# BIMM 143

## Structural Bioinformatics II

Lecture 12

**Barry Grant**

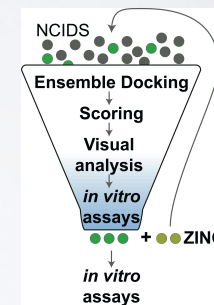UC San Diego

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

**Compound library**
(commercial, in-house, synthetic, natural)

**High throughput screening**
(HTS)

**Hit confirmation**

**Lead compounds**
(e.g., $\mu M\ K_d$)

**Lead optimization**
(Medicinal chemistry)

**Potent drug candidates**
($nM\ K_d$)

**Animal and clinical evaluation**

---

## COMPUTER-AIDED LIGAND DESIGN

Aims to reduce number of compounds synthesized and assayed

Lower costs
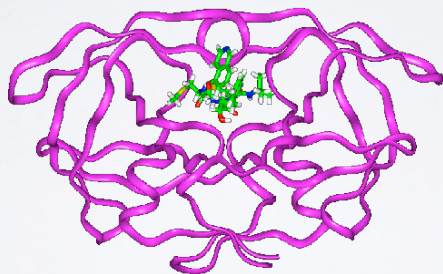
Reduce chemical waste

Facilitate faster progress

NCIDS

**Ensemble Docking**

**Scoring**

**Visual analysis**

*in vitro* **assays**

+ ZINC

*in vitro* **assays**

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
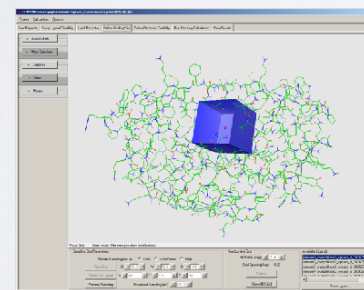


HIV Protease/KNI-272 complex

---

## PROTEIN-LIGAND DOCKING
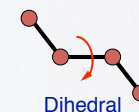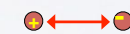
### Structure-Based Ligand Design

Docking software
Search for structure of lowest energy
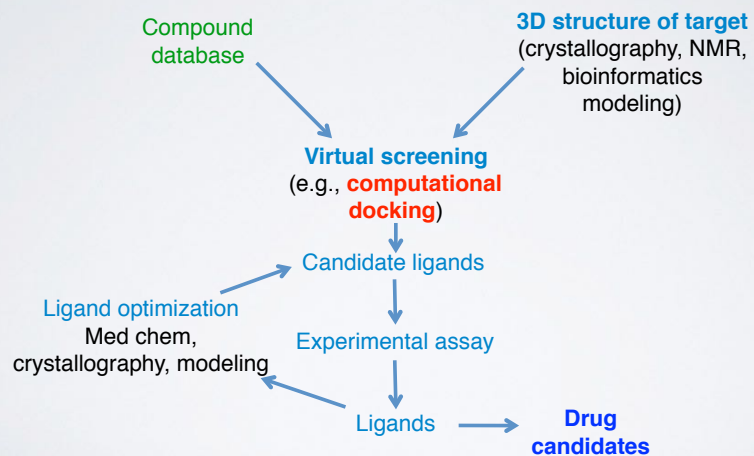
Potential function
Energy as function of structure



VDW

Screened Coulombic

Dihedral
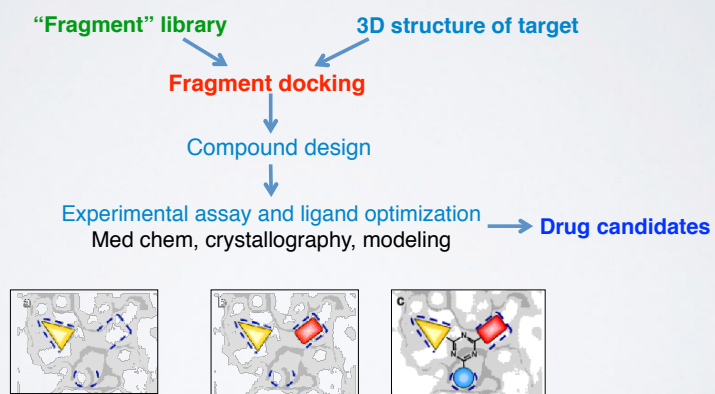
## STRUCTURE-BASED VIRTUAL SCREENING

Compound database

3D structure of target
(crystallography, NMR, bioinformatics modeling)

**Virtual screening**
(e.g., **computational docking**)

Candidate ligands

Ligand optimization
Med chem, crystallography, modeling

Experimental assay

Ligands → **Drug candidates**

## COMPOUND LIBRARIES



Commercial (in-house pharma)

Government (NIH)

Academia

## FRAGMENTAL STRUCTURE-BASED SCREENING

**"Fragment" library**

**3D structure of target**

**Fragment docking**

Compound design

Experimental assay and ligand optimization
Med chem, crystallography, modeling → **Drug candidates**



http://www.beilstein-institut.de/bozen2002/proceedings/Jhoti/jhoti.html
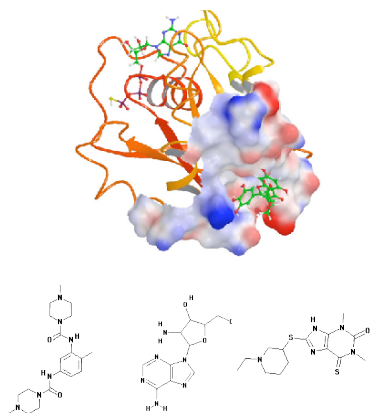
## Multiple non active-site pockets identified

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



Probe Occupancy
GTP
GDP
Residue No.

ethanol
isopropanol
acetone
cyclohexane
methylamine
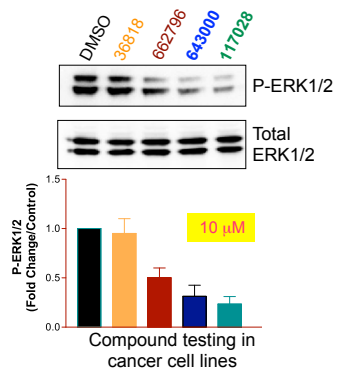benzene
phenol
acetamide

## 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.
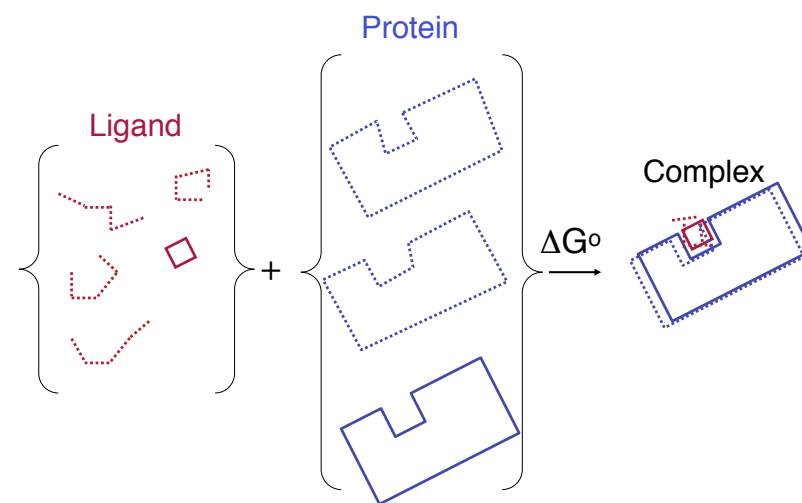
Ensemble computational docking

Compound effect on U251 cell line

DMSO  36818  662796  643000  117028

P-ERK1/2

Total ERK1/2

P-ERK1/2 (Fold Change/Control)

10 µM

Compound testing in cancer cell lines

**PLoS One (2011, 2012)**

---

## Proteins and Ligand are Flexible

Protein

Ligand

$\Delta G^o$

Complex

+

---

## 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**

## Slide 1

# Hand-on time!

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**

## Slide 2

### Proteins and Ligand are Flexible



Ligand + Protein $\xrightarrow{\Delta G^o}$ Complex

## Slide 3

HTTP://129.177.232.111:3848/PCA-APP/

HTTPS://DCMB-GRANT-SHINY.UMMS.MED.UMICH.EDU/PCA-APP/

HTTP://BIO3D.UCSD.EDU/PCA-APP/

## Slide 4

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

Compound Library          **Known Ligands**

**Molecular similarity**
Machine-learning
Etc.

Candidate ligands

Optimization
Med chem, crystallography, modeling

Assay

Actives → **Potent drug candidates**

---

## CHEMICAL SIMILARITY
## LIGAND-BASED DRUG-DISCOVERY

Compounds
(available/synthesizable)

Compare with known ligands

Different → Don't bother

Similar → Test experimentally

## CHEMICAL FINGERPRINTS
### BINARY STRUCTURE KEYS

Molecule 1 — phenyl, naphthyl, ketone, methyl, ethyl, aldehyde, alcohol, amide, carboxylate ... S-S bond, chlorine, fluorine

Molecule 2

$N_I = 2$

---

## CHEMICAL SIMILARITY FROM FINGERPRINTS

Molecule 1 — phenyl, naphthyl, ketone, methyl, ethyl, aldehyde, alcohol, amide, carboxylate ... S-S bond, chlorine, fluorine

Molecule 2

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

Bulky hydrophobe

5.0 ±0.3 Å

3.2 ±0.4 Å

+ 1

2.8 ±0.3 Å

Aromatic

---

## 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)

• 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



Descriptor 3
Descriptor 1
Descriptor 2

Point representing a compound in descriptor space

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



EMBL-EBI

---

## Drug properties



| Drug Type | Rule of Five | First in Class | Chirality | Prodrug | Oral | Parenteral | Topical | Black-Box Warning | Availability Type |

synthetic small molecule
natural product-derived
inorganic
polymer
monoclonal antibody
enzyme
peptide/ protein
oligonucleotide
oligosaccharide

racemic mixture
chirally pure

prescription only
over-the-counter
discontinued

Ingredient-related
(USANs, candidates and approved drugs)

Product-related
(approved drugs only)

EMBL-EBI

---

## 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



View cavities (and ligands) on structure

Details of sites identified

# Examples

## Quantifying the chemical beauty of drugs

G. Richard Bickerton, Gaia V. Paolini, Jérémy Besnard, Sorel Muresan & Andrew L. Hopkins

Affiliations | Contributions | Corresponding author

Citation | Reprints | Rights & permissions | Article metrics

### Abstract

Abstract · References · Author information · Supplementary information

Drug-likeness is a key consideration when selecting compounds during the early stages of drug discovery. However, evaluation of drug-likeness in absolute terms does not reflect adequately the whole spectrum of compound quality. More worryingly, widely used rules may inadvertently foster undesirable molecular property inflation as they permit the encroachment of rule-compliant compounds towards their boundaries. We propose a measure of drug-likeness based on the concept of desirability called the quantitative estimate of drug-likeness (QED). The empirical rationale of QED reflects the underlying distribution of molecular properties. QED is intuitive, transparent, straightforward to implement in many practical settings and allows compounds to be ranked by their relative merit. We extended the utility of QED by applying it to the problem of molecular target druggability assessment by prioritizing a large set of published bioactive compounds. The measure may also capture the abstract notion of aesthetics in medicinal chemistry.

Subject terms: Pharmacology · Theoretical chemistry

At a glance

# 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

# Examples

Abstract

The lack of success in target-based screening approaches to the discovery of antibacterial agents has led to resurgence 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.

---

# ARTICLE

doi:10.1038/nature11159

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

Eugen Lounkine, Michael J. Keiser, Steven Whitebread, Dmitri Mikhailov, Jacques Hamon, Jeremy L. Jenkins, Paul Lavan, Eckhard Weber, Allison K. Doak, Serge Côté, Brian K. Shoichet & Laszlo Urban

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 estrogen 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.

---

# 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 <u>intrinsically flexible</u> 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 <u>mapping of structure to function</u>**
  – **Molecular dynamics** (MD) and **normal mode analysis** (NMA) are two major methods for predicting and characterizing molecular motions and their properties

---

# MOLECULAR DYNAMICS SIMULATION



$\Delta t$

• Use force-field to find Potential energy between all atom pairs

• Move atoms to next state

• Repeat to generate trajectory

McCammon, Gelin & Karplus, *Nature* (1977)
[ See: https://www.youtube.com/watch?v=ui1ZysMFcKk ]

**Slide 1 (top-left):**

▷ Divide **time** into discrete (~1fs) **time steps** (**Δt**)
(for integrating equations of motion, see below)

$\longrightarrow t$

**Slide 2 (top-right):**

▷ Divide **time** into discrete (~1fs) **time steps** (**Δt**)
(for integrating equations of motion, see below)

$\longrightarrow t$

▷ 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*
$$E(\vec{R}) = \sum_{\text{bonded}} E_i(\vec{R}) + \sum_{\text{non-bonded}} E_i(\vec{R})$$
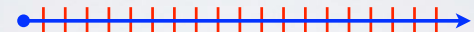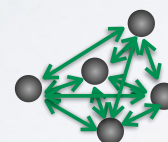
**Slide 3 (bottom-left):**

▷ Divide **time** into discrete (~1fs) **time steps** (**Δt**)
(for integrating equations of motion, see below)

$\longrightarrow t$

▷ 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)

$$\boldsymbol{v}\left(t + \frac{\Delta t}{2}\right) = \boldsymbol{v}\left(t - \frac{\Delta t}{2}\right) + \frac{\boldsymbol{F}(t)}{m}\Delta t$$
$$\boldsymbol{r}(t + \Delta t) = \boldsymbol{r}(t) + \boldsymbol{v}\left(t + \frac{\Delta t}{2}\right)\Delta t$$

**Slide 4 (bottom-right):**

## BASIC ANATOMY OF A MD SIMULATION

▷ Divide **time** into discrete (~1fs) **time steps** (**Δt**)
(for integrating equations of motion, see below)

$\longrightarrow t$

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

*Nucleic motion described classically*
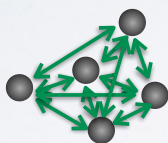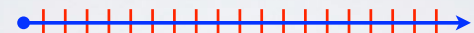$$m_i \frac{d^2}{dt^2} \vec{R}_i = -\vec{\nabla}_i E(\vec{R})$$

*Empirical force field*
$$E(\vec{R}) = \sum \ldots E_i(\vec{R})$$

▷ Use the forces to calculate **velocities** and move atoms to new **positions**
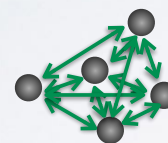(by integrating numerically via the "leapfrog" scheme)

$$\boldsymbol{v}\left(t + \frac{\Delta t}{2}\right) = \boldsymbol{v}\left(t - \frac{\Delta t}{2}\right) + \frac{\boldsymbol{F}(t)}{m}\Delta t$$
$$\boldsymbol{r}(t + \Delta t) = \boldsymbol{r}(t) + \boldsymbol{v}\left(t + \frac{\Delta t}{2}\right)\Delta t$$

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

# MD Prediction of Functional Motions



Accelerated MD simulation of nucleotide-free transducin alpha subunit

0.00 ns

"close"

0.00 ns

"open"

60.00 ns

Yao and Grant, Biophys J. (2013)

# Simulations Identify Key Residues Mediating Dynamic Activation



Yao … Grant, Journal of Biological Chemistry (2016)

# EXAMPLE APPLICATION OF MOLECULAR SIMULATIONS TO GPCRS



Binding

Cell Membrane

Activation

G protein coupling

GPCR

G protein

# PROTEINS JUMP BETWEEN MANY, HIERARCHICALLY ORDERED "CONFORMATIONAL SUBSTATES"



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

## MOLECULAR DYNAMICS IS VERY EXPENSIVE

**Example**: $F_1$-ATPase in water (183,674 atoms) for 1 nanosecond:

        => $10^6$ integration steps
        => $8.4 * 10^{11}$ floating point operations/step
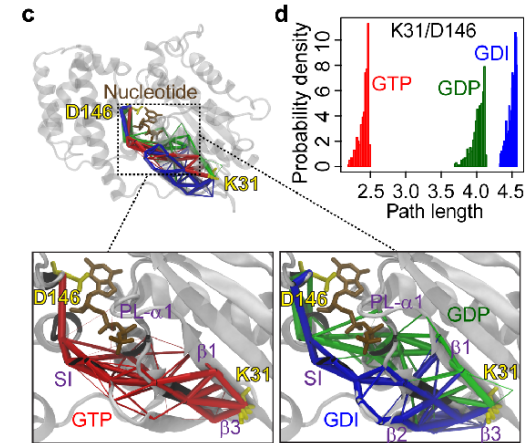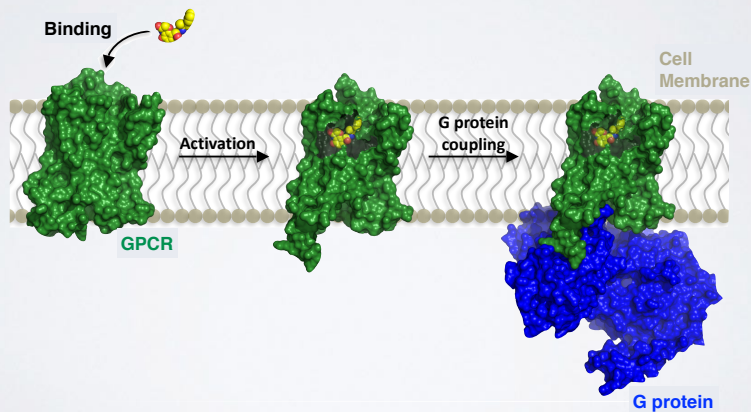          [n(n-1)/2 interactions]

          Total:    $8.4 * 10^{17}$ flop
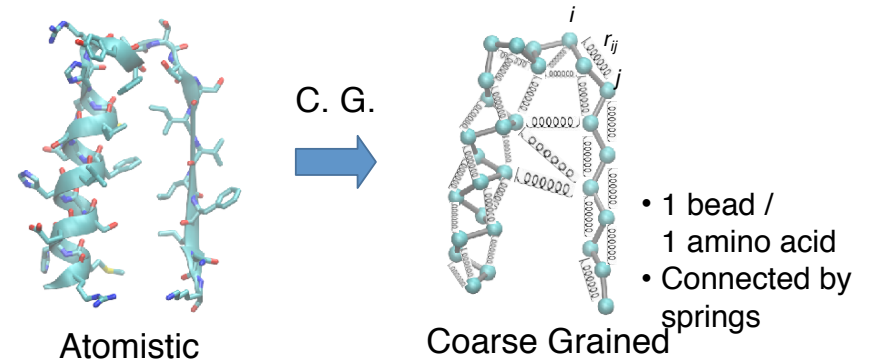          (on a 100 Gflop/s cpu:    **ca 25 years!**)

**… but performance has been improved by use of:**
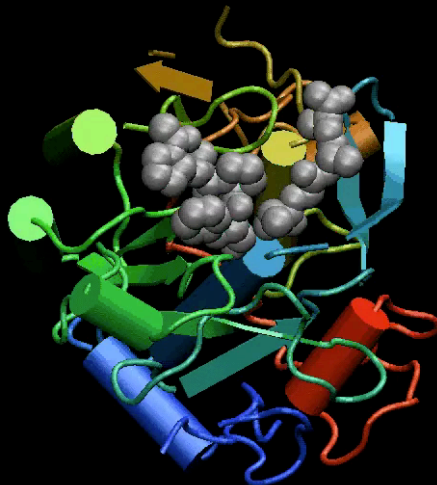   multiple time stepping         ca.  2.5 years
   fast multipole methods   ca.  1 year
   parallel computers         ca.  5 days
  modern GPUs           **ca.  1 day**
  **(Anton supercomputer     ca.  minutes)**

---

## 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.



C. G.

Atomistic        Coarse Grained

- 1 bead / 1 amino acid
- Connected by springs

---



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

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

Ilan Samish et al. Bioinformatics 2015;31:146-150

## INFORMING SYSTEMS BIOLOGY?



Genomes · DNA & RNA sequence · Literature and ontologies · Gene expression · Protein sequence · DNA & RNA structure · Protein structure · Protein families, motifs and domains · Chemical entities · Protein interactions · Pathways · Systems

## 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