# BGGN 213

## Structural Bioinformatics II

### Lecture 12

**Barry Grant**

UC San Diego
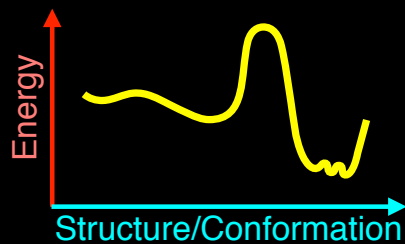
http://thegrantlab.org/bggn213

---

# Next Up:

- **Overview of structural bioinformatics**
  - Motivations, goals and challenges

- **Fundamentals of protein structure**
  - Structure composition, form and forces

- **Representing, interpreting & modeling protein structure**
  - Visualizing and interpreting protein structures
  - Analyzing protein structures
  - Modeling energy as a function of structure
  - Drug discovery & Predicting functional dynamics

---

# Key concept:

**Potential functions** describe a systems **energy** as a function of its **structure**



---

Two main approaches:
  (1). **Physics-Based**
  (2). **Knowledge-Based**

Two main approaches:
(1). **Physics-Based**
(2). **Knowledge-Based**

For physics based potentials
energy terms come from physical theory

$$V(R) = E_{\text{bonded}} + E_{\text{non.bonded}}$$

$$V(R) = E_{bonded} + E_{non.bonded}$$

Sum of bonded and non-bonded
atom-type and position based terms

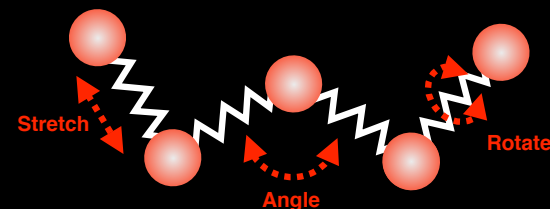$$V(R) = \boxed{E_{bonded}} + E_{non.bonded}$$

$E_{bonded}$ is itself a sum of three terms:

$$V(R) = \boxed{E_{bonded}} + E_{non.bonded}$$
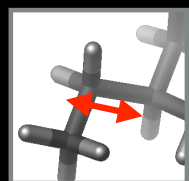
$E_{bonded}$ is itself a sum of three terms:

$$\boxed{E_{bond.stretch} + E_{bond.angle} + E_{bond.rotate}}$$

---

**Stretch**

**Angle**
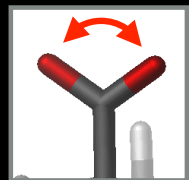
**Rotate**

$$V(R) = \boxed{E_{bonded}} + E_{non.bonded}$$

$E_{bonded}$ is itself a sum of three terms:

$$\boxed{E_{bond.stretch} + E_{bond.angle} + E_{bond.rotate}}$$

---

**Bond Stretch**
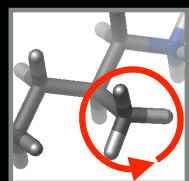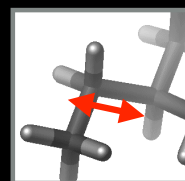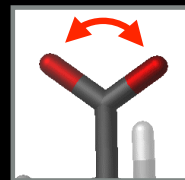
$$E_{bond.stretch}$$

**Bond Angle**

$$E_{bond.angle}$$

**Bond Rotate**

$$E_{bond.rotate}$$

---

**Bond Stretch**

$$\sum_{bonds} K_i^{bs}(b_i - b_o)$$

**Bond Angle**

$$\sum_{angles} K_i^{ba}(\theta_i - \theta_o)$$

**Bond Rotate**

$$\sum_{dihedrals} K_i^{br}[1 - cos(n_i\phi_i - \phi_o)]$$

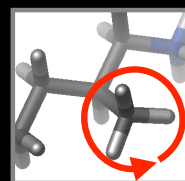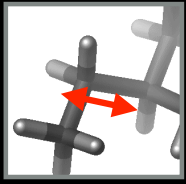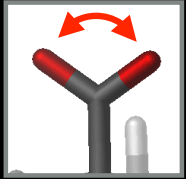**Bond Stretch**

$$\sum_{bonds} K_i^{bs}(b_i - b_o)$$

**Bond Angle**

$$\sum_{angles} K_i^{ba}(\theta_i - \theta_o)$$

**Bond Rotate**

$$\sum_{dihedrals} K_i^{br}[1 - cos(n_i\phi_i - \phi_o)]$$

$$V(R) = E_{bonded} + \boxed{E_{non.bonded}}$$

$E_{non.bonded}$ is a sum of two terms:

$$V(R) = E_{bonded} + \boxed{E_{non.bonded}}$$

$E_{non.bonded}$ is a sum of two terms:

$$E_{van.der.Waals} + E_{electrostatic}$$



**Non-bonded**

**Stretch**

**Rotate**

**Angle**

$$V(R) = E_{bonded} + \boxed{E_{non.bonded}}$$

$E_{non.bonded}$ is a sum of two terms:

$$E_{van.der.Waals} + E_{electrostatic}$$

Non-bonded

Stretch

Rotate

Angle

$$E_{electrostatic} = \sum_{pairs.i.j} \frac{q_i q_j}{\epsilon r_{ij}}$$

$$E_{van.der.Waals} = \sum_{pairs.i.j} \left[ \epsilon_{ij}\left(\frac{r_{o.ij}}{r_{ij}}\right)^{12} - 2\epsilon_{ij}\left(\frac{r_{o.ij}}{r_{ij}}\right)^{6} \right]$$

# Total potential energy
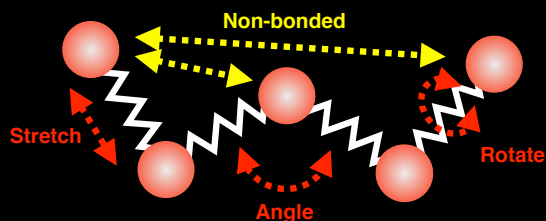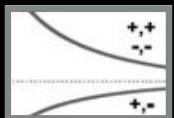
The potential energy can be given as a sum of terms for: Bond stretching, Bond angles, Bond rotations, van der Walls and Electrostatic interactions between atom pairs

$$V(R) = E_{bond.stretch}$$
$$\left. \begin{array}{l} +E_{bond.angle} \\ +E_{bond.rotate} \end{array} \right\} E_{bonded}$$
$$\left. \begin{array}{l} +E_{van.der.Waals} \\ +E_{electrostatic} \end{array} \right\} E_{non.bonded}$$

# Potential energy surface

Now we can calculate the potential energy surface that fully describes the energy of a molecular system as a function of its geometry



Energy (V)

Position (x)

# Potential energy surface

Now we can calculate the potential energy surface that fully describes the energy of a molecular system as a function of its geometry



Energy (V)

Position (x)

# Key concept:

Now we can calculate the potential energy surface that fully describes the energy of a molecular system as a function of its geometry



- The **forces** are the gradients of the energy

$$F(x) = -dV/dx$$

---

# Moving Over The Energy Surface

- **Energy Minimization** drops into local minimum

- **Molecular Dynamics** uses thermal energy to move smoothly over surface

- **Monte Carlo Moves** are random. Accept with probability:
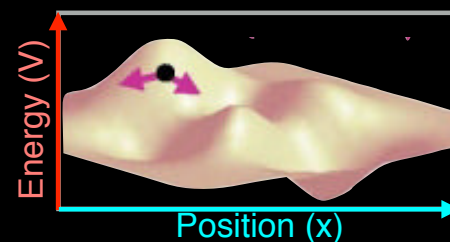
$$exp(-\Delta V/dx)$$



---

## PHYSICS-ORIENTED APPROACHES

**Weaknesses**
Fully physical detail becomes computationally intractable
Approximations are unavoidable
    (Quantum effects approximated classically, water may be treated crudely)
Parameterization still required

**Strengths**
Interpretable, provides guides to design
Broadly applicable, in principle at least
Clear pathways to improving accuracy

**Status**
Useful, widely adopted but far from perfect
Multiple groups working on fewer, better approxs
    Force fields, quantum
    entropy, water effects
Moore's law: hardware improving

---

## SIDE-NOTE: GPUS AND ANTON SUPERCOMPUTER



## SIDE-NOTE: GPUS AND ANTON SUPERCOMPUTER



## POTENTIAL FUNCTIONS DESCRIBE A SYSTEMS **ENERGY** AS A FUNCTION OF ITS **STRUCTURE**

Two main approaches:
(1). **Physics-Based**
(2). **Knowledge-Based**

## KNOWLEDGE-BASED DOCKING POTENTIALS



Histidine

Ligand carboxylate

2.69
3.01

Aromatic stacking

3.44
3.47
3.42

## ENERGY DETERMINES **PROBABILITY** (STABILITY)

Basic idea: Use probability as a proxy for energy



Boltzmann:

$$p(r) \propto e^{-E(r)/RT}$$

Inverse Boltzmann:

$$E(r) = -RT \ln\left[p(r)\right]$$

Example: ligand carboxylate O to protein histidine N

Find all protein-ligand structures in the PDB with a ligand carboxylate O
1. For each structure, histogram the distances from O to every histidine N
2. Sum the histograms over all structures to obtain $p(r_{O-N})$
3. Compute $E(r_{O-N})$ from $p(r_{O-N})$

---

## KNOWLEDGE-BASED POTENTIALS

**Weaknesses**
Accuracy limited by availability of data

**Strengths**
Relatively easy to implement
Computationally fast

**Status**
Useful, far from perfect
May be at point of diminishing returns
(not always clear how to make improvements)

---

# Computer Aided Drug Discovery

---

# Next Up:

- **Overview of structural bioinformatics**
  - Motivations, goals and challenges

- **Fundamentals of protein structure**
  - Structure composition, form and forces

- **Representing, interpreting & modeling protein structure**
  - Visualizing and interpreting protein structures
  - Analyzing protein structures
  - Modeling energy as a function of structure
  - Drug discovery & Predicting functional dynamics

**Slide 1: 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**

**Slide 2: 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

**Slide 3:**

Two main approaches:
(1). **Receptor/Target-Based**
(2). **Ligand/Drug-Based**

**Slide 4:**

Two main approaches:
(1). **Receptor/Target-Based**
(2). **Ligand/Drug-Based**

# SCENARIO I:
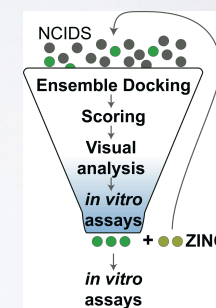## RECEPTOR-BASED DRUG DISCOVERY

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

## COMMON SIMPLIFICATIONS USED IN PHYSICS-BASED DOCKING

Quantum effects approximated classically

Protein often held rigid

Configurational entropy neglected

Influence of water treated crudely

---

# Hand-on time!

https://bioboot.github.io/bggn213_W19/lectures/#12

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

---

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



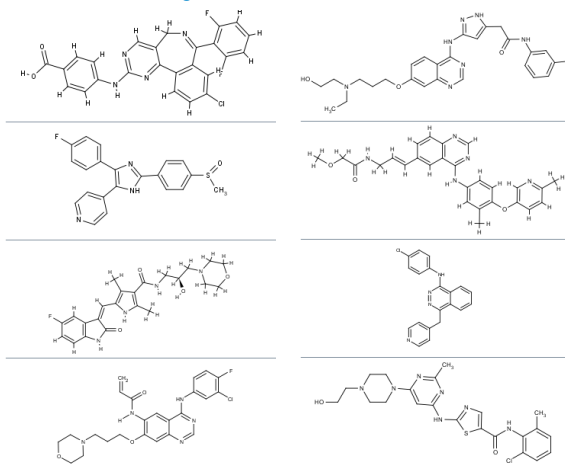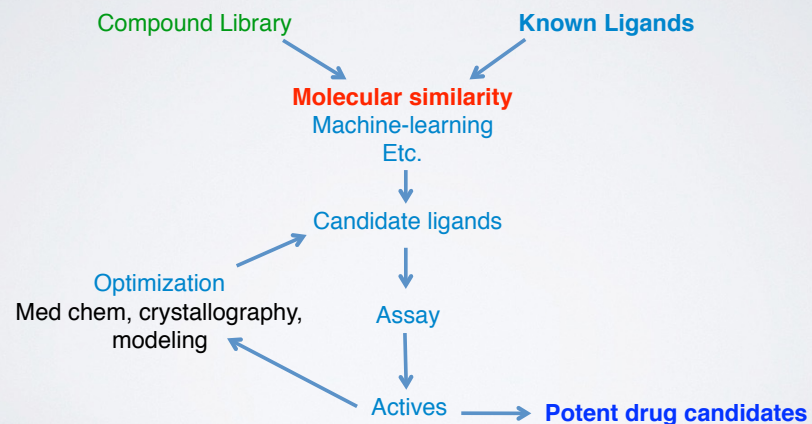Using knowledge of existing inhibitors to discover more

## 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

phenyl  naphthyl  ketone  methyl  ethyl  aldehyde  alcohol  amide  carboxylate  ...  S-S bond  chlorine  fluorine

Molecule 1

Molecule 2

# CHEMICAL SIMILARITY FROM FINGERPRINTS

phenyl naphthyl ketone methyl ethyl aldehyde alcohol amide carboxylate ... S-S bond chlorine fluorine

Molecule 1

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**

**Aromatic**

2.8 ±0.3 Å

---
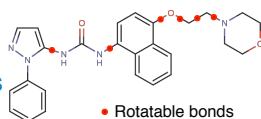
# Molecular Descriptors
## More abstract than chemical fingerprints

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 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.)

## Proteins and Ligand are Flexible

Ligand

Protein

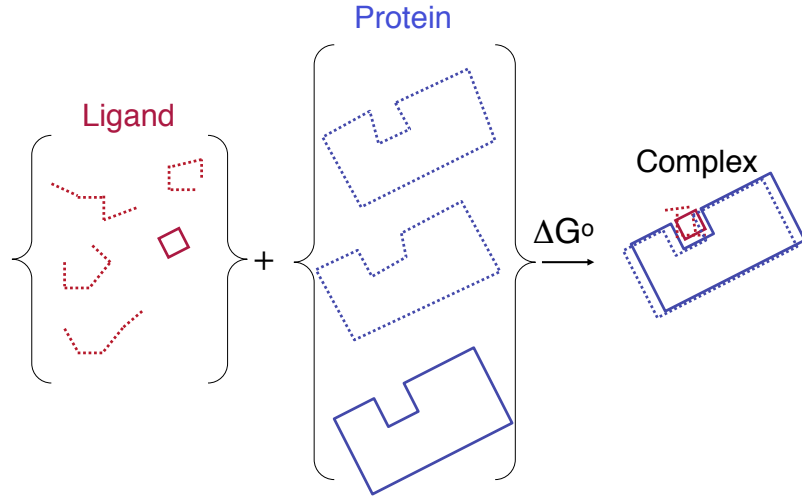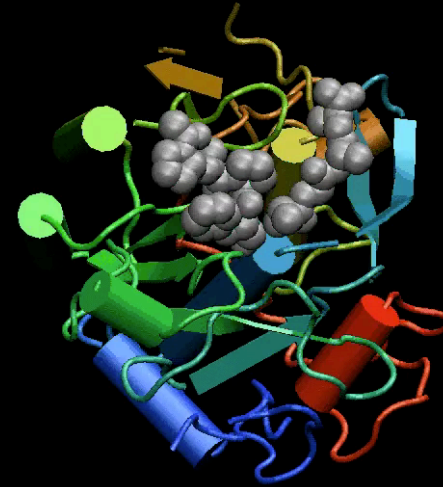Complex

$\Delta G^o$

---

**NMA** (Normal Mode Analysis) is a bioinformatics method to predict the intrinsic dynamics of biomolecules

https://bioboot.github.io/bggn213_W19/lectures/#12
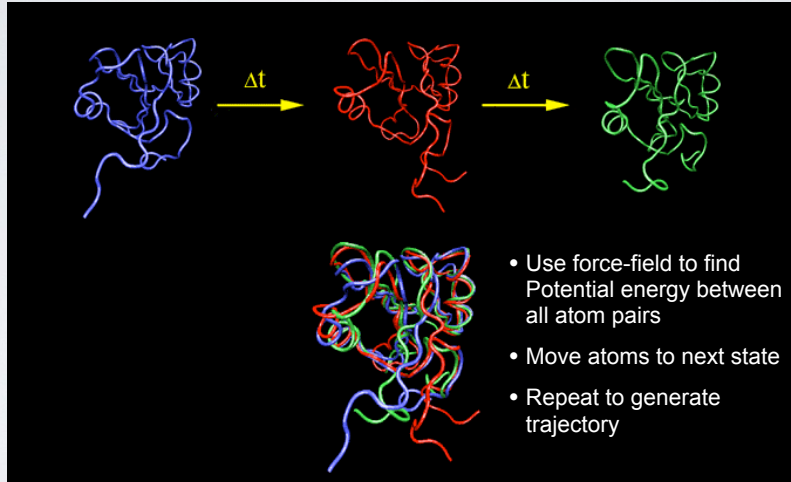
---

# Reference Slides

Molecular Dynamics (MD) and Normal Mode Analysis (NMA) Background and Cautionary Notes

[ Muddy Point Assessment ]

---

## 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



- 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)
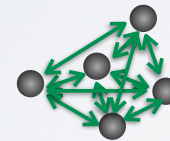[ See: https://www.youtube.com/watch?v=ui1ZysMFcKk ]

---

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



---

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



---

▷ Divide **time** into discrete (~1fs) **time steps** (**Δt**)
(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*
$$E(\vec{R}) = \sum_{\text{bonded}} E_i(\vec{R}) + \sum_{\text{non-bonded}} E_i(\vec{R})$$

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

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

$t$

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

*Nucleic motion described classically*

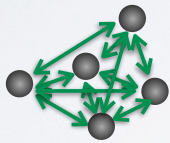$$m_i \frac{d^2}{dt^2}\vec{R}_i = -\vec{\nabla}_i E(\vec{R})$$

*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)

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

$$\boldsymbol{r}(t+\Delta t) = \boldsymbol{r}(t) + \boldsymbol{v}(t+\tfrac{\Delta t}{2})\Delta t$$

---

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

## BASIC ANATOMY OF A MD SIMULATION

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

$t$

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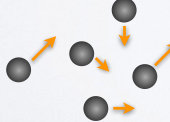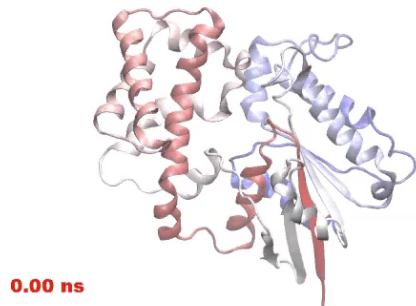(by evaluating **force-field** gradient)

*Nucleic motion described classically*

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

*Empirical force f...*

$$E(\vec{R}) = \sum \dots E_i(\vec{R})$$

▷ Use the for... ...te **velocities** and move atoms to new **positions**
...g numerically via the "leapfrog" scheme)

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

$$\boldsymbol{r}(t+\Delta t) = \boldsymbol{r}(t) + \boldsymbol{v}(t+\tfrac{\Delta t}{2})\Delta t$$

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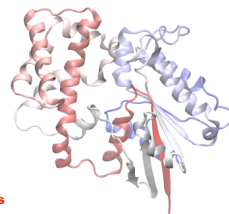
---

**Slide 3 (bottom-left):**

# MD Prediction of Functional Motions

Accelerated MD simulation of
nucleotide-free transducin alpha subunit

0.00 ns

Yao and Grant, Biophys J. (2013)

"close"

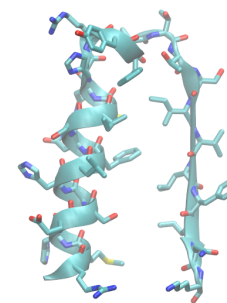0.00 ns

"open"

60.00 ns
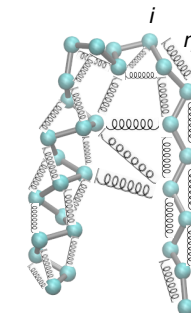
---

**Slide 4 (bottom-right):**
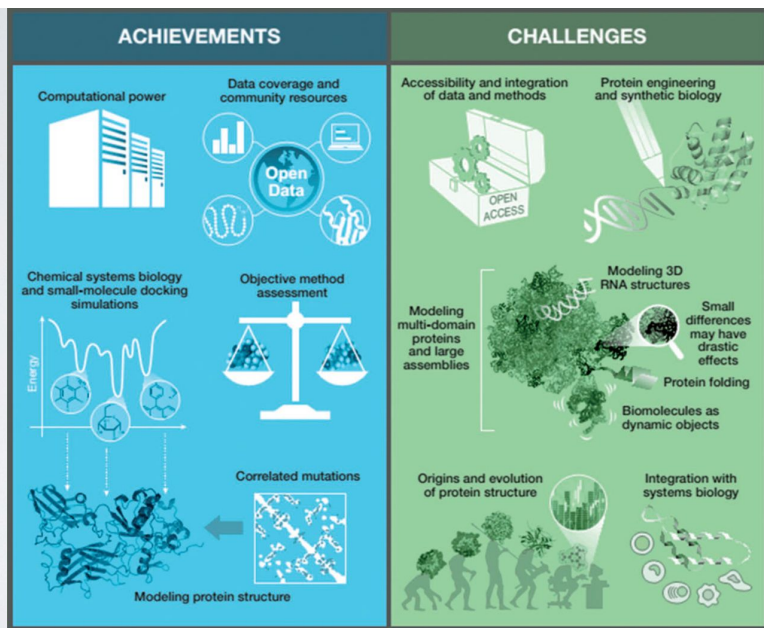
## COARSE GRAINING: **NORMAL MODE ANALYSIS**
(NMA)

• 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
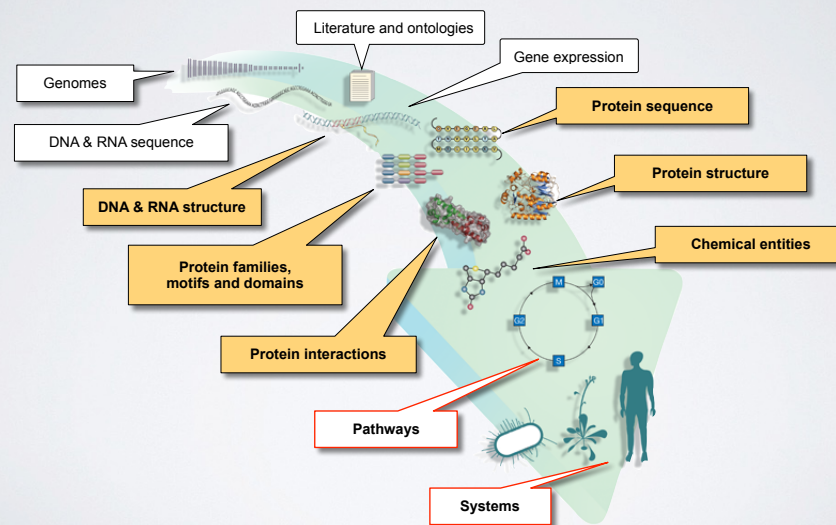even for large systems.

C. G.

$i$
$r_{ij}$
$j$

• 1 bead /
1 amino acid
• Connected by
springs

Atomistic

Coarse Grained

## Slide 1



ACHIEVEMENTS | CHALLENGES

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

## Slide 2

# INFORMING SYSTEMS BIOLOGY?



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

## Slide 3

# SUMMARY

- Structural bioinformatics is computer aided structural biology

- Described major motivations, goals and challenges of structural bioinformatics

- Reviewed the fundamentals of protein structure

- Explored how to use R to perform structural bioinformatics analysis!

- 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

[ Muddy Point Assessment ]

## Slide 4

# CAUTIONARY NOTES

- **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 parameters is an ongoing 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.