CODING A LIFE FULL OF ERRORS



What is Life? (biological and artificial)

0

- Self-replication.
- Emergence.
- Evolution.
- Non-equilibrium.
- Information.
- Geometry.
- Stochasticity.

- Viruses (bio and computer).
- Growth and form.
- Natural algorithms .
- Learning & robots.
- **Codes and Errors.**





Living information is carried (mostly) by molecules

"Living systems"

- I. Self-replicating information processors.
- II. Evolve collectively.
- III. Made of molecules.

- Generic properties of molecular codes subject to evolution?
- Information theory approach?



Challenges of molecular codes: rate and distortion

Distortion

- Noise, crowded milieu.
- Competing lookalikes.
- Weak recognition interactions $\sim k_B T$.
- Need diverse meanings.

"Synthesis of reliable organisms from unreliable components" (von Neumann, Automata Studies 1956)

Rate

How to construct the low-rate molecular codes at minimal cost of resources?

Rate-distortion theory (Shannon 1956)

Inside E. Coli, D. Goodsell

Codes are mappings, channels, representations, models...

• Code ϕ is a mapping between spaces, $\phi: S \to \mathcal{M}$.





- Molecular codes map\translate between molecular spaces\languages.
- Molecular spaces have inherent geometry\topology.
- Coding machinery affects organism's fitness.

Outline: molecular codes and errors

- → Living and artificial self-replication.
 - The main molecular codes of life (central dogma).
 - The translation machinery:
 - The genetic code, ϕ : codons \rightarrow amino-acids.
 - The ribosome and the problem of molecular recognition.



- Basic coding theory: geometrical aspects.
 - How codes cope with errors.
- Emergence and evolution of codes: rate-distortion.
- Accuracy vs. rate: proofreading schemes.



Coding and the problem of self-replication



Proposed demonstration of simple robot self-replication,

from advanced automation for space missions, NASA conference 1980.

Self-replication and accuracy in computers

Theory of Self-Reproducing Automata

Lectures on

PROBABILISTIC LOGICS AND THE SYNTHESIS OF RELIABLE ORGANISMS FROM UNRELIABLE COMPONENTS

JOHN VON NEUMANN edited and completed by Arthur W. Burks





Theory of (1966) Self-Reproducing Automata

Von Neumann's universal constructor

Self-reproducing machine: constructor + tape (1948/9).



Von Neumann's design allows open-ended evolution

Motivated by biological self-replication:

- Construction *universality*.
- Evolvability.

Key insight (before DNA) separation of information and function.

- Tape is read twice: for construction and when copied.
- How to design fast/accurate/compact constructor?
- Requires efficient and accurate coding...



Implementation by Nobili & Pesavento (1995)

Outline: molecular codes and errors

- Living and artificial self-replication.
- - The translation machinery:
 - − The genetic code, ϕ : codons → amino-acids.
 - The ribosome and the problem of molecular recognition.



- Basic coding theory: geometrical aspects.
 - How codes cope with errors.
- Emergence and evolution of codes.
- Accuracy vs. rate: proofreading schemes.



Dual spaces of DNA and proteins

DNA

- Building blocks :
 - 4 nucleic bases = $\{A, T, G, C\}$.



- Polymer: DNA double-helix.
- Inert information storage ("tape")

protein

• Building blocks:

20 amino acids.



- Polymer = protein.
- Functional molecules ("constructor")

RNA intermediates can be both tapes and machines



• Primordial "RNA world" :

RNA molecules are both information carriers (DNA) and executers (proteins).

The Central Dogma of molecular biology

Francis Crick: Ideas on Protein Synthesis (Oct. 1956)

The Doctrine of the Triad.

The Central Dogma: "Once information has got into a protein it can't get out again". Information here means the sequence of the amino acid residues, or other sequences related to it.

That is, we may be able to have



The central dogma graphs the main information channels between nucleotides and proteins

replication

translation

- Information from DNA sequence cannot be channeled back from protein to either protein or nucleic acid.
- 3 information carriers: DNA, RNA protein and 3×3 potential channels:
- 3 general channels (occur in most cells).
- 3 special channels (under "specific" conditions).
- 3 unknown transfers (no example (yet?)).

"special" information transfers

• **Reverse transcription (** $RNA \rightarrow DNA$):

Reverse transcriptase, in retroviruses (e.g. HIV) and eukaryotes (retrotransposons and telomeres).

• **RNA replication** (RNA \rightarrow RNA):

Many viruses replicate by *RNA-dependent RNA polymerases* (also used in eukaryotes for RNA silencing).

• **Direct translation** (DNA → protein):

demonstrated in extracts from E. coli which expressed proteins from foreign ssDNA templates.



RNA replication

Channels outside the dogma: Epigenetic information transfer

- Changes in methylation of DNA alter gene expression levels.
- Heritable change is called epigenetic.
- Effective information change but not DNA sequence.
- Gene "switched on"
- Active (open) chromatin
 Unmethylated cytosines (white circles)
- Acetylated histones

- Others:
- post-translational modification...

Gene "switched off"

- Silent (condensed) chromatin
- Methylated cytosines (red circles)
- Deacetylated histones





Post-translational modifications of proteins :

- Extends functionality by attaching other groups (e.g. acetate).
- Changes chemical nature of aa.
- Structural changes (disulfide bridges). 🖞
- Compensate for missing tRAS (Helicobacter pylori).
- Enzymes may remove amino acids or cut the peptide chain in the middle.





Universal constructors in the arts



The "replicator" (Star Trek)



The Santa Claus machine (A. Sward)

Universal constructors in the arts and in reality (?)



The Economist

FEBRUARY 12TH-18TH 2011

Europe loses the mobile-phone war Africa's new wealth Japan's tea party How to switch off the internet The shoe-thrower's index

Print me a Stradivarius The manufacturing technology that will change the world

Economist.com

This violin was made using an EOS laser-sintering 3D printer (and it plays beautifully)

Self-printing?

Outline: molecular codes and errors

- Living and artificial self-replication.
- The main molecular codes of life (central dogma).
- The translation machinery:
 - The genetic code, ϕ : codons \rightarrow amino-acids.
 - The ribosome and the problem of molecular recognition.



- Basic coding theory: geometrical aspects.
 - How codes cope with errors.
- Emergence and evolution of codes.
- Accuracy vs. rate: proofreading schemes.



The translation machinery is the main system of the living von Neumann's universal constructor

- Machinery parts = tRNA + synthetase + ribosome...
- The translation machinery conveys information from nucleotides to proteins.

Synthetases charge tRNAs according to the genetic code.



Ribosomes translate nucleic bases to amino acids

 Ribosomes are *large* molecular machines that synthesize proteins with mRNA blueprint and tRNAs that carry the genetic code.



Goodsell, The Machinery of Life

protein

1. Is the code $\phi(c)$ adapted to the noise problem?

Ribosome needs to recognize the correct tRNA





(ii) unbinding correct tRNAs: amino-acid = ϕ (codon)

2. How to construct fast\accurate\small molecular decoder ?

2.Decoding at the ribosome is a molecular recognition problem



• Central problem in biology and chemistry:

How to evolve molecules that recognize in a noisy environment?

(crowded, thermally fluctuating, weak interactions).

- How to estimate recognition performance ("fitness")?
- What are the relevant degrees-of-freedom? **Dimension**? **Scaling**?
- What is the role of conformational changes?

Ribosome sets physical limit on self-reproduction rate

Large fraction of cell mass is ribosomes.

- In self-reproduction each ribosome should self-reproduce.
- Sets lower bound on self-reproduction rate .

 $T \ge \frac{\text{mass}_{\text{ribo}}}{R_{\text{C}}} \approx \frac{10^4 \text{ amino-acids}}{20 \text{ amino-acids/sec}} = 500 \text{ sec}$

• "Fastest " growing bacteria (*Clostridium perfringens*): *T* ~ 500 sec.

Problem: how ribosome accuracy affects fitness depends on

- (i) Basic protein properties (mutations).
- (ii) Biological context (environment etc.).



Outline: molecular codes and errors

- Living and artificial self-replication.
- The main molecular codes of life (central dogma).
- The translation machinery:
 - The genetic code, ϕ : codons \rightarrow amino-acids.
 - The ribosome and the problem of molecular recognition.



- Basic coding theory: geometrical aspects.
 - How codes cope with errors.
- Emergence and evolution of codes.
- Accuracy vs. rate: proofreading schemes.



1. The genetic code maps DNA to protein

 Genetic code: maps 3-letter words in 4-letter DNA language (4³ = 64 codons) to protein language of 20 amino acids.

$$codon = b_1 b_2 b_3, \ b_i \in \{A, T, G, C\}.$$

 $\phi(codon) \rightarrow amino-acid.$

• Genetic code embeds the codon-graph (Hamming graph) into space of amino-acids ("digital to analog").



• Translation machinery, whose main component is the ribosome, facilitates the map.

Outline: molecular codes and errors

- Living and artificial self-replication.
- The main molecular codes of life (central dogma).
- The translation machinery:
 - The genetic code, ϕ : codons \rightarrow amino-acids.
 - The ribosome and the problem of molecular recognition.



- → Basic coding theory: geometrical aspects.
 - How codes cope with errors.
 - Emergence and evolution of codes.
 - Accuracy vs. rate: proofreading schemes.



Coping with unreliability of coding machinery

•	Error detecting code - parity checking.	signal	binary	parity
•	One check: Odd parity → mistake (e.g. 0111).	0	000	0
•	Retransmission	1	001	1
		2	010	1
•	Single error can be detected but not corrected.	3	011	0
		4	100	1
		5	101	0
•	The redundancy of the code:	6	110	0
	total # of bits n	7	111	1

$$R = \frac{1}{\text{\# of message bits}} = \frac{1}{n-1}$$

Error correction requires minimal redundancy

- *Error correcting code* can detect and correct errors.
- Multiple checks Locating errors by confluence.
- *Triplication* code send each message thrice (R = 3).
- What is the minimal number of checks *m*?
 - To locate *n* positions requires $2^m \ge n+1$.

$$R = \frac{n}{n-m} \ge \left(1 - \frac{\log_2(n+1)}{n}\right)^{-1} \simeq 1 + \frac{\log_2 n}{n}$$

Hamming's code reaches this limit.

	1001010111
0	1001101101
1	0010011110
0	0011111001
0	1101010001
1	1101001100

Rectangular code $R = 1 + \frac{2}{\sqrt{\pi}}$



Geometric view of error correction and detection

Messages are mapped between hypercubes

 $\phi: Y_n \to Y_{n+m} \quad (1001 \to 1001101)$

• Metric is the Hamming distance:

 $|x_i - x_j| = \#$ of different letters

• Sphere is $S = \{x \mid |x - x_0| \le r\}.$



1 Three-dimensional spheres about (0, 0, 0) and (1, 1, 1)



Error correction is packing hard spheres

- To correct *r* errors the spheres should be at least at distance d = 2r + 1.
- Correction: move to nearest sphere center.
- How many words can be encoded?
 Or how many spheres can be packed?



total volume \geq sphere volume \times # spheres $2^{n} \geq (n+1) \times 2^{n-m} \Longrightarrow 2^{m} \geq n+1$

spheres = # words =
$$2^{n-m} \le \frac{2^n}{n+1}$$

Shannon's channel coding theorem sets upper limit on the capacity of a noisy channel

- Noisy channel is defined by stochastic input\output $\phi(s|m)$.
- **Channel capacity** measures the input\output correlation

$$C = \max_{\phi(s)} I(S; M) = \max_{\phi(s)} \left\langle \log_2 \frac{\phi(s, m)}{\phi(s)\phi(m)} \right\rangle$$

- Channel rate $R = \lim_{n \to \infty} \frac{\log_2(\# \text{ words})}{n}$
- Shannon's coding theorem (1948/9): R=C.

• Proof: show that # hard spheres is
$$2^{nR} = 2^{nI(S,M)} = 2^{nC}$$
.

• Upper limit achieved only "recently" (turbo codes, LDPC).



The genetic code is a smooth mapping

UUU Phe	UCU Ser	UAU Tyr	UGU Cys
UUC Phe	UCC Ser	UAC Tyr	UGC Cys
UUA Leu	UCA Ser	UAA TER	UGA TER
UUG Leu	UCG Ser	UAG TER	UGG Trp
CUU Leu	CCU Pro	CAU His	CGU Arg
CUC Leu	CCC Pro	CAC His	CGC Arg
CUA Leu	CCA Pro	CAA GIn	CGA Arg
CUG Leu	CCG Pro	CAG GIn	CGG Arg
AUU lle	ACU Thr	AAU Asn	AGU Ser
AUC lle	ACC Thr	AAC Asn	AGC Ser
AUA lle	ACA Thr	AAA Lys	AGA Arg
AUG Met	ACG Thr	AAG Lys	AGG Arg
GUU Val	GCU Ala	GAU Asp	GGU Gly
GUC Val	GCC Ala	GAC Asp	GGC Gly
GUA Val	GCA Ala	GAA Glu	GGA Gly
GUG Val	GCG Ala	GAG Glu	GGG Gly



- Degenerate (20 out of 64) \rightarrow "spheres".
- Compactness of amino-acid regions.
- **Smooth** (similar "color" of neighbors).
- But not immune to one-letter errors ("soft" spheres).

Generic properties of molecular codes?


Gray code "smooths" the impact of errors

- Invented by Émile Baudot
 for telegraphy.
- Often used in $A \rightarrow D$ and $D \rightarrow A$ applications.
- Minimizes the number of changes between close by values → smooth code.
- Used in many modulation schemes.

(e.g. phase shift).



Figure 5-15.1 Eight-sector binary wheel





The smooth genetic code as a combinatorial game



"Marble packing"

(A) Max colors.

(B) Same\similar color of neighbors.

The genetic code maps codons to amino-acids

- Molecular code = map relating two sets of molecules.
- Spaces defined by similarity of molecules (size, polarity etc.)



Fitter codes have minimal distortion



$$\mathbf{D} = \left\langle c_{\alpha\omega} \right\rangle = \sum_{paths} p_{path} c_{\alpha\omega} = \operatorname{Tr} \left(e \cdot r \cdot d \cdot c \right)$$

- Distortion of noisy channel, **D** = average distortion of AA.
- *r* defines topology of codon space.
- *c* defines topology of amino-acid space.

Smooth codes minimize distortion



- Noise confuses close codons.
- Smooth code:
 - close codons = close amino-acids.
 - \rightarrow minimal distortion.

• Optimal code must balance contradicting needs for **smoothness** and **diversity**.



Channel rate is code's cost

- Diverse codes require high **specificity** = high binding energies ε .
- Cost ~ average binding energy $< \varepsilon >$.
- Binding prob. ~ Boltzmann: $\mathbf{E} \sim e^{\varepsilon/T}$.

$$I \sim \sum_{\alpha,i} e_{\alpha i} \ln e_{\alpha i} \sim \left\langle \mathcal{E}_{\alpha i} \right\rangle_e$$



• Cost *I* = Channel Rate (bits/message)

Rate-distortion theory of noisy information channels

- How well a mapping represents a signal?
- Example: quantization of continuous signal.
- The average distortion of a signal

$$\mathsf{D} = \left\langle c_{\alpha\omega} \right\rangle = \sum_{paths} p_{path} c_{\alpha\omega}$$



 Main theorem: there exists a *rate-distortion function* R(D) which is the minimal required rate R to achieve distortion D.



0.9

(Shannon, Kolmogorov 1956)

Code's fitness combines rate and distortion of map

Fitness = Gain x Distortion + Rate $H = \kappa D + I$

- **Gain** β increases with organism complexity and environment richness.
- Fitness **H** is "free energy" with inverse "temperature" κ.
- Evolution varies the gain κ .
- Population of self-replicators evolving according to code fitness *H*: mutation, selection, random drift.



Code emerges at a critical coding transition

• Low gain β : Cost too high

 \rightarrow no specificity \rightarrow **no code**.

• Code emerges when β increases:

channel starts to convey information $(I \neq 0)$.

- Continuous phase transition.
- Emergent code is smooth, low mode of **R**.





Distortion **Q**

Rate-distortion theory (Shannon 1956)

The emergent code is smooth

- Example: mapping between two cycles.
- $\kappa_c = 0.52$ $\kappa = 0.66$ $\kappa = 0.55$ $\kappa = 0.79$ meanings 0.9 0.8 0.7 symbols 0.6 $\kappa = 24.7$ 0.5 $\kappa = 1.26$ $\kappa = 1.52$ $\kappa = 6.13$ -0.4 0.3 0.2 0.1 0



Order parameter: deviation from random map

$$\delta e_{\alpha i} = e_{\alpha i} - e_{\alpha i}^{\mathrm{rand}}$$





Emergent code is a smooth mode of error-Laplacian

- Lowest excited modes of graph-Laplacian R .
- Single maximum for lowest excited modes (Courant).
- Every mode corresponds to amino-acid :

low modes = # amino-acids.

- \rightarrow single contiguous domain for each amino-acid.
- ightarrow Smoothness.



Probable errors define the graph and the topology of the genetic code

• Codon graph = codon vertices + 1-letter difference edges (mutations).





- Non-planar graph (many crossings).
- Genus γ = # holes of embedding manifold.

• Graph is holey : embedded in $\gamma = 41$

(lower limit is $\gamma = 25$)

Coloring number limits number of amino-acids

- Q: Minimal # colors suffices to color a map where neighboring countries have different colors?
- A: Coloring number, a topological invariant (function of genus):

$$chr(\gamma) = \left| \frac{1}{2} \left(7 + \sqrt{1 + 48\gamma} \right) \right|.$$

(Ringel & Youngs 1968)

 $max(\# amino-acids) = chr(\gamma)$



- From Courant 's theorem + "convexity" (tightness).
- Genetic code: $\gamma = 25-41 \rightarrow$ coloring number = 20-25 amino-acids

The genetic code coevolves with accuracy

 A path for evolution of codes: from early codes with higher codon degeneracy and fewer amino acids to lower degeneracy codes with more amino acids.

1 st	2 nd	3 rd	γ	chr #	
1	4	1	0	4	
2	4	1	1	7	
4	4	1	5	11	
4	4	2	13	16	
4	4	3	25	20	
4	4	4	41	25	



Part I: Summary

- The translation machinery:
 - The genetic code, ϕ : codons \rightarrow amino-acids.
- Genetic code is a smooth map that minimizes distortion.
- Model for emergence: phase transition in a noisy mapping.
- Free energy is rate-distortion function.
- Continuous transition.
- Topology governs emergent code.



Sources:

- Shannon, Mathematical Theory of Communication.
- Hamming, Coding and computation.
- von Neumann, In Automata Studies.
- Feynman, Lectures on Computation.
- Cover & Thomas, Elements of information theory.
- Berger T, Rate distortion theory.

Papers on coding: follow PITP link

CODING A LIFE FULL OF ERRORS



Ribosomes translate nucleic bases to amino acids

 Ribosomes are *large* molecular machines that synthesize proteins with mRNA blueprint and tRNAs that carry the genetic code.



Goodsell, The Machinery of Life

protein

Is the code $\phi(c)$ adapted to the noise problem?



Fig. 2. Pancreatic exocrine cell. Array of cisternae of the rough surfaced endoplasmic reticulum.

cs, cisternal space; cm, cytoplasmic matrix (cell sol); fr, free ribosomes; ar, attached ribosomes; mer, membrane of the endoplasmic reticulum.

x 50,000

George Palade (50s)

Ribosome needs to recognize the correct tRNA





(i) binding wrong tRNAs: amino-acid $\neq \phi(\text{codon})$ (ii) unbinding correct tRNAs: amino-acid $= \phi(\text{codon})$

How to construct fast\accurate\small molecular decoder ?

2.Decoding at the ribosome is a molecular recognition problem



• Central problem in biology and chemistry:

How to evolve molecules that recognize in a noisy environment?

(crowded, thermally fluctuating, weak interactions).

- How to estimate recognition performance ("fitness")?
- What are the relevant degrees-of-freedom? **Dimension**? **Scaling**?
- What is the role of conformational changes?

Ribosome sets physical limit on self-reproduction rate

Large fraction of cell mass is ribosomes.

- In self-reproduction each ribosome should self-reproduce.
- Sets lower bound on self-reproduction rate .

 $T \ge \frac{\text{mass}_{\text{ribo}}}{R_{\text{C}}} \approx \frac{10^4 \text{ amino-acids}}{20 \text{ amino-acids/sec}} = 500 \text{ sec}$

Problem: how ribosome accuracy affects fitness depends on

- (i) Basic protein properties (mutations).
- (ii) Biological context (environment etc.).



Ribosomes are complicated machines with many d.o.f.

Ribosomes are made of proteins and RNAs:

- $\sim 10^4$ nucleic bases in RNA.
- ~ 10^4 amino-acids in proteins.
- Total mass : ~ $3 \cdot 10^6$ a.u.
- High-res structure is known (Yonath et al.).

Within this known complexity:

- What are the relevant degrees-of-freedom?
- How does this machine operate?



(magenta – RNA, grey – protein, from Goodsell, *Nanotechnology*)

Decoding is determined by energy landscapes of correct and wrong tRNAs

- Decoding is multi-stage process.
- Kinetics involves *large* conformational changes.





Steady-state decoding rates (Arrhenius law, $k \propto e^{-\Delta G}$)

$$R_C \sim \frac{1}{e^{b_1} + e^{b_2} + e^{b_3}}$$

$$R_W \sim \frac{1}{e^{b_1} + e^{b_2} + e^{\delta + b_3}}$$

In Ehrenberg's notation

• Merge the first two barriers to get Michaelis-Menten kinetics: $e^{B} = e^{b_1} + e^{b_2}$.

$$E + S^{c(nc)} \underset{k_{d}^{c(nc)}}{\overset{k_{d}^{c(nc)}}{\longrightarrow}} E S^{c(nc)} \underset{k_{d}^{c}}{\overset{k_{d}^{c}}{\longrightarrow}} E + P^{c(nc)}$$

$$\int_{s_{2}}^{b_{1}} \overbrace{s_{2}}{\overset{k_{d}^{c}}{\longrightarrow}} \underbrace{f_{d}^{c}}{\overset{k_{d}^{c}}{\longrightarrow}} E + P^{c(nc)}$$

$$\int_{s_{2}}^{b_{1}} \overbrace{s_{2}}{\overset{k_{d}^{c}}{\longrightarrow}} \underbrace{f_{d}^{c}}{\overset{k_{d}^{c}}{\longrightarrow}} E + P^{c(nc)}$$

$$\frac{j^{c}}{s^{c}e} = \left(\frac{k_{cat}}{K_{m}}\right)^{c} = k_{a}^{c} \frac{k_{c}^{c}}{k_{c}^{c} + k_{d}^{c}} = k_{a}^{c} \frac{1}{1 + a} \iff R_{c} = k_{a}^{c} \frac{1}{1 + e^{b_{3} - B}} \propto \frac{1}{e^{B} + e^{b_{3}}}$$

$$\frac{j^{nc}}{s^{nc}e} = \left(\frac{k_{cat}}{K_{m}}\right)^{nc} = k_{a}^{nc} \frac{k_{c}^{nc}}{k_{c}^{nc} + k_{d}^{nc}} = k_{a}^{nc} \frac{1}{1 + d_{d}a} \iff R_{W} = k_{a}^{W} \frac{1}{1 + e^{b_{3} - B + \delta}} \propto \frac{1}{e^{B} + e^{b_{3} + \delta}}$$

$$A = \left(\frac{k_{cat}}{K_{m}}\right)^{c} / \left(\frac{k_{cat}}{K_{m}}\right)^{nc} = d_{a} \frac{1 + d_{d}a}{1 + a} = \frac{d_{a} + da}{1 + a} \iff \frac{R_{c}}{R_{W}} = \frac{e^{B} + e^{b_{3} + \delta}}{e^{B} + e^{b_{3}}} = \frac{1 + e^{b_{3} - B + \delta}}{1 + e^{b_{3} - B}}$$

$$a = \frac{\kappa_d}{k_c^c} = e^{b_3 - B}; \qquad d_a = \frac{\kappa_a}{k_a^{nc}} = 1;$$
$$d_d = \frac{k_d^{nc}}{k_c^{nc}} / \frac{k_d^c}{k_c^c} = e^{\delta}; \quad d = d_d d_a = e^{\delta};$$

1 (

$$\left(\frac{k_{cat}}{K_m}\right)^c = k_a^c \frac{1}{1+a} = k_a^c \frac{d-A}{d-1}$$
$$\iff R_C = k_a^C \frac{1}{1+e^{b_3-B}} = k_a^C \cdot \frac{e^{\delta} - \left(\frac{R_C}{R_W}\right)}{e^{\delta} - 1}$$

Ribosome kinetics exhibits large dimensionality reduction

- Effective dimension decreases by at least 3 orders of magnitude:
 - ~ 10⁴ structural parameters \rightarrow ~ 10 kinetic parameters (energy landscape).





 Generic phenomenon in biomolecules: many catalytic molecules (enzymes) can be described by a few kinetic parameters (transition state landscape).

What is the origin of dimensionality reduction?

- Hints:
- Protein function mainly involve the lowest modes of their vibrational spectra (hinges).
- Sectors: "Normal modes" of sequence evolution (Leibler & Ranganthan).

Transition states reduce the dimensionality of effective parameter space



Figure 2.1 Transition states occur at the peaks of the energy profile of a reaction, and intermediates occupy the troughs.

$$k_1 = \left(\frac{kT}{h}\right) \exp\left(\frac{-\Delta G^{\ddagger}}{RT}\right)$$

Reaction coordinate

FEBRUARY, 1935

JOURNAL OF CHEMICAL PHYSICS

VOLUME 3

The Activated Complex in Chemical Reactions

HENRY EYRING, Frick Chemical Laboratory, Princeton University (Received November 8, 1934)

leave no doubt of the proposed method of procedure in a particular case. We may write for the specific reaction rate constant for a reaction of any order

$$k_i = c(F_a/F_n)(\bar{p}/m^*) = c(F_a'/F_n)(kT/h)e^{-E_0/kT}$$
(10)

where F_a is simply the partition function (or Zustandsumme) for the activated state and F_n is the same quantity for the normal state. F_a' is the partition function for the activated complex for all the normal coordinates except the one in which decomposition is occurring. The partition function for this normal coordinate is included in the factor $(kT/h)e^{-E_0/kT}$. The other quantities have been defined.

Theory can be tested with measured rates

• The codon-specific stages are Codon recognition and GTP activation.



(U	(UUU) Cognate (CUC) Non-cognate						
k ₁	100-140 1/(µM·s)	100-140 1/(µM·s)					
<i>k</i> ₋₁	80-100 1/s	80-100 1/s					
<i>k</i> ₂	190 1/s	190 1/s					
k2	0.23 1/s	100 1/s					
k ₃	260 1/s	0.6 1/s					

(Rodnina's lab, Gottingen)

$$R_{C} \sim \frac{1}{e^{b_{1}} + e^{b_{2}} + e^{b_{3}}}$$
$$R_{W} \sim \frac{1}{e^{b_{1}} + e^{b_{2}} + e^{\delta + b_{3}}}$$

How to estimate recognition performance ("fitness")?

What is the actual dimension of the problem ?

Yonatan Savir

Recognition fitness has generic features

• "Fitness" *F* is often obscure and context-dependent:

→ look for generic properties of $F(R_C, R_W) = F(B, \delta, b_3)$.

• Only requirement: "biologically reasonable", $\frac{\partial F}{\partial R_C} \ge 0$, $\frac{\partial F}{\partial R_W} \le 0$.

$$R_C \sim \frac{1}{e^B + e^{b_3}}$$
$$R_W \sim \frac{1}{e^B + e^{\delta + b_3}}$$
$$(e^B = e^{b_1} + e^{b_2})$$



• Searching for optimum in (B, δ, b_3) space:

(i) $\frac{\partial F}{\partial \delta} \ge 0$: δ approaches biophysical limit.

(ii)
$$\left\{\frac{\partial F}{\partial B} = 0 \& \frac{\partial F}{\partial b_3} \ge 0\right\}$$
 or $\left\{\frac{\partial F}{\partial B} \le 0 \& \frac{\partial F}{\partial b_3} = 0\right\}$:

Optimization is essentially 1D (2 other parameters approach limit).

What is the optimal energy landscape of the ribosome?

• For example, distortion fitness from engineering (weight d is context-dependent) :

$$F = R_C - d \cdot R_W \propto \frac{1}{e^B + e^{b_3}} - \frac{d}{e^B + e^{b_3 + \delta}}$$

• 1D problem: optimum is along *b*₃.

(measured:
$$\Delta = b_3 - b_2$$
)

- What is the optimal b_3 (or Δ)?
- Is the ribosome optimal ?
- Role of conformational changes ?



Optimal design is a Max-Min strategy

- Weight *d* can vary.
- (i) For each d normalize F.
- (ii) "Worst case scenario": max(min(F)).





• Max-Min solution is "symmetric":

$$b_3 = -\frac{1}{2}\delta + B$$



Ribosome shows an energy barrier which is nearly optimal

- Measurements: $\Delta_C \approx -7 k_B T$, $\delta \approx 12 k_B T$, $\overline{B} = B b_2 \approx 1 k_B T$.
- <u>Prediction</u>: the optimal regime is **symmetric**, $\Delta_C < 0$, $\Delta_W > 0$.





• The ribosome is nearly optimal

(according to Max-Min prediction).

Decoding is optimal for all six measured tRNAs



• Except for UUC which encodes the same amino-acid

 ϕ (UUC) = ϕ (UUU)= phenylalanine.

Optimality is valid for wide range of fitness functions

• Ribosome optimal in wide region:

$$5 \cdot 10^{-6} = e^{-\delta} \le d \le e^{\delta} = 2 \cdot 10^5.$$

• **General feature**: any fitness function $F(R_C, R_W)$

exhibits optimum as long as both rates are "relevant".

$$e^{-\delta} < \left| \frac{\partial F}{\partial R_C} / \frac{\partial F}{\partial R_W} \right| < e^{\delta}$$

 $F = R_C - d \cdot R_W$





Theory predicts optimal regime of ribosomes for all organisms



- Optimal region in the space of all possible landscapes, $-\delta \leq \Delta_C \leq 0$.
- Mutations and antibiotics tend to push away of optimality.

What is the role of conformational changes?

• Energy barrier results from binding energy and deformation energy penalty:

$$\Delta_{C} = G_{\text{deform}} - G_{\text{bind}} = -\frac{1}{2}\delta + \overline{B}.$$

• Therefore

$$G_{\text{deform}} = G_{\text{bind}} - \frac{1}{2}\delta + \overline{B}.$$

• For any $G_{\text{bind}} \ge -\frac{1}{2}\delta + \overline{B} \approx 5 \text{ k}_{\text{B}}\text{T},$

 $G_{\text{deform}} \ge 0$ non-zero deformation is optimal for tRNA recognition.

Energy barrier that discerns the right target from competitors.


Recombination machinery recognizes homologous DNA

- Exchange between two *homologous* DNAs.
- Essential for:
 - Genome integrity (repair machinery).
 - Genetic diversity (crossover, sex).
- Task: Detect correct, homologous DNA target among many incorrect lookalikes.
- DNA stretches during recombination:

large deformation energy barrier.



Energy barriers for optimal recognition may be a general design principle of recognition systems with competition





Recombination optimizes extension energy of dsDNA.

Ribosome optimizes energy barriers of decoding

- **Conformational proofreading:** Design principle follows from optimization of information transfer function.
- May explain induced fit (Koshland 1958).
 Why molecules deform upon binding to target.



Relevant energy



Relevant energy

• Applies to any enzymatic kinetics in the presence of competition...

Open questions, future directions...

Understanding evolvable matter:

- What are the degrees-of-freedom underlying dimensional reduction? (Rubisco and other enzymes)
- Basic logic of molecular information channels (e.g. utilizing conformational changes, worst case scenario).

Translation machinery coevolved with proteins ightarrow

• Physics of the state of matter called "proteins"

(evolvable, mapped from DNA space, glassy dynamics).

Kinetic Proofreading

• The basic idea: iterations of irreversible discrimination step lead to exponential amplification.



$$Q_{0} \xrightarrow{K_{+}} Q_{1}$$

$$Q_{0} \xrightarrow{K_{+}} Q_{1}$$

$$Q_{0} \xrightarrow{K_{+}} Q_{2}$$

$$Q_{0} \xrightarrow{K_{+}} Q_{2}$$

$$Q_{N} \sim K^{-N} = \left(\frac{K_{+} + K_{-}}{K_{-}}\right)^{-N}$$

$$F = \frac{Q_{N}^{C}}{Q_{N}^{W}} = \left(\frac{K^{C}}{K^{W}}\right)^{-N} = \alpha^{N}$$

$$Q_{0} \xrightarrow{K_{+}} Q_{N-1}$$

$$K_{-}$$

QN

Hopfield (1974), Ninio (1975)

RecA dynamics exhibits multistage KPR



 $F = \alpha^{N}$

 $F = \alpha^{N^2/2}$

RecA filament strongly fluctuates

• Gradual depolymerization vs. polymerization jumps (similar to microtubules):

$$\frac{dp_n(t)}{dt} = -\kappa_{-}[p_n(t) - p_{n-1}(t)] + \kappa_{+}\left[\sum_{m=n+1}^{N} p_m(t) - np_n(t)\right]$$

• Continuum approximation $(p_n = P(n \text{ vacancies}), P_n = \sum_{m=n}^{N} p_m)$

$$\frac{\partial P(n,t)}{\partial t} = -\kappa_{-} \frac{\partial P(n,t)}{\partial n} - \kappa_{+} n P(n,t)$$

• Fluctuations "scan" the sequence and are therefore sensitive to mutations.



RecA dynamics is ultra-sensitive to DNA sequence

K.

ĸ

κ.

• At steady-state Gaussian amplification

$$P_s(n) = e^{-(\kappa_+/2\kappa_-)n^2}$$

• General result (Murugan, Huse & Leibler):

 $P \sim \exp(-\# \text{ loops}).$

• Can sense even single mutations:





Yonatan Savir

more: www.weizmann.ac.il\complex\tlusty

Early simulations of artificial "evolution"





				- E	1	15	1			5		1	2	1	1		- 1		-			- 1	_ 1			
				-	+	3	6	+	1.	2	7	2	-	-	-	7		1	-	-	+	+	-			
					-	L	12	⊢	11	5	1	-	-	-	4	21	-	-	-	+	_					
				5	1	Ľ	3	L	17	1		3	2	1	1		2	_			_	_				
			•	·]]		1	1		5	3	1	3	_	3	1		5	3		3		·				
				٦	T	5	3	1	15	-	3	1		5	3	1	3		3			5				
				13	11	11	1	tī	17		5	1	1	3	-	3	7		5	3		1				
				P ^a	17	1÷	⊢	Ě	17	1	5	<u>~</u>	i	-	-	-	-	1	1	-	7	÷t				
				F	12	1.	-	13	12	-	4		*	-	-	-	러	<u>+</u>	쒸	-+	쇍	+	_			
				L	12	12	1	3	1	<u>J</u> .	1		2	2	1	2	_	4	1	_	2	2	_			
				1	3		3	1		5	7	1	3		7	1		5	3	1	3	_	2			
				3	11	-	5	J	1	3		3	1		5	3	1	3		3	1	T	5			
	1			5	13	11	3	1	11	1		5	3	1	3		3	1	-	5	3	1	3			
				H	1ª	17	1	+	15	1	1	2		2	1	-		1	-	1	-	i	71			
				E	· L	12	Ľ	1	L.	5	Ľ.	4		-	-	_		픠		ě1		<u>x I</u>	-	ę., .		
										+										,						
			8																							
				З.									ċ													
-	-	1	_					-1		r÷r			Da	11	171	1	111	ħ		7-7	-T	-	71			
Ŧ	Ŧ	F	F	÷	П	7-	F	4	T	Ĥ	1,	5	, N	1	1	귀	2	À	Б	FI	Ŧ	+	FI		П	FI
Ŧ	Ŧ		E		H	Ŧ	F	4			5	2	¥.	1	1 3 5	2	2	1	1	Ð	Ŧ	+	Ę			E
Ŧ	Ŧ						5		5	3	5	3	1 2 3)	1 3 5 1	1	3 3 1				5		1		5	E
					5		5	2	5	3	5	2	1 2 3		1 3 5 3 1	3 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	35 311	1-31-1		1	5		1		5	
			5	37	5,	1	5 3	2	5 1 2 1	3	5	2	1 2 2 1	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	1 3 5 3 1 3	3 1 3	353113	1 3 1 1 3	1	1 3	5 3 7	5	1		5	
			5 3	31	5	2 1	5 3 1	3	5 3 3		5	1 1 3 5		11133	1 3 5 3 1 3	3 1 3 1 3 1	35313	1 3 . 1 3 5		1 3 3	5	5	3		5	
			5 21	31	5,		5 3 1 3	3	· · · · · · · · · · · · · · · · · · ·	1		3		1111511	1 3 5 3 1 3		35311311			1 3 3	5 1	5	3	5	5	
			5 3 1 3	<u> </u>	5		5 7 1 7	2	5 3 1 3 1 5 3 2 1 1 3 1	1		1 3 5 1		11113317	1 3 3 1 3	3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1	35 31 3 1 1			1 3 3 1 3		5 3 1	3	501	5 3	
					5 3 3 1 3 5 3	$\begin{array}{c} 5\\ \underline{j}\\ $	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\begin{array}{c c c c c c c c c c c c c c c c c c c $

Niels Aall Barricelli