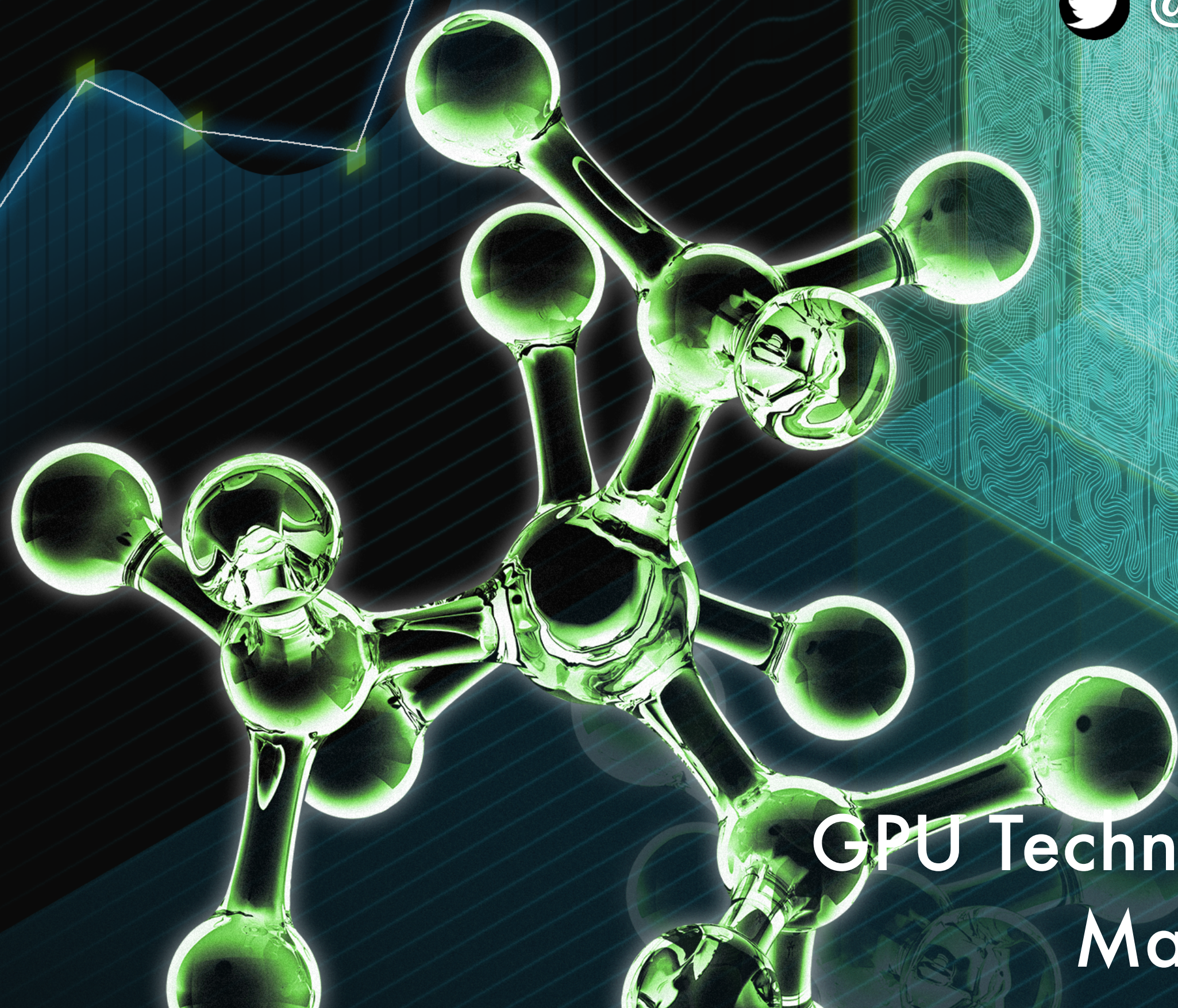# GPU-Accelerated Convolutional Neural Networks For Protein-Ligand Scoring

## David Koes

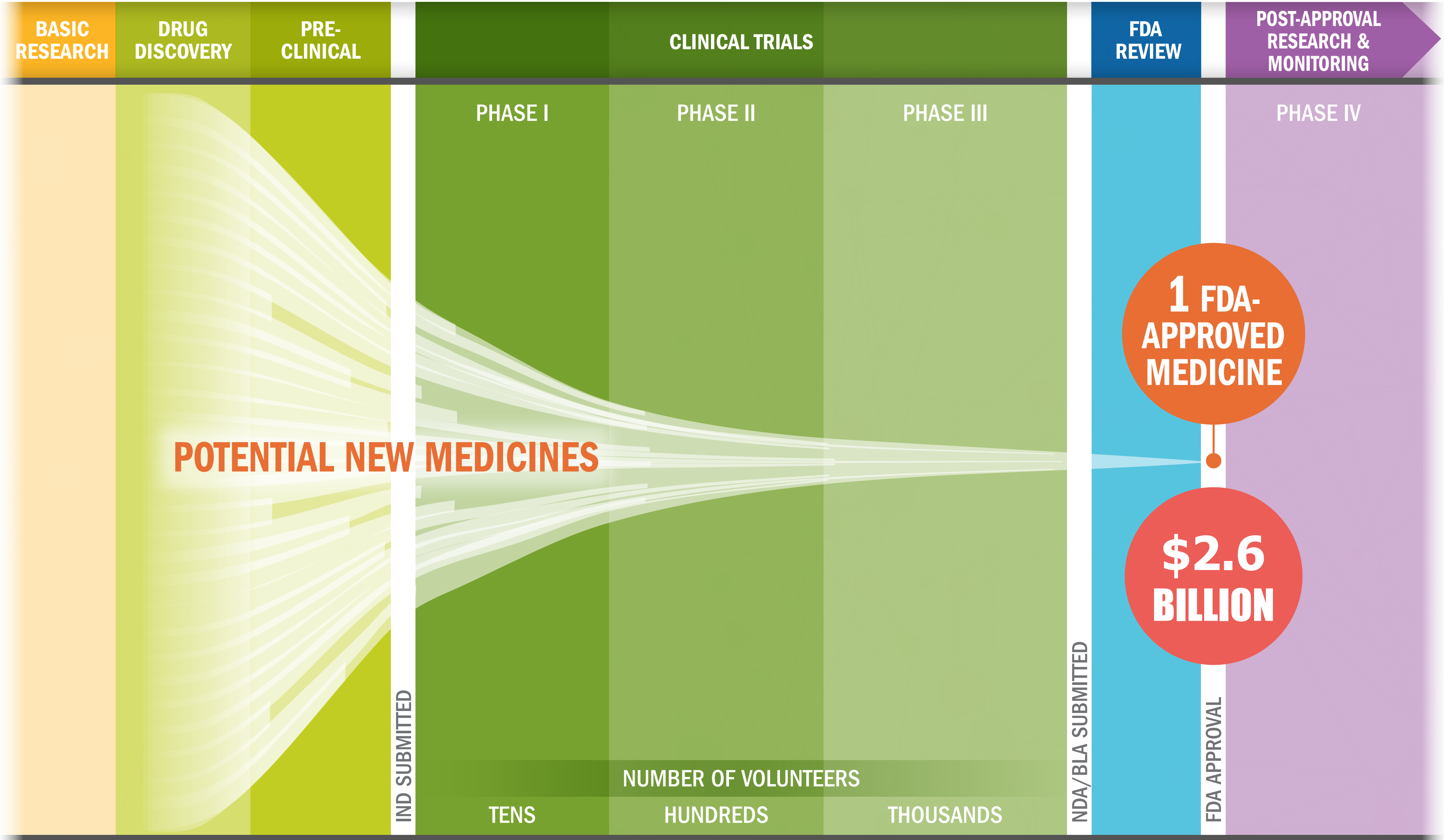@david_koes

GPU Technology Conference
May 8, 2017
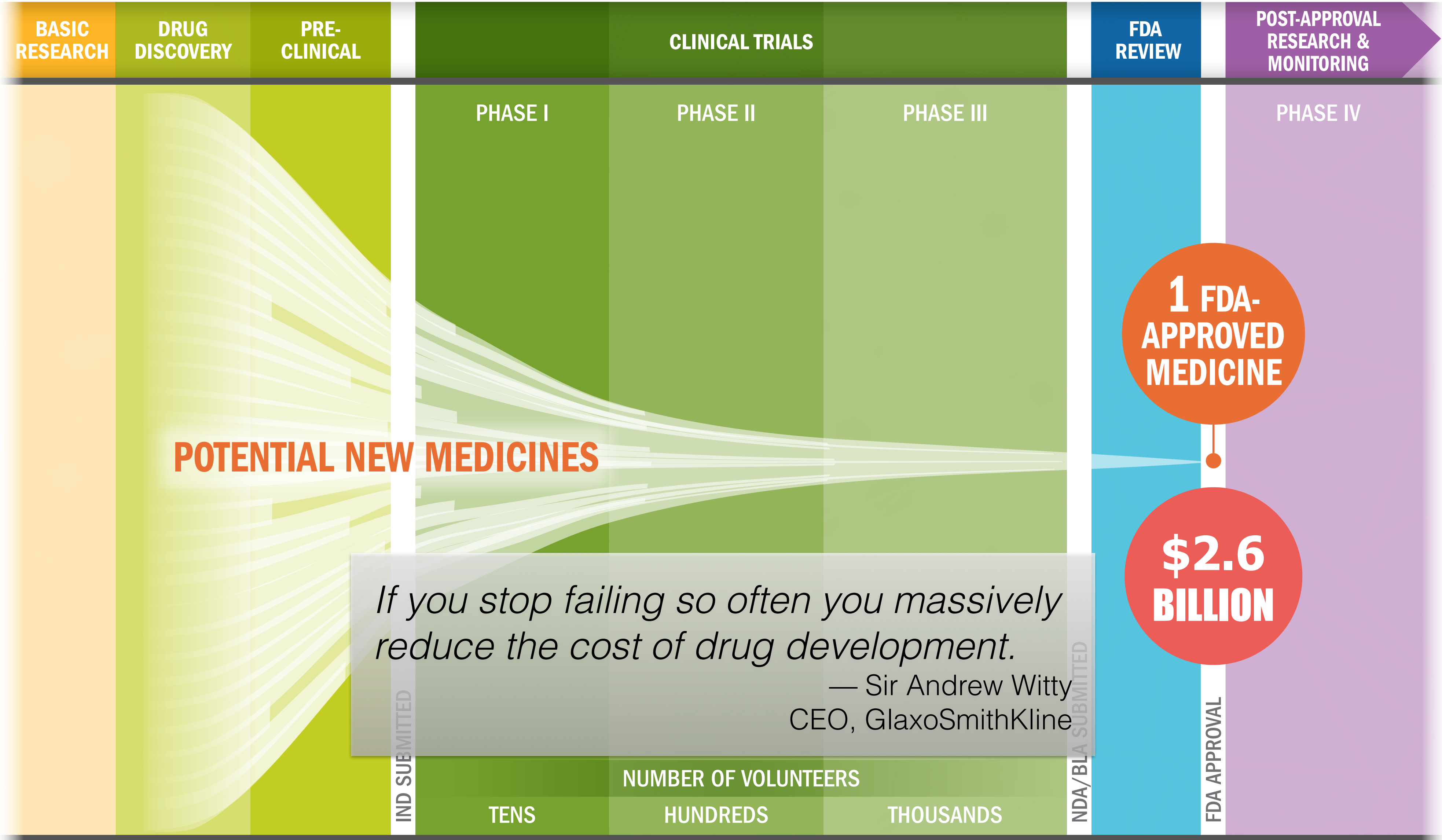
# THE BIOPHARMACEUTICAL RESEARCH AND DEVELOPMENT PROCESS



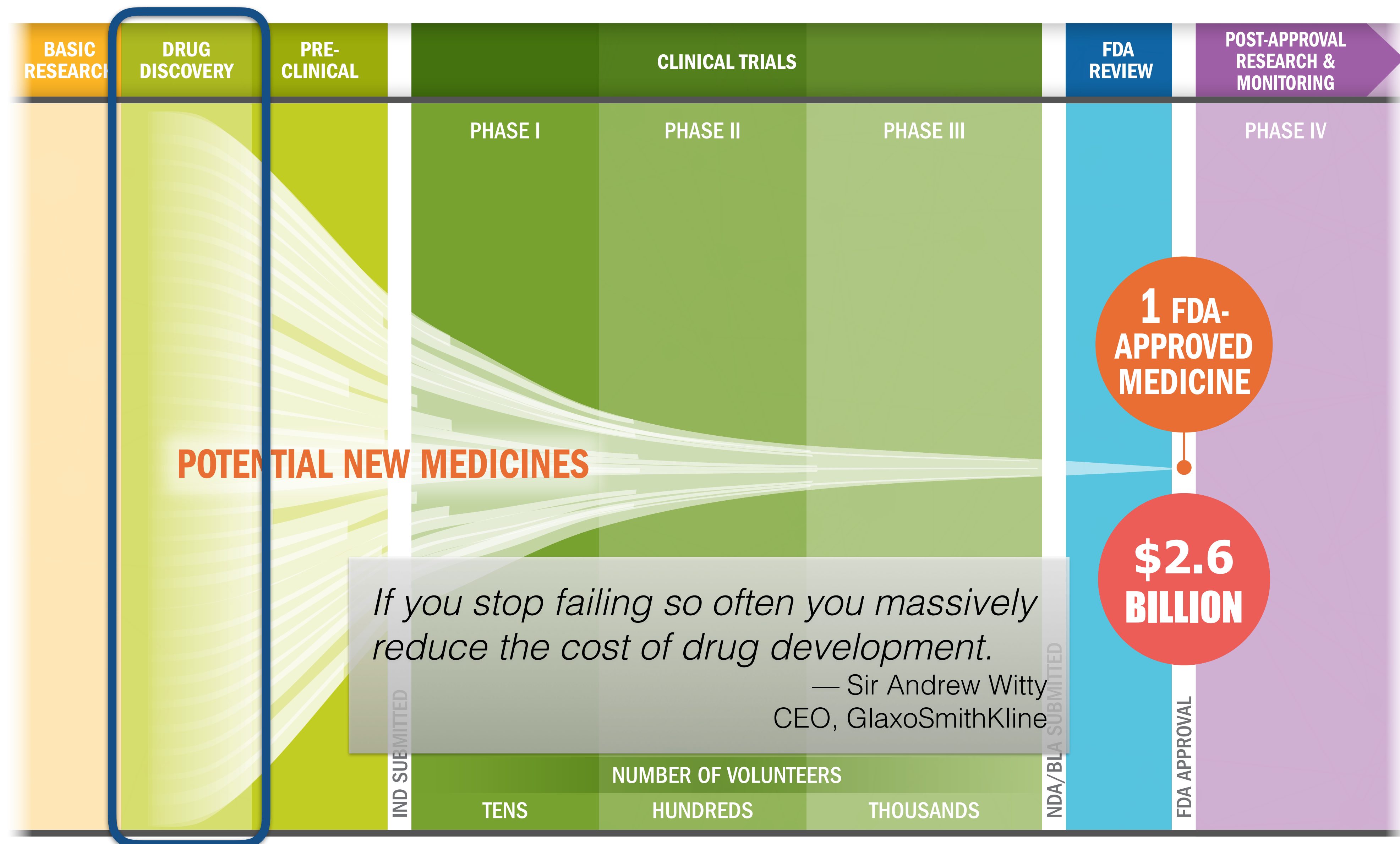Source: Pharmaceutical Research and Manufacturers of America (http://phrma.org)
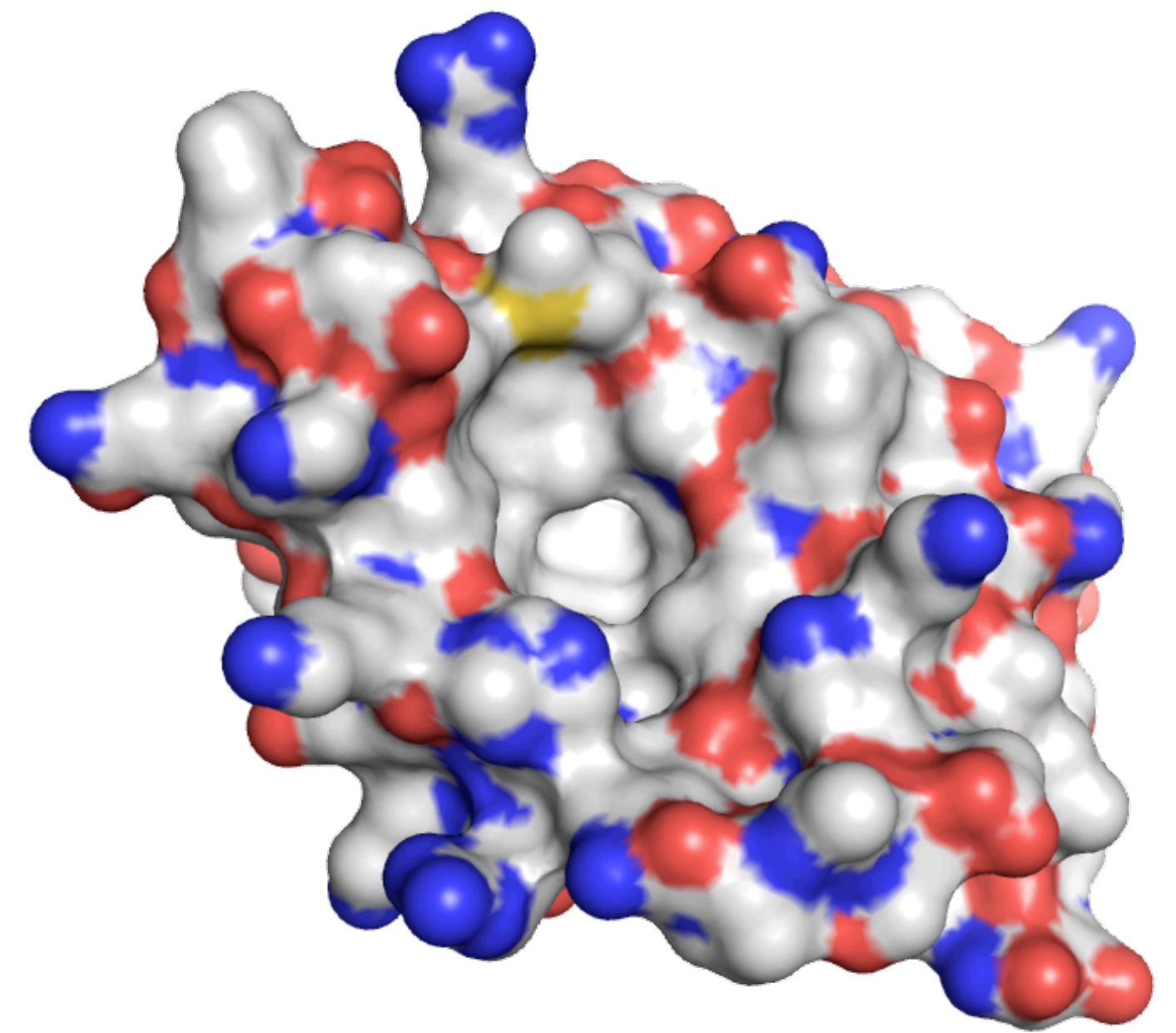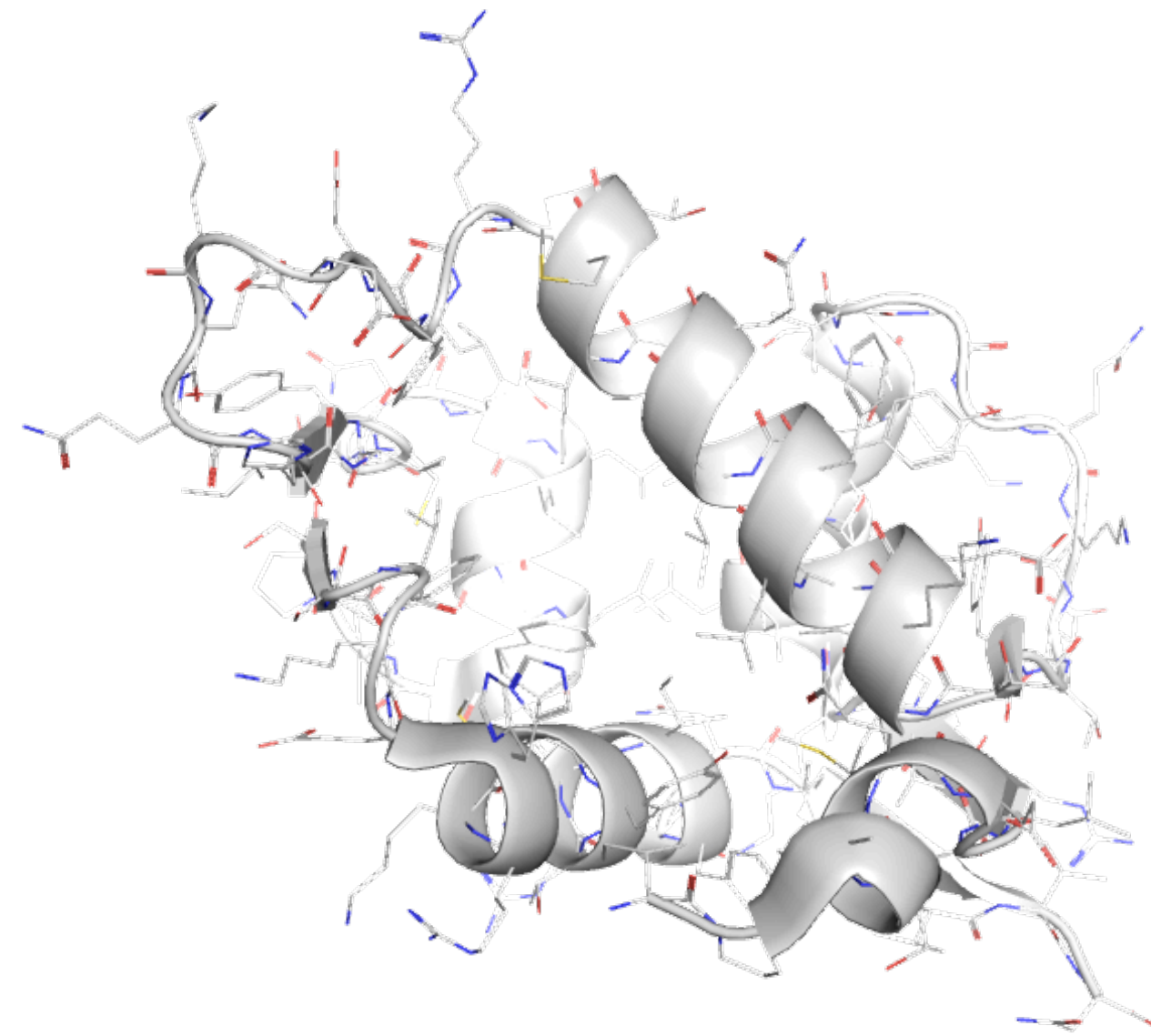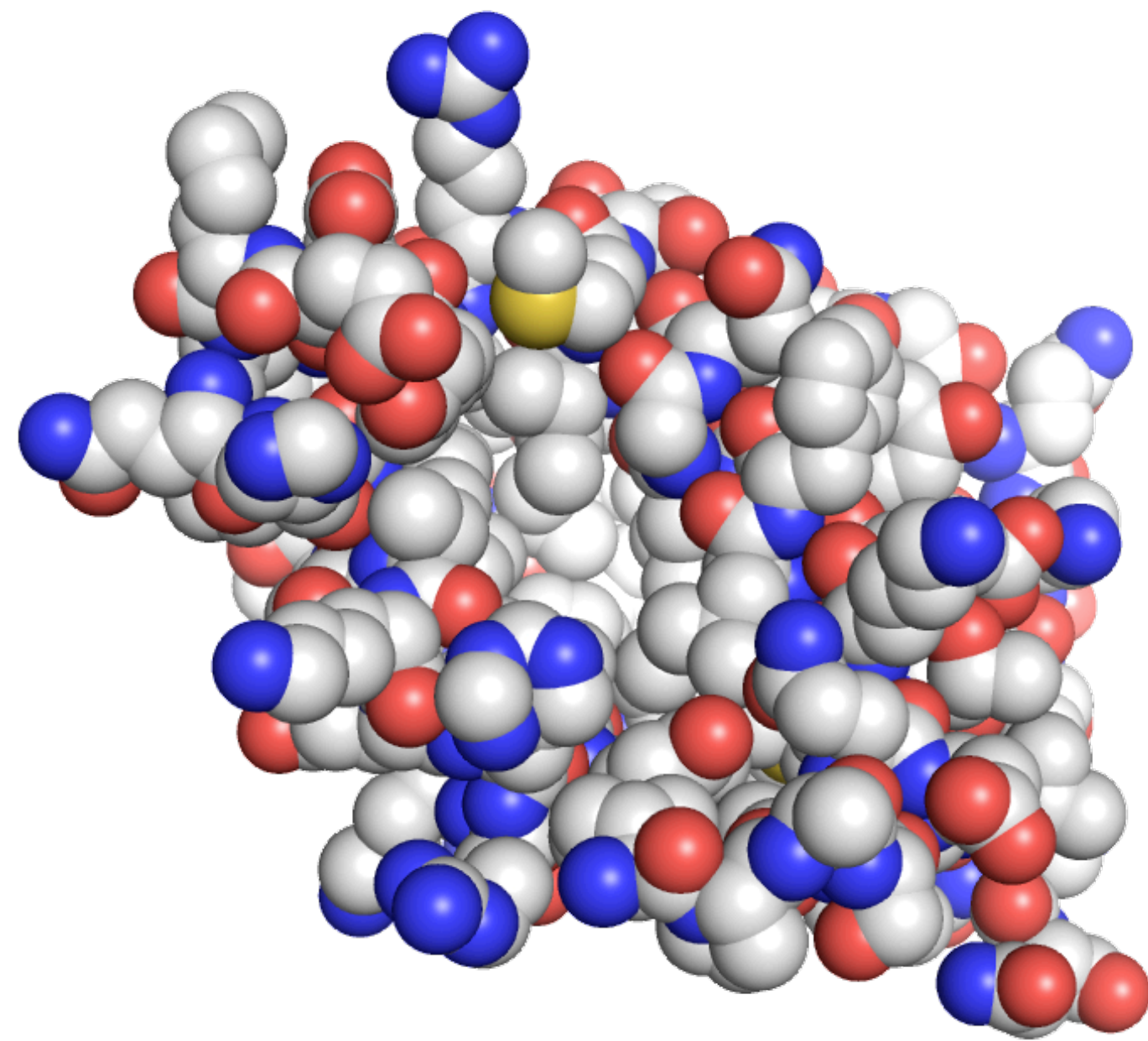
# THE BIOPHARMACEUTICAL RESEARCH AND DEVELOPMENT PROCESS



BASIC RESEARCH

DRUG DISCOVERY

PRE-CLINICAL

CLINICAL TRIALS

FDA REVIEW

POST-APPROVAL RESEARCH & MONITORING

PHASE I

PHASE II

PHASE III

PHASE IV

1 FDA-APPROVED MEDICINE

POTENTIAL NEW MEDICINES

$2.6 BILLION

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

IND SUBMITTED

NDA/BLA SUBMITTED

FDA APPROVAL

NUMBER OF VOLUNTEERS

TENS

HUNDREDS

THOUSANDS

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

# 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**

**1** FDA-APPROVED MEDICINE

**$2.6 BILLION**

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

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?

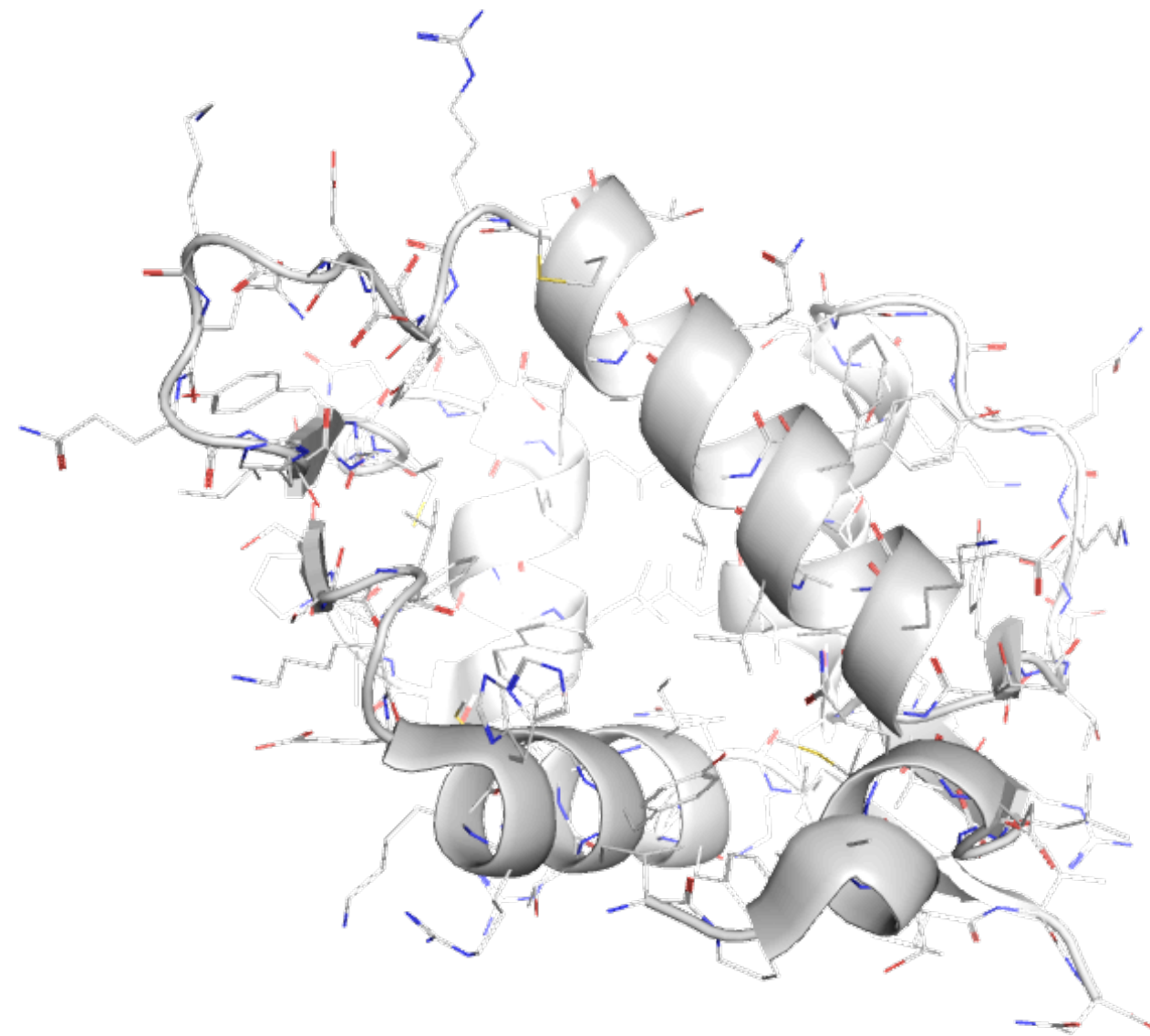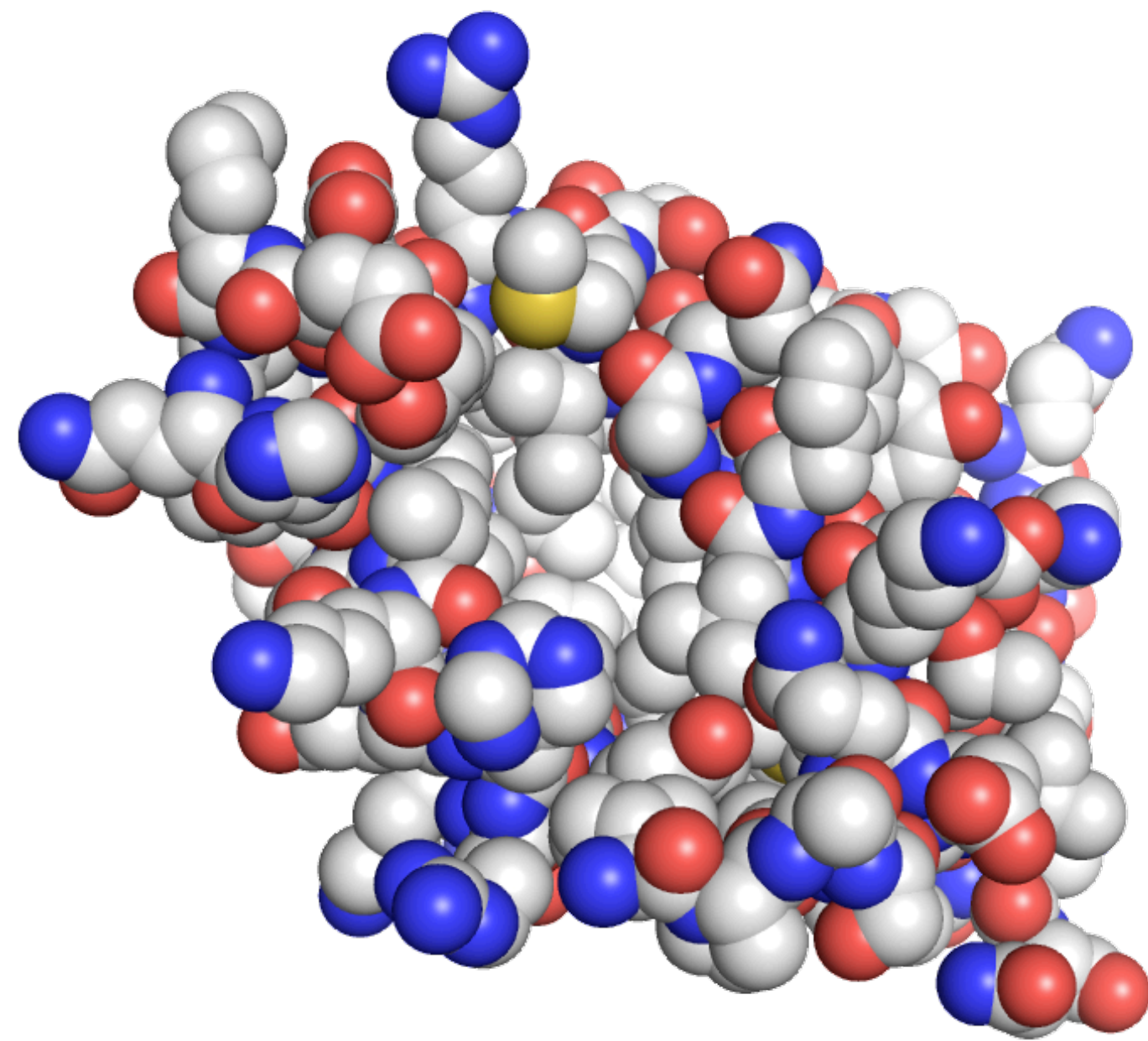3. Is what you want it to do the right thing?

# Protein Structures

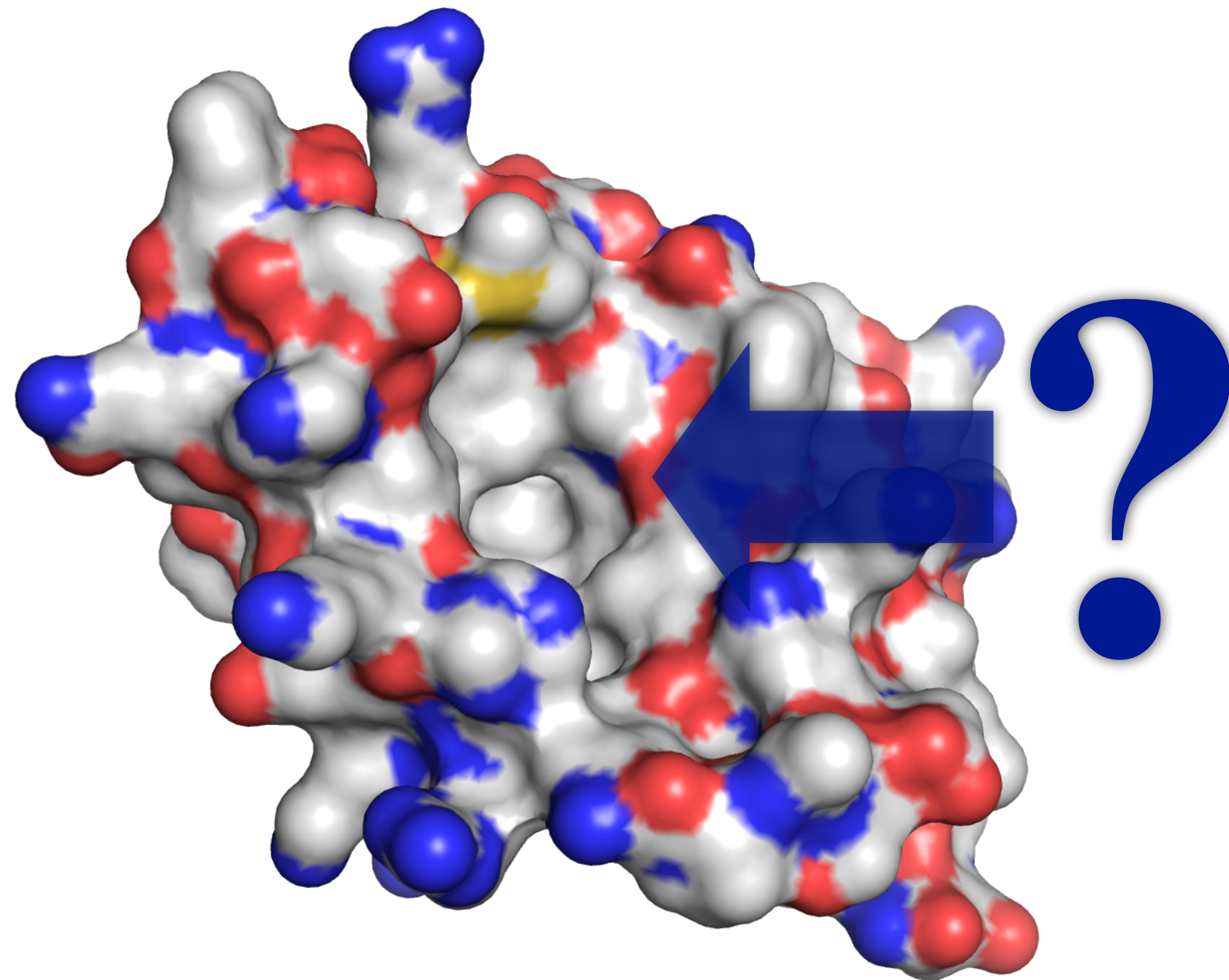sequence → **structure** → function

# Protein Structures

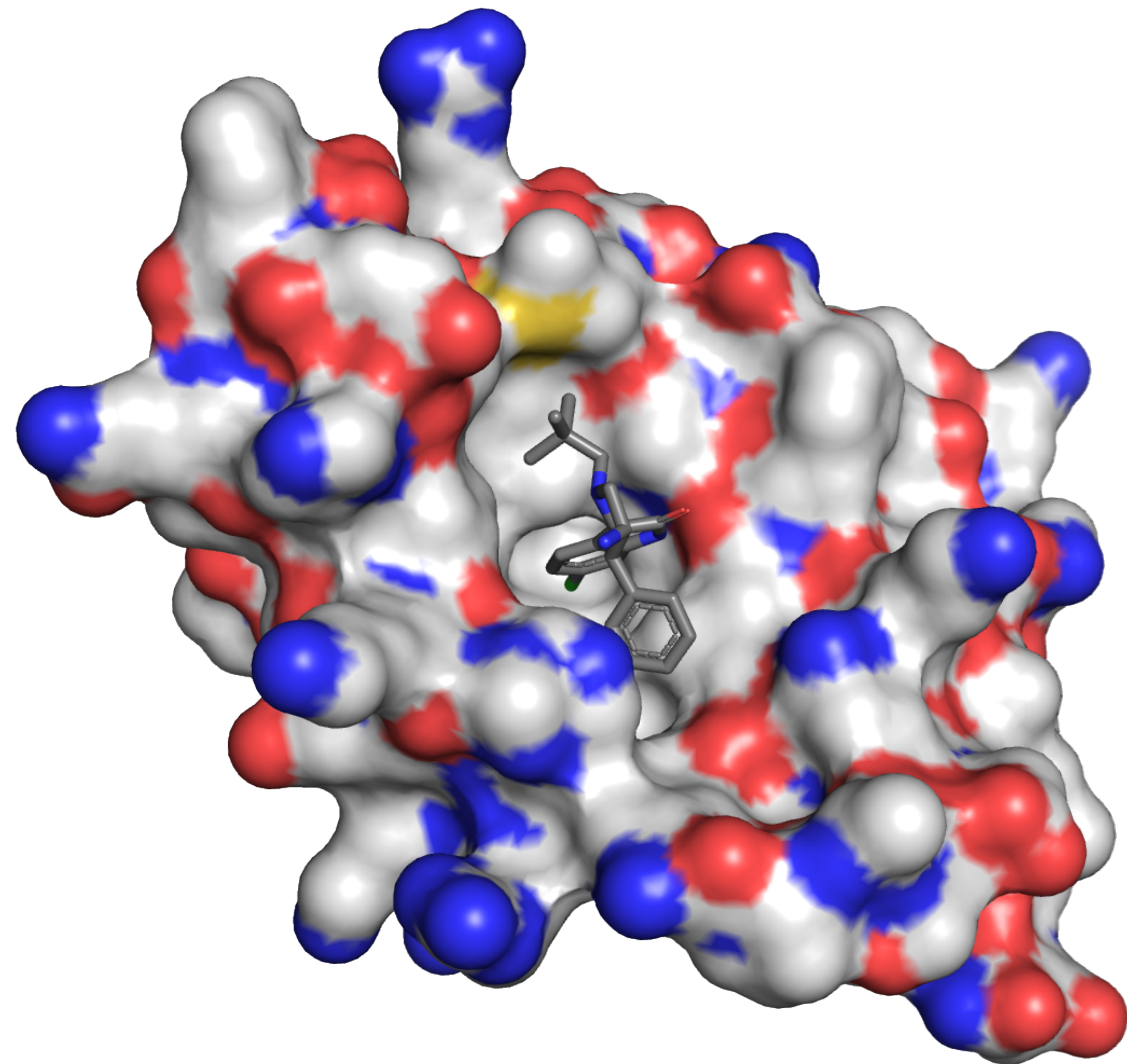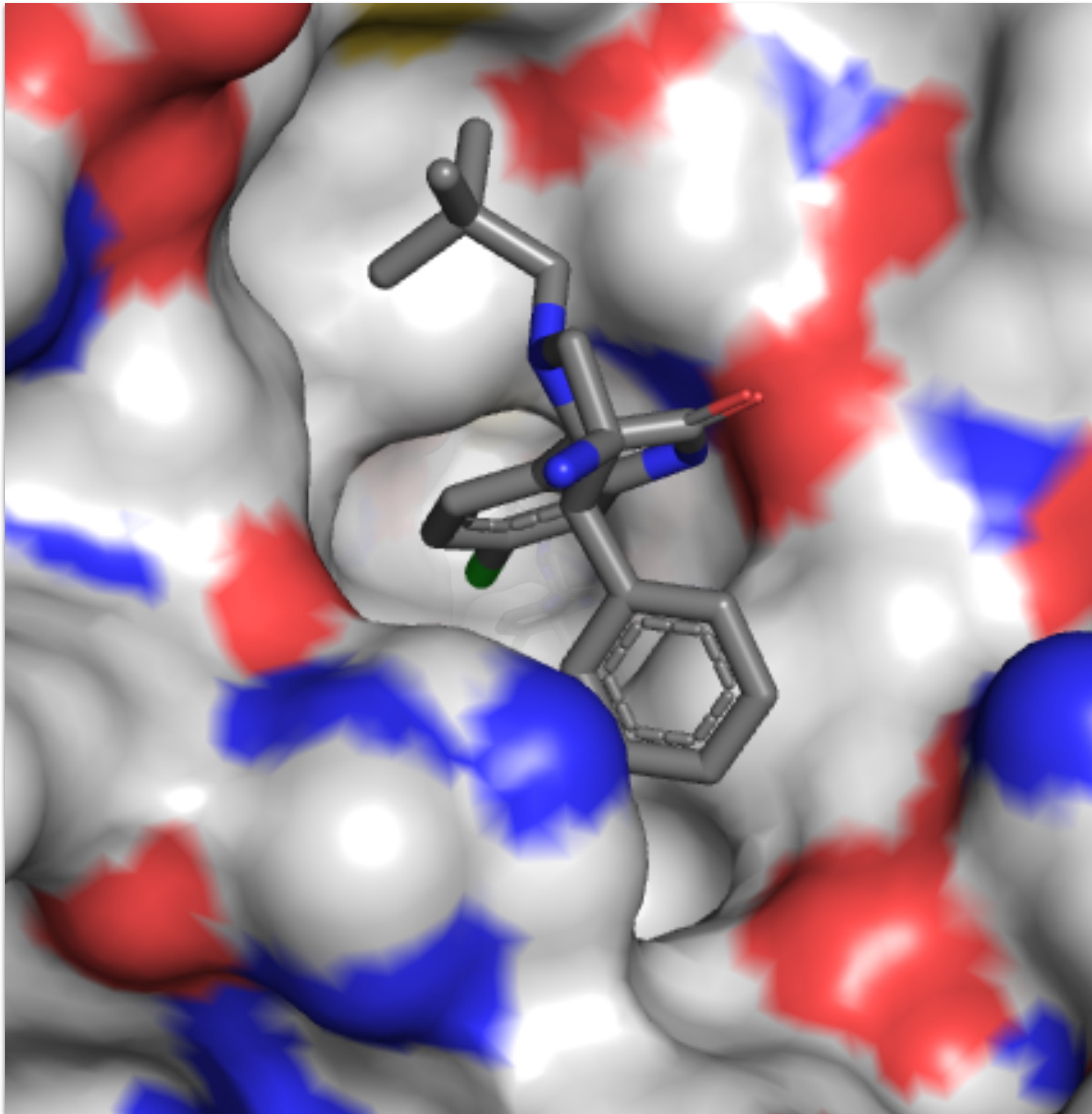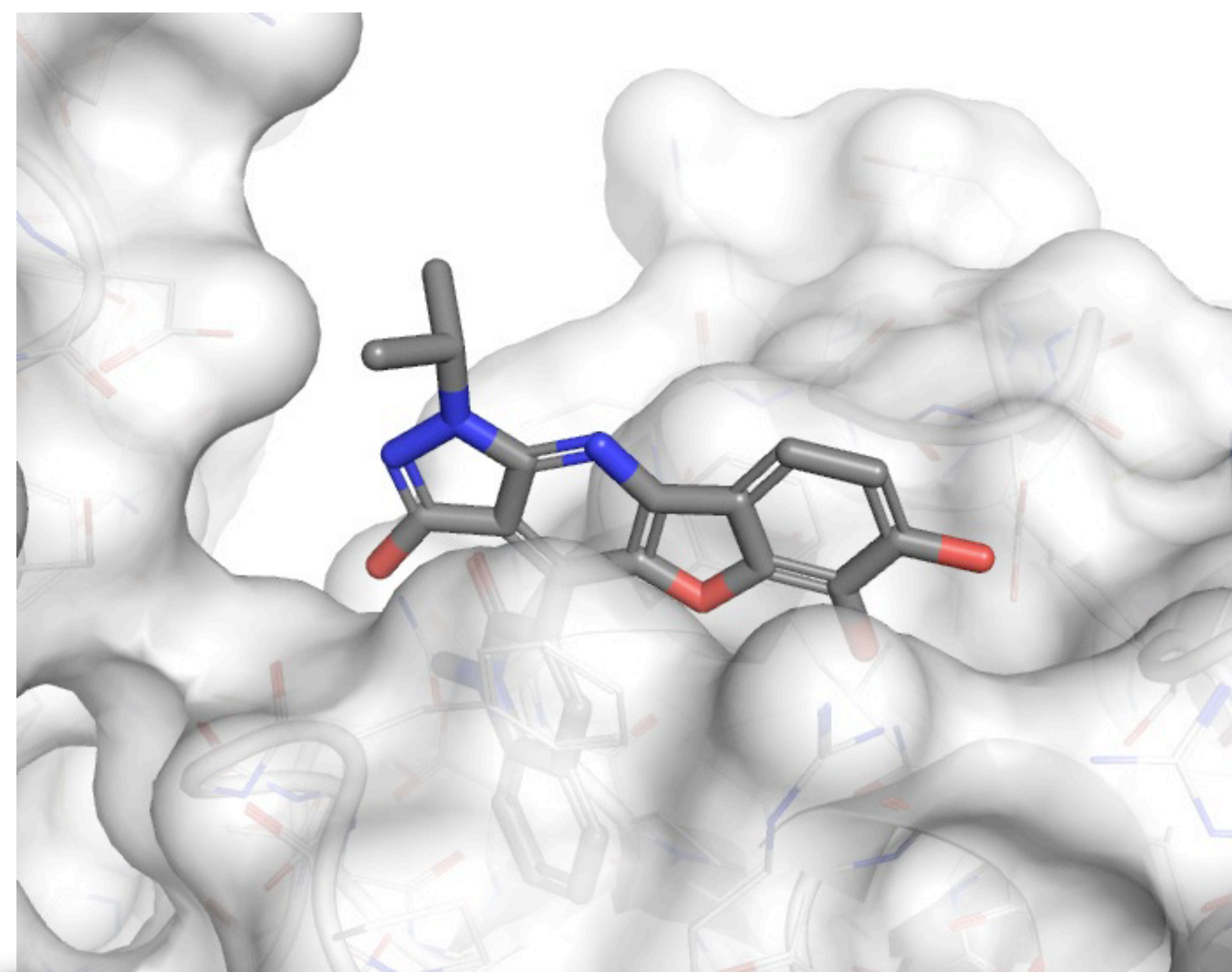sequence → **structure** → function

# Structure Based Drug Design



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

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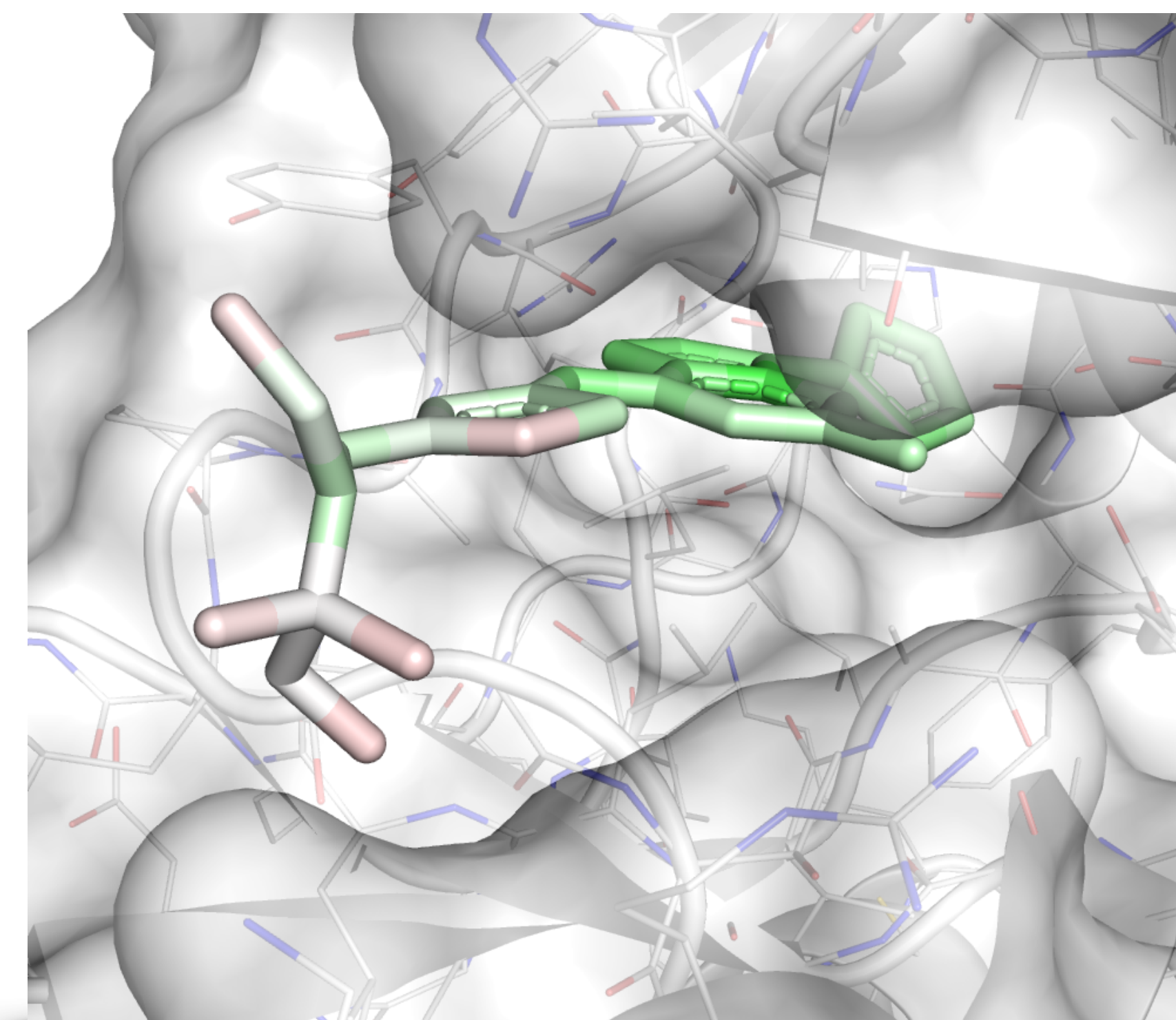**known structure** and **binding site**

# 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**



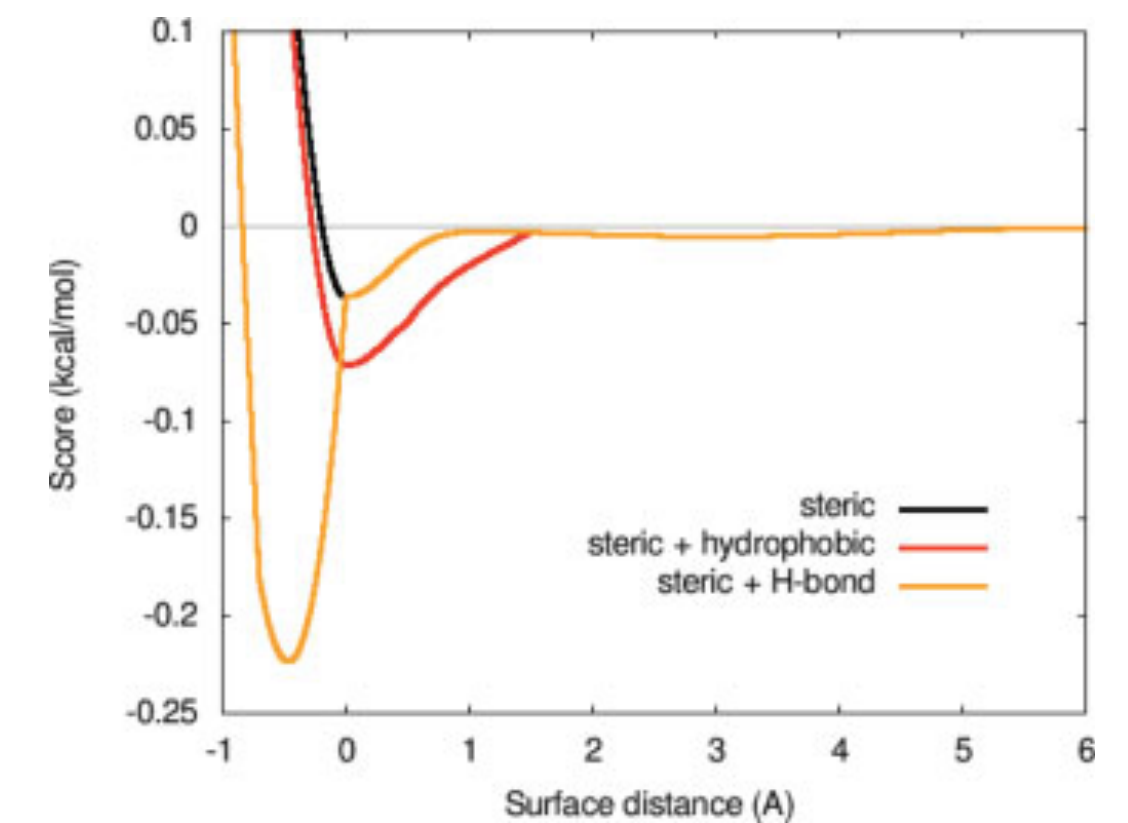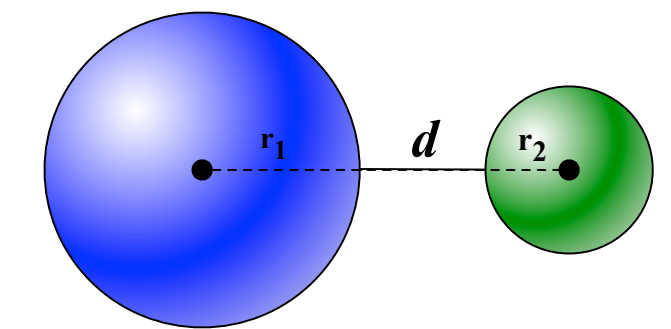Pose Prediction          Binding Discrimination          Affinity Prediction

# Structure Based Drug Design

**Virtual Screening**

**Lead Optimization**



Pose Prediction     Binding Discrimination     Affinity Prediction

# Protein-Ligand Scoring

## AutoDock Vina



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

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

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

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

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

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"

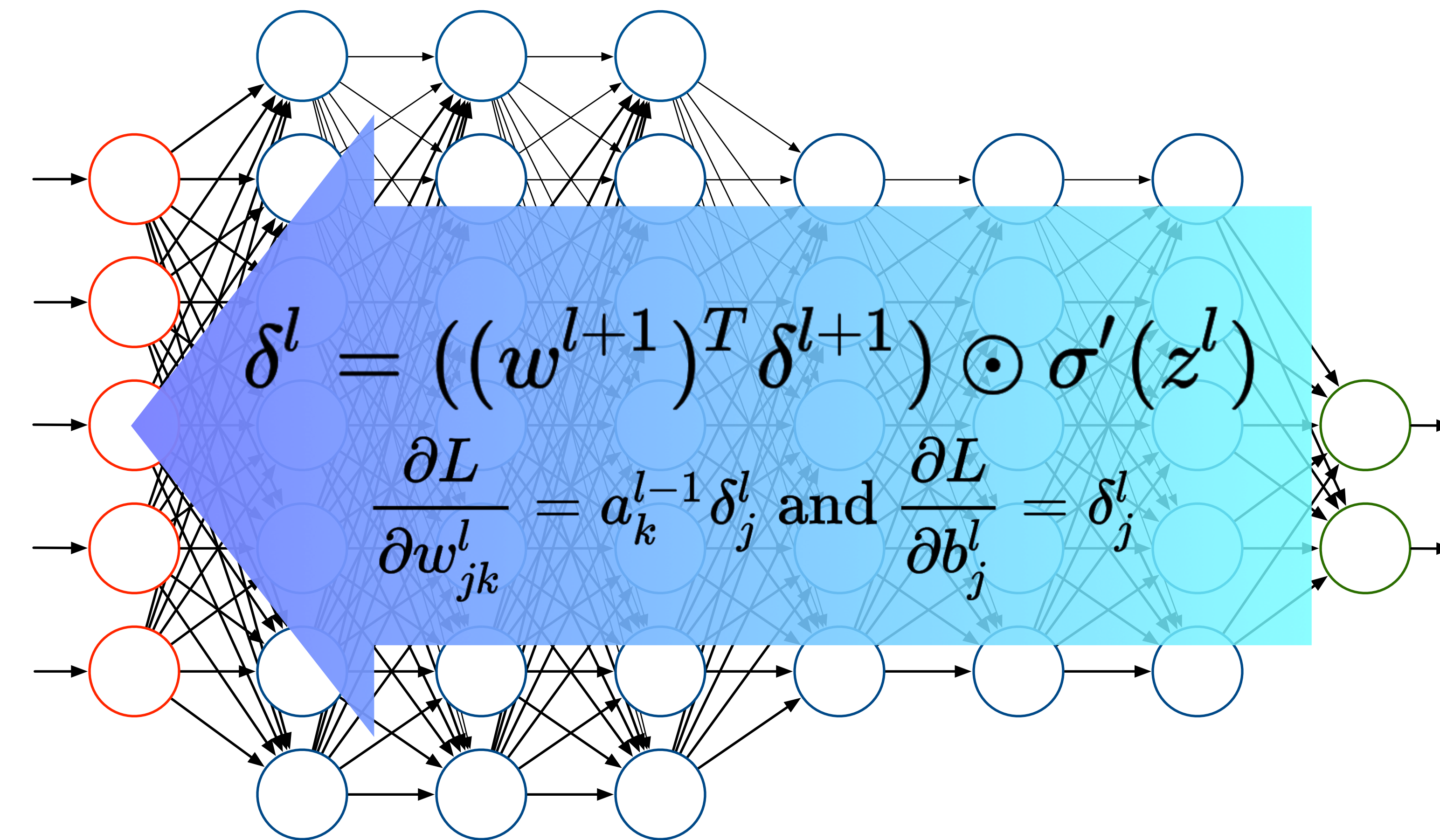- 231,655,275 bioactivities in PubChem
- 125,526 structures in the PDB
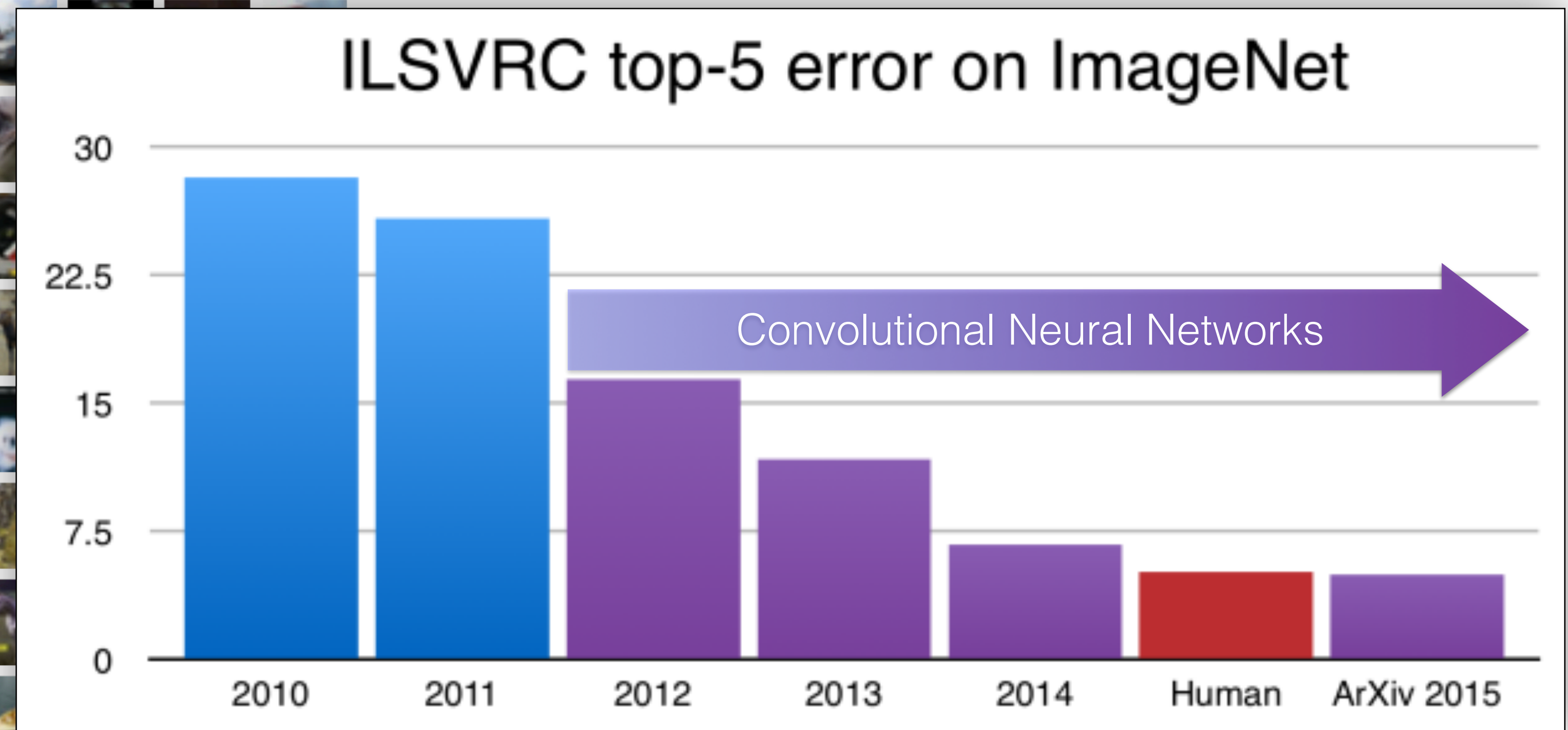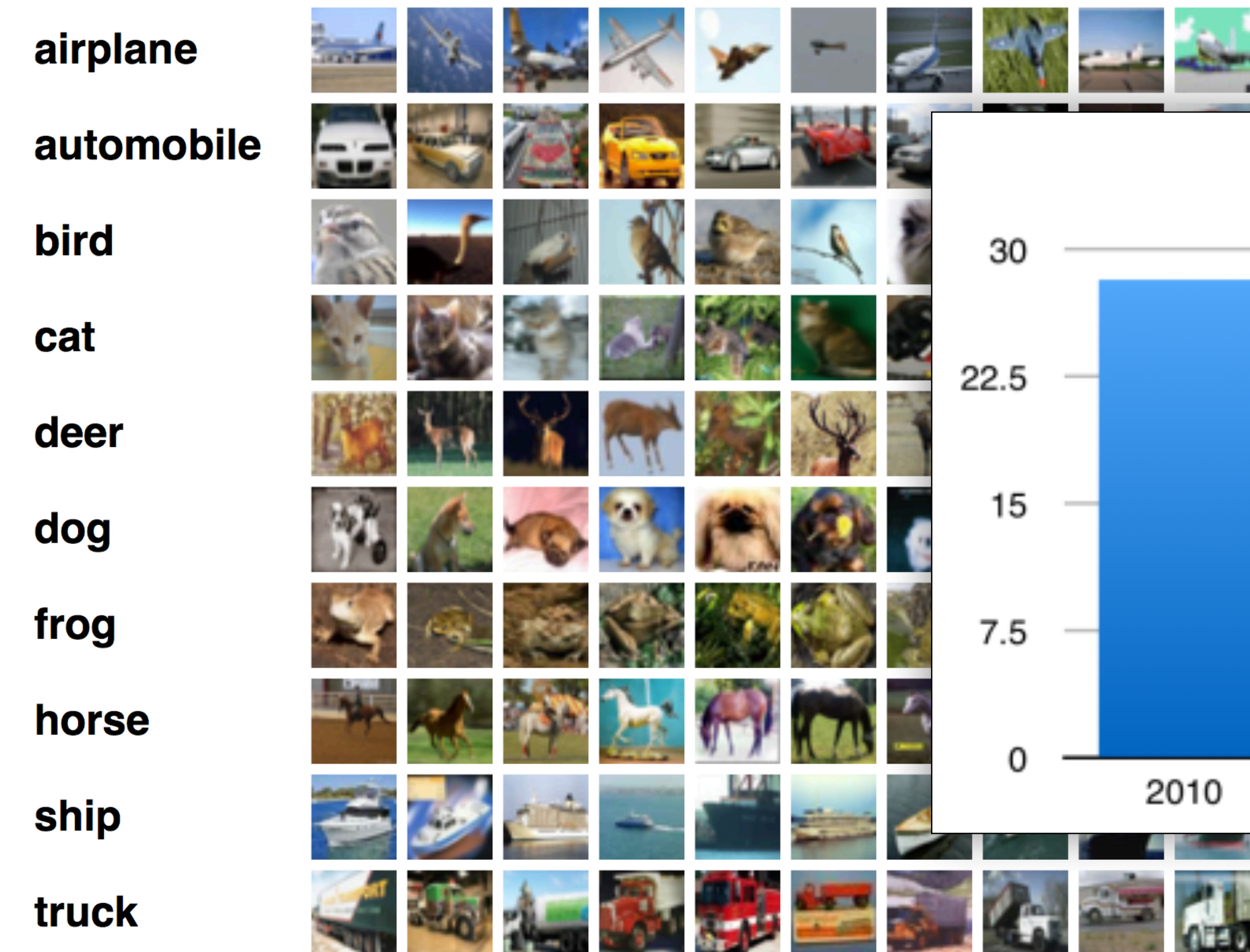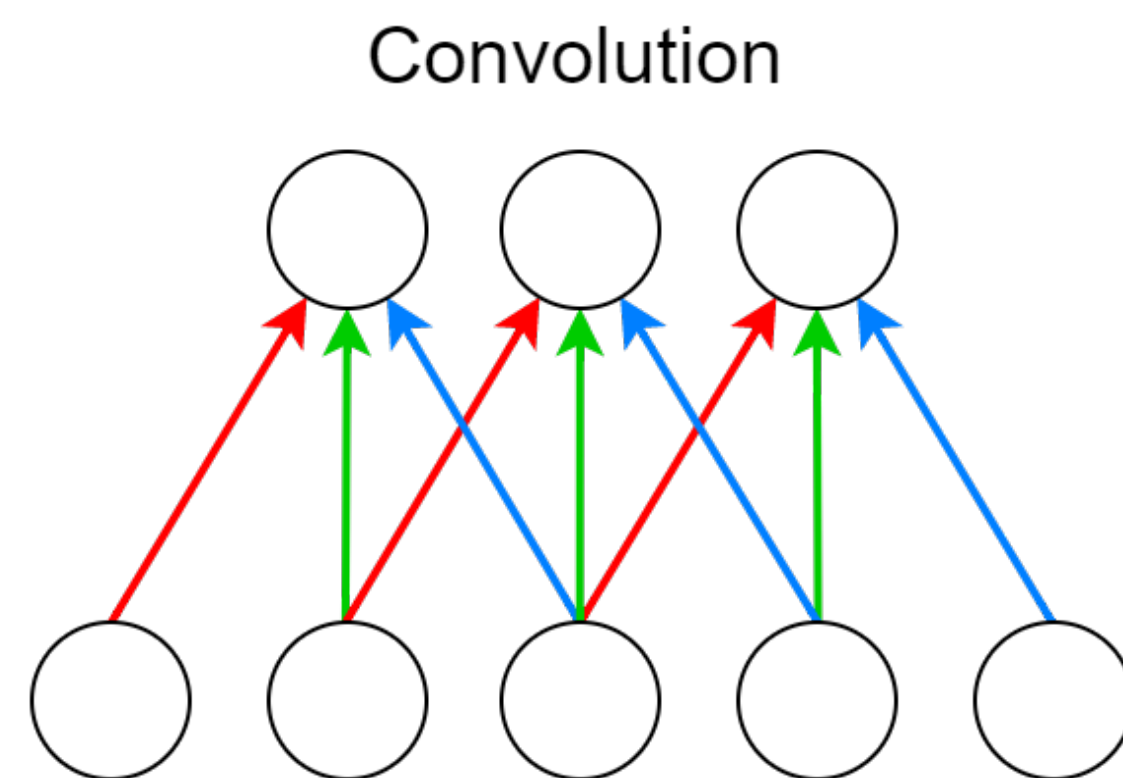- 16,179 annotated complexes in PDBbind

8

# Deep Learning

# Deep Learning



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

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

# 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

14

# 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
- <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 | | | |

# 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



20

# Model Optimization

# Model Optimization

# Cross-Validation Evaluation

# Pose Prediction (CSAR)

# Pose Prediction (CSAR)
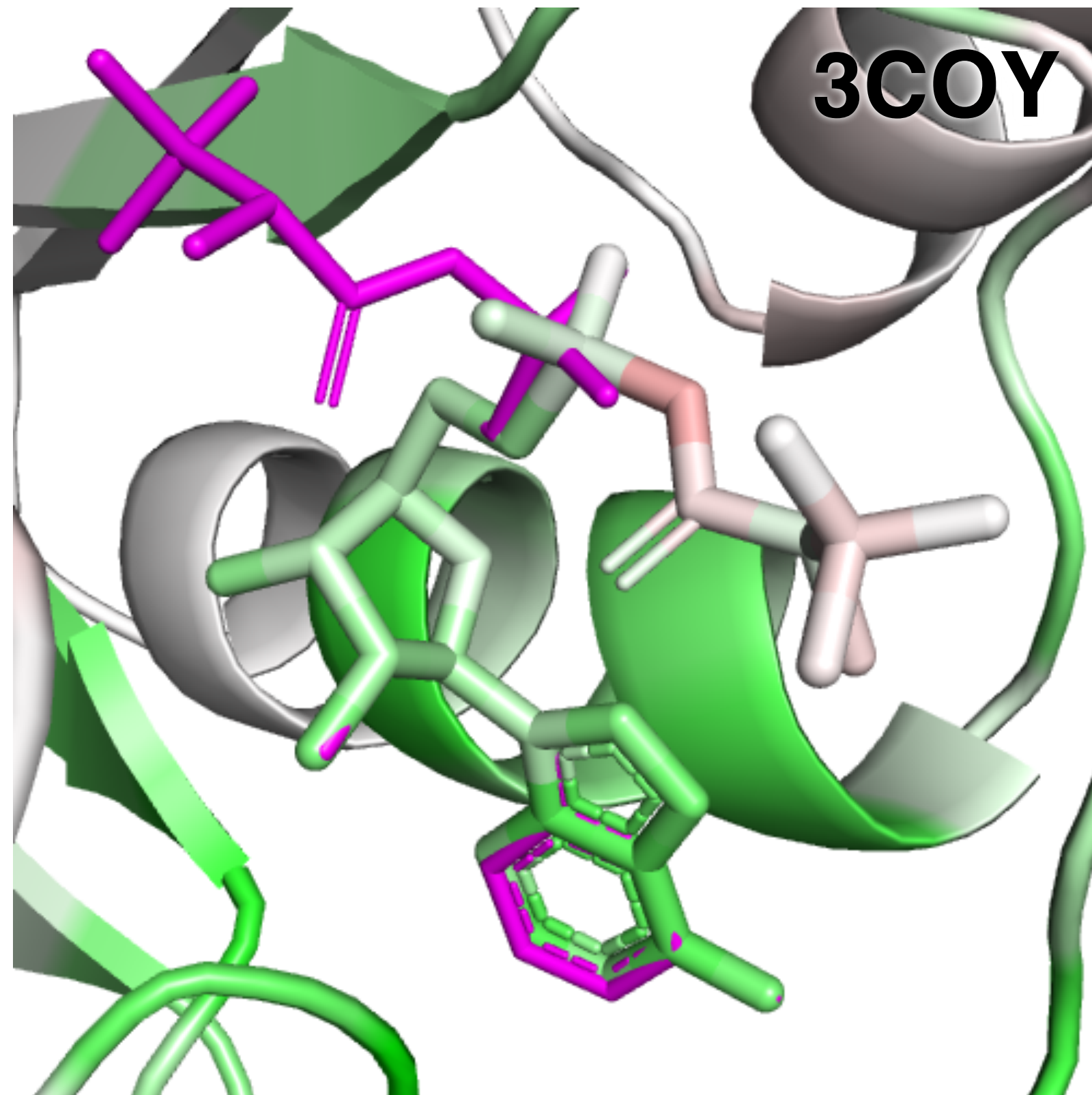


*inter*-target ranking

*intra*-target ranking

23

# Pose Prediction (PDBbind)

# Pose Prediction (PDBbind)



*inter*-target ranking

*intra*-target ranking

# Visualization



Delete single ligand atoms

Delete ligand fragments

Score

Average

Delete single residues

# Examples



Partially Aligned Poses

# Beyond Scoring

# Beyond Scoring



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

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

# Beyond Scoring

# Beyond Scoring

# Beyond Scoring

**2Q89**

Less Oxygen Here

More Oxygen Here

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

# The Future

Pose Selection

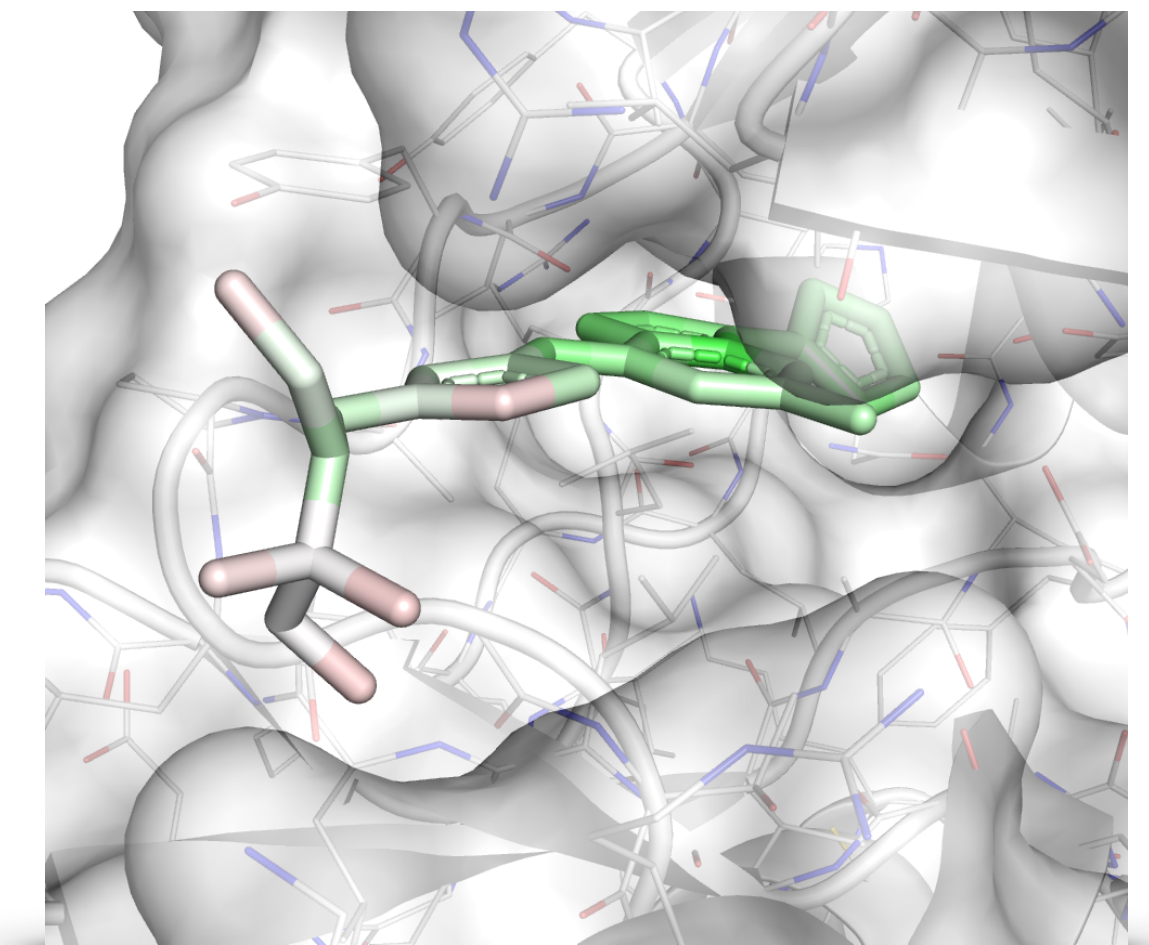**Iterative Training**

Pose *Generation*

**Iterative Training**

*Compound* Generation





**Virtual Screening**



**Lead Optimization**

# The Future

Pose Selection → **Iterative Training** → Pose *Generation* → **Iterative Training** → *Compound* Generation
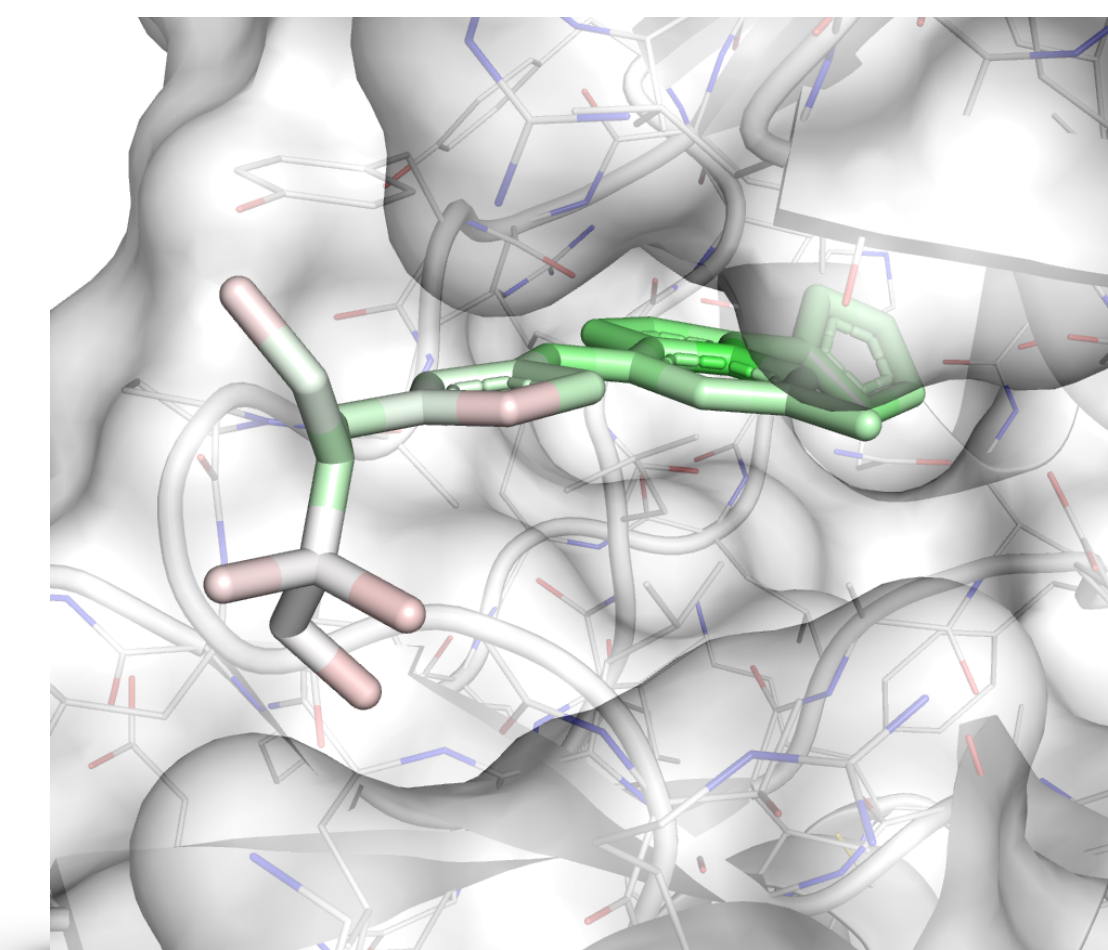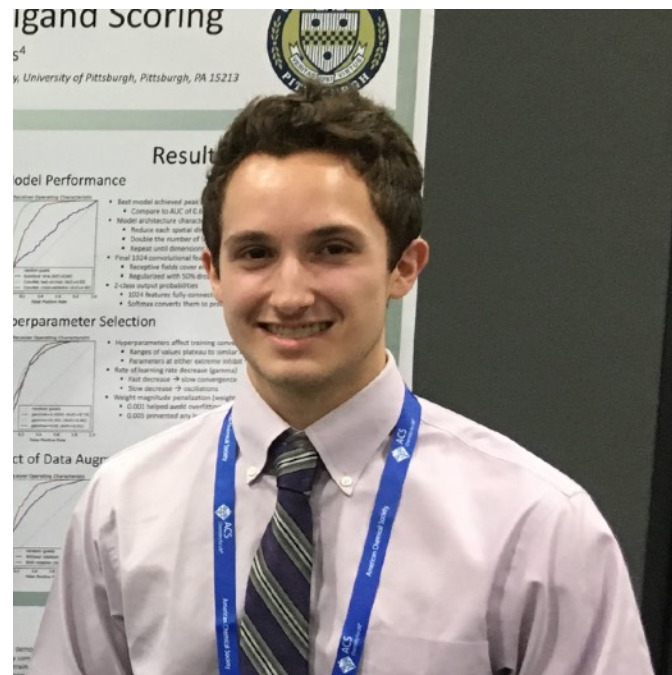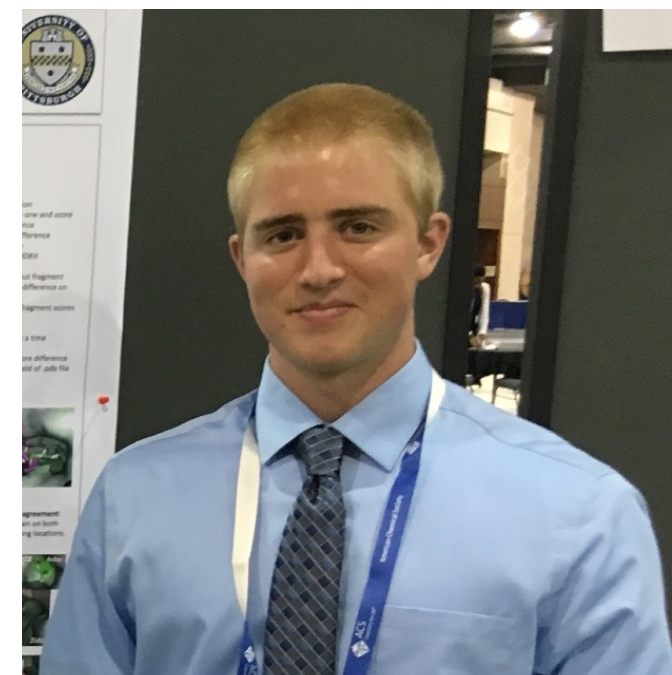


**Virtual Screening**

**Lead Optimization**

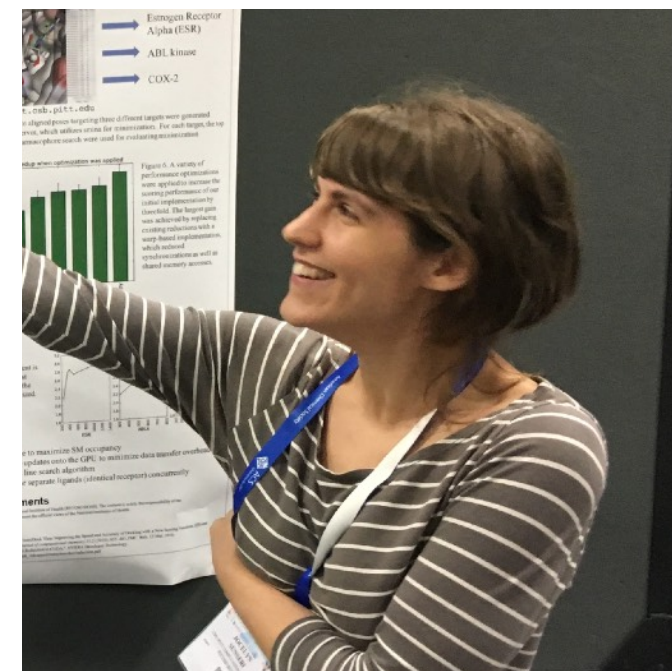# Acknowledgements



Matt Ragoza



Josh Hochuli



Elisa Idrobo



Jocelyn Sunseri

**Group Members**

Jocelyn Sunseri

Matt Ragoza

Josh Hochuli

Roosha Mandal

Alec Helbling

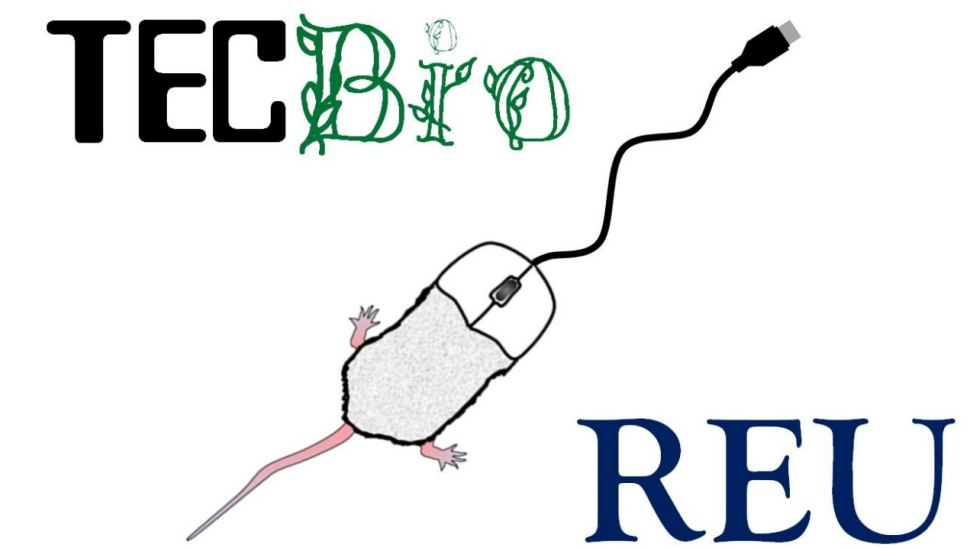Lily Turner

Aaron Zheng

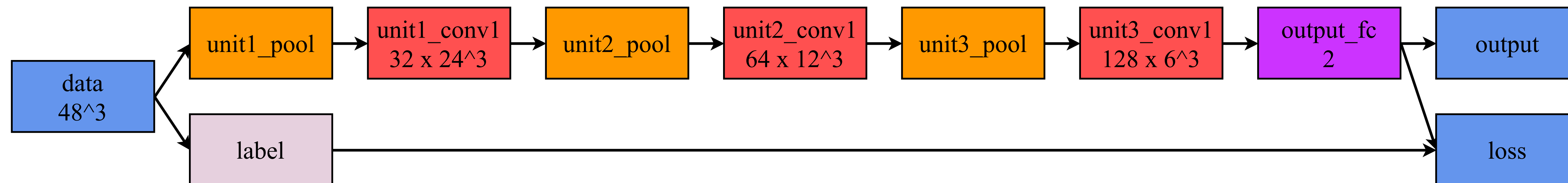Sara Amato

Lily Turner

Aaron Zheng

Gibran Biswas



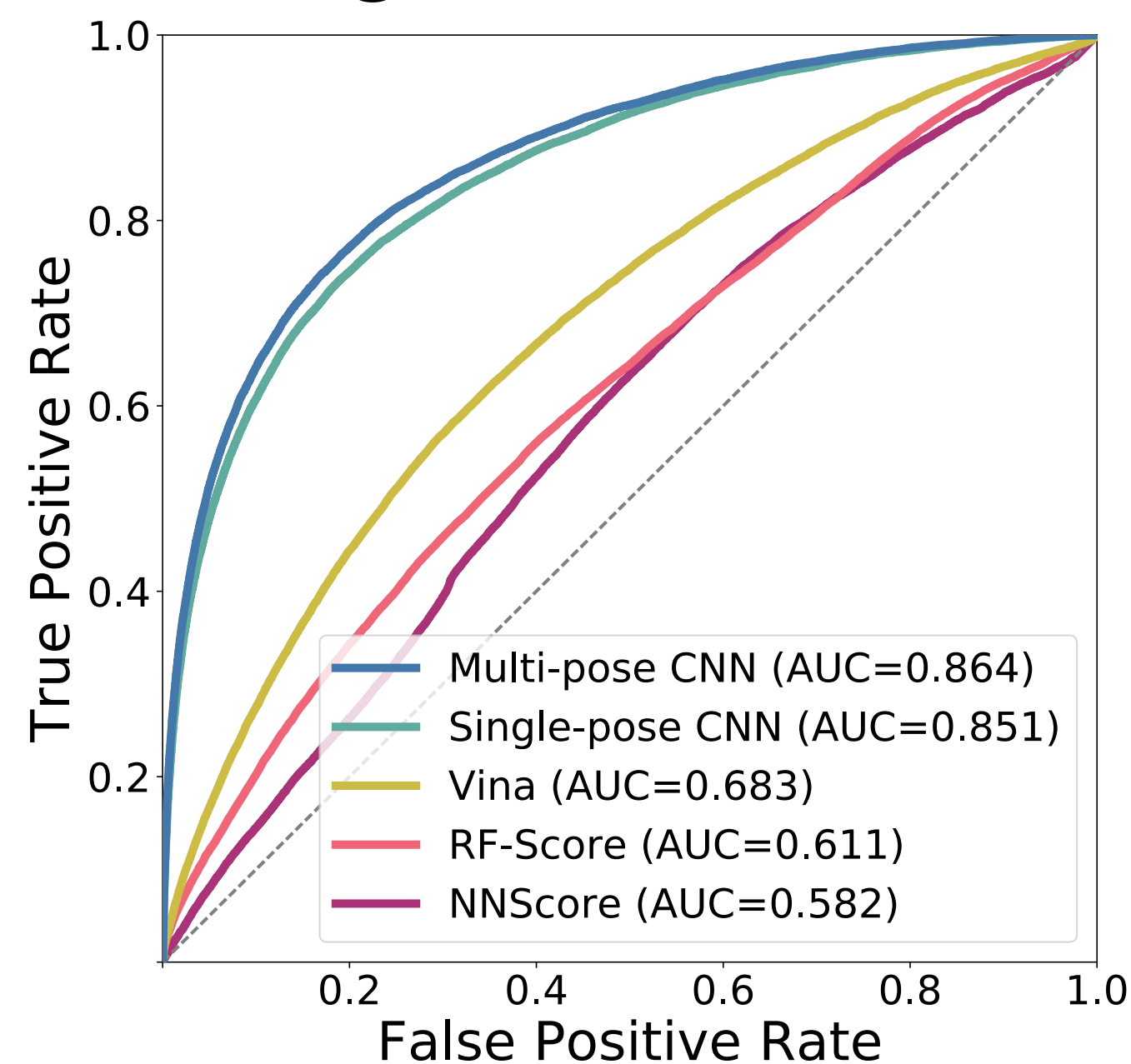Department of
Computational and
Systems Biology



TEC Bio

REU

# Questions?



## Binding Determination



## Affinity Prediction



R=0.687  RMS=2.186

## Relevance Propagation

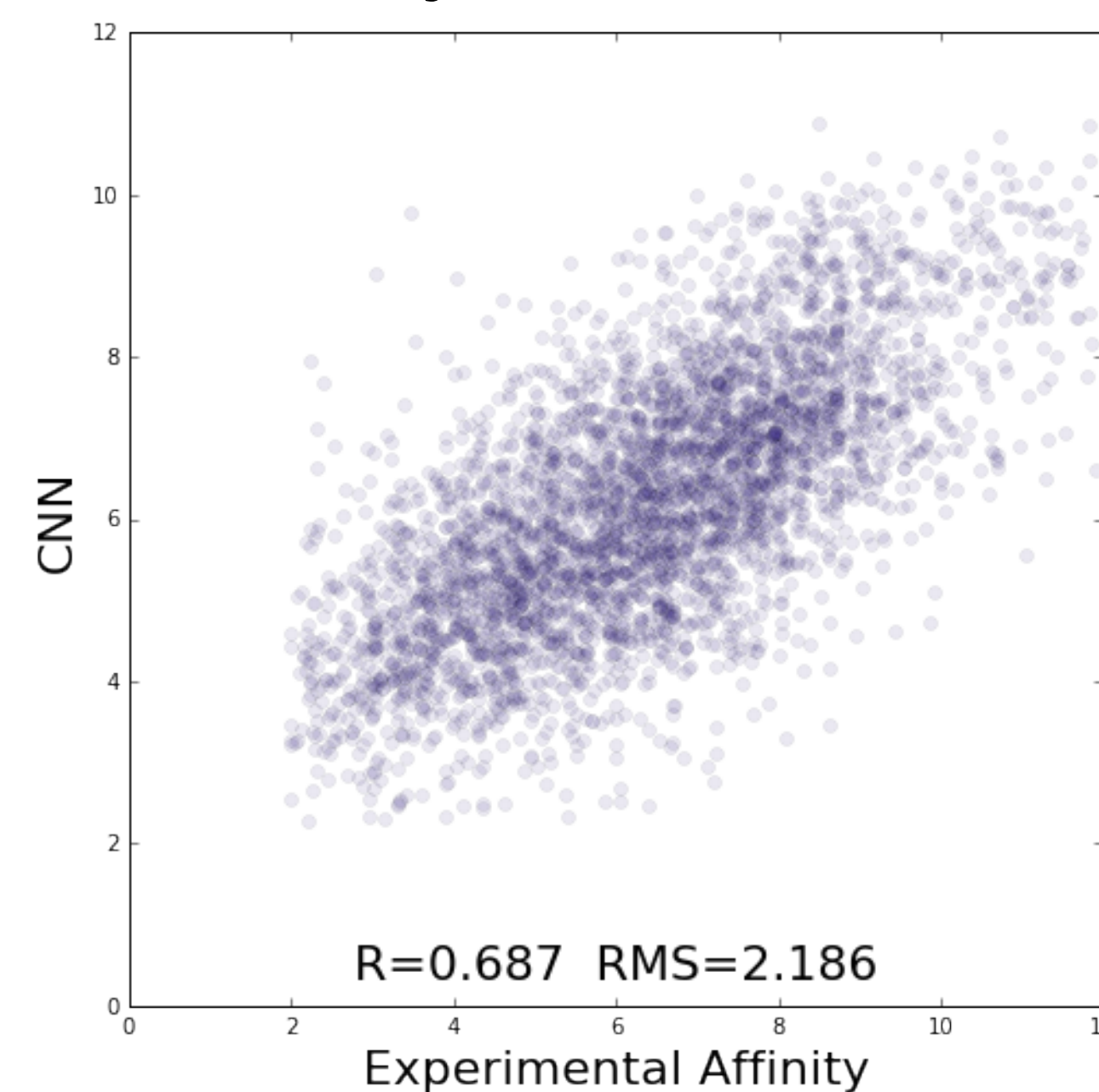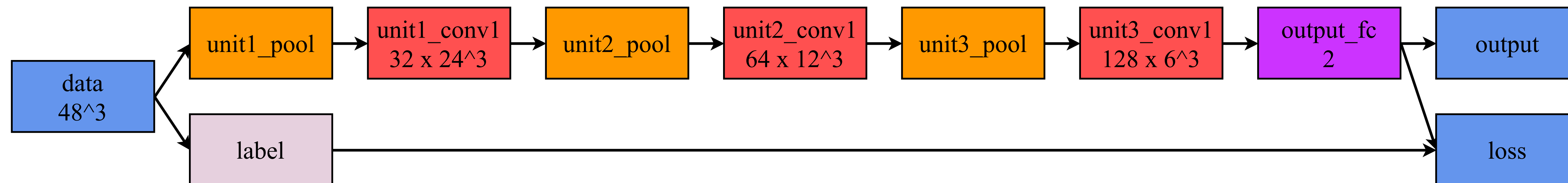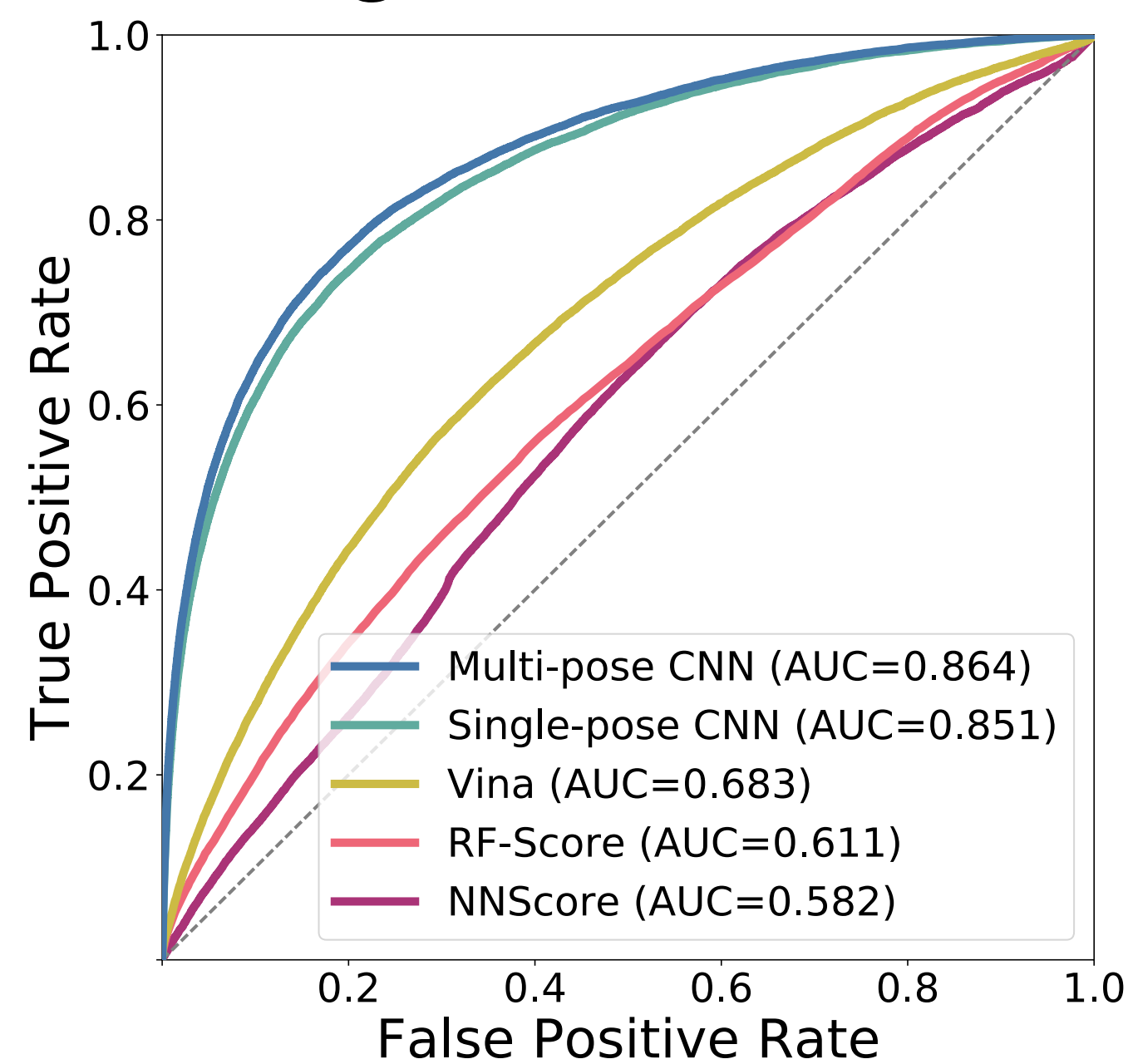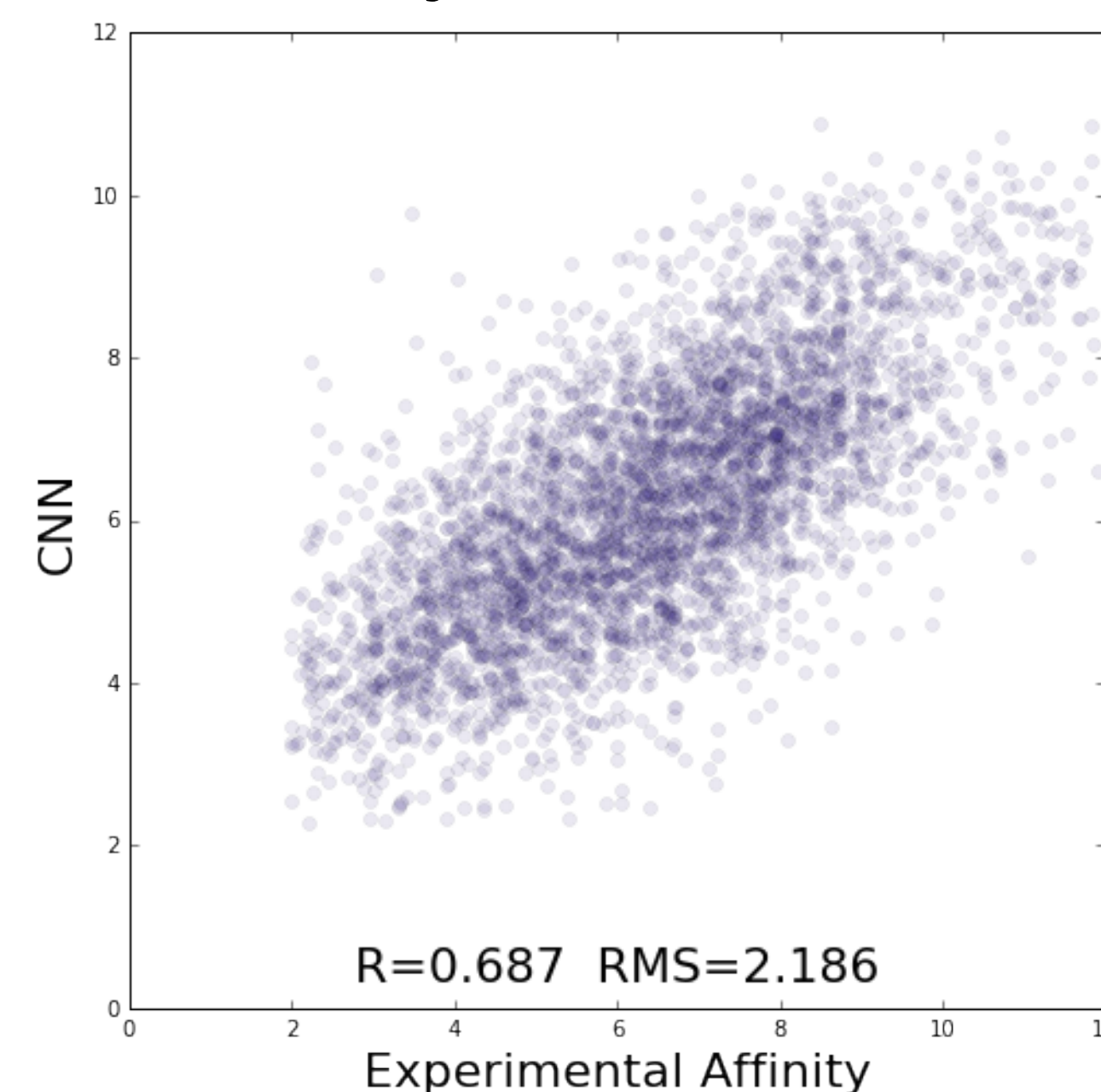# Questions?



## Binding Determination



## Affinity Prediction



## Relevance Propagation