



# BGGN 213

## Structural Bioinformatics II

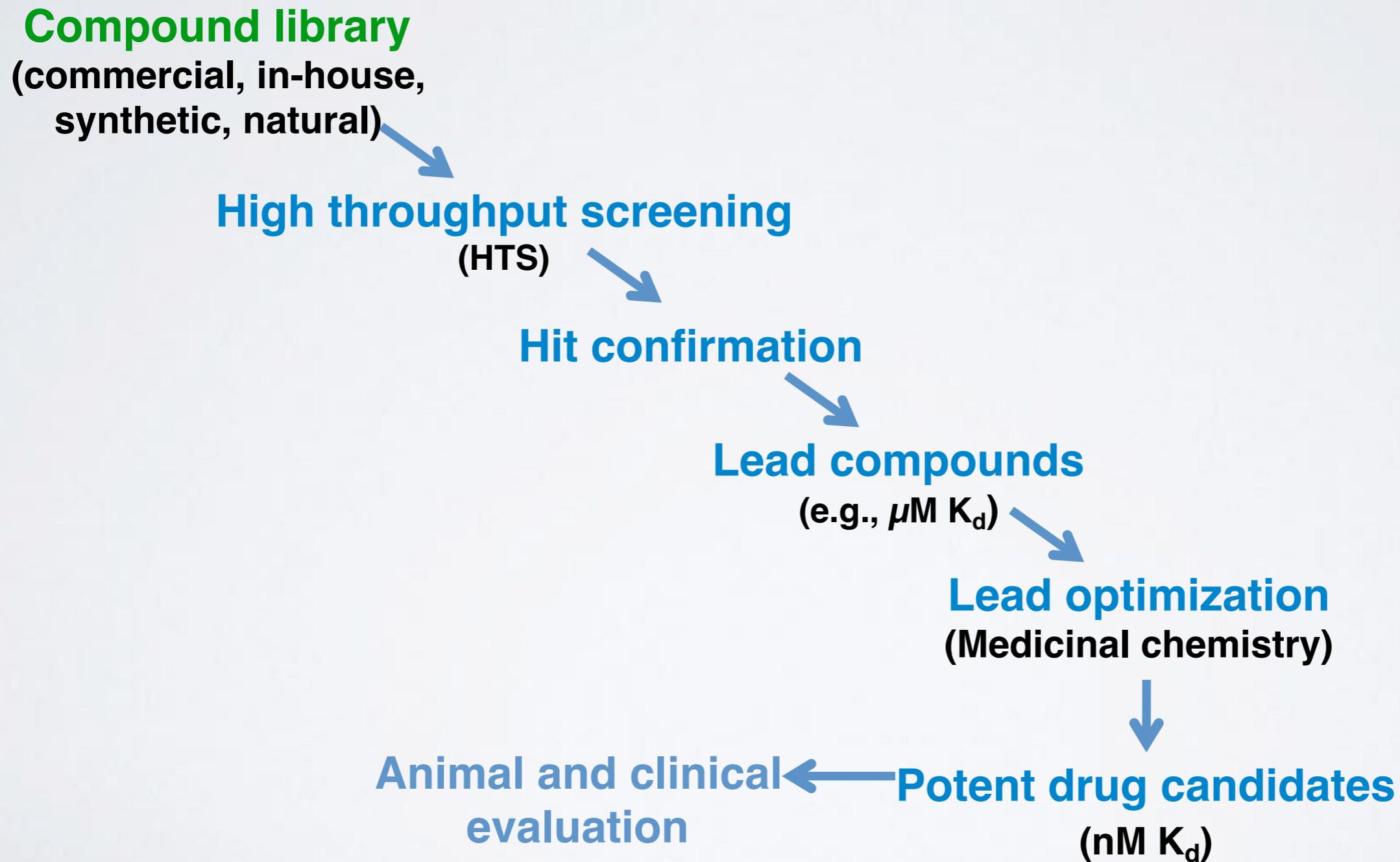
Barry Grant  
UC San Diego

<http://thegrantlab.org/bggn213>

# 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**
  - Predicting functional dynamics & drug discovery

# THE TRADITIONAL EMPIRICAL PATH TO DRUG DISCOVERY



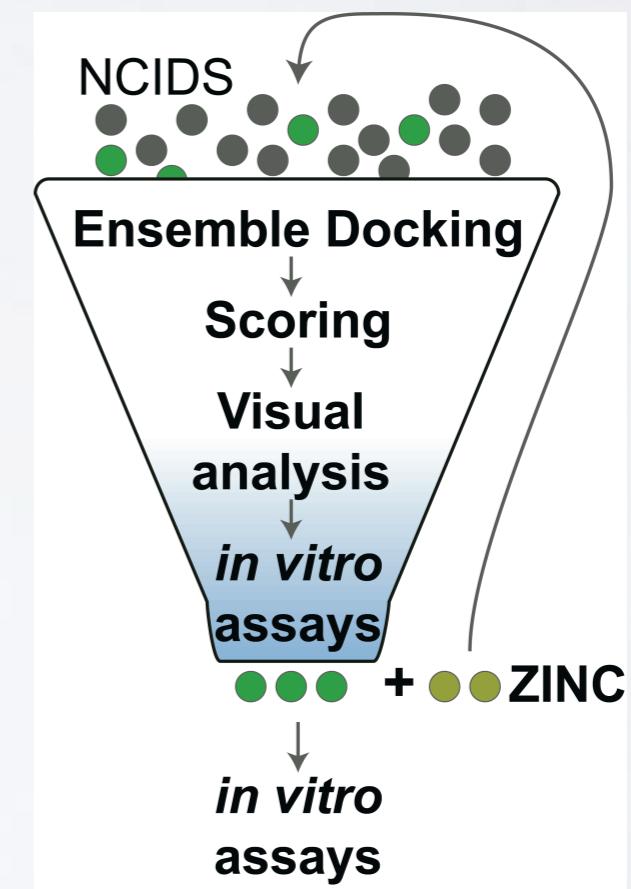
# COMPUTER-AIDED LIGAND DESIGN

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:

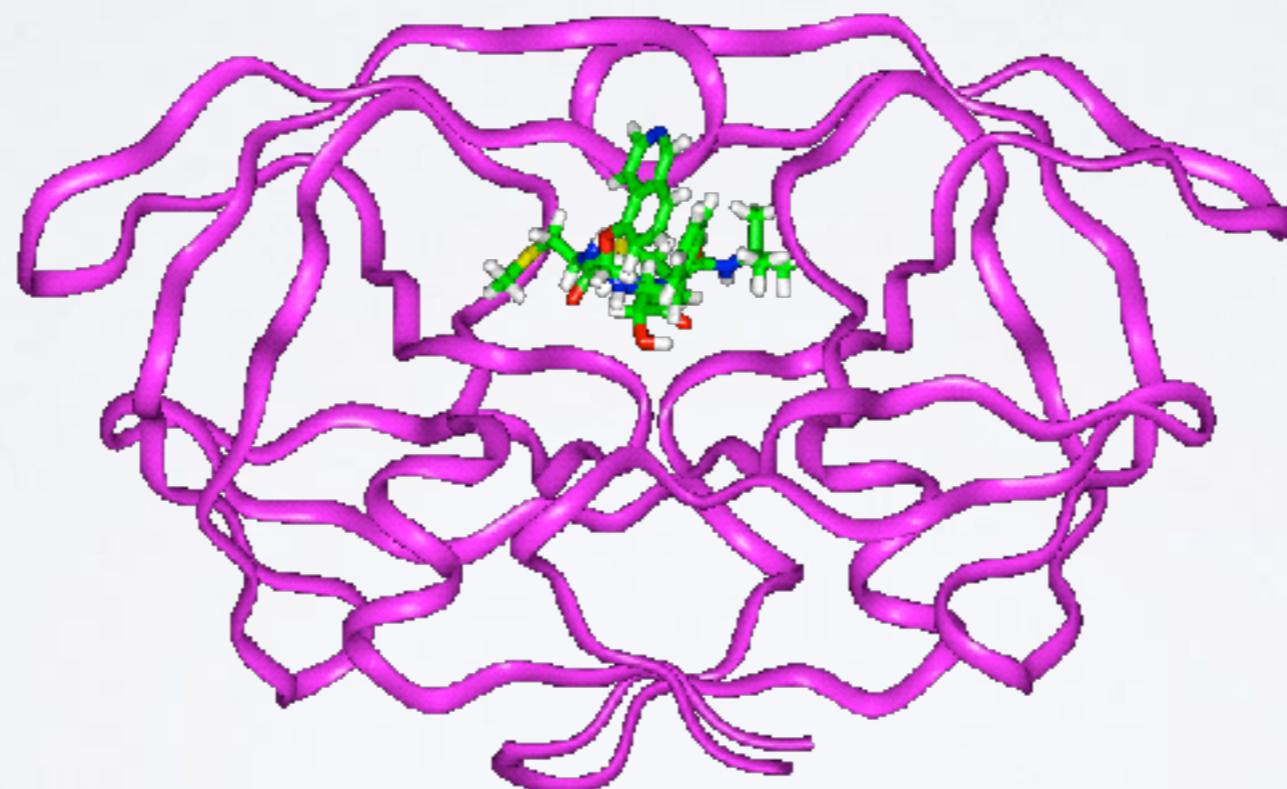
**(1). Receptor/Target-Based**

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

# **SCENARIO I:**

## RECEPTOR-BASED DRUG DISCOVERY

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



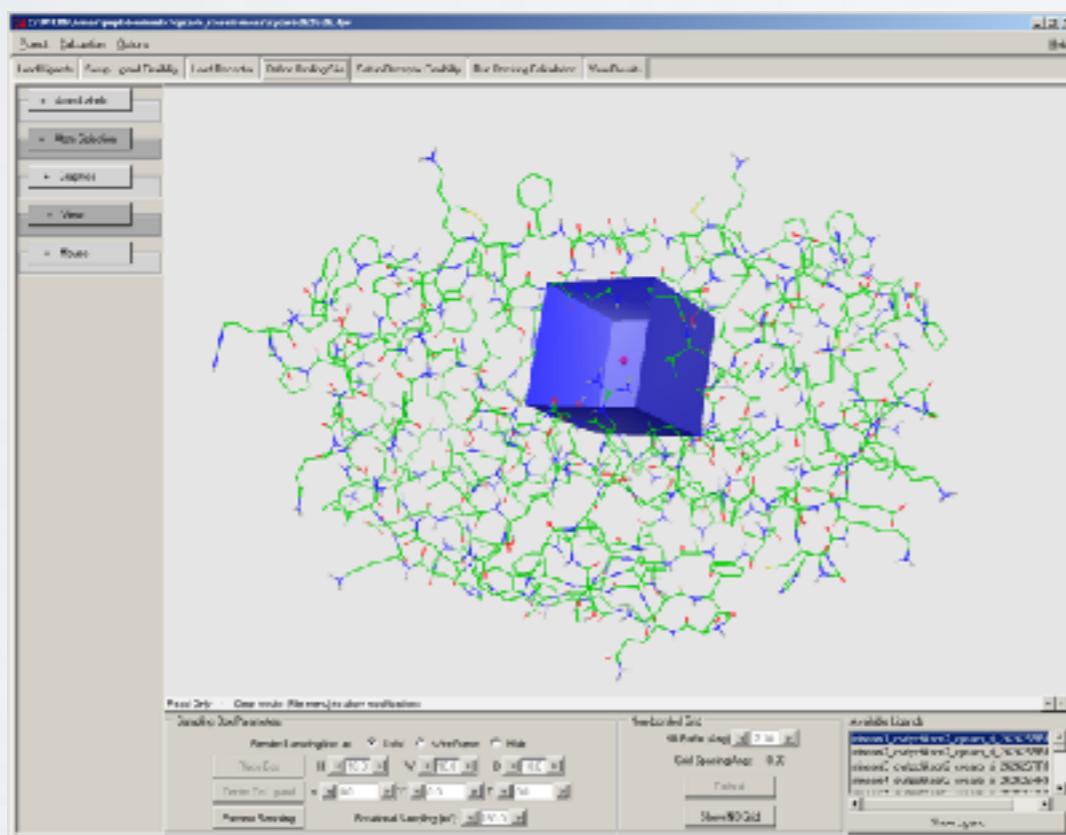
HIV Protease/KNI-272 complex

# PROTEIN-LIGAND DOCKING

## Structure-Based Ligand Design

Docking software

Search for structure of lowest energy



Potential function

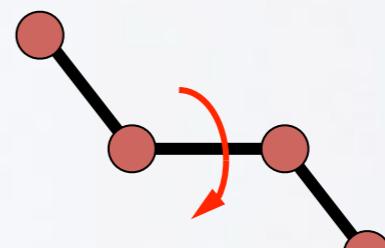
Energy as function of structure



VDW

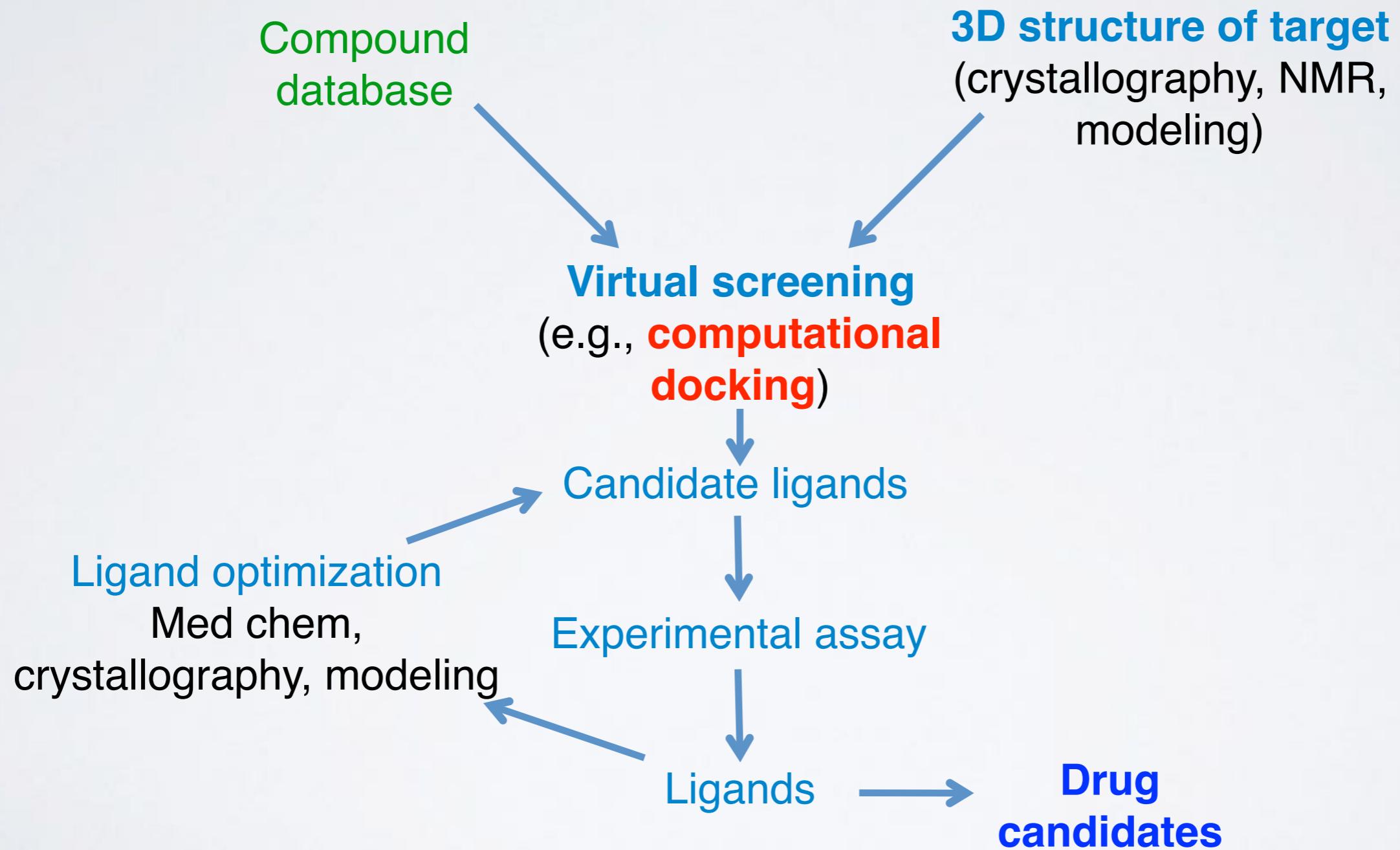


Screened Coulombic



Dihedral

# STRUCTURE-BASED VIRTUAL SCREENING



# COMPOUND LIBRARIES

The screenshot shows the Maybridge website. At the top, there's a search bar and a navigation menu with links like "HOME", "SCREENING SERVICES", "INDUSTRIAL SCREENING", "INDUSTRY", "INDUSTRY", and "CONTACT US". Below the header, there's a banner for "Maybridge HitFinder™". The main content area features a heading "Maybridge HitFinder™" and a sub-section "The pre-selected diverse screening library includes identifying potential drug leads easy, universal, and cost effective." It includes a section on "Screening quality data from your screens" and a "Ready to Screen" section with a grid of small molecule images.

The screenshot shows the NIH Molecular Libraries Small Molecule Repository website. The header includes the NIH logo and the text "NIH MOLECULAR LIBRARIES SMALL MOLECULE REPOSITORY". The main content area has a heading "A NIH Roadmap Initiative" and a "Welcome" section. It features a photograph of a scientist in a lab setting. The sidebar contains links for "Home", "MSMR Project", "MSMR Details", and "Submit Compounds". The footer includes the BioFocus logo and the text "BioFocus, a Galapagos company".

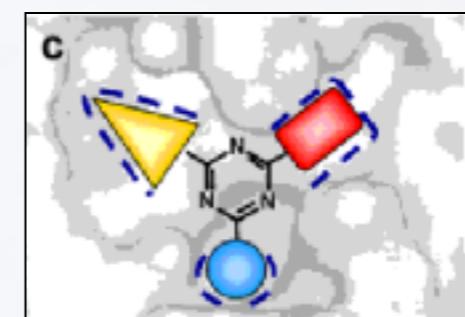
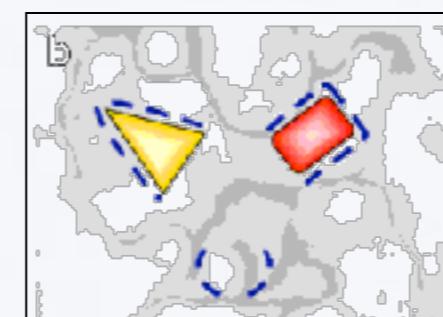
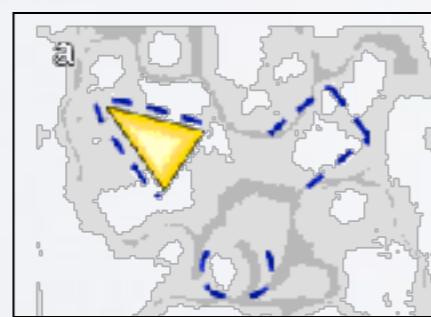
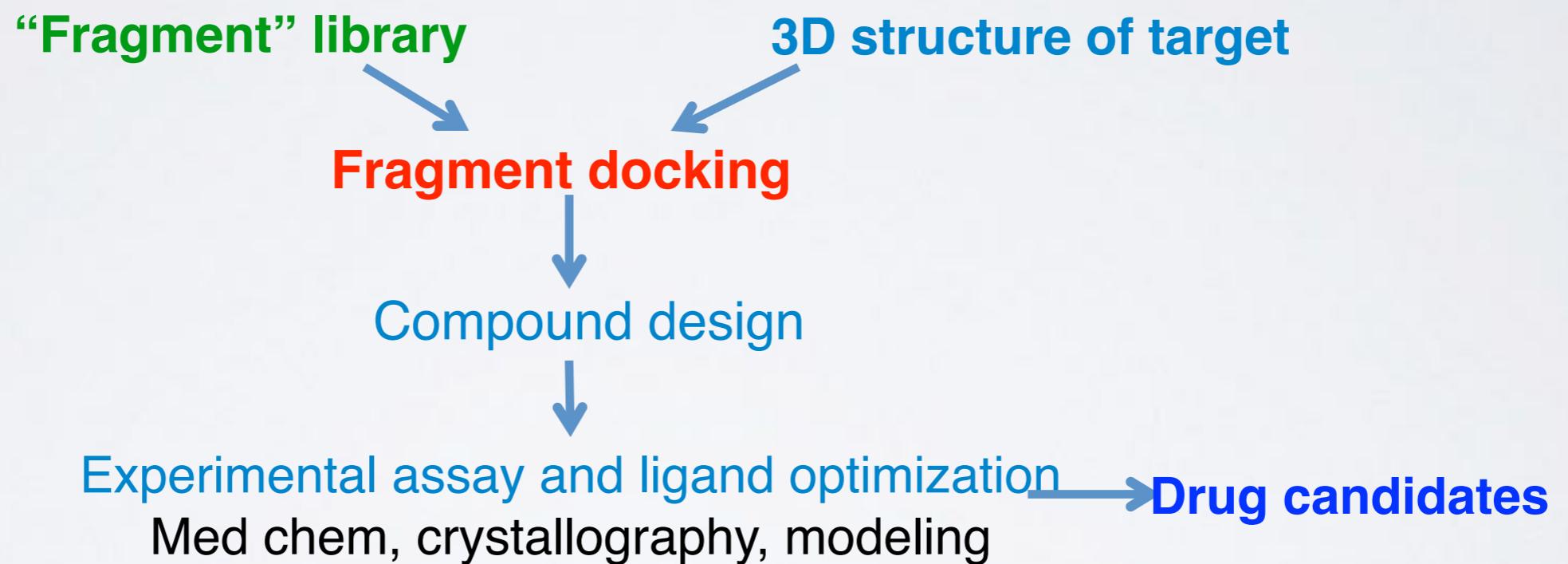
The screenshot shows the PMLSC website. The header includes the University of Pittsburgh logo and the text "University of Pittsburgh" and "Pittsburgh Molecular Libraries Screening Center". The main content area features a large image with the text "BIG DISCOVERIES FROM SMALL MOLECULES". The sidebar contains links for "HOME", "HISTORY", "PERSONNEL", "SCREENING TECHNOLOGY", "COMPONENTS", "RESEARCH & PUBLICATIONS", "LITERATURE", "ASSAY/PROTOLAB ASSAY PROTOCOLS", "PMLSC PROTO REPORTS", "LDR INDEX", "DATA ANALYSIS/INFORMATICS", "EDUCATIONAL ACTIVITIES", "MEMBERSHIPS", "LINKS", "CONTACTS", and "Corporate Search". The footer includes links for "Health Sciences & Pitt", "UPMC", "PSU", "School of Medicine", "Health Sciences Calendar", "Our News & Events", "Up or page | home | contact us", and "© 2008 by the Center for the Molecular Sciences, University of Pittsburgh. All rights reserved".

Commercial  
(in-house pharma)

Government (NIH)

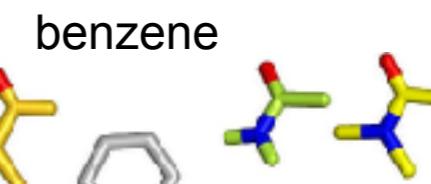
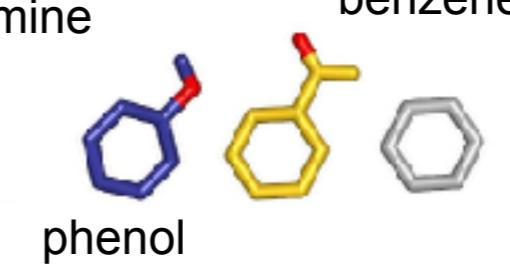
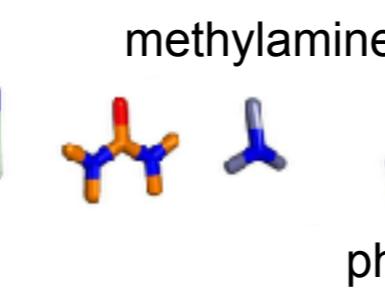
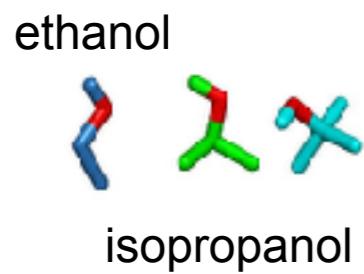
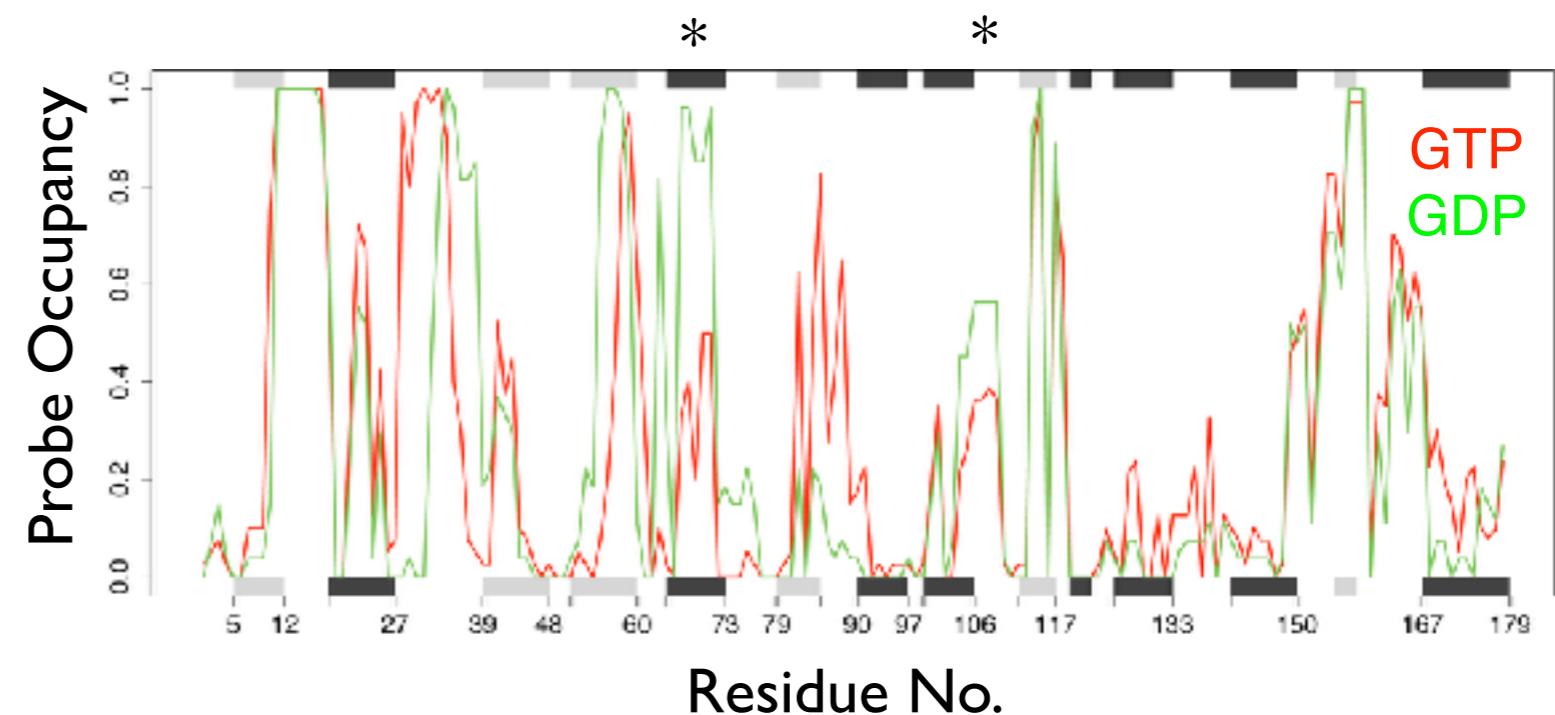
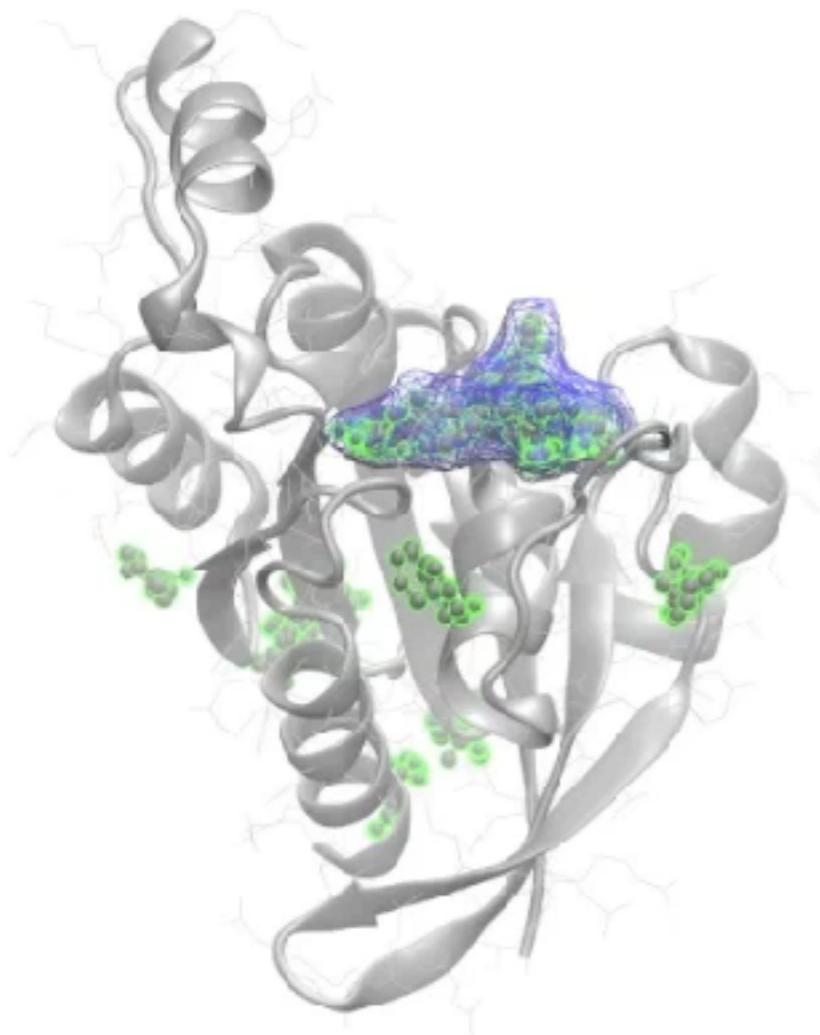
Academia

# FRAGMENTAL STRUCTURE-BASED SCREENING



# Multiple non active-site pockets identified

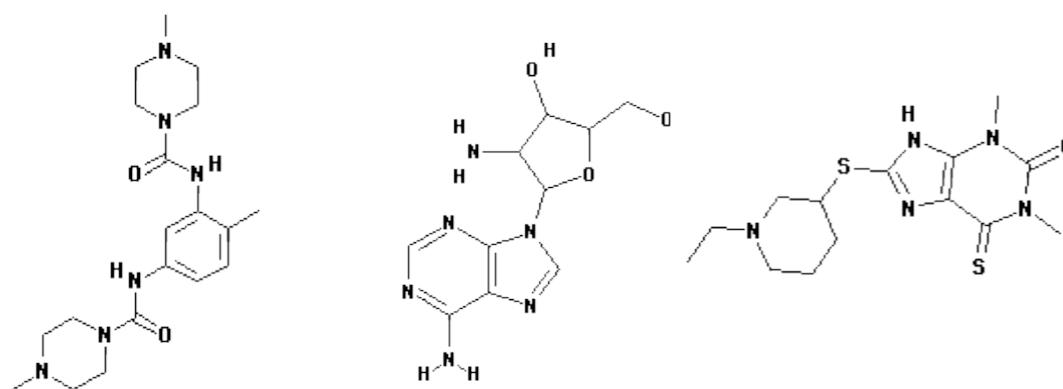
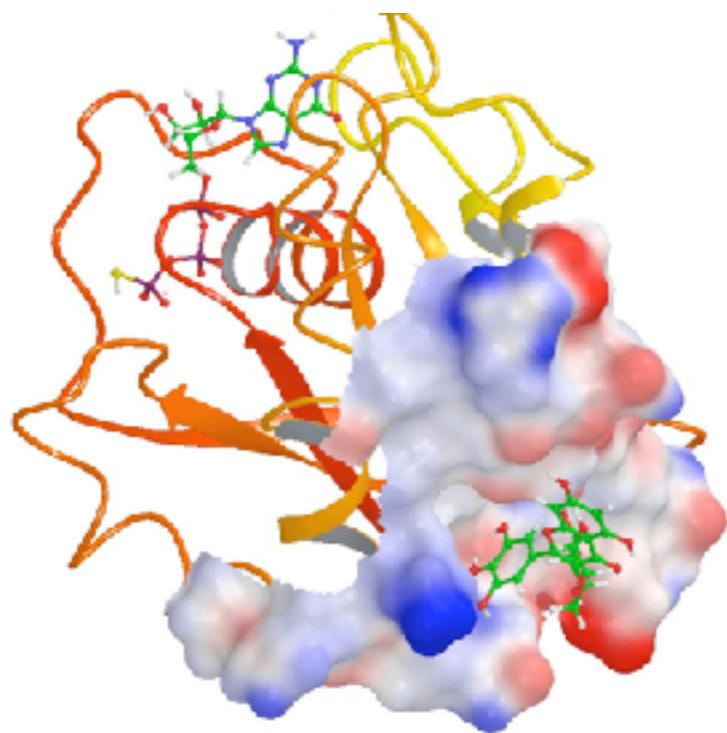
Small organic probe fragment affinities map multiple potential binding sites across the structural ensemble.



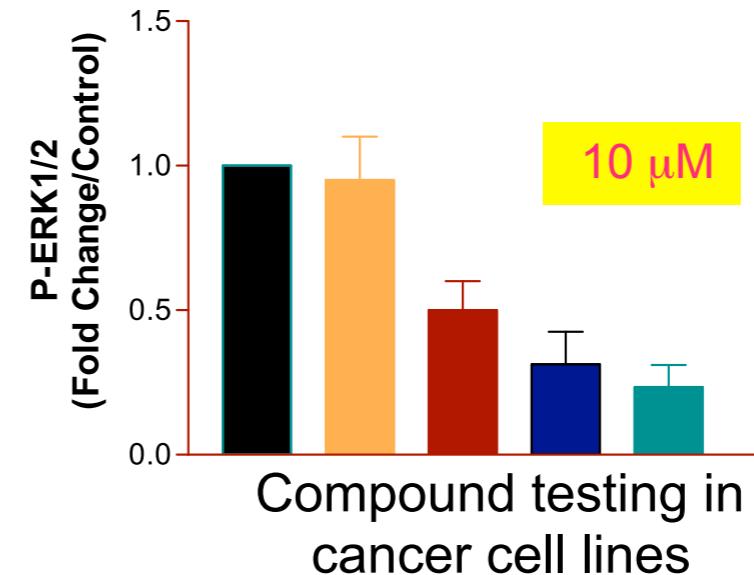
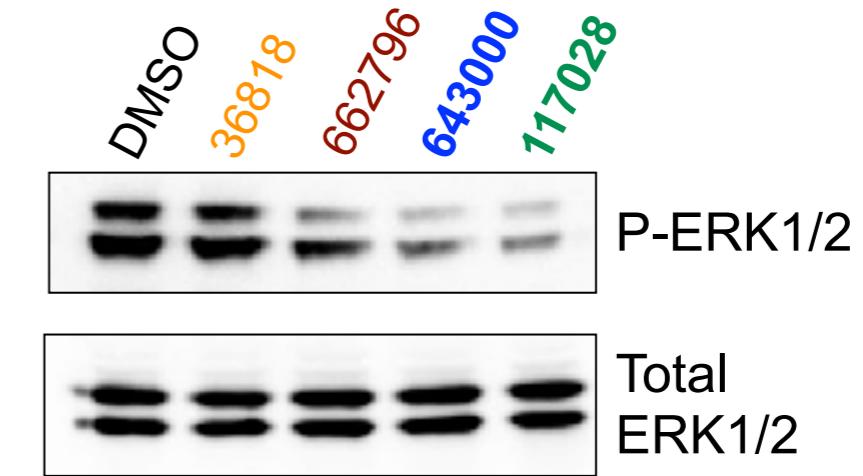
# Ensemble docking & candidate inhibitor testing

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

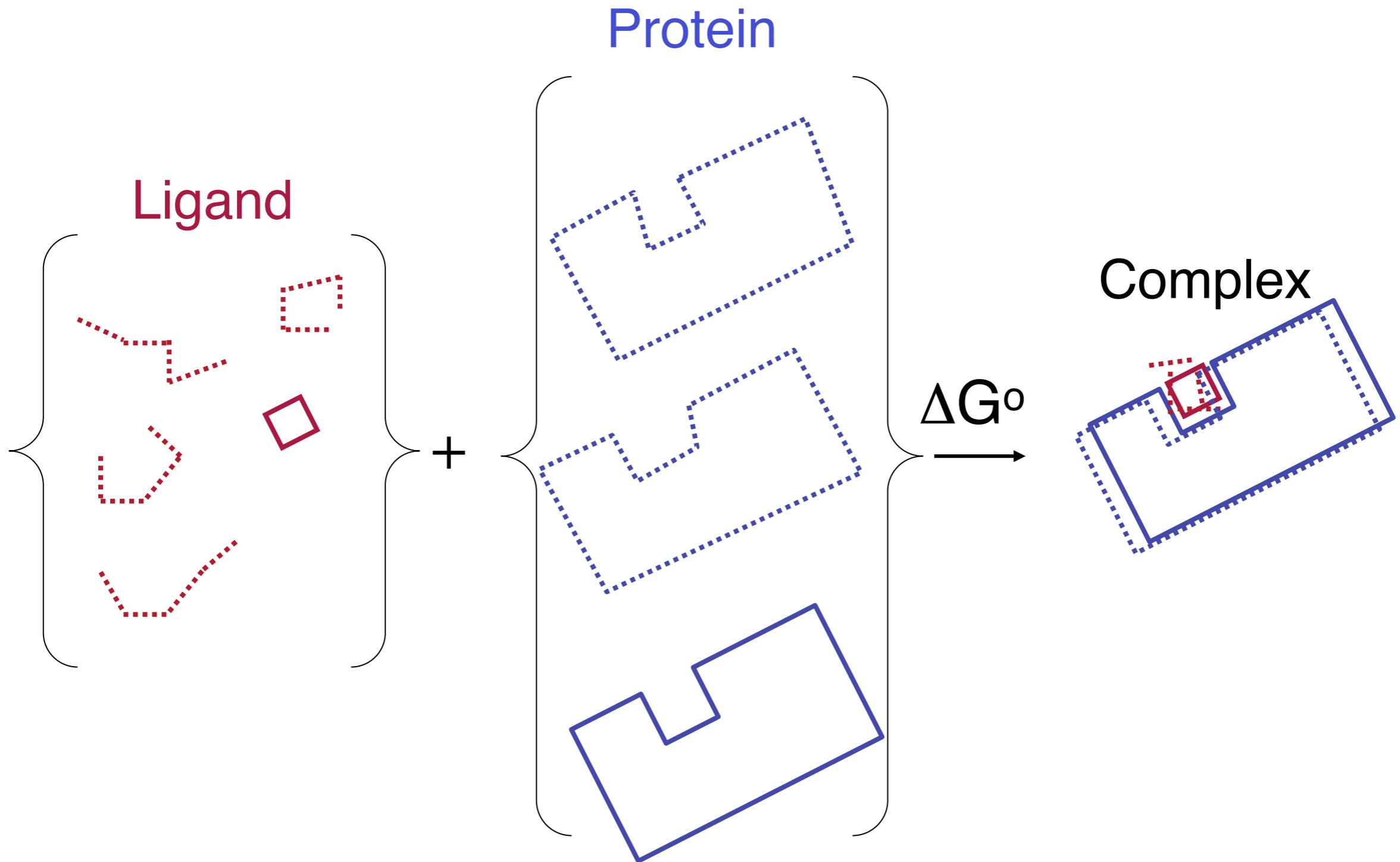
Ensemble computational docking



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

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

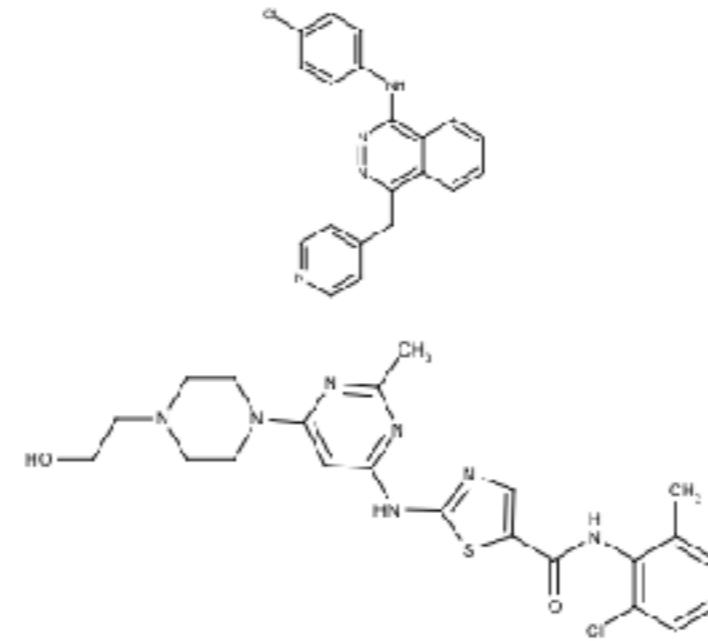
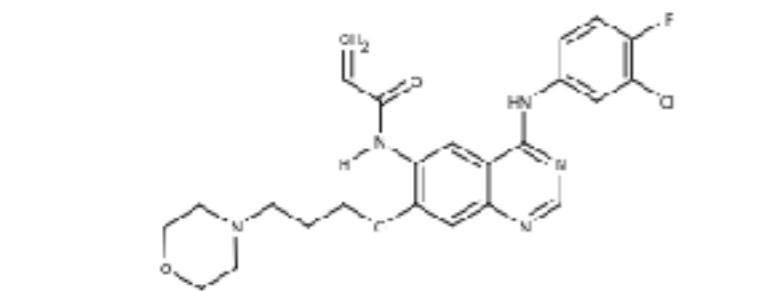
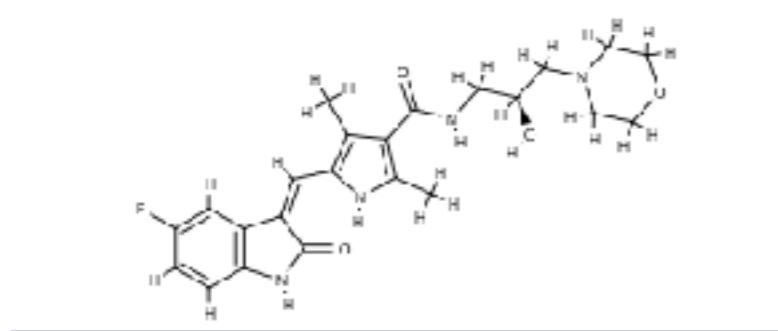
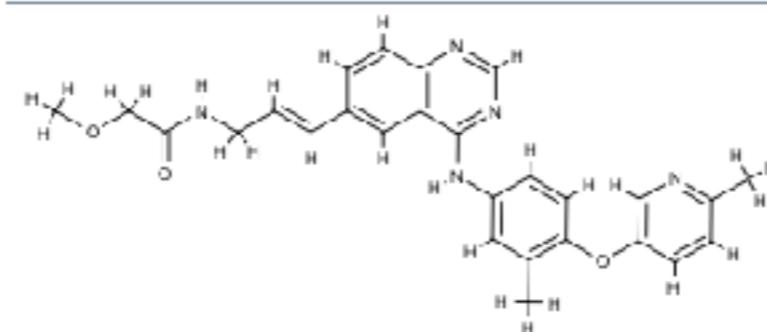
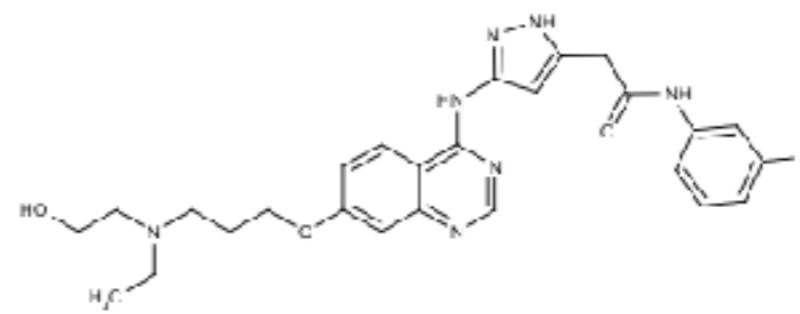
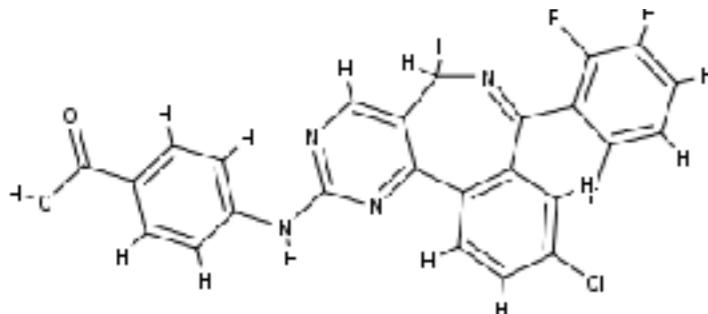
A company wants to work around another company's chemical patents

A high-affinity ligand is toxic, is not well-absorbed, etc.

# Scenario 2

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

e.g. MAP Kinase Inhibitors



Using knowledge of existing inhibitors to discover more

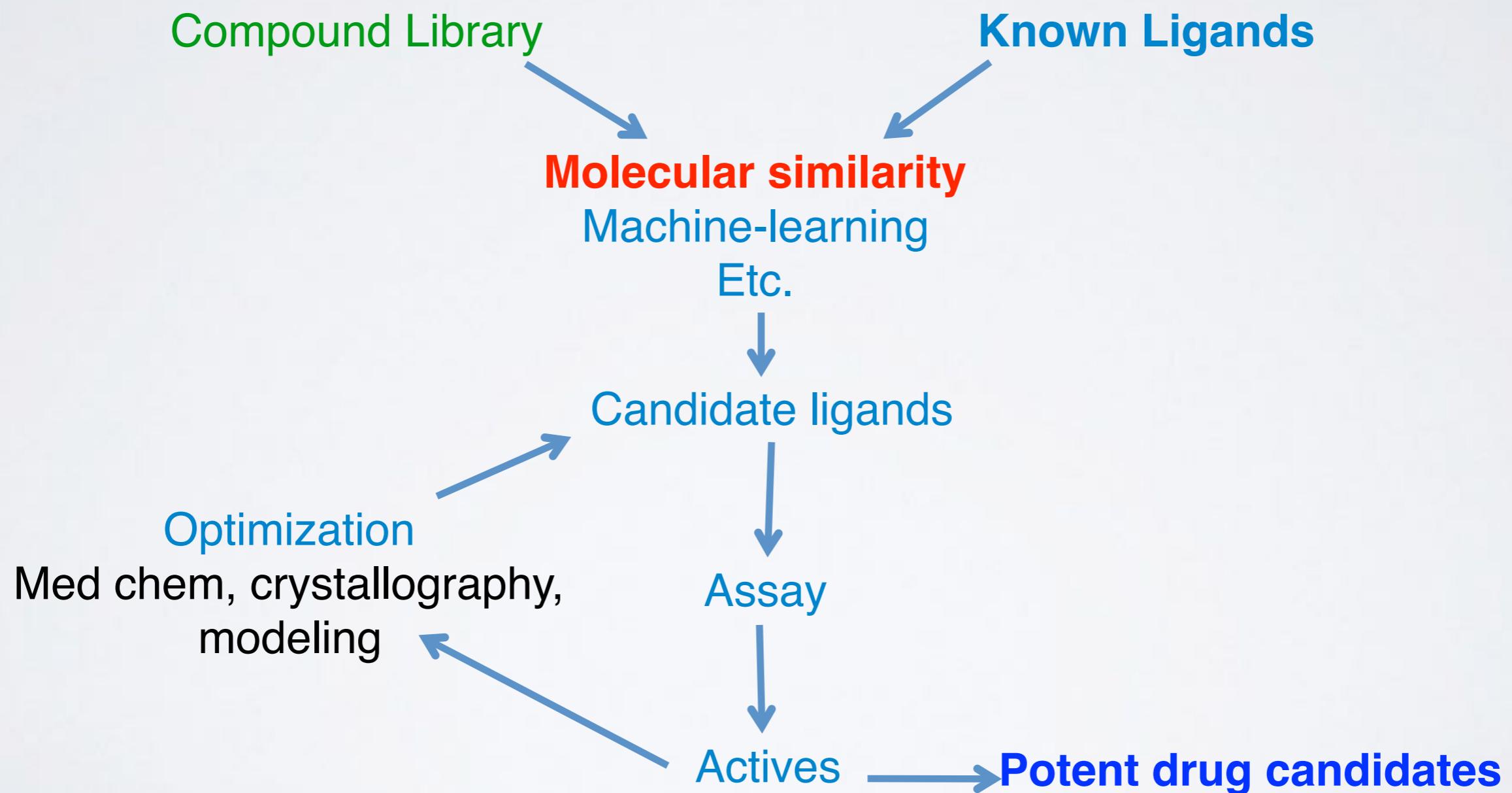
# Why Look for Another Ligand if You Already Have Some?

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

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

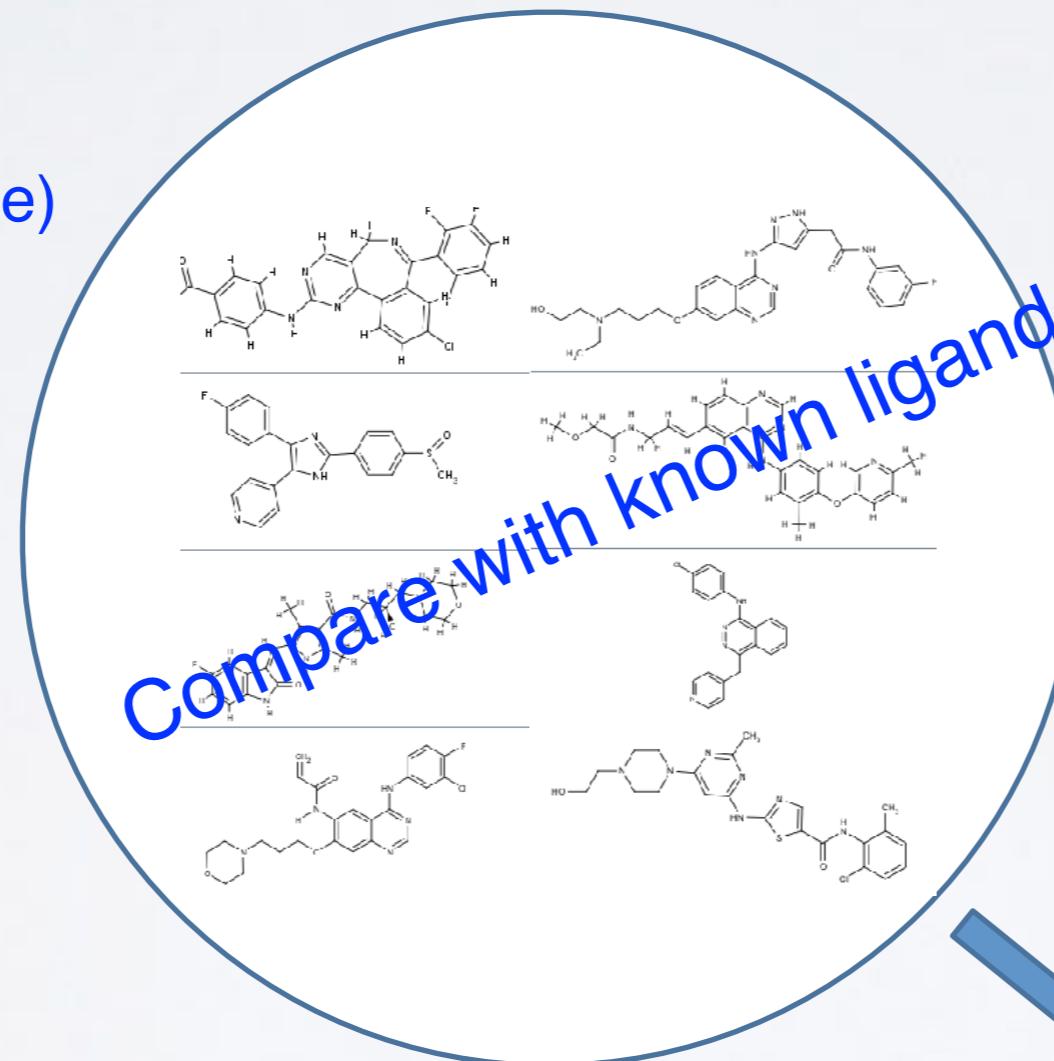
An high-affinity ligand is toxic, is not well-absorbed, etc.

# LIGAND-BASED VIRTUAL SCREENING



# CHEMICAL SIMILARITY LIGAND-BASED DRUG-DISCOVERY

Compounds  
(available/synthesizable)



Different

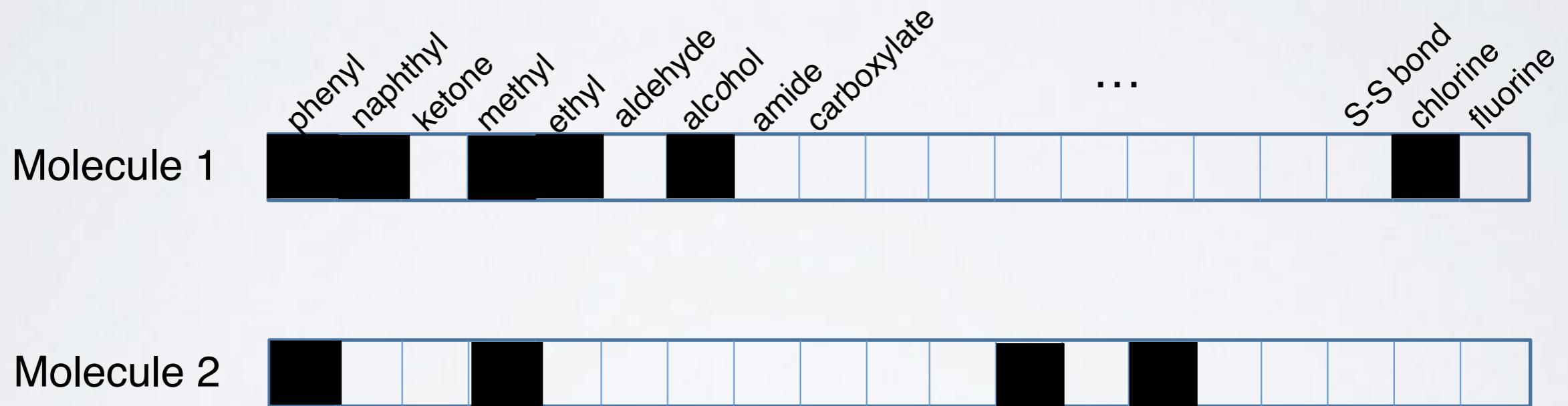
Don't bother

Similar

Test experimentally

# CHEMICAL FINGERPRINTS

## BINARY STRUCTURE KEYS

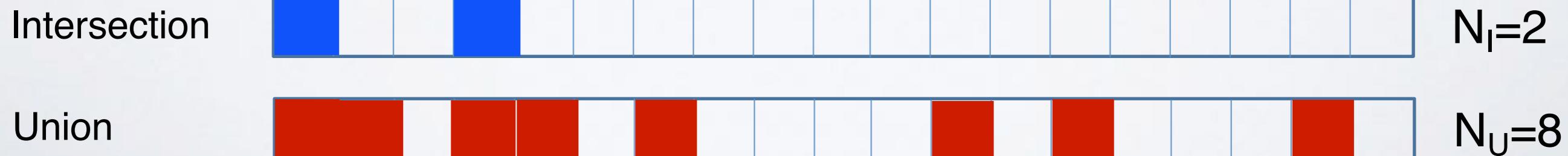


# CHEMICAL SIMILARITY FROM FINGERPRINTS



Tanimoto Similarity  
(or Jaccard Index),  $T$

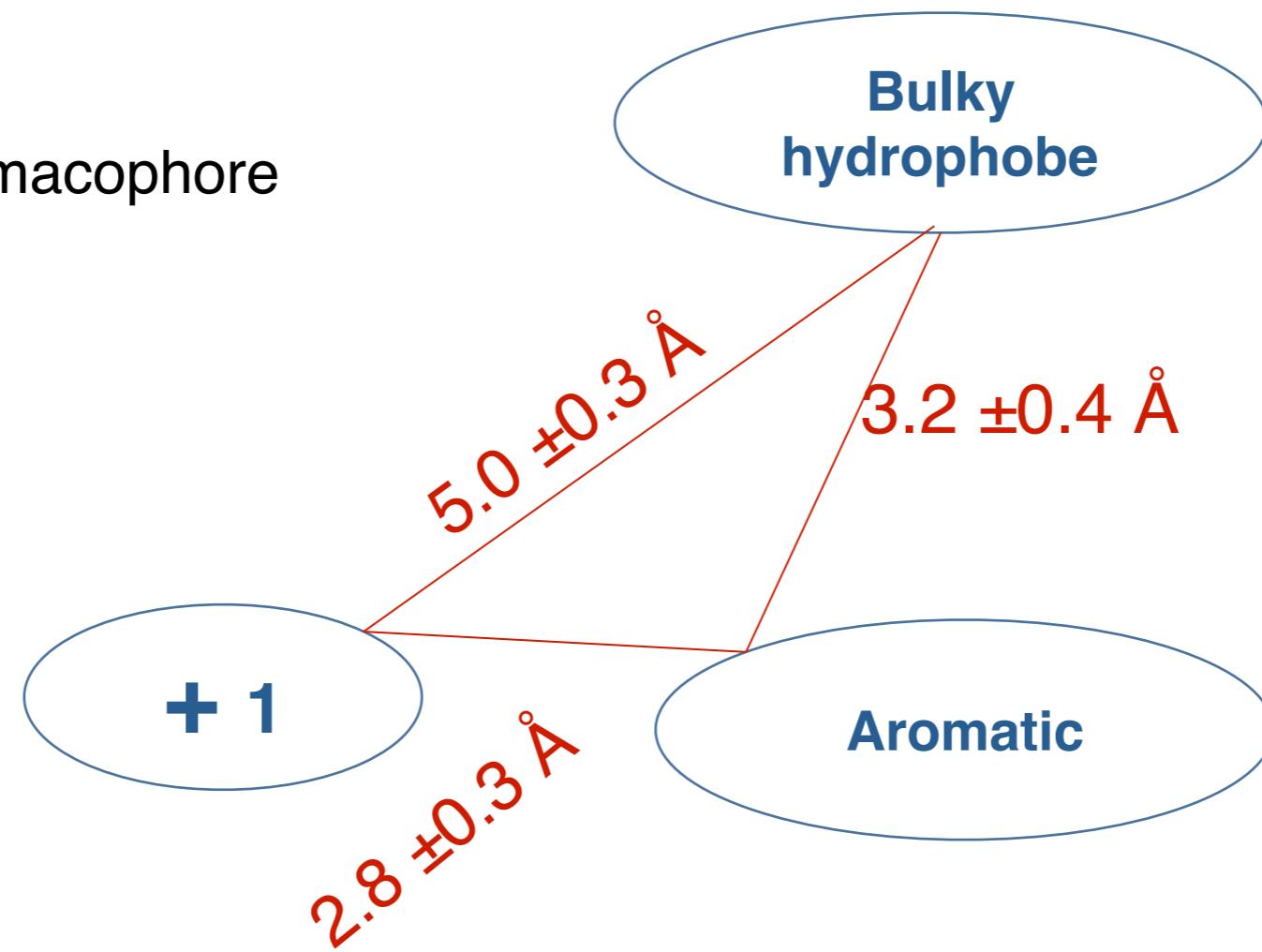
$$T \equiv \frac{N_I}{N_U} = 0.25$$



# Pharmacophore Models

Φάρμακο (drug) + Φορά (carry)

A 3-point pharmacophore



# Molecular Descriptors

## More abstract than chemical fingerprints

### Physical descriptors

molecular weight

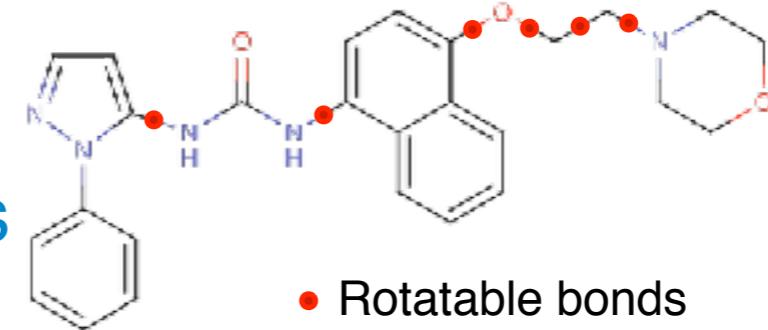
charge

dipole moment

number of H-bond donors/acceptors

number of rotatable bonds

hydrophobicity ( $\log P$  and  $c\log P$ )



• Rotatable bonds

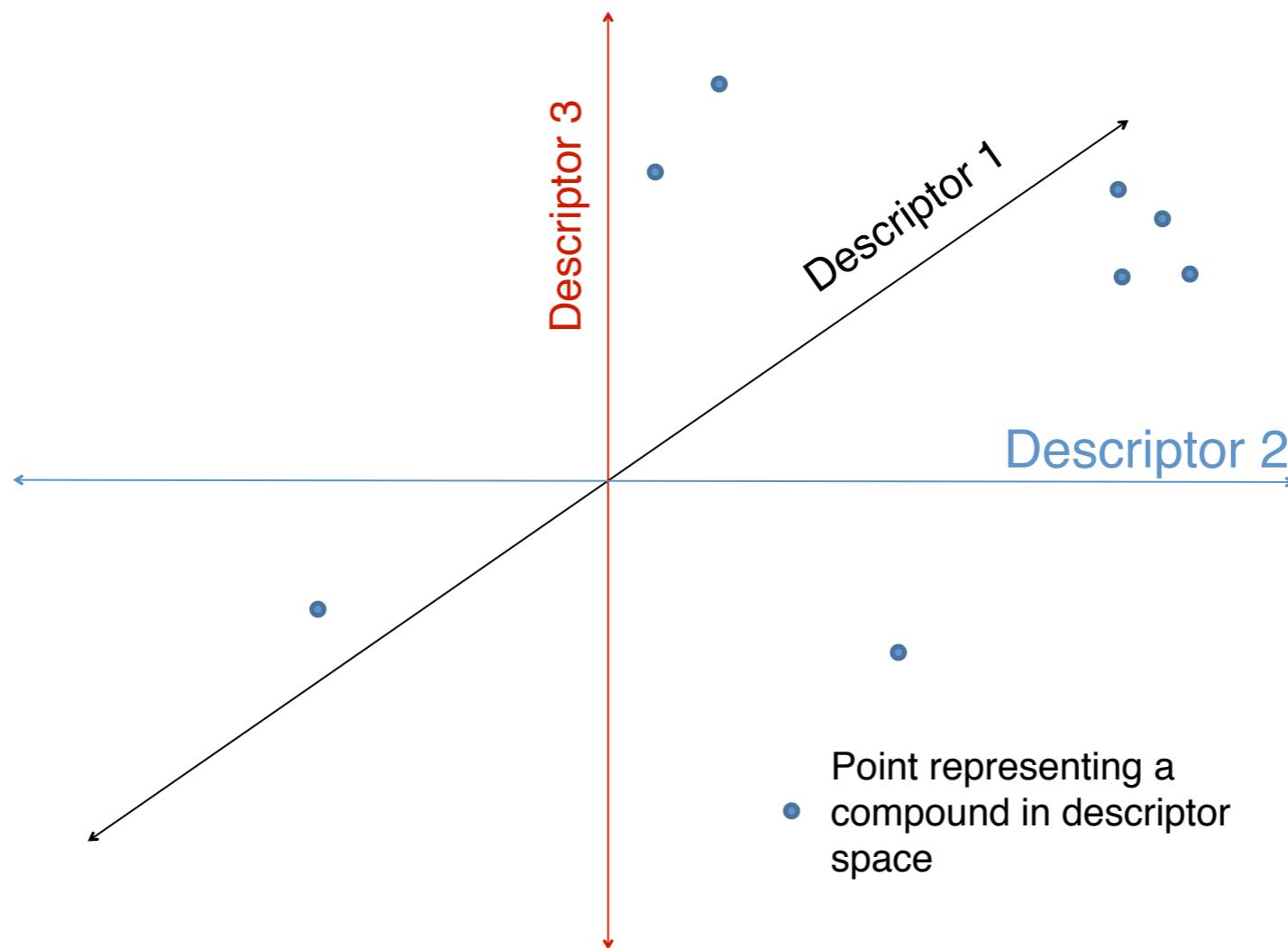
Topological  
branching index  
measures of linearity vs interconnectedness

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, support vector machines, random forest, etc.)

# CAUTIONARY NOTES

- “**Everything should be made as simple as it can be but not simpler**”

A model is **never perfect**. A model that is not quantitatively accurate in every respect does not preclude one from establishing results relevant to our understanding of biomolecules as long as the biophysics of the model are properly understood and explored.

- **Calibration of the parameters is an ongoing and imperfect process**

Questions and hypotheses should always be designed such that they do not depend crucially on the precise numbers used for the various parameters.

- **A computational model is rarely universally right or wrong**

A model may be accurate in some regards, inaccurate in others. These subtleties can only be uncovered by comparing to all available experimental data.

Do it Yourself!

# Hand-on time!

[https://bioboot.github.io/bggn213\\_f17/lectures/#12](https://bioboot.github.io/bggn213_f17/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**

# 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

# Summary

Overview of drug discovery

Computer-aided methods

Structure-based

Ligand-based

Interaction potentials

Physics-based

Knowledge-based (data driven)

Ligand-protein databases, machine-readable chemical formats

Ligand similarity and beyond

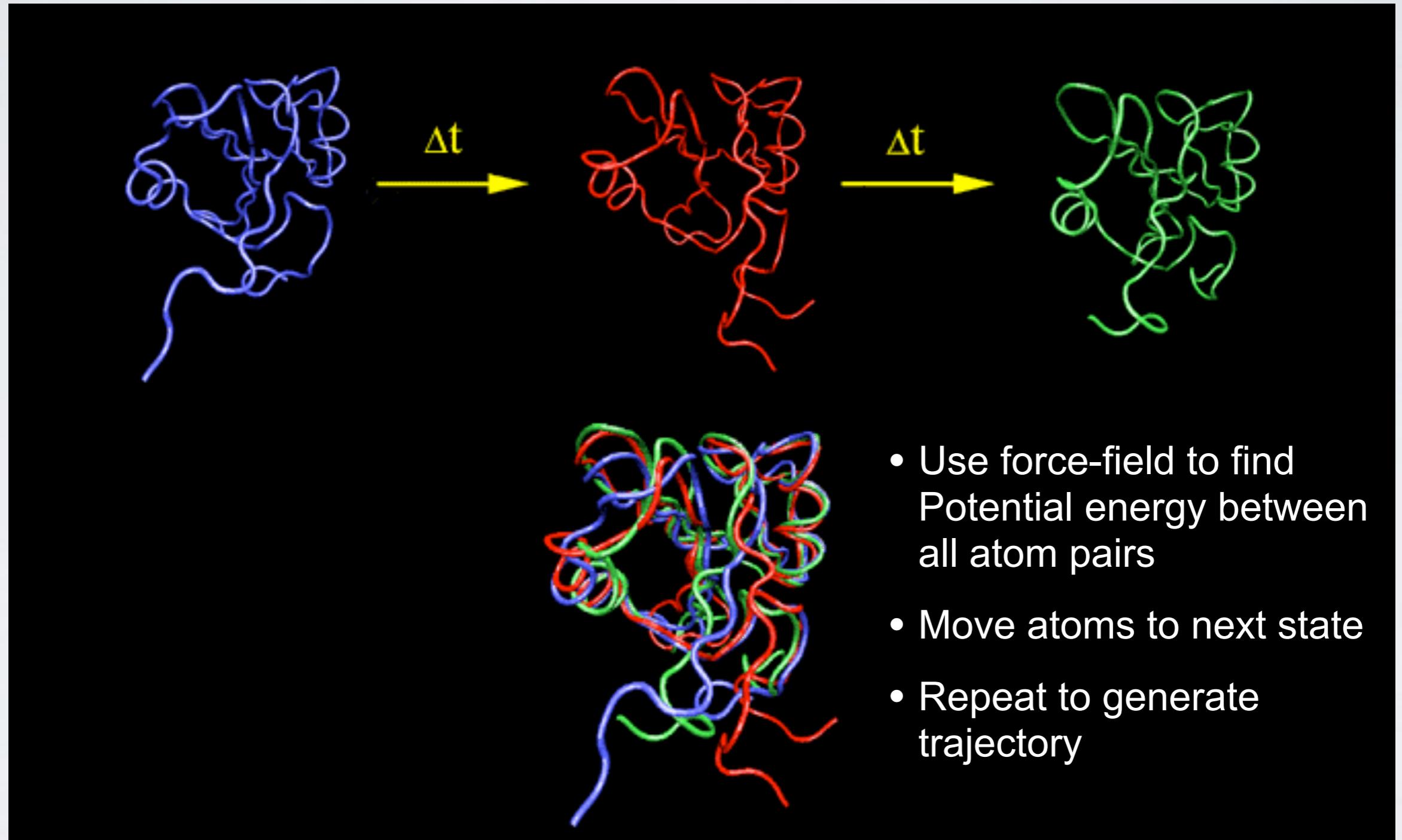
# NEXT UP:

- ▶ **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**
  - Predicting **functional dynamics & drug discovery**

# 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



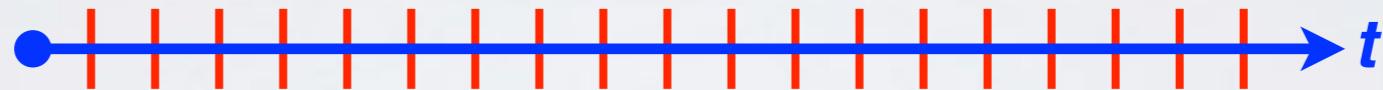
McCammon, Gelin & Karplus, *Nature* (1977)

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

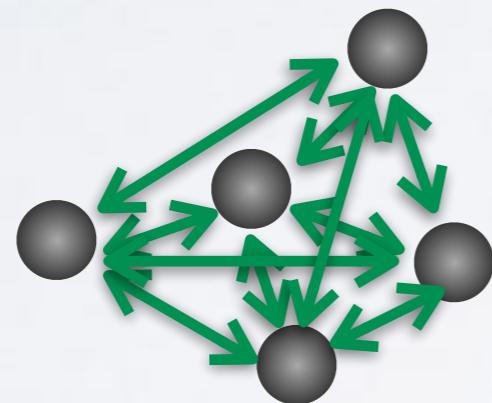
- ▶ Divide **time** into discrete ( $\sim 1\text{fs}$ ) **time steps ( $\Delta t$ )**  
(for integrating equations of motion, see below)



- ▶ Divide **time** into discrete ( $\sim 1\text{fs}$ ) **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*

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

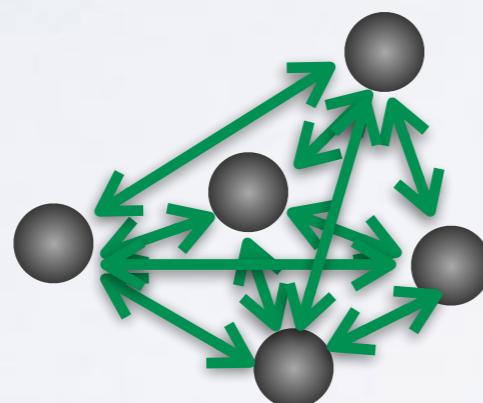
*Empirical force field*

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

- ▶ Divide **time** into discrete ( $\sim 1\text{fs}$ ) **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)



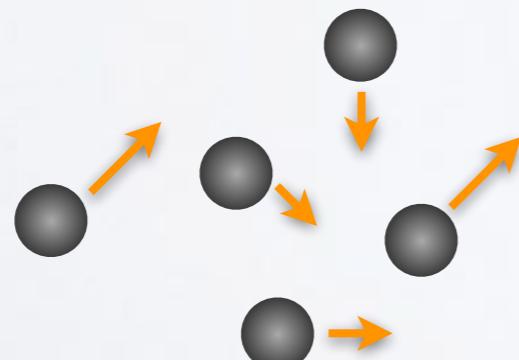
*Nucleic motion described classically*

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

*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)



$$\boxed{v(t + \frac{\Delta t}{2})} = v(t - \frac{\Delta t}{2}) + \frac{\mathbf{F}(t)}{m} \Delta t$$

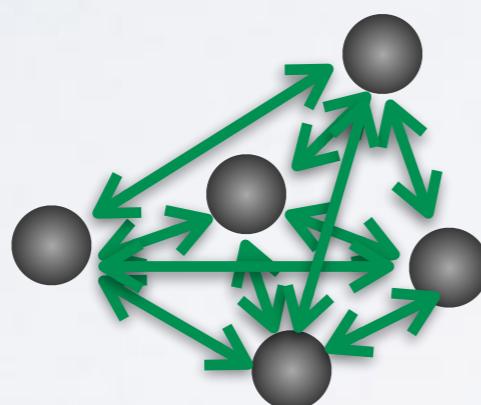
$$\mathbf{r}(t + \Delta t) = \mathbf{r}(t) + \boxed{v(t + \frac{\Delta t}{2})} \Delta t$$

# BASIC ANATOMY OF A MD SIMULATION

- Divide **time** into discrete ( $\sim 1\text{fs}$ ) **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)



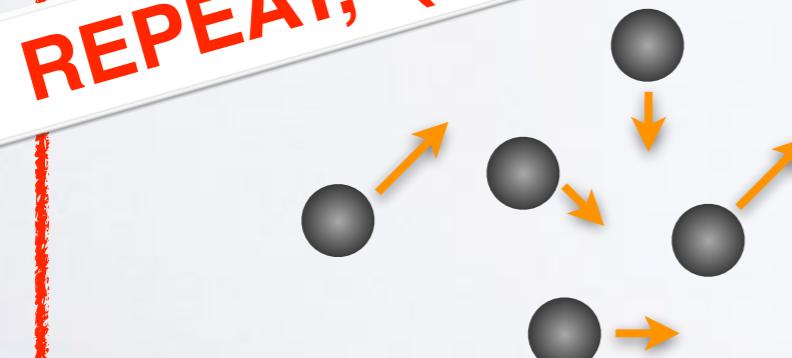
*Nucleic motion described classically*

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

*Empirical force function*

$$E(\vec{R}) = \sum_{i=1}^N \sum_{j \neq i, \text{non-bonded}} E_i(\vec{R})$$

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

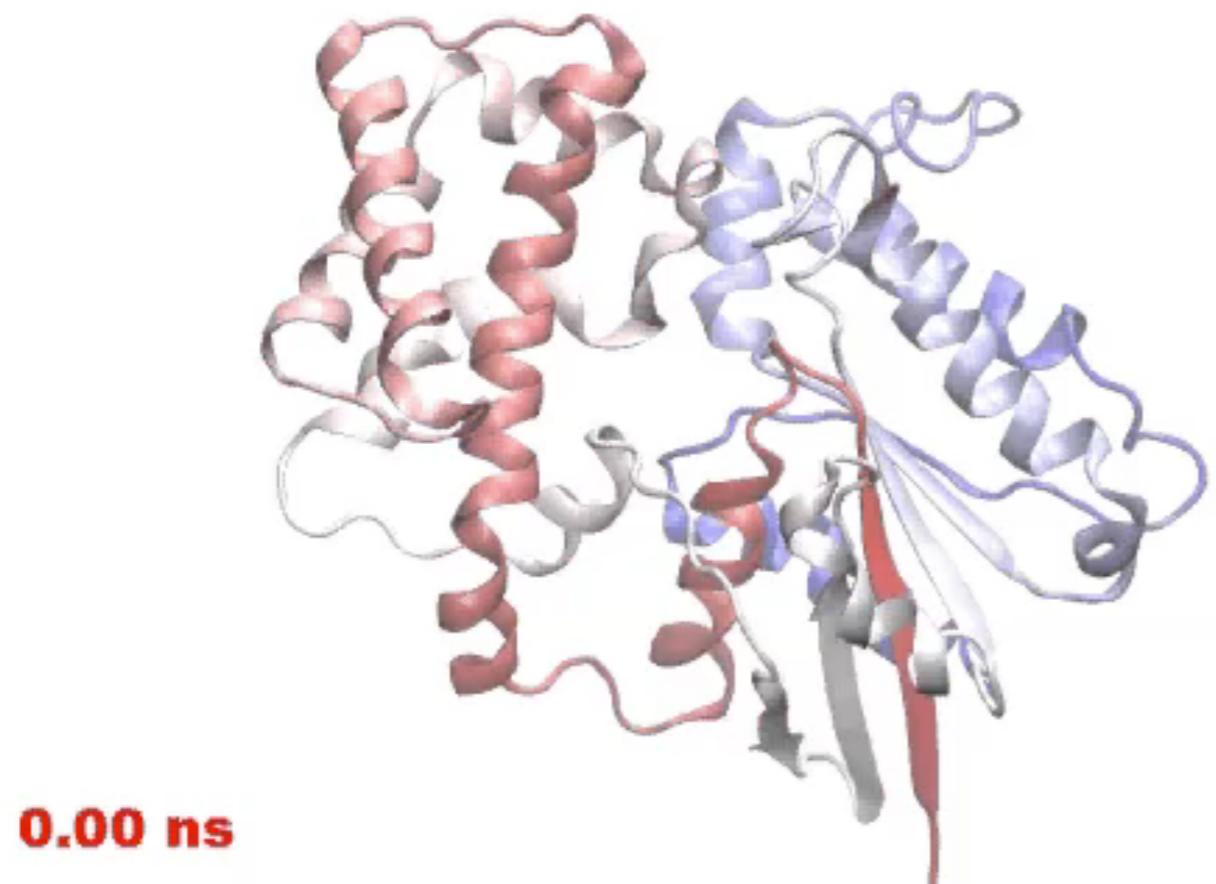


$$\begin{aligned} \mathbf{v}(t + \frac{\Delta t}{2}) &= \mathbf{v}(t - \frac{\Delta t}{2}) + \frac{\mathbf{F}(t)}{m} \Delta t \\ \mathbf{r}(t + \Delta t) &= \mathbf{r}(t) + \mathbf{v}(t + \frac{\Delta t}{2}) \Delta t \end{aligned}$$

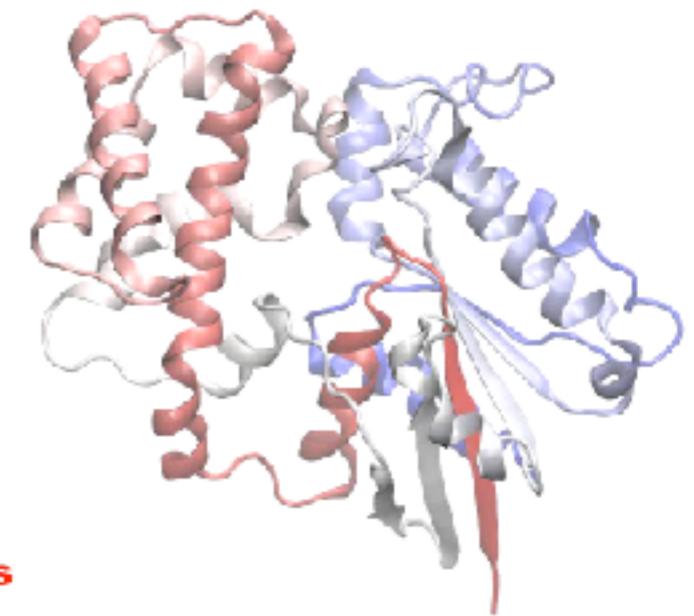
**REPEAT, (iterate many, many times... 1ms =  $10^{12}$  time steps)**

# MD Prediction of Functional Motions

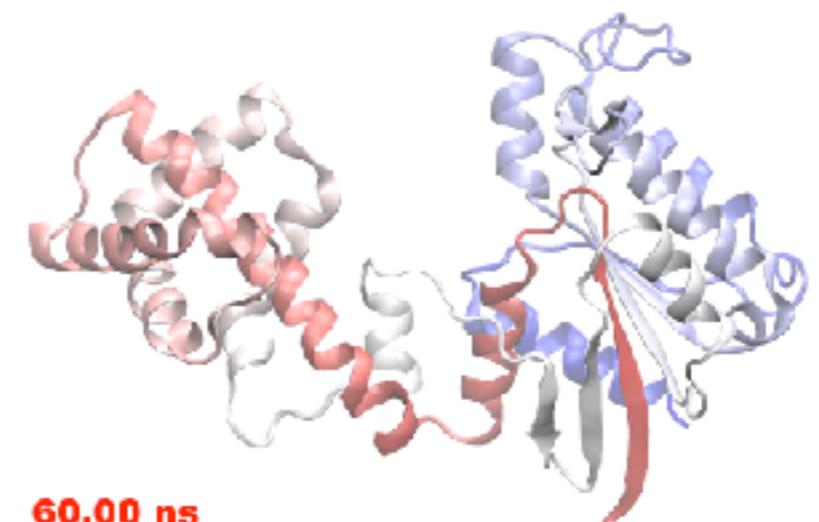
Accelerated MD simulation of  
nucleotide-free transducin alpha subunit



“close”

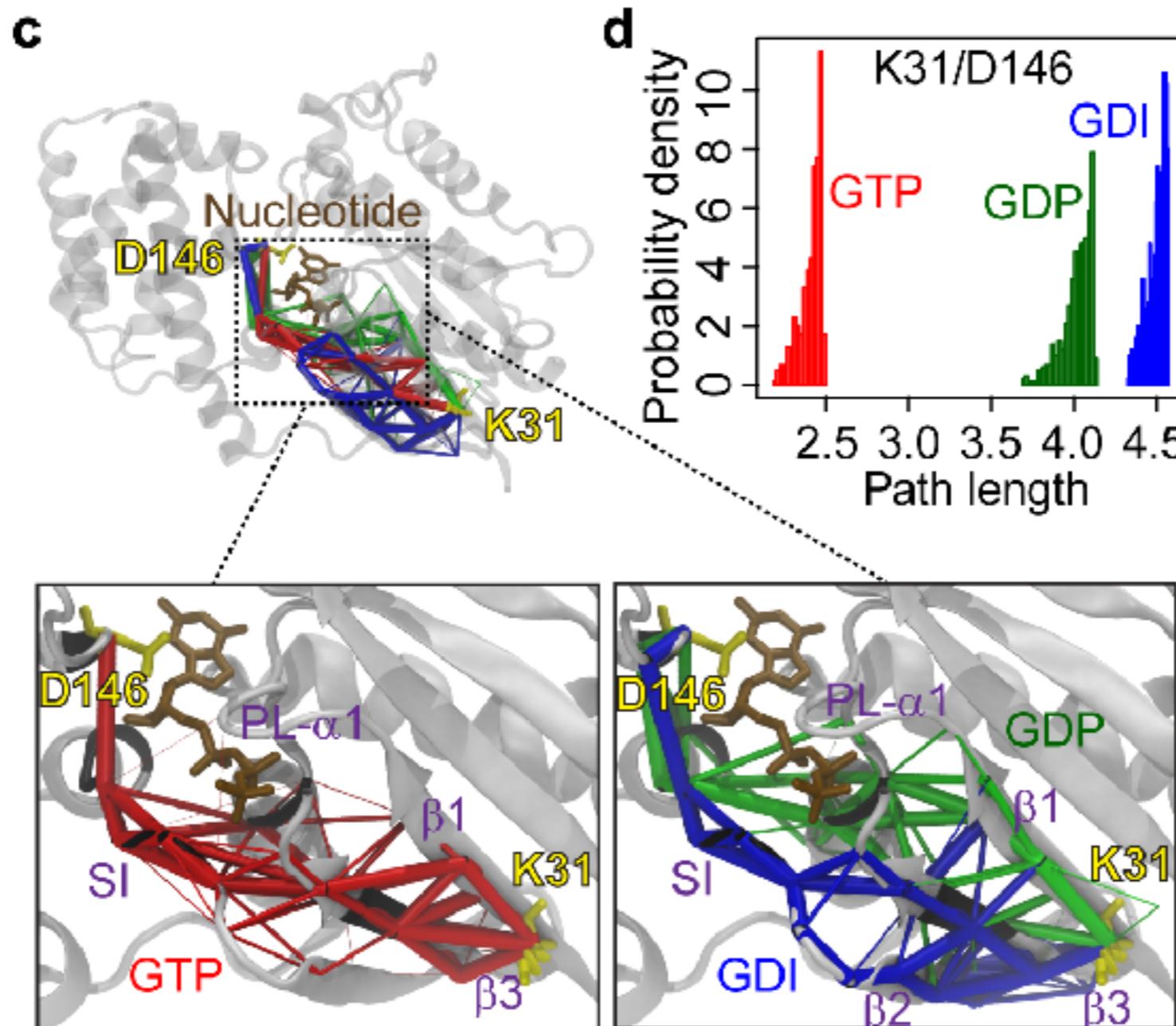


“open”

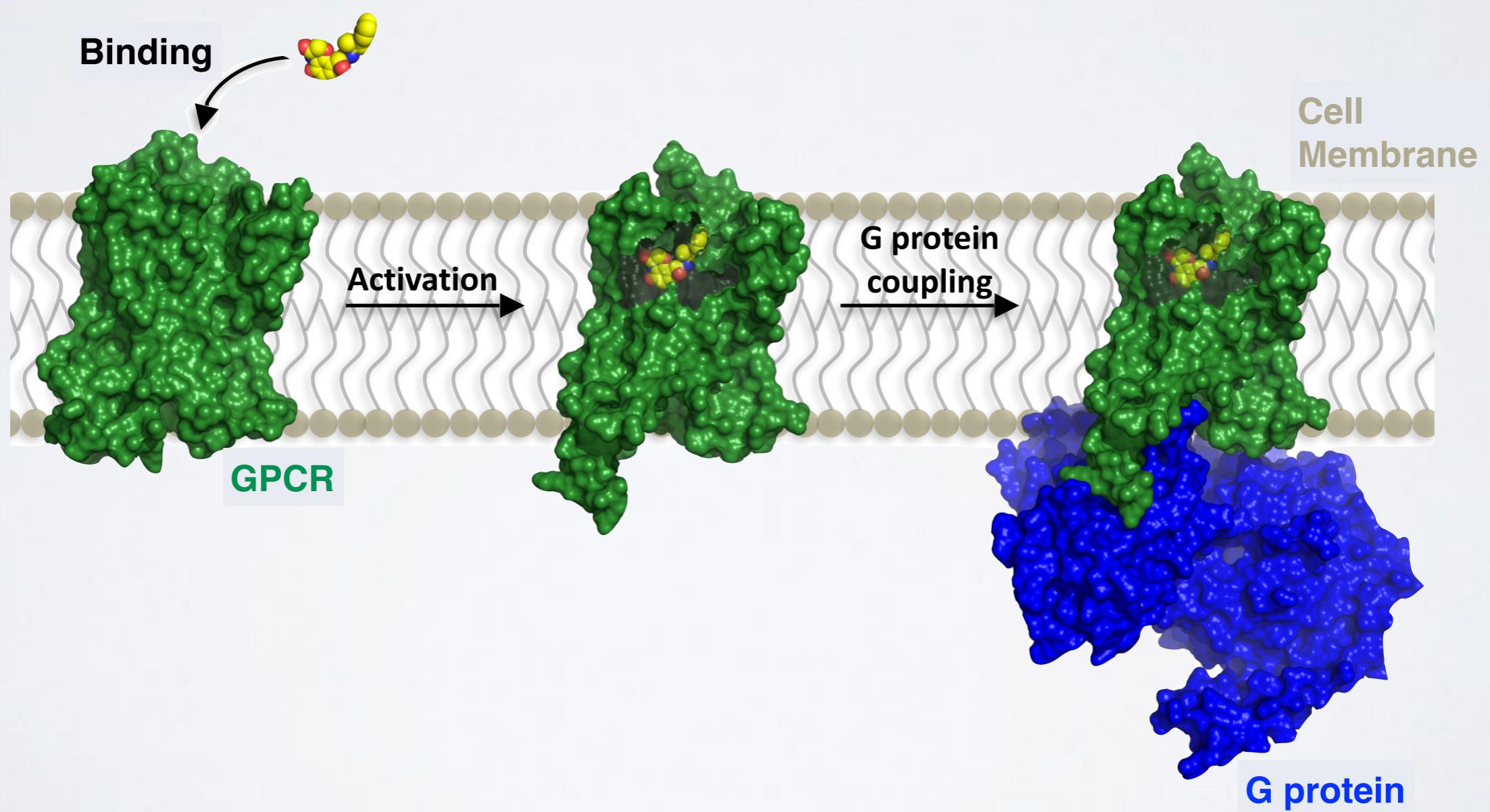


Yao and Grant, Biophys J. (2013)

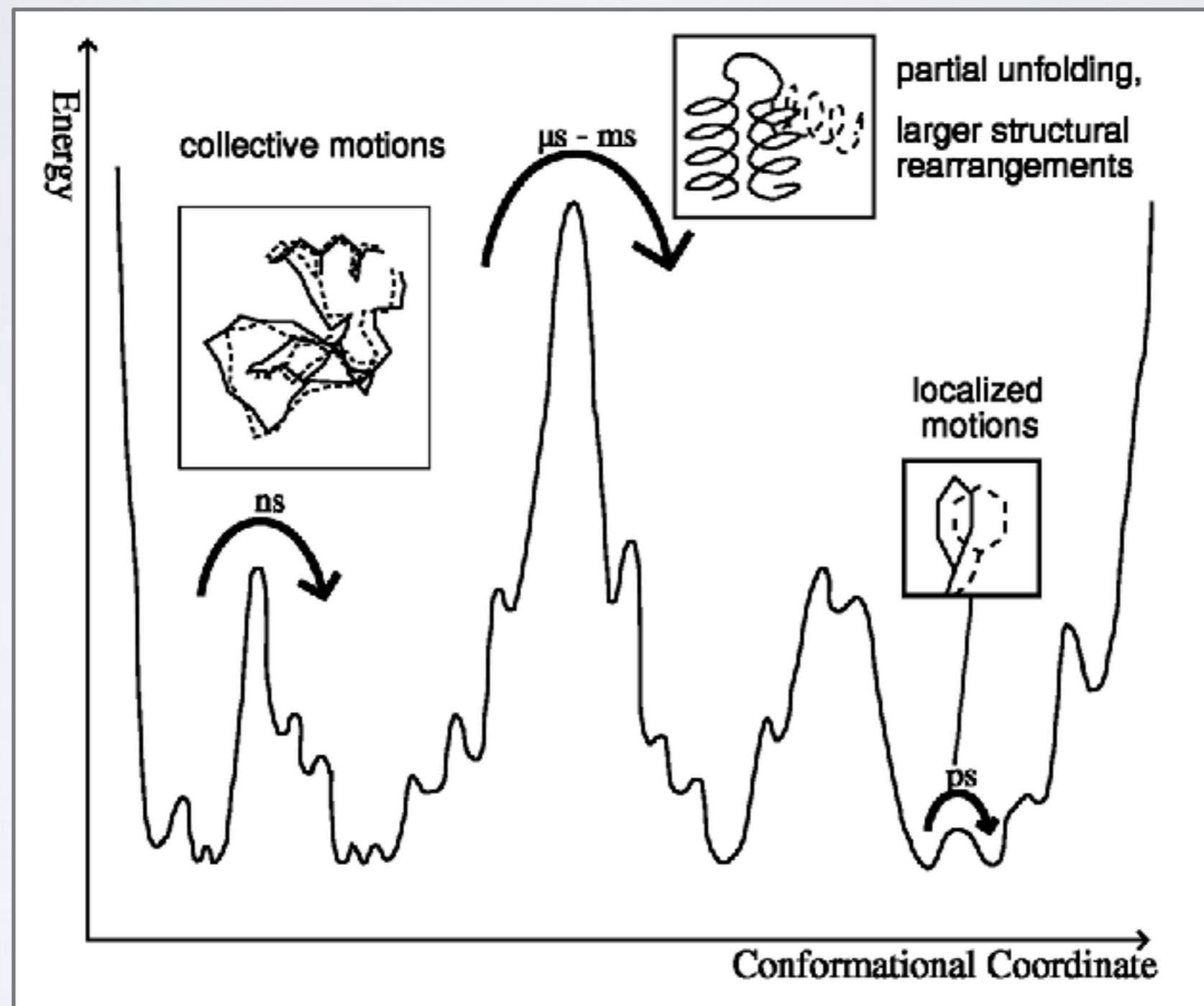
# Simulations Identify Key Residues Mediating Dynamic Activation



# 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

**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]

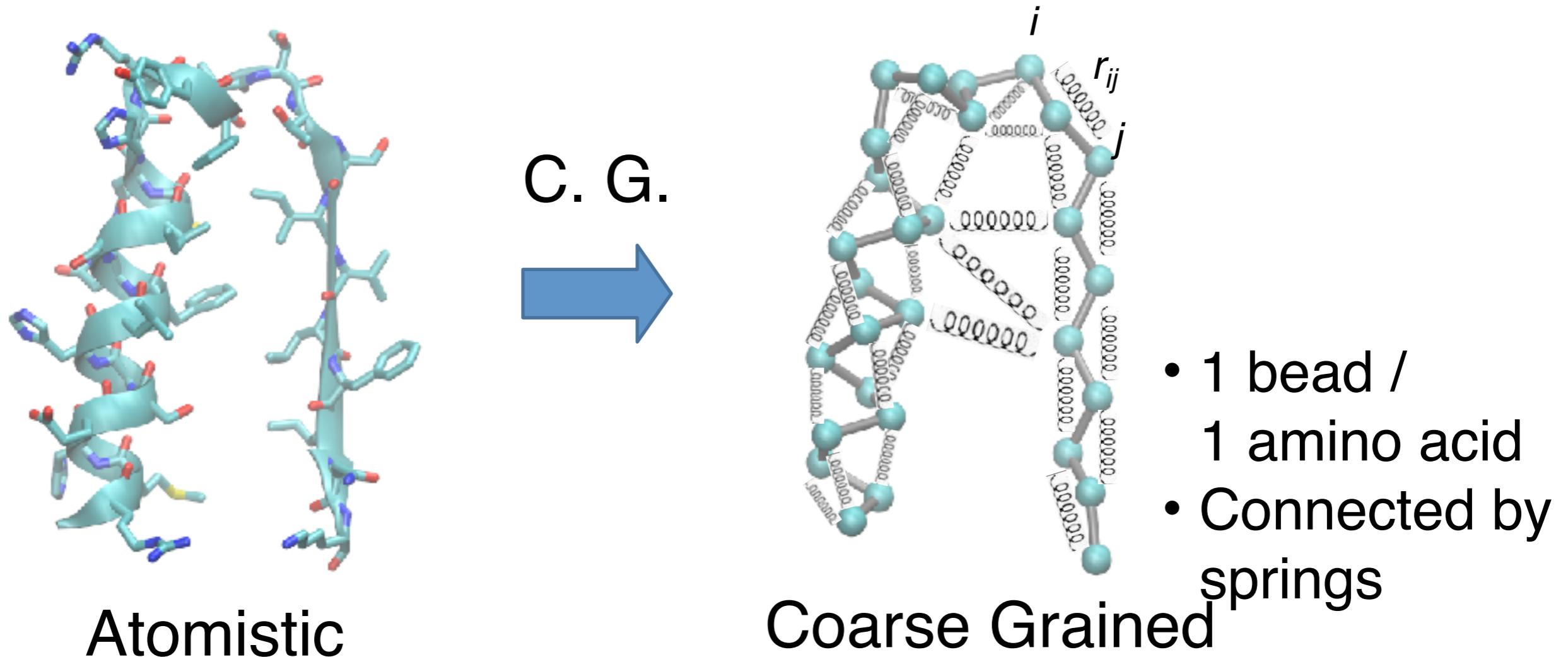
Total: 8.4 \* 10<sup>17</sup> flop  
(on a 100 Gflop/s cpu: **ca 25 years!**)

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

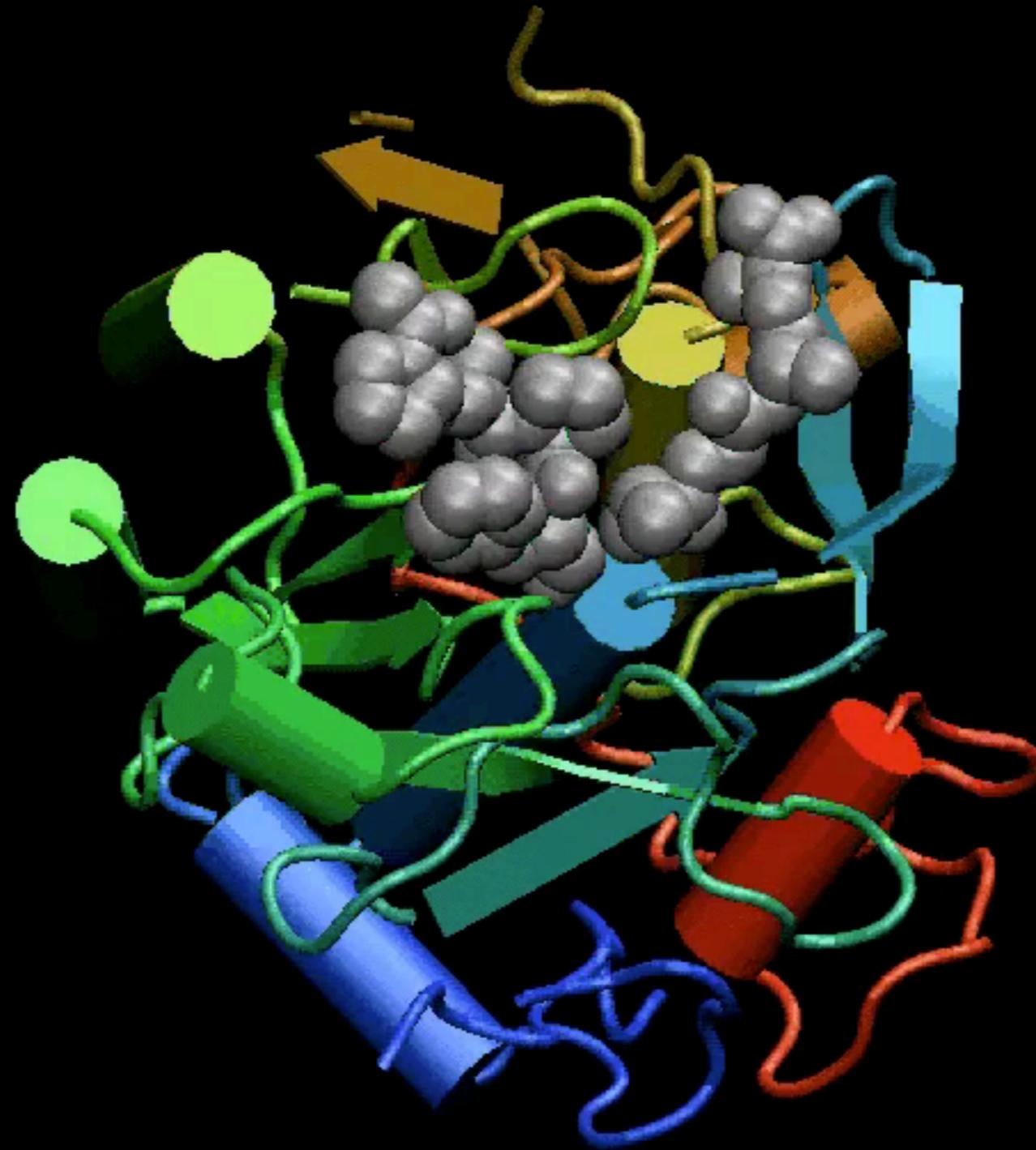
multiple time stepping	ca. 2.5 years
fast multipole methods	ca. 1 year
parallel computers	ca. 5 days
modern GPUs	ca. 1 day
<b>(Anton supercomputer</b>	<b>ca. minutes)</b>

# 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

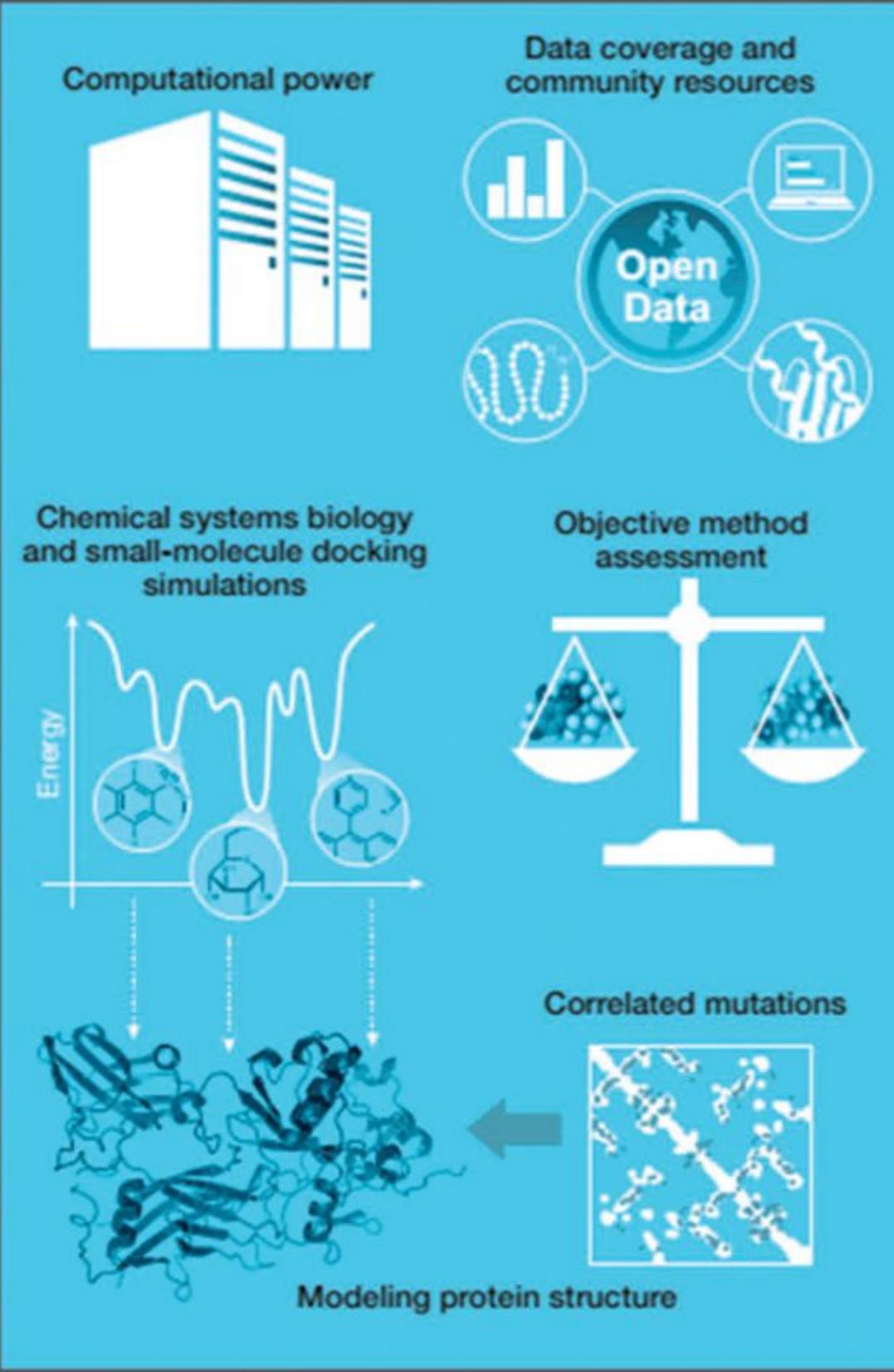
Do it Yourself!

# Hand-on time!

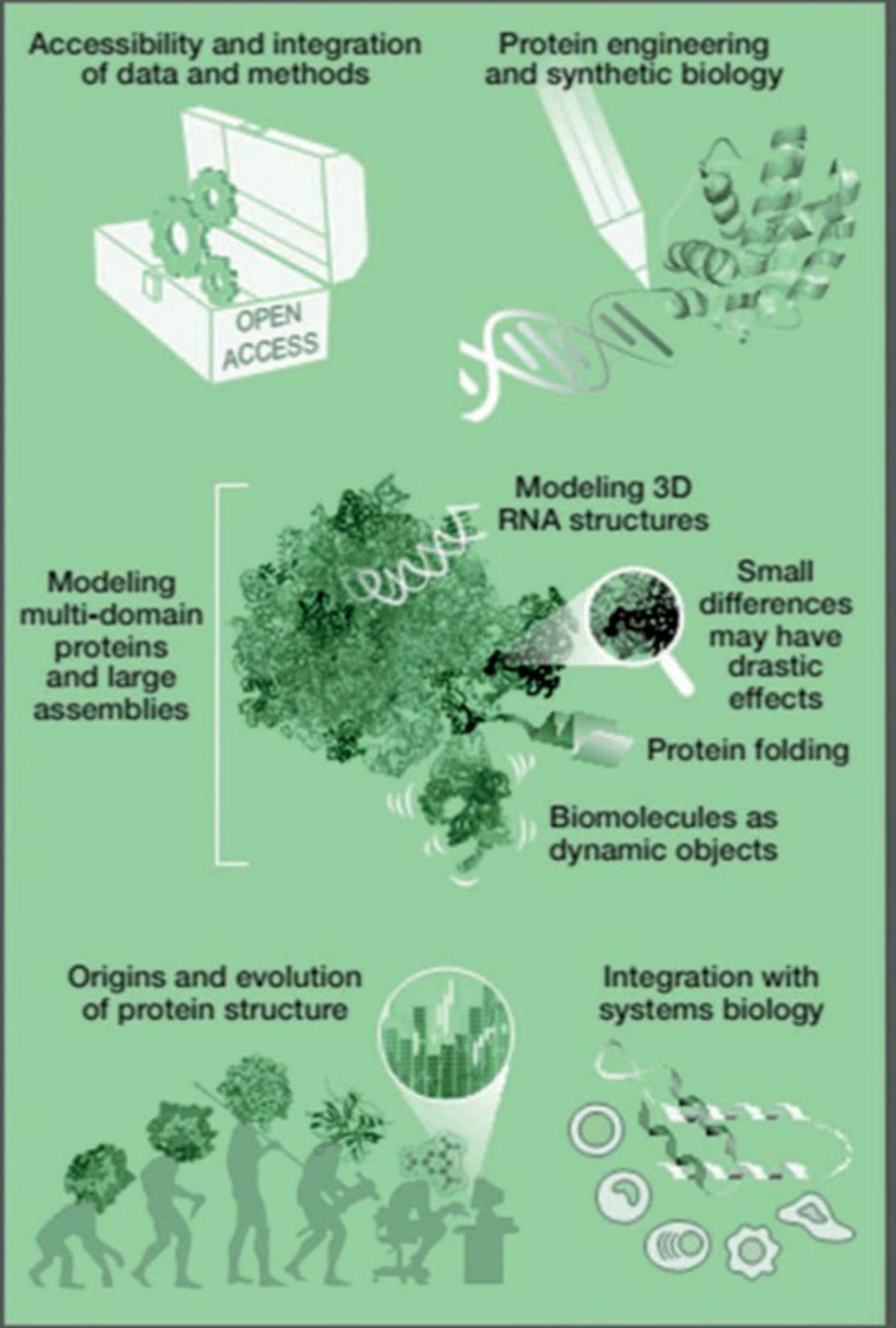
[https://bioboot.github.io/bggn213\\_f17/lectures/#12](https://bioboot.github.io/bggn213_f17/lectures/#12)

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

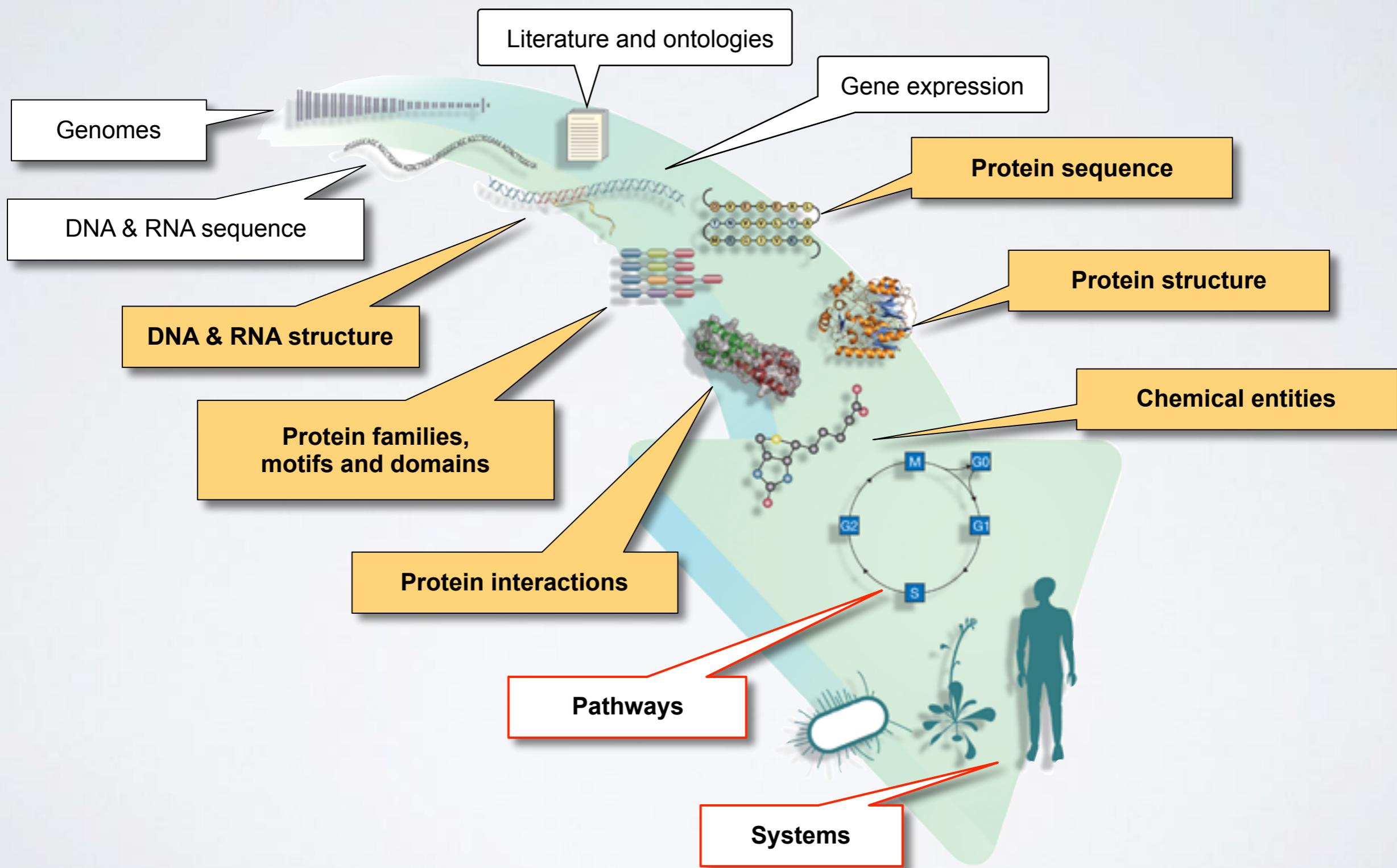
## ACHIEVEMENTS



## CHALLENGES



# 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