



**BIMM 143**  
**Structural Bioinformatics II**

Lecture 12

**Barry Grant**  
**UC San Diego**

<http://thegrantlab.org/bimm143>

# NEXT UP:

- ▶ **Overview of structural bioinformatics**

- Major motivations, goals and challenges

- ▶ **Fundamentals of protein structure**

- Composition, form, forces and dynamics

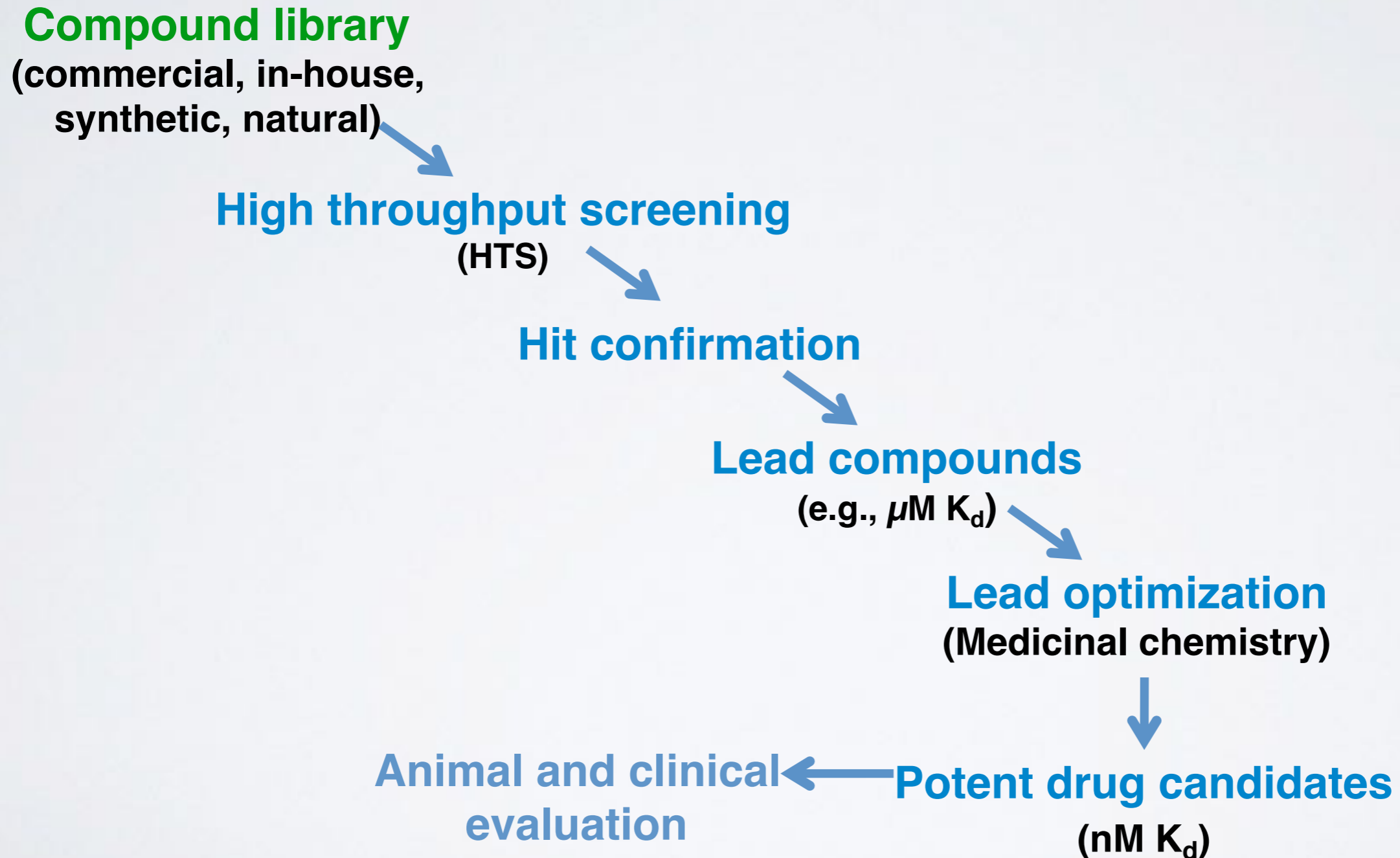
- ▶ **Representing and interpreting protein structure**

- Modeling energy as a function of structure

- ▶ **Example application areas**

- drug discovery & Predicting functional dynamics

# THE TRADITIONAL EMPIRICAL PATH TO DRUG DISCOVERY





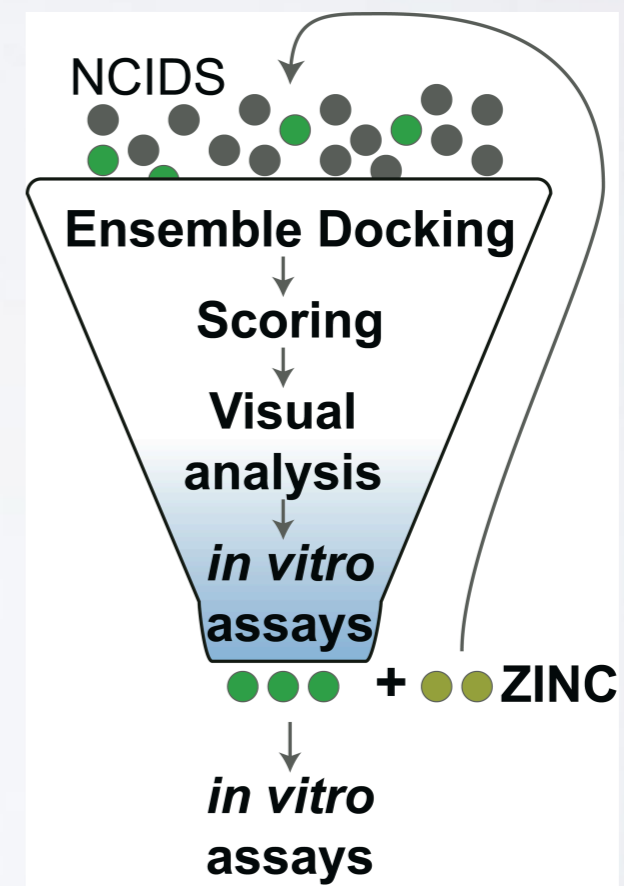
# COMPUTER-AIDED LIGAND DESIGN

Aims to reduce number of compounds synthesized and assayed

Lower costs

Reduce chemical waste

Facilitate faster progress





Two main approaches:

(1). **Receptor/Target-Based**

(2). **Ligand/Drug-Based**

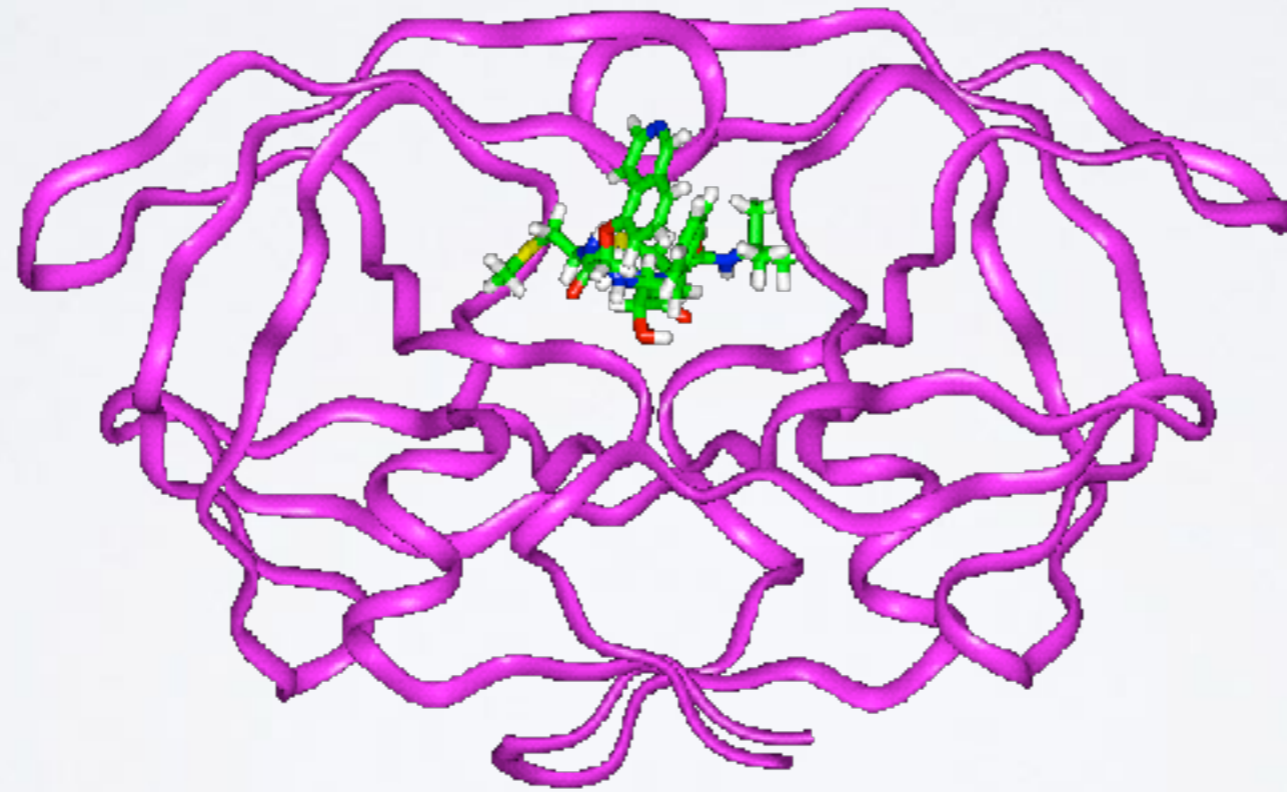
Two main approaches:

**(1). Receptor/Target-Based**

**(2). Ligand/Drug-Based**

# SCENARIO I: RECEPTOR-BASED DRUG DISCOVERY

Structure of Targeted Protein Known: **Structure-Based Drug Discovery**



HIV Protease/KNI-272 complex

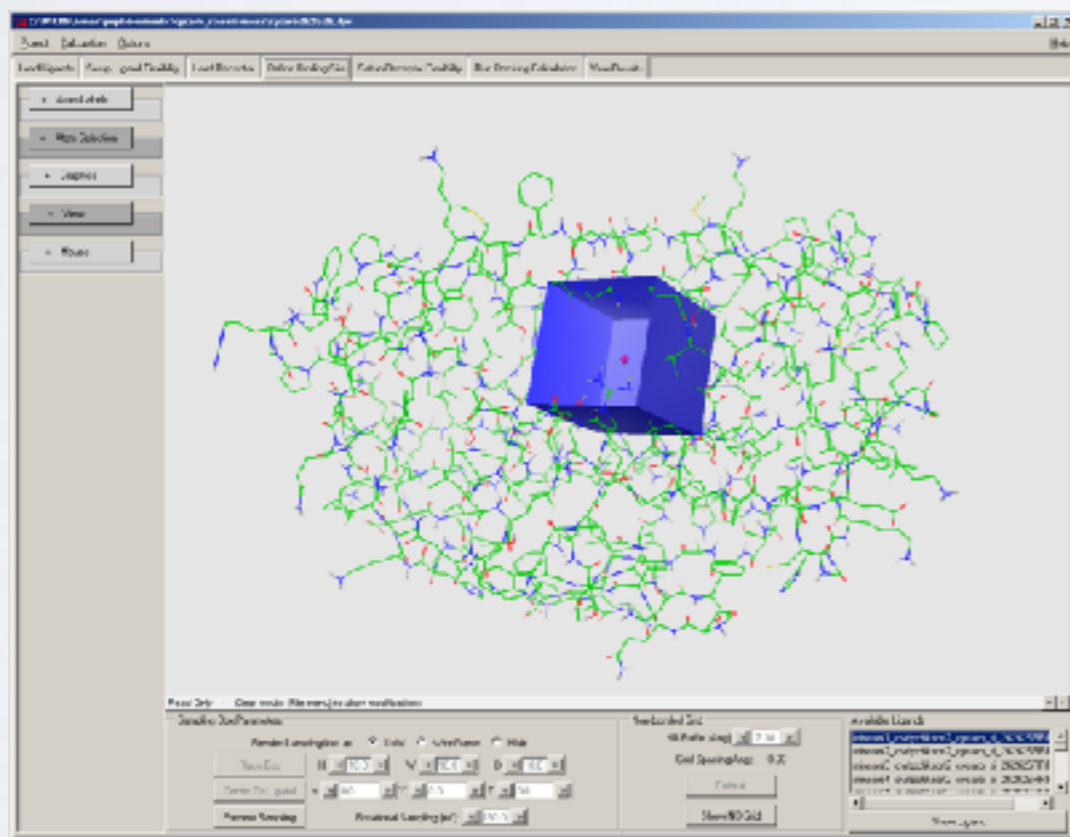


# PROTEIN-LIGAND DOCKING

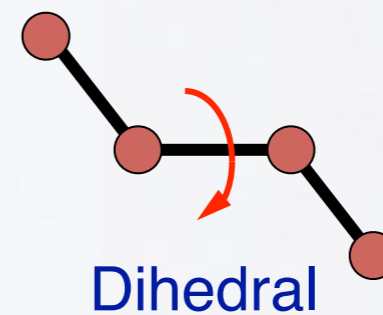
## Structure-Based Ligand Design

### Docking software

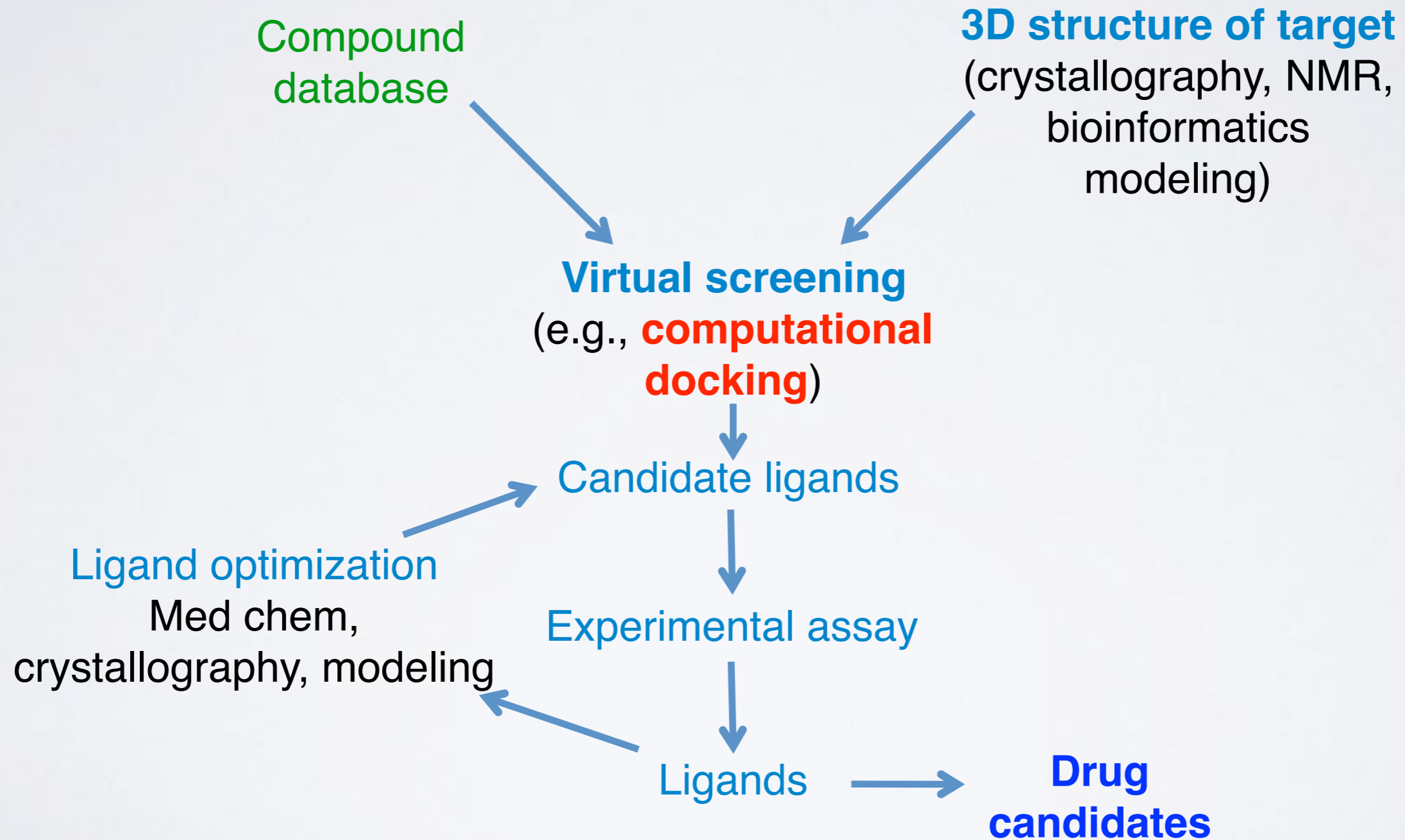
Search for structure of lowest energy



Potential function  
Energy as function of structure



# STRUCTURE-BASED VIRTUAL SCREENING



# COMPOUND LIBRARIES



The screenshot shows the Maybridge HiPlex website. At the top, there is a navigation bar with links for Home, Browse Stock, Shipping Options, About, and Contact Us. Below this is a search bar and a sidebar with navigation options like Home, About Us, and Contact Us. The main content area features the heading "Maybridge HiPlex™" and a sub-heading "This is a selected diverse screening library which identifies potential drug leads easy, universal, and cost effective." Below this, there are several bullet points describing the library's features, such as "The HiPlex™ collection comprises 1M+ high quality compounds representing the drug like diversity of the Maybridge Screening Collections" and "All screening compounds fit general guidelines for drug like molecules, and all have purity greater than 99%". There is also a section titled "Ready to Screen" with a list of bullet points detailing the library's capabilities, such as "Individualized dry film for easy storage and use" and "Ready to ship in 24 hours".

Commercial  
(in-house pharma)



The screenshot shows the NIH Molecular Libraries Small Molecule Repository website. At the top, there is a navigation bar with links for Home, MLEMR Project, Compound Identification, Quality Control, Sample Storage, Sample Access, and Information. Below this is a search bar and a sidebar with navigation options like Home, MLEMR Project, Compound Identification, Quality Control, Sample Storage, Sample Access, and Information. The main content area features the heading "NIH Molecular Libraries Small Molecule Repository" and a sub-heading "A Roadmap Initiative". Below this, there is a "Welcome" section with a paragraph describing the repository's mission: "NIH Molecular Libraries Small Molecule Repository collects samples for high throughput biological screening and distributes them to the NIH Molecular Libraries Probe Production Centers Network (MLPCN)." There is also a "Get Involved" section with a list of bullet points detailing the repository's capabilities, such as "The MLEMR is a key component of the Molecular Libraries Initiative, an NIH Roadmap project supporting the Pathways to Discovery in the 21st Century" and "The project is funded in whole with Federal funds from the National Institutes of Health, Department of Health and Human Services, under Contract No. HHS-N-279-2204-41001C." There is also a "Contact Us" section with a list of bullet points detailing the repository's contact information, such as "Contact: © 2007 Galapagos NV" and "BioFocus, a Galapagos company operates MLEMR in South San Francisco."

Government (NIH)

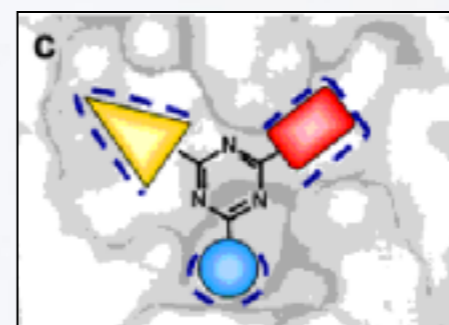
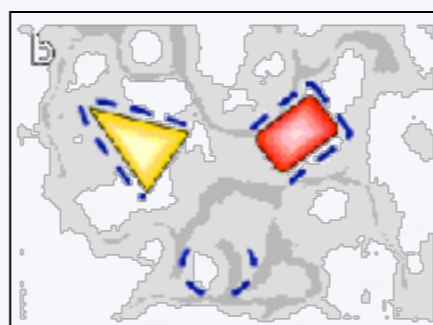
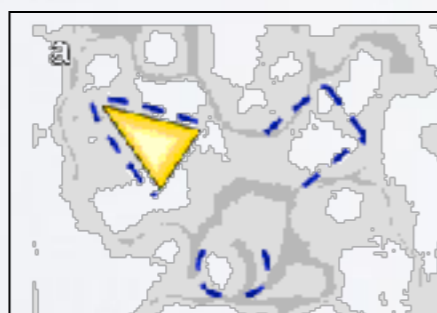
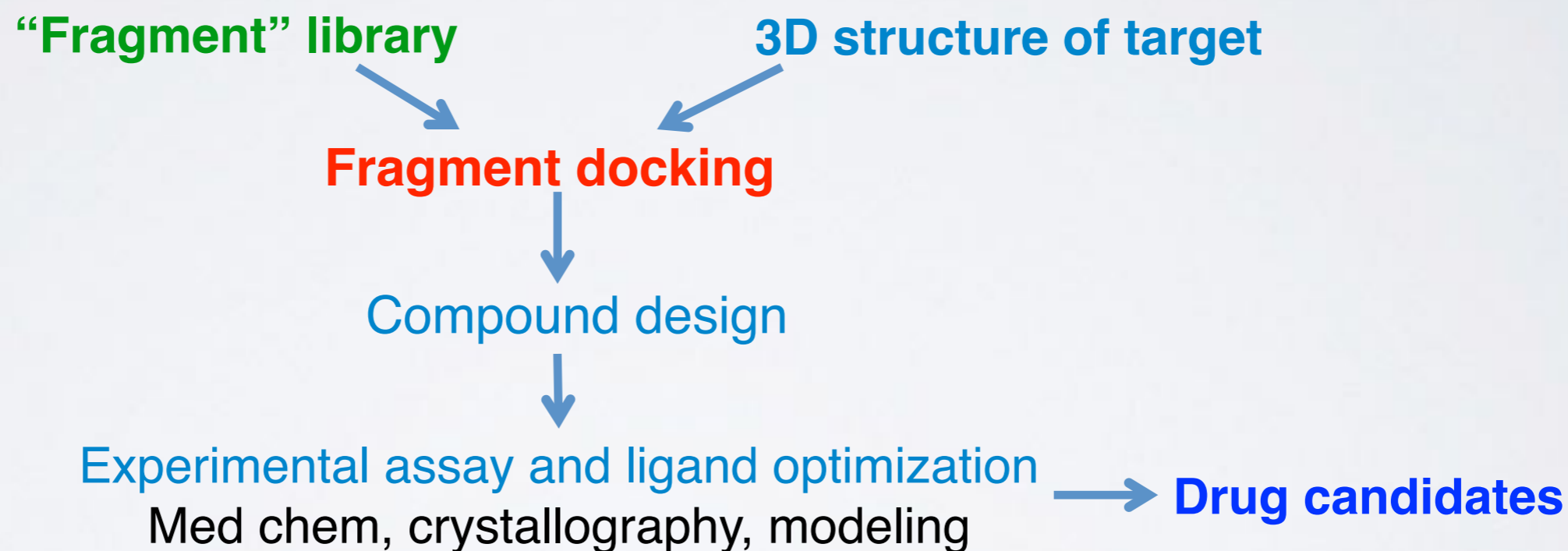


The screenshot shows the University of Pittsburgh Pittsburgh Molecular Libraries Screening Center website. At the top, there is a navigation bar with links for Home, About, Publications, Screening Technology, Compound Libraries, Resources & Applications, IT & Analytics, Approved PMLSC Assay Protocols, PMLSC Proc. Reports, Librarian, Data Analysis/Informatics, Educational Activities, Press Releases, Links, and Contacts. Below this is a search bar and a sidebar with navigation options like Home, About, Publications, Screening Technology, Compound Libraries, Resources & Applications, IT & Analytics, Approved PMLSC Assay Protocols, PMLSC Proc. Reports, Librarian, Data Analysis/Informatics, Educational Activities, Press Releases, Links, and Contacts. The main content area features the heading "PMLSC" and a sub-heading "BIG DISCOVERIES". Below this, there is a "Welcome" section with a paragraph describing the center's mission: "The Pittsburgh Molecular Library Screening Center (PMLSC) comprises investigators at the University of Pittsburgh and Carnegie Mellon University. The mission is to assist scientists and the National Institutes of Health to thoughtfully interrogate small molecule libraries using state-of-the-art High Throughput and High Content assays." There is also a "Get Involved" section with a list of bullet points detailing the center's capabilities, such as "The Pittsburgh Molecular Library Screening Center (PMLSC) comprises investigators at the University of Pittsburgh and Carnegie Mellon University" and "The mission is to assist scientists and the National Institutes of Health to thoughtfully interrogate small molecule libraries using state-of-the-art High Throughput and High Content assays."

Academia

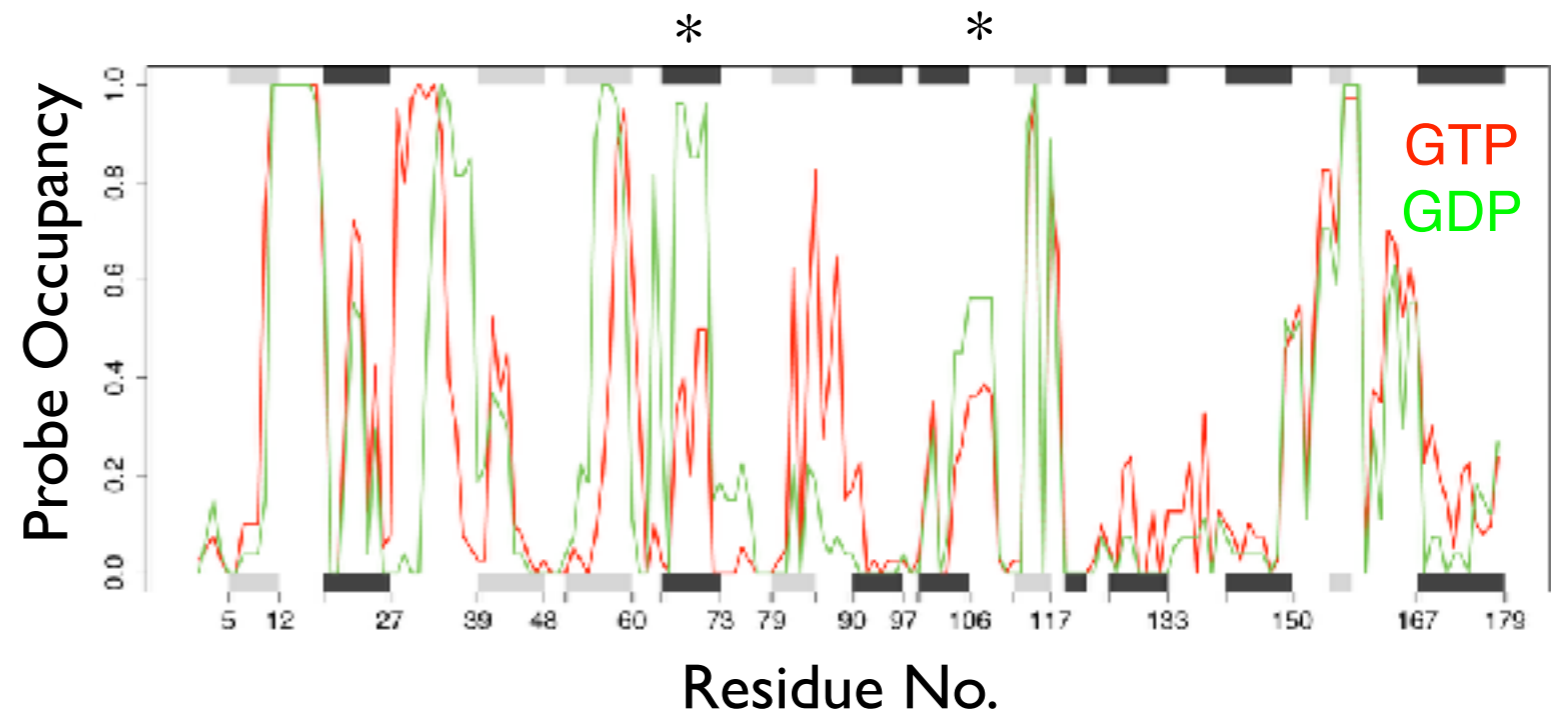
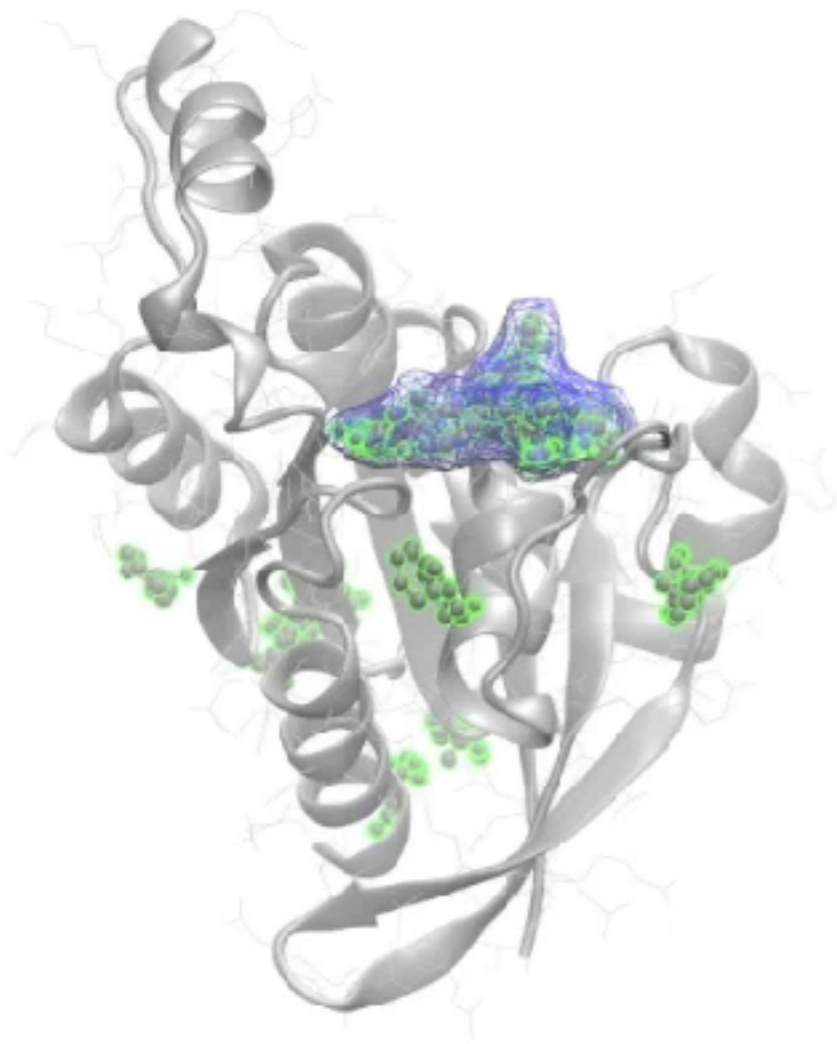


# FRAGMENTAL STRUCTURE-BASED SCREENING



# Multiple non active-site pockets identified

Small organic probe fragment affinities map multiple potential binding sites across the structural ensemble.



ethanol



isopropanol

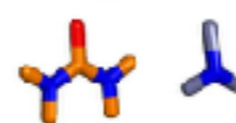
acetone



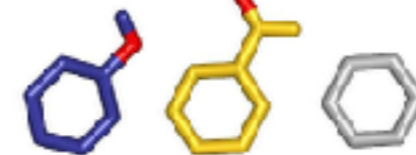
cyclohexane



methylamine



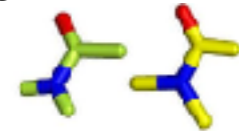
phenol



benzene



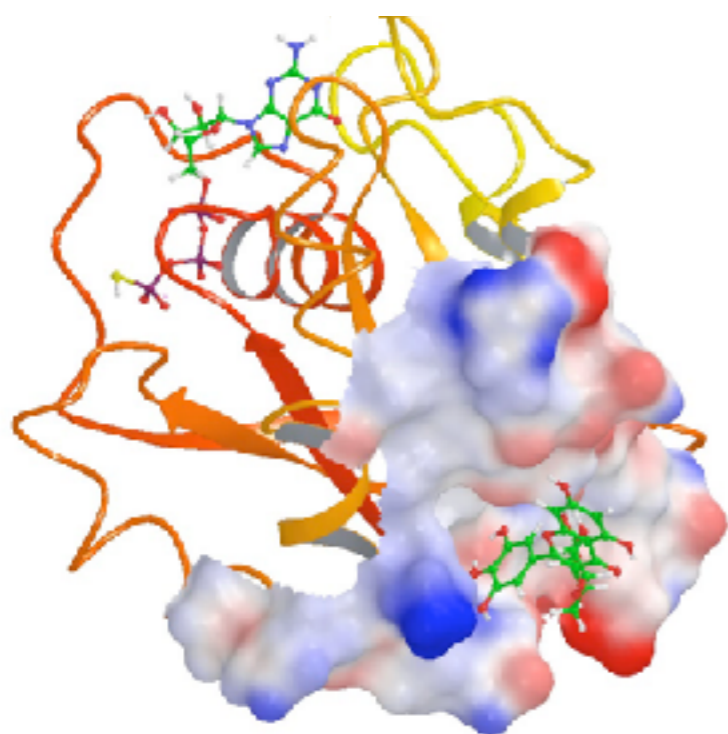
acetamide



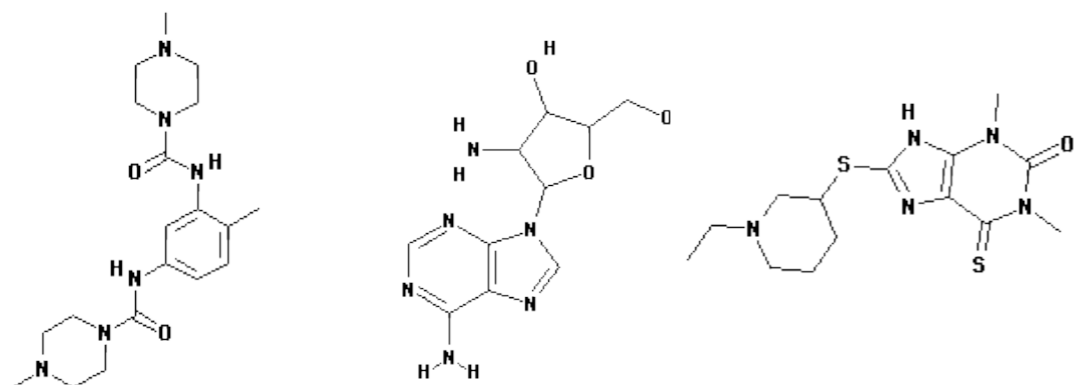
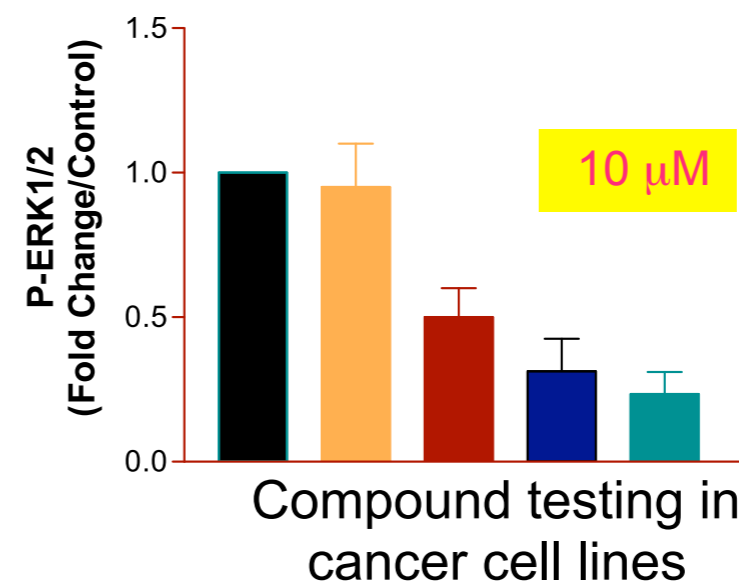
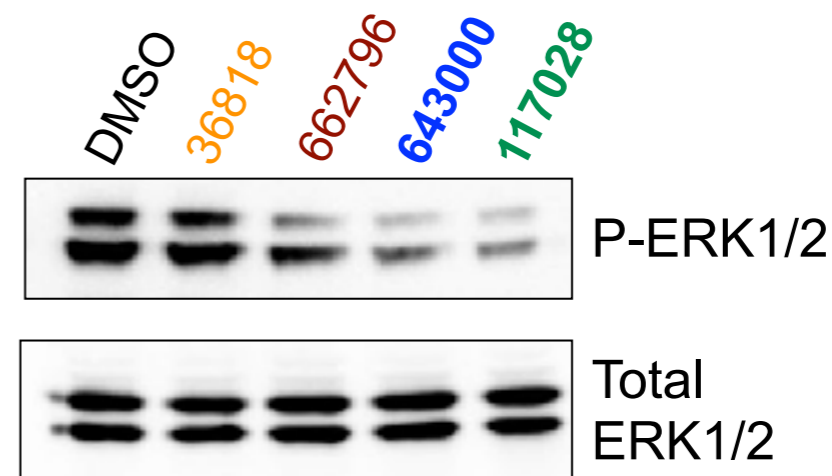
# Ensemble docking & candidate inhibitor testing

Top hits from ensemble docking against distal pockets were tested for inhibitory effects on basal ERK activity in glioblastoma cell lines.

Ensemble computational docking

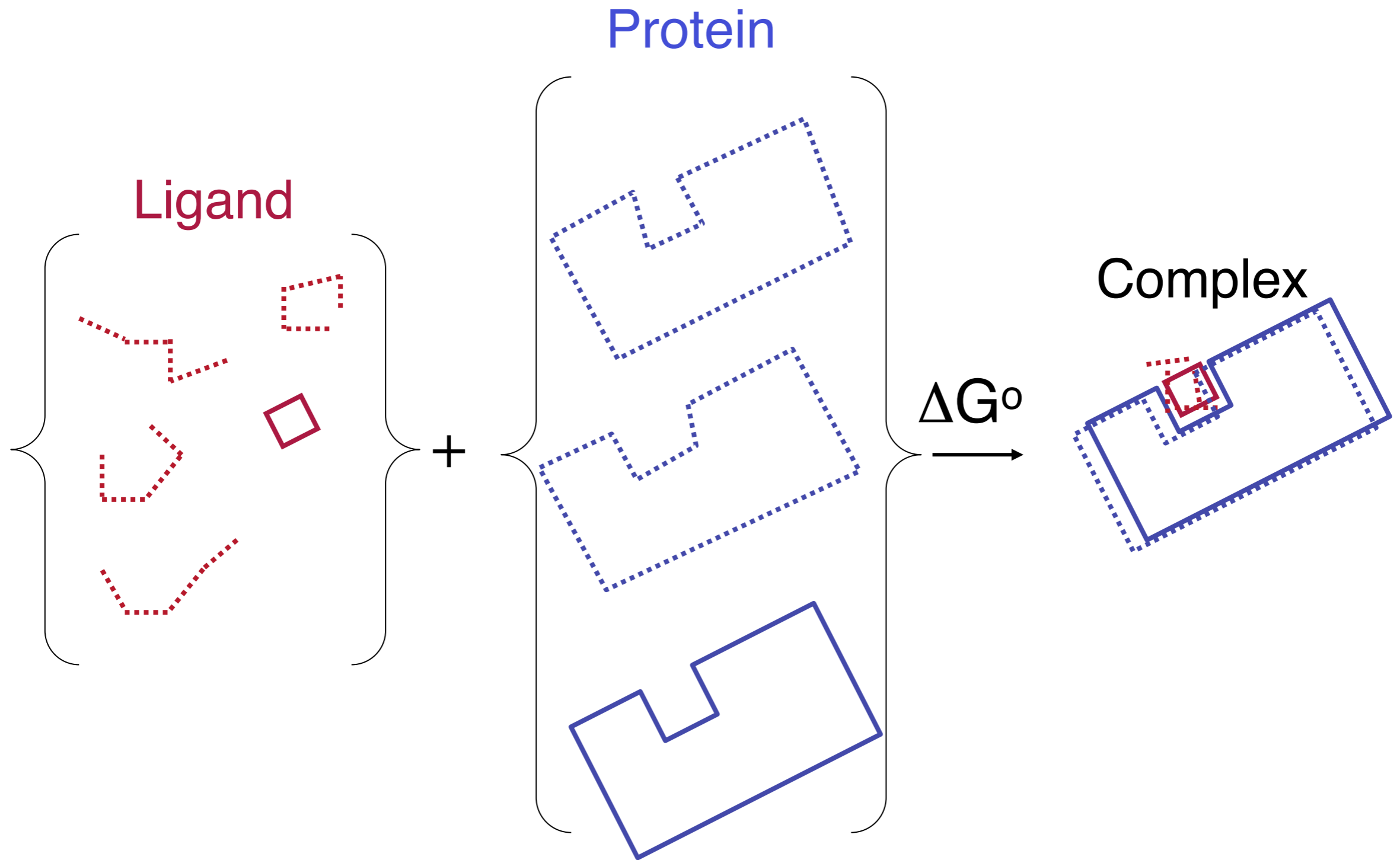


Compound effect on U251 cell line





# Proteins and Ligand are Flexible



# COMMON SIMPLIFICATIONS USED IN PHYSICS-BASED DOCKING

Quantum effects approximated classically

Protein often held rigid

Configurational entropy neglected

Influence of water treated crudely

Two main approaches:

**(1). Receptor/Target-Based**

**(2). Ligand/Drug-Based**

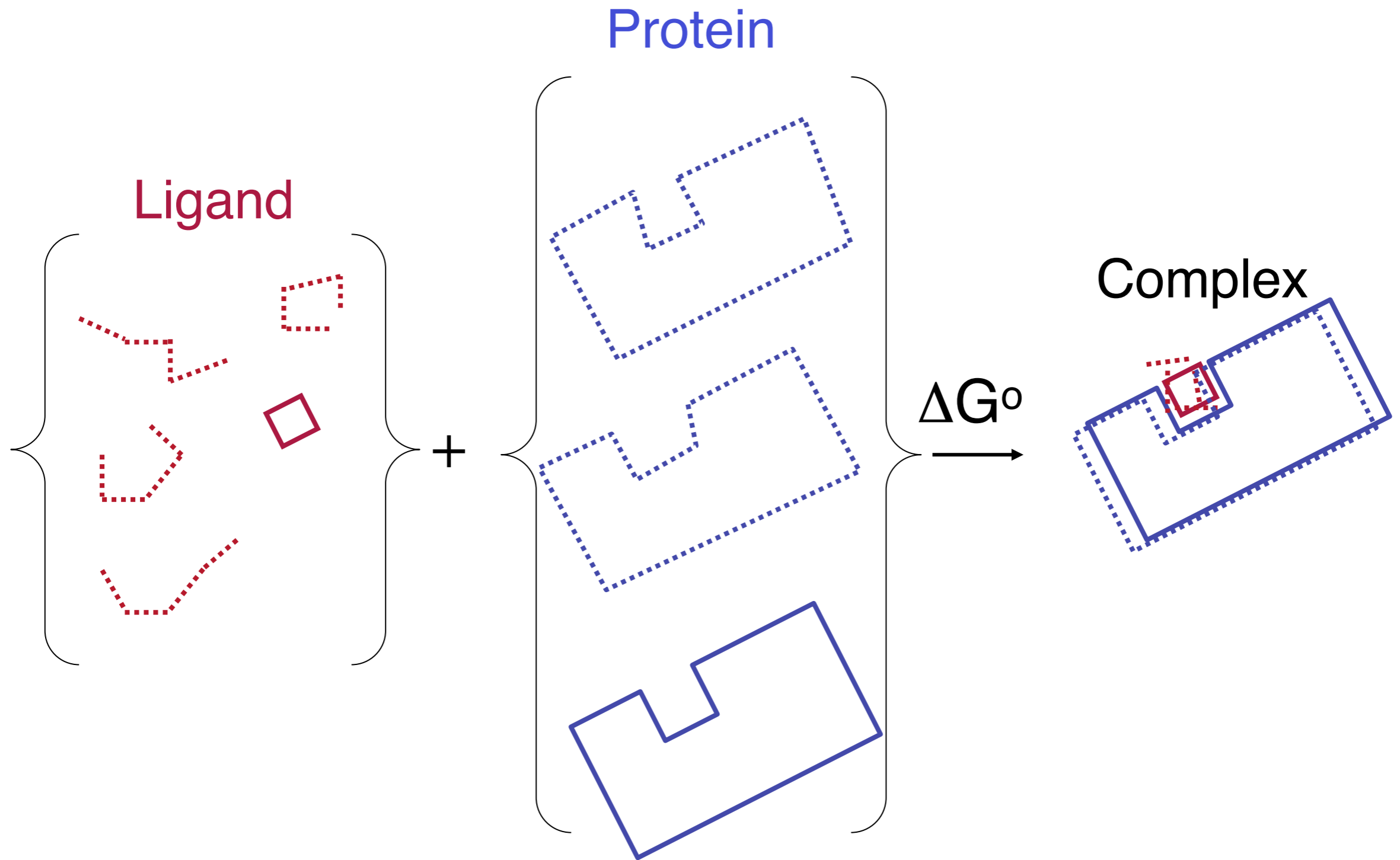
Do it Yourself!

# Hand-on time!

[https://bioboot.github.io/bimm143\\_W18/lectures/#12](https://bioboot.github.io/bimm143_W18/lectures/#12)

You can use the classroom computers or your own laptops. If you are using your laptops then you will need to install **VMD** and **MGLTools**

# Proteins and Ligand are Flexible





[HTTP://129.177.232.111:3848/PCA-APP/](http://129.177.232.111:3848/PCA-APP/)

[HTTPS://DCMB-GRANT-SHINY.UMMS.MED.UMICH.EDU/PCA-APP/](https://DCMB-GRANT-SHINY.UMMS.MED.UMICH.EDU/PCA-APP/)

[HTTP://BIO3D.UCSD.EDU/PCA-APP/](http://BIO3D.UCSD.EDU/PCA-APP/)

Two main approaches:

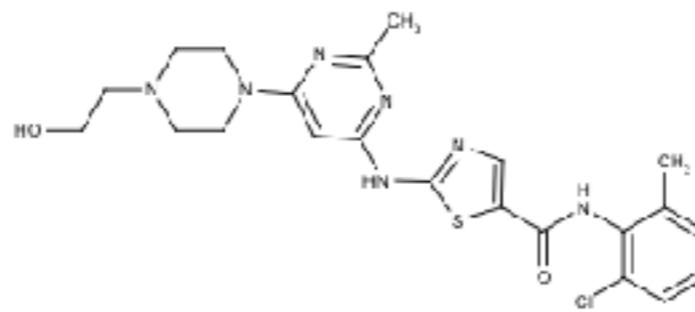
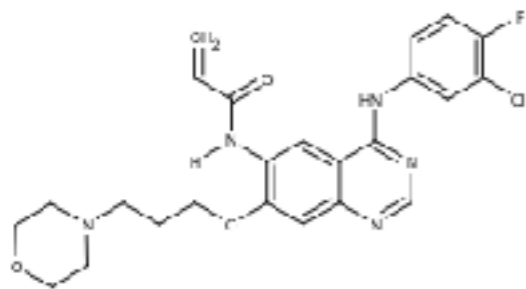
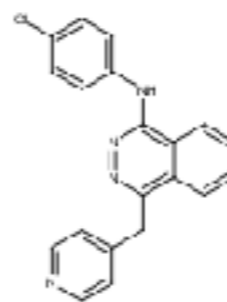
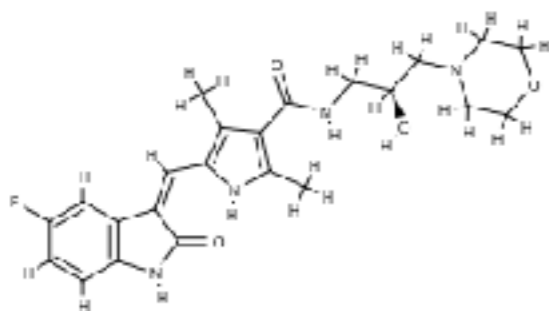
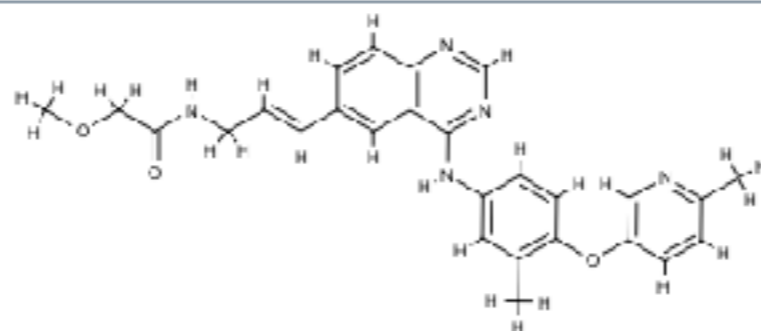
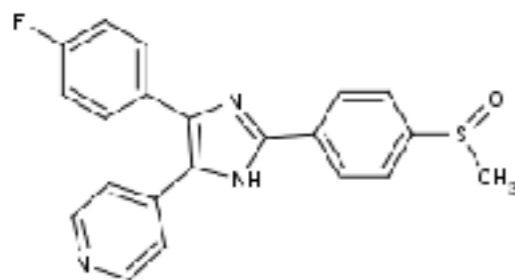
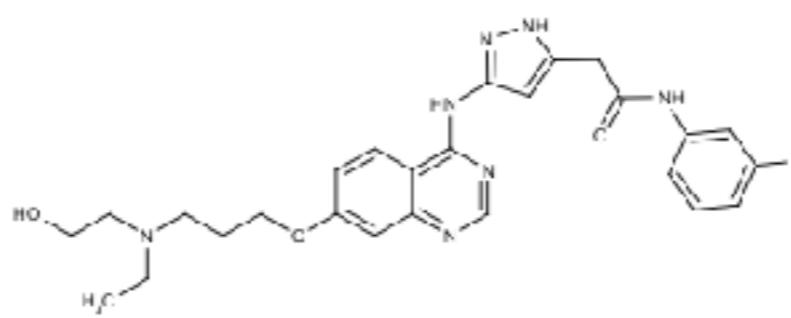
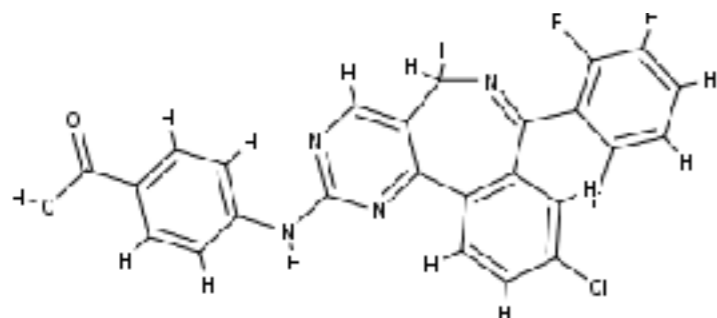
(1). **Receptor/Target-Based**

(2). **Ligand/Drug-Based**

# Scenario 2

## Structure of Targeted Protein Unknown: Ligand-Based Drug Discovery

e.g. MAP Kinase Inhibitors



Using knowledge of existing inhibitors to discover more

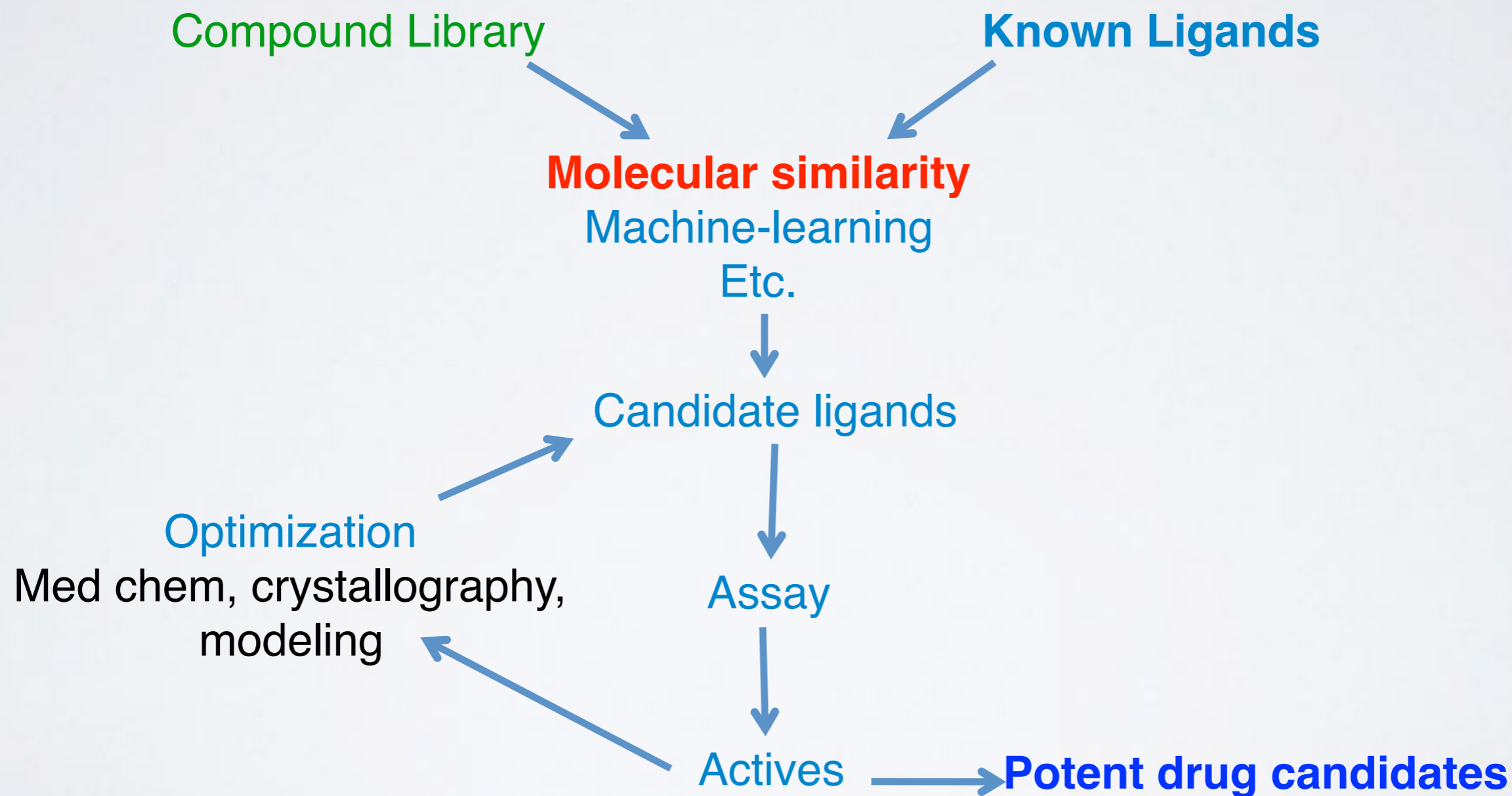
# Why Look for Another Ligand if You Already Have Some?

Experimental screening generated some ligands, but they don't bind tightly enough

A company wants to work around another company's chemical patents

An high-affinity ligand is toxic, is not well-absorbed, difficult to synthesize etc.

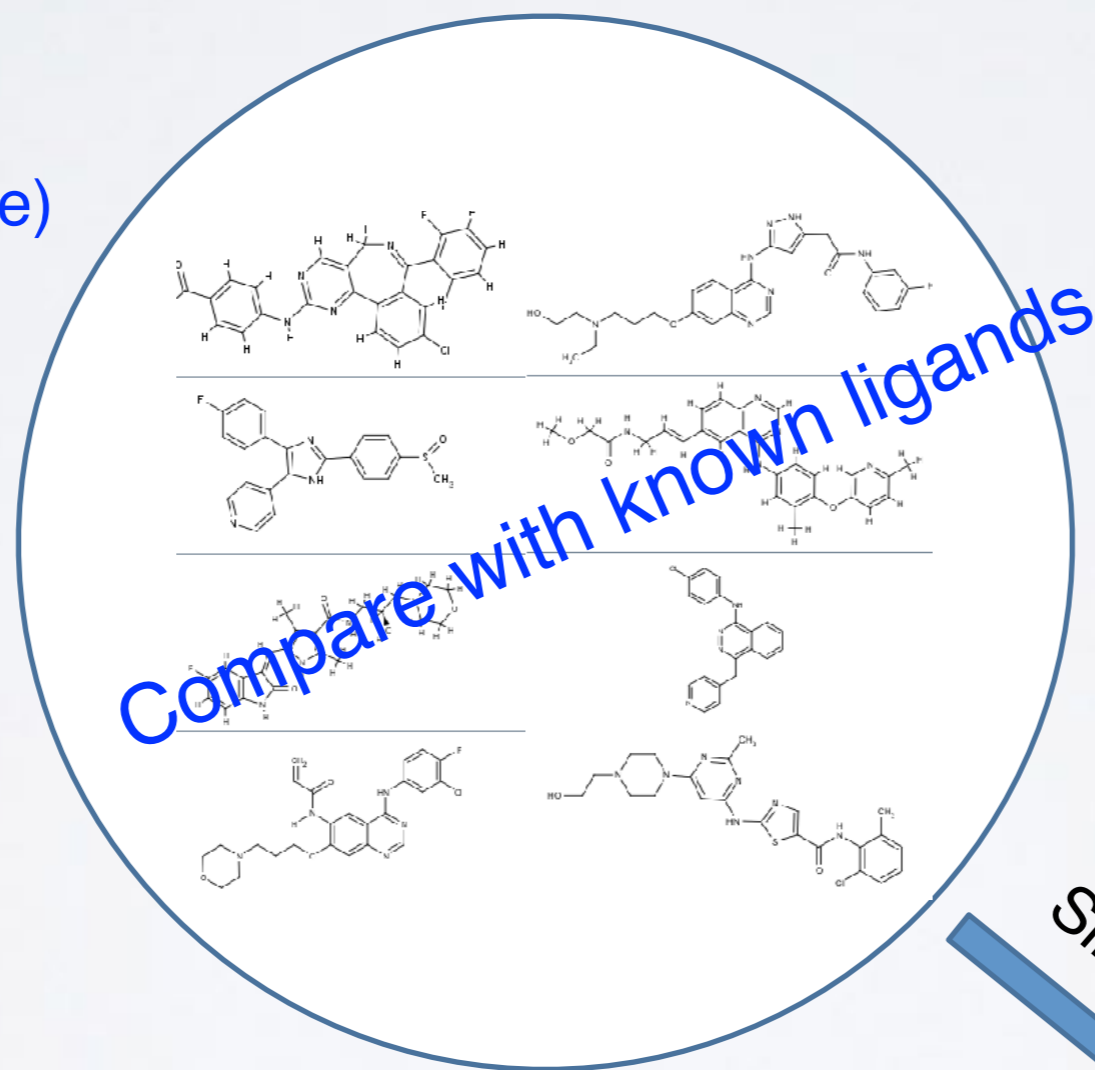
# LIGAND-BASED VIRTUAL SCREENING



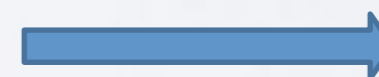


# CHEMICAL SIMILARITY LIGAND-BASED DRUG-DISCOVERY

Compounds  
(available/synthesizable)

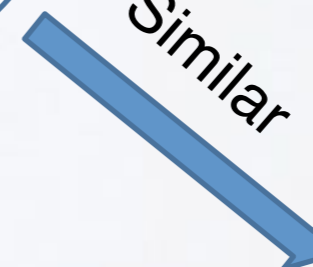


Different



Don't bother

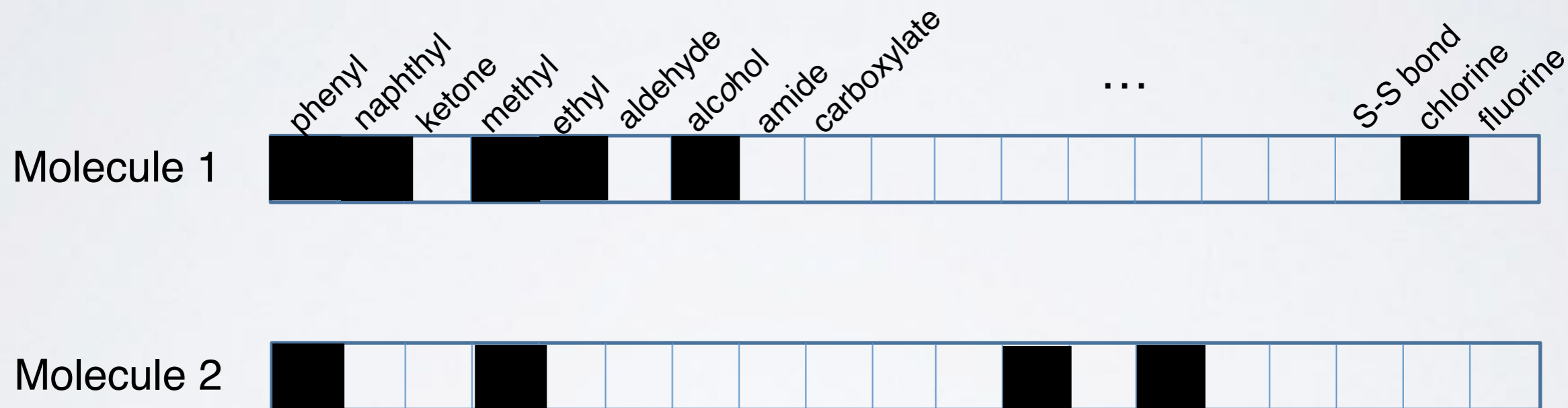
Similar



Test experimentally

# CHEMICAL FINGERPRINTS

## BINARY STRUCTURE KEYS



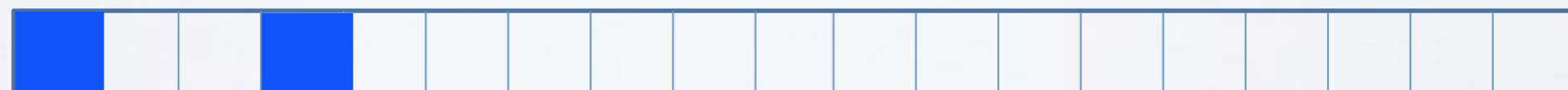
# CHEMICAL SIMILARITY FROM FINGERPRINTS



Tanimoto Similarity  
(or Jaccard Index),  $T$

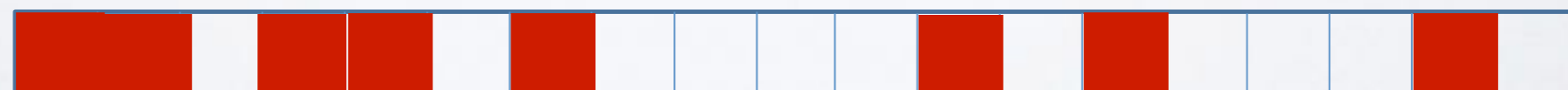
$$T \equiv \frac{N_I}{N_U} = 0.25$$

Intersection



$N_I=2$

Union

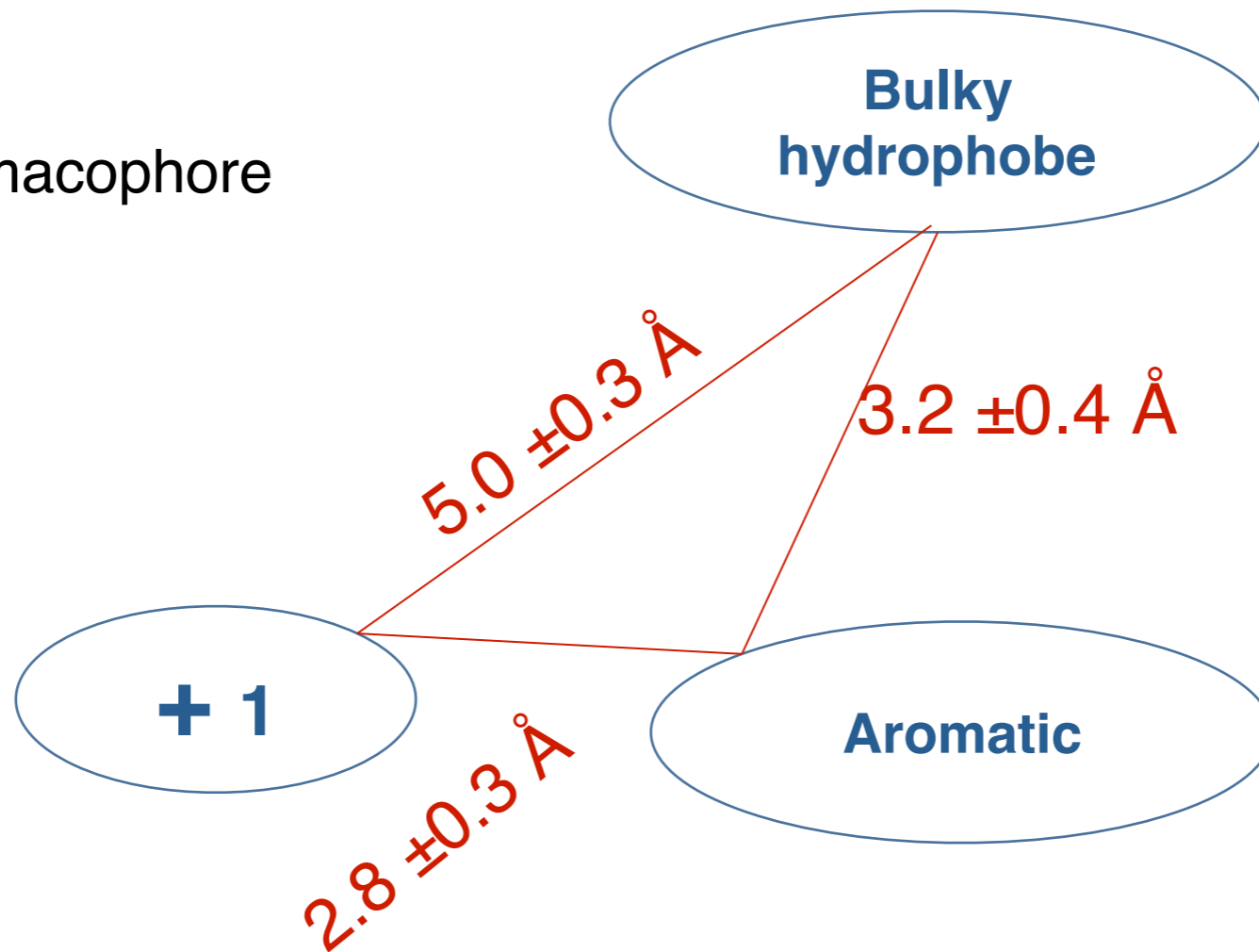


$N_U=8$

# Pharmacophore Models

Φάρμακο (drug) + Φορά (carry)

A 3-point pharmacophore



# Molecular Descriptors

More abstract than chemical fingerprints

## Physical descriptors

molecular weight

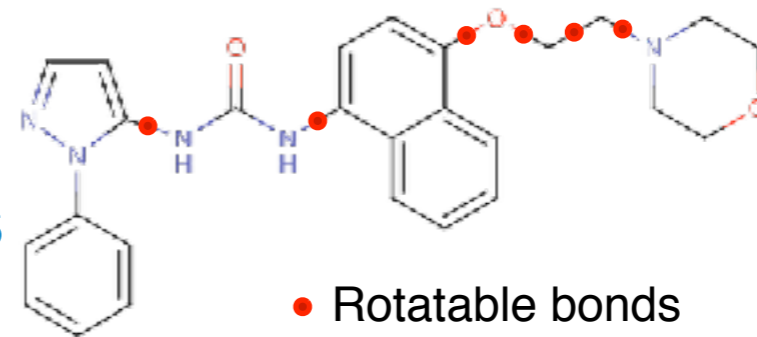
charge

dipole moment

number of H-bond donors/acceptors

number of rotatable bonds

hydrophobicity (log P and clogP)



## Topological

branching index

measures of linearity vs interconnectedness

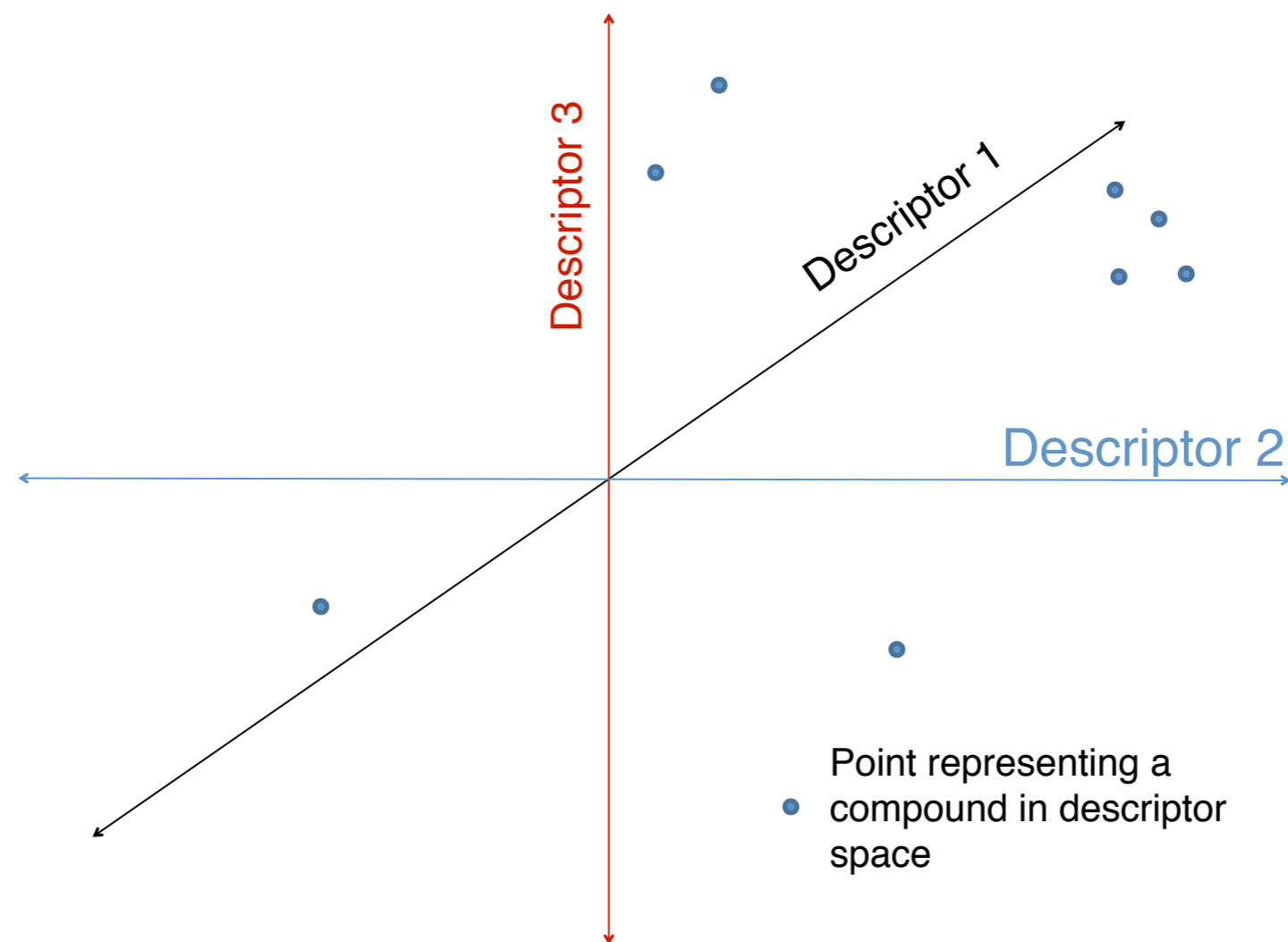
Etc. etc.



# A High-Dimensional “Chemical Space”

Each compound is at a point in an n-dimensional space

Compounds with similar properties are near each other



Apply **multivariate statistics** and **machine learning** for descriptor-selection. (e.g. partial least squares, PCA, support vector machines, random forest, deep learning etc.)

# Approved drugs and clinical candidates

- Catalogue approved drugs and clinical candidates from FDA Orange Book, and USAN applications
- Small molecules and biotherapeutics

ChEMBL wellcome trust



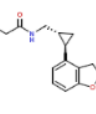

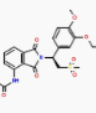

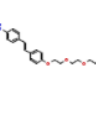

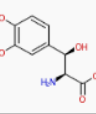

EBI > Databases > Small Molecules > ChEMBL Database > Home

Search ChEMBL... Compounds Targets Assays Documents [Activity Source Filter](#)

Ligand Search Target Search Browse Targets **Browse Drugs** Browse Drug Targets Drug Approvals About

Downloads... ▾

10 records per page Search:  Show / hide columns

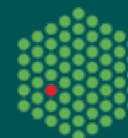
Parent Molecule	Synonyms	Phase	Research Codes	Applicants	USAN Stem	USAN Year	First Approval	ATC Code	Icon
 CHEMBL2108676	Elosulfase Alfa (INN, USAN)	4		Biomarin Pharmaceutical Inc.	-ase	2012	2014		
 CHEMBL2103822	Tasimelteon (FDA, INN, USAN)	4	BMS-214778 VEC-162	Vanda Pharmaceuticals Inc	-melteon	2007	2014		
 CHEMBL514800	Apremilast (FDA, INN, USAN)	4	CC-10004	Celgene Corp	-ast	2005	2014	L04AA32	
 CHEMBL1908906	Florbetaben F-18 (FDA) Florbetaben F18 (USAN)	4	BAY-949172 UNII-TLA7312TOI	Piramal Imaging Sa		2013	2014		
 CHEMBL1908906	Droxidopa (FDA, INN, USAN)	4	DOPS L-DOPS	Chelsea Therapeutics Inc	-dopa	2008	2014		

**ChEMBL Statistics**

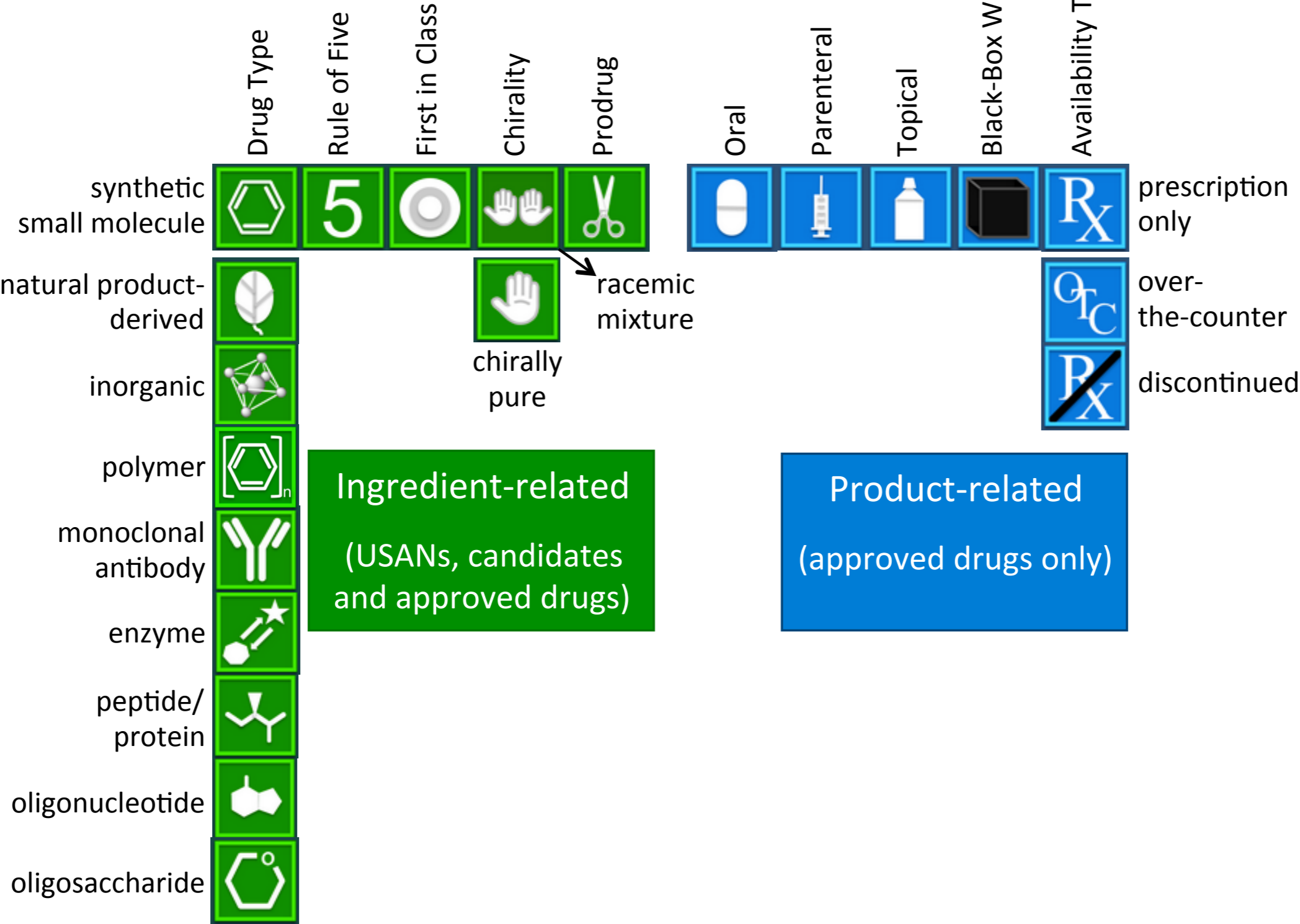
- DB: ChEMBL\_19
- Targets: 10,579
- Compound records: 1,638,394
- Distinct compounds: 1,411,786
- Activities: 12,843,338
- Publications: 57,156
- [Release Notes](#)

**ChEMBL Blog**

- [Another Confusion in the Literature - Trust but Verify](#)
- [Papers: Literature text mining and extensions to UniChem](#)



# Drug properties



# LIPINSKI'S RULE OF FIVE

Lipinski's rule of five states that, in general, an orally active drug has no more than one violation of the following criteria:

- Not more than 5 hydrogen bond donors (nitrogen or oxygen atoms with one or more hydrogen atoms)
- Not more than 10 hydrogen bond acceptors (nitrogen or oxygen atoms)
- A molecular mass less than 500 daltons
- An octanol-water partition coefficient  $\log P$  not greater than 5

# Rules for drug discovery success

- Set of approved drugs or medicinal chemistry compounds and their targets can be used to derive rules for drug discovery success (or failure):
  - What features make a successful drug target?
  - What features make a protein druggable by small molecules?
  - What features of a compound contribute to good oral bioavailability?
  - What chemical groups may be associated with toxicity?



# Druggability prediction

https://www.ebi.ac.uk/chembl/drugability/domain/32655

View cavities (and ligands) on structure

Details of sites identified

ChEMBL Blog

- Invitation to join the Teach-Discover-Treat Initiative
- Carbon and Oxygen - Simples

**Domain Details:**

PDB	1eeo <a href="#">SCOP</a>
Gene	P18031 <a href="#">View Protein Summary</a>
Description	Tyrosine-protein phosphatase non-receptor type 1
Fold	Phosphotyrosine protein phosphatases II
Superfamily	Phosphotyrosine protein phosphatases II
Family	Higher-molecular-weight phosphotyrosine protein phosphatases <a href="#">View Family Druggability</a>
Other PDB(s)	1eeo:A - px32655 ( Tractable: 1, Druggable: 0, Ensemble: -0.96 )

**Average Druggability Scores:**

Tractable	Druggable	Ensemble
0.97	0.02	-0.93

Tractable/Druggable ranges from low:0 to high:1. Ensemble ranges from low:-1 to high:+1.

**Site Druggability Details:**

Reset: <input type="radio"/>	Site 1	Site 2	Site 3	Site 4
Druggable	0.00	0.00	0.00	0.00
Confidence	0.73	0.96	0.96	0.96
Tractable	1.00	0.00	0.00	0.00
Confidence	0.92	0.86	0.83	0.86
Ensemble	-0.96	-0.99	-0.98	-0.99
Volume [Å <sup>3</sup> ]	1535.2	1318.36	1446.61	1454.2
Buried Surface [%]	71.3	65.25	72.27	64.08
Show Site	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Show Residues	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Ligand <input type="radio"/>	<p><b>PTR</b> <a href="#">CHEMBL286939</a></p>	-	-	-

Green :Druggable, Yellow :Tractable, Pink :Undruggable

Jmol

Use | Privacy | Cookies | EBI Funding | Contact EBI | © European Bioinformatics Institute 2012. EBI is an Outstation of the European Molecular Biology Laboratory.



# Target prediction models

- Active compounds from ChEMBL can be used to train target prediction models
- Variety of methods used
  - Multi-Category Naïve Bayesian Classifier (e.g., ChEMBL)
  - Chemical similarity between ligand sets (e.g., SEA)
  - 3D similarity between ligands (e.g., SwissTargetPrediction)
  - Protein and ligand descriptors (e.g., Proteochemometric models)
- Open source tools available for many methods
  - E.g., Scikit-learn with RDKit

Examples at: [https://github.com/chembl/mychembl/blob/master/ipython\\_notebooks](https://github.com/chembl/mychembl/blob/master/ipython_notebooks)





## OPEN ACCESS

**Citation:** Mugumbate G, Abrahams KA, Cox JAG, Papadatos G, van Westen G, Lelièvre J, et al. (2015) Mycobacterial Dihydrofolate Reductase Inhibitors Identified Using Chemogenomic Methods and *In Vitro* Validation. PLoS ONE 10(3): e0121492. doi:10.1371/journal.pone.0121492

**Academic Editor:** Anil Kumar Tyagi, University of Delhi, INDIA

**Received:** December 4, 2014

**Accepted:** February 1, 2015

**Published:** March 23, 2015

**Copyright:** © 2015 Mugumbate et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** All relevant data are within the paper and its Supporting Information files.

**Funding:** GM and GvW are grateful to EMBL and Marie Curie Actions for funding this work. GSB acknowledges support in the form of a Personal Research Chair from Mr. James Bardrick, a Royal Society Wolfson Research Merit Award, The Wellcome Trust (681569/Z/06/Z), and the Medical Research Council (MR/K012118/1). The research also received funding from the European Union's 7th framework programme (FP7-2007–2013) under grant agreement ORCID no. 261378. The funders had no

## RESEARCH ARTICLE

# Mycobacterial Dihydrofolate Reductase Inhibitors Identified Using Chemogenomic Methods and *In Vitro* Validation

Grace Mugumbate<sup>1</sup>, Katherine A. Abrahams<sup>2</sup>, Jonathan A. G. Cox<sup>2</sup>, George Papadatos<sup>1</sup>, Gerard van Westen<sup>1</sup>, Joël Lelièvre<sup>3</sup>, Szymon T. Calus<sup>2</sup>, Nicholas J. Loman<sup>2</sup>, Lluís Balcells<sup>3</sup>, David Barros<sup>3</sup>, John P. Overington<sup>1\*</sup>, Gurdyal S. Besra<sup>2\*</sup>

**1** European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Wellcome Trust Genome Campus, Hinxton, Cambridge, United Kingdom, **2** Institute of Microbiology and Infection (IMI), School of Biosciences, University of Birmingham, Edgbaston, Birmingham, United Kingdom, **3** Diseases of the Developing World, GlaxoSmithKline, Severo Ochoa 2, 28760 Tres Cantos, Madrid, Spain

\* [jpo@ebi.ac.uk](mailto:jpo@ebi.ac.uk) (JPO); [g.besra@bham.ac.uk](mailto:g.besra@bham.ac.uk) (GSB)

## Abstract

The lack of success in target-based screening approaches to the discovery of antibacterial agents has led to reemergence of phenotypic screening as a successful approach of identifying bioactive, antibacterial compounds. A challenge though with this route is then to identify the molecular target(s) and mechanism of action of the hits. This target identification, or deorphanization step, is often essential in further optimization and validation studies. Direct experimental identification of the molecular target of a screening hit is often complex, precisely because the properties and specificity of the hit are not yet optimized against that target, and so many false positives are often obtained. An alternative is to use computational, predictive, approaches to hypothesize a mechanism of action, which can then be validated in a more directed and efficient manner. Specifically here we present experimental validation of an *in silico* prediction from a large-scale screen performed against *Mycobacterium tuberculosis* (*Mtb*), the causative agent of tuberculosis. The **two** potent anti-tubercular compounds studied in this case, belonging to the tetrahydro-1,3,5-triazin-2-amine (THT) family, were predicted and confirmed to be an inhibitor of dihydrofolate reductase (DHFR), a known essential *Mtb* gene, and already clinically validated as a drug target. Given the large number of similar screening data sets shared amongst the community, this *in vitro* validation of these target predictions gives weight to computational approaches to establish the mechanism of action (MoA) of novel screening hit.

## Introduction

The human pathogen, *Mycobacterium tuberculosis* (*Mtb*) is the causative agent of tuberculosis (TB), an infectious disease that is widespread, infecting around one third of the world's population [1]. The discovery of streptomycin in 1943, and the subsequent discovery and

# Large-scale prediction and testing of drug activity on side-effect targets

Eugen Lounkine<sup>1\*</sup>, Michael J. Keiser<sup>2,3\*</sup>, Steven Whitebread<sup>1</sup>, Dmitri Mikhailov<sup>1</sup>, Jacques Hamon<sup>4</sup>, Jeremy L. Jenkins<sup>1</sup>, Paul Lavan<sup>4</sup>, Eckhard Weber<sup>4</sup>, Allison K. Doak<sup>3</sup>, Serge Côté<sup>4</sup>, Brian K. Shoicher<sup>3</sup> & Laszlo Urban<sup>4</sup>

Discovering the unintended 'off-targets' that predict adverse drug reactions is daunting by empirical methods alone. Drugs can act on several protein targets, some of which can be unrelated by conventional molecular metrics, and hundreds of proteins have been implicated in side effects. Here we use a computational strategy to predict the activity of 656 marketed drugs on 73 unintended 'side-effect' targets. Approximately half of the predictions were confirmed, either from proprietary databases unknown to the method or by new experimental assays. Affinities for these new off-targets ranged from 1 nM to 30 μM. To explore relevance, we developed an association metric to prioritize those new off-targets that explained side effects better than any known target of a given drug, creating a drug-target-adverse drug reaction network. Among these new associations was the prediction that the abdominal pain side effect of the synthetic oestrogen chlorotrianisene was mediated through its newly discovered inhibition of the enzyme cyclooxygenase-1. The clinical relevance of this inhibition was borne out in whole human blood platelet aggregation assays. This approach may have wide application to de-risking toxicological liabilities in drug discovery.

Adverse drug reactions (ADRs) can limit the use of otherwise effective drugs. Next to lack of efficacy, they are the leading cause for attrition in clinical trials of new drugs<sup>1–3</sup> and are more prominent still in the failure of molecules to advance from pre-clinical research into human trials<sup>4</sup>. Some ADRs are caused by modulation of the primary target of a drug<sup>5</sup>, others result from non-specific interactions of reactive metabolites<sup>6</sup>. In many cases, however, ADRs are caused by unintended activity at off-targets. Notorious examples of off-target toxicity include that of the appetite suppressant fenfluramine-phenentermine (fen-phen), which was withdrawn from the market after numerous patient deaths. These owed to the activation of the 5-hydroxytryptamine-2B (5-HT<sub>2B</sub>) receptor by one of its metabolites, norfenfluramine, leading to proliferative valvular heart disease<sup>7</sup>. Similarly, well-known drugs, such as the antihistamine terfenadine, have been withdrawn because they caused arrhythmias and death, which have been attributed to their off-target inhibition of the human *ether-à-go-go*-related gene potassium channel (hERG, also known as KCNH2)<sup>8,9</sup>. Prediction of unknown off-target drug interactions might prevent such disastrous drug toxicities, which are often detected only after fatalities in the clinic, and might allow safer molecules to be prioritized for pre-clinical development. Methods to systematically predict off-targets, and associate these with side effects, have thus attracted intense interest<sup>10–16</sup>, frequently in the form of either chemical genomics<sup>17,18</sup> or informatics<sup>19–26</sup> approaches.

Whereas the informatics methods have never been tested systematically on a large scale, in principle they can be deployed against thousands of targets. Here we present a large-scale, prospective evaluation of safety target prediction using one such method, the similarity ensemble approach (SEA)<sup>25–27</sup>. SEA calculates whether a molecule will bind to a target based on the chemical features it shares with those of known ligands, using a statistical model to control for random similarity. Because SEA relies only on chemical similarity, it can be applied systematically and, for those targets that have known ligands,

comprehensively. For 656 drugs approved for human use (Supplementary Table 1), targets were predicted from among 73 proteins (Supplementary Table 2 and Methods) with established association of ADRs<sup>22,28</sup>, for which assays were available at Novartis. Encouragingly, many of the predictions were confirmed, often at pharmacologically relevant concentrations. This motivated us to develop a guilt-by-association metric that linked the new targets to the ADRs of those drugs for which they are the primary or well-known off-targets, creating a drug-target-ADR network. The applicability and the limitations of this approach will be considered.

## Testing the predictions

The 656 drugs were computationally screened for their likelihood to bind to 73 targets (Supplementary Table 2) using SEA<sup>25–27</sup>. The targets belong to the Novartis *in vitro* safety panels based on their association with ADRs<sup>22,28</sup>. Here we insisted that they also be described in the ChEMBL database<sup>29</sup>, enabling correspondence with SEA predictions (Supplementary Table 2). ChEMBL annotates more than 285,000 ligands modulating more than 1,500 different human targets with affinities better than 30 μM. SEA calculated the similarity of each drug versus each set of ligands for the 73 targets, comparing the overall set similarity to a model of such expected at random. For instance, the sodium channel blocker aprindine loosely resembled the set of histamine H<sub>1</sub> ligands; although no single H<sub>1</sub> ligand was strongly similar to the drug (Table 1), the overall similarity of the set was much greater than expected at random, leading to a highly significant SEA expectation value (*E* value) of  $5 \times 10^{-26}$  between aprindine and H<sub>1</sub> receptor ligands. Only 1,644 of the more than 47,000 possible drug-target pairs had significant *E* values. Of these, 403 were already known in ChEMBL and so were trivially confirmed; we do not consider these further. Of the remaining 1,241 predictions, 348 (28%) were unknown to ChEMBL, but could be found in proprietary ligand-target databases that were unavailable to SEA (see Methods). The remaining

<sup>1</sup>Novartis Institutes for Biomedical Research, Cambridge, Massachusetts 02139, USA. <sup>2</sup>SeaChange Pharmaceuticals Inc, 409 Illinois Street, San Francisco, California 94158, USA. <sup>3</sup>Department of Pharmaceutical Chemistry, University of California, San Francisco, 1700 4th Street, Byers Hall Suite 508D, California 94158-2550, USA. <sup>4</sup>Novartis Institutes for Biomedical Research, 4056 Basel, Switzerland.

\*These authors contributed equally to this work.

# NEXT UP:

- ▶ **Overview of structural bioinformatics**

- Major motivations, goals and challenges

- ▶ **Fundamentals of protein structure**

- Composition, form, forces and dynamics

- ▶ **Representing and interpreting protein structure**

- Modeling energy as a function of structure

- ▶ **Example application areas**

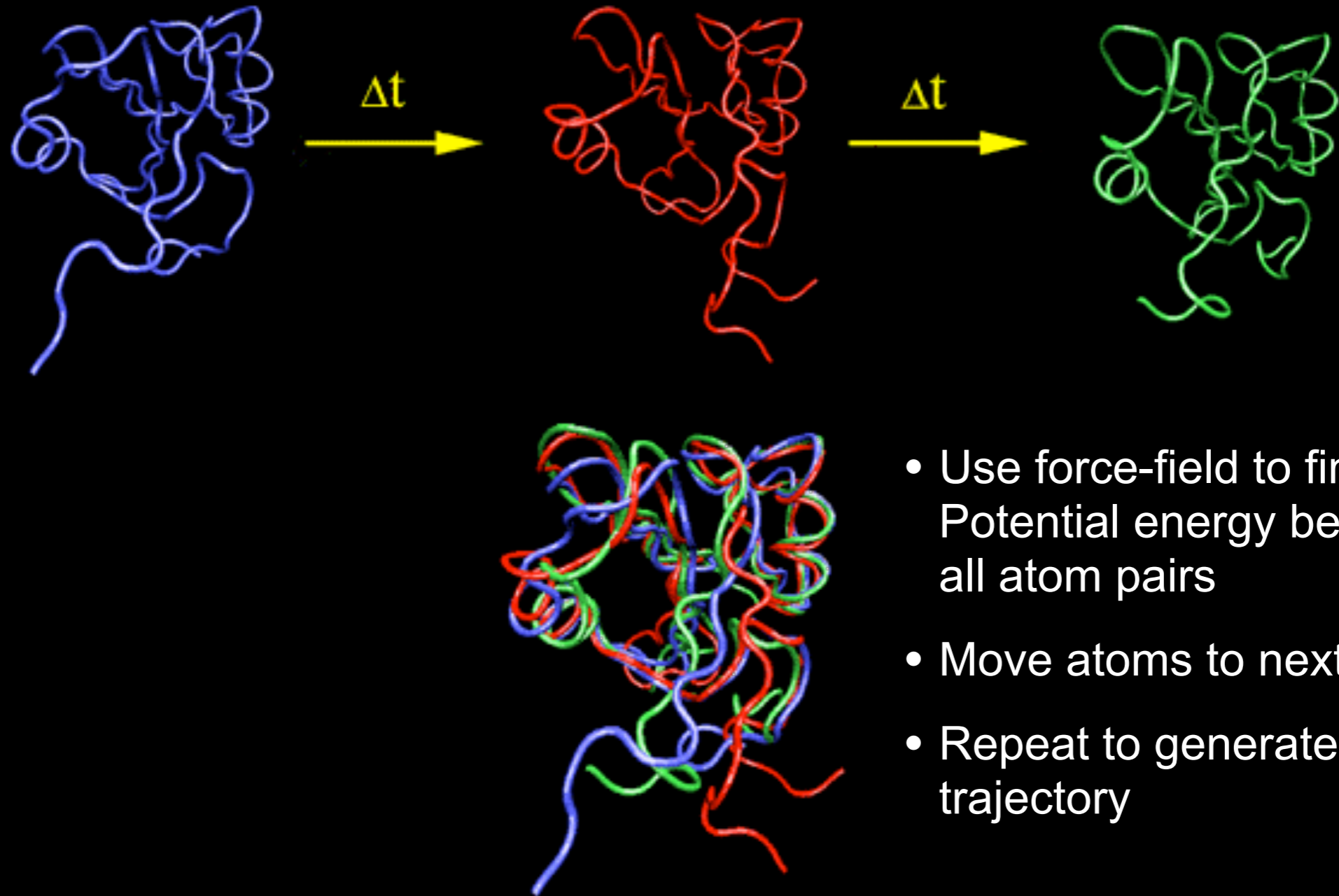
- Drug discovery & predicting functional dynamics

# PREDICTING FUNCTIONAL DYNAMICS

- Proteins are intrinsically flexible molecules with internal motions that are often intimately coupled to their biochemical function
  - E.g. ligand and substrate binding, conformational activation, allosteric regulation, etc.
- Thus knowledge of dynamics can provide a deeper understanding of the mapping of structure to function
  - Molecular dynamics (MD) and normal mode analysis (NMA) are two major methods for predicting and characterizing molecular motions and their properties



# MOLECULAR DYNAMICS SIMULATION



McCammon, Gelin & Karplus, *Nature* (1977)

[ See: <https://www.youtube.com/watch?v=ui1ZysMFcKk> ]

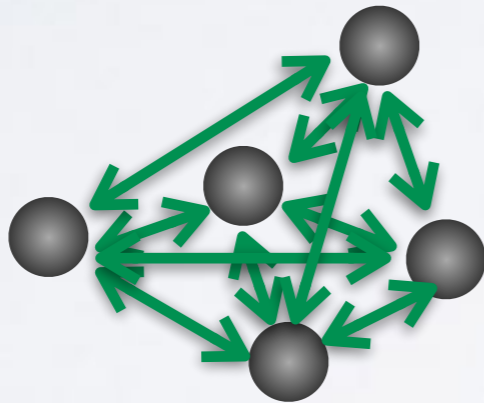
- ▶ Divide **time** into discrete ( $\sim 1$ fs) **time steps** ( $\Delta t$ )  
(for integrating equations of motion, see below)



- ▶ Divide **time** into discrete ( $\sim 1$ fs) **time steps** ( $\Delta t$ )  
(for integrating equations of motion, see below)



- ▶ At each time step calculate pair-wise atomic **forces** ( $F(t)$ )  
(by evaluating **force-field** gradient)



*Nucleic motion described classically*

$$m_i \frac{d^2 \vec{R}_i}{dt^2} = -\vec{\nabla}_i E(\vec{R})$$

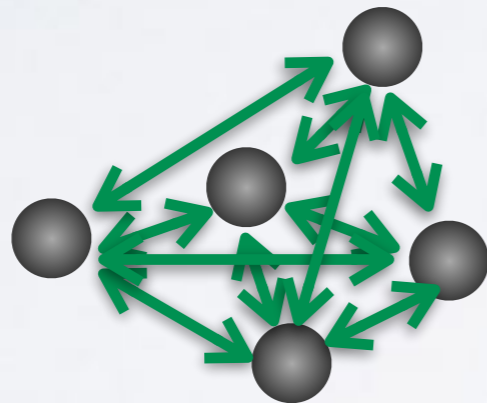
*Empirical force field*

$$E(\vec{R}) = \sum_{\text{bonded}} E_i(\vec{R}) + \sum_{\text{non-bonded}} E_i(\vec{R})$$

- ▶ Divide **time** into discrete ( $\sim 1$ fs) **time steps** ( $\Delta t$ )  
(for integrating equations of motion, see below)



- ▶ At each time step calculate pair-wise atomic **forces** ( $F(t)$ )  
(by evaluating **force-field** gradient)



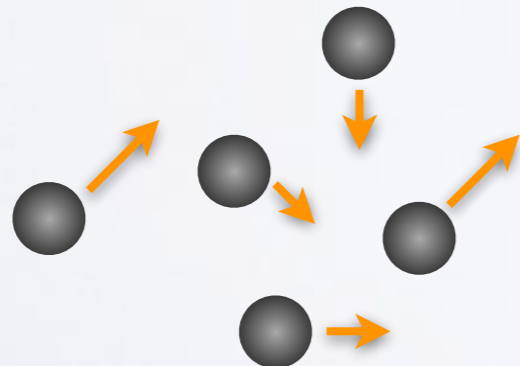
*Nucleic motion described classically*

$$m_i \frac{d^2 \vec{R}_i}{dt^2} = -\vec{\nabla}_i E(\vec{R})$$

*Empirical force field*

$$E(\vec{R}) = \sum_{\text{bonded}} E_i(\vec{R}) + \sum_{\text{non-bonded}} E_i(\vec{R})$$

- ▶ Use the forces to calculate **velocities** and move atoms to new **positions**  
(by integrating numerically via the “leapfrog” scheme)



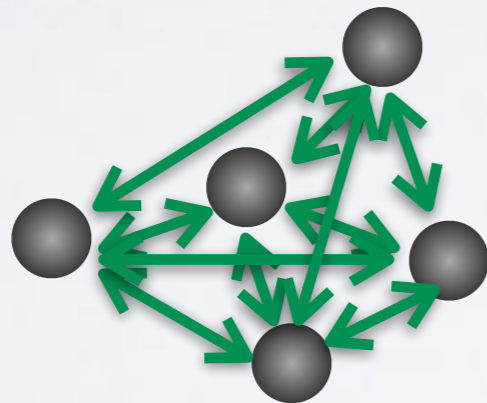
$$\begin{aligned} \mathbf{v}\left(t + \frac{\Delta t}{2}\right) &= \mathbf{v}\left(t - \frac{\Delta t}{2}\right) + \frac{\mathbf{F}(t)}{m} \Delta t \\ \mathbf{r}(t + \Delta t) &= \mathbf{r}(t) + \mathbf{v}\left(t + \frac{\Delta t}{2}\right) \Delta t \end{aligned}$$

# BASIC ANATOMY OF A MD SIMULATION

- ▶ Divide **time** into discrete ( $\sim 1$ fs) **time steps** ( $\Delta t$ )  
(for integrating equations of motion, see below)



- ▶ At each time step calculate pair-wise atomic **forces** ( $F(t)$ )  
(by evaluating **force-field** gradient)



*Nucleic motion described classically*

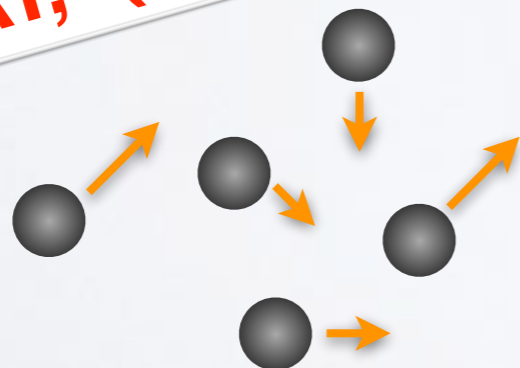
$$m_i \frac{d^2 \vec{R}_i}{dt^2} = -\vec{\nabla}_i E(\vec{R})$$

*Empirical force field*

$$E(\vec{R}) = \sum_{\text{non-bonded}} E_i(\vec{R})$$

- ▶ Use the forces to calculate **velocities** and move atoms to new **positions**  
(the integration is done numerically via the "leapfrog" scheme)

**REPEAT, (iterate many, many times... 1ms = 10<sup>12</sup> time steps)**

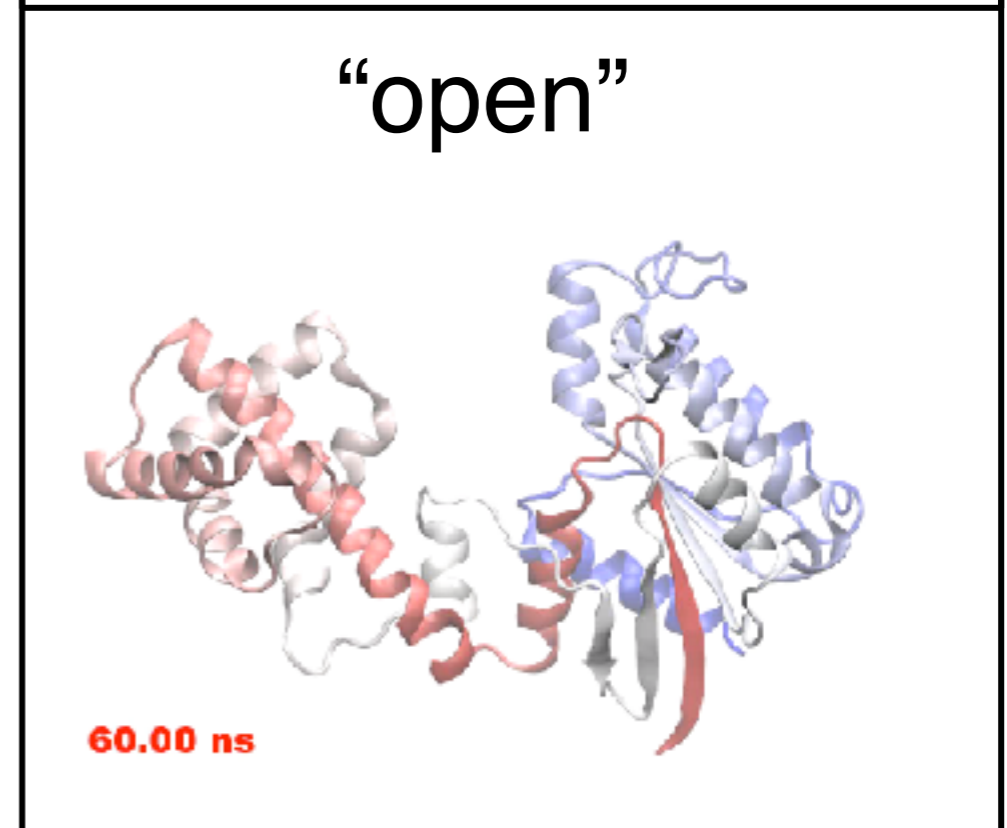
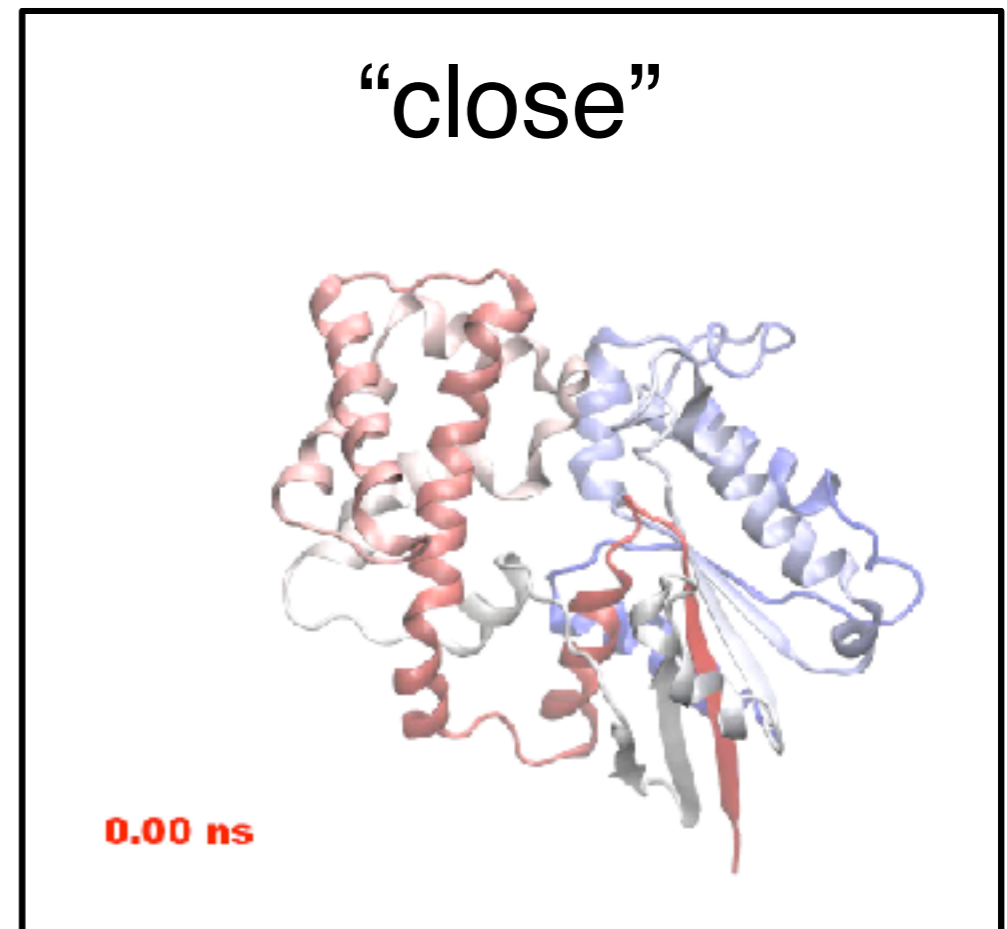
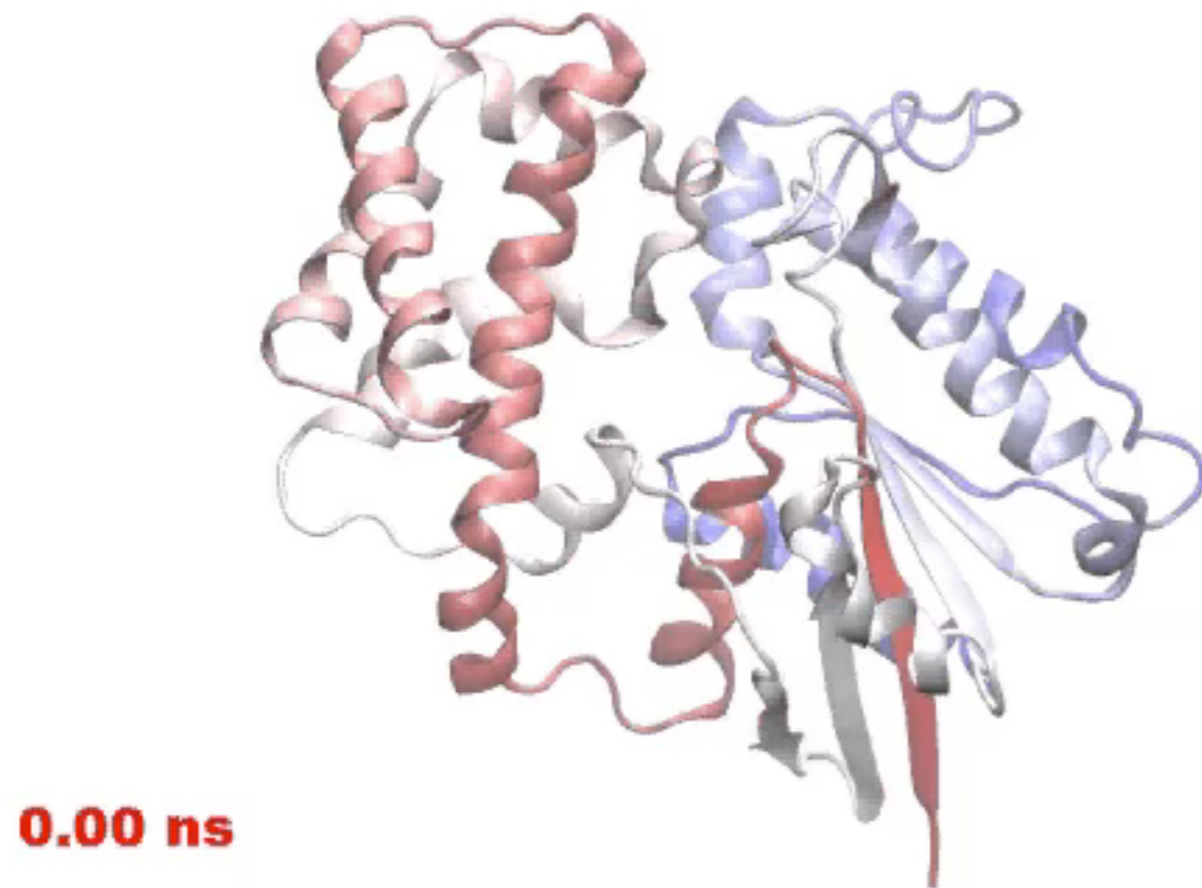


$$\begin{aligned} \mathbf{v}(t + \frac{\Delta t}{2}) &= \mathbf{v}(t - \frac{\Delta t}{2}) + \frac{\mathbf{F}(t)}{m} \Delta t \\ \mathbf{r}(t + \Delta t) &= \mathbf{r}(t) + \mathbf{v}(t + \frac{\Delta t}{2}) \Delta t \end{aligned}$$



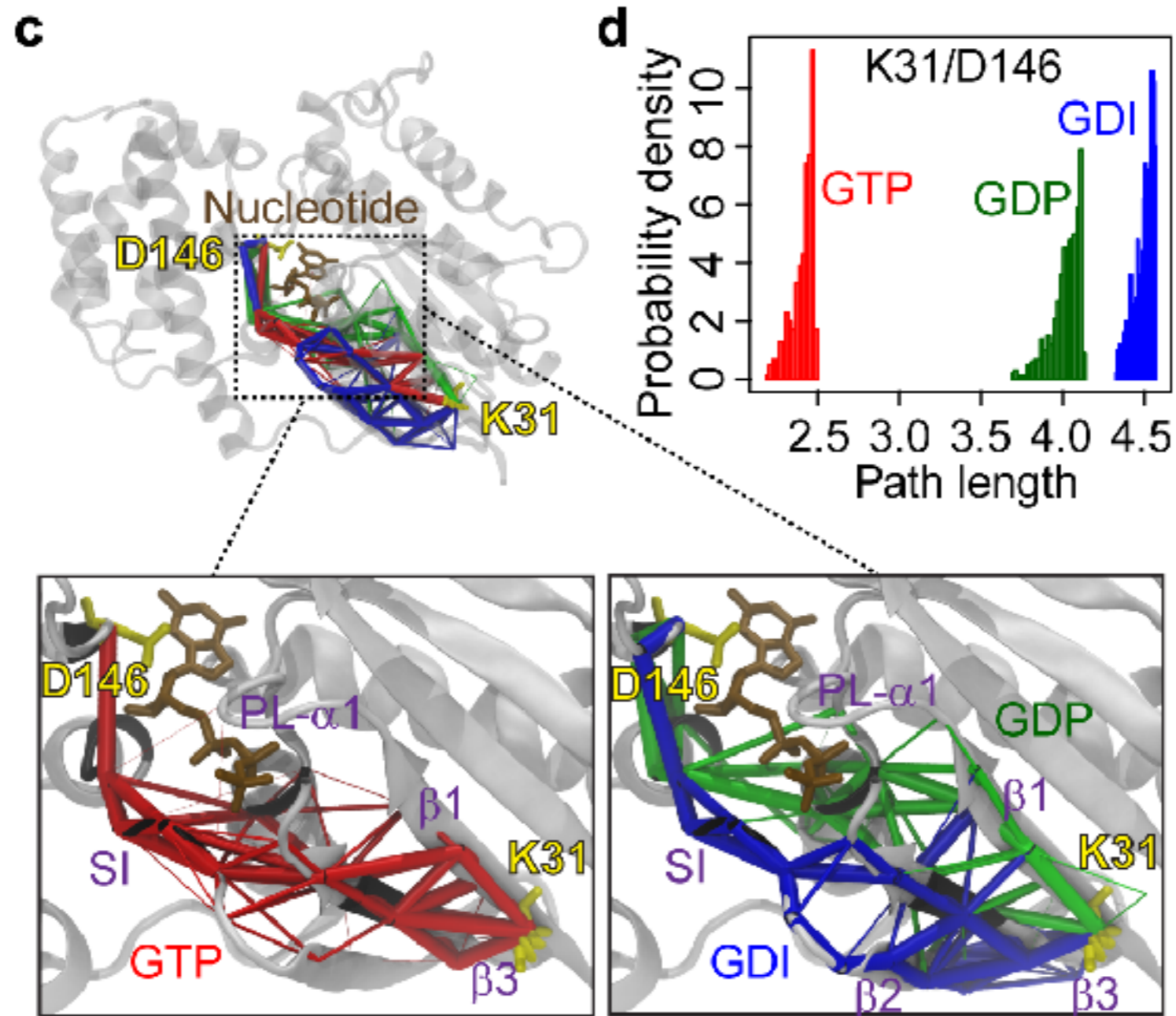
# MD Prediction of Functional Motions

Accelerated MD simulation of  
nucleotide-free transducin alpha subunit



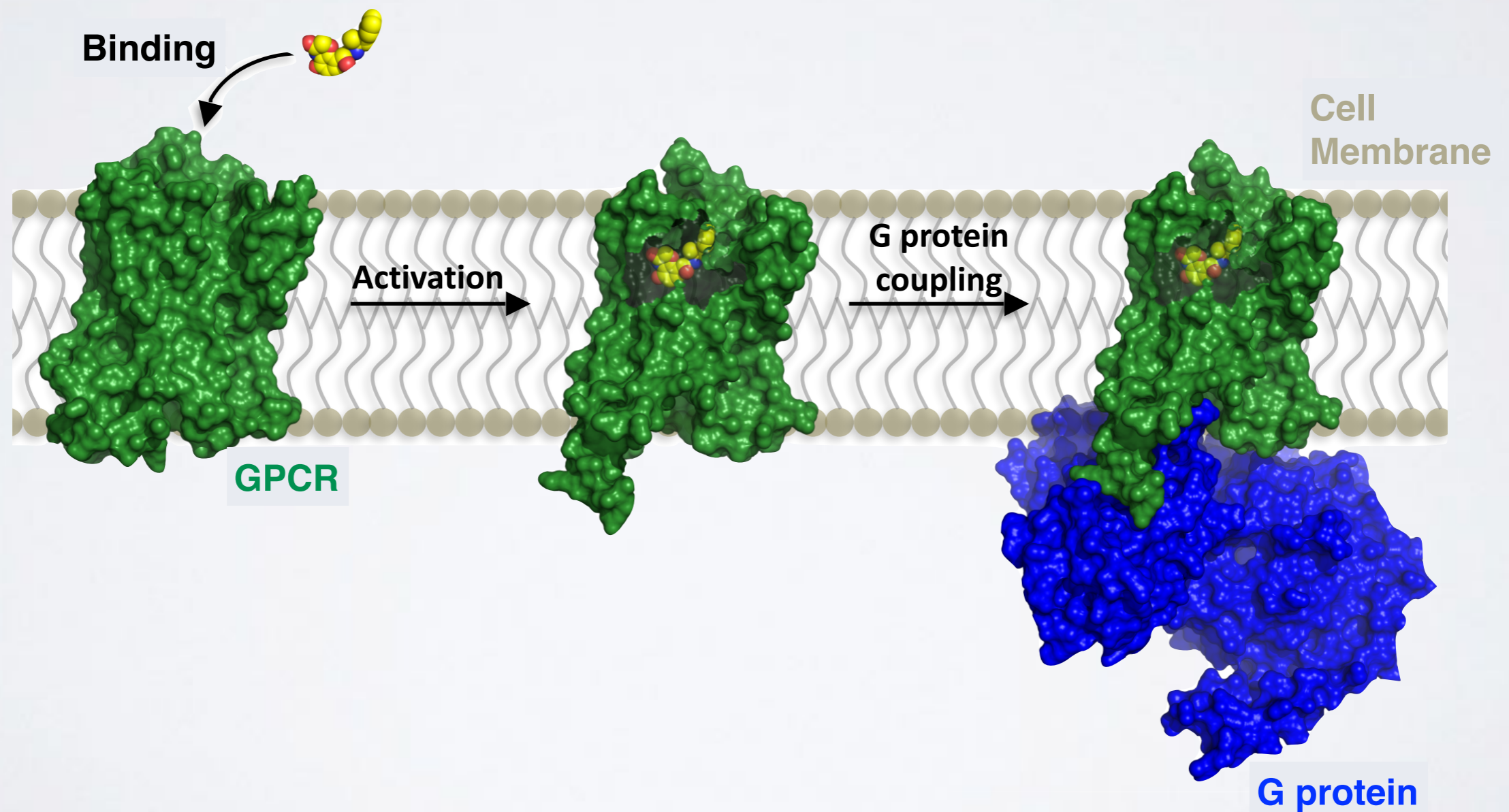
Yao and Grant, Biophys J. (2013)

# Simulations Identify Key Residues Mediating Dynamic Activation

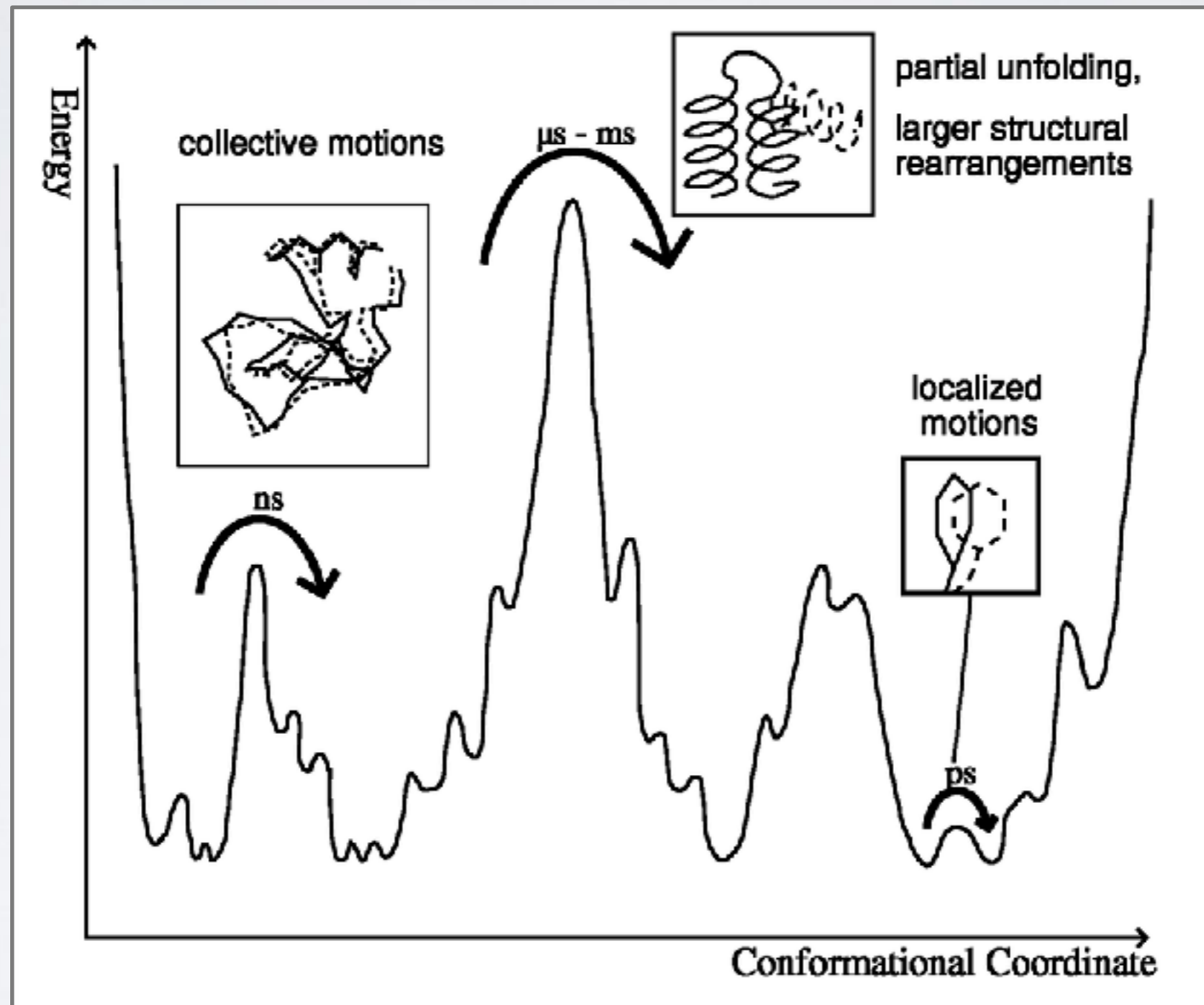




# EXAMPLE APPLICATION OF MOLECULAR SIMULATIONS TO GPCRS



# PROTEINS JUMP BETWEEN MANY, HIERARCHICALLY ORDERED “CONFORMATIONAL SUBSTATES”



H. Frauenfelder et al., *Science* **229** (1985) 337

# MOLECULAR DYNAMICS IS VERY EXPENSIVE

**Example:** F<sub>1</sub>-ATPase in water (183,674 atoms) for 1 nanosecond:

=> 10<sup>6</sup> integration steps

=> 8.4 \* 10<sup>11</sup> floating point operations/step

[n(n-1)/2 interactions]

Total: 8.4 \* 10<sup>17</sup> flop

(on a 100 Gflop/s cpu: **ca 25 years!**)

**... but performance has been improved by use of:**

multiple time stepping ca. 2.5 years

fast multipole methods ca. 1 year

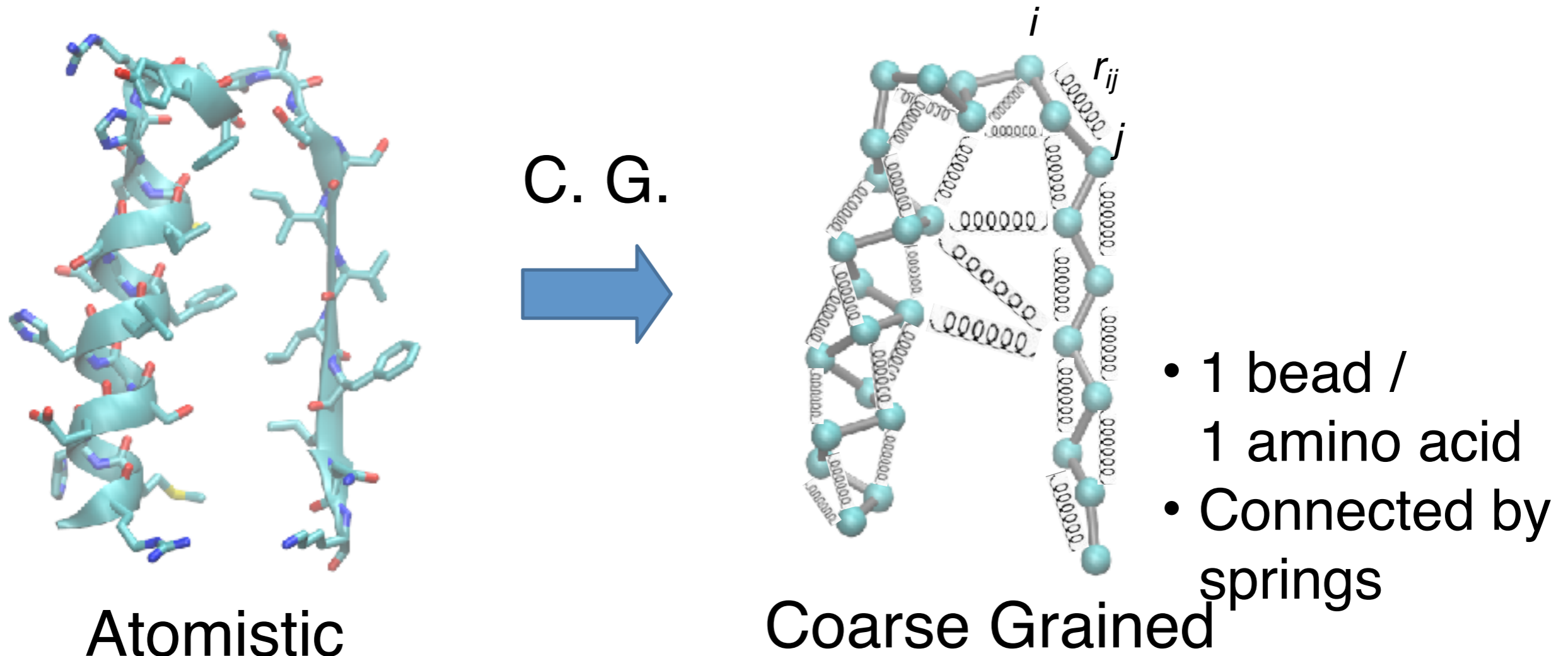
parallel computers ca. 5 days

modern GPUs **ca. 1 day**

**(Anton supercomputer ca. minutes)**

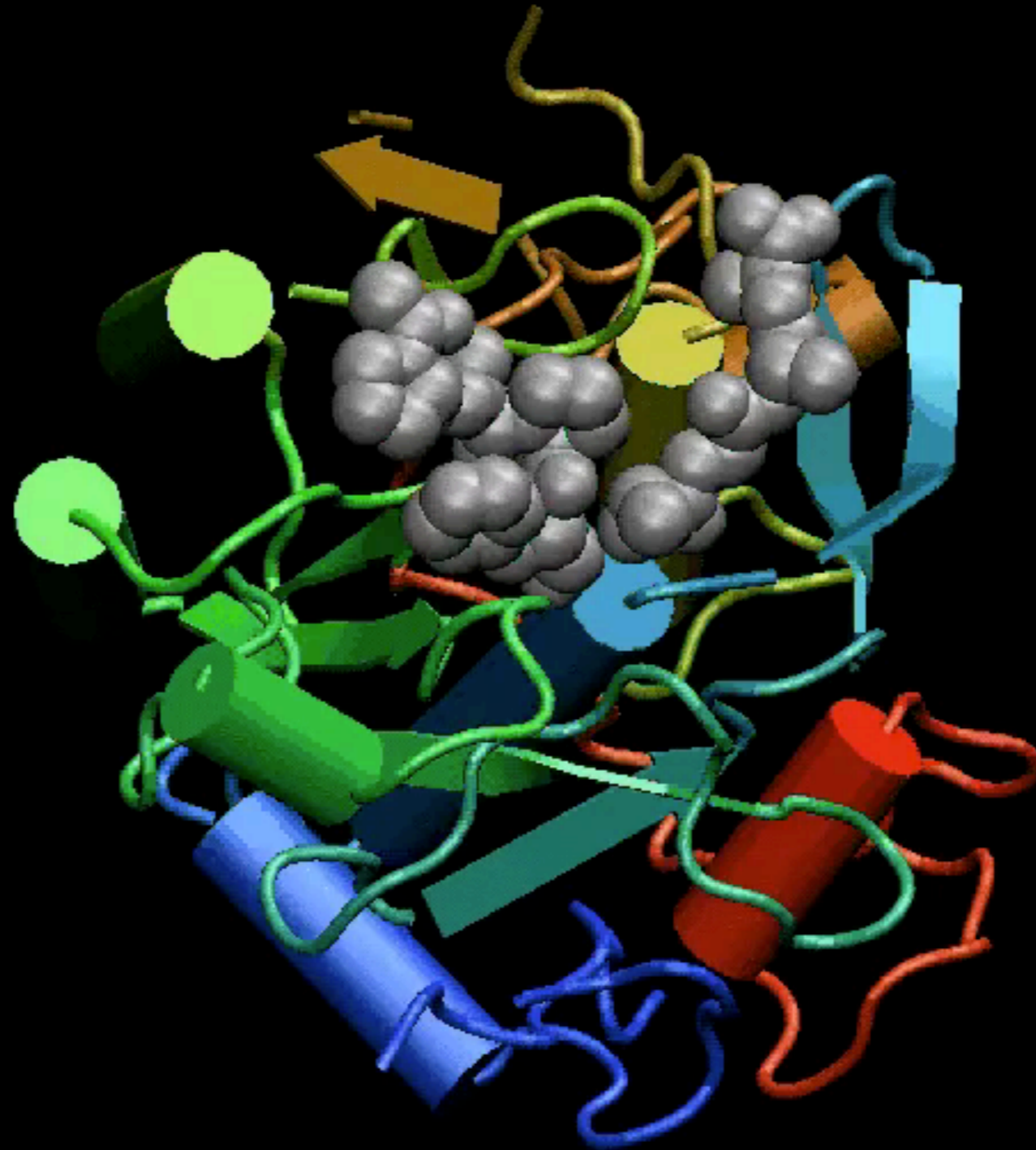
# COARSE GRAINING: **NORMAL MODE ANALYSIS** (NMA)

- MD is still time-consuming for large systems
- Elastic network model NMA (ENM-NMA) is an example of a lower resolution approach that finishes in seconds even for large systems.





NMA models the protein as a network of elastic strings



Proteinase K

Do it Yourself!

# Hand-on time!

[https://bioboot.github.io/bimm143\\_W18/lectures/#12](https://bioboot.github.io/bimm143_W18/lectures/#12)

Focus on **section 3** & **4** exploring **PCA** and **NMA apps**

## ACHIEVEMENTS

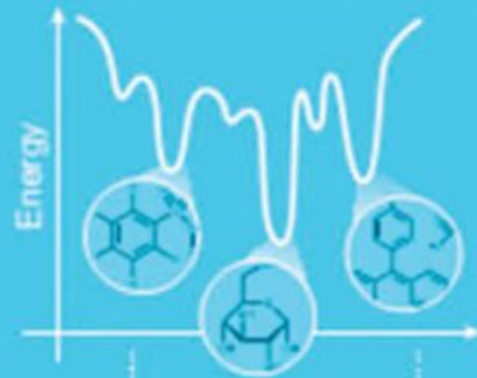
Computational power



Data coverage and community resources



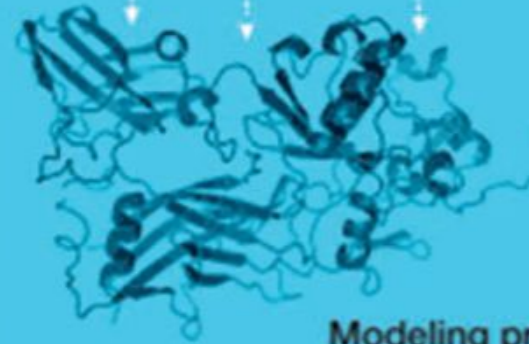
Chemical systems biology and small-molecule docking simulations



Objective method assessment



Correlated mutations



Modeling protein structure

## CHALLENGES

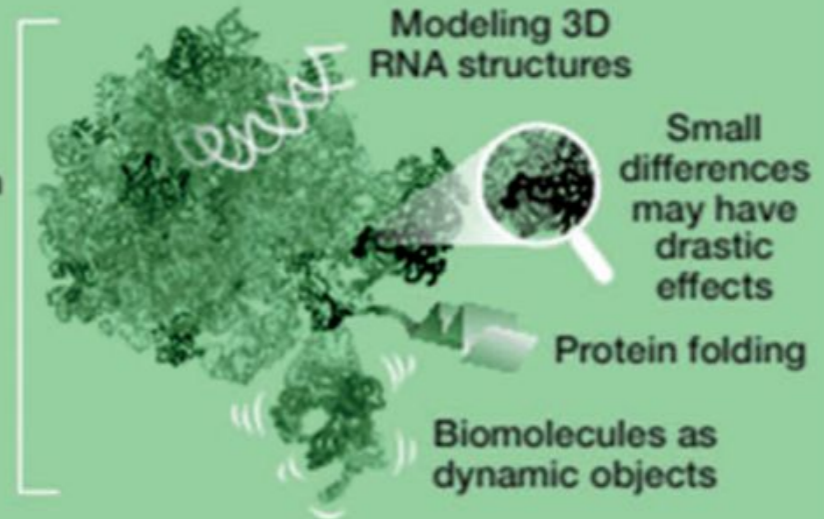
Accessibility and integration of data and methods



Protein engineering and synthetic biology



Modeling multi-domain proteins and large assemblies



Origins and evolution of protein structure

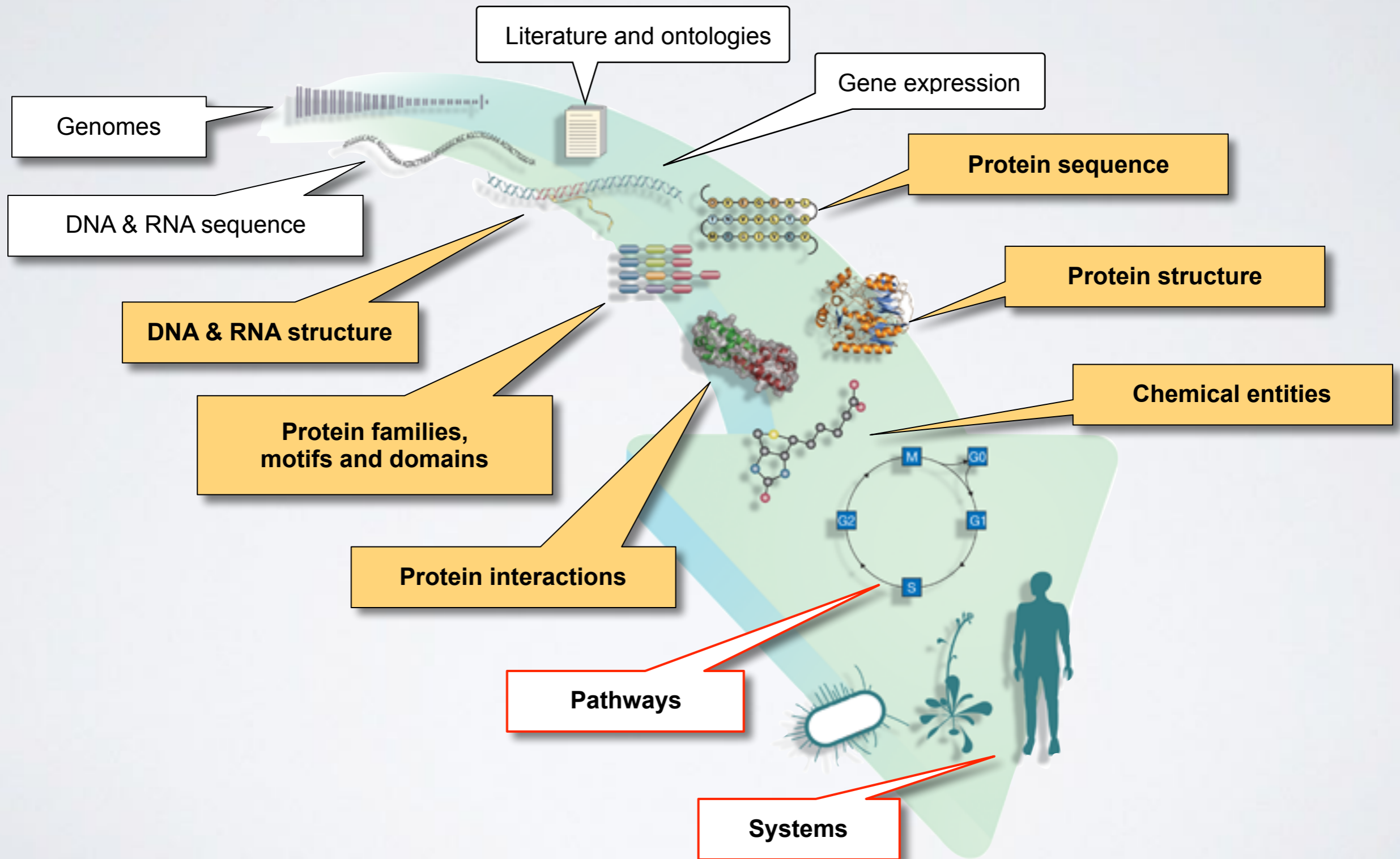


Integration with systems biology





# INFORMING SYSTEMS BIOLOGY?



# SUMMARY

- Structural bioinformatics is computer aided structural biology
- Described major motivations, goals and challenges of structural bioinformatics
- Reviewed the fundamentals of protein structure
- Introduced both physics and knowledge based modeling approaches for describing the structure, energetics and dynamics of proteins computationally
- Introduced both structure and ligand based bioinformatics approaches for drug discovery and design