

# Computational tools in pharmacology: Merging artificial and organic intelligence



Shannon T. Smith  
Ph.D. Candidate – Dr. Jens Meiler  
6 May 2022

# Disclaimers:

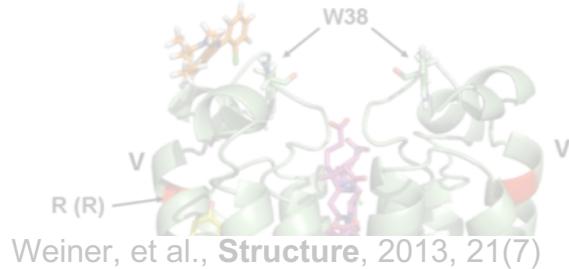
- Some slides were generously given by Dr. Jens Meiler
- Others were adapted from Dr. Carrie Jones's Modern Drug Discovery class
  - Drs. Dave Weaver, Craig Lindsley, Alex Waterson, Jeff Conn, Colleen Niswender



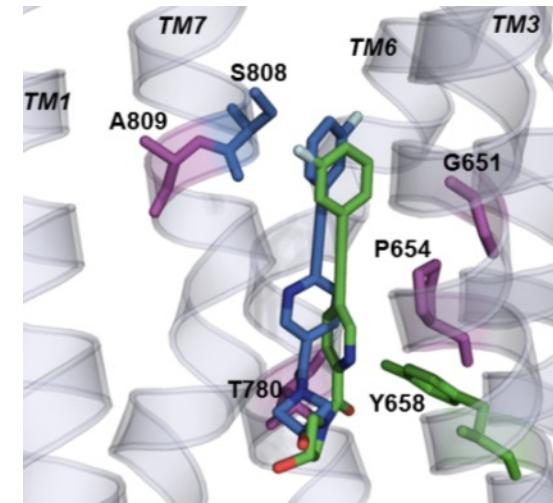
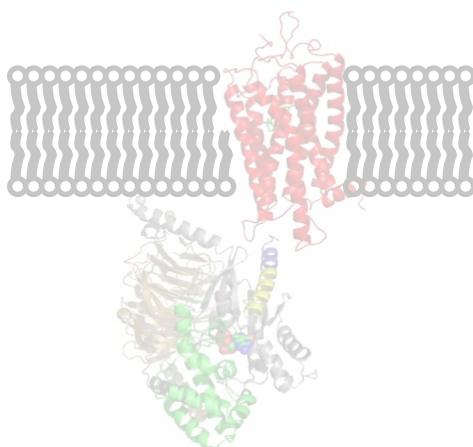
Olivia Marie Smith  
Born May 3, 2022



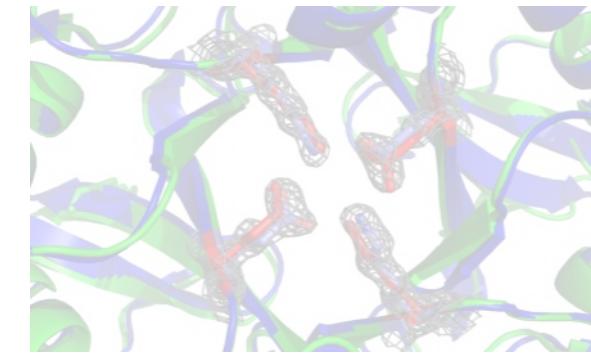
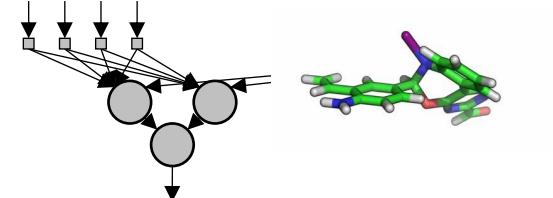
# Computational Structural and Chemical Biology in the Meiler Lab



*Protein structure prediction  
de novo and from limited  
experimental data*



*Merging ligand- and structure-  
based computer-aided drug  
discovery*



*Design of large protein  
scaffolds, antibodies, and  
protein/ligand interfaces*

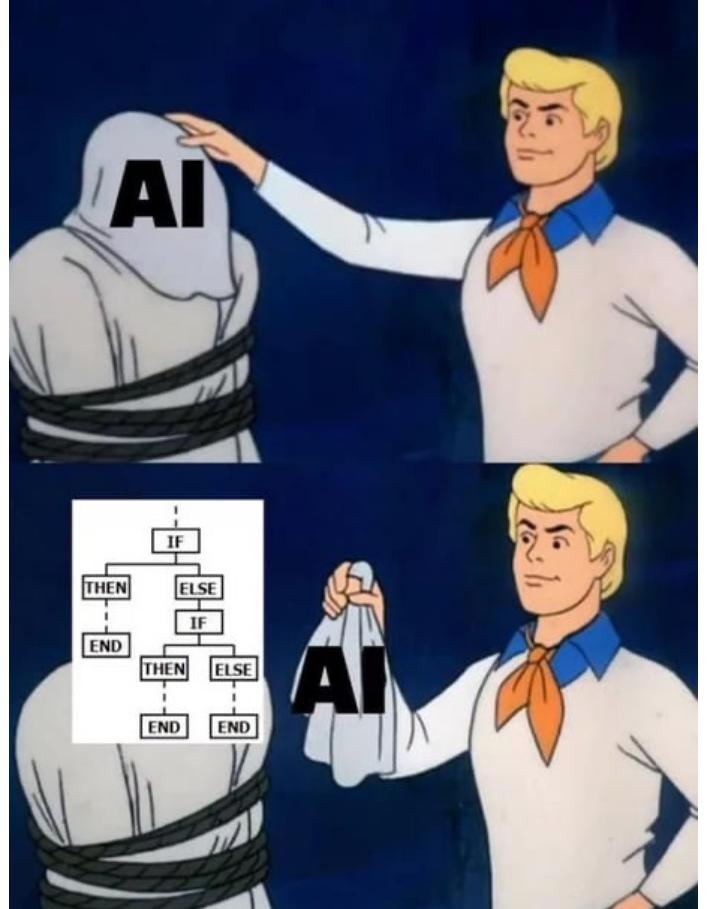


# Goals for this talk:

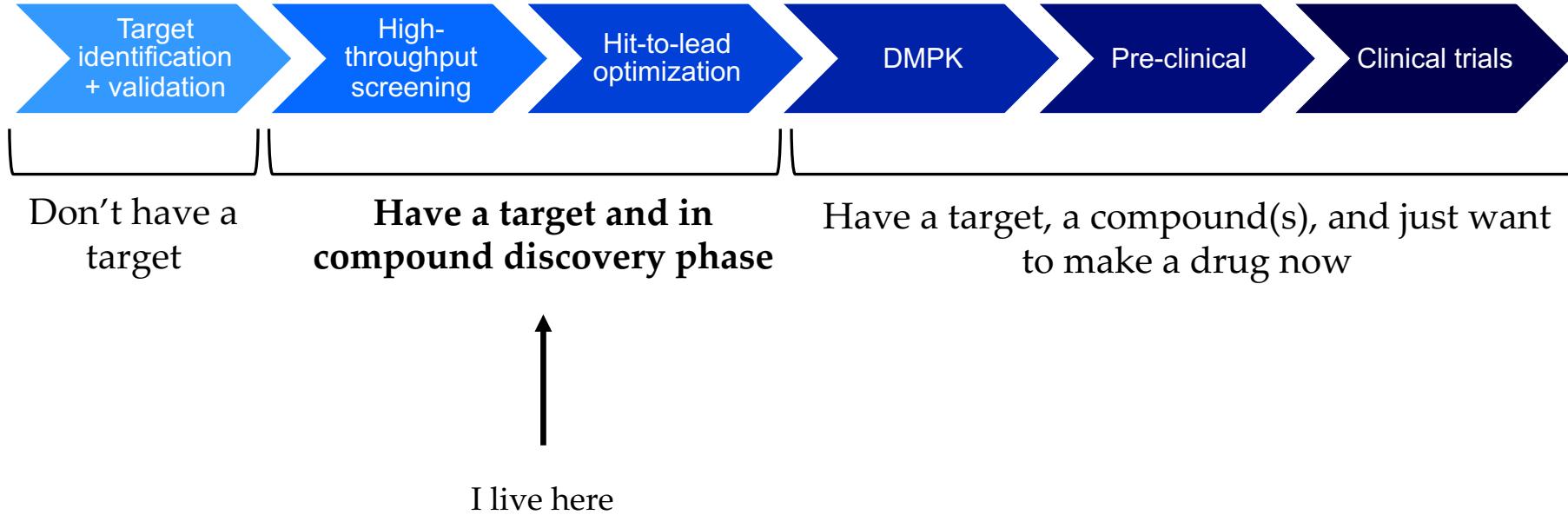
- Obtain a basic understanding for how computation, including AI, is currently applied in drug discovery at different stages and how you can apply this to your own work
- Demystifying AI in drug discovery
- Outline the advantages and disadvantages, provide realistic expectations

This talk is **NOT** about:

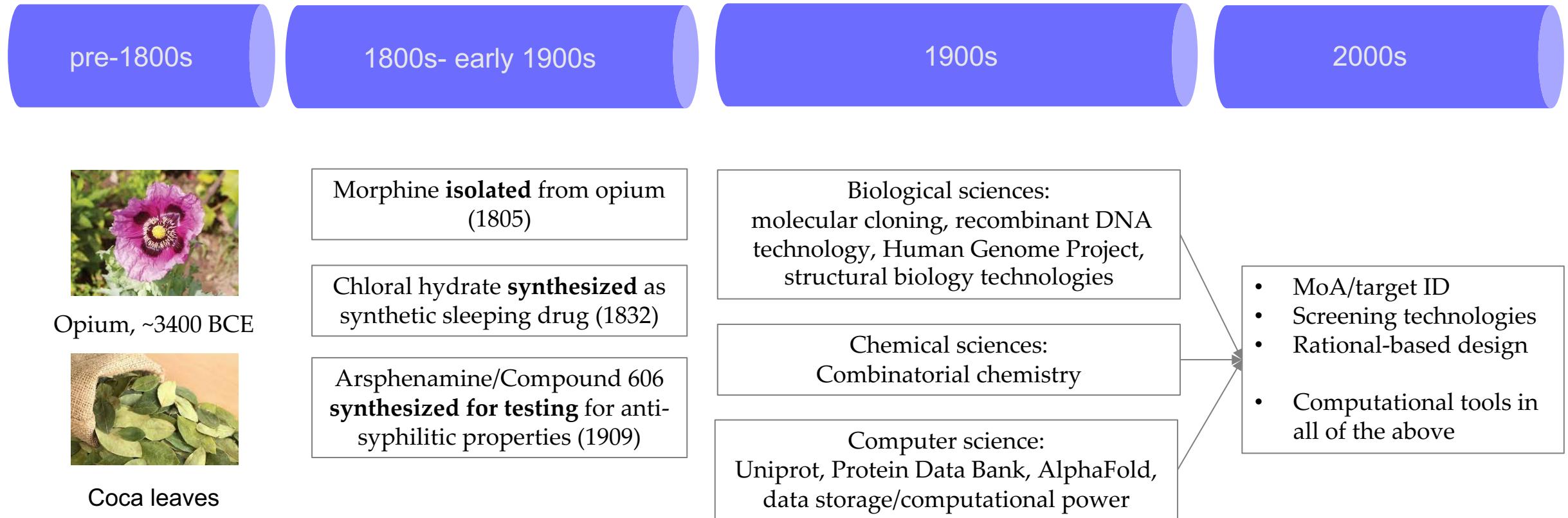
- Explaining the details of algorithms
- Me convincing you that AI will solve all your problems (though it can help)



Which stage are you working in?

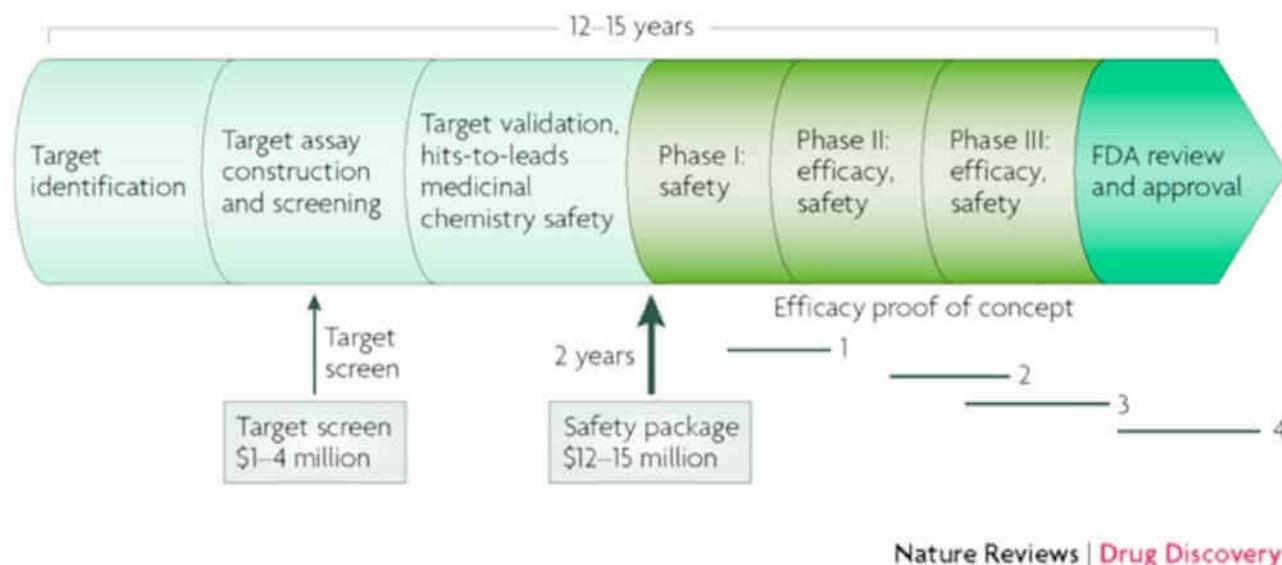


# Drug discovery: a crude history



# Why use AI in drug discovery?

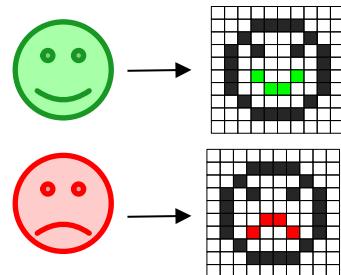
- Because we need to make drugs faster and cheaper



# Machine Learning (ML) is great for pattern recognition

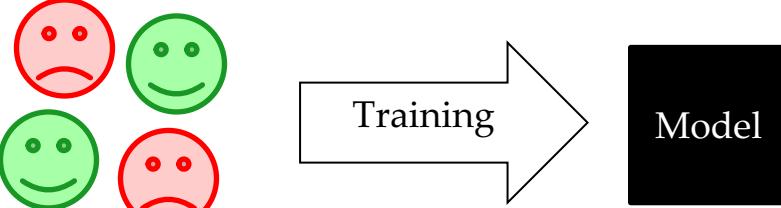
## Step 1: Data acquisition

Make faces machine-readable)

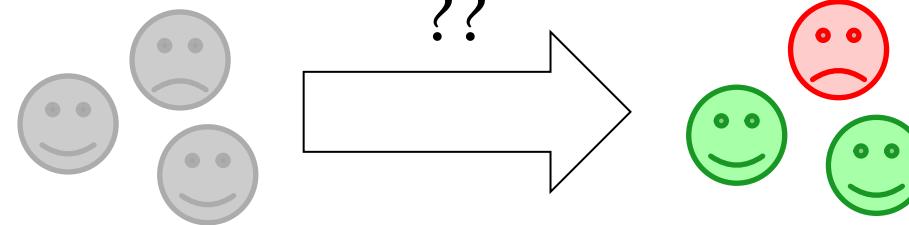


## Step 2: Training

Given known smiles and frowns, computer figures out which features define each.

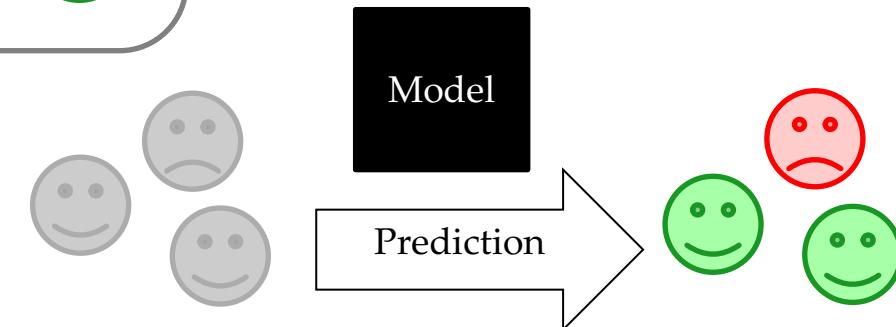


What determines if a face is a smile or frown?



## Step 4: Prediction

Given faces with unknown expressions, can we predict smile or frown?



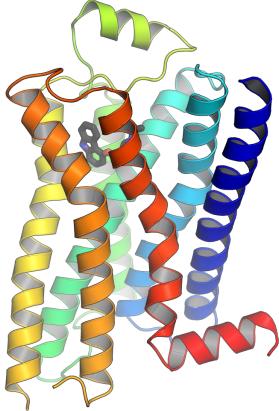
## Step 3: Validation

How well do we perform on a known test set?

		Actual	
		Smile	Frown
Predicted	Smile	9	1
	Frown	0	10

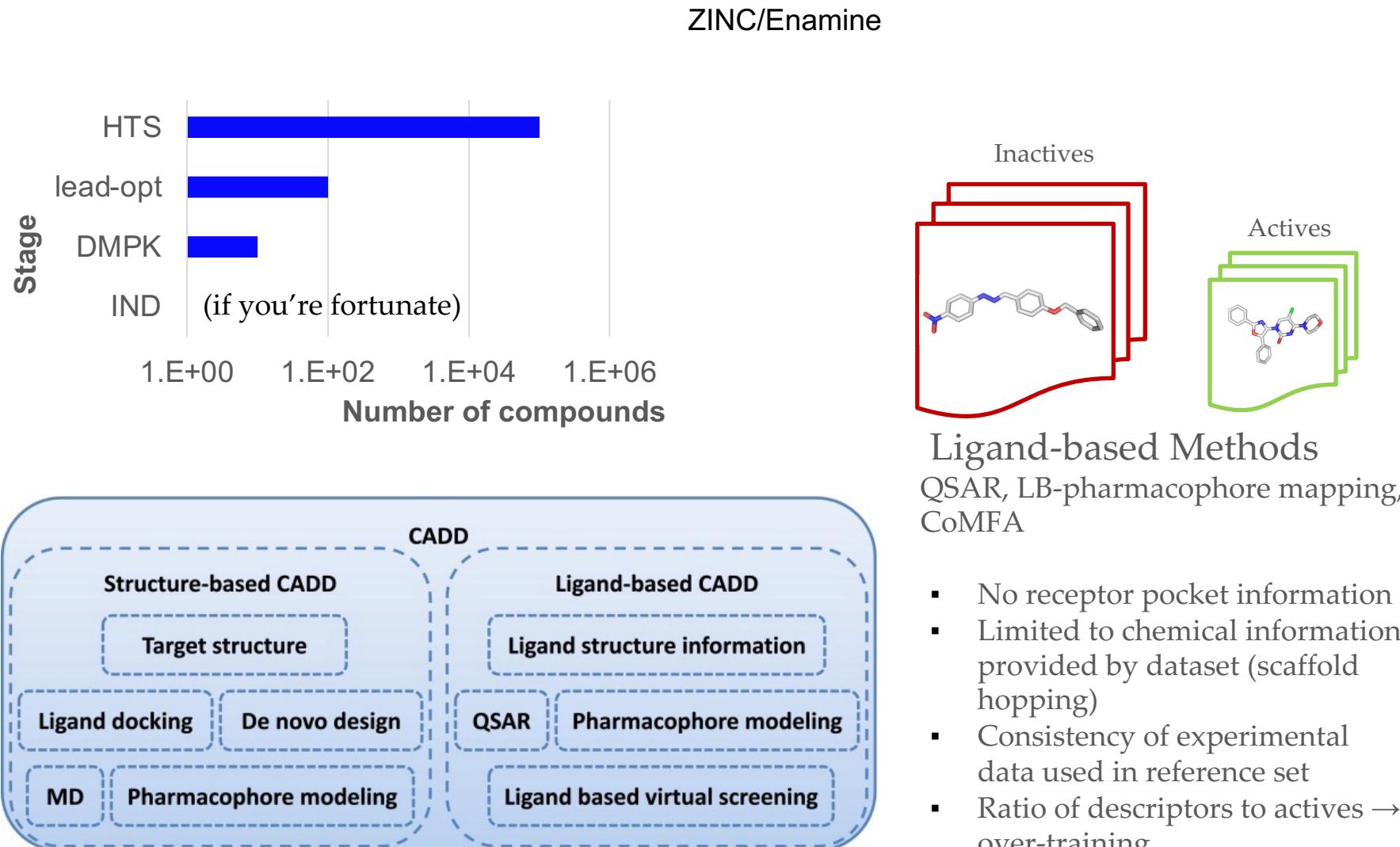
$$\text{Accuracy} = (9+10)/20 = 0.95$$

# Computer-aided drug discovery from a distance



Structure-based Methods  
Docking, Molecular Dynamics, SB-pharmacophore mapping

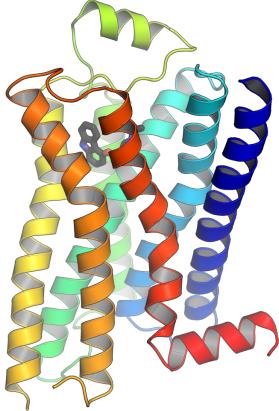
- Need high-resolution structure
- Understand mode of binding, key residues
- Rely on scoring functions
- Poorly ranks binding affinities of multiple compounds
- Slow/more expensive



- No receptor pocket information
- Limited to chemical information provided by dataset (scaffold hopping)
- Consistency of experimental data used in reference set
- Ratio of descriptors to actives → over-training

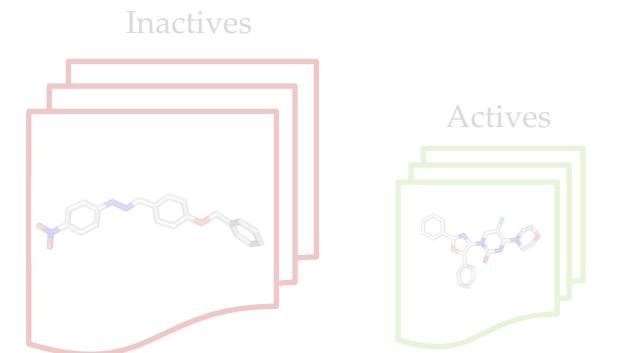
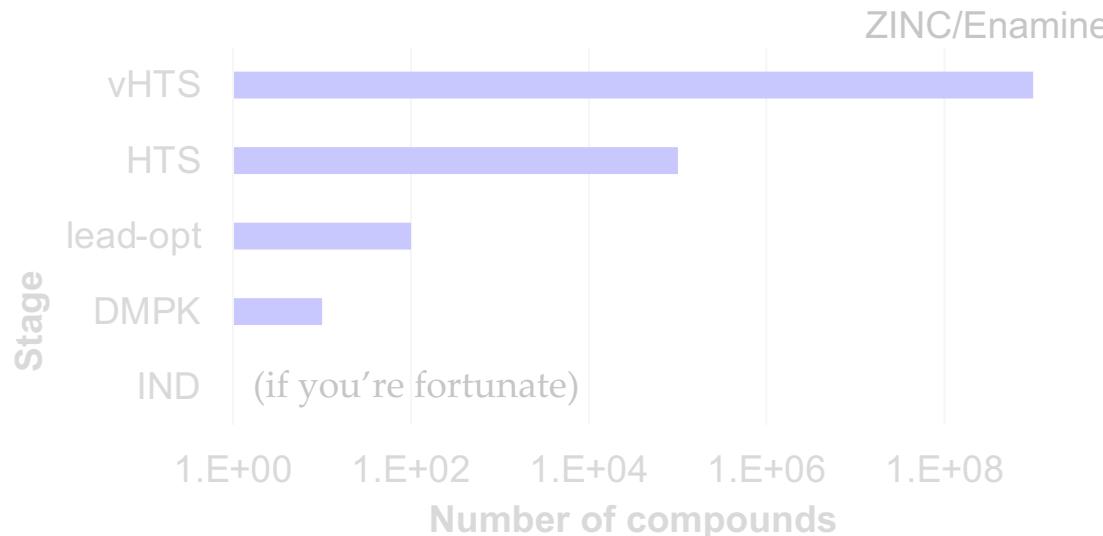


# Computer-aided drug discovery from a distance

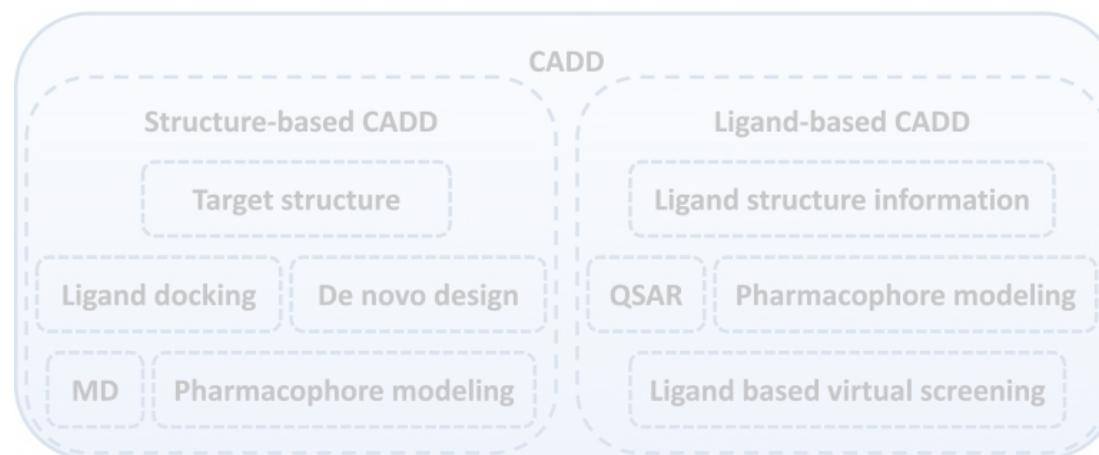


Structure-based Methods  
Docking, Molecular Dynamics, SB-pharmacophore mapping

- Need high-resolution structure
- Understand mode of binding, key residues
- Rely on scoring functions
- Poorly ranks binding affinities of multiple compounds
- Slow/more expensive



Ligand-based Methods  
QSAR, LB-pharmacophore mapping, CoMFA

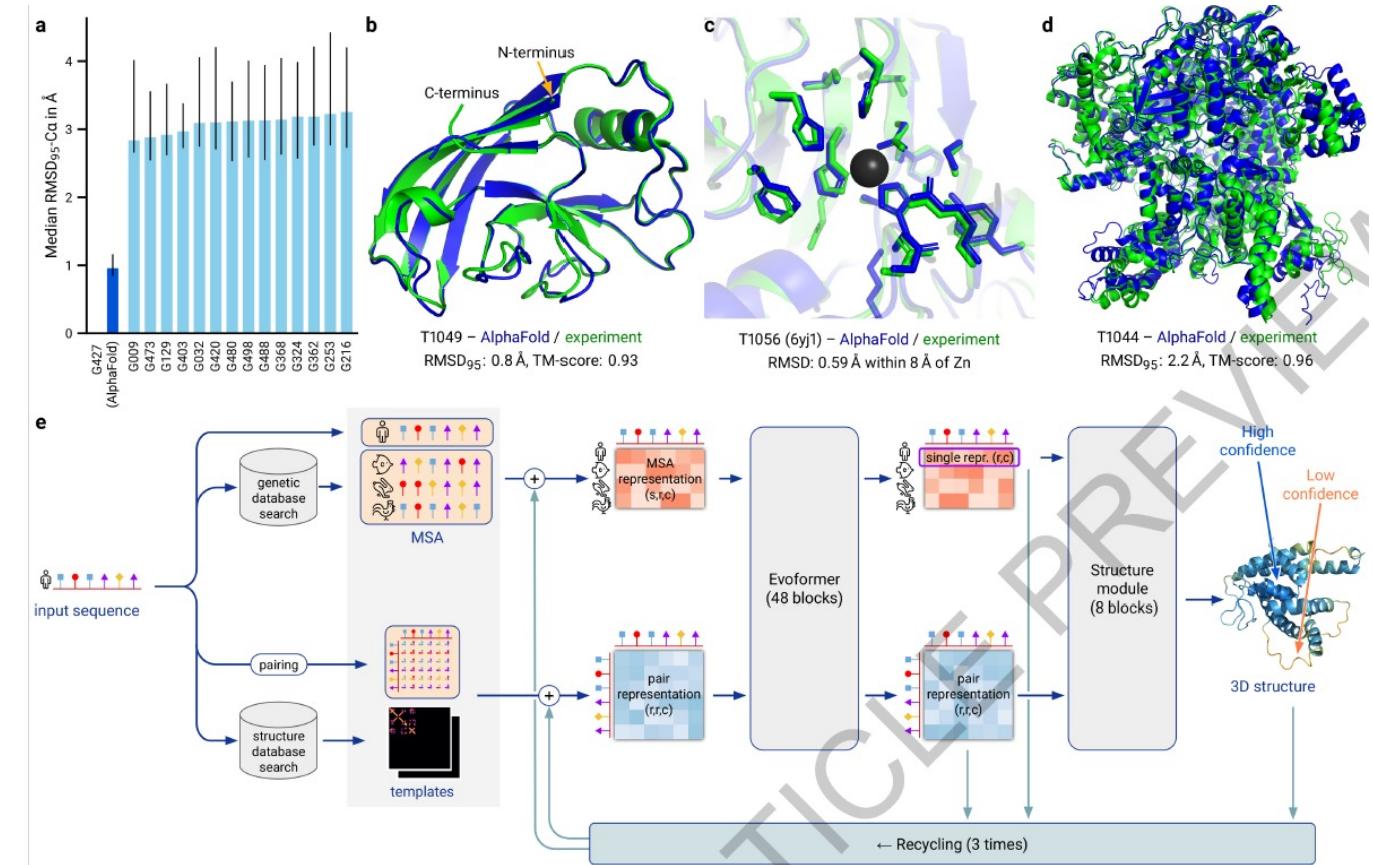


- No receptor pocket information
- Limited to chemical information provided by dataset (scaffold hopping)
- Consistency of experimental data used in reference set
- Ratio of descriptors to actives → over-training

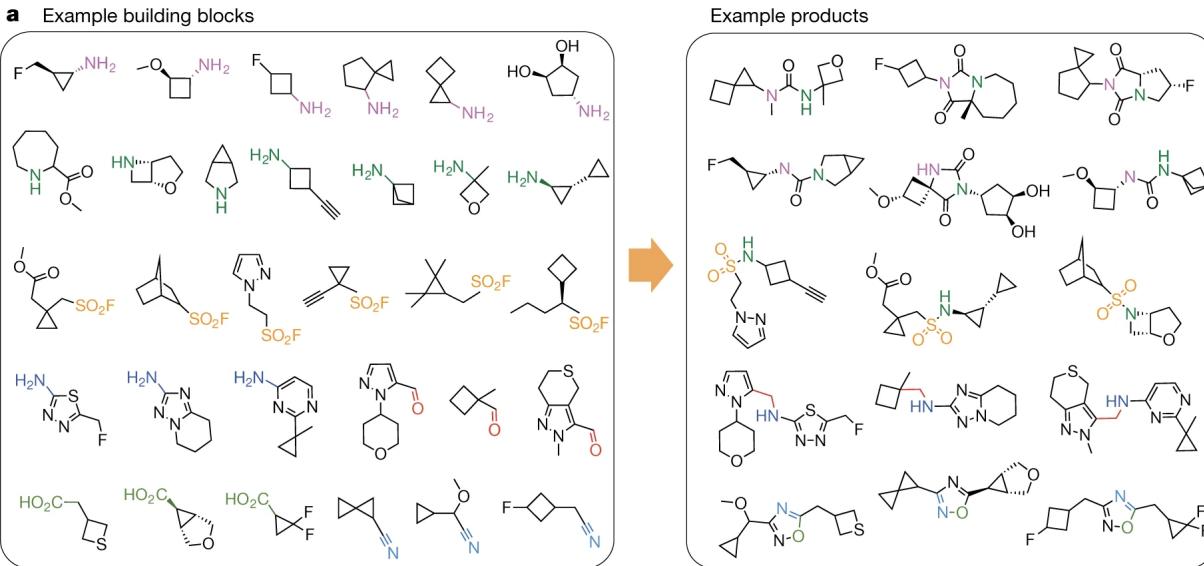


# AlphaFold: Highly accurate protein structure prediction

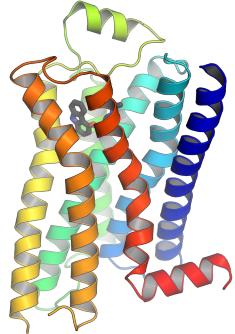
- DeepMind released AlphaFold algorithm and publicly available database of predicted structure of "entire" human proteome and ~20 model organisms
- This was one big data problem "solved" with AI
- Accurate with membrane proteins also



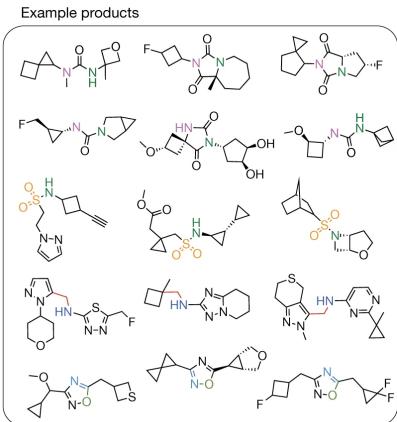
# Virtual synthesis library contains ~1 billion compounds to screen



# Molecular docking of ultra-large libraries identifies binders of AmpC and D<sub>4</sub>



+



## AmpC

Dock 99 Million Compounds



Cluster/select compounds to test



44/55 synthesized



Hits ranged from  
 $K_i = 1.3\text{-}400 \mu\text{M}$



Optimized to 77 nM

## D<sub>4</sub> dopamine receptor

Dock 138 Million Compounds



Cluster/select compounds to test



549/589 synthesized



81 hits  
 $K_i < 10 \mu\text{M}$

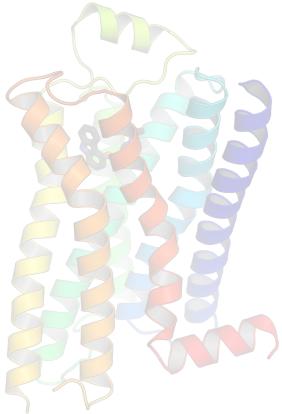


180 pM hit

There are efforts to incorporate ML into these structure-based methods

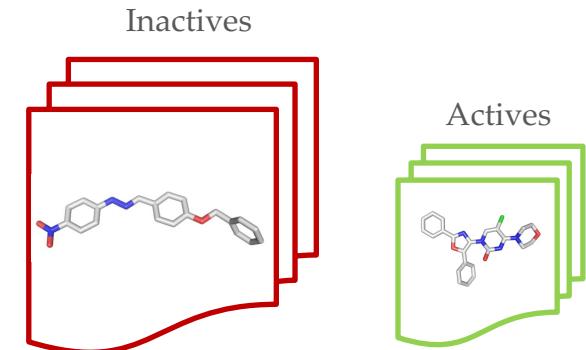
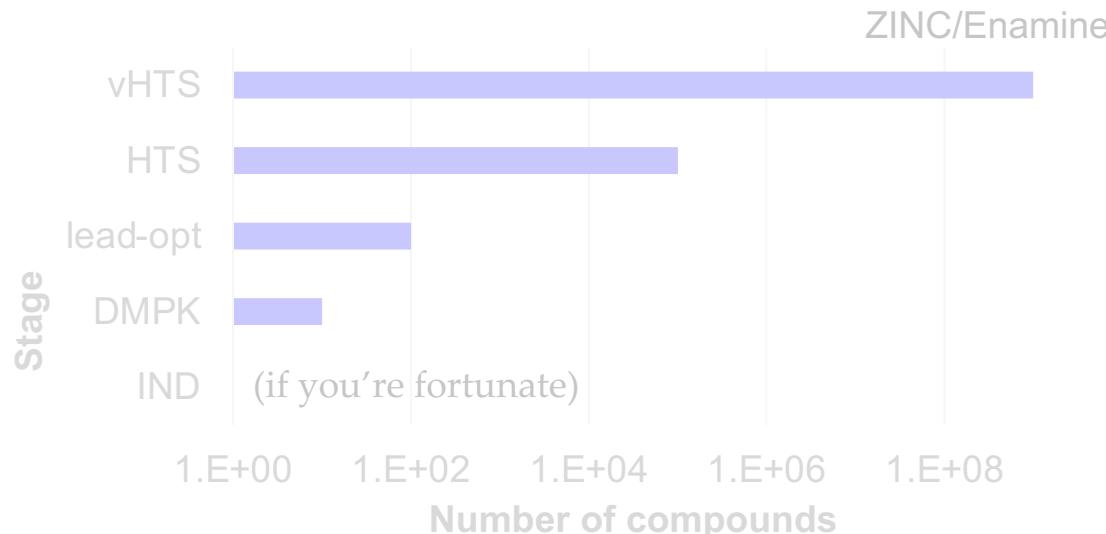
- Learning contact prediction
- de novo ligand design
- MD speed-up with ML

# Computer-aided drug discovery from a distance

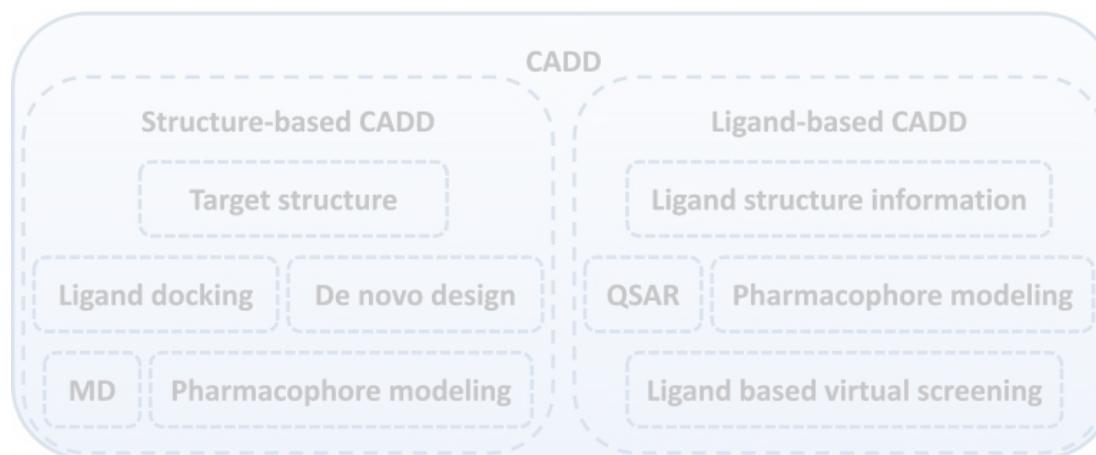


Structure-based Methods  
Docking, Molecular Dynamics, SB-pharmacophore mapping

- Need high-resolution structure
- Understand mode of binding, key residues
- Rely on scoring functions
- Poorly ranks binding affinities of multiple compounds
- Slow/more expensive



Ligand-based Methods  
QSAR, LB-pharmacophore mapping, CoMFA



- No receptor pocket information
- Limited to chemical information provided by dataset (scaffold hopping)
- Consistency of experimental data used in reference set
- Ratio of descriptors to actives → over-training



# Journal of Medicinal Chemistry

© Copyright 1964 by the American Chemical Society

VOLUME 7, NUMBER 4

JULY 6, 1964

## What is the Quantitative SAR?

response = average + effect of R<sub>1</sub> substituent + effect of R<sub>2</sub> substituent

$$LD_{50} = \mu + a[H] + a[CH_3] + b[N(CH_3)_2] + b[N(C_2H_5)_2]$$

where  $\mu$  = over-all average

$a[H]$  = contribution of H substituent at position R<sub>1</sub>

$a[CH_3]$  = contribution of CH<sub>3</sub> substituent at position R<sub>1</sub>

$b[N(CH_3)_2]$  = contribution of N(CH<sub>3</sub>)<sub>2</sub> substituent at position R<sub>2</sub>

$b[N(C_2H_5)_2]$  = contribution of N(C<sub>2</sub>H<sub>5</sub>)<sub>2</sub> substituent at position R<sub>2</sub>



Assuming

1. R<sub>1</sub> and R<sub>2</sub> are independent
2. Symmetric contributions to R group total

Multiple Linear Regression problem

- R<sub>1</sub> = H contributes +0.245 to LD<sub>50</sub>
- R<sub>1</sub> = CH<sub>3</sub> contributes -0.245 to LD<sub>50</sub>
- R<sub>2</sub> = N(CH<sub>3</sub>)<sub>2</sub> contributes +0.425 to LD<sub>50</sub>
- R<sub>2</sub> = N(C<sub>2</sub>H<sub>5</sub>)<sub>2</sub> contributes -0.425 to LD<sub>50</sub>

	R <sub>2</sub>	R <sub>1</sub>		LD <sub>50</sub> values (mg/10g)
		H	CH <sub>3</sub>	
	N(CH <sub>3</sub> ) <sub>2</sub>	2.13	1.64	
	N(C <sub>2</sub> H <sub>5</sub> ) <sub>2</sub>	1.28	0.85	

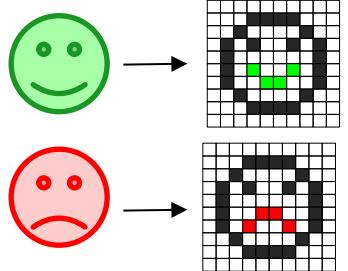
## What is the SAR (qualitatively)?

- R<sub>1</sub>: larger group contributes to lower LD<sub>50</sub> values
- R<sub>2</sub>: larger group contributes to lower LD<sub>50</sub> values
- Assuming R<sub>1</sub> and R<sub>2</sub> are independent

"The suggested mathematical models do not compensate for the three dimensionality of compounds, pH, pK<sub>a</sub>, or other similar physical properties. Perhaps, in time, these can be built into the models for better estimation."



# Chemical “descriptors” allow us to describe models with numeric values



- 1-D (scalar) descriptors**
- MW
  - # of C/N/O/H atoms
  - Formal charge

## 2-D descriptors

- TPSA
- # of HBA/HBD
- # of rotatable bonds
- logP
- # of rings (by ring size)
- 2-D autocorrelations:  
Property distribution by bond distance

## 3-D descriptors

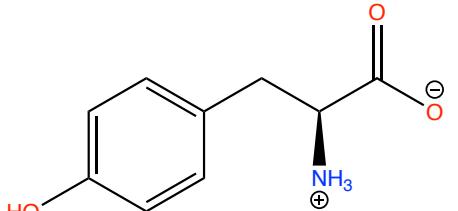
- Volume
- Density
- Radius of gyration
- Principal moment of inertia
- 3-D autocorrelations:  
Property distribution by Euclidian distance

## 4-D descriptors

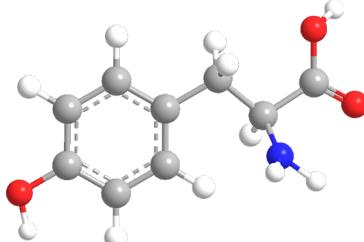
- Determined by conformational ensemble of molecule
- Grid occupancy of lower-dimensional descriptors



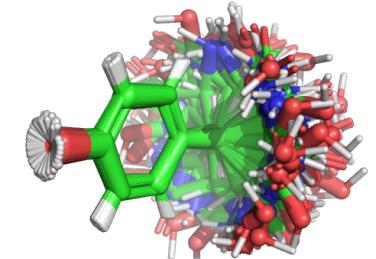
Number of atom types



Connectivity (+ stereochemistry)



Connectivity with spatial info

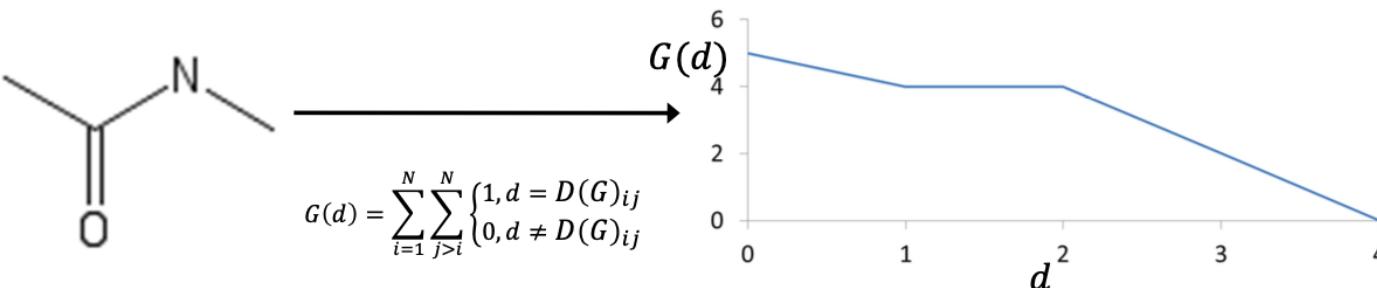


Dynamic information



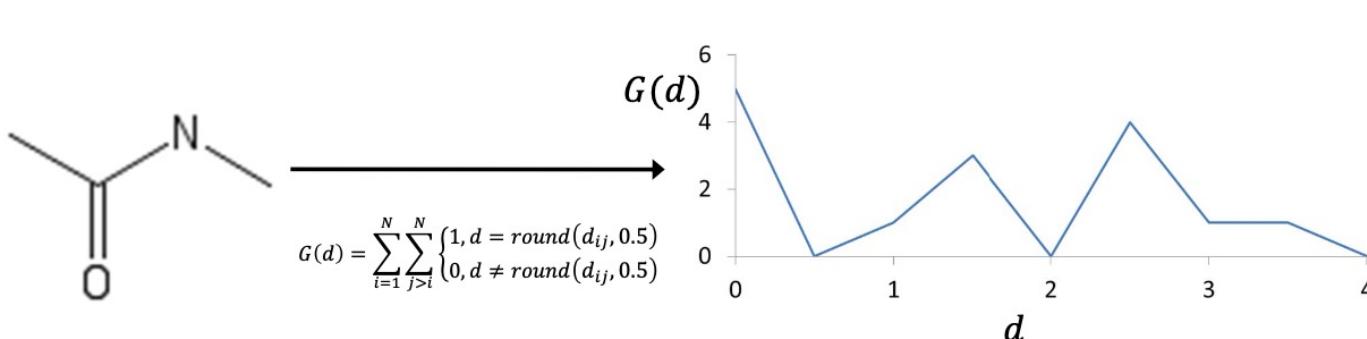
# 2D and 3D autocorrelation descriptors

Distribution of number of bonds



	C	C	N	C	O
C	0	1	2	3	2
C	1	0	1	2	1
N	2	1	0	1	2
C	3	2	1	0	3
O	2	1	2	3	0

Distribution of atom pair distances



	C	C	N	C	O
C	0.00	1.54	2.5*	3.5*	2.5*
C	1.54	0.00	1.47	2.5*	1.23
N	2.5*	1.47	0.00	1.47	2.5*
C	3.5*	2.5*	1.47	0.00	3.0*
O	2.5*	1.23	2.5*	3.0*	0.00

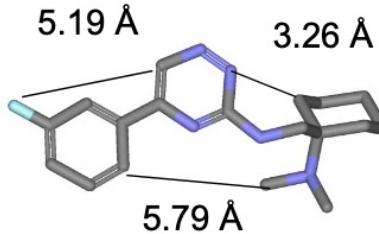
\* - estimates



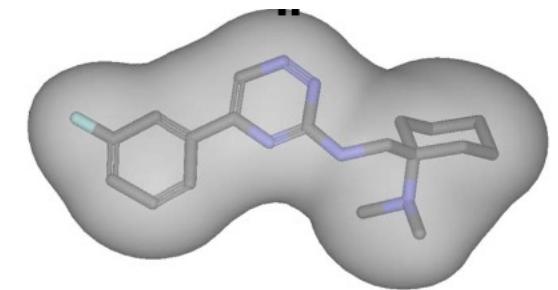
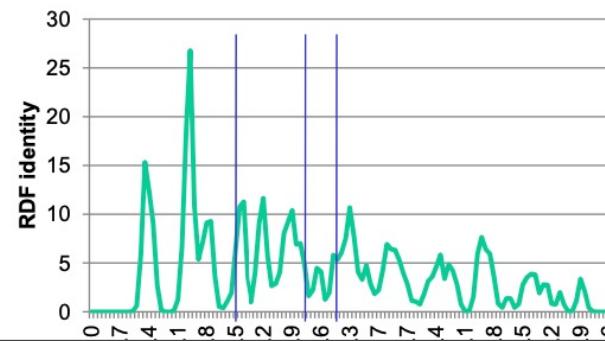
# Radial distribution functions describe 3D distribution of descriptors

What is the 3D distribution of...

Atoms?



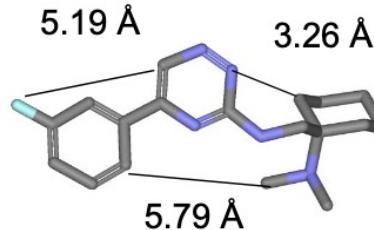
$$g(d) = \sum_{\text{atom pairs}} e^{-B(d-d_{ij})^2}$$



Lone pair electrons?

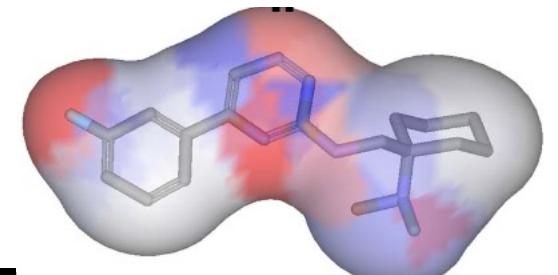
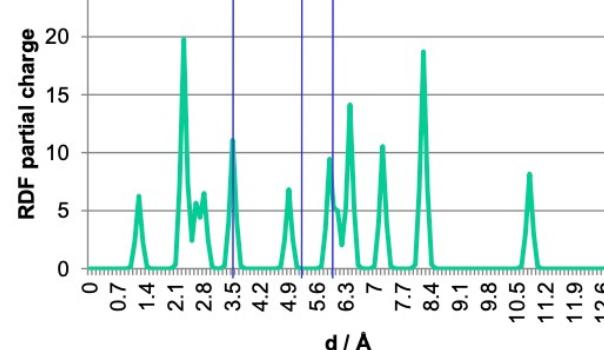
- Bottom line:
1. These descriptors can get complicated
  2. These are really useful descriptors...why do you think that is?

Partial charges?

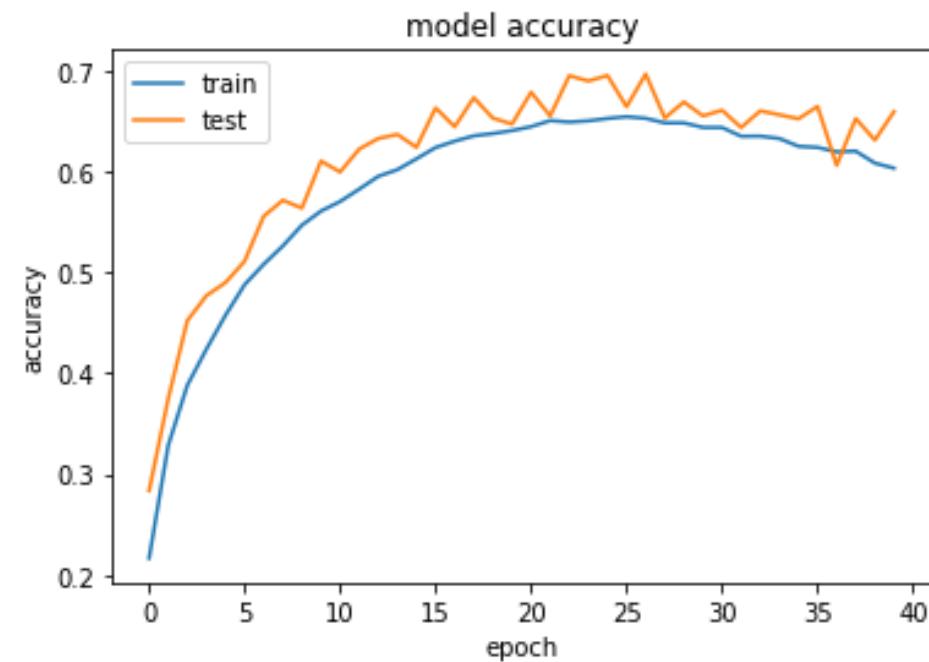
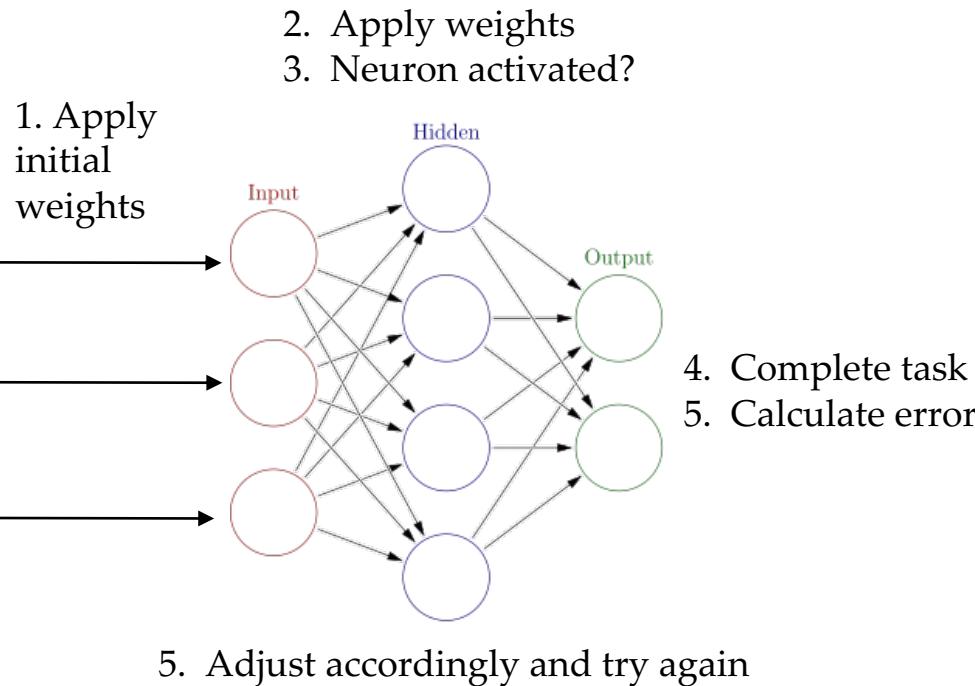
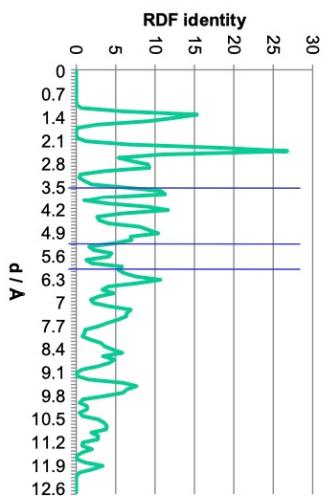
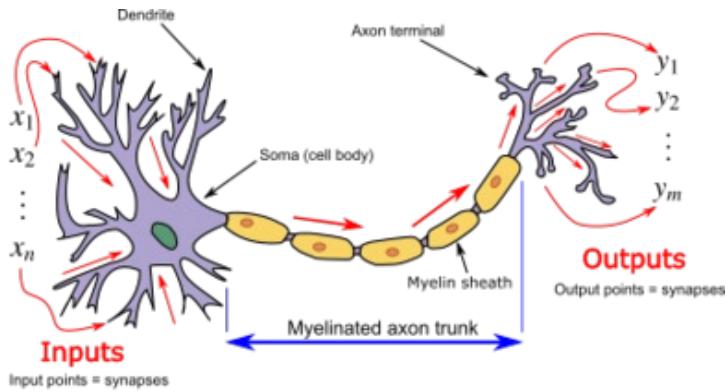


Gilson (2003)

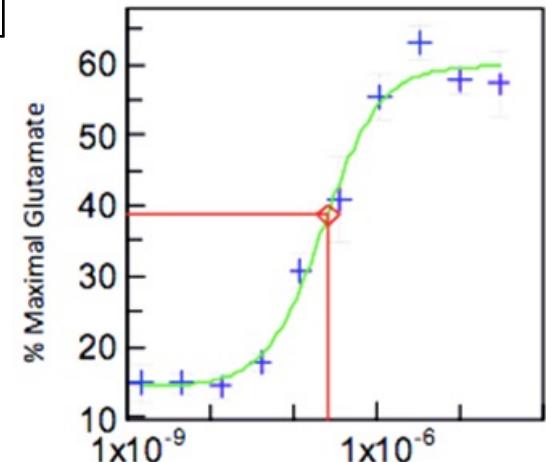
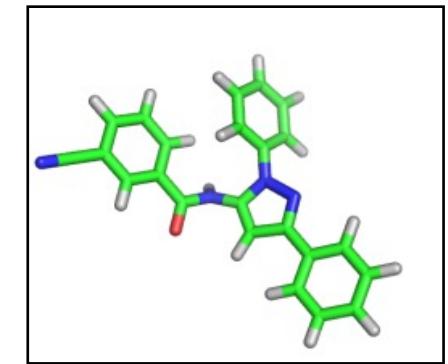
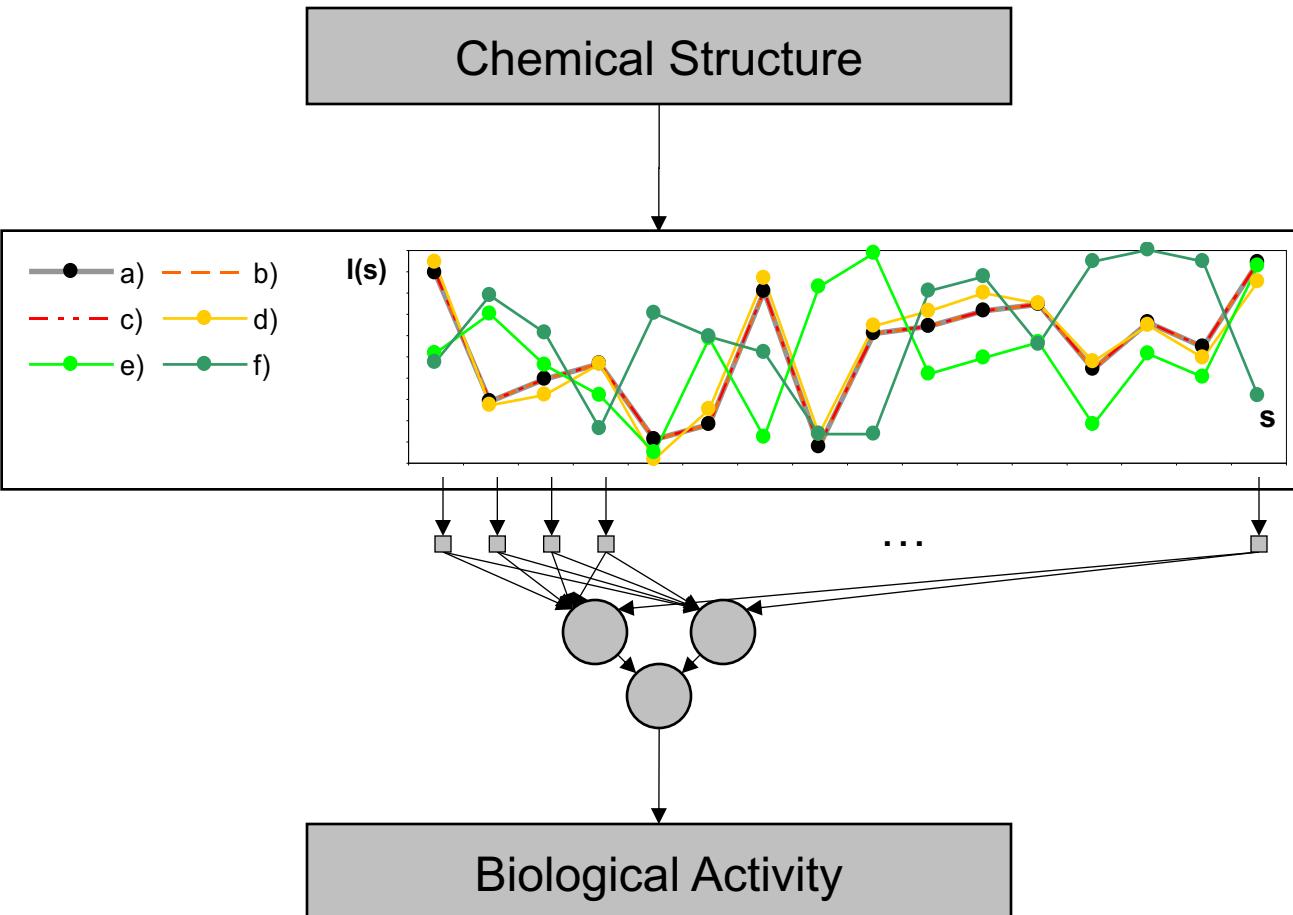
$$g(d) = \sum_{\substack{\text{atom pairs} \\ i,j}} A_i A_j e^{-B(d-d_{ij})^2}$$



# Artificial Neural Networks (ANNs)



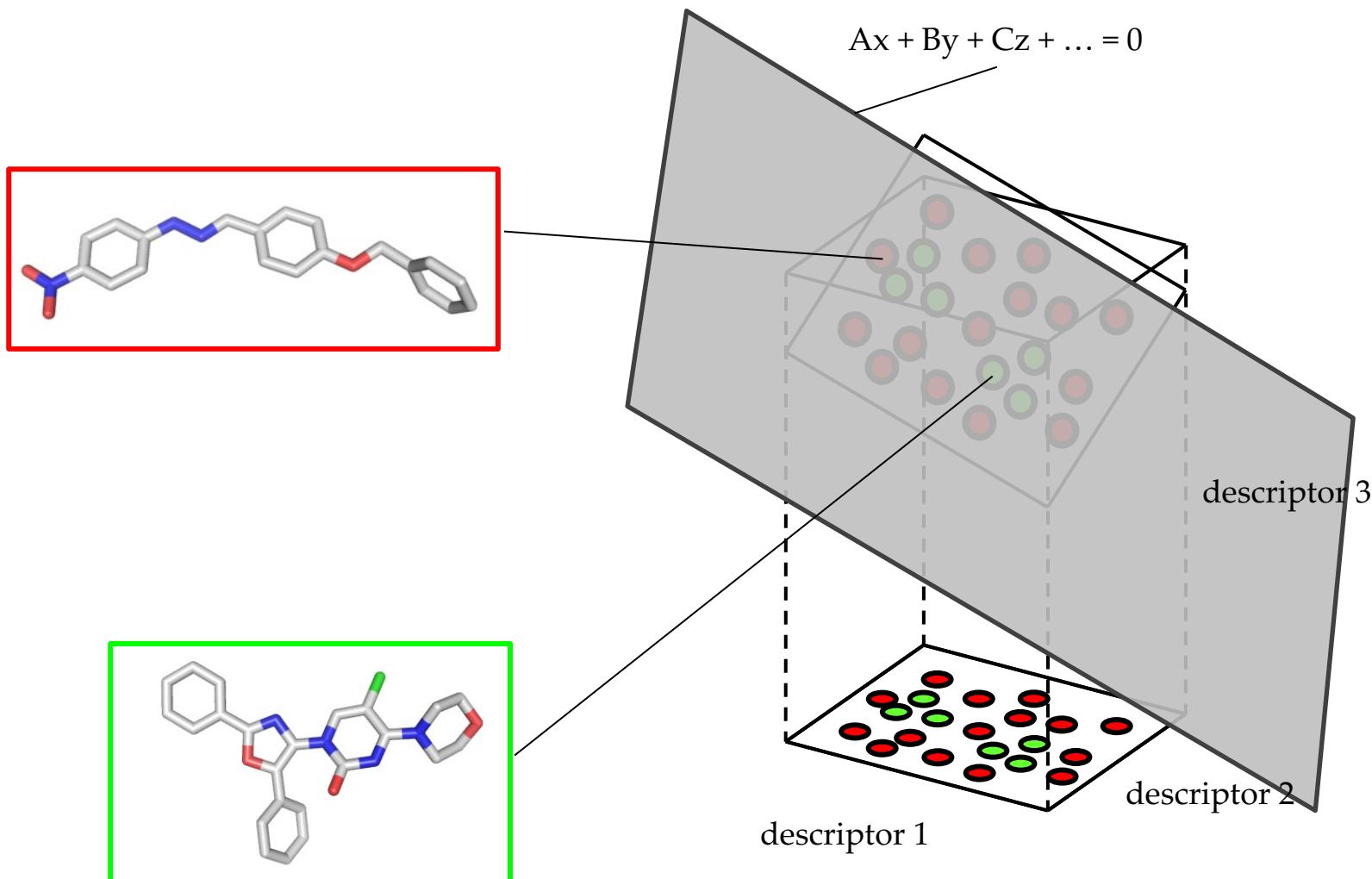
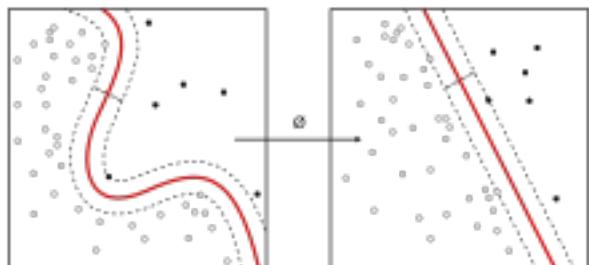
# Artificial Neural Networks (ANNs) predicted biological activity from compound descriptors



# Mapping descriptor space into hyperspace

Which descriptors, and their weights, are important to separate actives and inactives?

- Classic support vector machine (SVM) problem
- Works better with small sample sizes
- Easier interpretability of results: you can figure out most important features



# Optimizing descriptors for given test set

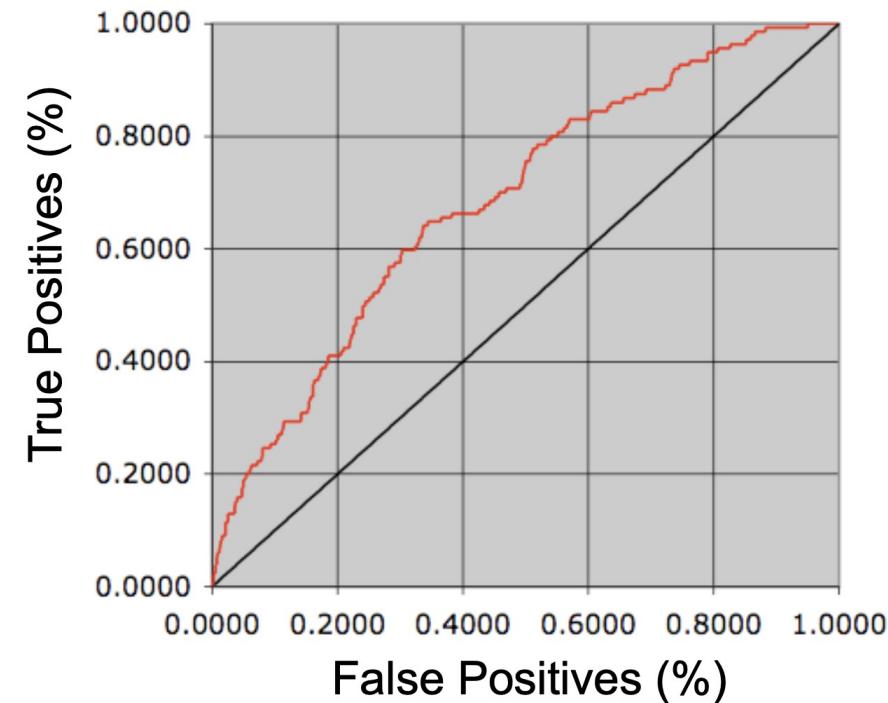
The more generalizable, the better! Less over-fitting.

Number descriptors

| 8 |

- Molecular Weight
- Number H bond donors
- Number H bond acceptors
- XlogP
- Polar surface area
- Mean molecular polarizability
- Molecular dipole moment
- Aqueous solubility

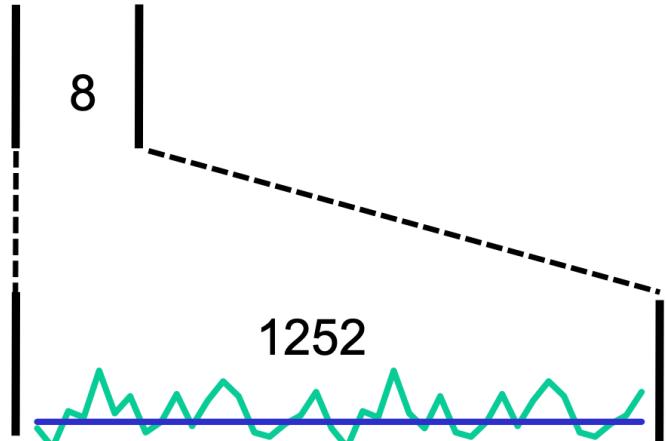
vHTS Training Optimization (ROC curves)



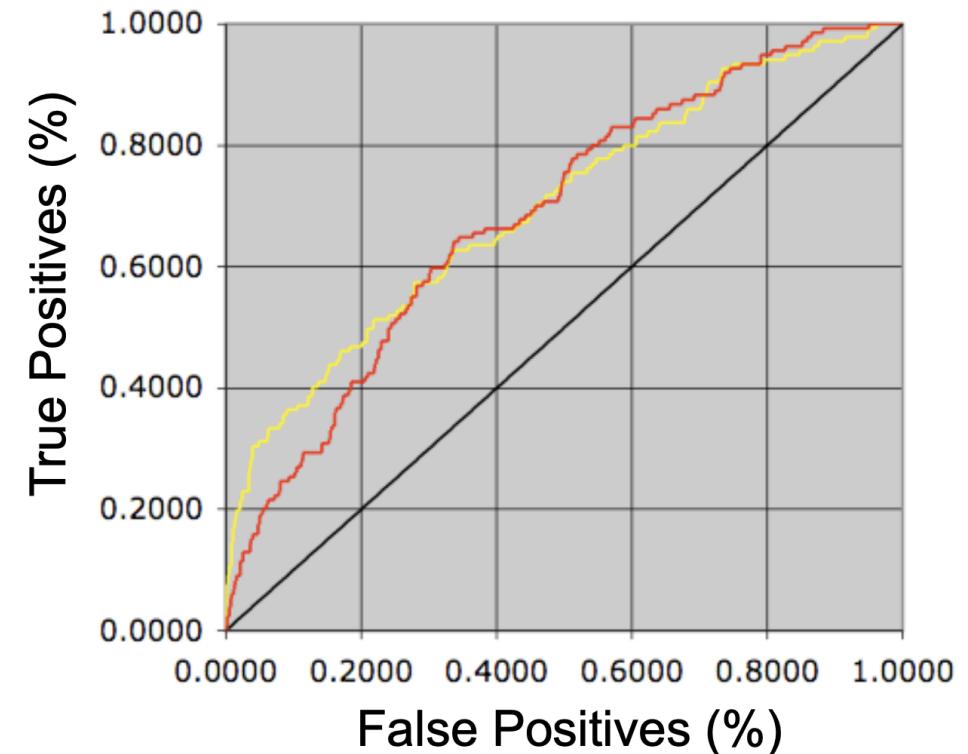
# Optimizing descriptors for given test set

The more generalizable, the better! Less over-fitting.

Number descriptors



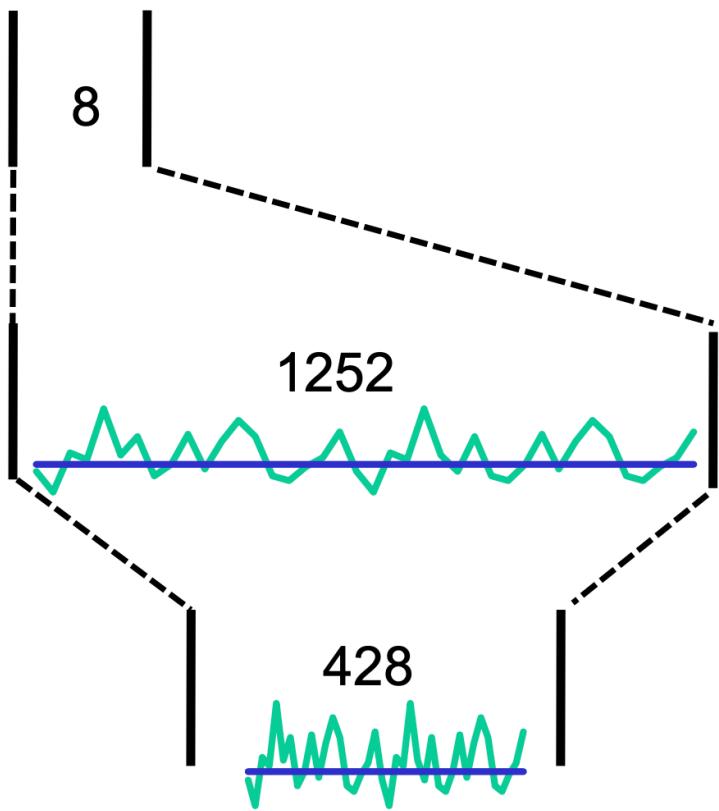
vHTS Training Optimization (ROC curves)



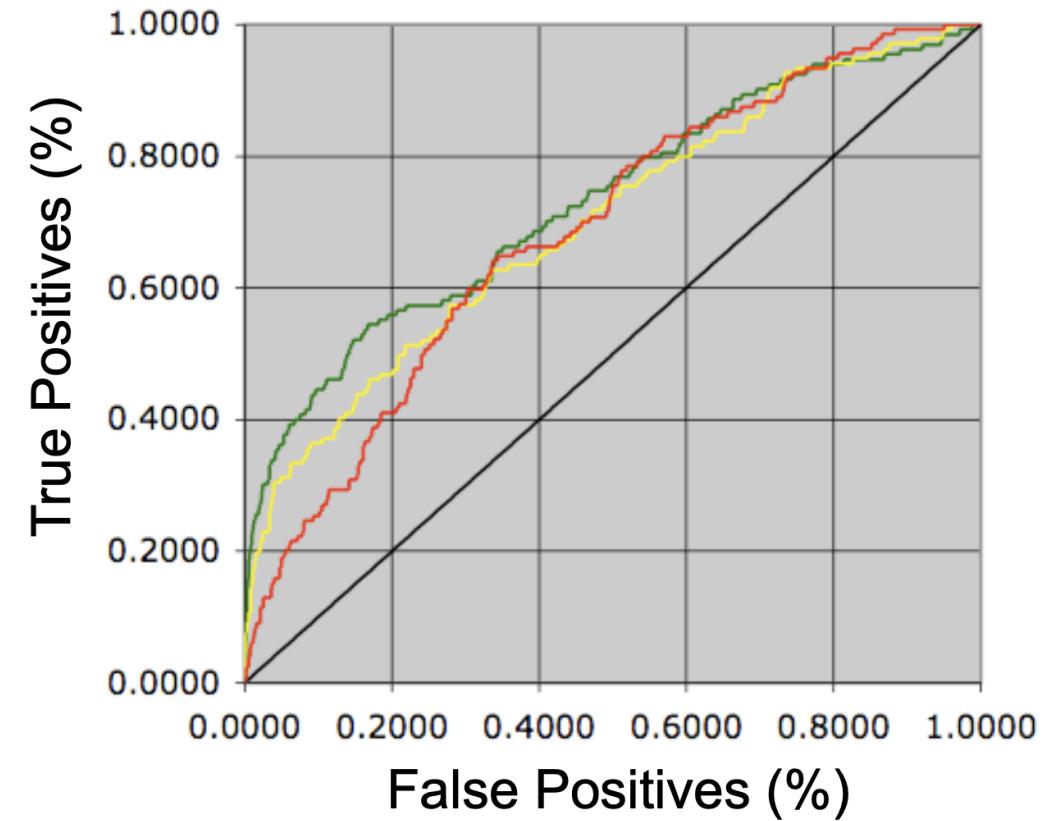
# Optimizing descriptors for given test set

The more generalizable, the better! Less over-fitting.

Number descriptors



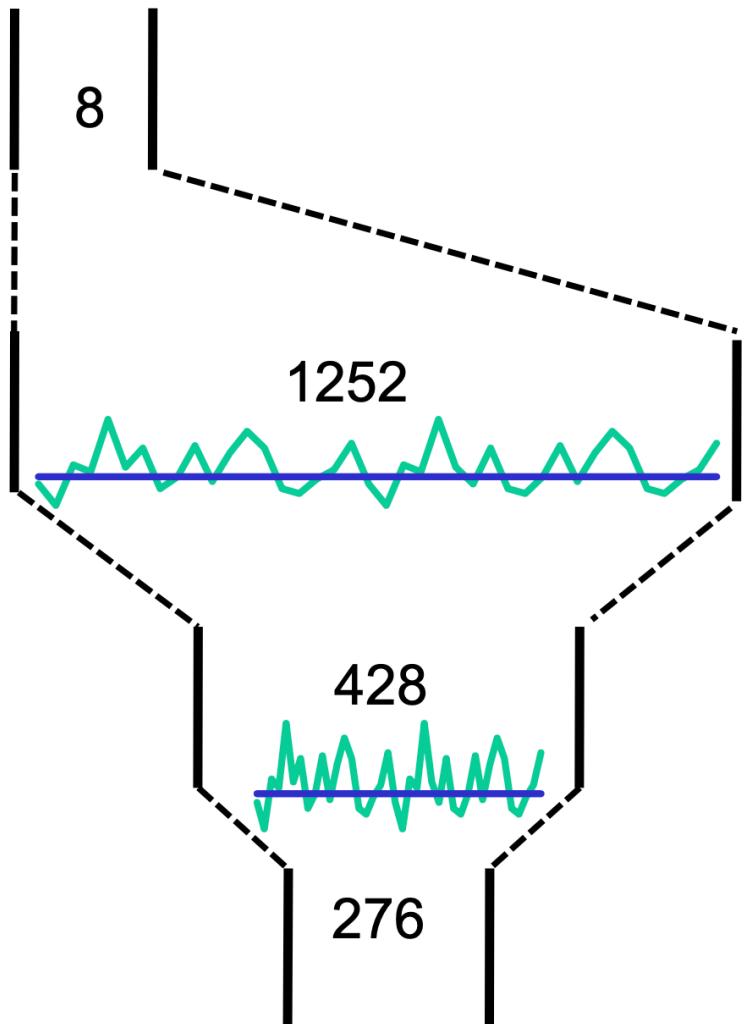
vHTS Training Optimization (ROC curves)



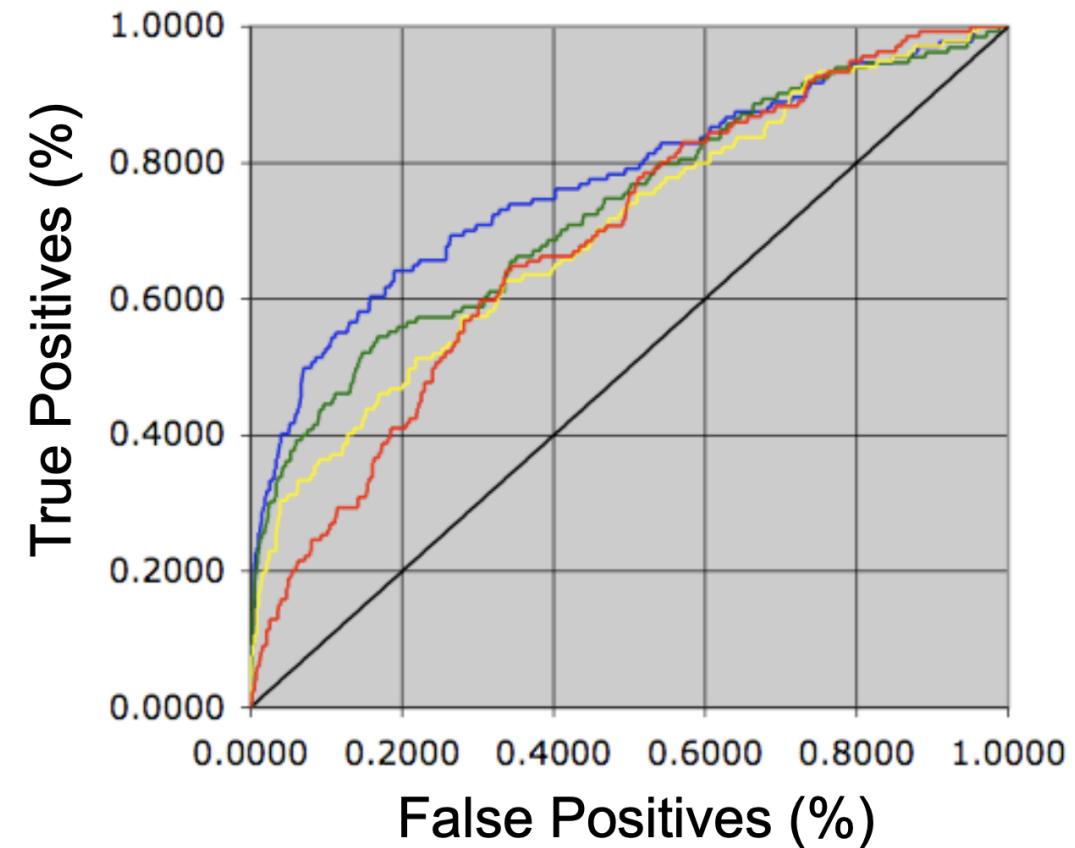
# Optimizing descriptors for given test set

The more generalizable, the better! Less over-fitting.

Number descriptors

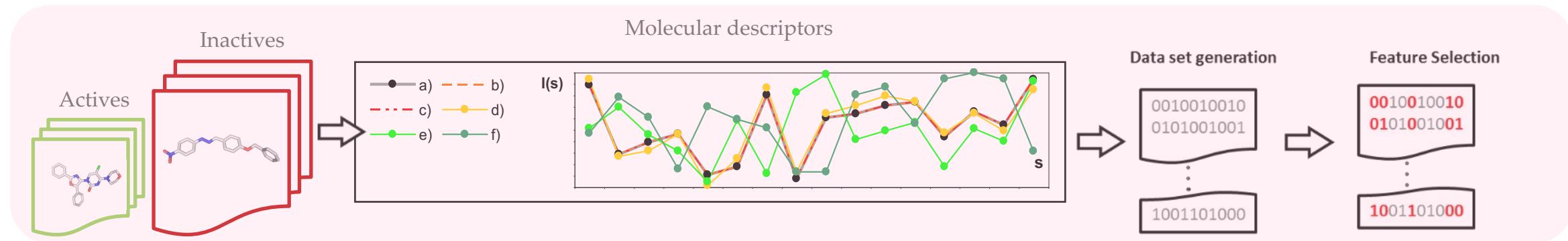


vHTS Training Optimization (ROC curves)

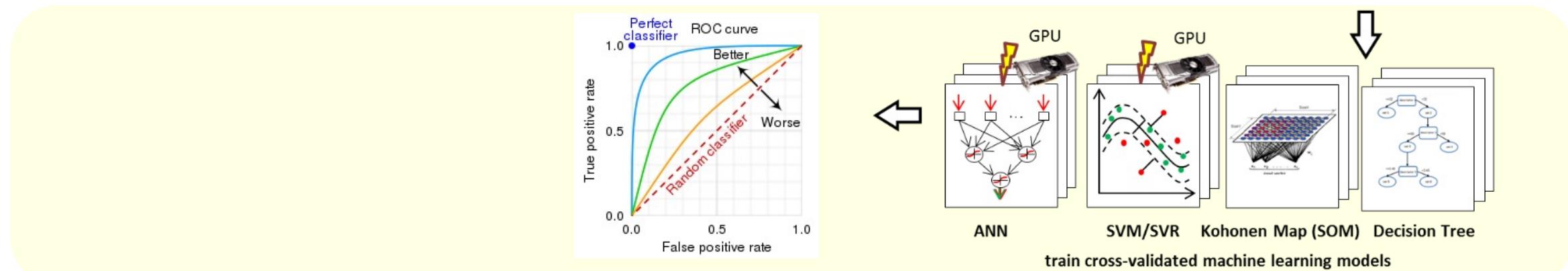


# Putting all this together → LB-vHTS pipeline

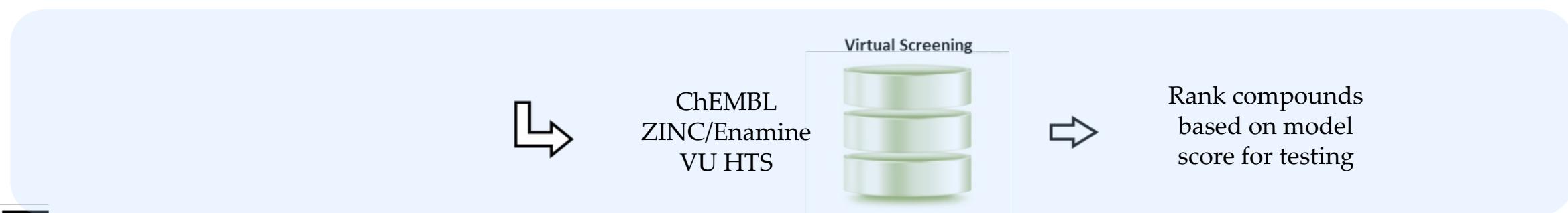
Dataset acquisition



Training and validation



Prediction



# ANN QSAR identifies mGluR<sub>5</sub> allosteric modulators

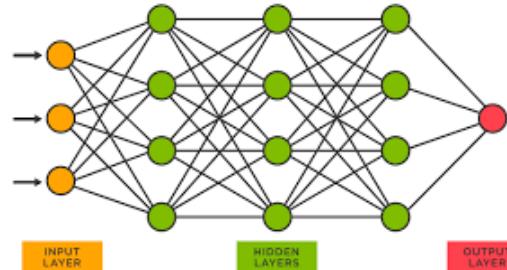
150,000 compounds were experimentally tested for allosteric modulation of mGluR<sub>5</sub> measuring receptor-induced intracellular release of calcium



- 1,387 (0.94%) compounds were verified as mGluR<sub>5</sub> PAMs
- 345 (0.23%) compounds were verified as mGluR<sub>5</sub> NAMs

Niswender, C. M.; Johnson, K. A.; Luo, Q.; Ayala, J. E.; Kim, C.; Conn, P. J.; Weaver, C. D. *Mol Pharmacol* 2008, 73, 1213-24.

~750,000 compounds (ChemBridge) virtually screened with ANN QSAR for allosteric modulation of mGluR<sub>5</sub>



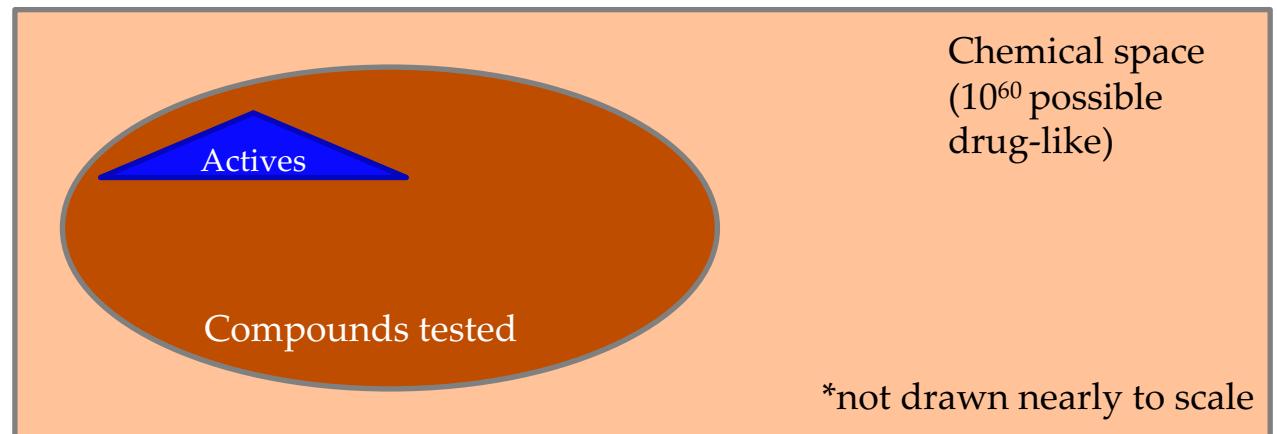
- 749 compounds predicted EC50 < 10uM
- 12 novel NAM scaffolds

$$\text{Enrichment} = 3.6\% / 0.23\% = 16$$

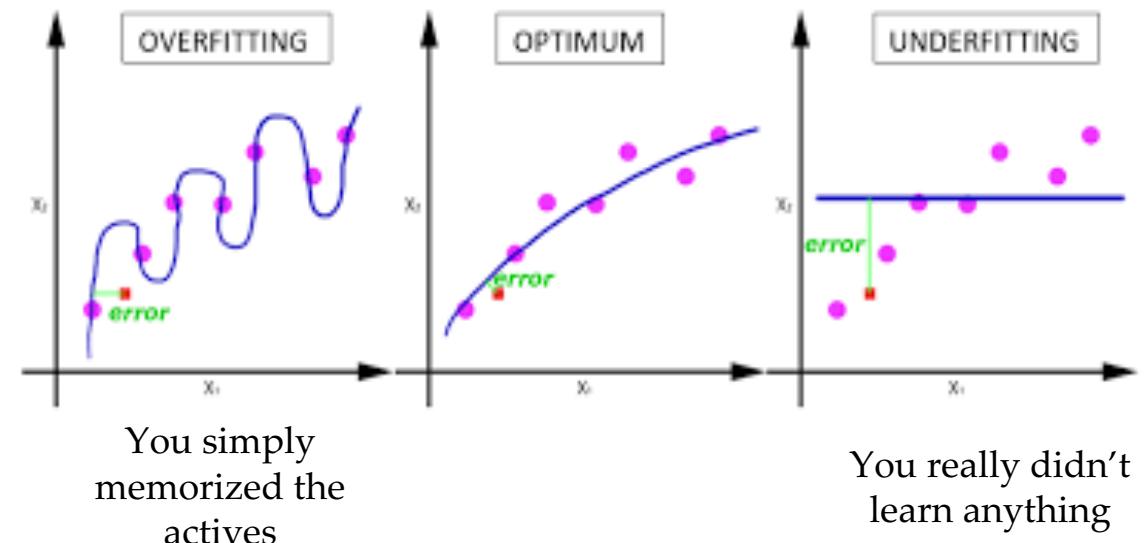
R. Mueller, E. S. Dawson, J. Meiler, A. L. Rodriguez, B. A. Chauder, B. S. Bates, A. S. Felts, J. P. Lamb, U. N. Menon, S. B. Jadhav, A. S. Kane, C. K. Jones, K. J. Gregory, C. M. Niswender, P. J. Conn, C. M. Olsen, D. G. Winder, K. A. Emmitte and C. W. Lindsley. *ChemMedChem*; 2012; Vol. 7 (3): p. 406-14.

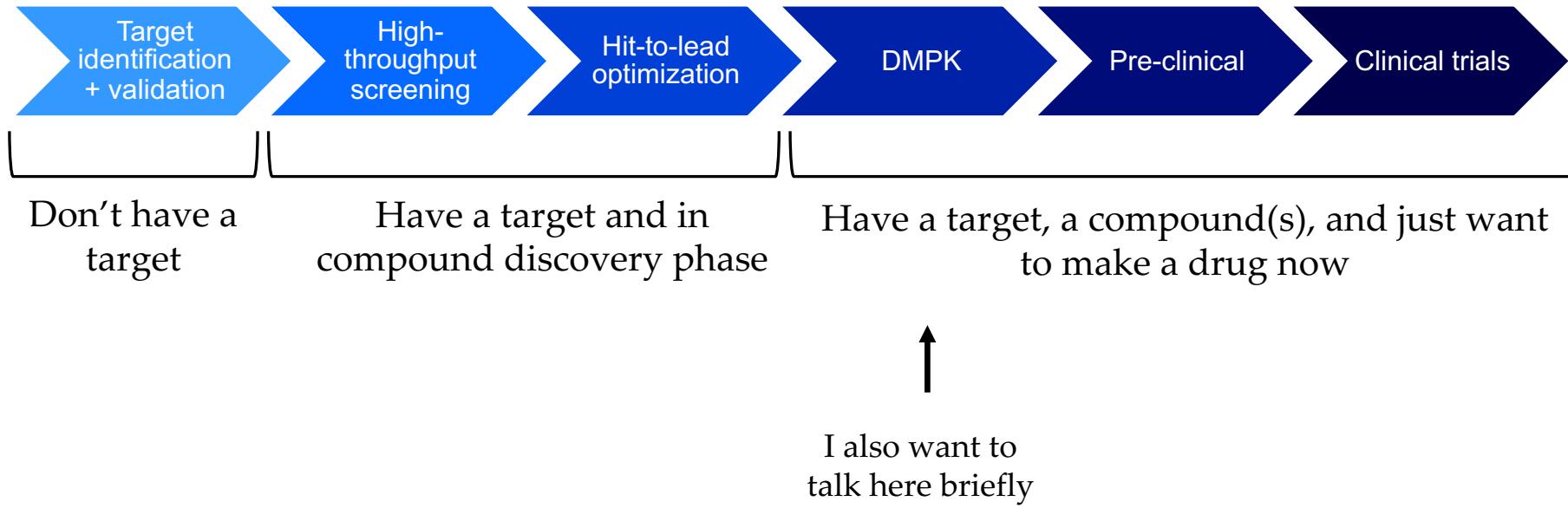
# A couple things to be weary of with QSAR

- Less confidence in predictions for molecules further in chemical space from your test set
  - We have information of only a tiny fraction of available compound space
  - Nobody likes to report inactive molecules
  - And the actives are usually chemically similar



- Heavily biased dataset composition
  - Many descriptors and only a few outcomes
  - Over-fitting if not careful
- Can be a black-box and tough to interpret





# “It’s relatively easy to discovery a potent ligand – it’s damn tough to discover a drug” – E.H. Cordes

90% of drug candidates that get IND fail in CTs

Causes for attrition	
Animal tox/safety	20%
ADME	41%
Lack of efficacy	27%
Marketing consideration	7%
Formulation	4%
Miscellaneous	20%

	Antibiotic	CNS	Cardiovascular	Respiratory
Failure due to <b>safety</b> (Phase 1)	85-100%	20-30%	30-40%	50%
Failure due to <b>efficacy</b> (Phase 2/3)	5-15%	70%	60%	50%

Failure is industry-wide

Failure is indication-wide

Reasons for the failures vary between indications, but NOT between companies

**What do we use now?**

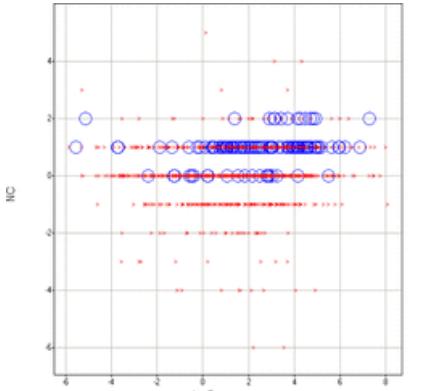
Lipinski's Rule of 5 (sort of)

Can we figure out which molecular features don't pass tox?

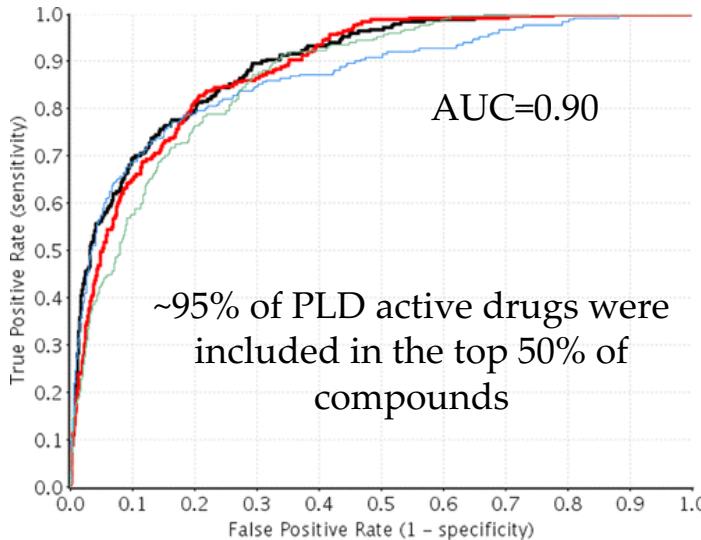
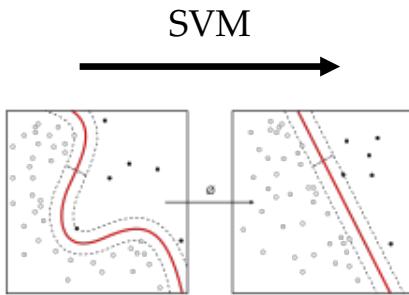


# Drug-induced phospholipidosis (PLD) prediction using SVMs

- PLD is a phospholipid storage disorder characterized by excessive accumulation of intracellular phospholipids in a variety of tissue types
- Many drugs can accumulate in tissues by forming complexes with the polar phospholipids → “drug-induced PLD”
- Potentially toxic and can halt drug development in later stages
- **Can we identify these compounds beforehand? Let's fail earlier!**



Commonly-used descriptors  
for this task do not separate  
PLD inducing drugs

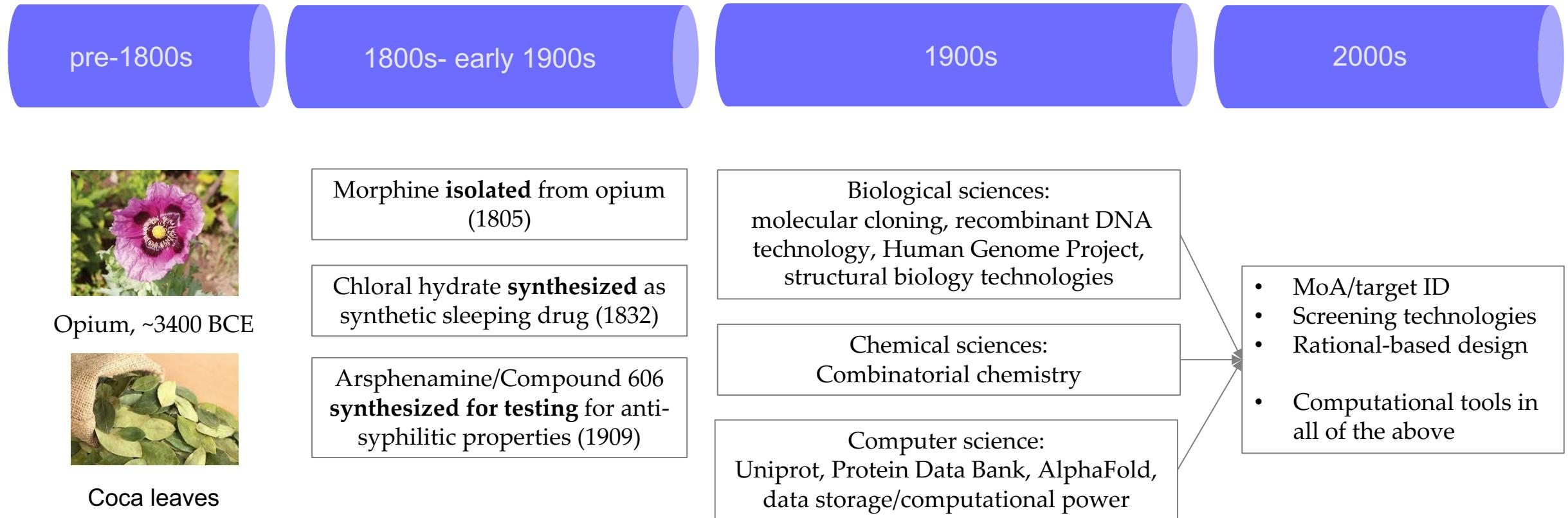


## Most important features:

1. N16 (a positively charged nitrogen atom in a saturated ring)
2. C3 (an aromatic carbon atom with no substitution and adjacent to two aromatic carbon atoms)
3. H2 (a hydrogen bonded to an aromatic carbon)
4. M12 (a number of aromatic rings)
5. S5 (sulfur in a ring bonded to two aromatic atoms)

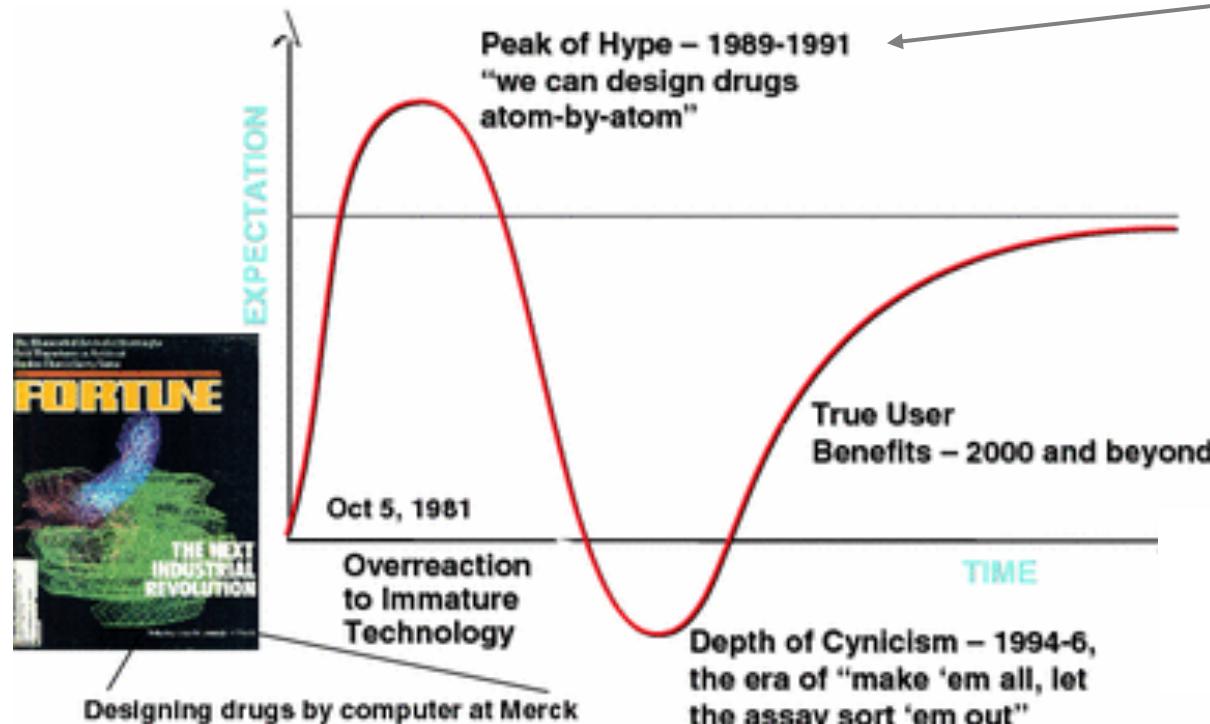


# Drug discovery: a crude history



# "Computer-aided drug design: the next 20 years" (2007)

Paper based on talk given at 2007 ACS meeting:  
"The aim in this Perspectives is to provoke  
some thinking about what may lie ahead."



Peter Goodford (1989)

"We must...

- F 1. deal properly with water
- C 2. remember that conformation depends on structure and environment
- B 3. combine theory and experiment
- D 4. predict solubility
- A 5. improve homology modeling"

- (mid-1990s) HIV protease inhibitors via SB-CADD
- (1995): de novo ligand design of binders to dihydrofolate reductase and thymidylate synthase
- (1999) tirofiban by pharmacophore vHTS



# This is becoming more accessible

- Python sklearn is very well-documented and freely accessible—go try it!



# Dr. Fiona Marshall Visit

## **Merck/MSD**

Senior VP Head of Discovery, Preclinical and  
Translational Medicine

2018-2021:

- Senior VP Head of Discovery, Preclinical & Translational Medicine
- VP Head of Neuroscience and Head of MSD UK Discovery Research
- VP Head of Discovery Research MRL UK

## **Sosei Group**

2015-2018: Executive VP and Chief Scientific Officer

## **Heptares Therapeutics**

2006-2018: Chief Scientific Officer and Founder



**Tuesday, May 17<sup>th</sup> from 2:00-5:00pm**

**CPB Retreat Speaker on May 18<sup>th</sup>  
9:00 in first floor class ESB**

**Please talk to me or email for details  
(shannon.t.smith.1@vanderbilt.edu)**



# Acknowledgements

## Meiler Lab

Jens Meiler

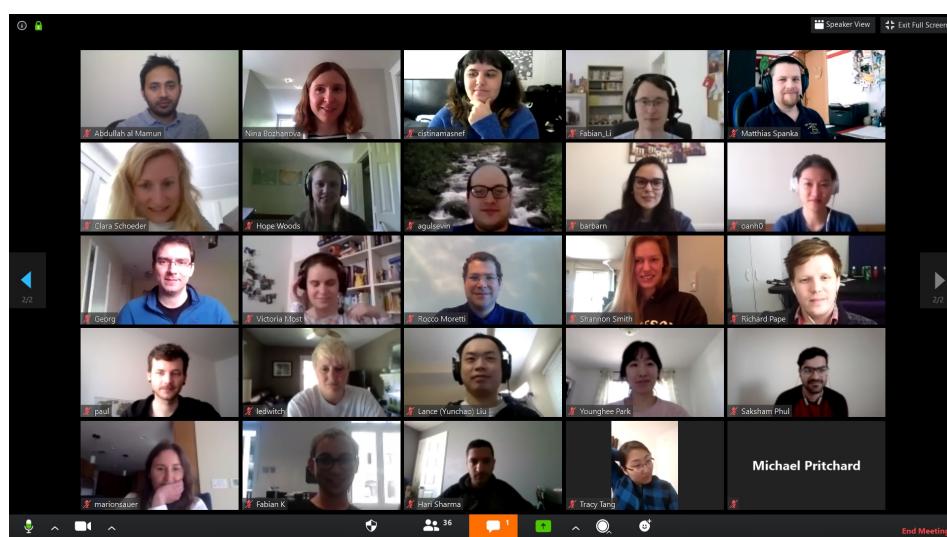
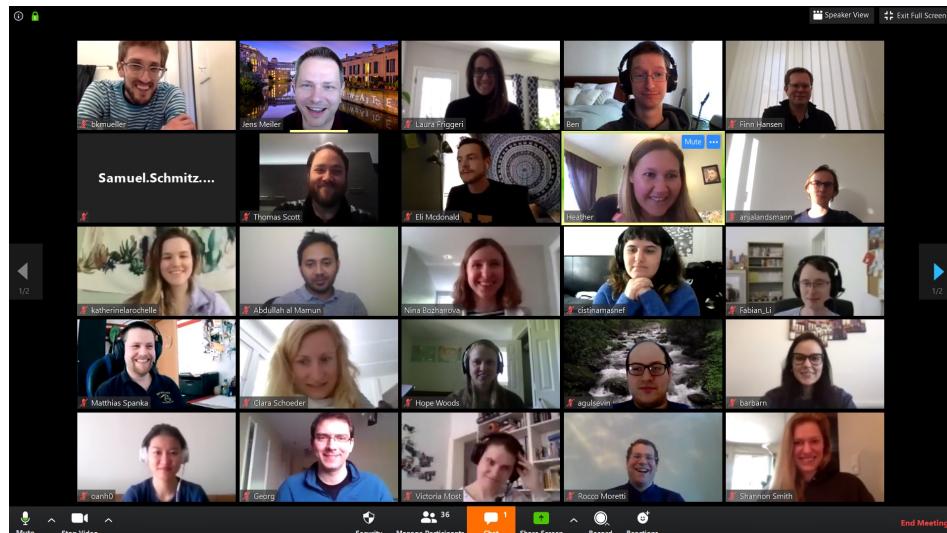
Ben Brown

Rocco Moretti

Oanh Vu

Jeffrey Mendenhall

## Funding



# You can throw the kitchen sink at your problem

- Python sklearn is very well-documented and freely accessible—go try it!

```
from sklearn.neural_network import MLPClassifier  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.svm import SVC  
from sklearn.gaussian_process import GaussianProcessClassifier  
from sklearn.gaussian_process.kernels import RBF  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier  
from sklearn.naive_bayes import GaussianNB  
from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis  
  
Classifiers = [  
    KNeighborsClassifier(3),  
    SVC(kernel="linear", C=0.025),  
    SVC(gamma=2, C=1),  
    GaussianProcessClassifier(1.0 * RBF(1.0)),  
    DecisionTreeClassifier(max_depth=5),  
    RandomForestClassifier(max_depth=5, n_estimators=20, max_features='auto'),  
    MLPClassifier(alpha=1, max_iter=1000),  
    AdaBoostClassifier(),  
    GaussianNB(),  
    QuadraticDiscriminantAnalysis()]
```

