## Protein-Ligand Scoring with Convolutional Neural Networks

American Chemical Society Meeting April 3, 2017

## David Koes

>@david\_koes

## Structure Based Drug Design Lead Optimization **Virtual Screening**



### Pose Prediction



**Binding Discrimination** 

### Affinity Prediction



## Structure Based Drug Design Lead Optimization **Virtual Screening**



### Pose Prediction



**Binding Discrimination** 

### Affinity Prediction



## **Protein-Ligand Scoring**



### AutoDock Vina



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





### 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

## Can we do better?







## Deep Learning



Ar hast - a compater program that can beat a champion Go player materi



SATEALARD TRAASPARIATS

WIESGANS







limator.





## Deep Learning



At host - a computer program that can beat a champion Go player Met of











## Image Recognition









## **Convolutional Neural Networks**

## **CNNs for Protein-Ligand Scoring**



### Pose Prediction

### Binding Discrimination

Affinity Prediction





## **CNNs for Protein-Ligand Scoring**



- Training

### Input representation

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



## Atom Types

### Ligand

AliphaticCarbonXSHydrophobe AliphaticCarbonXSNonHydrophobe AromaticCarbonXSHydrophobe AromaticCarbonXSNonHydrophobe Bromine Chlorine Fluorine lodine Nitrogen NitrogenXSAcceptor NitrogenXSDonor NitrogenXSDonorAcceptor Oxygen OxygenXSAcceptor OxygenXSDonorAcceptor Phosphorus Sulfur SulfurAcceptor

### Receptor

AliphaticCarbonXSHydrophobe AliphaticCarbonXSNonHydrophobe AromaticCarbonXSHydrophobe AromaticCarbonXSNonHydrophohe 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</li>
- generate poses with Vina
  - 745 <2Å RMSD (actives)
  - 3251 >4Å RMSD (decoys)



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)



## Training Data

### **Binding Discrimination**



102 targets

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

### **Affinity Prediction**



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





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

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



## Model Evaluation

### Clustered Cross-validation















## Data Augmentation







## Data Augmentation







### **University of Pittsburgh**



## 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



## **Cross-Validation Evaluation**



## Pose Prediction (CSAR)







## Pose Prediction (CSAR)







## Pose Prediction (PDBbind)



## Binding Determination



### E D

### 102 targets

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







## Binding Determination





### **University of Pittsburgh**











Partially Aligned Poses Combined 2:1 Training Set





Preentation	COMP-Divisio
Presentation	COMP-Divisio.

6:00pm-8:00pm Apr 4

COMP 290: Visualization of convolutional neural network scoring of protein-ligand binding

8:00pm-10:00pm Apr 3

COMP 290: Visualization of convolutional neural network scoring of protein-ligand binding











## **Beyond Scoring**







## **Beyond Scoring**







### More Oxygen Here

2Q89



### Less Oxygen Here



![](_page_35_Picture_7.jpeg)

### More Oxygen Here

2Q89

![](_page_36_Picture_4.jpeg)

### Less Oxygen Here

$\partial L$ _	$\nabla \partial L$	$\partial x_j$	$\partial dist_{A,x_{jata}}$	
$\overline{\partial A_{\{x,y,z\}}}$ –	$\sum_{j} \overline{\partial x_{j}}$	$\overline{\partial dist_{A,x_j}}$	$\partial A_{\{x,y,z\}}^{48^{3}}$	

![](_page_36_Picture_7.jpeg)

![](_page_36_Picture_8.jpeg)

![](_page_36_Picture_9.jpeg)

![](_page_37_Picture_0.jpeg)

![](_page_38_Picture_0.jpeg)

![](_page_39_Picture_0.jpeg)

![](_page_39_Picture_1.jpeg)

![](_page_40_Picture_0.jpeg)

![](_page_40_Picture_1.jpeg)

## CNN Summary

Pose Prediction (Selection)

- consistently better than Vina at *inter*-target ranking
- consistently worse than Vina at *intra*-target ranking

Binding Determination (Virtual Screening)

- Generally better than Vina, **but**
- the model is pose-insensitive

Combined Training

- Get (mostly) best of both worlds
- ... including affinity prediction

![](_page_41_Figure_13.jpeg)

You get what you train for...

![](_page_41_Figure_16.jpeg)

![](_page_41_Picture_17.jpeg)

![](_page_41_Figure_18.jpeg)

### ...but you can train for what you want

![](_page_41_Picture_20.jpeg)

![](_page_41_Picture_21.jpeg)

![](_page_41_Picture_22.jpeg)

## CNN Summary

Pose Prediction (Selection)

- consistently better than Vina at *inter*-target ranking
- consistently worse than Vina at *intra*-target ranking

Binding Determination (Virtual Screening)

- Generally better than Vina, **but**
- the model is pose-insensitive

Combined Training

- Get (mostly) best of both worlds
- ... including affinity prediction

![](_page_42_Figure_13.jpeg)

You get what you train for...

![](_page_42_Figure_16.jpeg)

![](_page_42_Picture_17.jpeg)

![](_page_42_Figure_18.jpeg)

### ...but you can train for what you want

![](_page_42_Picture_20.jpeg)

![](_page_42_Picture_21.jpeg)

![](_page_42_Picture_22.jpeg)

## Acknowledgements

![](_page_43_Picture_2.jpeg)

### Matt Ragoza

![](_page_43_Picture_4.jpeg)

![](_page_43_Picture_5.jpeg)

Josh Hochuli

![](_page_43_Picture_7.jpeg)

### Elisa Idrobo Jocelyn Sunseri

![](_page_43_Picture_9.jpeg)

### **Group Members**

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

**Gibran Biswas** 

National Institute of **General Medical Sciences** R01GM108340

![](_page_43_Picture_15.jpeg)

Department of Computational and Systems Biology

![](_page_43_Picture_17.jpeg)

![](_page_43_Picture_18.jpeg)

# ③@david\_koes ③ github.com/gnina http://bits.csb.pitt.edu

![](_page_44_Picture_1.jpeg)

# ③@david\_koes ③ github.com/gnina http://bits.csb.pitt.edu

![](_page_45_Picture_1.jpeg)