Protein-Ligand Scoring with Convolutional Neural Networks

D3R Workshop San Diego February 22, 2018

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?





Machine Learning





Computational and Systems Biology

Noce $\rightarrow y$ Prediction





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











Image Recognition

airplane	in 1	X		X	*	+	2	-7	-	
automobile					-		11			
bird	S	ſ	12			A	-	30	_	
cat	-		-	du)		1		2.5		
deer	6	48	X	RT		Y	1	2.0		
dog	17	6	-	9).	1		-	15		
frog	-	19			2			7.5	-	
horse	-	The second	P	2	(m)	TAN	-	0		
ship		Carlo	1	+	144	-		Ū		201
truck								1		





https://devblogs.nvidia.com

Convolutional Neural Networks

Convolutional Filters

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

-1	0	1	-1	-1	-1
-1	0	1	-1	8	-1
-1	0	1	-1	-1	-1

CNNs for Protein-Ligand Scoring

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**.

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

Pooling Max 2x2

| 2×| 2×| 2×32

Convolution 3×3×3

Rectified Linear Unit

| 2×| 2×| 2×64

Pooling Max 2×2

6×6×6×64

Convolution 3x3x3

Rectified Linear Unit

6×6×6×128

Fully Connected	Pseudo-Huber Loss
Fully Connected	Softmax+Logistic Loss

Pose Score

Model

Results

Trained on PDBbind refined; tested on CSAR

Beyond Scoring

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

Iterative Refinement

Iterative Refinement

Iterative Refinement

Docking vina/smina/gnina

best

poses

Sampling

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

Refinement

D3R Results

Grand Challenge 3 - CatS_stage2

Affinity Ranking - Kendall's Tau

Receipt ID

Green circle indicates your predictions (requires login)

Grand Challenge 3 - JAK2_SC2

Affinity Ranking - Kendall's Tau

Receipt ID

Green circle indicates your predictions (requires login)

Grand Challenge 3 - p38a

Affinity Ranking - Kendall's Tau

Receipt ID

Green circle indicates your predictions (requires login)

Grand Challenge 3 - TIE2

Affinity Ranking - Kendall's Tau

Receipt ID

Green circle indicates your predictions (requires login)

Grand Challenge 3 - VEGFR2

Affinity Ranking - Kendall's Tau

Receipt ID

Green circle indicates your predictions (requires login)

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
veqfr2	0.366	0.484	0.434	0.448	0.414
iak2	0.428	0.338	0.39	0.27	0.106
iak2 sub3	0.68	0.369	-0.372	0.159	-0.633
tio2	0.648	0.835	0 136	-0.078	0 561
abl1	0.634	0.745	0.005	0 182	0 713

Receipt ID

GC3: Pose Prediction

Cathepsin Phase 1

Affinity Ranking (Stage 2) - Kendall's Tau

Receipt ID

Green circle indicates your predictions (requires login)

Affinity Ranking (Stage 2) - 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

National Institute of **General Medical Sciences** R01GM108340

Matt Ragoza Josh Hochuli

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

O github.com/gnina

http://bits.csb.pitt.edu

@david_koes

Prospective Case Study: TIGIT

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

Cancer Cell Article

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,*} ¹Department of Cancer Immunology ²Department of Protein Chemistry ³Department of Bioinformatics and Computational Biology ⁴Department of Oncology Biomarker Development ⁵Department of Translational Oncology ⁶Department of Translational Immunology ⁷Department 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

University of Pittsburgh

Computational and Systems Biology

Filter Visualization

