

Course Outline

Atomic-level modeling of proteins / macromolecules
Protein structure; Energy functions + molecular conformation; MD simulation; protein structure prediction; protein design; ligand docking

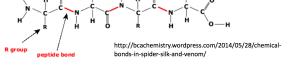
Coarser-level modeling and imaging-based methods

Fourier transforms and convolution; Image analysis, X-ray crystallography; Single-particle electron microscopy, Microscopy, Diffusion and single-molecule simulation, Genome structure

ATOMIC-LEVEL MODELING

Protein structure. Proteins are constantly jiggling around; as a result, they have many states (think molten water). There's also a cell membrane (made of lipids). Each protein can be made structures; but they tend to be similar. Usually we talk about the average structure (typically what is measured by X-ray crystallography in PDB). Surrounding molecules play a key role in determining structure.

2D Protein structure. Proteins are chains of amino acids. Amino acids are identical except for their sidechains. Proteins have regular backbones with differing side chains.



3D Protein structure - basic interactions Bond length stretching, bond angle bending, torsional angle twisting, electrostatic interaction

3D Protein structure - complex interactions Hydrogen bonds, hydrophobic effect, particularly important.

ENERGY FUNCTIONS + MOLECULAR CONFORMATION

Potential energy functions Potential energy function $U(x)$ specifies total potential energy of a system of atoms; as a function of all their positions x . For a system with n_A atoms, x is a vector of length $3n_A$ (x_1, y_1, z_1 coordinates for every atom). More generally, include not only atoms in a protein but also surrounding atoms.

Force calculation. Take derivatives of potential energy to obtain force.

Example energy function.

$$\begin{aligned} U = & \sum_{\text{bonds}} k_b (h-b)^2 & \text{Bond lengths ("Stretch")} \\ & + \sum_{\text{angles}} k_\theta (\theta-\theta_0)^2 & \text{Bond angles ("Bend")} \\ & + \sum_{\text{torsions}} k_\phi [1 + \cos(n\phi - \phi_0)] & \text{Torsional/dihedral angles} \\ & + \sum_{\text{electrostatics}} \frac{q_i q_j}{r_{ij}} & \text{Electrostatic} \\ & + \sum_{\text{Vander Waals}} \frac{A_{ij}}{r_{ij}^6} - \frac{B_{ij}}{r_{ij}^{12}} & \text{Non-bonded terms} \end{aligned}$$

Boltzmann distribution. Relates potential energy of an arrangement of atoms to the probability of observing that arrangement of atoms (At equilibrium).

$$p(x) \propto \exp(-U(x)/(k_B T))$$

- At low temperature, $P(C) > P(A)$
- At high temperature, $P(A) > P(C)$

All temperature differences in energy transfer calculations are relative to the energy difference between the probability distributions of the two conformations. Since reference C is arbitrary, we can choose it to be zero.

Macrostates. We typically care most about probability that protein atoms will be in some approximate arrangement, with any arrangement of surrounding atoms. Macrostates are "wells" of the potential energy function. To calculate the probability of a well, sum the probabilities of all specific atomic arrangements it contains.

$$P(A) = \int_{x \in A} P(x) \propto \int_{x \in A} \exp(-U(x)/k_B T) dx$$

Free energies. The free energy G_A of a macrostate A satisfies $P(A) = \exp(-G_A/(k_B T))$. Free energies are clearly defined only for macrostates. However, in protein structure prediction, design, and docking; often useful to define a "free energy function" that approximates a free energy function for some neighborhood of each arrangement of protein atoms. To predict protein structure - minimize free energy, not potential energy. The term energy function is used for both potential energy and free energy functions.

MOLECULAR DYNAMICS

MD overview. An MD simulation predicts how atoms move based on physical models of their interactions. Closest to the physics; aims to predict real dynamics of system. Can capture structural changes in proteins; protein-ligand binding, or protein folding.

Basic MD algorithm. Step through time (very short steps). At each step, calculate force acting on every atom using a mechanics formulae based on forces. Update atom positions and velocities using Newton's second law. $\frac{dx}{dt} = v; \frac{dv}{dt} = F(x)$

Note - this is an approximation, because we are using classical physics rather than quantum mechanics. But quantum mechanical calculations can be used to parametrize force fields.

Why MD is computationally intensive. One needs to take millions to trillions of timesteps to get to timescales on which events of interest take place. Computing the forces at each step involves substantial computation. Particularly for non-bonded forces! Which at best every pair of atoms.

Sampling. Given enough time, an MD simulation will sample the full Boltzmann distribution of the system. Hence - if one took a snapshot from the simulation after a long period of time; then the probability of the atoms being in a particular arrangement is given by the Boltzmann distribution.

One can also sample the Boltzmann distribution in other ways, including Monte Carlo sampling with the Metropolis criterion - Metropolis Monte Carlo: generates moves at random. Accept any move that decreases the energy. Accept moves that increase the energy by ΔU with probability $\exp(-\Delta U/(k_B T))$. - If one decreases the temperature over time, this becomes a minimization method (simulated annealing).

PROTEIN STRUCTURE PREDICTION.

Overview. The goal: given the amino acid sequence of a protein, predict its average three-dimensional structure. In theory, one could do this by MD simulation, but that isn't practical. Practical methods for protein structure prediction take advantage of existing data on protein structure (and sequence)

Two approaches. Template-based modeling (homology modeling) - Used when one can identify one or more likely homologs of known structure (usually the case) • Ab initio structure prediction - Used when one cannot identify any likely homolog of known structure - Even ab initio approaches usually take advantage of available structural data, but in more subtle ways

Template based modeling. Workflow. • User provides a query sequence with unknown structure • Search the PDB for proteins with similar sequence and known structure. Pick the best match (the template). • Build a model based on that template - One can also build a model based on multiple templates, where different templates are used for different parts of the protein.

Structure is more conserved than sequence. Proteins with similar sequences tend to be homologous, meaning that they evolved from a common ancestor - The fold of the protein (i.e., its overall structure) tends to be conserved during evolution • The tendency is very strong. Even proteins with 15% sequence identity usually have similar structures. - During evolution, sequence changes are more frequent than structure changes

Ab initio modeling - typically Rosetta. Search for structures that minimizes a energy function - This energy function is knowledge-based (informed, in particular, by statistics of the PDB), and it approximates a free energy function - Use a knowledge-based search strategy - Rosetta uses a Monte Carlo search method involving "fragment assembly," in which it tries replacing structures of small fragments of the protein with fragment structures found in the PDB

Protein design. • Goal: given a desired approximate three-dimensional structure (or structural characteristics), find an amino acid sequence that will fold to that structure - In principle, we could do this by searching over sequences and doing ab initio structure prediction for each possible sequence, but that's not practical

Simplifying the problem. - Instead of predicting the structure for each sequence considered, just focus on the desired structure, and find the sequence that minimizes its energy - Energy is generally measured by a knowledge-based free energy function - Consider a discrete set of rotamers for each amino acid side chain - Minimize simultaneously over identities and rotamers of amino acids - Assume the backbone is fixed - Or give it a bit of "wiggle room"

Heuristic but effective. These simplifications mean that protein structure prediction methodologies are highly heuristic, but they're proven surprisingly effective - The minimization problem itself is also usually solved with heuristic methods (e.g., Metropolis Monte Carlo)

LIGAND DOCKING.

Goals: - Given a ligand known to bind a particular protein, determine its binding pose (i.e., location, orientation, and internal conformation of the bound ligand) - Determine how tightly a ligand binds a given protein

Contour map of electron density:

Solving for molecular structure. Step 1: Initial phasing - Come up with an approximate solution for the structure (and thus an approximate set of phases), often using a homologous protein as a template - Step 2: Refinement - Search for perturbations that improve the fit to the experimental data (the diffraction pattern), often using simulated annealing - Restrict the search to "realistic" molecular structures, usually using a molecular mechanics force field

SINGLE PARTICLE EM.

Overview. We want the structure of a "particle": a molecule (e.g., protein) or a well-defined complex composed of many molecules (e.g., a ribosome) - We spread identical particles out on a film, and image them using an electron microscope - The images are two-dimensional, each representing a projection of the 3D shape (density) of a particle. Each particle is positioned at a different, unknown angle. - Given enough 2D images of particles, we can computationally reconstruct the 3D shape of a particle

Electron beam

Particles

Images

Predicts...

- The pose of the molecule in the binding site
- The binding affinity or a score that reflects the strength of binding

Binding affinity. Binding affinity quantifies the binding strength of a ligand to a protein (or other target) - Conceptual definition: we want the protein and the ligand (with no other ligands) to have a higher binding affinity than the protein and a ligand bound to it - Affinity can be expressed as the difference ΔG in free energy of the bound state and the unbound state, or as the concentration of unbound ligand molecules at which half the protein molecules will have a ligand bound

Docking is heuristic. In principle, we could estimate binding affinity by measuring the fraction of time the ligand is bound in an MD simulation, but this isn't practical

Docking methodology. Ligand docking is a fast, heuristic approach with two key components - A scoring function that very roughly approximates the binding affinity of a ligand to a protein given a binding pose - A search method that searches for the best-scoring binding pose for a given ligand - Most ligand docking methods assume that - The protein is rigid - The approximate binding site is known - The first one is to search for ligands that will bind to a binding site in the target - In reality, proteins, protein mobility, and water molecules all play a role in determining binding affinity - Docking is approximate but useful - The scoring function is used instead of energy function to emphasize the highly approximate nature of the scoring function

Does this apply to RNA? Some of these concepts apply with little or no modification to RNA and other biomolecules - Energy functions and their relationship to molecular conformation - Molecular dynamics simulation - Ligand docking - In other cases, the basic ideas apply, but the techniques are different - (RNA) structure prediction - (RNA) design

COASER MODELING + IMAGING

Fourier Transform. Writing functions as sums of sinusoids

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: 1.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude: -0.3 Phase: -0.9

Magnitude: 1.9 Phase: -0.9

Magnitude: 0.27 Phase: -1.4

Magnitude: 0.39 Phase: 2.8

Given a function defined on an interval of length L , we can write it as a sum of sinusoids with the following frequencies/periods: - Frequencies: $0, 1/L, 2/L, 3/L, \dots$ - Periods: constant, term $L_1/L, L_2/L, \dots$ Each of these sinusoidal terms has a magnitude (scale factor) and a phase (shift).

We can thus express the original function as a series of magnitude and phase coefficients as a complex number - The Fourier transform maps the function to this set of complex numbers, providing an alternative representation of the function. - This also works for functions of 2 or 3 variables (e.g., images) - Fourier transforms can be computed efficiently using the Fast Fourier Transform (FFT) algorithm

Convolution. A weighted moving average. To convolve one function with another, we computed a weighted moving average of one function using the other function to specify the weights

Original function

Sum of sinusoids below

Magnitude:

Take away from Rosetta energy functions: The coarse-grained Rosetta energy function is essentially entirely knowledge-based – Based on statistics compiled from the PDB + Many of the terms are of the form $-\log_e(P(A))$, where $P(A)$ is the probability of some event A – This is essentially the free energy of event A. Recall definition of free energy:

$$G_A = -k_B T \log_e(P(A)) \quad P(A) = \exp\left(-\frac{G_A}{k_B T}\right)$$

Simplifying assumptions of Rosetta • Still makes simplifying assumptions: – Do not explicitly represent solvent (e.g., water) – Assume all bond lengths and bond angles are fixed – Functional forms are a hybrid between molecular mechanics force fields and the (coarse-grained) Rosetta energy function – Partly physics-based, partly knowledge-based

PE functions or FF functions? • The energy functions of previous lectures were potential energy functions. • One can also attempt to construct a free energy function, where the energy associated with a conformation is the free energy of the set of “similar” conformations (for some definition of “similar”) • The Rosetta energy functions are sometimes described as potential energy functions, but they are closer to approximate free energy functions – This means that searching for the “minimum” energy is more valid – Nevertheless, typical protocol is to repeat the search process many times, cluster the results, and report the largest cluster as the solution. This rewards wider and deeper wells.

Rosetta fragment assembly • Uses a large database of 3-residue and 9-residue fragments, taken from structures in the PDB • Monte Carlo sampling algorithm: Step 1: Start with the fragment in its intended conformation. Step 2: Randomly select a 3-residue or 9-residue section. Step 3: Find a move in the library whose sequence resembles it. Step 4: Consider a move in which the backbone dihedrals of the selected section are replaced by those of the fragment. Calculate the effect on the entire protein structure. Step 5: Evaluate the Rosetta energy function before and after the move. Step 6: Use the Metropolis criterion to accept or reject the move. Step 7: Return to step 2 • The search algorithm adds some bells and whistles

Rosetta refinement • Refinement is performed using the Rosetta altatom energy function, after building in side chains • Refinement involves a combination of Monte Carlo moves and energy minimization – The Monte Carlo moves are designed to perturb the structure much more gently than those used in the coarse search – May still involve the use of fragments

Best structure prediction method • Currently, it's probably I-TASSER – <http://zhanglab.mech.umd.edu/I-TASSER/> • I-TASSER is template-based, but it uses threading, meaning that when selecting a template it maps the query sequence onto the template structure and evaluates the quality of the fit – This allows it to incorporate homologous I-TASSER-like multiple alignments – It incorporates a surprisingly large number of different components and strategies, including an ab initio prediction module – It runs many algorithms in parallel and then looks for a consensus between the results • Example: at least seven different threading algorithms – Ineligible but effective

Structure prediction game FoldIt, Eterna (RNA design), allow players to fold / with the goal of achieving minimum energy conformations.

Comparing structures of a protein Comparing structures of a protein • The most common measure of similarity between two structures for a given protein is root mean squared distance/deviation (RMSD), defined as

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - w_i)^2}.$$

where x gives the coordinates for one structure and w the coordinates for the other • We generally want to align the structures, which can be done by finding the rigid-body rotation and translation of one structure that will minimize its RMSD from the other – The relative measure of similarity is RMSD after alignment

GENE STRUCTURE

Epigenetics • Epigenetics is the study of how the same genome produces different functions. – DNA can be chemically marked – DNA can be left accessible or packaged away – RNA is also regulated in many different ways • Gene regulation is closely tied to genome organization.

DNA is stiff

- DNA has repeating negative charge – it is a polyelectrolyte
- The negative charges repel, making DNA quite inflexible at the scale of 10s of base pairs
- Ions in solution affect DNA flexibility – Monovalent ions screen charge – Divalent cations (Mg^{2+}) promote compaction – Neutral salt increases stiffness
- How do you compact DNA? By wrapping it around something with positive charges.

Nucleosome

- Around 150bp of DNA winds around an octamer of eight histone proteins
- Together, the DNA and the histone proteins are called a nucleosome
- Histones have many positively charged amino acids
- This wrapping compacts the genome linearly by about 7x
- Stringing multiple nucleosomes together gives chromatin – DNA and its associated proteins.

Chromatin fiber

The chromatin fiber

- Repeating nucleosomes form the 10nm fiber, like beads on a string
- In a test tube, chromatin condenses to form a thick fiber about 30nm in diameter



Structure of 30nm fiber

- Two main models: solenoid and zigzag
- Use coarse-grained simulations to examine possibilities
- Use metropolis Monte Carlo
- Parameters: linker length, linker angle, turn angle, nucleosome thickness
- Consider various physical measurements: thickness, mass density, spring behavior

Nucleosomes / gene regulation

Nucleosomes affect gene expression in two important ways:

1. Presence of a nucleosome at the beginning of a gene inhibits its transcription
 2. Chromatin can be *accessible* or *compacted*. Compaction inhibits transcription
- This is determined by chemical modifications on histone tails

Chromosome conformation capture

- Under a microscope, all DNA sequences “look the same,” so you can’t match your image to genomic position
- There are ways to mark specific DNA sequences (fluorescence in situ hybridization), but very low throughput
- We can use DNA sequencing to map self-contacts of the folded 3D genome
- This method is called **chromosome conformation capture (3C)**

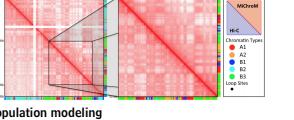
3D genome reconstruction

Given a map of self-contacts, how can you reconstruct the structures that produced it?

1. Knowledge-based: 3C or other data can be used as constraints
2. Mechanism-based: propose a mechanism that organizes the genome and simulate it
3. Mechanisms that organize the genome are still quite mysterious. Need computation to tackle this problem.

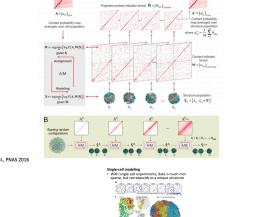
Minimal Chromatin Model

- **Minimal Chromatin Model** (MiChrom)
- Knowledge-based: use 1D compartment info
- Beads in the same compartment prefer to interact
- Train energy function on one chromosome and apply it to the whole genome



Population modeling

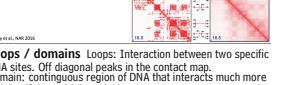
- Assignment step: decide which contacts are most important for each individual structure
- Modeling step: optimize the structure based on those constraints



Single cell modeling

Binding of multivalents protein

- Binding of multivalent proteins that prefers to bind to one “type” of chromatin, and binds to multiple sites
- This causes clustering that generates compartments.
- Mechanism-based method



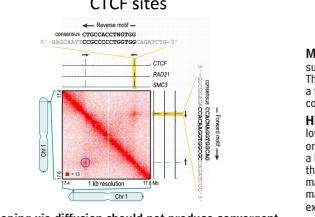
Loop formation correlates with nearby gene activation.

Anchored at convergent CTCF sites

Loops are anchored at convergent CTCF sites

Loops are anchored at convergent CTCF sites

Looping via diffusion should not produce convergent sites.



Values should sum to 1, so overall brightness of the image remains constant

Why 1/9? Values should sum to 1, so overall brightness of the image remains constant

Median filter A median filter ranks the pixels in a square surrounding the pixel by the average of a set of pixels

– For example, a 3x3 median filter replaces each pixel with the average of a 3x3 square of pixels

– This is equivalent to convolving the image with a 3x3 matrix:

Overview • Why design proteins? • Overall approach: Simplifying the protein design problem • Protein design methodology – Designing the backbone – Select sidechain rotamers: the core optimization problem – Optional: giving the backbone wiggles room – Optional: negative design – Optional: complementary experimental methods • Examples of recent designs • How well does protein design work?

Applications of protein design Sample applications • Designing enzymes (proteins that catalyze chemical reactions) – Required for production of industrial chemicals and drugs • Designing proteins that bind specifically to other proteins – Proteins that detect TNT – Proteins that detect SARS-CoV-2 spike antibody design • Designing sensors (protein that bind to and detect the presence of small molecules—for example, by lighting up or changing color) – Calcium sensors used to detect neuronal activity in imaging studies • Proteins that detect TNT or other explosives, for mine detection • Making a more stable variant of an existing protein – Or a water-soluble variant of a membrane protein

Direct approach to protein design • Computationally intractable – We’ll need to use ab initio structure prediction – Ab initio structure prediction for even one sequence is a complex problem, involving many degrees of freedom with N residues • May not be good enough – Ab initio structure prediction is far from perfect, in part because energy functions are imperfect – Given an energy function, what we really want is to maximize the probability of the desired structure (compared to all other possible folded and unfolded structures) – We could do this by sampling the full Boltzmann distribution for each candidate sequence ... but that’s even more computationally intensive!

How to simplify protein design • We can dramatically simplify this problem by making a few assumptions 1. Assume the backbone geometry is fixed 2. Assume each amino acid can only take on one of a few possible conformations 3. Assume that what we want to do is to maximize the energy driven from the completely unfolded state to the target geometry – In other words, simply ignore all the other possible folded structures that we want to avoid

Simplified problem • At each position on the backbone, choose a rotamer (an amino acid type and a side-chain geometry) to minimize overall energy – We assume the energy is a free energy. The Rosetta all-atom force field (physics-based/knowledge-based hybrid) is a common choice.

– For each amino acid sequence, energy is measured relative to the unfolded state. In practice, the reference energy of each amino acid is subtracted off, corresponding roughly to how much that amino acid favors folded states – You’re not responsible for this – Assume that energy can be expressed as a sum of terms that depend on one rotamer or two rotamers each. This is the case for the Rosetta force fields (and for most molecular mechanics force fields as well).

Want to minimize total free energy:

$$E_T = \sum_i [E_i(r_i) + \sum_{j \neq i} E_{ij}(r_i, r_j)]$$

Designing backbone • The first step of most protein design protocols is to select one or more target backbone structures. – These are usually chosen by hand, but in some cases, many different structures are selected, because some won’t work. (Apparently proteins can only adopt a limited set of backbone structures, but we don’t have a great description of what that set is.) • Methods to do this: – Use an experimentally determined backbone structure – Assemble secondary structural elements by hand – Use a fragment assembly program like Rosetta, selecting fragment combinations that fit some approximate desired structure

Optimization methods for sidechain rotamers

– Heuristic methods: Not guaranteed to find optimal solution, but fast and efficient to implement – Monte Carlo moves may be as simple as randomly choosing a position, then randomly choosing a new rotamer at that position • May decrease temperature over time (simulated annealing)

Flexible backbone design – Better to have “wiggle room” in backbone – This requires optimizing simultaneously over rotamers and backbone geometry – Often addressed through a Monte Carlo search procedure that alternates between local tweaks to backbone dihedrals and changes to side-chain rotamers

Negative design • In negative design, we identify a few structures that we want to avoid, and we try to keep their energies high during the design process. – This can help, but we cannot explicitly avoid all possible incorrect structures without making the problem much more complicated. So the overall approach is still heuristic.

Experimental methods Computational protein design is often coupled with experimental validation of the predicted structures • For example, computational designs can often be improved by directed evolution. Directed evolution involves introducing random mutations to proteins and picking out the best ones – Usually this is done in living cells, with the fittest cells (i.e., those containing the “best” version of the protein) selected by some measure

Median filter

• The simplest way is to use a “mean filter”

– Replace each pixel by the average of a set of pixels surrounding it

– For example, a 3x3 mean filter replaces each pixel with the average of a 3x3 square of pixels

– This is equivalent to convolving the image with a 3x3 matrix:

$$\begin{bmatrix} \frac{1}{9} & \frac{1}{9} & \frac{1}{9} \\ \frac{1}{9} & \frac{1}{9} & \frac{1}{9} \\ \frac{1}{9} & \frac{1}{9} & \frac{1}{9} \end{bmatrix} = \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

Median filter A median filter ranks the pixels in a square surrounding the pixel by the average of a set of pixels

– Under certain conditions, entire proteins will pack into a regular grid (a “lattice”) (e.g., insulin crystals)

Electron density • The electron density corresponding to the 3D structure of a molecule gives the probability of finding an electron at each point in space • X-rays bounce off electrons they hit

What causes diffraction patterns? • Short answer: interference of light – The bright spots are places where light interferes constructively. Elsewhere it tends to interfere destructively (cancel out).

Diffraction pattern; electron density • It turns out that the diffraction pattern is the Fourier transform of the electron density – Both the electron density and the diffraction pattern are functions of three dimensions (i.e., defined at every point in a 3D volume) – Each diffraction pattern corresponds to one sinusoidal component of the electron density – The Fourier transform gives a magnitude and a phase for each sinusoid, but it’s only practical to measure the magnitude, not the phase • Brightness of the spot gives the magnitude

Initial phasing • The most common method for initial phasing is molecular replacement – Start with a computational model of the protein structure (often the structure of a homologous protein)

– Search over the possible ways that a protein with this structure could be packed into a crystal, and find the one that gives the best fit to the data

Phase refinement • Once we have an initial model, we can search for perturbations to that model that improve the fit to the experimental data – This is usually done through a Monte Carlo search (via simulated annealing). – One useful refinement technique is “realistic” protein structures using a molecular mechanics force field – This dramatically improves the accuracy of the results • The idea was introduced by Axel Brunger, now on the Stanford faculty

– Typically done with cross validations.

Misc notes on crystallography • Protein crystals contain water – Often half the crystal is water (filling all the empty spaces between copies of the protein) – Usually only a few water molecules are visible in the structure, because the rest are too mobile

– You can’t determine hydrogen positions by x-ray crystallography – But one can model them in computationally – Some high-quality, published crystal structures have turned out to be completely incorrect due to computational problems/errors

Single particle EM (cryo-EM)

– Single-particle EM has come a long way for decades, it has improved dramatically in the last five years due to: – Improved electron cameras • Until recently electrons were detected either by photographic film, or by scintillator-based digital cameras which converted electrons to photons for detection • New “direct-electron detectors” can detect electrons directly, substantially improving image resolution and quality – Better computational reconstruction techniques • Single-particle EM is turning into much wider use, and may challenge crystallography as the dominant way to determine experimental structures

Single particle EM vs. X-ray crystallography • Single-particle EM’s major advantage over crystallography is that it does not require formation of a crystal – Particularly advantageous for large complexes, which are usually difficult to crystallize – Also avoids structural artifacts due to packing in a crystal lattice. In EM, particles are in a more natural environment

– On the other hand – Single-particle EM’s resolution is typically lower than X-ray crystallography. Resolving the structures of small proteins in EM images is difficult, because the structures from different orientations look similar (i.e., “a blob”)

– Bottom line: single-particle EM is particularly advantageous for large complexes, because – Large complexes tend to be harder to crystallize – The computational reconstruction problem in single-particle EM is usually easier to solve for large particles than for small ones

Single particle EM uses transmission electron microscopy • In transmission electron microscopy, a beam of electrons pass through a thin sample before forming an image

Beating diffraction limit. Decrease wavelength • Higher-frequency radiation (e.g., x-rays) has shorter wavelengths and thus allows higher resolution – It also damages the sample more • It’s possible to image with electrons, which have a much shorter wavelength (.1 nm) – Electron microscopy can thus achieve much higher resolution (below 4 Å) – Disadvantages: can’t use living cells, and molecules of interest won’t glow

Super resolution microscopy • A number of recently developed techniques achieve resolution well beyond the diffraction limit – This requires violating some of the assumptions of the theory – I’ll briefly describe the most popular of these techniques, known alternately as STORM (stochastic optical reconstruction microscopy) or PALM (photoactivation localization microscopy)

STORM If we have only a few fluorophores in an image, we can use them very accurately – That by itself is not very useful, but if we have many, we can turn them on at a time, identifying their locations in each image, and then use that information (computationally) across many images, we can build a composite image of very high resolution

Fick's law • Suppose that particles are uniformly distributed in the y and z dimensions, and vary only in x . Let c represent concentration (a function of x) • Define the flux J as the rate at which particles diffuse across a boundary – Fick's 1st law:

$$J = -D \frac{\partial c}{\partial x}.$$

Concentration

- New concentration and flux as a function of position x and time t
- The concentration at a particular position goes down with time if there is more flux away from that position than there is coming to it in that position (in other words, if the net flux is increasing as one moves in the positive x direction)

$$\frac{\partial c}{\partial t} = -D \frac{\partial^2 c}{\partial x^2}$$

$$\frac{\partial c}{\partial t} = D \frac{\partial^2 c}{\partial x^2}$$

(Can be generalized to multiple dimensions).

Continuum approach • Advantage: faster • Disadvantage: less accurate for small numbers of molecules • Unlike the particle-based approach, the continuum approach is deterministic

Why is docking useful? • Virtual screening: Identifying drug candidates by considering large numbers of possible ligands

• Lead optimization: Modifying a drug candidate to improve its properties • Once a candidate is found, docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how binding energy is modified

• Docking can help validate the binding of the ligand would affect how