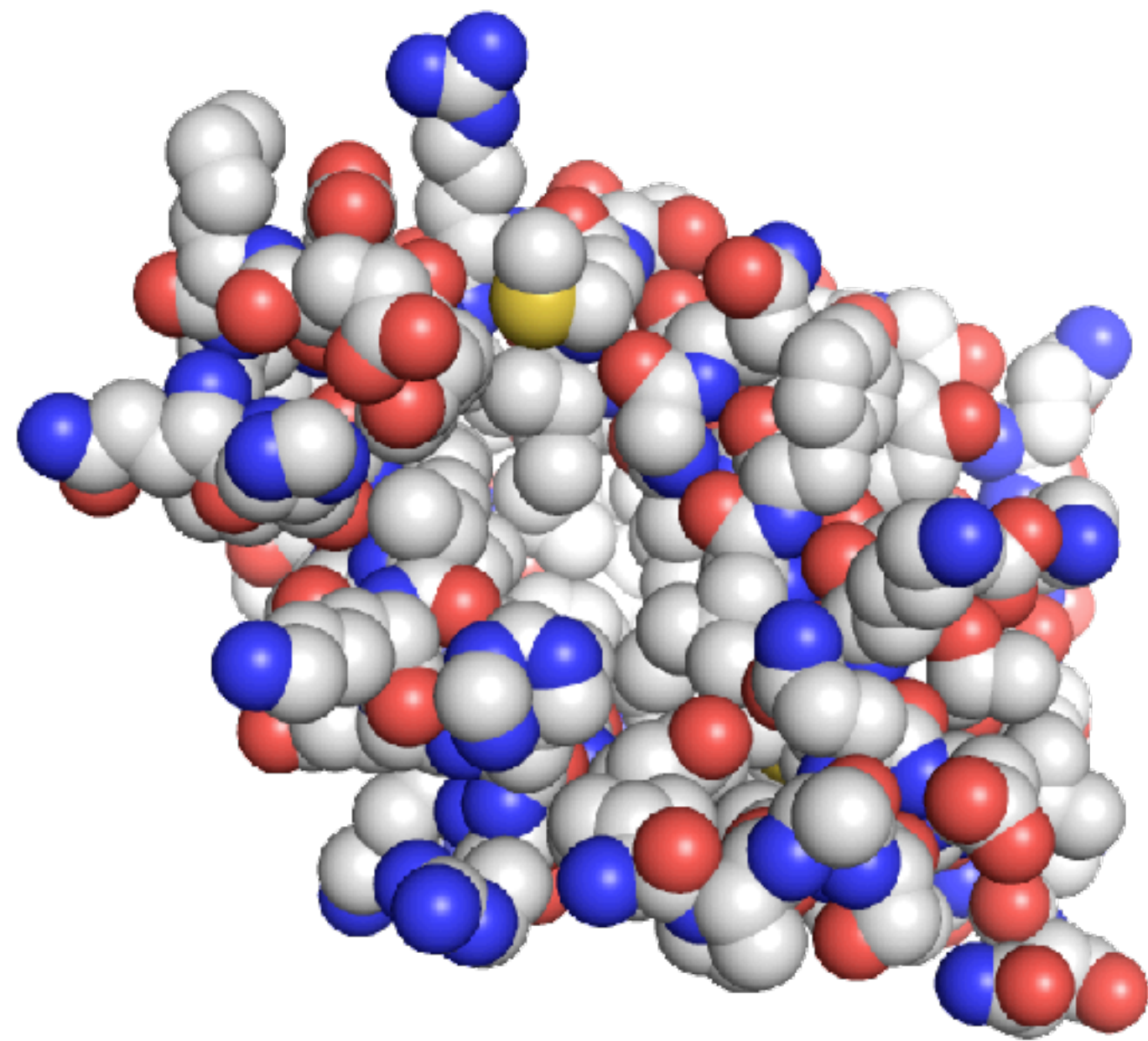3. Is what you want it to do the right thing?

# Protein Structures
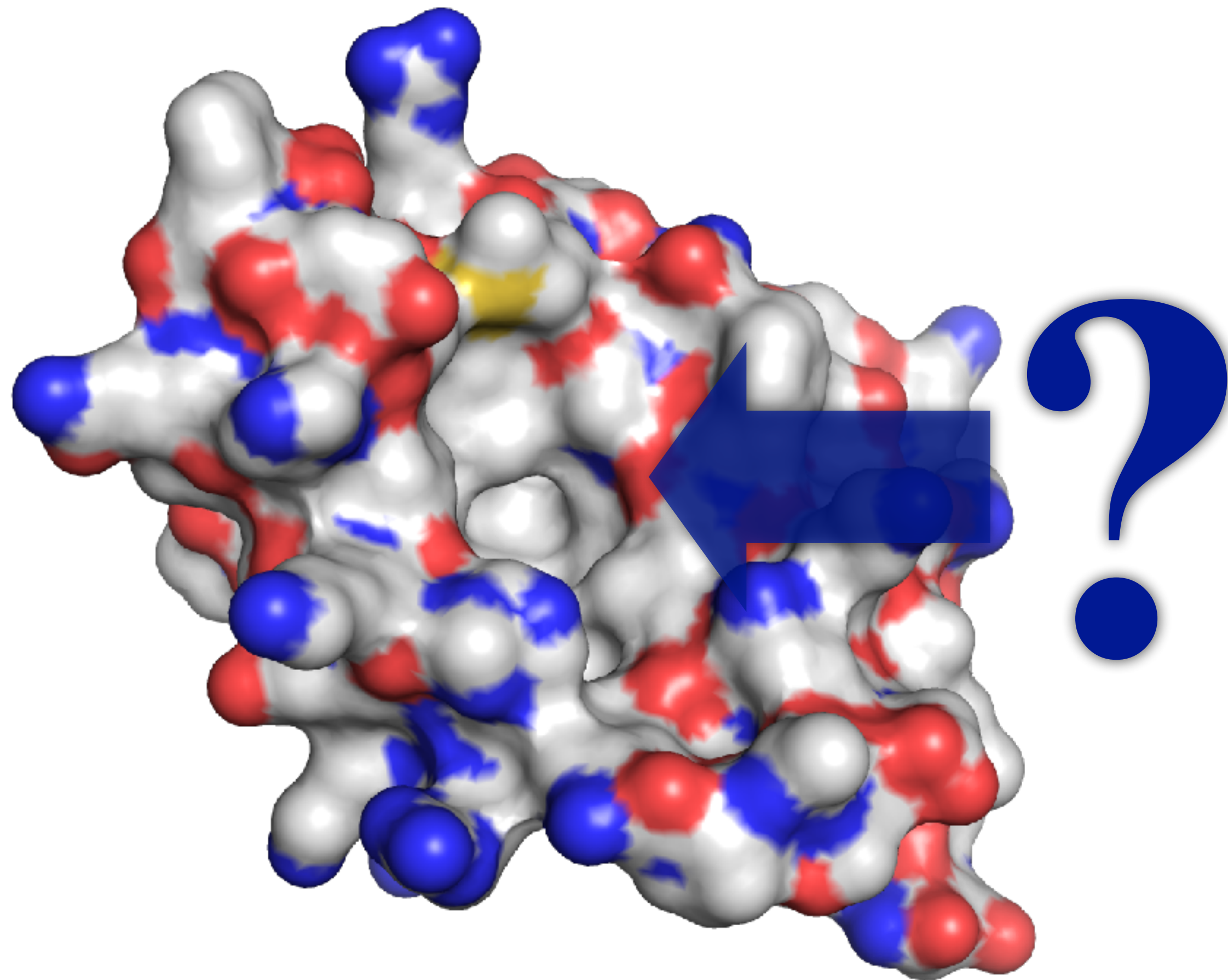
sequence → **structure** → function

# Protein Structures

sequence → **structure** → function
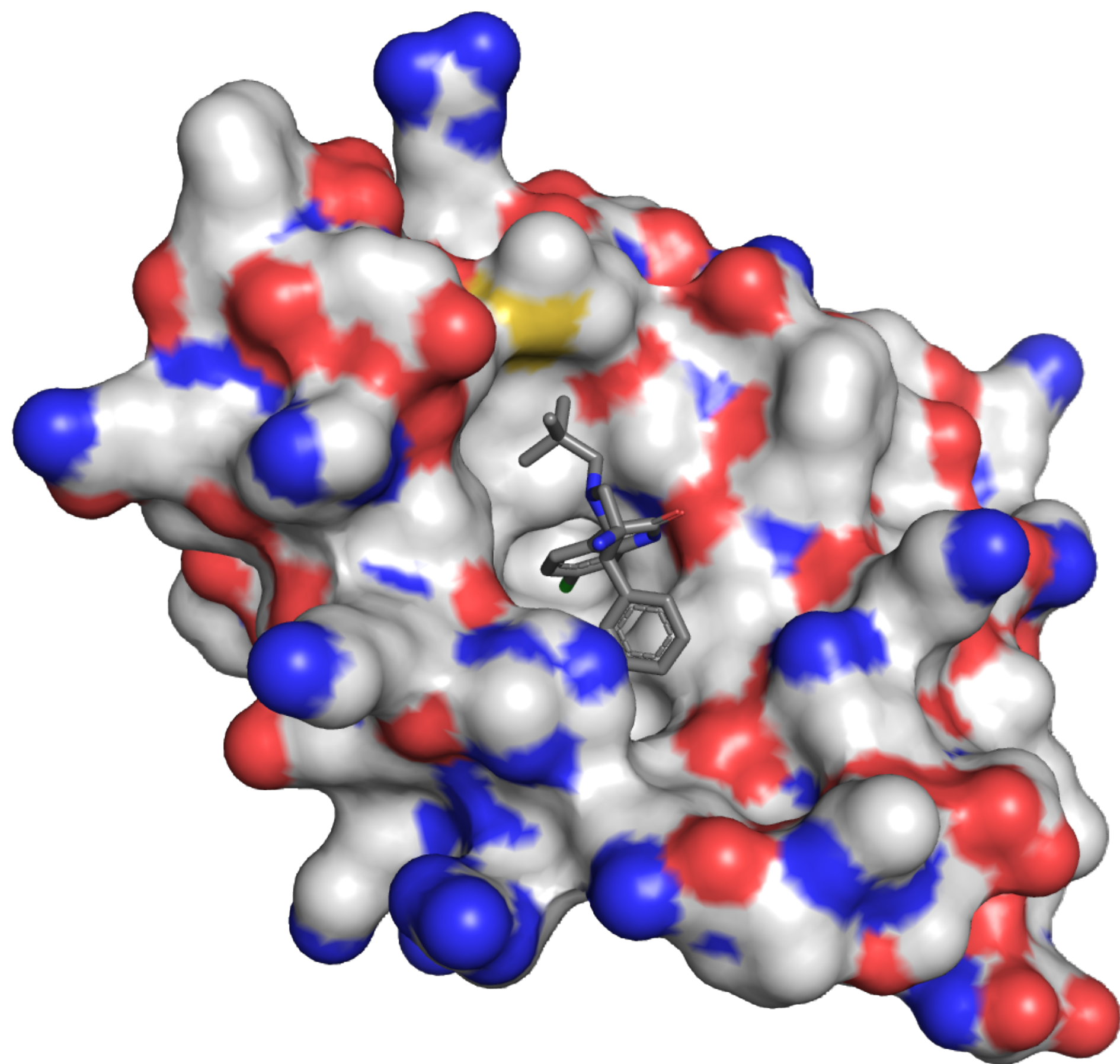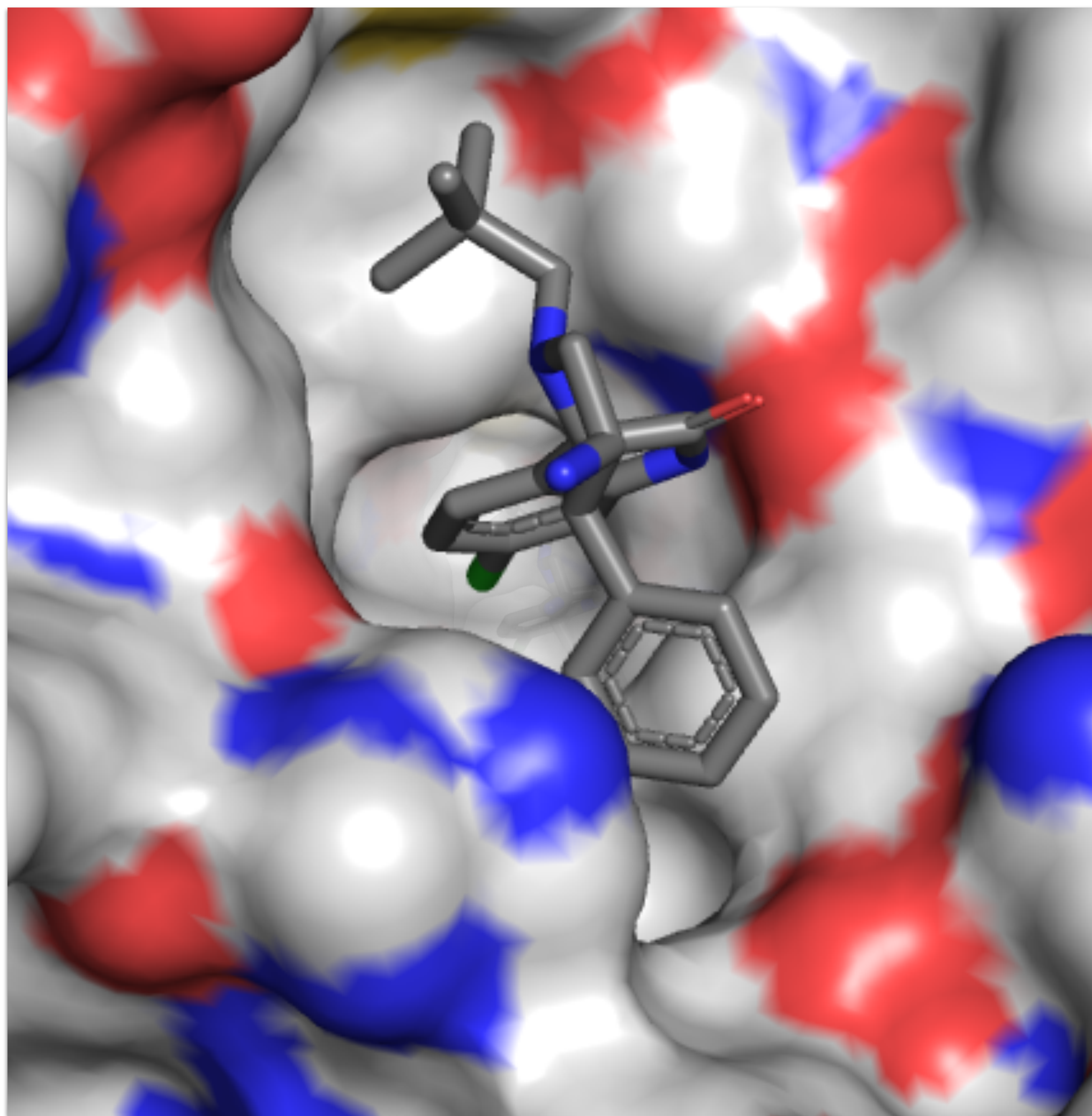
# Structure Based Drug Design



Unlike ligand based approaches, **generalizes to new targets**

Requires **molecular target** with **known structure** and **binding site**

# Structure Based Drug Design
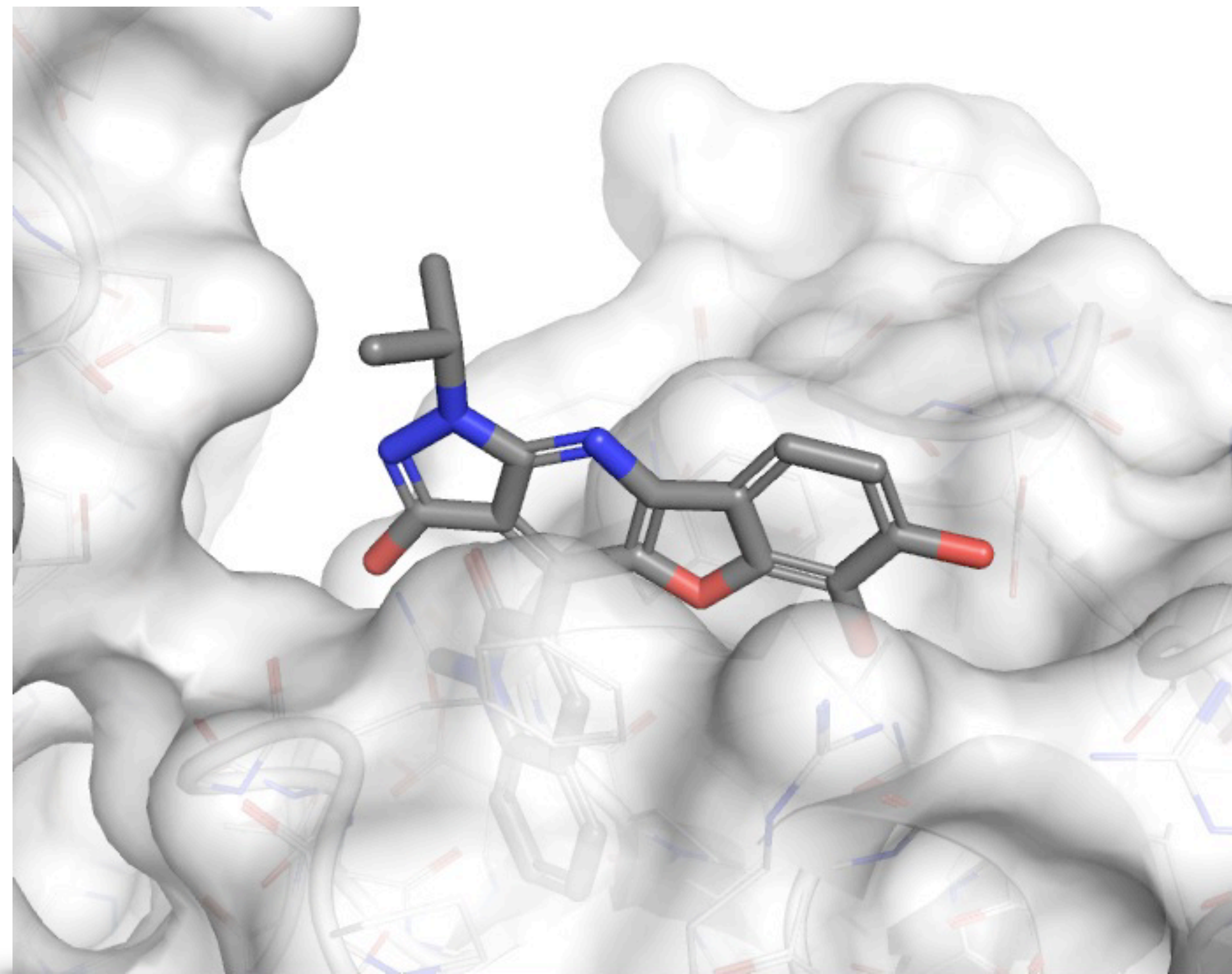


Unlike ligand based approaches, **generalizes to new targets**

Requires **molecular target** with **known structure** and **binding site**

# Structure Based Drug Design



Unlike ligand based approaches, **generalizes to new targets**
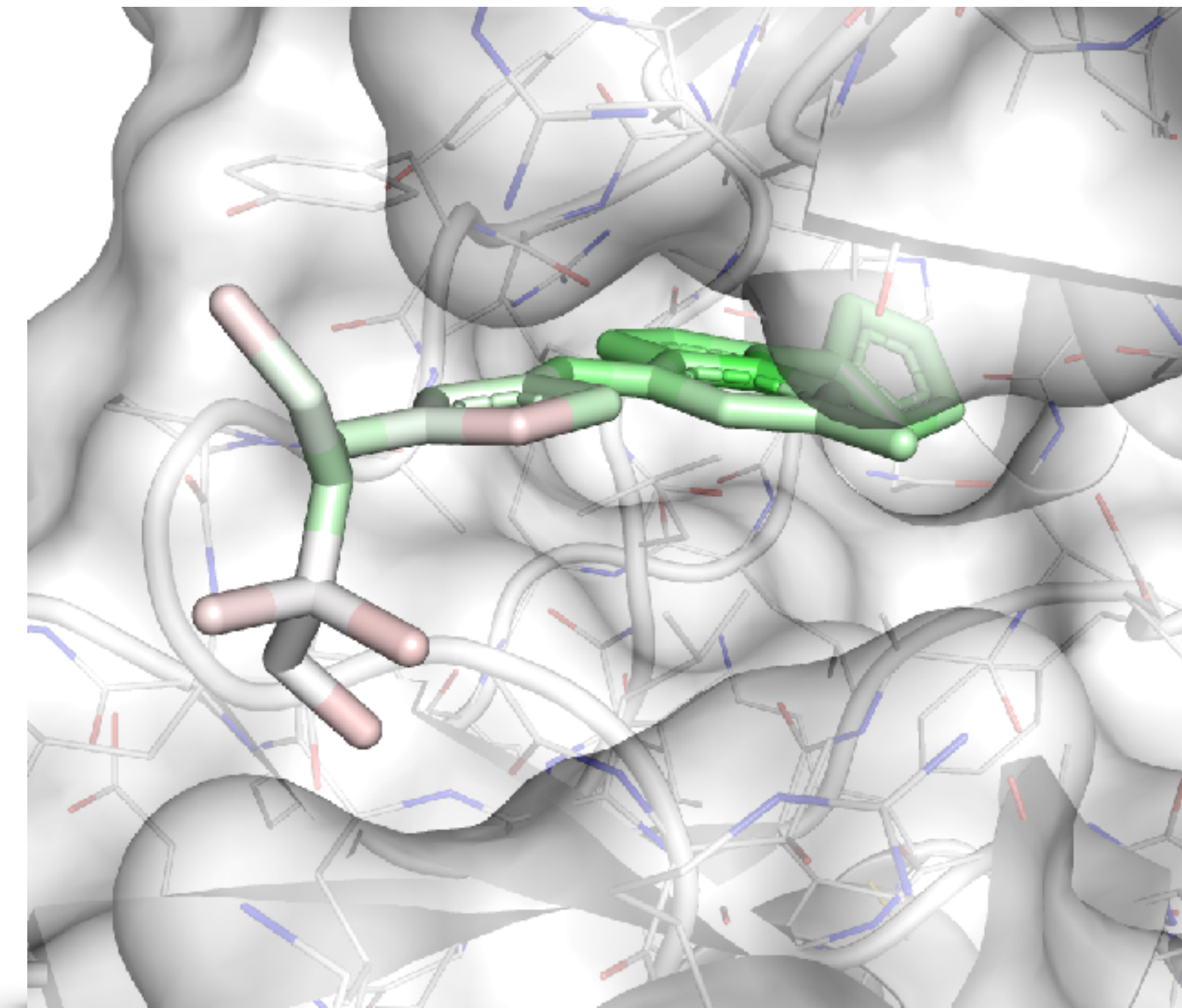
Requires **molecular target** with **known structure** and **binding site**

# Structure Based Drug Design

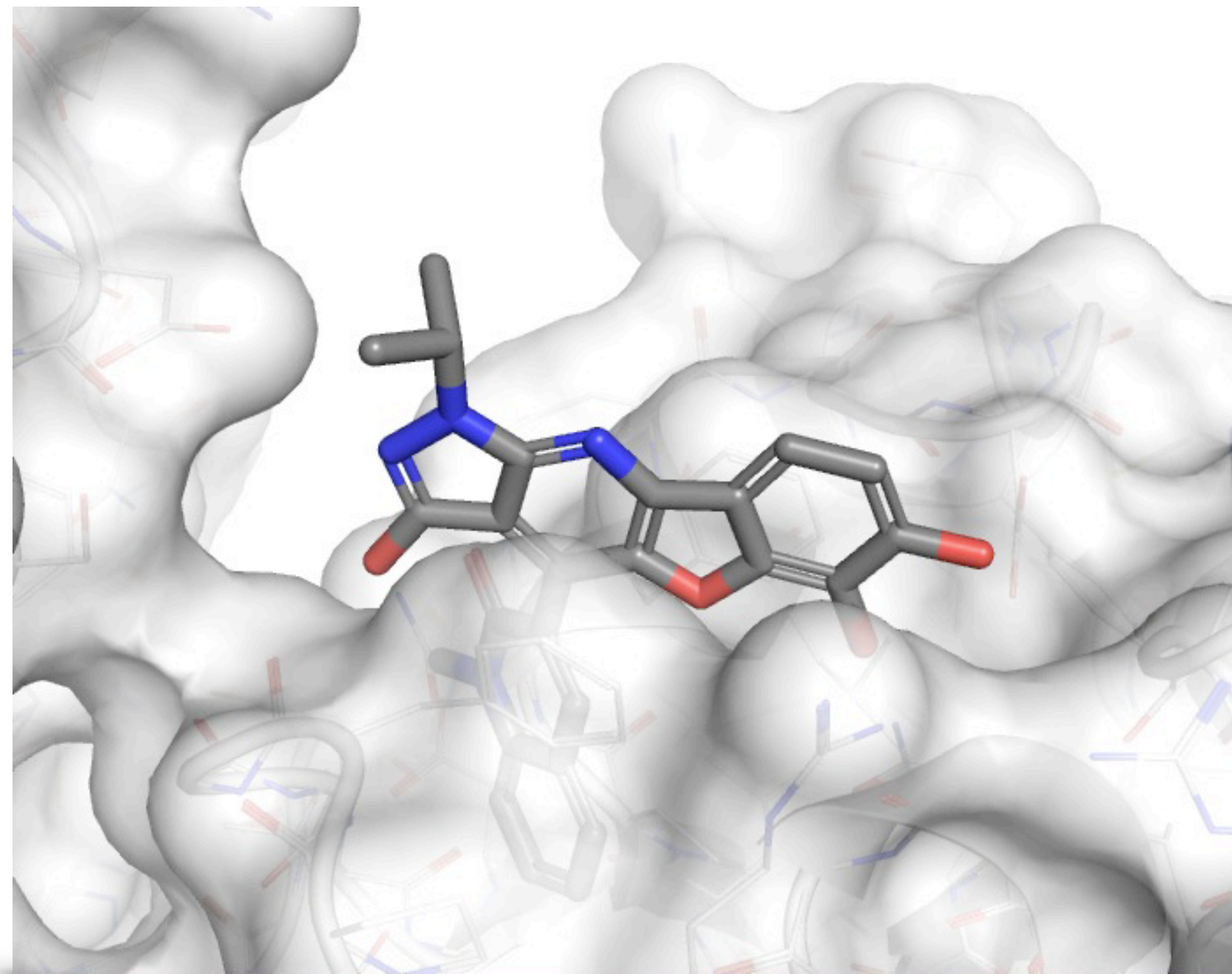**Virtual Screening**                    **Lead Optimization**



Pose Prediction        Binding Discrimination        Affinity Prediction
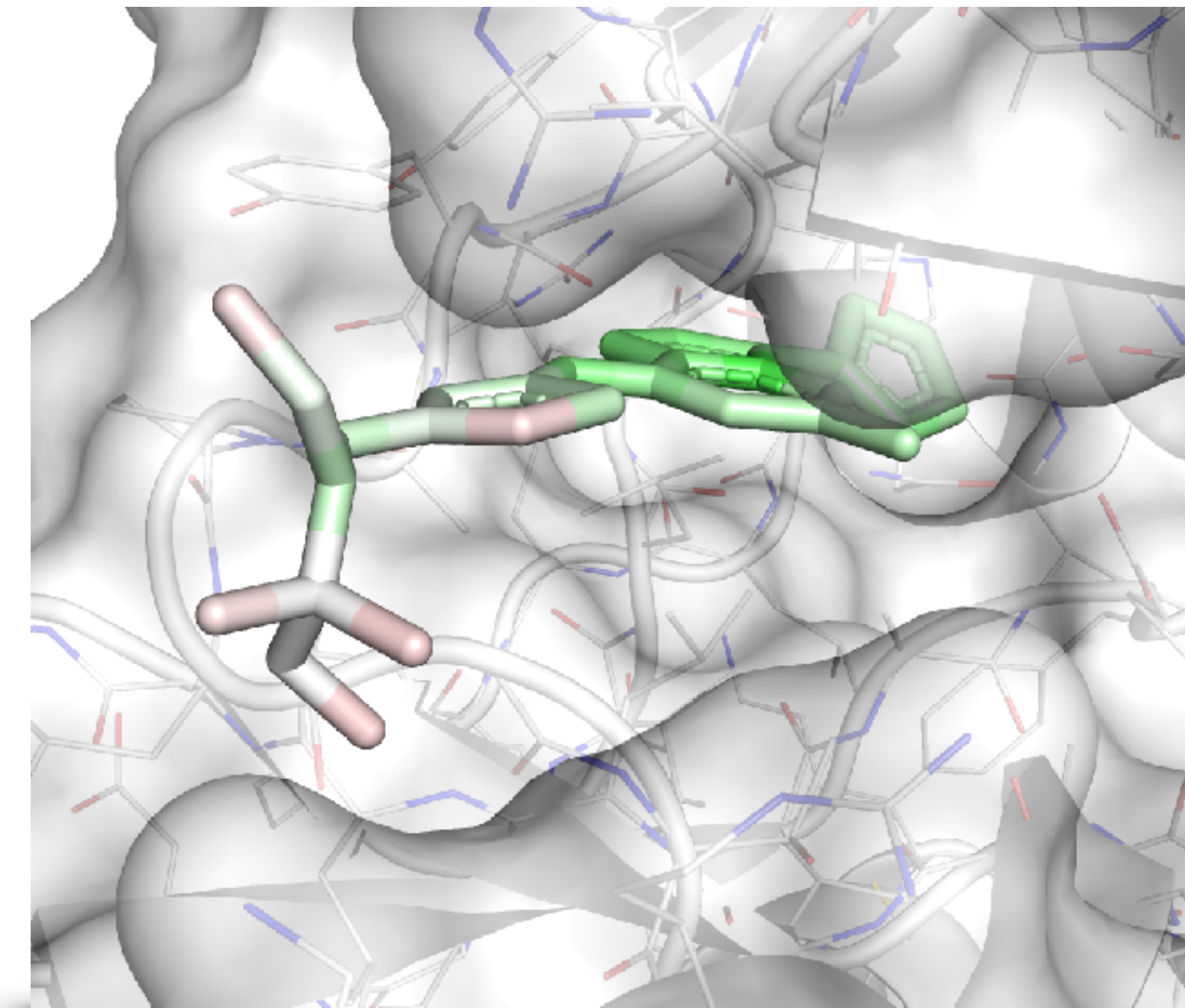
# Structure Based Drug Design

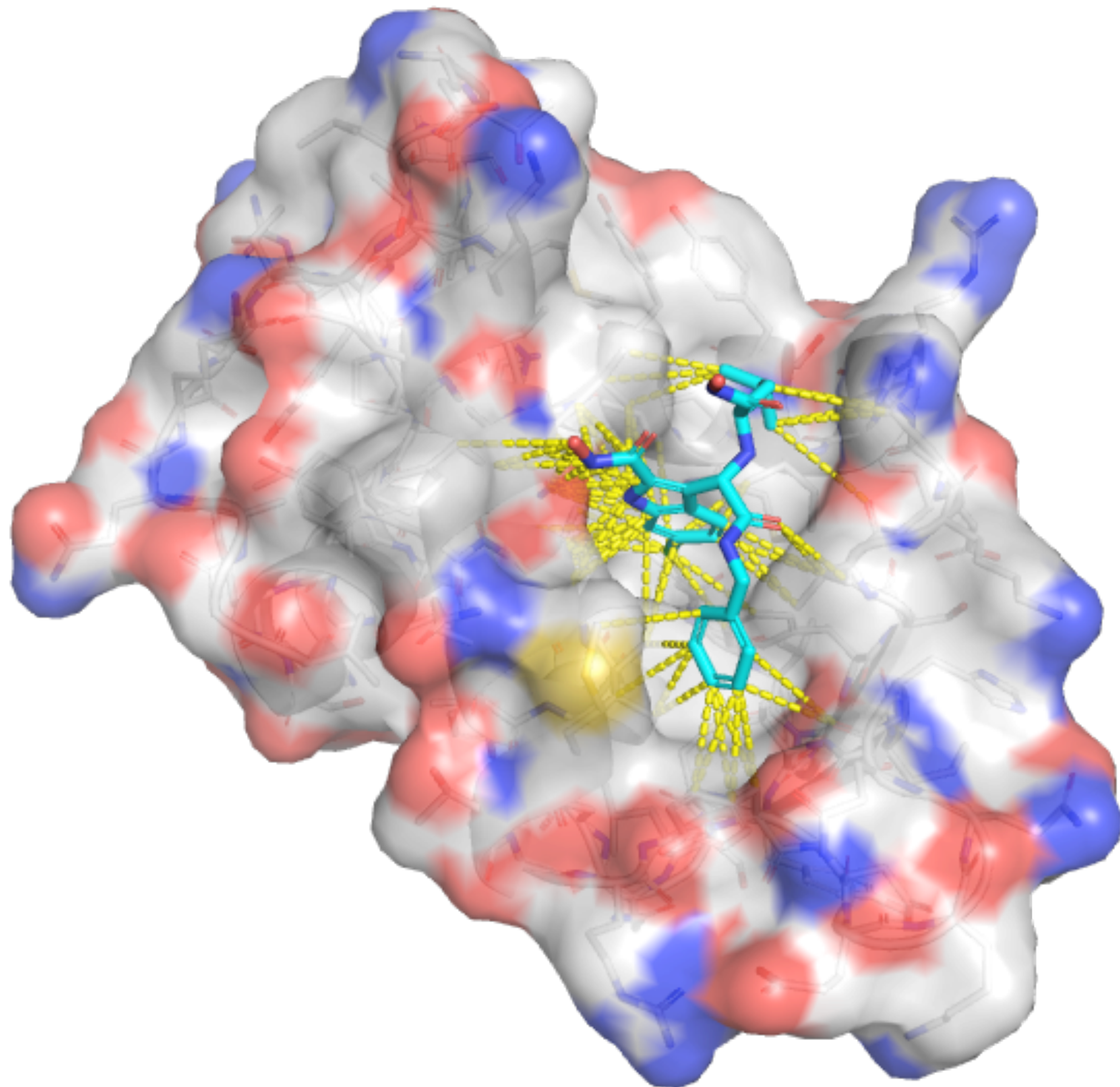**Virtual Screening**                                    **Lead Optimization**



Pose Prediction            Binding Discrimination        Affinity Prediction
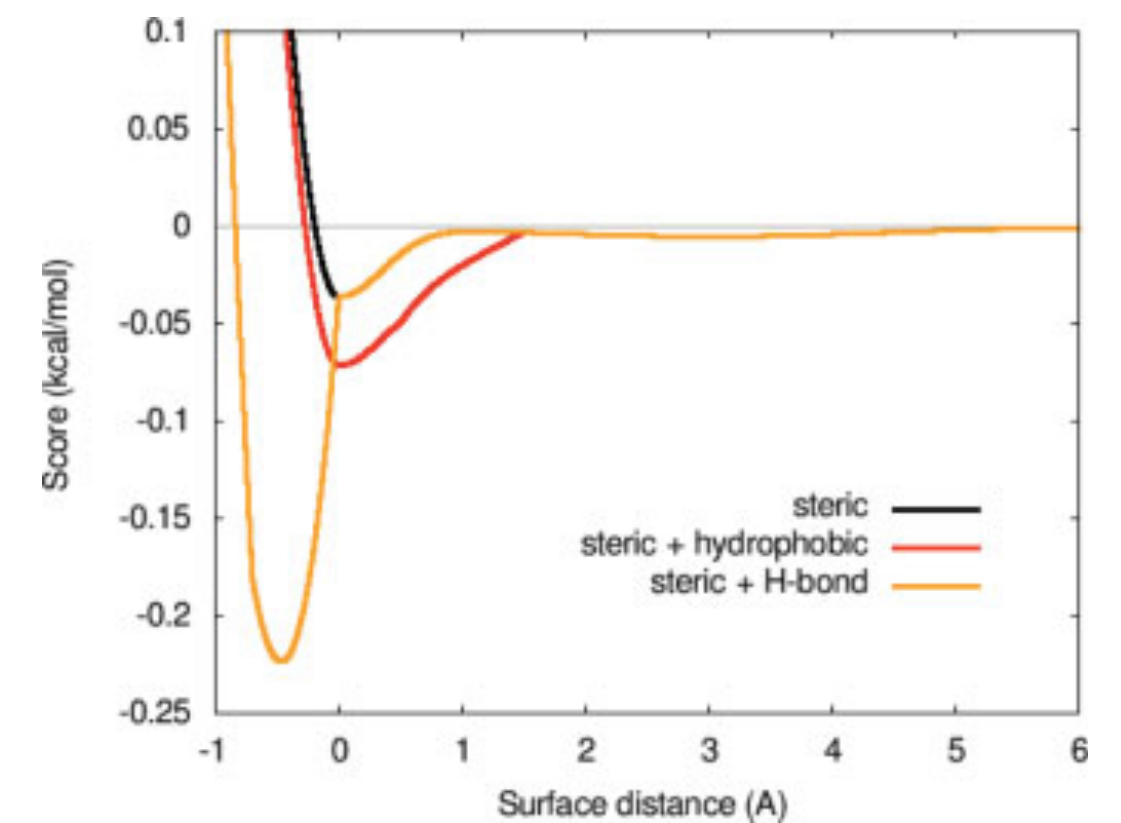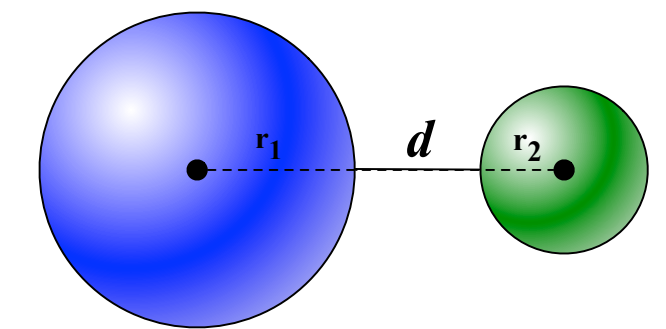
# Protein-Ligand Scoring

## AutoDock Vina

$$\text{gauss}_1(d) = w_{\text{guass}_1} e^{-(d/0.5)^2}$$

$$\text{gauss}_2(d) = w_{\text{guass}_2} e^{-((d-3)/2)^2}$$

$$\text{repulsion}(d) = \begin{cases} w_{\text{repulsion}} d^2 & d < 0 \\ 0 & d \geq 0 \end{cases}$$

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

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

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

7

# Can we do better?

Accurate pose prediction, binding discrimination, **and** affinity prediction without sacrificing performance?
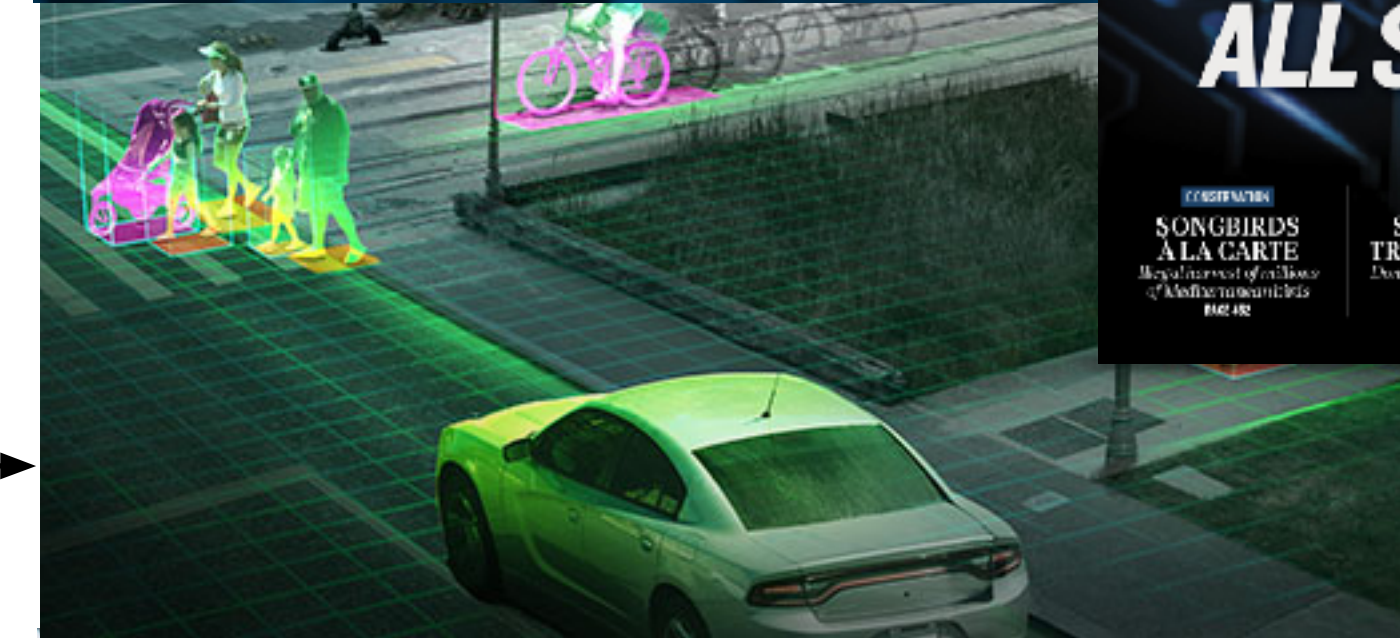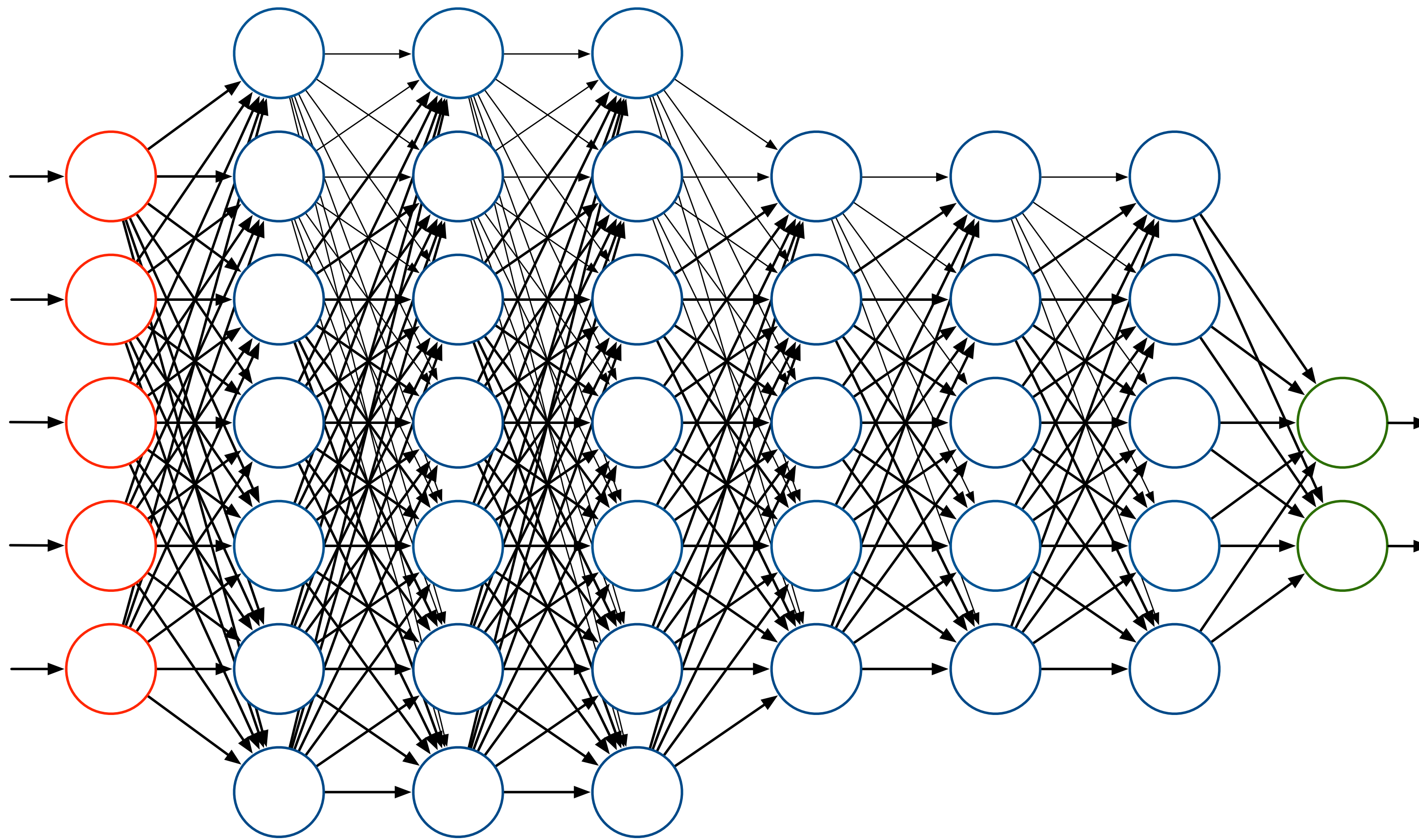
# Can we do better?

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

**Key Idea:** Leverage "big data"
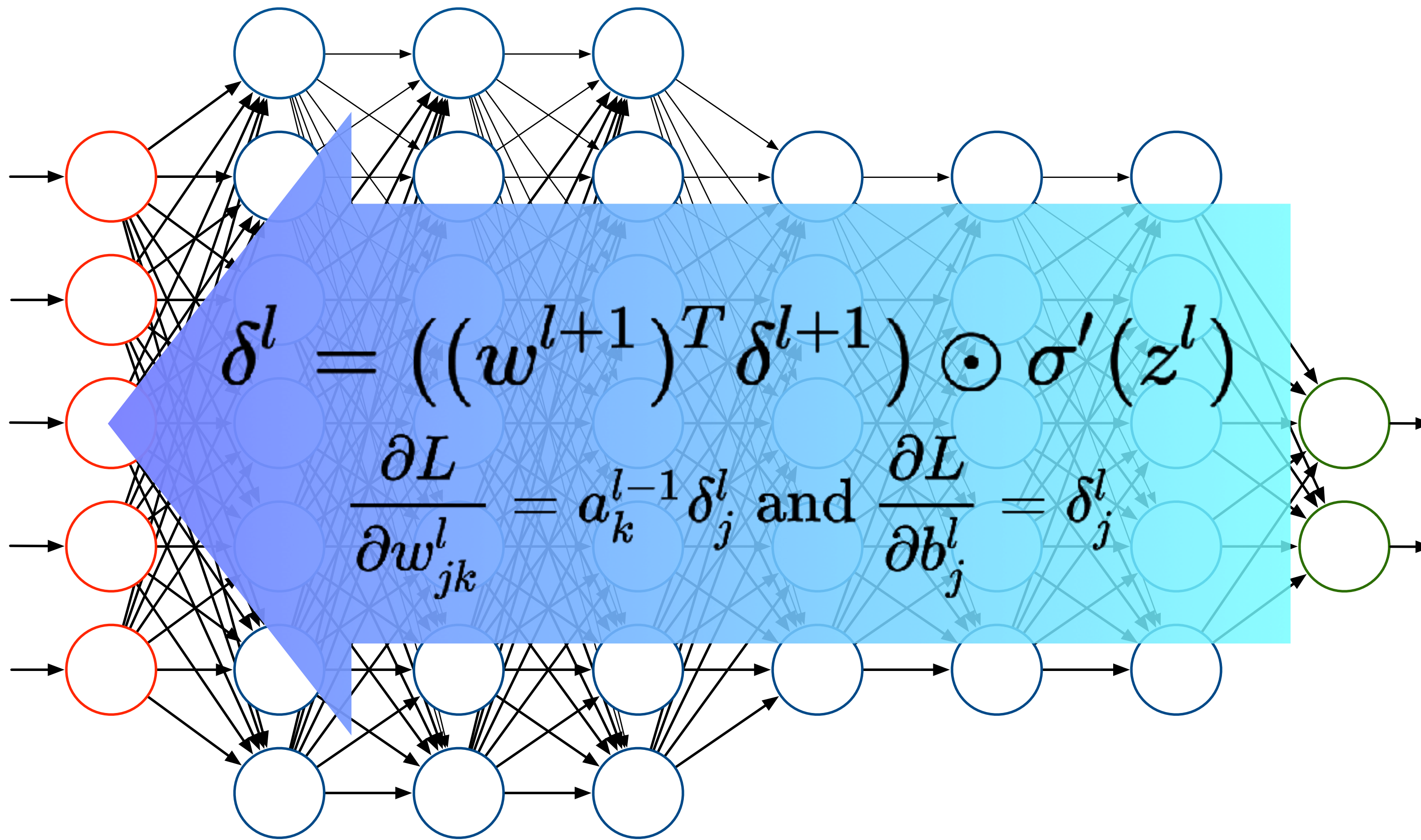  • 231,655,275 bioactivities in PubChem
  • 125,526 structures in the PDB
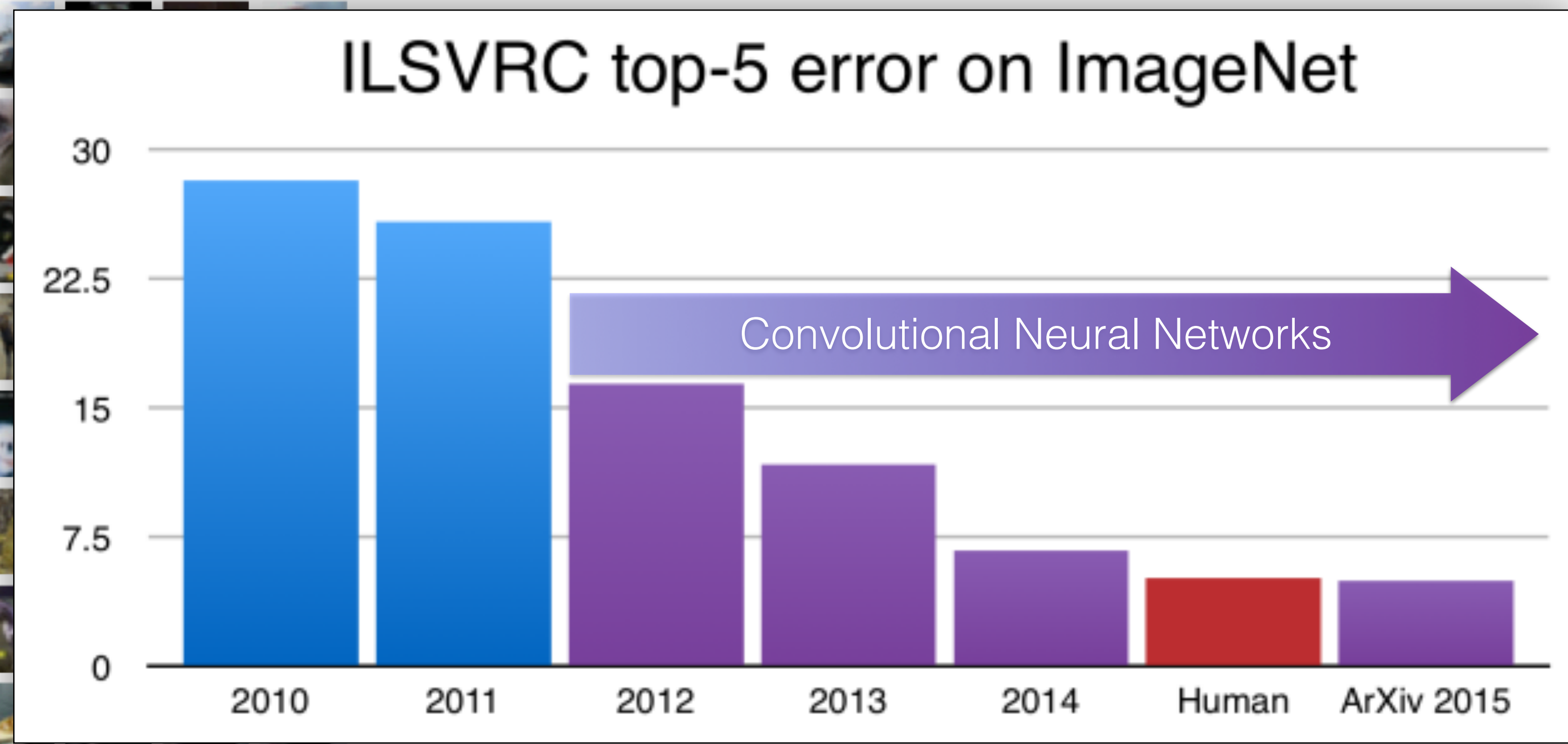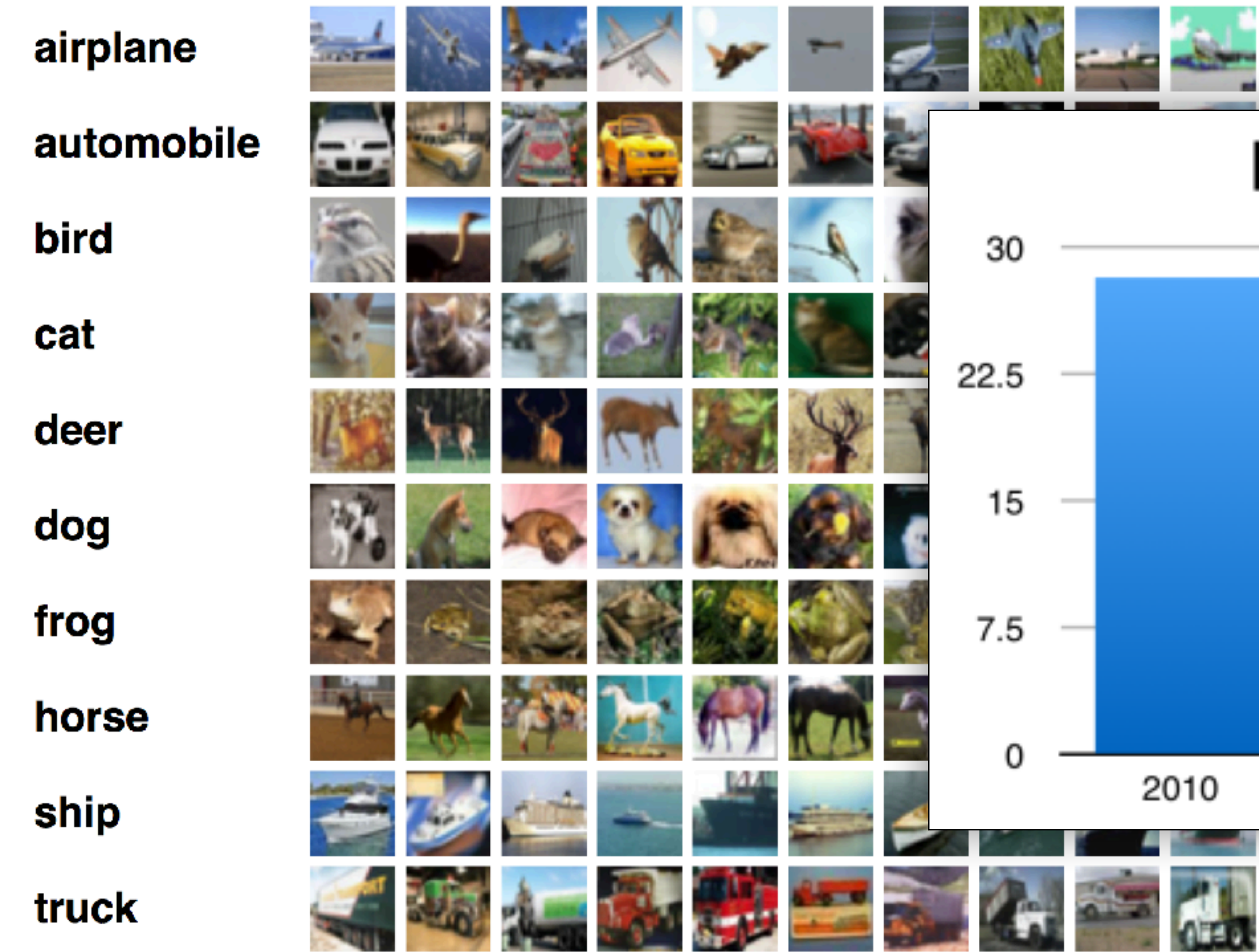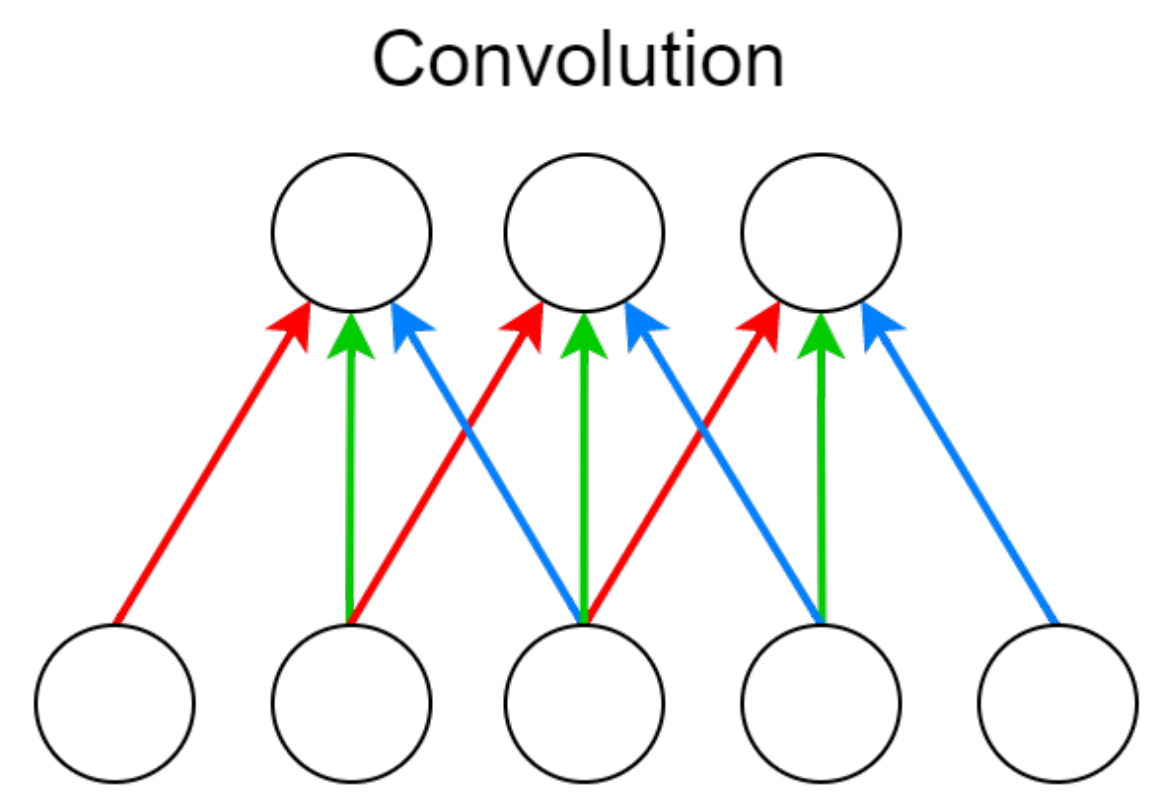  • 16,179 annotated complexes in PDBbind

# Deep Learning

# 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

11

# CNNs for Protein-Ligand Scoring



**CNN**

Pose Prediction

Binding
Discrimination

Affinity Prediction

# CNNs for Protein-Ligand Scoring



**CNN**

Pose Prediction

Binding
Discrimination

Affinity Prediction

# 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

$$A(d,r) = \begin{cases} e^{-\frac{2d^2}{r^2}} & 0 \le d < r \\[2ex] \frac{4}{e^2 r^2}d^2 - \frac{12}{e^2 r}d + \frac{9}{e^2} & r \le d < 1.5r \\[2ex] 0 & d \ge 1.5r \end{cases}$$



Gaussian

# Atom Types



| Ligand | Receptor |
|---|---|
| AliphaticCarbonXSHydrophobe | AliphaticCarbonXSHydrophobe |
| AliphaticCarbonXSNonHydrophobe | AliphaticCarbonXSNonHydrophobe |
| AromaticCarbonXSHydrophobe | AromaticCarbonXSHydrophobe |
| AromaticCarbonXSNonHydrophobe | AromaticCarbonXSNonHydrophobe |
| Bromine | Calcium |
| Chlorine | Iron |
| Fluorine | Magnesium |
| Iodine | Nitrogen |
| Nitrogen | NitrogenXSAcceptor |
| NitrogenXSAcceptor | NitrogenXSDonor |
| NitrogenXSDonor | NitrogenXSDonorAcceptor |
| NitrogenXSDonorAcceptor | OxygenXSAcceptor |
| Oxygen | OxygenXSDonorAcceptor |
| OxygenXSAcceptor | Phosphorus |
| OxygenXSDonorAcceptor | Sulfur |
| Phosphorus | Zinc |
| Sulfur | |
| SulfurAcceptor | |

15

# Training Data

## Pose Prediction





337 protein-ligand complexes
- curated for electron density
- diverse targets
- <10µM affinity
- **generate poses** with Vina
  - 745  <2Å RMSD (actives)
  - 3251 >4Å RMSD (decoys)

12,484 protein-ligand complexes
- diverse targets
- wide range of affinities
- **generate poses** with AutoDock Vina
- include minimized crystal pose
  - 24,727  <2Å RMSD (actives)
  - 244,192 >4Å RMSD (decoys)

# Model Evaluation

**CSAR**: >90% similar targets kept in same fold

**PDBbind**: >80% similar targets kept in same fold

# 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 $O(N_{atoms})$ check



For example, consider subgrid 5:

| Atom mask: | 1 | 1 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| Exclusive scan: | 0 | 1 | 2 | 2 | 2 |
| Final indices: | 0 | 1 | | | |

18

# 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

# Cross-Validation Evaluation

# Pose Prediction (CSAR)

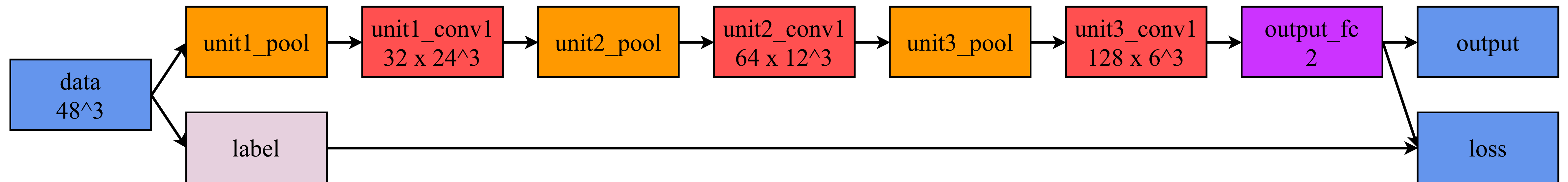# Pose Prediction (CSAR)



*inter*-target ranking

*intra*-target ranking

# Pose Prediction (PDBbind)

# Pose Prediction (PDBbind)



*inter*-target ranking

*intra*-target ranking

# Visualization

# Examples



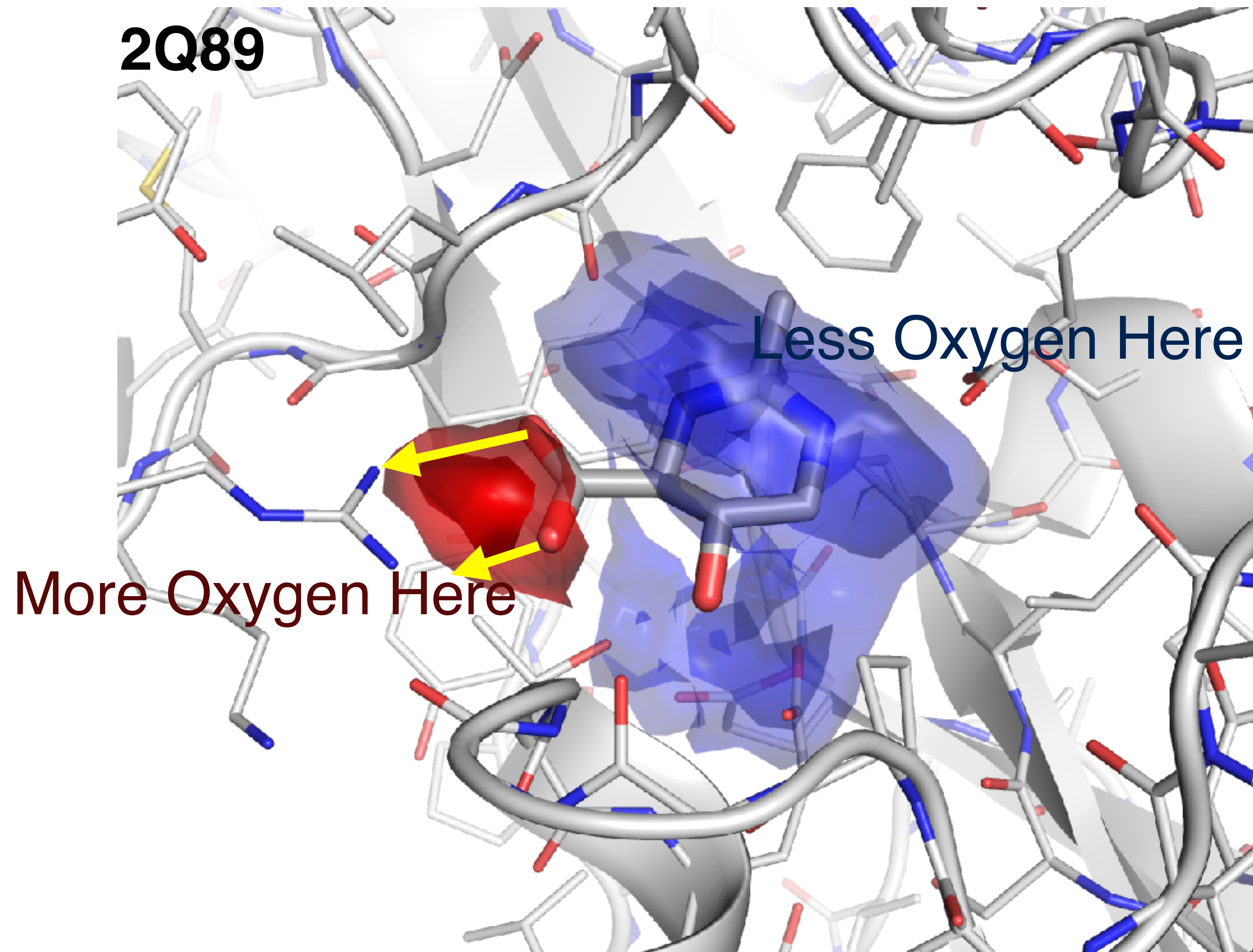3COY     2QMJ     3OZT

Partially Aligned Poses

# Beyond Scoring

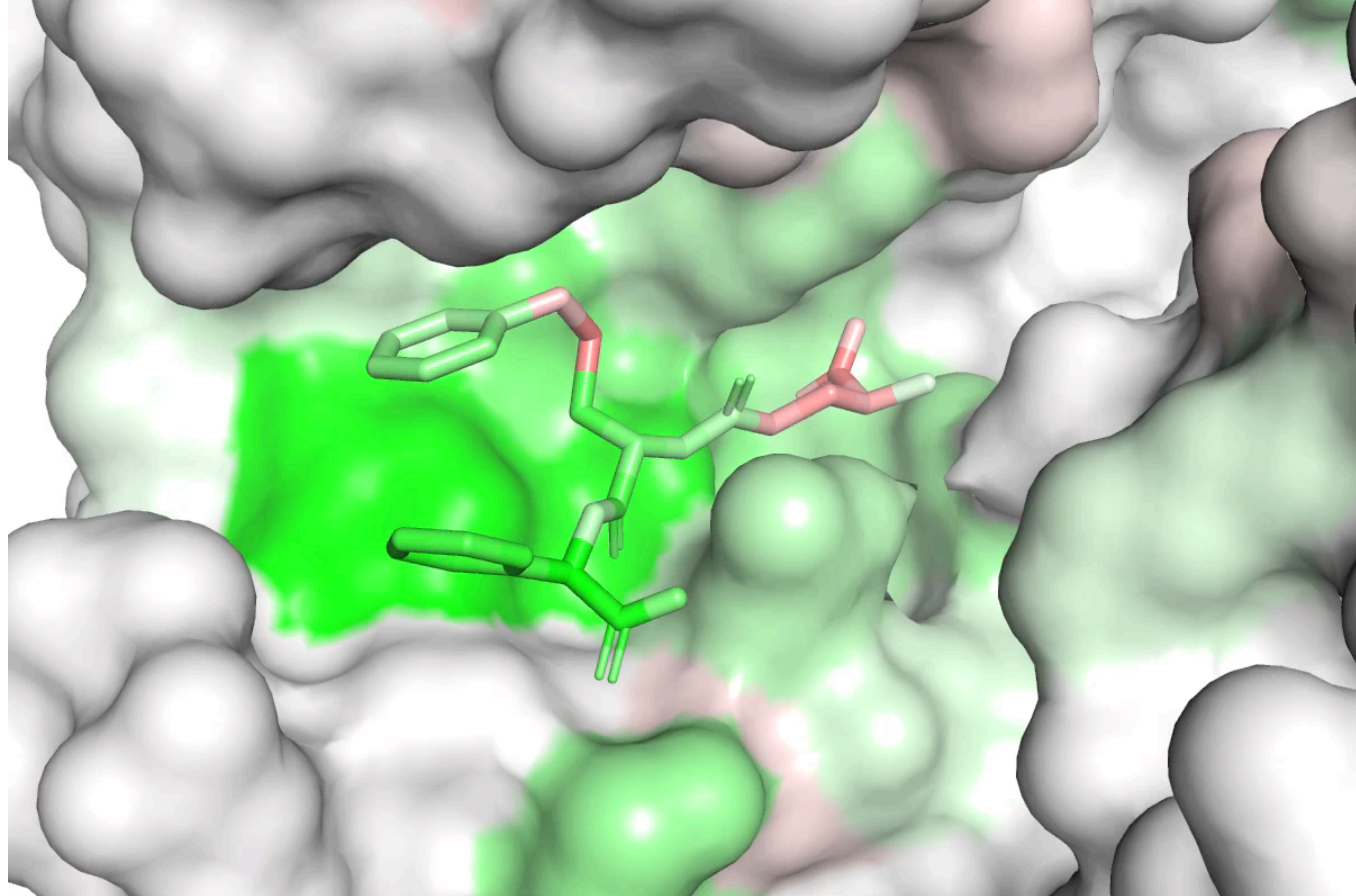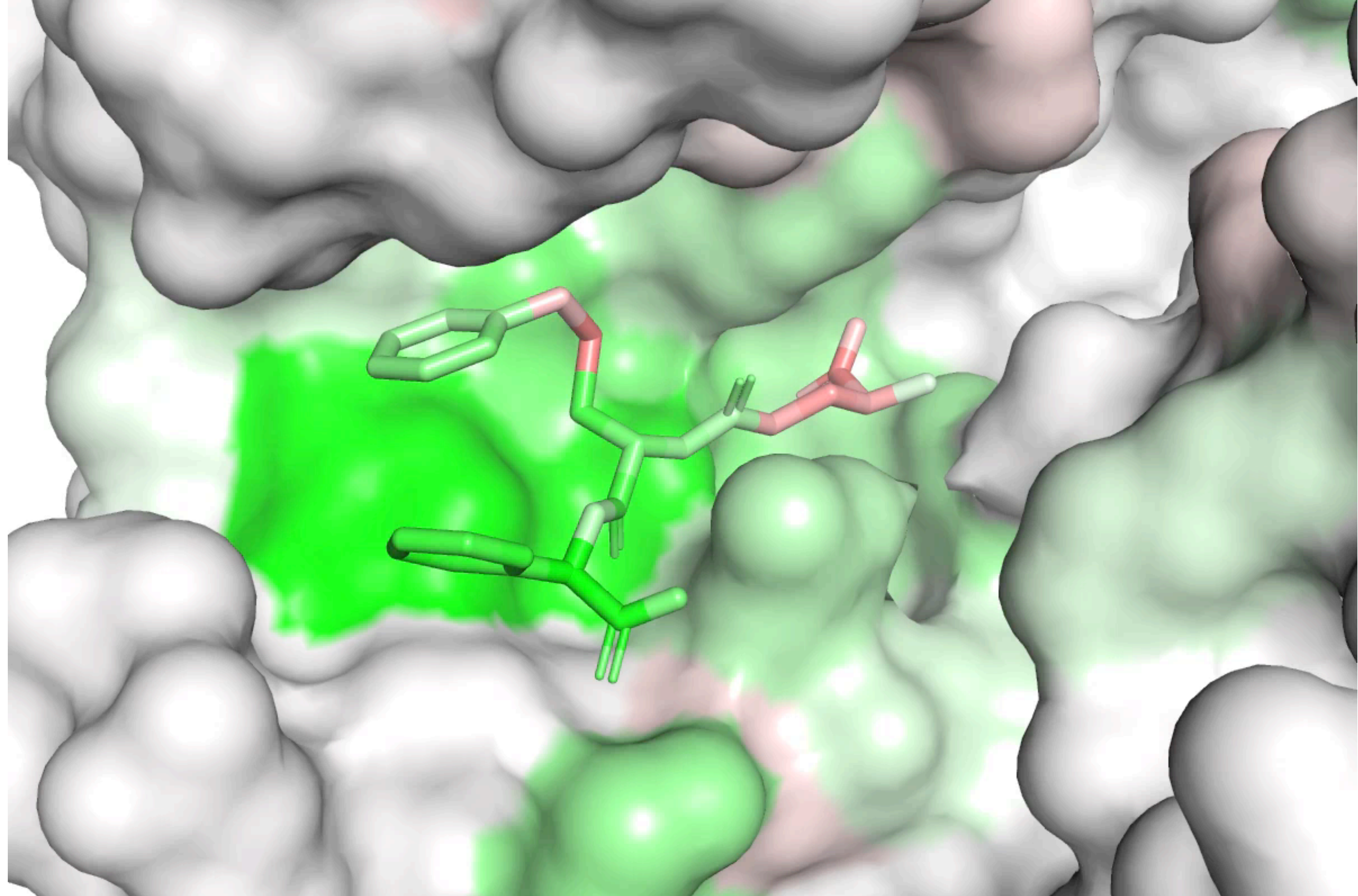# Beyond Scoring

# Beyond Scoring

# Beyond Scoring

# Beyond Scoring



**2Q89**

Less Oxygen Here

More Oxygen Here

$$\frac{\partial L}{\partial A} = \sum_{i \in G_A^{48^{\wedge}3}} \frac{\partial L}{\partial G_i} \frac{\partial G_i}{\partial D} \frac{\partial D}{\partial A}$$
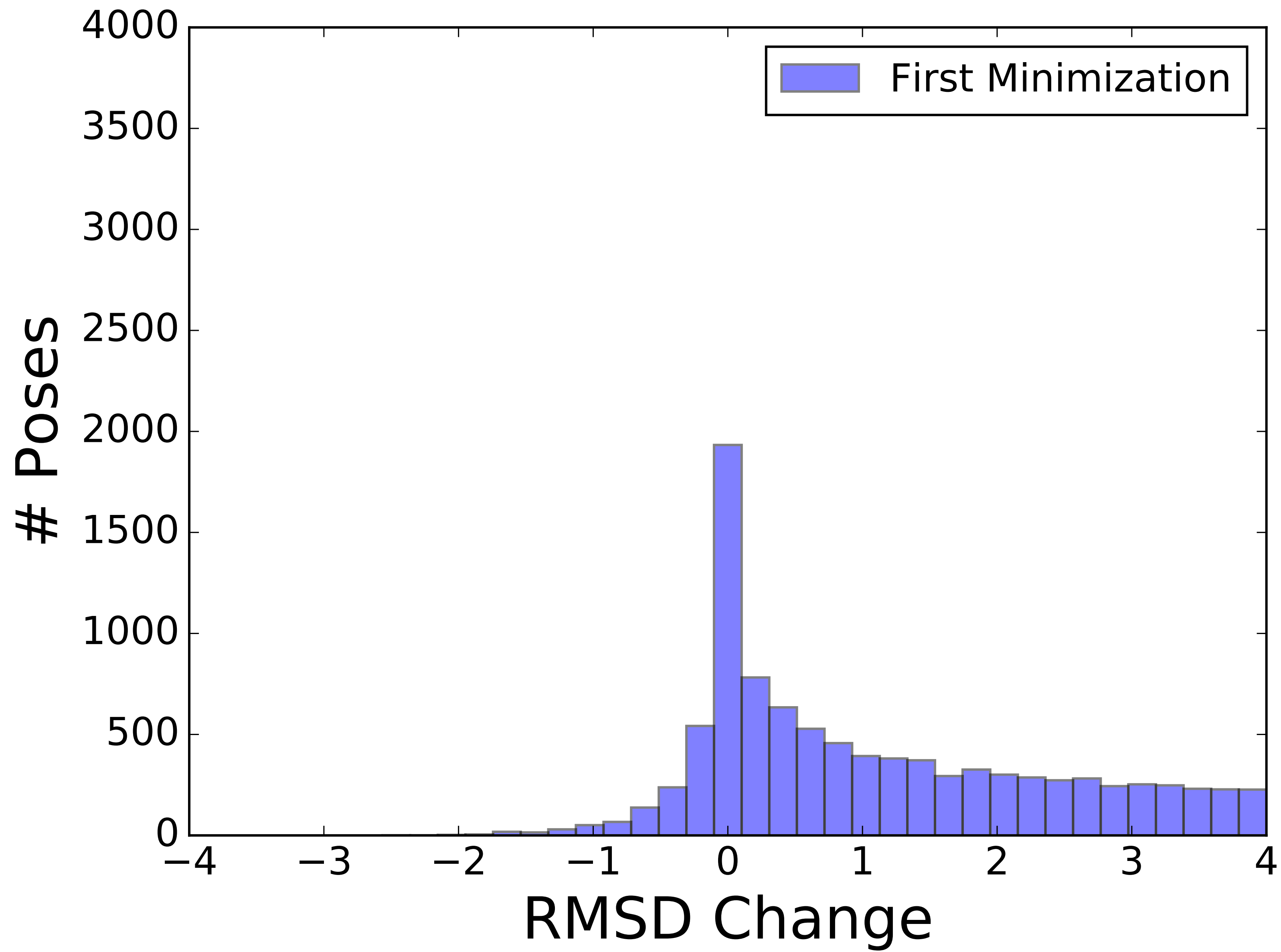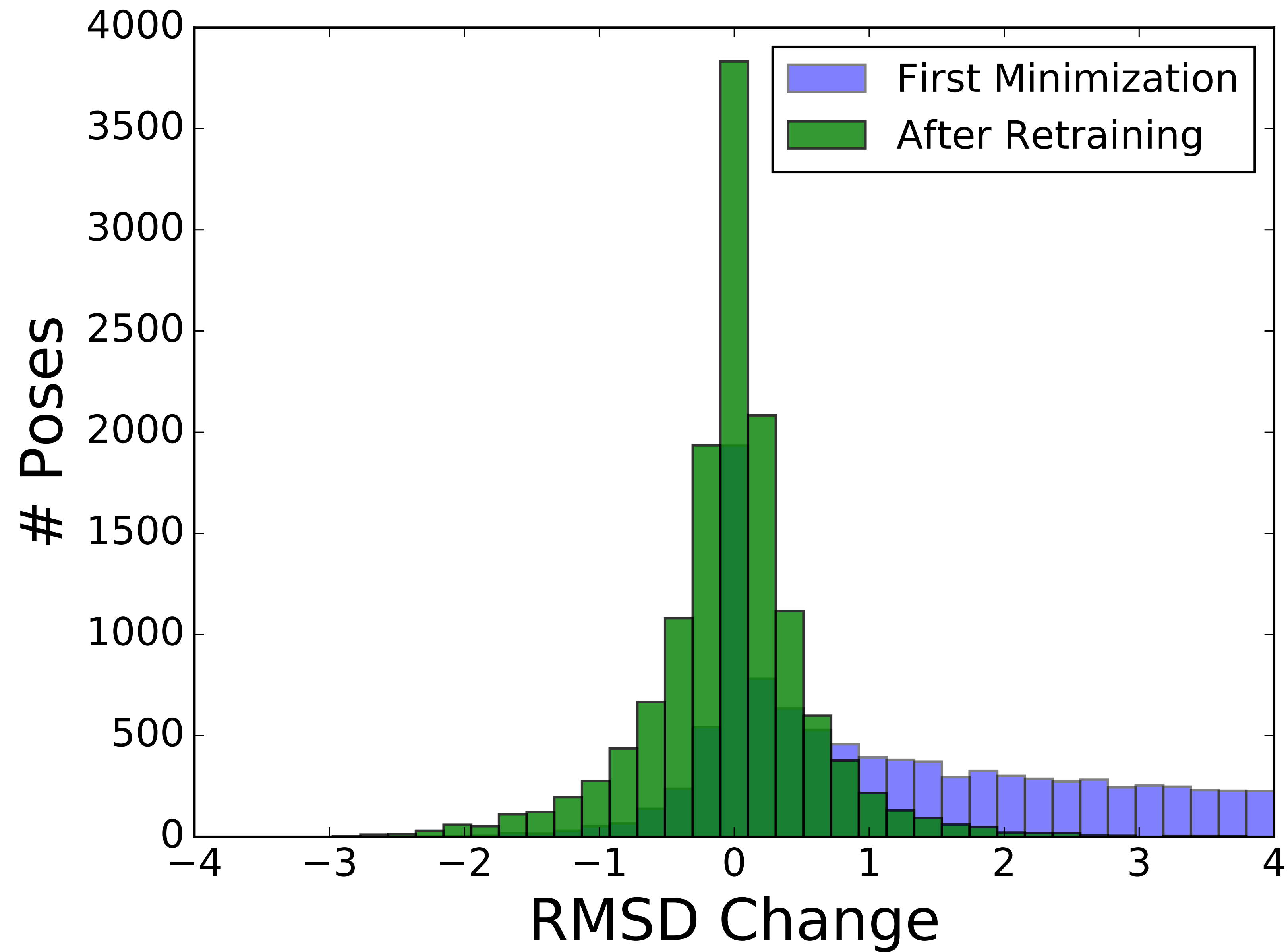
data

unit1_pool

label

28

# The Future

Pose Selection

**Iterative Training**

Pose *Generation*

**Iterative Training**
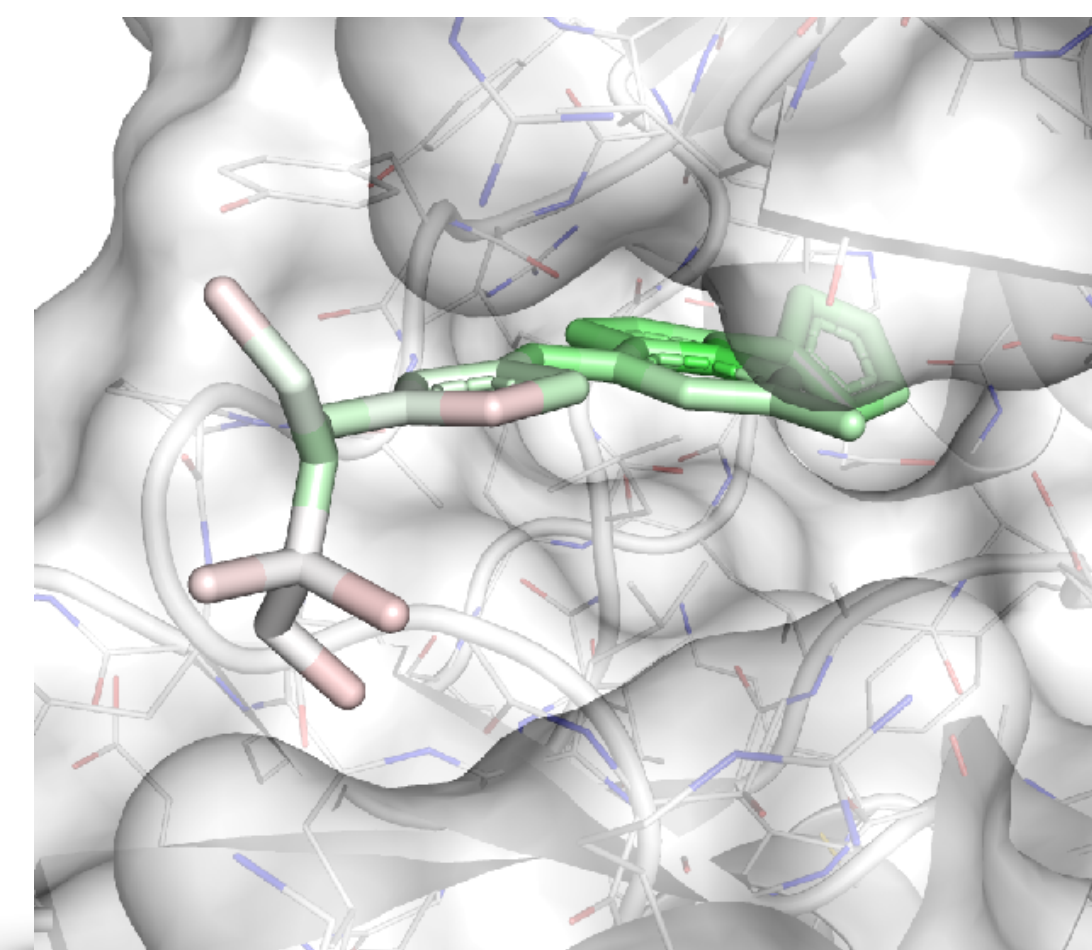
*Compound* Generation



**Virtual Screening**

**Lead Optimization**

# The Future

Pose Selection

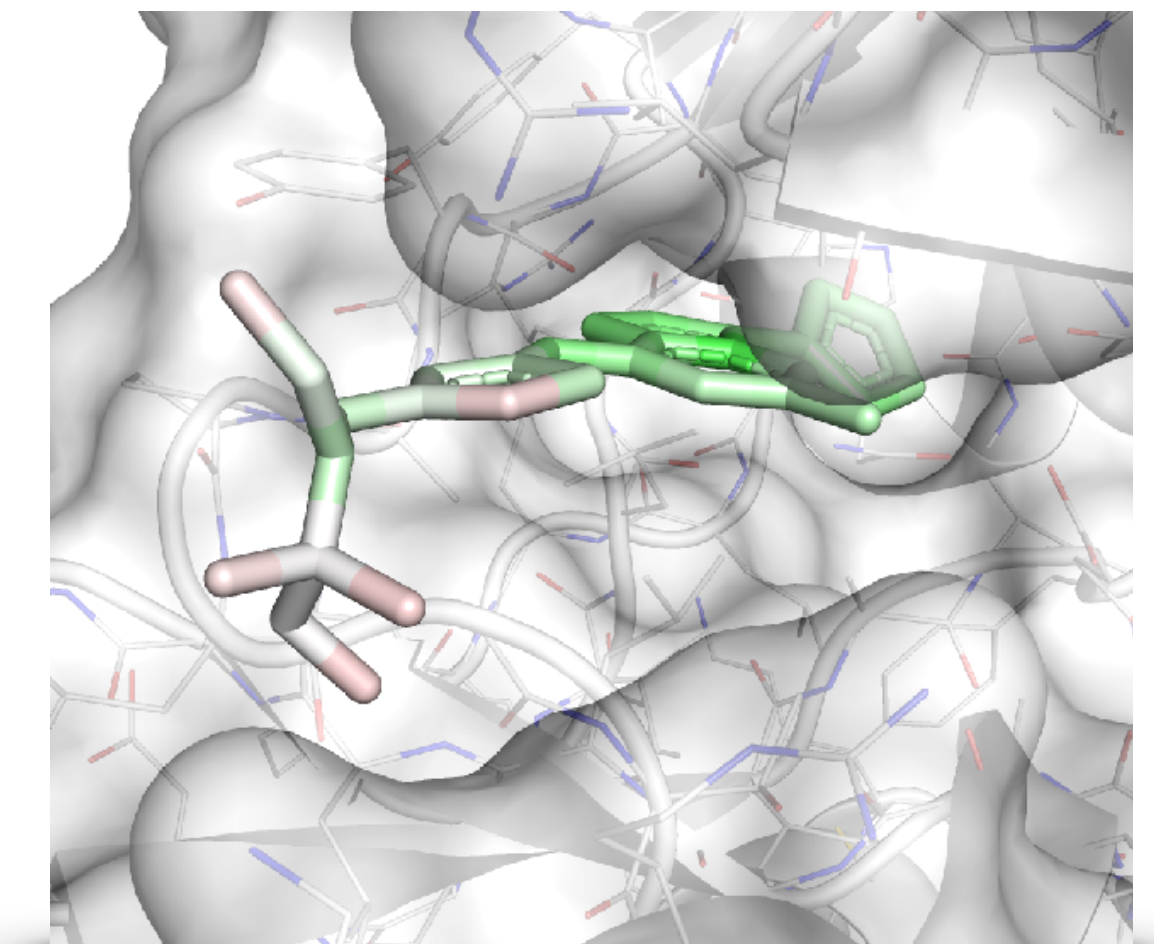**Iterative Training**

Pose *Generation*

**Iterative Training**

*Compound* Generation



**Virtual Screening**
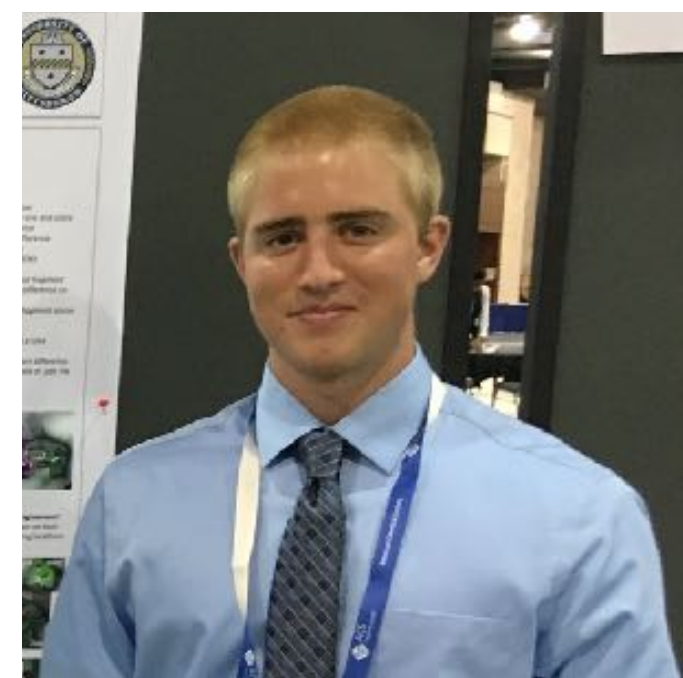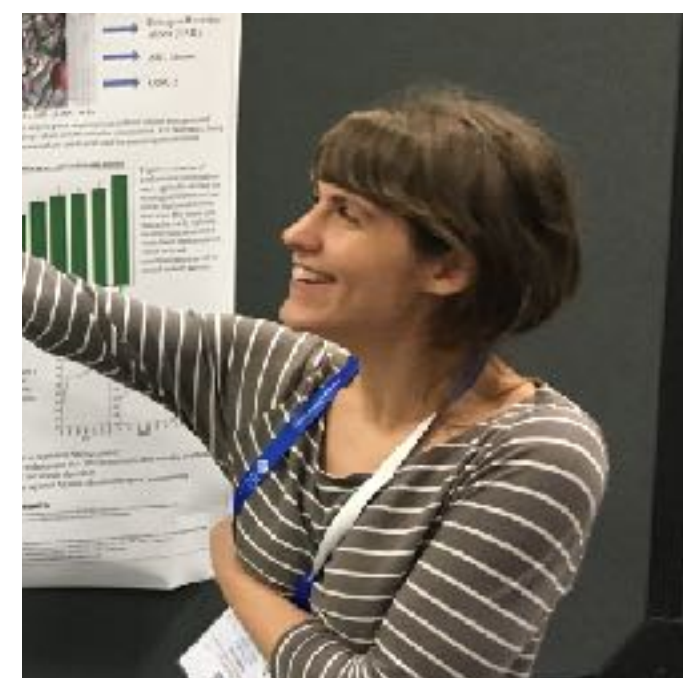


**Lead Optimization**
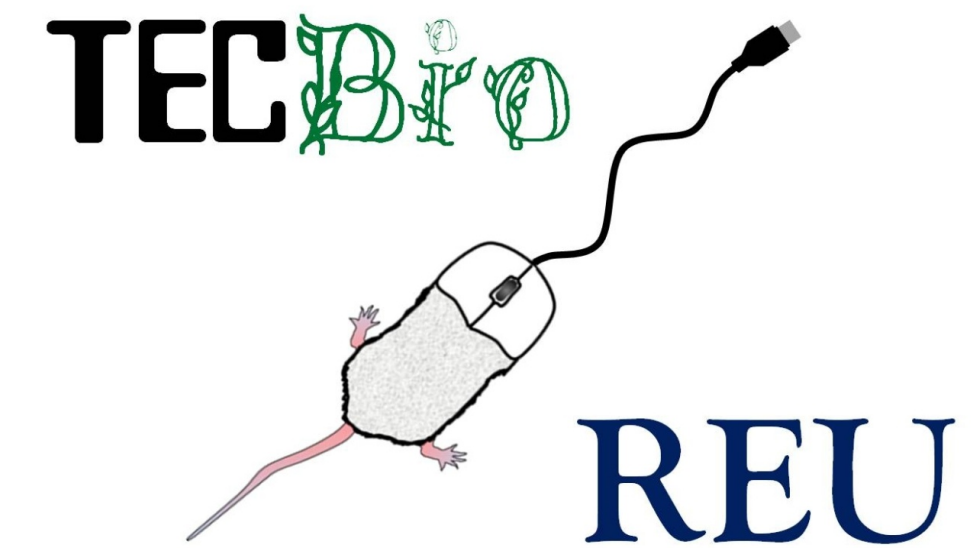
# Acknowledgements


Matt Ragoza


Josh Hochuli


Elisa Idrobo


Jocelyn Sunseri

**Group Members**
Jocelyn Sunseri
Matt Ragoza
Josh Hochuli
Roosha Mandal
Alec Helbling
Lily Turner
Aaron Zheng
Sara Amato
Lily Turner
Aaron Zheng
Gibran Biswas



Department of
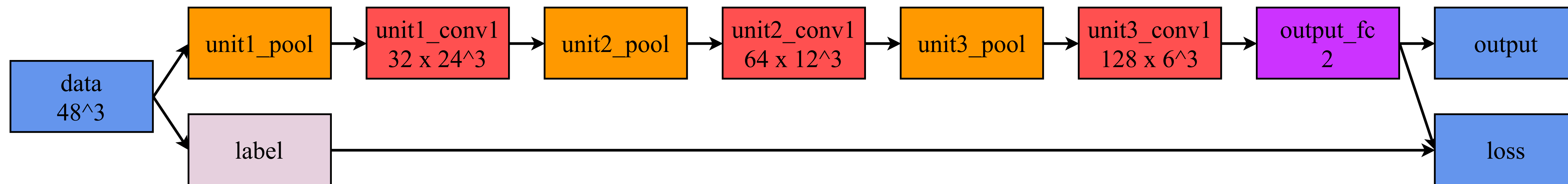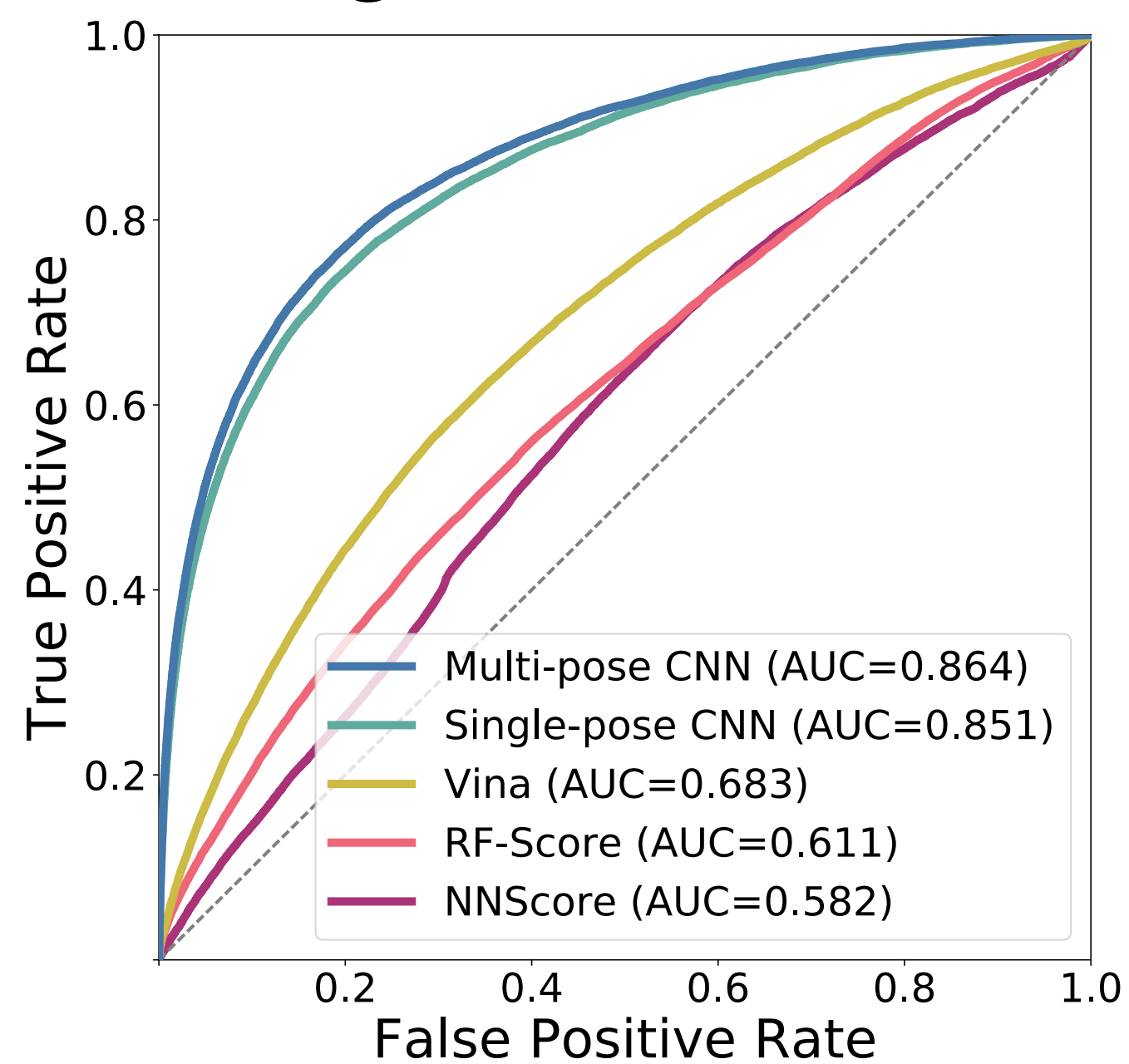Computational and
Systems Biology

# Questions?



## Binding Determination



## Affinity Prediction



R=0.687  RMS=2.186

## Relevance Propagation

# Questions?



data
48^3

unit1_pool

unit1_conv1
32 x 24^3

unit2_pool

unit2_conv1
64 x 12^3

unit3_pool

unit3_conv1
128 x 6^3

output_fc
2

output

label

loss

## Binding Determination



Multi-pose CNN (AUC=0.864)
Single-pose CNN (AUC=0.851)
Vina (AUC=0.683)
RF-Score (AUC=0.611)
NNScore (AUC=0.582)

## Affinity Prediction



R=0.687  RMS=2.186

## Relevance Propagation



33