# Molecular Evolution

Andrew Torda summer semester 2012, Struktur & Simulation

Why ?

- applications not possible with detailed models

Ingredients in this set of lectures

- models for proteins
  - simple representation - lattices, simple energy functions
- Boltzmann relation and partition function
  - ability to calculate probability of conformations
- definitions from earlier lectures – foldability / stability

Aim

- from very few assumptions
- simulation which reproduces physical properties

# Why lattice models ?

- Earlier – building models
  - how much detail - rather arbitrary
- theme of these lectures
  - model to mimic physics
    - simulate the behaviour of a system
  - minimal model
    - what observables can I reproduce with which features in the model ?

- here – minimal models
  - one does not need serious chemistry to reproduce protein properties
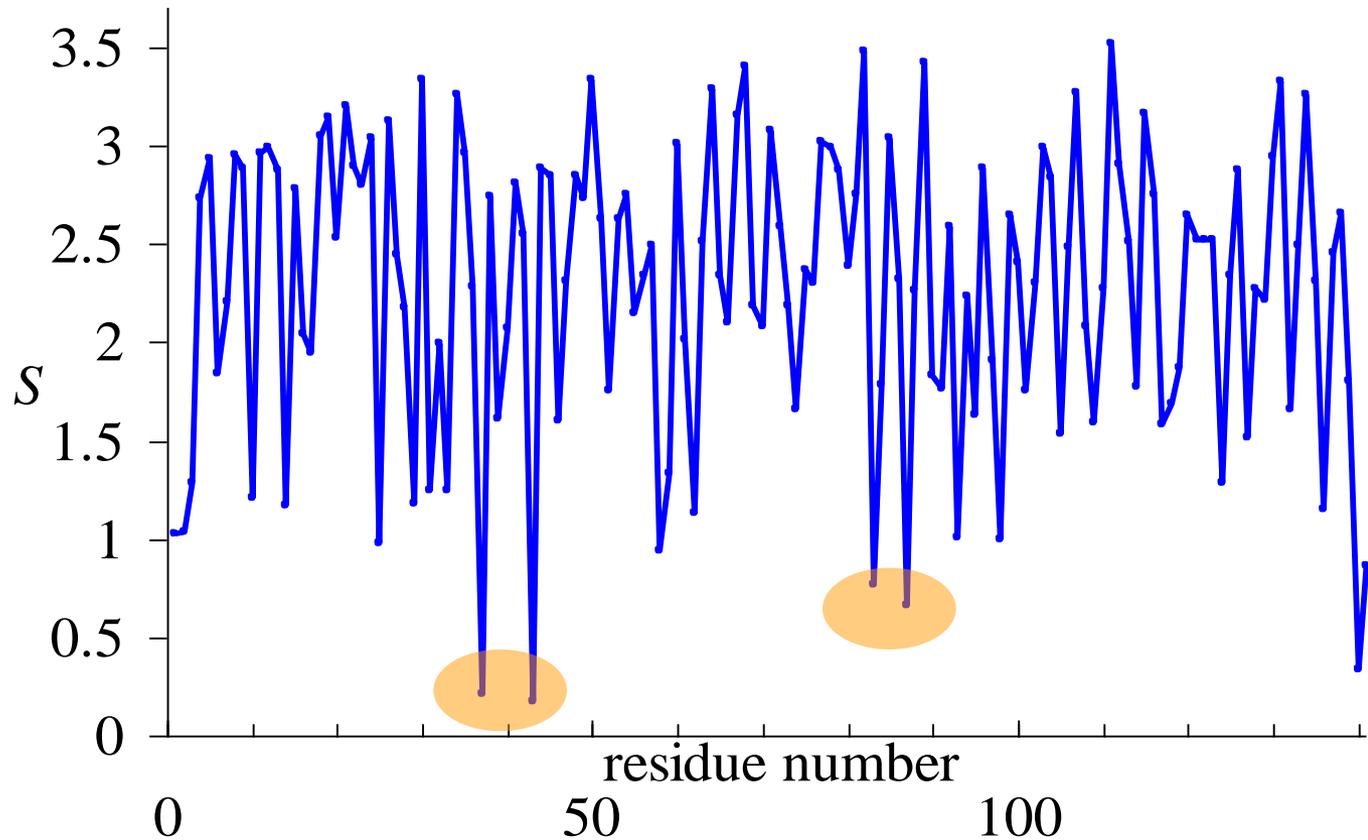  - evolutionary pressure may not be real

# Plan

- generalities
  - sources of evolutionary pressure
- example of unexpected evolutionary pressures (Darwinian)
- neutral networks
- another molecular consequence (not so Darwinian)


- calculations which can only be done on idealised systems


- two example papers – in stine

# Evolution observables

- In the real world, not much
  - phenotypes
    - blue eyes, brown eyes (macroscopic)
    - different proteins (molecular)
  - genotypes (with more effort)
  - population properties
- consequence ?
  - mostly look at evolution in terms of pressure on phenotypes
  - classic adaptive Darwinism

- first:
  - a property to be explained later

# Haemoglobin conservation
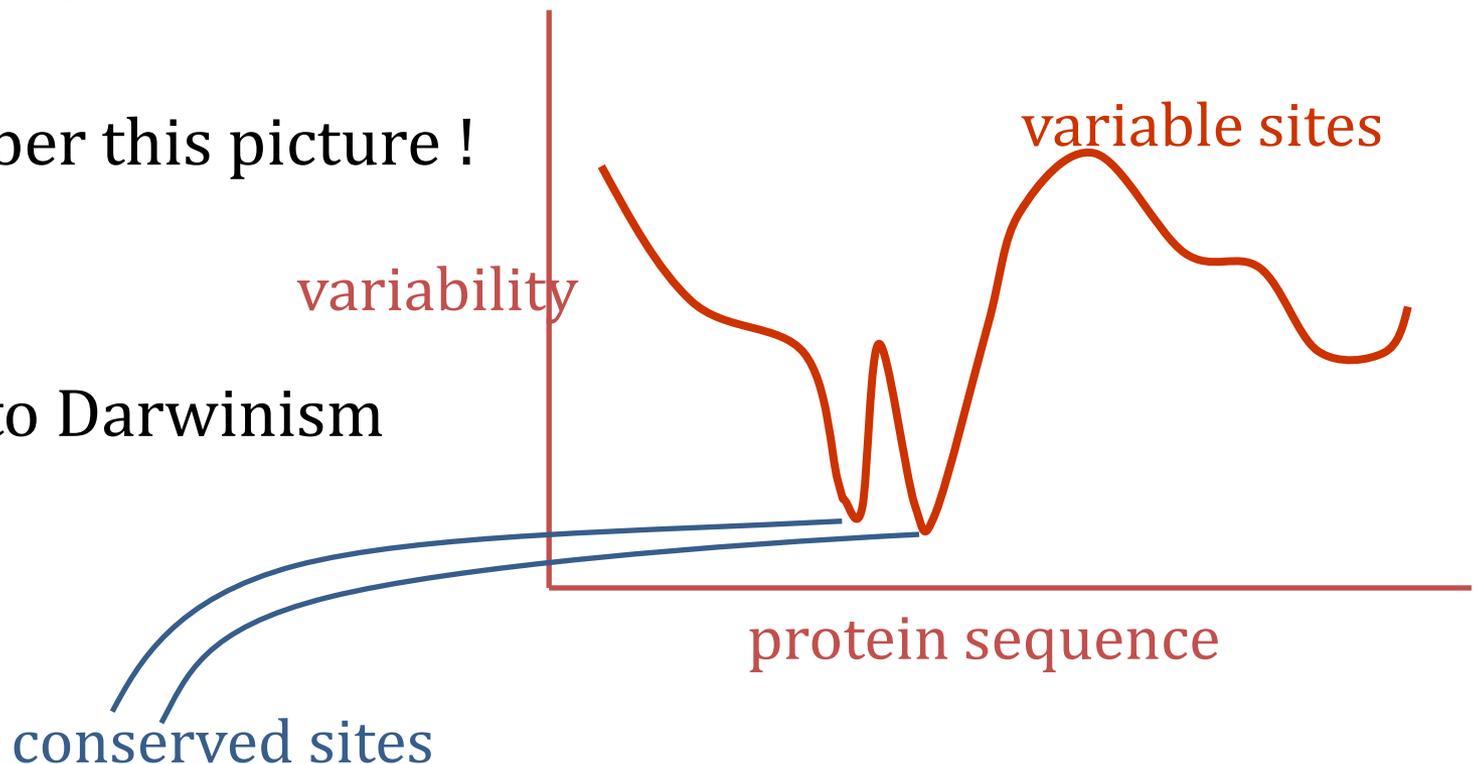
- look at residues 37, 43, 83 and 87



- 4 residues (maybe more) stand out as conserved

# Sequence variability

- take family of related sequences
  - see how conserved / variable they are
- variable sites
  - are they unimportant ?

- remember this picture !

- return to Darwinism

variability

variable sites

protein sequence

conserved sites

# Adaptive Darwinism

- I see a fish which lives behind a rock and eats seaweed
- A mouse is just the right size to squeeze through the hole in my wall
- Voltaire (1694-1778)

```
Master Pangloss taught ..
dass in dieser besten aller möglichen Welten, ...
    „Es ist erwiesen" sagte er, „dass die Dinge nicht
    anders sein können: denn da Alles zu einem Zweck
    geschaffen worden, ist Alles notwendigerweise zum
    denkbar besten Zweck in der Welt. Bemerken Sie wohl,
    dass die Nasen geschaffen wurden, um den Brillen als
    Unterlage zu dienen, und so tragen wir denn auch
    Brillen"
```

- Two aspects
  - adaptation to glasses (evolution is directed)
  - best of all possible worlds (we / the world are optimised)

# Classic Darwinism – molecular level

Obvious pressures

- function
- stability

Less obvious, but simple

- folding

# Stability (first version)

- I must be stable at room temperature
- proteins in us must evolve to be stable under different conditions (organelles)
- extreme examples – bacteria
  - thermophiles, acidophiles, halophiles, …

- proteins are not really very stable

# Function

More difficult to explain

- how does a sugar enzyme change to a muscle protein ?
- almost must have redundant copies of function
  - if one is broken, you do not die

Consequence

- we are not "optimal"

Experiment

- make "knockout" animals to look at function
- results are often not clear
  - prion proteins (verrückte Kuh Krankheit / Mäuse)

# Folding

Subtle phenotype

- we cannot look at a population and see it
- we can simulate it

intuitively plausible



stable but will not fold

native

configurations

native

configurations

# less obvious

- composition of proteins
  - trp costs far more energy to make than gly or ala
  - this is an observable phenotype
- DNA base pair composition
- …

# Other evolutionary pressures

- is it good to be resistant to mutation ?
  - what if a gamma ray hits me and my children die ?
- more formally
  - a sequence (protein) is more likely to propagate if
    - it can be changed
    - it keeps functioning
- can this be modelled ?

Plan :
- be Darwinian
- (later) show why it is probabilistic (not Darwinian)

# Simulating mutation resistance

Lattice simulations

- 25 residues, 2 dimensional, compact, 5×5 lattice
- 20 residue types
- 1081 conformations
- remember we can calculate $Z$ and stability
- for any sequence can say
  - will this sequence fold or not ? $\Delta G_{fold}$
    - how different is lowest energy to other energies
- too big to check all sequences

Example calculation

- look at differences with and without evolution

*Taverna, DM and Goldstein RA, J. Mol. Biol. 315, 479-484 (2002)
Why are proteins so robust to site mutations ?

# Example evolution calculation

Evolution simulation

- apply mutations infrequently / randomly
- sequence must maintain
  - same structure
  - foldability
- for each member of population
  - check lowest energy configuration
    - if it has changed – sequence dies
  - check $\Delta G_{fold}$
    - if sequence is not foldable – dies
  - of remaining sequences, randomly pick for reproduction

# Comparing populations

Take a sequence which folds

- copy 3 000 times – initial population

| | |
|---|---|
| initial population | |
| ↓ 30 000 generations | forget (equilibration) |
| diverse population | |
| ↓ 30 000 generations | keep and sample |

generate sequences randomly

evolved sequences

random sequences

# Properties to look at

- How often does a mutation make a protein more stable ?
- How often does
  - a stable protein become more stable ? (not often)
  - an unstable protein become more stable ? (must be higher)
- Do the fractions differ between
  - random sequences (right hand side previous Folien)
  - evolved sequences (left hand side)
- For some protein we know $\Delta G$
- From simulation look at proteins with some $\Delta G$
  - after mutation get new $\Delta G$
  - look at large number of mutations, get probability $P(\Delta \Delta G > 0)$ of becoming even less stable

# What do you expect ?

- Evolved sequences must be more stable than random ones
- Will they also be more resistant to mutations ?
  - if they were not, they would die

# Simulation results

- Take a sequence and have a look
  - when it mutated and survived
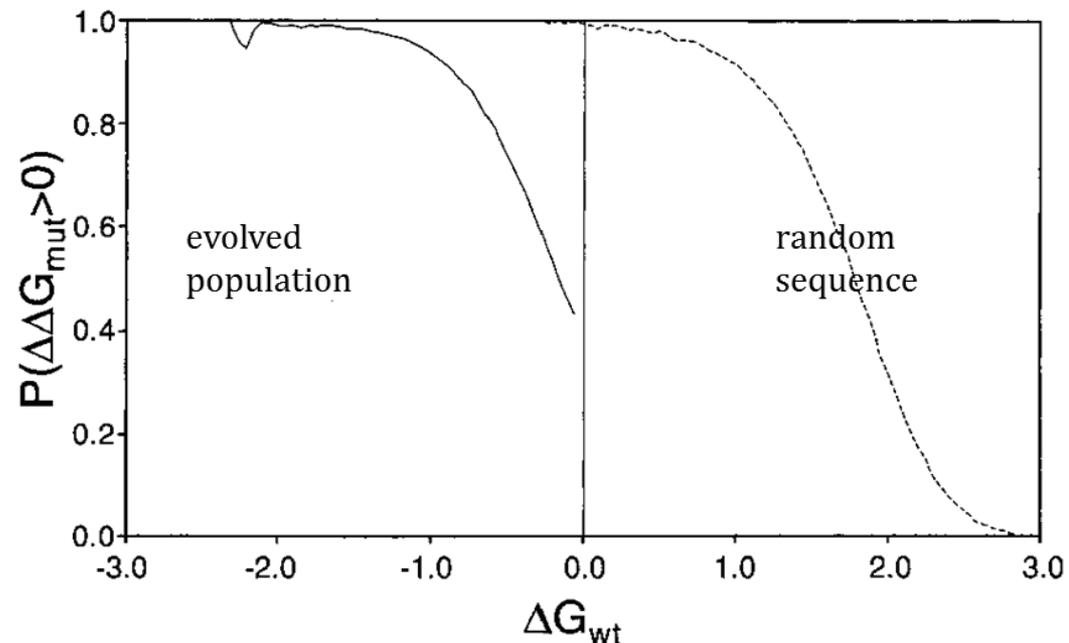    - how often did it become less stable $P(\Delta \Delta G > 0)$ ?



evolved population

random sequence

very stable proteins $\leftrightarrow$ less stable

stability before mutation

# Interpreting results

random sequence

- unstable $\Delta G > 0$
  - not easy to make more stable
- stable ? $\Delta G < 0$
  - all mutations make it worse

evolved sequence

- very stable ?
  - cannot make better
- marginally stable ?
  - mutations often OK

# Results explanation

Without explicitly adding idea
- evolution makes
  - more stable proteins
  - proteins which survive mutations



Agree with experiment ?
- small amount of the time
  - mutations have no effect
  - make protein more stable than natural version

# Sequence variability interpretation

- Typical part of sequence analysis
- look at collection of related sequences and see how conserved they are (conservation, profiles, sequence entropy, ..)

- Why are some sites so well conserved ?
  - function ?
- Why do some sites vary ?
  - old view: they do not matter
  - this paper
    - this is a consequence of evolution
  - if they are important and fragile, you die

variable sites

variability

protein sequence

conserved sites

# Subtle evolutionary pressure ?

Is this an evolutionary pressure ?
- seems like a good idea to not die when mutated
- authors argue that the reason is different
- neutral evolution ...

## so far

- very simple lattice model reproduces
- stability, evolutionary pressures

- now .. background for neutral evolution

# Neutral evolution

Classical view (selective adaptation) explains life
- we are always trying to adapt to each other, environment ...
- there is some diversity when there is no cost (blue / brown eyes)

Alternative
- most mutations have no effect (neutral)
- if they far outnumber the selected mutations, they will dominate

Macroscopic
- brown eyes versus blue – not so surprising
- microscopic / molecular ?

Neutral evolution
- consequences ?
- predictions ?
- predictions at molecular level / simulations

# Background of neutral evolution

At molecular level

- DNA level (obvious)
  - 64 codons / 20 amino acids / much redundancy
    - CUG / CUC both ile (+ many more)
  - lots of mutations have no (not much) effect
- Protein
  - bit less clear
  - we can change amino acids and
    - preserve structure
    - often function
- Net effect
  - we can make many mutations
  - some do not affect the protein
  - some protein effects are very small

# Simulating at the molecular level

Basic idea

- take a population  (maybe $10^3$ or as big as possible)
    - make random changes
    - look at consequences
    - kill or reproduce molecules

Most popular

- RNA
    - for a given mutation, can guess at secondary structure
- Proteins
    - lots of lattice calculations

# Simulation machinery



HP model in two dimensions
- length 18
  - one can look at all sequences
  - all conformations
  - ... for any sequence
    - can find minimum energy structure
  - for any structure
    - we can find all sequences which have this as minimum energy

Bornberg-Bauer, E (1997) Biophys J. 73, 2393-2403
How are model protein structures distributed in sequence space ?

# Calculations

Find popular structures
- which is best for many sequences
- collect these sequences
  - neutral set

Neutral mutations
- which of these sequences are connected by a point mutation?
- example
  - `HPHP`**`H`**`HH..` and `HPHP`**`P`**`HH..` have same ground state
  - they are connected by one change
  - this change does not cost anything in evolution
    - it is "neutral"
  - in pictures...

# Neutral mutations

- look at sites which can be changed
  - many possible sequences
- can one mutate each to every other ?
  - `HPHP`**`HH`**`H..` and `HPHP`**`PP`**`H` are not connected
- what can we say about the connected sequences ?
  - form connected sets

✔ sites where neutral mutations were found

- `HPHP`**`HH`**`H` and `HPHP`**`PP`**`H` may be a set, but not connected

Bornberg-Bauer, E (1997) Biophys J. 73, 2393-2403
How are model protein structures distributed in sequence space ?

# Connected and non-connected sets

- each dot is one protein sequence/structure



neutral set with two connected sets

neutral set and connected set

# Neutral networks

- Sequences which can turn into each other are "neutral network"
- How big are the neutral sets ?
  - about ¼ have more than 5 sequences
  - most popular has 48 sequences
  - lots of very rare structures
- Are these sets fully connected ?
(can anyone eventually mutate into anyone else) ?
  - about 80 % of time

# Evolutionary consequences

- a population can quickly spread over a huge number of accessible sequences
- immense variation at molecular level is possible

- Can one hop between different connected networks ?
  - in this model – not so easily ( ≥ 2 mutations)

More interesting consequences
- some structures are hard to find by random moves
- some are very popular
- what does this say about mutation study ?

# Mutation resistance revisited

Earlier slides

- it seems as if proteins evolve in order to be resistant to mutations (sounds Darwinian)
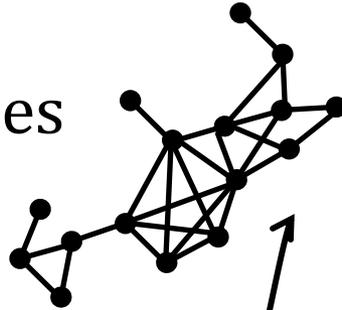
- Alternative
    - think of sequence space
    - a group of related sequences are a cluster in this space

residue 3

residue 2

A   C   D   …   W   Y

residue 1

# Networks, probabilities, mutation resistance

huge network
    1000's sequences

small network

- mutate to here
  - seems mutation resistant
    - lots of possibilities to mutate and maintain structure
  - more likely to be found (more sequences)
- mutate here ? likely to die

# Darwinian versus neutral evolution

- Crux of these lectures
- Darwinian evolution – what you see is
  - most fit (selection pressure)
- Neutral evolution – what you see is
  - whatever is most likely to occur

- Relevance to mutation resistance
  - Darwinian
    - useful trait that will be selected for
  - Neutral
    - larger neutral networks
      - by definition – mutation tolerant
      - because they are larger, more likely to be found
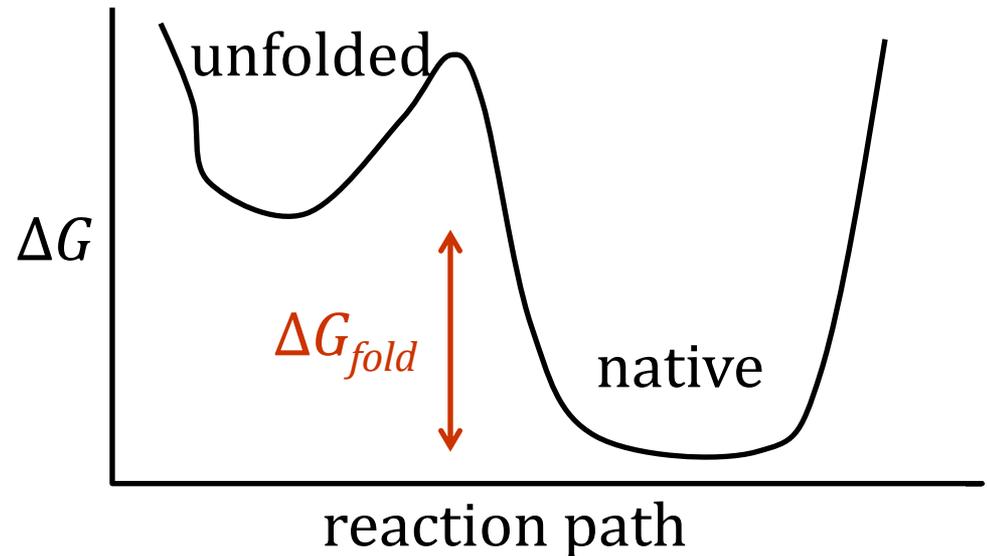
# Simple models and lattice models

- details (numbers) are not so vital
- problem is made tractable by use of simple model

# Protein stability

more work from same group[*]

Most proteins are NOT very stable ($5 - 10$ kcal mol$^{-1}$)

- claims:
  - less stable, more flexible
  - easier to have chemical function



*Taverna, DM, Goldstein, RA, 2002, Proteins, 46, 105-109, Why are proteins marginally stable ?

# Another model calculation

- 5×5 lattice 1081 conformations
- 20 amino acid types
- cannot visit all sequences, can visit all structures
- use a definition of foldable

$$\Delta G_{folding} = E_f + kTln\left(Z - \exp\left(\frac{-E_f}{kT}\right)\right)$$

3 simulations
1. long walk of one sequence
2. population
3. random sequences

# Sidetrack for arguments

Goldstein's formula

- $p_f$ probability of folded state $\quad p_f = \dfrac{\exp\left(\frac{-E_f}{kT}\right)}{Z}$

- $p_u$ probability of unfolded state
  - probability all states (1)−minus probability of folded

$$p_u = \dfrac{\sum_i \exp\left(\frac{-E_i}{kT}\right) - \exp\left(\frac{-E_f}{kT}\right)}{Z}$$

$$\dfrac{p_f}{p_u} = \dfrac{\exp\left(\frac{-E_f}{kT}\right)}{\sum_i \exp\left(\frac{-E_i}{kT}\right) - \exp\left(\frac{-E_f}{kT}\right)}$$

$$= \dfrac{\exp\left(\frac{-E_f}{kT}\right)}{Z - \exp\left(\frac{-E_f}{kT}\right)}$$

# Getting free energy expression

$$\Delta G = -kT \ln\left(\frac{p_f}{p_u}\right)$$

$$= kT \ln\left(\frac{\exp\left(\frac{-E_f}{kT}\right)}{Z - \exp\left(\frac{-E_f}{kT}\right)}\right)$$

$$= -kT \ln \exp\left(\frac{-E_f}{kT}\right) - kT \ln\left(Z - \exp\left(\frac{-E_f}{kT}\right)\right)$$

$$= E_f + kT \ln\left(Z - \exp\left(\frac{-E_f}{kT}\right)\right)$$

# Simulation (long walk)

- Take viable sequence
- mutate
  - if (foldable)
    - keep
  - else
    - retain old sequence

# Simulation (population)

- Take 3 000 identical sequences
- mutate

- calculate $\Delta G_{folding}$ for all members
- kill (remove) non-folders
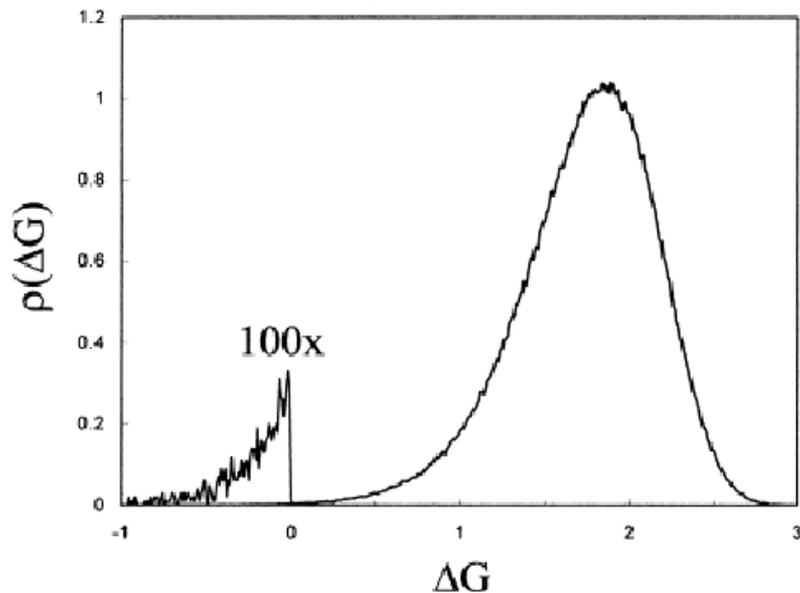- copy random survivors to keep population at 3 000

# Stability of results

What is the result
- from random sequences ? (left)
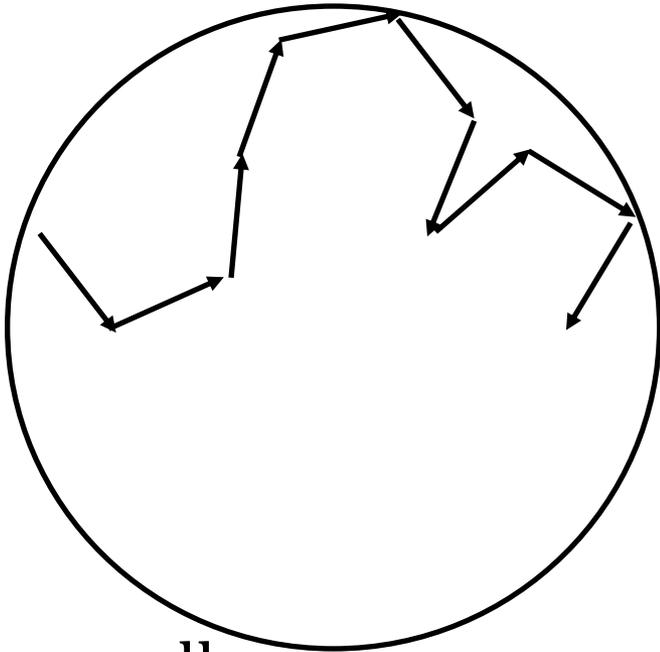- from a long walk (right A)
- from a population (right B)

Sequences become more stable
- but barely so

# Where does the population result come from ?

Proteins die if they are unstable

- the population moves to folding sequences (this is selected)

- there is no force to make them more stable


- high dimensional object arguments / population phenomena
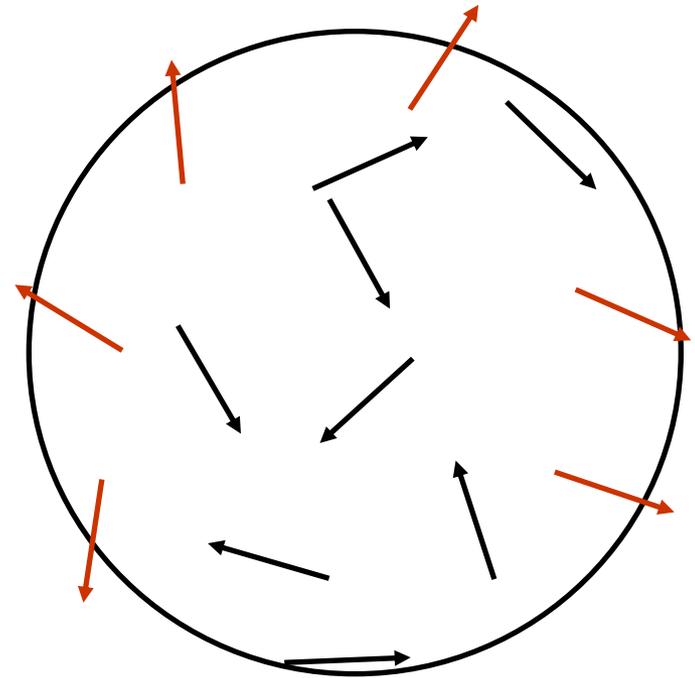
  - explain the population result

# Walk versus Population

- high dimensional objects
  - high proportion near to surface

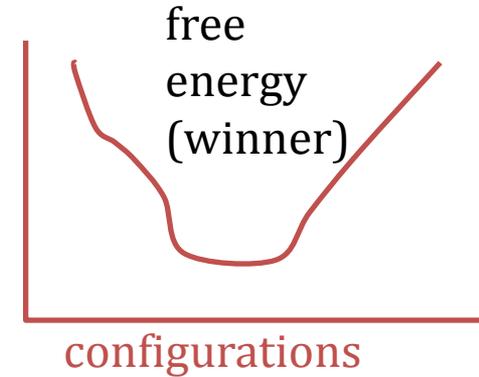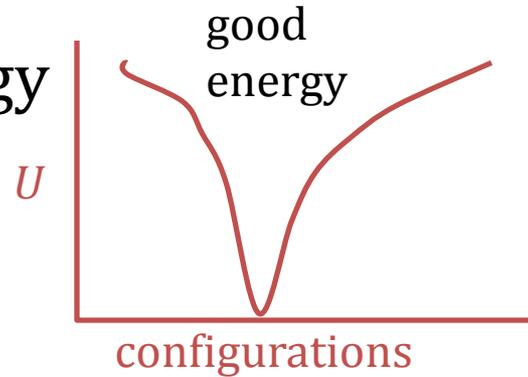long walk

- sequences bounce around
  near surface

population

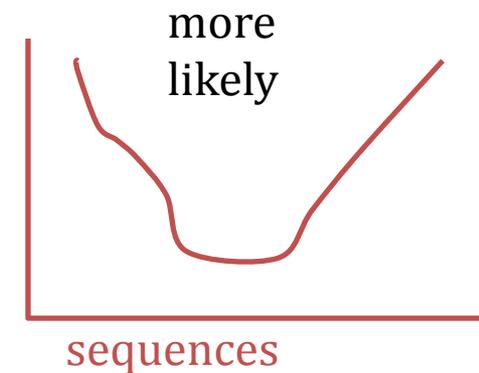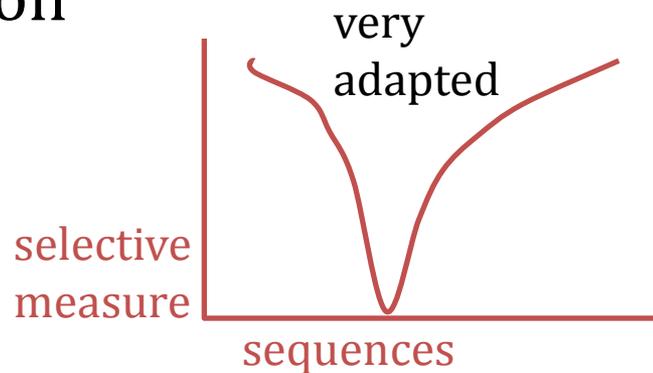- sequences near surface removed, others reproduce

- Population acts as if there is a sink removing most unstable proteins

- Results give marginally stable proteins
  - no mention of function
  - arguments purely statistical

# Analogy evolution and free energy

energy /free energy minima

good energy

$U$

configurations

free energy (winner)

configurations

evolutionary version

very adapted

selective measure

sequences

more likely

sequences

- evolution is adaptive, but subject to statistical effects
- statistical effects may look like evolutionary pressures (mutation resistance, stability)

# Summary

- Neutral evolution began in the late 1960's
  - nicest evidence from simple simulations

- Molecular models can be applied in unexpected places

- We interpret the world in terms of observables (numbers, colours, stability, ...)
  - this may be over-interpretation

- First lattice lectures
  - one can do Monte Carlo simulations
- Now
  - there are other types of simulation