

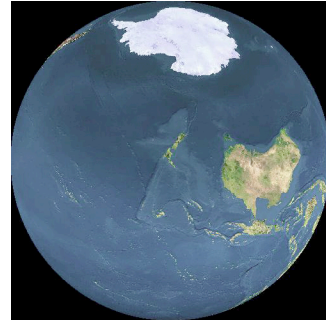
Genes, bioinformatics and dynamics

Peter R Wills

Department of Physics
University of Auckland

Supported by the Alexander von Humboldt Foundation

A different perspective



Metaphysics of science

The Greeks' question:

What remains the same through change?

Post-Einsteinian answer:

Energy/matter as *Urstoff*. Conserved quantity (undifferentiated).

The "gene"

Like an atom. An irreducible element of hereditary. A material particle, but an information carrier.

But what does a gene *do*?

The gene as a causal agent - a programmer.

The Century of the Gene Evelyn Fox Keller (2000)

- 1900 De Vries, Correns & Tschermak
'rediscover' Mendel's rules
- 1906 Bateson coins the term 'genetics'
- 1909 Johannsen uses the word 'gene' [unit factors, elements, allelomorphs in gametes]

The century of the gene

- 1943 Avery, McCleod & McCarty identify DNA as carrier of specificity in bacteria
- 1953 Watson & Crick's DNA structure
- 2000 Sequence of human genome

Schrödinger, 1944

"How are we to understand that [Hapsburger Lippe] has remained unperturbed by the disordering tendency of the heat motion for centuries?"

The genetic material as an "aperiod crystal", an idea of Delbrück's (1935)

Biological information – the old paradigm

Crick, 1958

Central dogma: DNA → RNA → protein

The Central Dogma ... states that once 'information' has passed into protein *it cannot get out again*. ... Information here means the *precise* determination of sequence, either of the bases in the nucleic acid or of amino acid residues in the protein.

Sequence hypothesis

The specificity of a piece of nucleic acid is expressed solely by the sequence of its bases, and this sequence is a (simple) code for the amino acid sequence of a particular protein.

A fuller picture

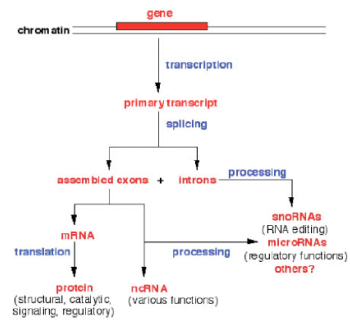
DNA → RNA → protein → function

The real genotype-phenotype mapping occurs in the last (implicit) step of biological information flow.

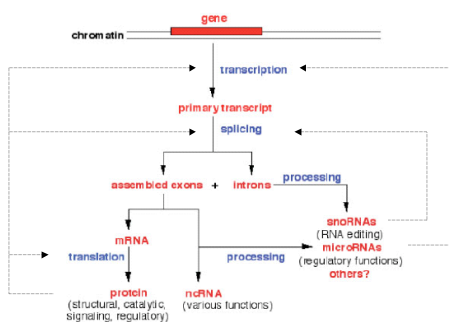
It is governed by:

- folding processes
- structure-function relationship

The new ncRNA paradigm



The new ncRNA paradigm



Bioinformatics

“We have coined the term **Bioinformatics** for the study of informatic processes in biotic systems. Our Bioinformatic approach typically involves spatial, multi-levelled models with many interacting entities whose behaviour is determined by local information.” (1978)

Prof. dr. Paulien Hogeweg
[Bioinformatics group](#), Utrecht University

Bioinformatics

The study of the

- generation
- transfer
- coding

of biological information

Importance of dynamics

Informational stability
in biological systems
depends on
dynamic stability.

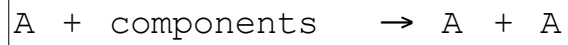
Importance of dynamics

Informational (in)stability
in biological systems
depends on
(thermo)dynamic (in)stability.

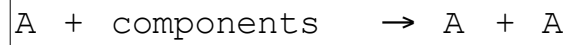
Informational (in)stability

hereditary transfer (mutation/selection)
gene expression (regulation)
translation (maintenance/control)
signalling (cascading)

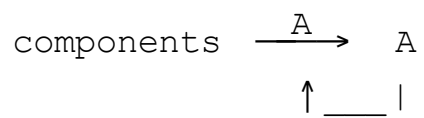
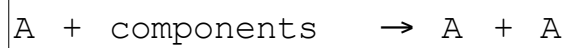
Replication/autocatalysis



Replication/autocatalysis



Replication/autocatalysis

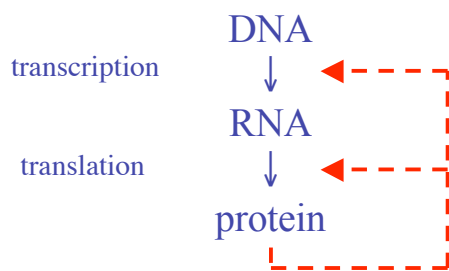


Replication/autocatalysis

Simplest dynamics for information transfer

nucleic acid replication; quasi-species selection

Autocatalysis in protein synthesis



Autocatalysis in protein synthesis -the information bootstrap problem

If (prebiotic) proteins are produced with more or less random sequences, how can the special “assignment catalysts” needed for a code be selected?

Restatement of the problem

You cannot make protein catalysts with **high assignment specificity** unless you can construct specific structures (amino acid sequences) accurately. For that you need **high assignment specificity**.

Why don't errors in translation lead to the catastrophic collapse of the process?

Binary nucleic acid to protein assignments (**code**).

x = codon; **y** = amino acid

K → **k** **L** → **k**

K → **l** **L** → **l**

Coding self-organization

In straightforward simulations of translation, the proteins that perform the two coding assignments ($\mathbf{K} \rightarrow \mathbf{k}$, $\mathbf{L} \rightarrow \mathbf{l}$) are selected and the proteins that perform the twelve “unwanted” assignments ($\mathbf{K} \rightarrow \mathbf{l}$, $\mathbf{L} \rightarrow \mathbf{k}$) die out.

HOW??

RNA-sequence-dependent protein synthesis

LLKKKKLKK | KLKKKKKKL



kllkllkll

kllkkkllk

Which assignments do these proteins catalyse?

“Reflexive” information

We are assuming that proteins are produced from RNA in a collinear fashion (codon frame is ‘read’ sequentially).

In simulations, the information provided is such that when it is translated by proteins that perform coding assignments, those very same coding-assignment catalysts are produced

– the process is then autocatalytic.

(But how is reflexive information selected?)

Catalytic centres

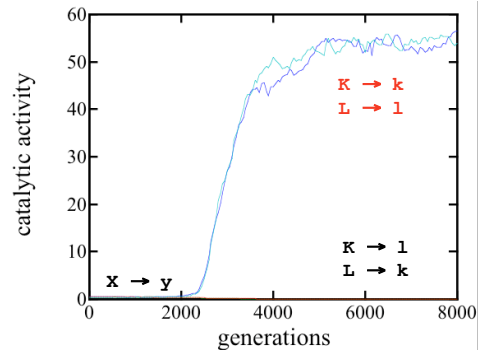
$\mathbf{K} \rightarrow \mathbf{k}$	llkkkkllkllk
$\mathbf{K} \rightarrow \mathbf{l}$	kkkllllkkll
$\mathbf{L} \rightarrow \mathbf{k}$	lkkllkllkkll
$\mathbf{L} \rightarrow \mathbf{l}$	kllkkkkllllk

In a random population of proteins, all assignments (of codons to amino acids) will be catalysed at equal rates and the further random sequences will be synthesised.

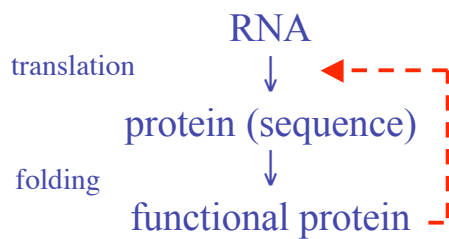
BUT

In the presence of reflexive information, this situation may dynamically unstable, (depending on the way in which catalytic function varies with protein sequence).

Selection of coding catalysts



Coding self-organisation



Problem 1

Where does the reflexive information come from? How can reflexive information be preserved from loss through mutation?

Ans. The coevolution of reflexive information and coding can be stabilised, including protection against parasitism, by co-localisation of nucleic acid information carriers and protein assignment catalysts.

Füchslin & McCaskill, 2002

Problem 2

Qu. Isn't it impossibly improbable that the reflexive information needed for a code as complicated as the one we have (64 codons, 20 amino acids) could be self-selecting because of dynamic instability in a mutually generating population of random nucleic acids and proteins?

Ans. Yes!

What then?

Nucleic acid and protein decomposition

$$K = \{A, B\}; L = \{B, C\}$$

$$k = \{a, b\}; l = \{c, d\}$$

$$X = \{A, B, C, D\}$$

$$y = \{a, b, c, d\}$$

For example ...

$K = \text{GYN, etc.}$ $L = \text{YRN, etc.}$
 $k = \text{hydrophobic}$ $l = \text{hydrophilic}$

$$K = \{A, B\}; L = \{C, D\}$$

$$k = \{a, b\}; l = \{c, d\}$$

$A = \text{GCN}; B = \text{other GYN, etc.}$
 $a = \text{ala}; b = \text{other hydrophobic}$

Binary nucleic acid to protein assignments (code).

$x = \text{codon}; y = \text{amino acid}$

$$K \rightarrow k \quad L \rightarrow k$$

$$K \rightarrow l \quad L \rightarrow l$$

Code decomposition

$A \rightarrow a$	$B \rightarrow a$	$C \rightarrow a$	$D \rightarrow a$
$K \rightarrow k$		$L \rightarrow k$	
$A \rightarrow b$	$B \rightarrow b$	$C \rightarrow b$	$D \rightarrow b$
$B \rightarrow c$	$B \rightarrow c$	$C \rightarrow c$	$D \rightarrow c$
$K \rightarrow l$		$L \rightarrow l$	
$B \rightarrow d$	$B \rightarrow d$	$C \rightarrow d$	$D \rightarrow d$

Coding assignment set

$A \rightarrow a$	$B \rightarrow a$	$C \rightarrow a$	$D \rightarrow a$
$A \rightarrow b$	$B \rightarrow b$	$C \rightarrow b$	$D \rightarrow b$
$B \rightarrow c$	$B \rightarrow c$	$C \rightarrow c$	$D \rightarrow c$
$B \rightarrow d$	$B \rightarrow d$	$C \rightarrow d$	$D \rightarrow d$

Catalytic centres (binary assignments)

$K \rightarrow k$	11kkkk1kk11k
$K \rightarrow l$	kkk11111kk11
$L \rightarrow k$	1kk1kk111kk1
$L \rightarrow l$	k1kkkkk1111k

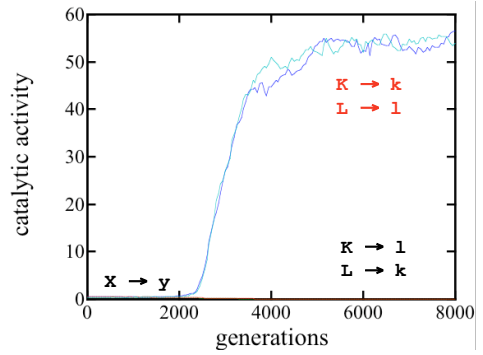
Decomposition

$K \rightarrow k$	11kkkk1kk11k
↓	
$A \rightarrow a$	d d b b a b c a a d c a
$A \rightarrow b$	c d b b b b c a b c d a
$B \rightarrow a$	d c a a a a d a a c c b
$B \rightarrow b$	c d a b a b d b a d c a

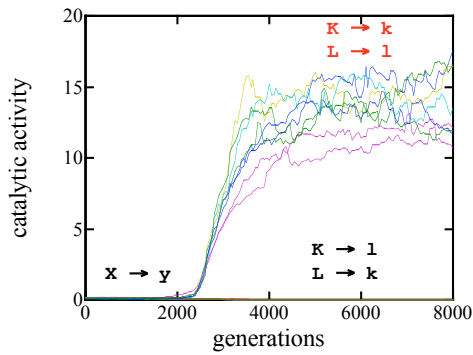
Catalytic centres (quarternary assignments)

A → a	d d b b a b c a a d c a	A → c	b a a c d d d d b a d c
A → b	c d b b b b c a b c d a	A → d	a b a c d d d c b b d c
B → a	d c a a a a d a a c c b	B → c	b b b c d c c c b a c c
B → b	c d a b a b d b a d c a	B → d	b a a d c d d c b b d c
C → a	d b a d a b c d d b b c	C → c	b c b a a b b b c c d b
C → b	d a b d a b d c d a b d	C → d	b c a b a b a a d d c b
D → a	d a a c b b c c c a a d	D → c	b d b b a b b a c c d b
D → b	c a b d b b d c d b a c	D → d	a c a b b a b a d d c a

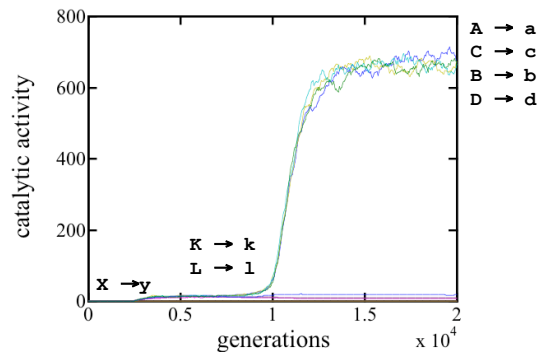
Selection of binary coding catalysts

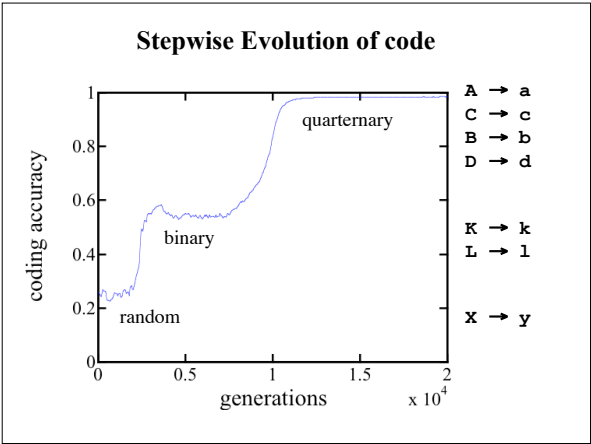
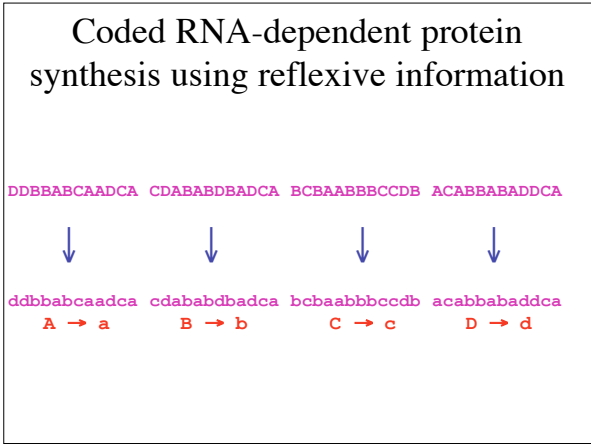
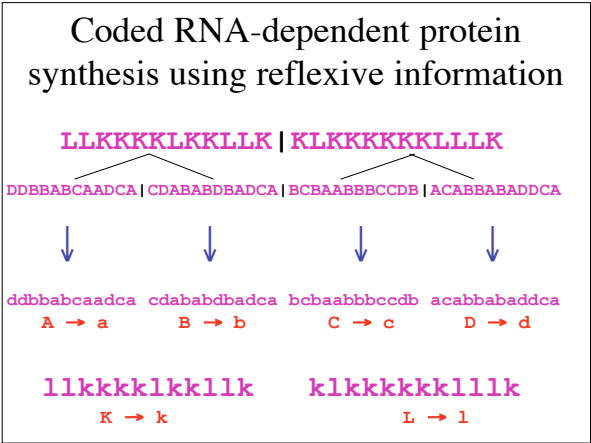
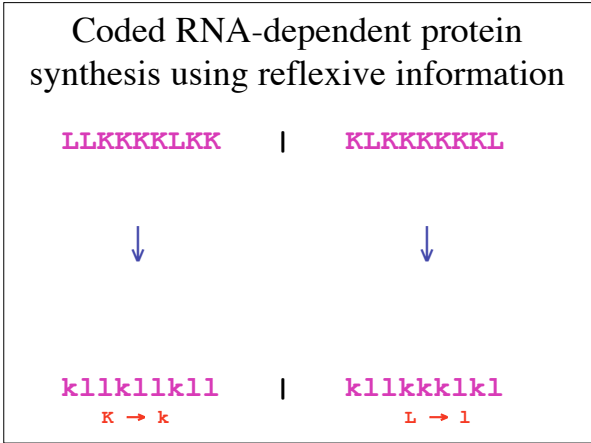


Selection of binary code catalysts

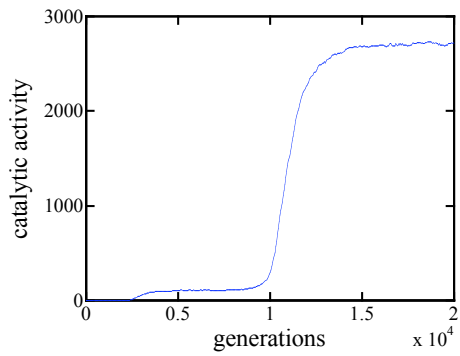


Transition to quarternary code

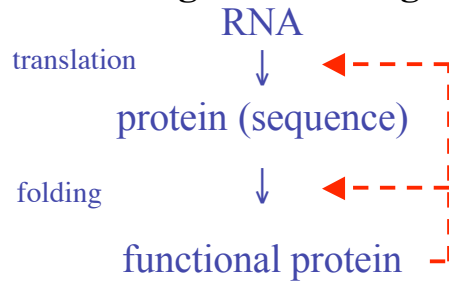




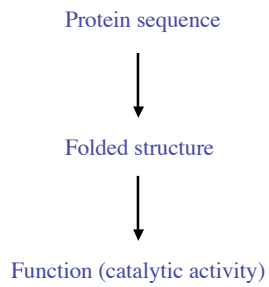
Thermodynamic driving force



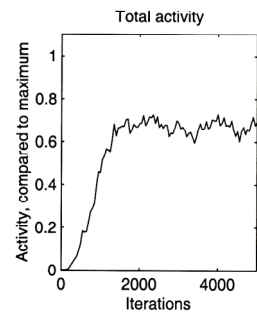
Coding self-organisation with guided folding



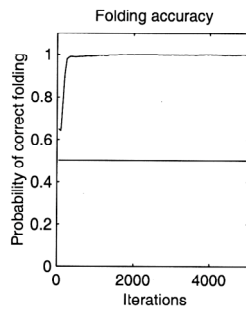
Folding process



Coding self-organisation



Simultaneous folding self-organisation



Susceptibility to parasites

Cooperative autocatalytic systems exhibit susceptibility to parasites that rely on the functionality of the whole system in order to reproduce but which contribute nothing to the operation of the system.

Folding acts as a kind of “scrambling” mechanism that requires parasites to become more and more specialised in order to reproduce (c.f. prions).

Dynamic determination

The biological “meaning” of polymeric sequence information depends on:

- coding
- folding

The operation of these processes depends on the maintenance of a state of self-organization. Thus the outcome of these processes is maintained and determined dynamically. It is not genomically encoded.

Contingency of biological function

Biological function is contingent on

- genetic information
- thermodynamic self-organization

This latter aspect has been generally ignored in bioinformatic studies:

– genetic phylogenies are established on the assumption that there is an unchanging “interpreter” that maps genotype onto phenotype

The real challenge for post-genomic bioinformatics!

To discover the how the molecular biological
“interpreter” has evolved.