Topic Course

Gene and protein evolution

Lecture 6

Winter 2016 Department of Molecular Genetics University of Toronto

Hue Sun Chan

Theory of protein sequence space organization and the dynamics of molecular evolution

- Sequence-space topology, neutral nets, and their biophysical basis
- Simple exact models of the mapping between protein sequences and structures; the hydrophobic-polar model
- Superfunnels and mutational robustness, correlation between thermodynamic and mutational stabilities
- Recombination, local conformational preference, and autonomous folding units
- Promiscuous functions and conformational switches
- Escape from adaptive conflicts: mutational robustness and subfunctionalization

Reference: Sikosek & Chan (2014) Biophysics of protein evolution and evolutionary protein biophysics. *J Royal Society Interface* **11**:20140419 (http://rsif.royalsocietypublishing.org/content/11/100/20140419.full).

Models of structural evolution of proteins are based on a presumed mapping between protein sequences and structures and a presumed relationship between structure and function



The NK model

The spin-glasses inspired **NK** model defines a combinatorial space consisting of every string **s** of length **N**, each site of the string is chosen from an alphabet of **A** alternative states – e.g., the 20 amino acids. A scalar value (which may be identified with "fitness" **F**) is assigned to each string, where

$$F(s) = \frac{1}{N} \sum_{i=1}^{N} F^{(K)}{}_{i}(s_{i}; s_{i1}, \dots, s_{iK})$$

If a distance metric is defined between strings (e.g. Hamming distance, i.e., number of amino acid substitutions), the resulting construct is a fitness landscape.

Example of a model sequence space: A *four-dimensional* Boolean *hypercube*, in which each of the 16 vertices represents one of the possible strings of four 0 or 1 values. Here each such string is interpreted as a specific tetrapeptide with two possible types of amino acid at each position.

Kauffman & Levin (1987) Towards a general theory of adaptive walks on rugged landscapes. *J Theor Biol* **128**:11-45. Kauffman & MacReady (1995) Search strategies for applied molecular evolution. *J Theor Biol* **173**:427-440.

Figure from the following references for the NK model:



Example of genotype-phenotype map in the modeling of evolution

Schematic representation of the genotype–phenotype map used in the analysis of Draghi et al.* Each circle corresponds to a genotype; colours denote phenotypes. A proportion q of mutations are neutral (solid lines) and the remaining mutations are non-neutral (dashed lines). A non-neutral mutation changes an individual's phenotype to one of the *K* accessible alternatives that form the individual's phenotypic neighbourhood. When *K* is smaller than the total number of alternative phenotypes in the landscape, *P*, individuals may have different phenotypic neighbourhoods. The central pair of adjacent genotypes shown here express the same phenotype, but they have different phenotypic neighbourhoods.

Figure from: Draghi, Parsons, Wagner & Plotkin (2010) Mutational robustness can facilitate adaptation. *Nature* **463**:353-355.

Concepts of genotype-phenotype (sequence-structure) map and neutral net have long been used in theories of protein evolution



Figure from: Harms & Thornton (2013) Evolutionary biochemistry: Revealing the historical and physical causes of protein properties. *Nature Review Genetics* 14:559-571.

Although theories based on *ad hoc* genotype-phenotype maps or maps that were not derived from explicit-chain protein models (e.g. the *NK* model) have led to much advance, it has become increasingly clear that sequence-structure maps based on explicit-chain biophysics are essential for gaining insights into the evolution of protein folds.



Figure from: Harms & Thornton (2013) Evolutionary biochemistry: Revealing the historical and physical causes of protein properties. *Nature Review Genetics* 14:559-571.

Conceptual Models of Protein Evolution

An atomistic sequence-structure map is currently not feasible to construct, because such a feat amounts to first solving the protein folding problem for all possible amino acid sequences. In view of this fundamental limitation in theory and computation, highly simplified *explicit*-chain models that aim to capture the essential biophysics of protein folding are being used to develop theories of protein evolution.

Hydrophobic Interaction is a Major Driving Force for Protein Folding

Kauzmann (1959), Dill (1990), etc.



A Simple Exact Biophysical Model: The Hydrophobic-Polar (HP) Model



Lau & Dill, Macromolecules (1989); Chan & Dill, J Chem Phys (1991) • Figure from: Chan & Dill, Physics Today (1993)

The Hydrophobic-Polar (HP) Lattice Protein Model

- Hydrophobic-hydrophobic contacts are favorable
- Allows for exact (exhaustive) enumeration of sequences and conformations
- Provides an exact sequence-to-structure mapping for the study of evolution [Chan & Bornberg-Bauer, Appl Bioinformatics (2002) ; Xia & Levitt, Curr Opin Struct Biol 14, 202 (2004)]



What properties of real proteins can simple HP model capture?



• Not every property; e.g., HP lattice proteins undergo much less cooperative folding than real proteins. But the HP model is apparently quite reasonable as a model for the mapping between protein sequences and their ground-state (lowest-free-energy) structures.



HP pattern is an important determinant of protein structure

Sequence Pattern Statistics of HP Model Proteins are similar to that of Real Proteins



Figures and analysis taken from: Irbäck & Sandelin, Biophys J 79, 2252 (2000)

Early applications of the HP model to the study of evolution

Lipman & Wilbur* recognized that the twodimensional (2D) HP model can provide a biophysical assessment of J. Maynard-Smith's idea (1970) that 'if evolution by natural selection is to occur, functional proteins must form a continuous network which can be traversed by unit mutational steps without passing through non-functional intermediates' [Nature 225:563-564]. The authors pointed out that Maynard-Smith's concept is related to S. Wright's more general notion of an adaptive landscape (1932), and noted further that in modeling adaptive landscapes, the 2D HP model sequence-structure map offers a useful, more physical alternative to NK approaches that treated this map as random. Lipman & Wilbur designate HP sequences that do not have a unique ground-state contact pattern as nonviable and thus mutations to these non-encoding sequences as "lethal". Within this modeling construct, Maynard-Smith's condition was found to be satisfied by the 2D HP model.



• Importance of neutral mutations: Any path from phenotype A to C requires at least one neutral mutation in phenotype B (figure from Lipman & Wilbur, 1991*).

Lau & Dill (1990) Theory for protein mutability and biogenesis. PNAS 87:638-642.

Chan & Dill (1991) Sequence space soup of proteins and copolymers. J Chem Phys 95:3775-3787;
* Lipman & Wilbur (1991) Modeling neutral and selective evolution of protein folding. Proc R Soc B-Biol Sci 245:7-11.

Superfunnels: Correlation between Thermodynamic and Mutational Stabilities



Bornberg-Bauer & Chan (1999) *PNAS* **96**:10689-10694; Wroe, Bornberg-Bauer & Chan (2005) *Biophys J* **88**:118-131.

Master Equations may be used to model evolutionary population dynamics in the limit of an infinite population. Finitepopulation effects may be modeled by Monte Carlo simulations.



The steady-state population of a sequence is governed by the balance between gains from mutational influx from adjacent viable sequences (------>) and losses from viable mutations (----->) of the given sequence. Hence a sequence with more viable neighbors tend to accumulate a higher steady-state evolutionary population.

Bornberg-Bauer & Chan (1999) PNAS **96**:10689-10694; Chan & Bornberg-Bauer (2002) Applied Bioinformatics **1**:121-144.

Interplay of network topology and fitness effects on evolutionary population

Mutational robustness alone may not be sufficient to account for the experimentally observed concentration of evolutionary populations at prototype-like sequences. Selection for native stability and kinetic stability is evident.



Sikosek & Chan, J R Soc Interface (2014)



Neutral Nets and Supernets in Sequence Space

Other 2D lattice Protein Models

Models that restrict conformational variation to maximally compact conformations are computationally more tractable but are less physical Example:



The most frequently encoded structure in the standard 2D HP model of chain length n = 25, shown by one of the 326 sequences that encode uniquely for it [Irbäck & Troein (2002) *J Biol Phys* **28:**1-15]. All 5,768,299,665 square-lattice n = 25 conformations are considered to arrive at this ground-state structure.



The most frequently encoded structure in the modified HP model of Xia and Levitt (2002) [99:10382-10387] that restricts conformational variation to the 1,081 maximally compact conformations confined to a 5×5 square, shown by one of the 67,615 sequences encoding uniquely for this structure.



Size and Structure of Neutral Nets and Superfunnels Depend on Model Interaction Scheme



Another 2D lattice model example of superfunnel





Recombinatoric exploration of folded structures

more efficient than single-point mutations



Results from: Cui, Wong, Bornberg-Bauer & Chan (2002) Recombinatoric exploration of novel folded structures: A heteropolymer-based model of protein evolutionary landscapes. PNAS 99:809-814.



 \bigcirc = neutral net; 897 neutral nets (n = 18; 4,553 sequences total) in the big net

Cui, Wong, Bornberg-Bauer & Chan, PNAS (2002)



Sequence and structure exploration facilitated by crossovers



Local Signal for Sequence Uniqueness



• Model proteins prefer certain local sequence patterns.

• Consistent with a subsequent experiment on β -lactamase showing that for a given number of amino acid substitutions, recombined variants are much more likely to retain function than variants generated by random point mutations [Drummond, Silberg, Meyer, Wilke & Arnold, *PNAS* (2005)].

Cui, Wong, Bornberg-Bauer & Chan, PNAS (2002)



Cui, Wong, Bornberg-Bauer & Chan, PNAS (2002)

Recombinations are more effective than point substitutions in pushing evolutionary sequence population towards the prototype sequence



Results and figure from: Xia & Levitt (2002) Roles of mutation and recombination in the evolution of protein thermodynamics. PNAS **99:**10382-10387.

Latent Evolutionary Potentials:

Selection of excited-state (promiscuous) functions can speed up evolution dramatically





Model population dynamics with excited-state selections

• While the population drifts along network **A**, the changing sequence similarity indicates that it is getting closer to the target C but structurally the vast majority of the sequences still stay on the neutral network for structure A.

The upper and lower dotted curves give the maximum and minimum similarity of the sequences in the population to the target sequence.

• Evolution of the average pairwise Hamming distance between all pairs of sequences in the population shows a reduction in genotypic variation immediately before the steps of adaptation, while at the same time structural distance from the originating structure increases.

Number of sequences (out of 1000), which code uniquely for structure A (solid curve), B (dotted), and C (dashed) as

functions of generation.

Results and figure from: Wroe, Chan & Bornberg-Bauer (2007) HFSP J 1:79-87.



The simple 2D HP model provides a biophysical rationalization for the intriguing, and otherwise puzzling experimental observation that adaptation to new requirements can proceed while the "old," phenotypically dominant function is maintained along a series of seemingly neutral mutations.



Evolutionary transitions between neutral nets. **Top:** possible development of the relative fitness during adaptation of a population which was initially optimal for one function (\times , *left*) and then selected for another function (target: \Box , *right*). **Bottom:** an interpretation

in terms of structural phenotypes according to Maynard-Smith's idea of a continuous phenotype space (there is one direct transition from one representative of one structural phenotype to the other) and the superfunnel paradigm (genotypes depicted in the center have more neutral neighbors and code for thermodynamically most stable structures).

Tawfik and coworkers characterized ~300 variants of the enzyme PON1 that are apparently neutral, or close to neutral, with respect to PON1's levels of expression and native lactonase activity. Their activities with promiscuous substrates and ligands indicated significant changes in adaptive potentials.



"Inasmuch as the selection of nonnative functions and structures is operative, the attraction of any superfunnel towards its prototype may extend to proteins that are not yet part of its neutral network." – Sikosek et al, PNAS (2012)



Results and figure from: Wroe, Chan & Bornberg-Bauer (2007) HFSP J 1:79-87.

Conformational switches in real proteins

• Arc repressor homodimer: Double mutant L12 \leftrightarrow N11: $\beta \rightarrow \alpha$; single mutant N11L *bi-stable* [Cordes *et al* & Sauer, *Science* **284**, 325 (1999); *Nature Struct Biol* **7**, 1129 (2000)].

• One-mutation switch between the human serum albumin-binding domain (G_A) and the IgG-binding domain (G_B) of *Streptococcus* protein G [Alexander, He, Chen, Orban &



Bryan, *PNAS* **106**, 21149 (2009)]. NMR structures were determined for $G_A95 \& G_B95$ (the 56aa sequences are 95% identical). Three substitutions in G_A95 or G_B95 shift the equilibrium from >99% in one structure to >99% in the other.

G _A , G _B sequence			αl 000000000000000000000000000000000000		α2 <u>00000000</u>	α: 0000000	α3 2000000000000	
variants:		1	10	20	30	40	50	
	G _A 77	TTYKLIL	N <mark>l</mark> kQ <mark>a</mark> kee	AI <mark>K</mark> E <mark>L</mark> VDA <mark>G</mark>	IAEKY <mark>I</mark> KL <mark>I</mark> ANA	KTVEGVWT <mark>L</mark> I	KDE <mark>ILKA</mark> TVTE	
	G _A 88	TTYKLIL	NLKQAKEE	AI <mark>K</mark> E <mark>L</mark> VDAG	IAEKY <mark>I</mark> KL <mark>I</mark> ANA	KTVEGVWT <mark>L</mark> I	KDEI <mark>L</mark> TFTVTE	
	G _A 91	TTYKLIL	NLKQAKEE	AIKE <mark>L</mark> VDAG	TAEKY <mark>I</mark> KLIANA	KTVEGVWT <mark>L</mark> I	KDEI <mark>L</mark> TFTVTE	
	G_A95	TTYKLIL	NLKQAKEE	AIKE <mark>L</mark> VDAG	TAEKY <mark>I</mark> KLIANA	KTVEGVWT <mark>L</mark> I	KDEIKTFTVTE	
	$G_A 98$	TTYKLIL	NLKQAKEE	AIKELVDAG	TAEKYFKLIANA	KTVEGVWT <mark>L</mark> I	KDEIKTFTVTE	
	G-98	TTYKLIL	NIKOAKEE	AIKELVDAG	TAEKYFKLIANA	KTVEGVWT <mark>Y</mark>	VDEIKTFTVTE	
	G ₈ 95	TTYKLIL	NLKQAKEE	AIKE <mark>A</mark> VDAG	TAEKY <mark>F</mark> KLIANA	KTVEGVWT <mark>Y</mark>	KDEIKTFTVTE	
	G _B 91	TTYKLIL	NLKQAKEE	AIKE <mark>A</mark> VDAG	TAEKY <mark>F</mark> KLIANA	KTVEGVWT <mark>Y</mark> I	KDEI <mark>K</mark> TFTVTE	
	G _B 88b	TTYKLIL	NLKQAKEE	AI <mark>TE</mark> AVDAG	TAEKY <mark>F</mark> KL <mark>Y</mark> ANA	KTVEGVWT <mark>Y</mark> I	KDE I <mark>K</mark> TFTVTE	
	G _B 77	TTYKLIL	N <mark>G</mark> KQ <mark>l</mark> kee	AI <mark>T</mark> E <mark>A</mark> VDA <mark>A</mark>	TAEKY <mark>F</mark> KL <mark>Y</mark> ANA	KTVEGVWT <mark>Y</mark> I	KDE <mark>TKTF</mark> TVTE	
Figures from: Alexander et al (2009))	β1	→	β2	αl	β3	β4	



Model calculations we performed on the experimental G_A/G_B series of sequences indicate that the stability of the excited (latent) state gradually increases (i.e., with decreasing free energy) as the switch point is approached.

Sikosek, Bornberg-Bauer & Chan, PLoS Comput Biol (2012)





Experimentally designed bi-stable proteins and mutationinduced conformational switches



Bouvignies *et al.*, *Nature* (2011)

 bi-stable
 Cordes et al., Nature

 N11L
 L12N
 Cordes et al., Nature

 Arc repressor
 double mutant
 Struct Biol (2000)



GA 27 mut. L45Y

chameleon 🚫 P22 Cro 9 mut. 14 mut. 2 Cro

Alexander et al., PNAS (2009)

Meier et al.,

Curr Biol (2007)

Anderson *et al.*, *Protein Eng Des Sel* (2011)



Role of excited-state selection in Escape from Adaptive Conflict



The "Escape from Adaptive Conflict" Perspective focuses on adaptation before gene duplication [Hittinger & Carrol, *Nature* **449**, 677 (2007); Des Marais & Rausher, *Nature* **454**, 762 (2008)]



Sikosek et al., PNAS (2012); Sikosek et al., PLoS Comput Biol (2012); Sikosek & Chan, J R Soc Interface (2014)

Escape from Adaptive Conflict Follows from Weak Functional Trade-Offs and Mutational Robustness.

• Biophysics-based network connections.

• Evolutionary dynamics under mutations and gene duplications computed using both an analytical master equation and stochastic Monte Carlo simulations.

• Fitness is proportional to the stability (concentration) of the functional structures up to a certain optimum concentration above which fitness does not increase further with concentration.

• The optimum concentration corresponds to a measure of selection pressure.

Sikosek, Chan & Bornberg-Bauer (2012) *PNAS* **109**:14888-14893.



Master Equation of Evolutionary Population Dynamics with Single-Point Mutations and Gene Duplications

Now we need to consider two types of genotypes: (i) those with a single gene and (ii) those with two genes.



(ii) Genotypes with two genes:

$$P_{ij}(q+1) = \left[-2\mu n P_{ij}(q) + \mu \sum_{r=1}^{A_i} P_{\nu_i(r)j}(q) + \mu \sum_{s=1}^{A_j} P_{i\nu_j(s)}(q) + \mu^2 \sum_{r=1}^{A_i} \sum_{s=1}^{A_j} P_{\nu_i(r)\nu_j(s)}(q) + \delta_{ij}\mu_d P_i(q) + P_{ij}(q)\right] \frac{\mathcal{N}(q)W_{ij}}{\tilde{W}(q)}$$

The formalism can also be applied to Monte Carlo simulations of finite evolutionary populations.

Sikosek, Chan & Bornberg-Bauer (2012) Escape from adaptive conflict follows from weak functional trade-offs and mutational robustness. PNAS 109:14888-14893.

Neofunctionalization follows from strong selection pressures (strong trade-offs), whereas *Subfunctionalization* results from intermediate selection pressures (weak trade-offs)



Sikosek, Chan & Bornberg-Bauer (2012) PNAS 109:14888-14893.

Escape from Adaptive Conflict (EAC) in the real world

Experiments showing that a reconstructed common ancestor of the fluorescent proteins in corals that emit either red or green light can emit light of both colors are indicative of EAC.



Subfunctionalization (SUBF) as a Consequence of Neutral Network Topology

A strong tendency toward mutationally robust genotypes (i.e., those with many sequence-space neighbors): a sequence-space "entropy" effect.

Because of this "entropic" driving force, SUBF can be nonadaptive.



Sikosek, Chan & Bornberg-Bauer (2012) PNAS 109:14888-14893.



Examples of 3D lattice protein models of evolution

• Maximally compact 27mers configured on the simple cubic lattice and restricted to a 3x3x3 cube, with a 20-letter alphabet



• Useful and versatile construct, but non-maximally compact conformations – including open unfolded conformations – are *not* considered for the sake of computational tractability. This limitation has a significant impact on the physicality of the energetics of these 3D lattice protein models as for the 2D lattice protein models.

Number of 27mer conformations in a 3×3×3 cube = 103,346

Total number of 27mer unrestricted conformations on the simple cubic lattice = 11,447,808,041,780,409

Reference: Sikosek & Chan (2014) *J R Soc Interface* **11**:20140419

Possible biophysical basis for the anti-correlation between protein expression level and evolutionary rate





Another probable biophysical constraint underpinning the anti-correlation between expression level and evolutionary rate is the need for a protein to avoid *misinteraction* with other proteins



simulation of evolution. • C_i is the concentration of free molecules of protein i, C_{ii} is the concentration of the protein complex composed of a protein *i* and a protein *j*, *D_i* is the total concentration of protein *i* in the cell.

hypothesis does not fully explain the E-R (Expression-Rate) anti-correlation, especially for residues on protein surface. They propose that natural selection against protein-protein misinteraction, which wastes functional molecules and is potentially toxic, constrains the evolution of surface residues. Because highly expressed proteins are under stronger pressures to avoid misinteraction, surface residues are expected to show an E-R anticorrelation. Their findings indicate a pluralistic origin of the E-R anti-correlation.



Results and figures from: *Yang, Liao, Zhuang & Zhang (2012) Protein misinteraction avoidance causes highly expressed proteins to evolve slowly. PNAS 109:E831-E840.

Summary

- Biophysical sequence-to-structure mappings based upon simple explicit-chain protein models are versatile conceptual tools for addressing general principles of evolution.
- The superfunnel paradigm underscores a fundamental positive correlation between thermodynamic stability of the native structure and the mutational stability of a protein.
- Excited-state selection of promiscuous function can speed up evolution considerably.
- Subfunctionalization after duplication of a bi-stable gene with dual functions can be driven by sequence-space topology (i.e., mutational robustness) in an essentially nonadaptive manner.