

# Fragment binding prediction using unsupervised learning of ligand substructure binding sites

Grace Tang  
Altman Lab  
July 20, 2013  
Berlin, Germany

# Structure Based Virtual Screening

Databases ..... Virtual Screening ..... End Goal

**ZINC<sup>1</sup>**

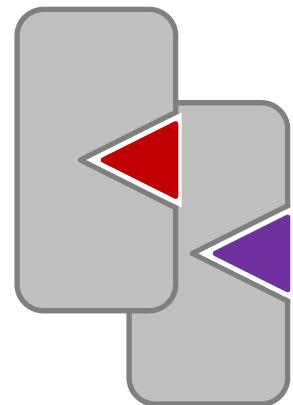
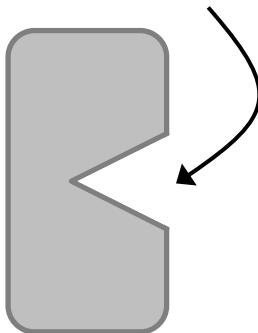
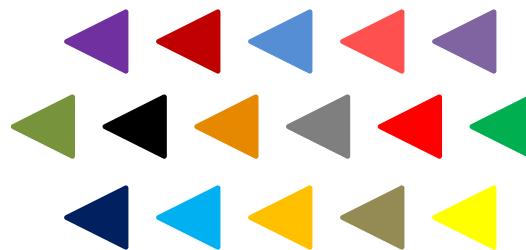
21 million  
purchasable  
compounds

**PubChem<sup>2</sup>**

48 million  
compounds

**GDB-13<sup>3</sup>**

970 million  
drug-like  
small  
molecules



1. Irwin, J.J., et al., ZINC: A Free Tool to Discover Chemistry for Biology. *J Chem Inf Model*, 2012.

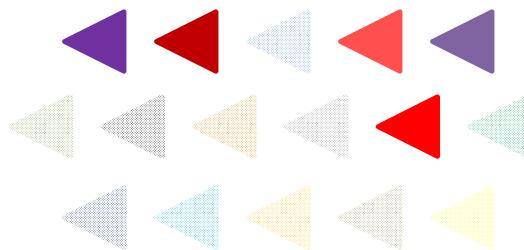
2. Bolton E., et al., PubChem: Integrated Platform of Small Molecules and Biological Activities. *Annual Reports in Computational Chemistry*, 2008. 2

# Adding Prior Knowledge

Databases ..... Virtual Screening ..... End Goal

**ZINC<sup>1</sup>**

21 million  
purchasable  
compounds

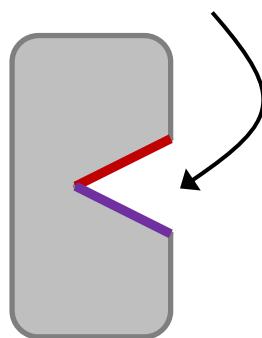


**PubChem<sup>2</sup>**

48 million  
compounds

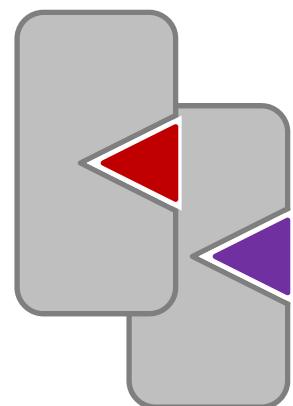
**GDB-13<sup>3</sup>**

970 million  
drug-like  
small  
molecules



contains fragment: X

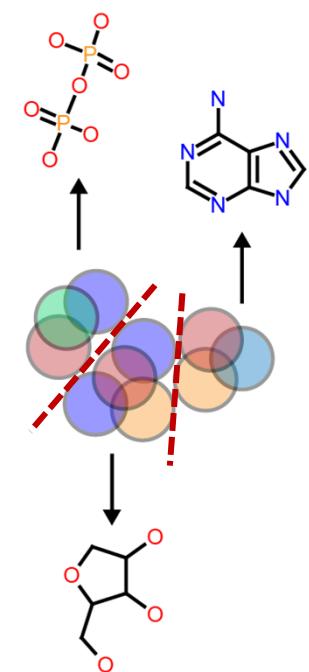
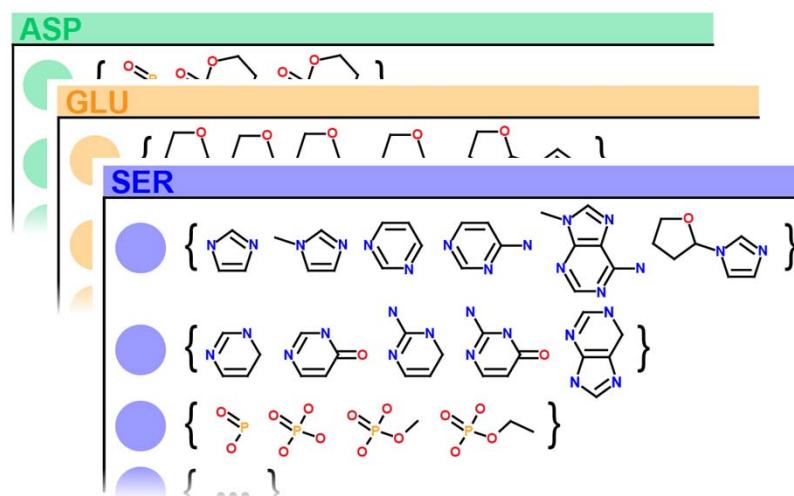
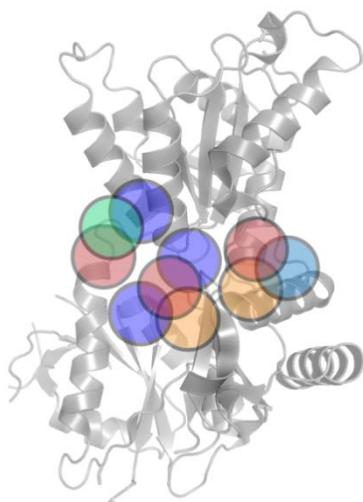
contains fragment: Y



3. Blum, L.C. and J.L. Reymond, 970 million druglike small molecules for virtual screening in the chemical universe database GDB-13. J Am Chem Soc, 2009. **131**(25): p. 8732-3.

# Method Overview

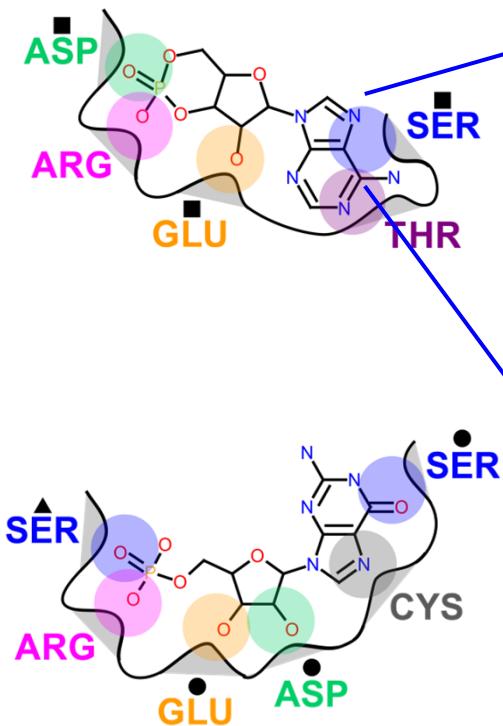
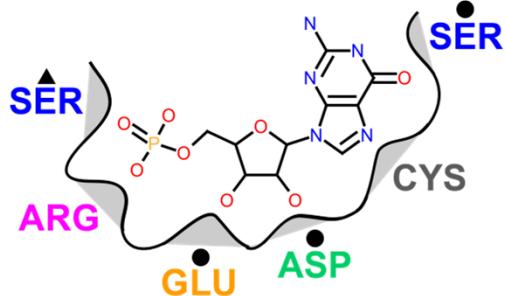
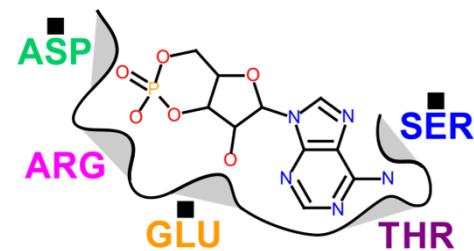
Query Protein ..... Knowledge Base ..... Fragment Predictions



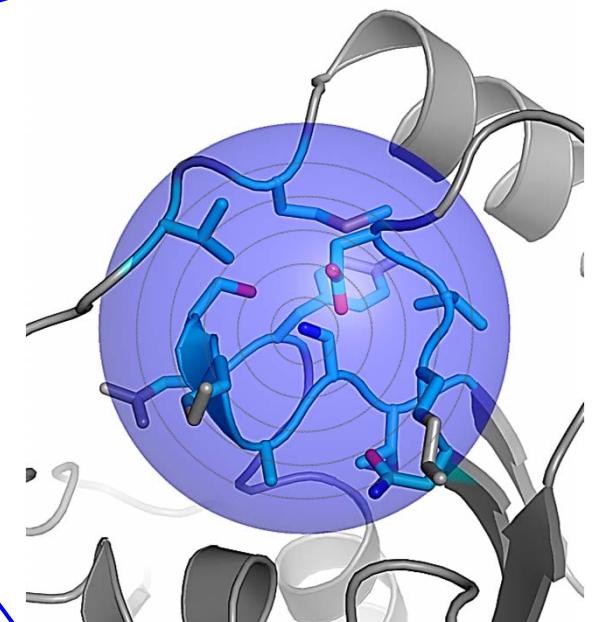
# Derivation of Knowledge Base

Protein-Ligand  
Complexes<sup>1</sup>

Structural  
Information<sup>2</sup>



FEATURE



Atom Type  
Atom Element  
Residue Name  
Residue Class

Partial Charge  
Hydrophobicity  
Aromatic  
etc.

1. Berman, H.M., et al., *The Protein Data Bank*. Nucleic Acids Res, 2000. **28**(1): p. 235-42.

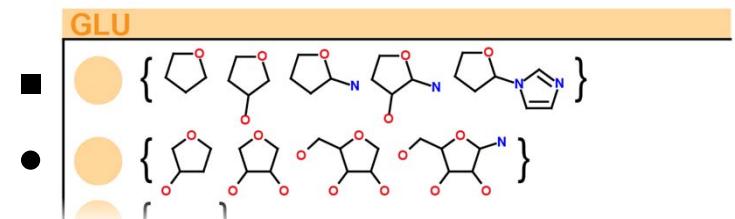
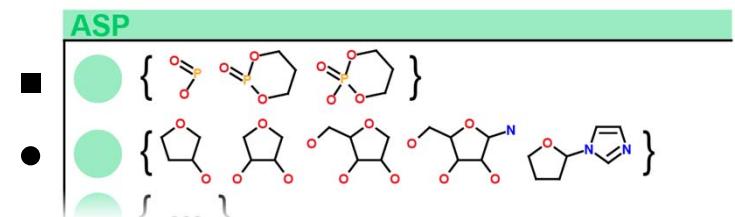
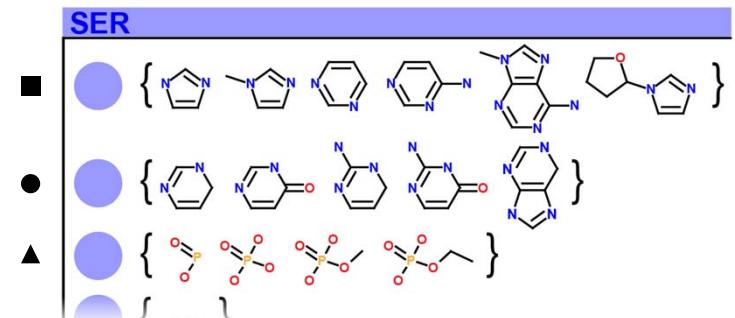
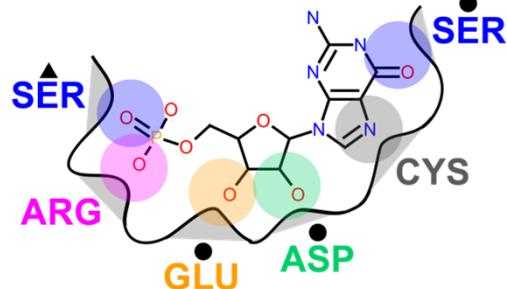
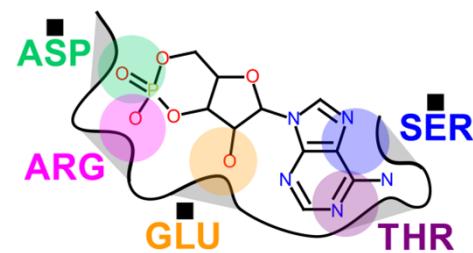
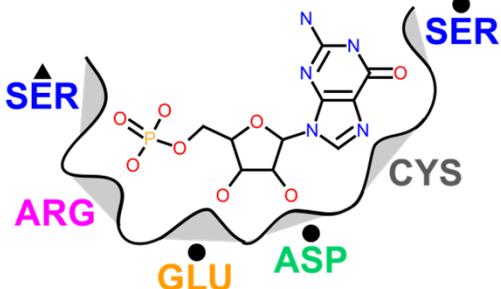
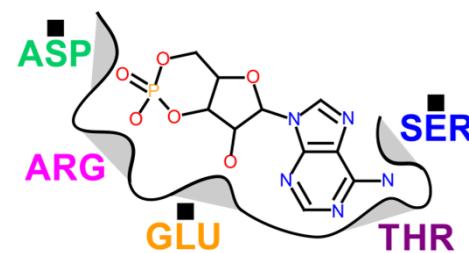
2. Halperin, I., et al., *The FEATURE framework for protein function annotation*. BMC Genomics, 2008. **9 Suppl 2**: p. S2.

# Derivation of Knowledge Base

Protein-Ligand  
Complexes<sup>1</sup>

Fragment  
Information<sup>2</sup>

Knowledge Base

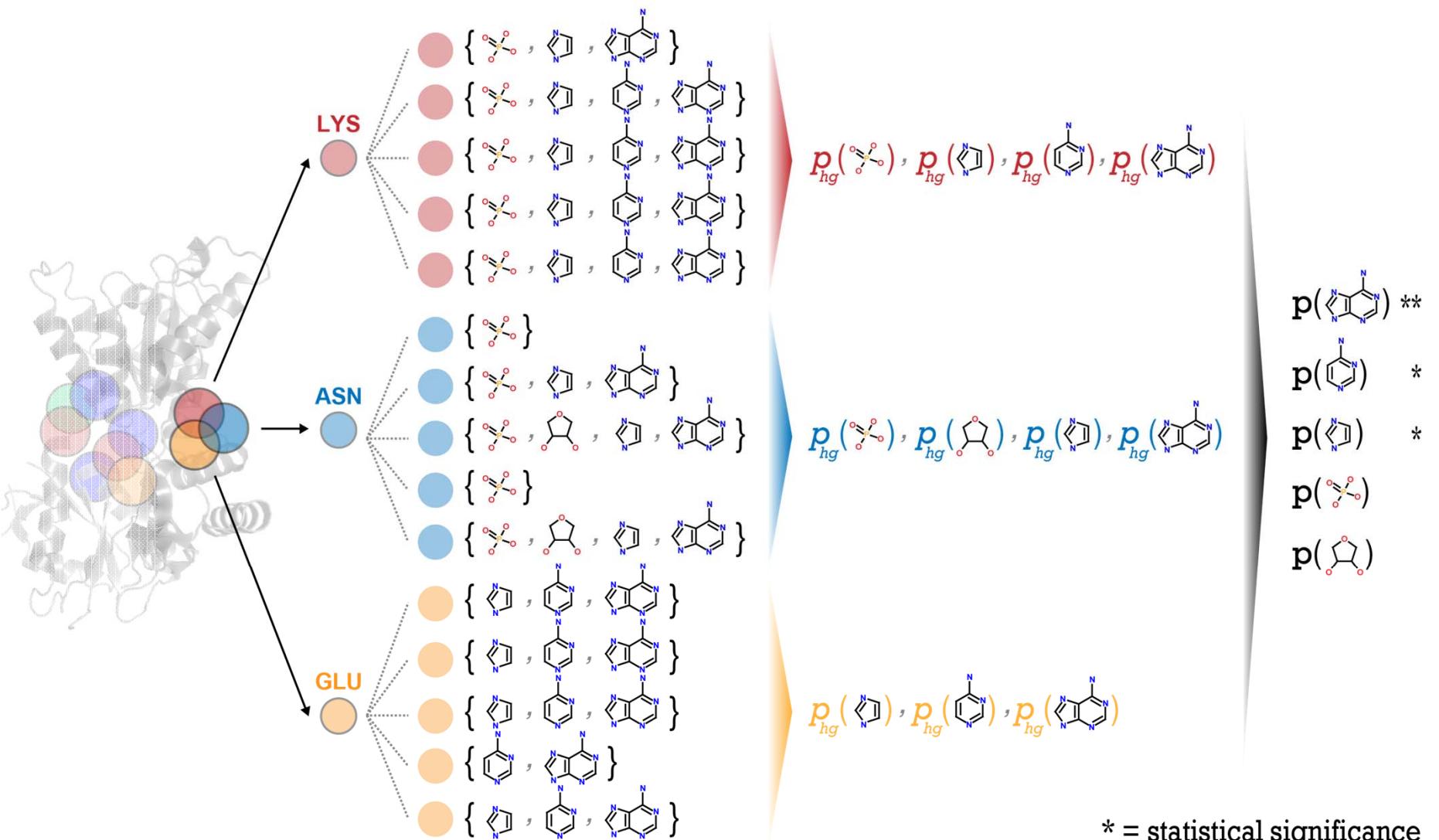


1. Berman, H.M., et al., *The Protein Data Bank*. Nucleic Acids Res, 2000. **28**(1): p. 235-42.

2. Rahman, S.A., et al., *Small Molecule Subgraph Detector (SMSD) toolkit*. J Cheminform, 2009. **1**(1): p. 12.

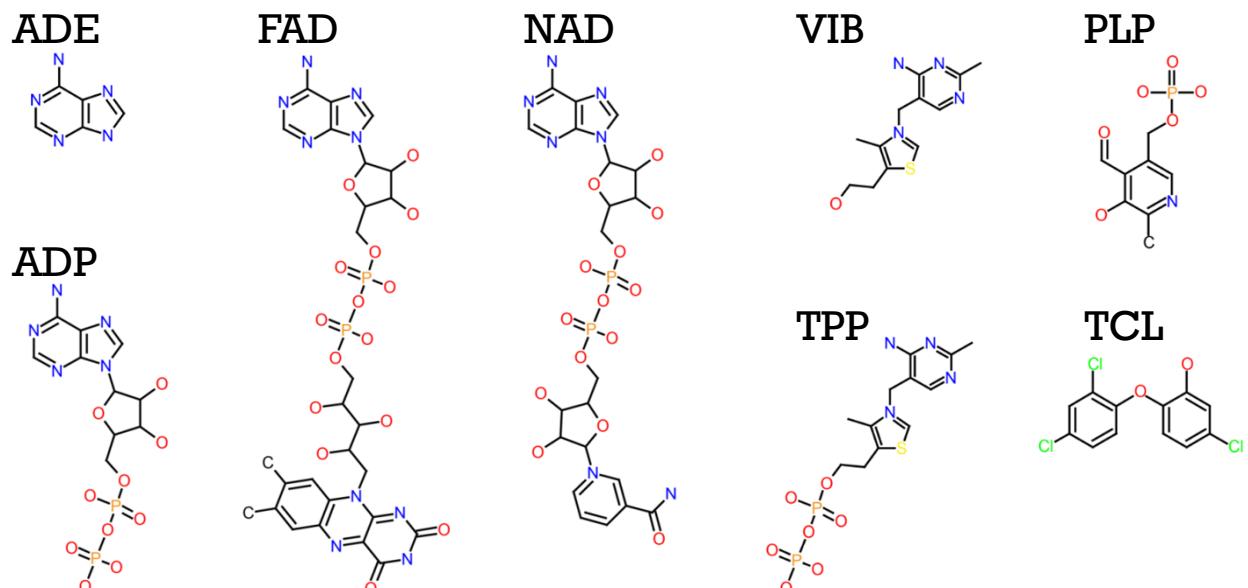
# Method

Query Protein ..... 5 Nearest Non-Homologous Neighbors ..... Hypergeometric p-value ..... Fisher's p-value



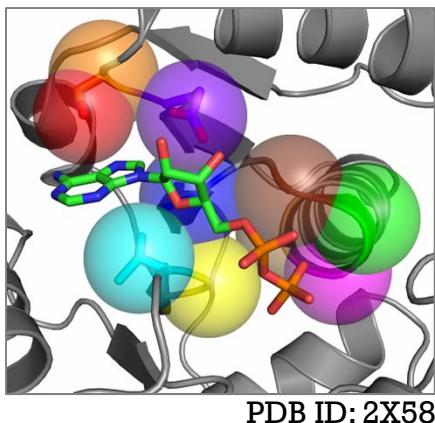
# Validation Ligands

ADE	adenine	n = 123
ADP	adenosine-5'-diphosphate	n = 2640
FAD	flavin-adenine dinucleotide	n = 2769
NAD	nicotinamide-adenine dinucleotide	n = 2309
VIB	thiamin, vitamin B1	n = 19
TPP	thiamine diphosphate	n = 217
PLP	pyridoxal-5'-phosphate	n = 1227
TCL	triclosan	n = 88

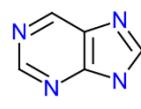


# Validation

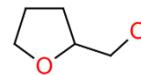
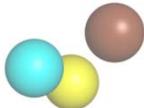
## Protein – ADP Complex



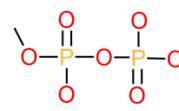
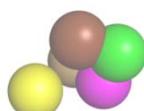
## Predictions and Moieties



adenine ring 1 & 2



ribose



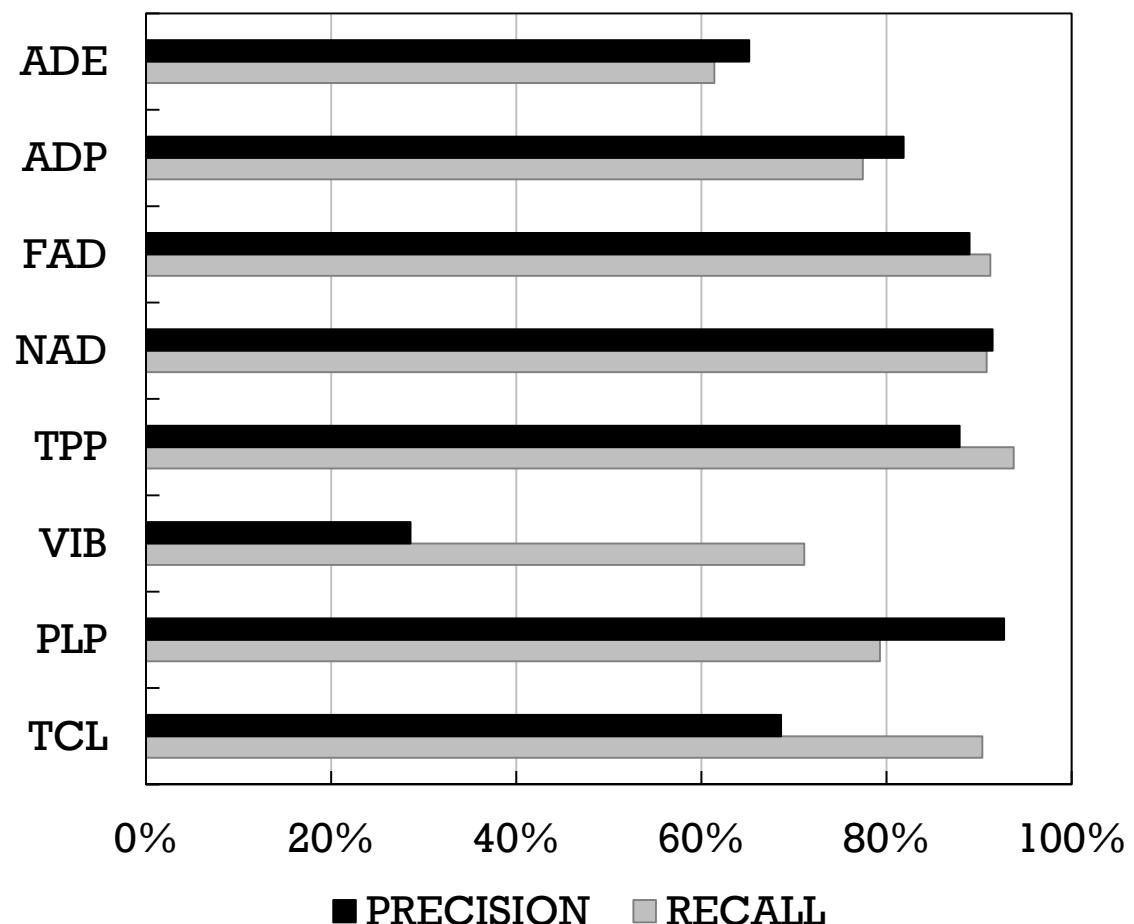
phosphate 1 & 2

Precision:

$$\frac{\text{\# Correct Predictions}}{\text{\# Total Predictions}}$$

Recall:

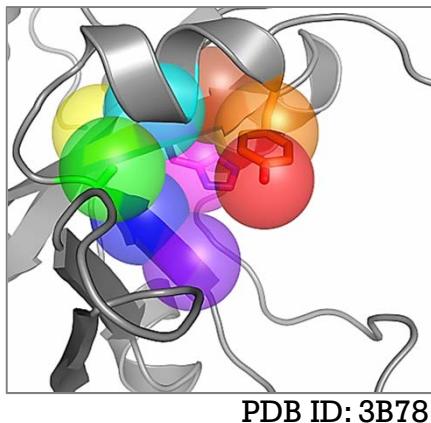
$$\frac{\text{\# Predicted Moieties}}{\text{\# Bound Moieties}}$$



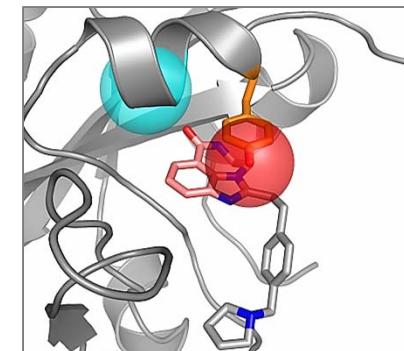
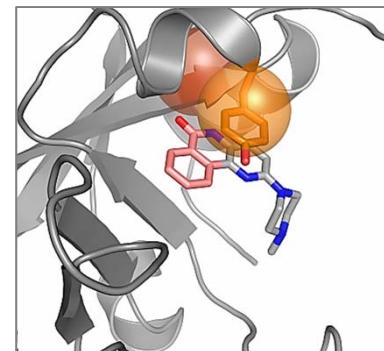
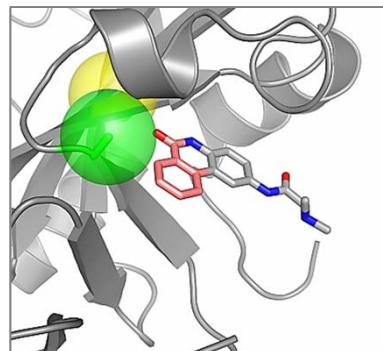
# Exotoxin A

Exotoxin A from *Pseudomonas aeruginosa* is an ADP-ribosyltransferase that inactivates eukaryotic ribosomal elongation factor 2, preventing protein synthesis and triggering cell necrosis.

Query



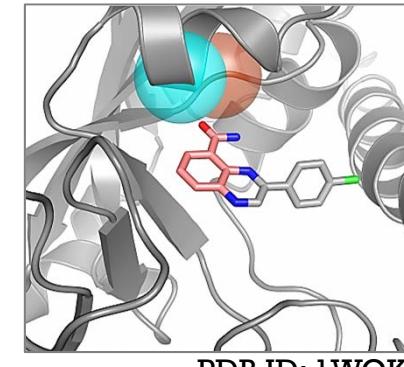
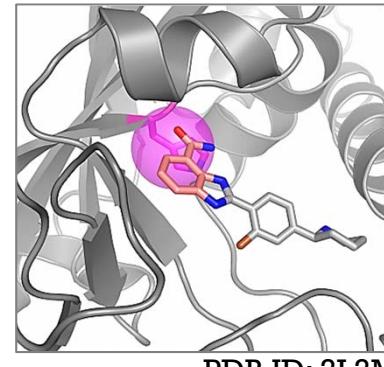
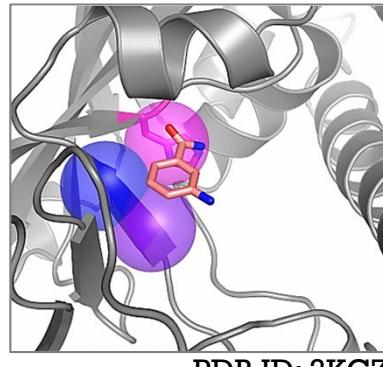
Example Nearest Microenvironment Neighbors



Predicted Fragment

CID: 450318

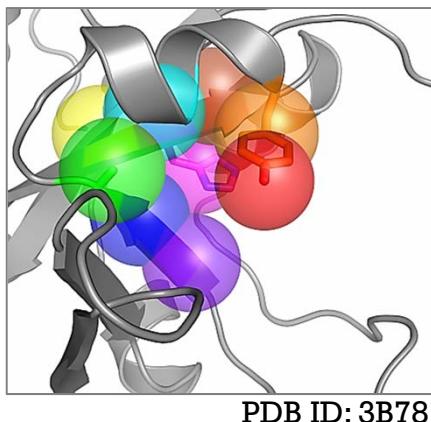
( $p$ -value  $1 \times 10^{-28}$ )



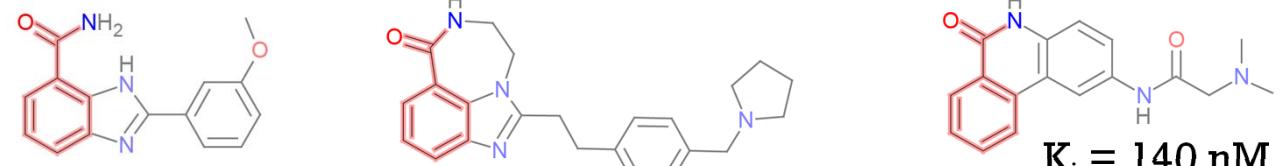
# Exotoxin A

Exotoxin A from *Pseudomonas aeruginosa* is an ADP-ribosyltransferase that inactivates eukaryotic ribosomal elongation factor 2, preventing protein synthesis and triggering cell necrosis.

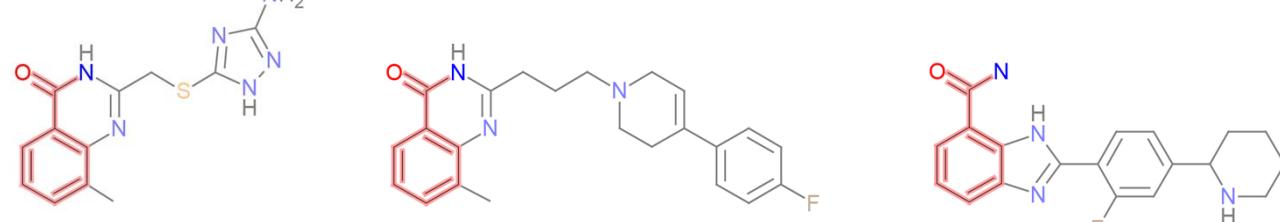
Query



Ligands Bound by Nearest Microenvironment Neighbors



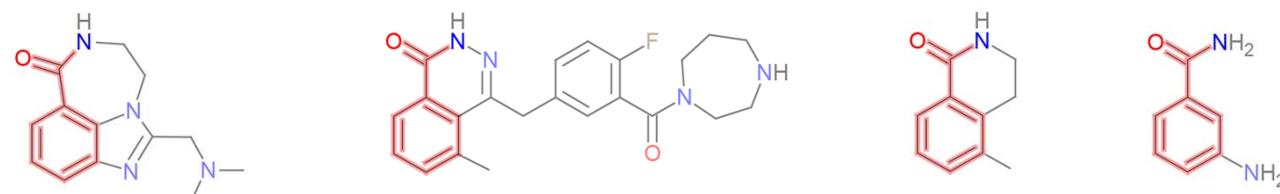
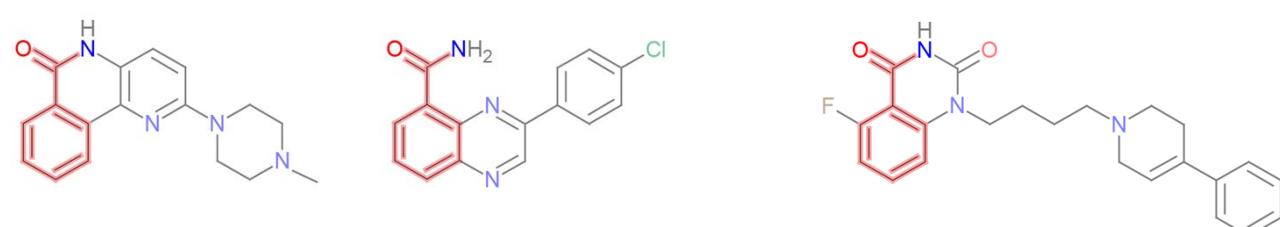
$K_i = 140 \text{ nM}$



Predicted Fragment

CID: 450318

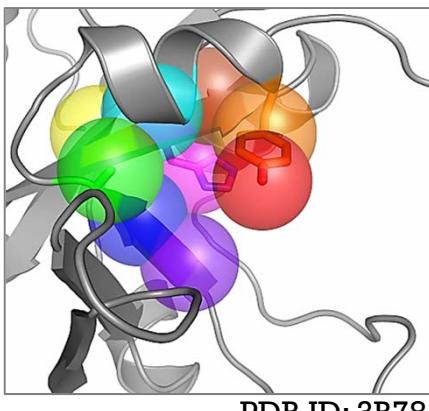
( $p$ -value  $1 \times 10^{-28}$ )



# Exotoxin A

Exotoxin A from *Pseudomonas aeruginosa* is an ADP-ribosyltransferase that inactivates eukaryotic ribosomal elongation factor 2, preventing protein synthesis and triggering cell necrosis.

Query



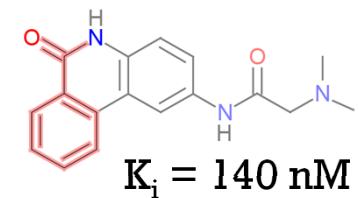
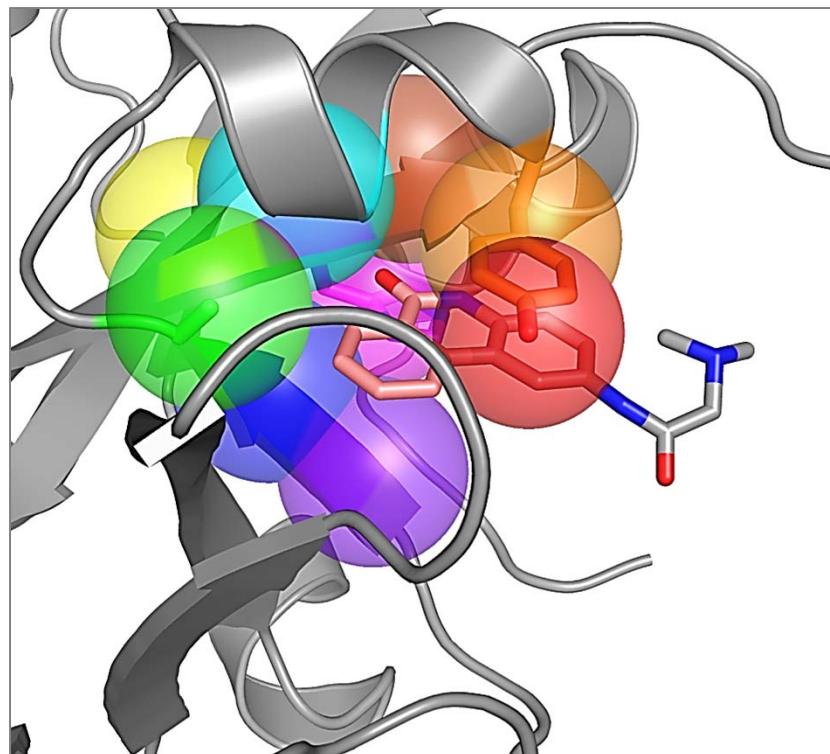
Predicted Fragment

CID: 450318

( $p$ -value  $1 \times 10^{-28}$ )



Structural validation of predicted fragment

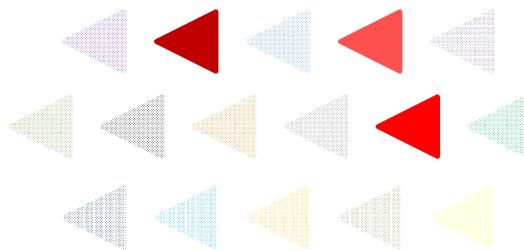


# Adding Fragment Knowledge

Databases ..... Virtual Screening ..... End Goal

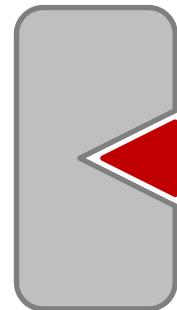
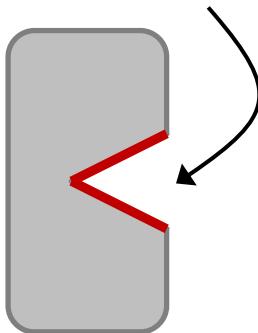
## ZINC

21 million  
purchasable  
compounds



## PubChem

48 million  
compounds



## GDB-13

970 million  
drug-like  
small  
molecules

contains fragment:



# Acknowledgements

- Russ Altman
- Altman Lab
- Funding Sources



- Travel funding to ISMB/ECCB 2013 was generously provided by ISCB



# Thank You!

gwtang@stanford.edu

3DSIG:  
Poster ... 046

ISMB/ECCB:  
Poster ... L081

## Fragment binding prediction using unsupervised learning of ligand substructure binding sites

Grace W. Tang<sup>1</sup> and Russ B. Altman<sup>1,2</sup>

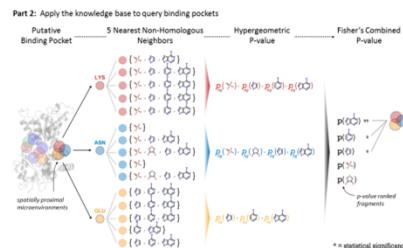
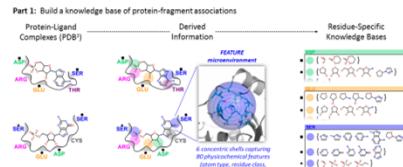
Departments of <sup>1</sup>Bioengineering and <sup>2</sup>Genetics, Stanford University

### Motivation

In structure-based drug discovery, virtual screening plays an important role in the identification of hit and lead molecules for a protein target of interest. Docking algorithms are popular for this task but have limited throughput in order to account for target and ligand flexibility. These algorithms cannot handle the enormous size of small molecule databases, and most of them lack a measure of method complexity and accuracy. Identification of database subsets pre-filtered for molecules with specific interactions with a protein target facilitates computational screening techniques. As bioactive compounds against a target frequently share chemical substructures, we propose an unsupervised machine learning approach to predict small molecule fragments compatible with a target protein structure. We take advantage of the availability of structural data for proteins whose bound ligands share substructures to enhance our understanding of fragment binding to develop a fragment predictor. Our method requires as input a 3D structure of the target but requires no prior knowledge of active compounds, enabling broad usage.

### Method

For all protein residues involved in ligand binding, we collect their local structural microenvironment using FEATURE and annotate them with the ligand fragments they bind<sup>1</sup>. This serves as the knowledge base of protein-fragment interactions. Comparison to the knowledge base enables retrieval of fragments statistically preferred by the microenvironments of a target structure.

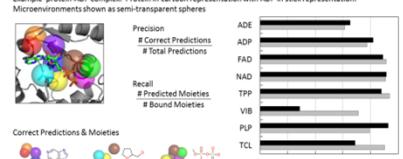


### Performance

Results on a dataset of proteins binding a variety of ligands show strong ability to rediscover fragments corresponding to the ligand bound, validating the knowledge base and method.

	n	
ADE	adenine	123
ADP	adenosine-5'-diphosphate	2640
FAD	flavin-adenine dinucleotide	2769
NAD	nicotinamide-adenine dinucleotide	2309
VIB	thiamine-vitamin B1	237
TPP	thiamine diphosphate	227
PLP	pyridoxal-5'-phosphate	2227
TCL	tricloso	88

Example: protein-ADP complex. Protein in cartoon representation with ADP in stick representation. Microenvironments shown as semi-transparent spheres.



### References

1. Johnson, L., Olson, J. S., Shi, L., & Altman, R. B. The FEATURE framework for protein function annotation: modeling new functions, improving performance, and extending to novel applications. *BMC Bioinform*. 9, 132 (2008).
2. Hwang, K. S., Kim, D. Y., & Cho, K. J. A fast algorithm for small molecule subgraph detection (SMID) toolkit. *J. Cheminf*. 3, 11 (2011).
3. Bernier, M. M. et al. The Protein Data Bank. *Biochem Biophys Res Commun*. 299, 295–298 (2002).

This work is supported by NIMHD5602 and Travel Funding provided by

### Inhibitor Fragment Predictions

For multiple protein targets, we identify high scoring fragments that are substructures of known inhibitors. Our method therefore predicts fragments suitable for preprocessing small molecule databases to enrich for bioactive compounds for a given protein target.

Exotoxin A from *Pseudomonas aeruginosa* is an ADP-ribosyltransferase that inactivates eukaryotic ribosomal elongation factor 2, preventing protein synthesis and triggering cell necrosis. Shown are nearest neighbors for the microenvironments predicting benzamide.

