

The 12th KIAS protein folding winter school High1 resort, Jan 21-25, 2013 Eighteenth-century science, following the Newtonian revolution, has been characterized as developing the sciences of organized simplicity, 19th century science, via statistical mechanics, as focusing on disorganized complexity, and 20th- and 21st-century science as confronting organized complexity.

Nowhere is this confrontation so stark as in biology. Nowhere are new conceptual tools so deeply needed.

"Origin of Order" Stuart Kaufman





























The goal of this lecture: get a global view of Biology



Outline

- Biology by numbers
- Genome
- Gene expression, genetic switch
- Foldings and functions of RNA

I. Biology by numbers

Physical Biology of the Cell by R. Philips et al.

http://bionumbers.hms.harvard.edu



Figure 1-21 Molecular Biology of the Cell 5/e (© Garland Science 2008)



bacteria (prokaryote)

eukaryote

 $V_{E.Coli} \approx 1 \mu m^3 = 1 fL$ $m_{E.Coli} \approx 1 pg$ $\rho_{E.Coli} \approx 1g / mL \approx \rho_{water}$

(A) pili
pili
(B) flagella
(C)
$$(1 + 1)$$

 $(2 + 1)$
 $(2 + 1)$
 $(2 + 1)$
 $(3 + 1)$
 $(4 + 1)$
 $(5 + 1)$
 $(6 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7 + 1)$
 $(7$

Figure 2.1 Physical Biology of the Cell (© Garland Science 2009)

 $1 \, molecule \, / \, 1fL = \frac{1mol}{6 \times 10^{23} \times 10^{-15} \, L} \approx 2nM$

Water is 70 % of the cell mass $(m_{E.Coli} = 1pg)$

Dry mass of the cell (30% of 1pg) = 0.3 pg

Half of the dry mass (=0.15 pg) = proteins

1 Da = mass of a hydrogen atom = $\frac{1}{6 \times 10^{23}}g = 1.6 \times 10^{-24}g$

1 amino acid =
$$100$$
 Da

average protein size (E. coli) = 300 a.a. \rightarrow 30,000 Da

 There are 2×10^6 proteins in cytoplasm

$$c_{prot} = 2 \times 10^{6} / 1 \mu m^{3}$$

$$d_{prot-prot} = c^{-1/3} = \left(\frac{\left(10^{3} nm\right)^{3}}{2 \times 10^{6}}\right)^{1/3} = (500 nm)^{1/3} \approx 8 nm$$





RIBOSOME

20% of the protein complement of a cell = ribosomal proteins

Ribosome (70S) = Large subunit + small subunit Large subunit (50S) = 23S r-RNA + r-proteins Small subunit (30S) = 16S r-RNA + r-proteins

S : Svedberg constant (sedimentation constant)A heavier particle sediments faster in the centrifugation, thus have a larger S value.

ribosome (2.5MDa)
$$\begin{cases} r - RNA : 2/3 \text{ of the mass} \\ r - protein : 1/3 \text{ of the mass} \end{cases}$$



$$m_{r-protein} = 830,000 Da$$

$$M_{r-protein} = 20\% \times M_{protein}$$

$$N_{ribosome} = \frac{M_{r-prot}}{m_{r-prot}} = \frac{0.2 \times 0.15 \, pg}{830,000 \, Da} = 20,000$$

$$\approx 19,000$$

yeast cell : model system to study a single eukaryote cell



Important model organisms in Biology

- Bacteriophage: Gene structure, gene regulation
- *E. coli* : Genetic network (lac operon), flagella
- Yeast (*S. cerevisiae*) :
 - Cell cycle, genetics, cell biology
 - quick and easy to grow
- Flies (*Drosophila melanogaster*)
 - Easy to grow, rapid generations
 - Easy mutation induction, many observable mutation
 - Giant chromosome in salivary glands that can be examined under a light microscope
 - Molecular genetics, Population genetics, Developmental biology
- C. elegans
 - Excellent model of the genetic control of development and physiology
 - The first multicellular organism whose genome was completely sequenced.
- Mice
 - Genome size and organization is very similar to human's
 - Small, easy to breed, ideal for lab study.
 - Transgenic, knock-out mouse.
- Human
 - cf) HeLa (Henrietta Lacks) cell
 - A cell type in an immortal cell line used in scientific research
 - Oldest and most commonly used human cell line.
 - Incredibly hardy, which makes them useful for medical research
 - The line derived from cervical cancer cells of a patient taken on 1951
 - No aging, no programmed cell death (apoptosis)
- and others (Giant squid)

Bio"polymers"

- Cells are made of biopolymers (DNA, RNA, proteins, cytoskeletal filaments,)
- Length : $L = N \times a$
- Size : $R_g \sim N^{\nu}$
- Flexibility (persistence length)

$$\left\langle \vec{u}(s) \cdot \vec{u}(s') \right\rangle = e^{-|s-s'|/l_p}$$

$$\left\langle R^2 \right\rangle = \int_0^L ds \int_0^L ds' e^{-|s-s'|/l_p} = 2Ll_p \left[1 - \frac{l_p}{L} \left(1 - e^{-L/l_p} \right) \right]$$

Length (contour length)

 $L = N \times a$

Proteins

cf. Largest protein size : titin 33423 aa









Size





18

Estimating the rate of biological processes

800 bp/sec

Replication $\frac{dN_{bp}}{dt} \approx \frac{N_{bp}}{\tau_{cell}} \approx \frac{5 \times 10^6 bp}{3000 \text{ sec}} \approx 2000 bp / \text{sec}$ (or 1000bp/sec per DNA replication complex) Biochemical study found 250-1000bp/sec range Let's take it 800 bp/sec.

Transcription

40nt / sec/ RNAP

40 nt/sec

typical length of transcript : $3nt / a.a \times 300a.a \approx 1,000nt$

The time to produce a transcript using an RNAP ~ 25 sec

Protein synthesis
$$\frac{dN_{protein}}{dt} \approx \frac{N_{protein}}{\tau_{cell}} \approx \frac{3 \times 10^{6} \text{ proteins}}{3000 \text{ sec}} \approx 1000 \text{ proteins / sec}$$

15 aa/sec
 $\rightarrow \frac{(3 \times 10^{6}) \times 300 a.a.}{3000 \text{ sec}} \approx 3 \times 10^{5} a.a / \text{sec}$
protein synthesis rate per ribosome
= 3,000,000 proteins/3000 sec /20,000 ribosomes
= 0.05 proteins/sec/ribosome $\rightarrow 15$ aa/sec/ribosome

Energy scale relevant for biomolecules

Biomolecules are "marginally" stable; thus their dynamics are susceptible to thermal fluctuations, external noises.

 $1k_BT = 4.14 \text{ pN} \times \text{nm} = 4.14 \times 10^{-21} \text{ J} = 0.59 \text{ kcal/mol at 300 K}$ $1eV \approx 38.7k_BT \sim 40k_BT$

The characteristics of "marginal stability" allow biomolecules to adapt their conformational states into, if any, an alternative form.

However, for some molecules or as the size of molecules grows the thermal fluctuation itself is not sufficient to trigger the conformational transitions within a biologically relevant time scale. In this case, various cofactors, NTPs, or interaction with other molecular machines are used as free energy sources (active process).

Binding Free Energy

 $AB \rightleftharpoons A + B$

Thermal energy scale (at 300K) $k_B T = 4.14 \times 10^{-21} J = 4.14 \, pN \cdot nm$ $\doteq 0.59 \, kcal \, / \, mol$

$$\Delta G = k_B T \log K_d = k_B T \log \frac{[A][B]}{[AB]} \approx k_B T \log 10^{-9} = -2.3 \times 9k_B T \approx -20k_B T$$

Typical protein-protein dissociation free energy $\approx 20k_BT(K_d = 1nM)[1molecule / cell] \\\approx 14k_BT(K_d = 1\mu M)[10^3 molecules / cell] \\\approx 7k_BT(K_d = 1mM)[10^6 molecules / cell]$

Insulin dimer : 7 kcal/mol $(K_d = 10^{-5} M)$ Trypsin-PTI : 18 kcal/mol $(K_d = 10^{-13} M)$ Haemoglobin $\alpha\beta$ dimer : > 11 kcal/mol $(K_d << 10^{-8} M)$

Thermodynamics of protein folding





Kinetics of biopolymers



cf. Another Arrhenius-like expression, Eyring equation (or transition state theory (TST)) derived by theoretical chemists H. Eyring and M. Polanyi

$$k_{TST} = \kappa \frac{k_B T}{h} e^{-E^{\dagger}/k_B T}$$

is often used in estimating the height of kinetic barrier. However, it is of note that this equation should only be applied for chemical reactions associated with bond breaking (or formation) dynamics of simple organic (or inorganic) compounds in gas phase. The prefactor $(k_BT/h)^{-1} \approx 0.2$ psec is the vibrational frequency of a critical molecular bond at the transition state to be ruptured by chemical dynamics. Blind use of this prefactor overestimate kinetic barriers. To describe dynamics in condensed media, including biopolymer folding and biomolecular assembly, it is more appropriate to use Kramers' equation, which correctly predicts the effect of solvent viscosity to the reaction rate.

$$k_{KR} = \frac{\omega_0 \omega_b}{2\pi\gamma} e^{-\Delta G^{\dagger}/k_B T}$$

As estimated from the fits to folding time of RNA and proteins, the prefactor of Kramers' rate for RNA and proteins are \sim (0.1-1) µsec. These values indeed agree with the recent experimental measurements of **"speed limit" of protein folding**, which were performed using barrierless downhill folders (Gruebele), and computational and experimental reports of **transition path time** by D. E. Shaw and W.A. Eaton and coworkers, respectively.

Transition path time $\tau_0 \sim 1 \mu sec$ is size-independent. Fast and slow folding proteins with vastly different topologies have almost the same value.

How Fast-Folding Proteins Fold

Kresten Lindorff-Larsen,¹*† Stefano Piana,¹*† Ron O. Dror,¹ David E. Shaw^{1,2}† SCIENCE VOL 334 28 OCTOBER 2011



517

Single-Molecule Fluorescence Experiments Determine Protein Folding Transition Path Times

Hoi Sung Chung,* Kevin McHale, John M. Louis, William A. Eaton*

SCIENCE VOL 335 24 FEBRUARY 2012





2. Genome

A. Chromosome organization inside nucleus B. Composition of Genome

A. Chromosome organization inside nucleus $L=(3\times10^{9} \text{ bp})\times(0.34\text{ nm/bp})\approx1\text{ m}$ \rightarrow nucleus of Iµm size



 $U_{\text{int}}(\{\vec{r}_i\}) = \sum_{i=1}^{N} \frac{k_r}{2} (r_{i,i+1} - a)^2$

Gaussian chain (random flight chain) Freely jointed chain $(k_r \gg 1)$

$$U_{\text{int}}(\{\vec{r}_i\}) = \sum_{i=1}^{N} \frac{k_r}{2} (r_{i,i+1} - a)^2 + \sum_{i=1}^{N-1} \frac{k_{\theta}}{2} (\theta_i - \theta_0)^2$$
Semiflexible chain,
$$\left(\doteq \sum_{i=1}^{N} \frac{k_r}{2} (r_{i,i+1} - a)^2 - \sum_{i=1}^{N-2} \frac{k_{\theta}}{2} \hat{r}_{i,i+1} \cdot \hat{r}_{i+1,i+2} \right)$$
Worm-like chain
$$(k_r \gg 1, \theta_0 = \pi)$$

$$U_{\text{int}}(\{\vec{r}_i\}) = \sum_{i=1}^{N} \frac{k_r}{2} (r_{i,i+1} - a)^2 + \sum_{i=1}^{N-1} \sum_{j=i+1}^{N} \varepsilon \left(\frac{\sigma}{r_{i,j}}\right)^{12}$$

Self-avoiding chain

$$\left\langle R^2 \right\rangle = \left\langle \sum_{i}^{N} r_i \cdot \sum_{j}^{N} r_j \right\rangle = \sum_{i=1}^{N} \left\langle r_i^2 \right\rangle + \sum_{i \neq j}^{N} \left\langle r_i \cdot r_j \right\rangle = Na^2$$

$$\therefore R \sim N^{1/2}$$

$$\left\langle R^2_{R} \right\rangle = \frac{1}{N} \sum_{i=1}^{N} (\vec{R}_i - \vec{R}_c)^2 = \frac{1}{N} \sum_{i=1}^{N} (\vec{R}_i - \frac{1}{N} \sum_{i=1}^{N} \vec{R}_i)^2$$

$$= \frac{1}{2N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} (\vec{R}_i - \vec{R}_j)^2 = \frac{1}{N^2} \sum_{i=1}^{N-1} \sum_{j=i+1}^{N} |i-j|^2 a^2 = \frac{Na^2}{6}$$

$$= \left\langle R^2 \right\rangle / 6$$

Size scaling of SAW: de Gennes' argument

$$F = E - TS = \frac{T}{2}c^2R^d + T\frac{3R^2}{2\langle R^2 \rangle} \approx \frac{T}{2}\left(\frac{N^2}{R^d}\right) + \frac{3T}{2}\frac{R^2}{Na^2}$$

$$\frac{\partial F}{\partial R} = 0; \quad N^2 R^{-d-1} \sim N^{-1} R$$

$$\therefore R \sim N^{\nu} \quad \text{where} \quad \nu = \frac{3}{d+2} = \begin{cases} 1, \quad d=1\\ 3/4, \quad d=2\\ 3/5, \quad d=3 \end{cases}$$

cf. polymer melts

$$\left\langle R_g^2 \right\rangle \sim N$$

Gaussian statistics

Flexible polymer at θ -condition

$$U_{\text{int}}(\{\vec{r}_i\}) = \sum_{i=1}^{N} \frac{k_r}{2} (r_{i,i+1} - a)^2$$

$\Delta F = \Delta E - T \Delta S$





Flexible polymer in good solvent (SAW)

$$U_{\text{int}}(\{\vec{r}_i\}) = \sum_{i=1}^{N} \frac{k_r}{2} (r_{i,i+1} - a)^2 + \varepsilon_l \sum_{i=1}^{N-1} \sum_{j=i+1}^{N} \left(\frac{\sigma}{r_{i,j}}\right)^{12}$$



Collapse of flexible polymer in poor solvent

$$U_{\text{int}}(\{\vec{r}_i\}) = \sum_{i=1}^{N} \frac{k_r}{2} (r_{i,i+1} - a)^2 + \varepsilon_h \sum_{i=1}^{N-1} \sum_{j=i+1}^{N} \left[\left(\frac{\sigma}{r_{i,j}}\right)^{12} - 2\left(\frac{\sigma}{r_{i,j}}\right)^6 \right]$$



DNA compaction



Euchromatin : Loosely packed, Transcriptionally active Heterochromatin : Densely packed, Transcriptionally inactive




Visualization of chromosome territories using FISH (fluorescence in situ hybridization)





Restriction enzymes

- Type I enzymes cleave at sites remote from recognition site; require both ATP and S-adenosyl-L-methionine to function; multifunctional protein with both restriction and methylase activities.
- Type II enzymes cleave within or at short specific distances from recognition site; most require magnesium; single function (restriction) enzymes independent of methylase.
- Type III enzymes cleave at sites a short distance from recognition site; require ATP (but do not hydrolyse it); Sadenosyl-L-methionine stimulates reaction but is not required; exist as part of a complex with a modification methylase.
- Type IV enzymes target modified DNA, e.g. methylated, hydroxymethylated and glucosyl-hydroxymethylated DNA





Hi-C method

Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome Erez Lieberman-Aiden *et al. Science* **326**, 289 (2009); DOI: 10.1126/science.1181369



(IMb resolution)

Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome Erez Lieberman-Aiden *et al. Science* **326**, 289 (2009); DOI: 10.1126/science.1181369



Human chromosomes contact freq.



B

Human chromosomes



 $R_g \sim N^{1/2}$ for polymer melts at an equilibrium, suggesting that each chain in the melts obey statistics of an ideal chain. Therefore, the contact probability $P(r + r_2)$ for given chain contour should be

$$P(\vec{r}_{1},\vec{r}_{2}) = \left(\frac{3}{2\pi\langle r_{12}^{2}\rangle}\right)^{3/2} \exp\left[-\frac{3(\vec{r}_{1}-\vec{r}_{2})^{2}}{2\langle r_{12}^{2}\rangle}\right] - \left[\frac{1}{2}\left(\frac{1}{2}\right)^{2}\right]$$



Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome Erez Lieberman-Aiden *et al. Science* **326**, 289 (2009); DOI: 10.1126/science.1181369

Fig S25. Self-similarity of Fractal Globule



Fractal Globule

B. Composition of Genome



-ome

- Genome: the entirety of an organism's hereditary information, which includes both the genes and the non-coding sequences of the DNA/RNA.
- Transcriptome: the set of all RNA molecules, including mRNA, rRNA, tRNA, and other ncRNA produced in one or a population of cells (connect the genome to gene function).
- Proteome : the entire set of proteins expressed by a genome

Central Dogma





Figure 1 | **The central dogma of gene expression.** In the typical process of eukaryotic gene expression, a gene is transcribed from DNA to pre-mRNA. mRNA is then produced from pre-mRNA by RNA processing, which includes the capping, splicing and polyadenylation of the transcript. It is then transported from the nucleus to the cytoplasm for translation. TSS, transcription start site; TTS, transcription termination site.



Exon-Intron boundary consensus sequences





cf) group I, II intron Self-splicing ribozyme



Locating the Genes in a Genome Sequence

- ORF scanning
 - Open Reading Frames (ORF)

Initiation codon: ATG, Termination codon: TAA, TAG, TGA





Genes and Signals, © 2002 by Cold Spring Harbor Laboratory Press, Chapter 1, Figure 10

Genes and Signals, © 2002 by Cold Spring Harbor Laboratory Press, Chapter 1, Figure 13. (Modified, with permission, from Blackwell Science http://www.blacksci.co.uk/.)

Locating the Genes in a Genome Sequence (ORF scanning)

If DNA has a random sequence each termination codon (TAA, TAG, TGA) appear once every 4³=64 bp;
GC content is greater than 50 %;
Frequency of termination codon from random sequence would be 1/200 - 1/100;
Thus, the shortest length of a putative gene should be ~> 50-100 codons.

• For higher eukaryote, ORFs are not continuous but often split by introns; thus ORF scanning makes use of the followings:

(I) Codon bias

- (2) Exon-Intron boundary consensus sequences
 - (introns in prokaryotes are extremely rare)
- (3) Upstream regulatory sequences
- (4) Upstream CpG islands (vertebrate)
- (5) Comparative genomics (use homology search)
- (6) Databases of cDNA sequences

Locating the Genes for Functional RNAs

- signatures of intramolecular base pairing
- tRNA, Sequence of one of the Escherichia coli tRNA^{leu} genes







Codon bias

CODON	USAGE	IN	Е.	COLI	GENES ¹

	Codon	Amino	9 % 3	Ratio ⁴	Codon	Amino	%	Ratio	Codon	Amino	%	Ratio	Codon	Amino	%	Ratio	
		ac id ²				ac id				acid				acid			
U	ບບບ	Phe (F)	1.9	0.51	UCU	Sei (8)	1.1	0.19	UAU	Туз (Ү)	1.6	0.53	UGU	Cys (C)	0.4	0.43	U
	UUC	Phe (F)	1.8	0.49	UCC	Ser (8)	1.0	0.17	UAC	Tyı (Y)	1.4	0.47	UGC	Cys(C)	0.6	0.57	C
	UUA	Leu (L)	1.0	0.11	UCA	Ser (8)	0.7	0.12	UAA	STOP	0.2	0.62	UGA	STOP	0.1	0.30	Α
	UUG	Leu (L)	1.1	0.11	UCG	Ser (8)	0.8	0.13	UAG	STOP	0.03	0.09	UGG	Tip (V)	1.4	1.00	G
C	CUU	Leu (L)	1.0	0.10	CCU	Pro(P)	0.7	0.16	CAU	His (H)	1.2	0.52	CGU	A1g (R)	2.4	0.42	U
	CUC	Leu (L)	0.9	0.10	CCC	Pro(P)	0.4	0.10	CAC	His (H)	1.1	0.48	CGC	Aig (R)	2.2	0.37	С
	CUA	Leu (L)	0.3	0.03	CCA	Pro(P)	0.8	0.20	CAA	Gln (Q)	1.3	0.31	CGA	Aig (R)	0.3	0.05	Α
	CUG	Leu (L)	5.2	0.55	CCG	P10(P)	2.4	0.55	CAG	Gln (Q)	2.9	0.69	CGG	Aig (R)	0.5	0.08	G
Α	AUU	Ile (I)	2.7	0.47	ACU	Thu (T)	1.2	0.21	AAU	Asn (N)	1.6	0.39	AGU	Ser (8)	0.7	0.13	U
	AUC	Ile (I)	2.7	0.46	ACC	Thu (T)	2.4	0.43	AAC	Asn (N)	2.6	0.61	AGC	Ser (8)	1.5	0.27	C
	AUA	Ile (I)	0.4	0.07	ACA	Thu (T)	0.1	0.30	AAA	Lys (K)	3.8	0.76	AGA	Aig (R)	0.2	0.04	Α
	AUG	Met (M)	2.6	1.00	ACG	Thu (T)	1.3	0.23	AAG	Lys (K)	1.2	0.24	AGG	Aig (R)	0.2	0.03	G
G	GUU	Val(∀)	2.0	0.29	GCU	Ala (A)	1.8	0.19	GAU	Asp (D)	3.3	0.59	GGU	Gly (G)	2.8	0.38	U
	GUC	Val(♥)	1.4	0.20	GCC	Ala (A)	2.3	0.25	GAC	Asp (D)	2.3	0.41	GGC	Gly (G)	3.0	0.40	С
	GUA	Val (V)	1.2	0.17	GCA	Ala (A)	2.1	0.22	GAA	Glu(E)	4.4	0.70	GGA	Gly (G)	0.7	0.09	Α
	GUG	Val(♥)	2.4	0.34	GCG	Ala (A)	3.2	0.34	GAG	Glu(E)	1.9	0.30	GGG	Gly (G)	0.9	0.13	G
	U					С			Α			G					

¹ The data shown in this table is from the Arabidopsis Research Companion on the World Wide Web (//weeds/mgh.harvard.edu). Codon

frequencies for many other bacteria can be found at http://morgan.angis.su.oz.au/Angis/Tables.html.

 2 The letter in parenthesis represents the one-letter code for the amino acid.

³ % represents the average frequency this codon is used per 100 codons.

⁴ Ratio represents the abundance of that codon relative to all of the codons for that particular amino acid.

Annotation of genome sequences



Distribution of exon in genome?



Some interesting facts about human genome

- The smallest protein-coding gene in the human genome is only 500-nt long and has no introns. It encodes a histone protein.
- The largest human gene encodes the protein dystrophin, which is missing or non-functional in the disease muscular dystrophy. This gene is 2.5 million nucleotides in length and it takes over 16 hours to produce a single transcript. However, more than 99 percent of the gene made up of its 79 introns.
- Most of the big differences between human and chimpanzee DNA lie in regions that do not code for genes, according to a new study. Instead, they may contain DNA sequences that control how gene-coding regions are activated and read. - "The differences between chimps and humans are not in our proteins, but in how we use them,"

Sequence variation

- Single Nucleotide Polymorphism (SNP)
- Insertion or deletions of one or more bases
- Repeat length polymorphism, rearrangement

Nonredundant SNP = 1,419,190 in human genome sequence Average SNP density \approx one SNP/1.91kb

Alternative Splicing



...and many different proteins

Note that our genome and chimpanzee's are 99 % identical. However, we humans are very different from them. Gene regulation through transcriptome and proteome is extremely important for species variation !



miRNAs

- MicroRNAs are evolutionarily conserved, small (~22 nucleotide) noncoding RNAs that are encoded within the genomes of almost all eukaryotes, from plants to mammals.
- Derive from dedicated genes, exons and introns of coding genes.
- Base pairing to partially complementary sequences in the 3' untranslated regions (3'UTRs) of target mRNAs.
- miRNA mediate mRNA repression by recruiting the <u>miRNA-induced silencing</u> <u>complex</u> (miRISC)
- Computer-assisted estimates predict ~1000 miRNAs in the human genome.
- miRNAs have multiple targets and thus might regulate ~30% of the proteincoding genome.
- miRNAs function in a variety of biological processes, including tissue differentiation and organ development, control of cell proliferation, apoptosis, fat metabolism and insulin secretion.
- miRNAs are more than a fine-tuning agent, should be upgraded to be a partner of TF.

miRNA biogenesis pathways

Nature Cell Biology 11, 228 - 234 (2009)



Biogenesis pathway of siRNA (exogeneous)







Dicer



NSMB 19:436 (2012)

Argonaute (AGO) proteins

- Bilobal architecture consisting of four evolutionarily conserved domains:
 - N-terminal and Piwi-Argonaute-PAZ domains


Regulatory role of microRNAs

Drosophila embryo



zebrafish embryos



http://www.mirbase.org/index.shtml

Stem-loop sequence hsa-mir-17

Accession	MI000071	
Symbol	HGNC:MIR17	
Description	Homo sapiens miR-17 stem-loop	
Gene family	MIPF0000001; mir-17	
Community annotation	This text is a summary paragraph taken from the <u>Wikipedia</u> entry entitled <u>mir-17 microRNA precursor family</u> . miRBat The miR-17 microRNA precursor family are a group of related small non-coding RM family, includes miR-20, miR-91, and miR-103. miRNAs are transcribed as ~70 nuc product. In this case the mature sequence comes from the 3' arm of the precursor mRNA. A screen of 17 miRNAs that have been predicted to regulate a number of b these patients were noncarriers of BRCA1 or BRCA2 mutations, lending the possible Show Wikipedia entry View @ Wikipedia Edit Wikipedia entry	se NA cle ilit
Stem-loop	ga -ca a g g - au guca auaaugu aagugcuu ca ugcag uag ug a u cagu uauuacg uucacgga gu acguc auc ac g gg aug a g - u gu Get sequence	
Deep sequencing	35854 reads, 48 experiments	

GUCAGAANAAUGUCAAAGUGCUUACAGUGCAGGUAGUGAUAUGUGCAUCUACUGCAGUGAAGGCACUUGUAGCAUUAUGGUGAC

ure sequence hsa-miR-17-5p			
Accession	MIMAT0000070		
revious IDs	hsa-miR-17-5p;hsa-miR-17		
Sequence	14 - caaagugcuuacagugcagguag - 36 Get sequence		
Deep sequencing	31906 reads, 42 experiments		
Evidence	experimental; cloned [2,5-8], Northern [4]		
Validated targets	TARBASE: hsa-miR-17-5p		
Predicted targets	DIANA-MICROT: <u>hsa-miR-17-5p</u> MICRORNA.ORG: <u>hsa-miR-17-5p</u> MIRDB: <u>hsa-miR-17-5p</u> RNA22-HSA: <u>hsa-miR-17-5p</u> TARGETMINER: <u>hsa-miR-17-5p</u> TARGETSCAN-VERT: <u>hsa-miR-17</u> PICTAR-VERT: <u>hsa-miR-17-3p</u> PICTAR-VERT: <u>hsa-miR-17-5p</u>		

Matu

Mature sequence hsa-miR-17-3p			
Accession	MIMAT0000071		
Previous IDs	hsa-miR-17-3p;hsa-miR-17*		
Sequence	51 - acugcagugaaggcacuuguag - 72 Get sequence		
Deep sequencing	3947 reads, 33 experiments		
Evidence	experimental; cloned [1,5,7-8], Northern [1]		
Validated targets	TARBASE: <u>hsa-miR-17-3p</u>		
Predicted targets	DIANA-MICROT: <u>hsa-miR-17-3p</u> MICRORNA.ORG: <u>hsa-miR-17-3p</u> MIRDB: <u>hsa-miR-17-3p</u> RNA22-HSA: <u>hsa-miR-17-3p</u> TARGETMINER: <u>hsa-miR-17-3p</u> PICTAR-VERT: <u>hsa-miR-17-3p</u> PICTAR-VERT: <u>hsa-miR-17-5p</u>		





piRNA (Piwi interacting RNA)

- The largest class of small ncRNA molecules expressed in animal cells.
- forms RNA-protein complexes through interactions with piwi proteins
- , which leads to epigenetic and post-transcriptional gene silencing of retrotransposons and other genetic elements in germ line cells.
- Distinct from microRNA or siRNA in size (26-31 nt)
- has highly complex and heterogeneous population revealed by recent DNA sequencing.
- Biogenesis of piRNAs not fully understood yet.

Reorganization of Genome

Homologous recombination

Site-specific recombination



minisatellite : 10-60 bp around core sequence(--GGGCAGGANG--) microsatellite : 2-13 nt long short tandem repeats

CRISPRs (Clustered Regularly Interspaced Short Palindromic Repeats)

Short sequence tags from invading genetic elements are actively incorporated into the host's CRISPR locus to be transcribed and processed into a set of small RNAs that guide the destruction of foreign genetic material.

- Immune system of bacteria and archaea against predators
- Confers resistance to exogenous genetic elements (plasmids, phages).
- Short segments of foreign DNA (spacers) incorporated into genome between CRISPR repeats and serve as a 'memory' of past exposures.
- CRISPR spacers are used to recognize and silence exogenous genetic elements in a manner analogous to RNAi in eukaryotic organisms.



Epigenetics



- Chemical modifications of the DNA (methylation) or the histones (deacetylation) alter the chromatin structure without changing the DNA sequence
- Inappropriate gene silencing (or activation), brought about by epigenetic modification, cause a number of human diseases including cancers.

Epígenetíc landscape



S. Huang / Progress in Biophysics and Molecular Biology 110 (2012) 69-86

Evolution

- Natural selection Survival of the fittest
- Mutation Malicious / Beneficial mutation





The *E. coli* long-term evolution experiment is an ongoing study in experimental evolution led by Richard Lenski that has been tracking genetic changes in 12 initially identical populations of asexual Esherichia coli bacteria since 24 February 1988. The populations reached the milestone of 50,000 generations in February 2010.

Since the experiment's inception, Lenski and his colleagues have reported a wide array of genetic changes; some evolutionary adaptation have occurred in all 12 populations, while others have only appeared in one or a few populations. One particularly striking adaption was the evolution of a strain of *E. coli* that was able to use citric acid as a carbon source in an aerobic environment.



Growth in cell size of bacteria in the Lenski experiment

Growth rate (=Fitness)



Number of consecutive divisions of the old-pole mother cell

Current Biology 20, 1099–1103,

3. Foldings and Functions of RNA



Figure 6-21a Molecular Biology of the Cell 5/e (© Garland Science 2008)





RNA folding



Primary to Secondary to 3D structure formation

 Characteristic base-pairing rule for RNA (stacking, loop, bulge, internal multiloop free energy,)
 →1970s (Tinoco, Uhlenbeck, Crothers, Porschke...)



Statistics of RNA structures in PDB It is about time...

Α



RNA motifs for tertiary structure formation formation



Ribose zipper, kissing hairpin, Coaxial stack,



Prediction of RNA native state (RNA folding problem) is difficult

Building block of RNA : Chemically similar 4-nt (instead of chemically dissimilar 20 aa)
 Polyelectrolyte nature (charged phosphate group)

3. Uniformity of hydrophilic backbone along with the lack of diversity in the bases makes RNA closer to "homopolymer" than polypeptide chains -> Make the stability gap smaller and alternate structure possible.



See <u>http://www.major.iric.ca/MC-Pipeline/</u> (NCM, nucleotide cyclic motifs).

RNA energy landscape is rugged One dimensional chain + Heterogeneous interactions → Topological / Energetic frustration





Thirumalai and Hyeon, Biochemistry (2005)

$$f_N(t) = \Phi_{max} - \Phi e^{-k_{fast}t} - \sum_i A_i e^{-k_i t}.$$

The value of Φ depends on point mutation, ionic environment, and initial condition ...



A 0.3

0.2

0.1

0

0.8

0.6

0.4

0

U,

0

00000

5

80 120

Fraction fast (

0.4

Misfolded

Native

Misfolded

Native

N.

>250 mM

Na⁺

600

L Icommitment

I_{commitment}

150-250 mM Na*

40

400

[Na⁺], mM

200

rap

Time, s

Fraction native

в

Fraction correct (
)

<150 mM Na* U

150-250

mM Na*

>250 mM Na*

<150 mM Na* U -

Initial condition dependent folding routes of T. Ribozyme

Varying the Na⁺ concentration in the preincubation from 20 **620** to mM increased the fraction that folded correctly from 0.2 to ~ 0.8 (Fig. 1B) even though folding was always under the same conditions (5 mM Mg²⁺, 20 mM Na⁺). Also, the fraction of the correct-folders that folded fast (1 s⁻¹) increased from ~0.45 to ~0.9.



Point mutation



Pan et al. (1997) J. Mol. Biol.

Φ=0.06-0.1 for Tetrahymena ribozyme. But, a point mutation U273A in P3 increases Φ

Rugged folding landscape of RNA and functions

- RNA consisting of four chemically similar building blocks (A,C, G, U)
 - \rightarrow RNA are closer to homopolymers
 - → rugged landscape
 - \rightarrow adopt alternative structures
 - → kinetic partitioning mechanism

Guo & Thirumalai (1995) Biopolymers Pan, Thirumalai, Woodson (1997) JMB Thirumalai & Hyeon (2005) Biochemistry



 Foldings and functions of RNA are associated with conformational adaptation (CBA↔NBA) in response to cellular signals (metal ions, metabolites, proteins, nucleic acids,...)



divalent metal-ion induced secondary structure rearrangement





Functions of RNA are often dictated by the conformational switches between the alternative states



Al-Hashimi & Walter (2008) COSB







DEAD-box proteins assist secondary structure rearrangement by acting as RNA strand separators (or helicases).



"DEAD-box 'Helicase' Proteins as Chaperones of RNA Folding"

Conclusions

- RNAs are more primitive and homopolymer-like than proteins.
- RNA folding landscapes are rugged.
- Kinetic partitioning mechanism
- Ionic environment is crucial for RNA folding and function
- RNA chaperone is essential to facilitate RNA folding

4. Gene expression, Genetic switch





$$\frac{dm}{dt} = \alpha_m - \beta_m m$$

$$m^* = \frac{\alpha_m}{\beta_m} \approx (1 - 10)nM$$

$$\frac{dp}{dt} = \alpha_p m - \beta_p p$$

 $p^* = \frac{\alpha_m \alpha_p}{\beta_m \beta_p} = (10^2 - 10^3) nM$





For Ecoli, mRNA degrades much faster than proteins, and reaches the SS values. $\beta_m \gg \beta_p$

$$\frac{dm}{dt} \approx 0; \quad m^* = \frac{\alpha_m}{\beta_m} g_R(R)$$
$$\frac{dR}{dt} = \alpha_p m^* - \beta_p R$$

Since $g_R(R) \sim 1$ for $R << K_R$ and $g_R(R) \sim 0$ for $R >> K_R$ it is expected that R(t) increase slowly when $R > K_R$ and the steady state value of R is smaller than unrepressed case.



In 2D





For Ecoli, mRNA degrades much faster than proteins, and reaches the SS values. $\beta_m \gg \beta_p$

$$m_1^*(R_2) = \frac{\alpha_m}{\beta_m} g_R(R_2) \qquad m_2^*(R_1) = \frac{\alpha_m}{\beta_m} g_R(R_1)$$

$$\frac{dR_1}{dt} = \frac{\alpha_p \alpha_m}{\beta_m} g_R(R_2) - \beta_p R_1 \xrightarrow{ss} R_1^* = \frac{\alpha_m \alpha_p}{\beta_m \beta_p} g_R(R_2^*)$$

$$\frac{dR_2}{dt} = \frac{\alpha_p \alpha_m}{\beta_m} g_R(R_1) - \beta_p R_2 \xrightarrow{ss} R_2^* = \frac{\alpha_m \alpha_p}{\beta_m \beta_p} g_R(R_1^*)$$



At crossing points

 $\frac{dR_1}{dt} = \frac{dR_2}{dt} = 0$

 $R_1^* = \frac{\alpha_m \alpha_p}{\beta_m \beta_p} g_R(R_2^*) \quad ... \text{ to assess the stability}$ $\Rightarrow R_1^* \qquad \text{consider } \delta R_i = R_i - R_i^*$

 $\frac{dR_1}{dt} = \frac{\alpha_p \alpha_m}{\beta_m} g_R(R_2) - \beta_p R_1 = F(R_1, R_2)$ $\frac{dR_2}{dt} = \frac{\alpha_p \alpha_m}{\beta_m} g_R(R_1) - \beta_p R_2 = G(R_1, R_2)$

$$\longrightarrow \left(\begin{array}{c} \delta R_1 \\ \delta R_2 \end{array}\right) = \left(\begin{array}{c} X_1 \\ Y_1 \end{array}\right) e^{\lambda_1 t} + \left(\begin{array}{c} X_2 \\ Y_2 \end{array}\right) e^{\lambda_2 t}$$

$$\frac{d}{dt} \begin{pmatrix} \delta R_1 \\ \delta R_2 \end{pmatrix} = \begin{pmatrix} \frac{dF}{dR_1} & \frac{dF}{dR_2} \\ \frac{dG}{dR_1} & \frac{dG}{dR_2} \end{pmatrix}_{R_1^*, R_2^*} \begin{pmatrix} \delta R_1 \\ \delta R_2 \end{pmatrix}$$
$$\begin{pmatrix} -\beta_p - \lambda & \frac{\alpha_p \alpha_m}{\beta_m} g'_R(R_2^*) \end{pmatrix} \begin{pmatrix} X \end{pmatrix} = 0$$

$$\frac{\alpha_p \alpha_m}{\beta_m} g'_R(R_1^*) \qquad -\beta_p - \lambda \qquad \int_{R_1^*, R_2^*} \left(\begin{array}{c} Y \end{array} \right)^{=0}$$
$$\lambda_{1,2} = -\beta_p \pm \frac{\alpha_m \alpha_p}{\beta_m} \sqrt{g_R'(R_1^*)g_R'(R_2^*)}$$

Then, at each fixed point $g_{R}'(R_{1}^{*}) = g_{R}'(R_{2}^{*}) \approx 0 \rightarrow \lambda_{1,2} < 0 \text{ stable fixed point}$ $g_{R}'(R_{1}^{*}) = g_{R}'(R_{2}^{*}) \approx 1 \rightarrow \lambda_{1} > 0, \lambda_{2} < 0 \text{ unstable}$



Construction of a genetic toggle switch in *Escherichia coli*

Timothy S. Gardner*†, Charles R. Cantor* & James J. Collins*†

* Department of Biomedical Engineering, † Center for BioDynamics and ‡ Center for Advanced Biotechnology, Boston University, 44 Cummington Street, Boston, Massachusetts 02215, USA





$$\frac{\mathrm{d}U}{\mathrm{d}t} = \frac{\alpha_1}{1 + V^\beta} - U$$
$$\frac{\mathrm{d}V}{\mathrm{d}t} = \frac{\alpha_2}{1 + U^\gamma} - V$$



0

Small Regulatory RNAs May Sharpen Spatial Expression Patterns

Erel Levine[®], Peter McHale[®], Herbert Levine[®]

Center for Theoretical Biological Physics, University of California San Diego, La Jolla, California, United States of America

















Look at the mean and variance about a spatially *inhomogeneous* steady-state.



Patterning mechanism remains robust in very small populations