# Darwin and Evolutionary Dynamics
## 150 Years After the ‚Origin of Species'

Peter Schuster

Institut für Theoretische Chemie, Universität Wien, Austria

and

The Santa Fe Institute, Santa Fe, New Mexico, USA

Evolution of Genomes and Origin of Species

Ohio State University, Columbus, 10.11.2008

Web-Page for further information:

http://www.tbi.univie.ac.at/~pks

1. Charles Darwins pathbreaking thoughts

2. Evolution without cellular life

3. Chemical kinetics of molecular evolution

4. Neutrality in replication

5. Modeling optimization of molecules

6. Complexity of biology

1. **Charles Darwins pathbreaking thoughts**

2. Evolution without cellular life

3. Chemical kinetics of molecular evolution

4. Neutrality in replication

5. Modeling optimization of molecules

6. Complexity of biology

Populations adapt to their environments through multiplication, variation, and selection – Darwins natural selection.

All forms of (terrestrial) life descend from one common ancestor – phylogeny and the tree of life.
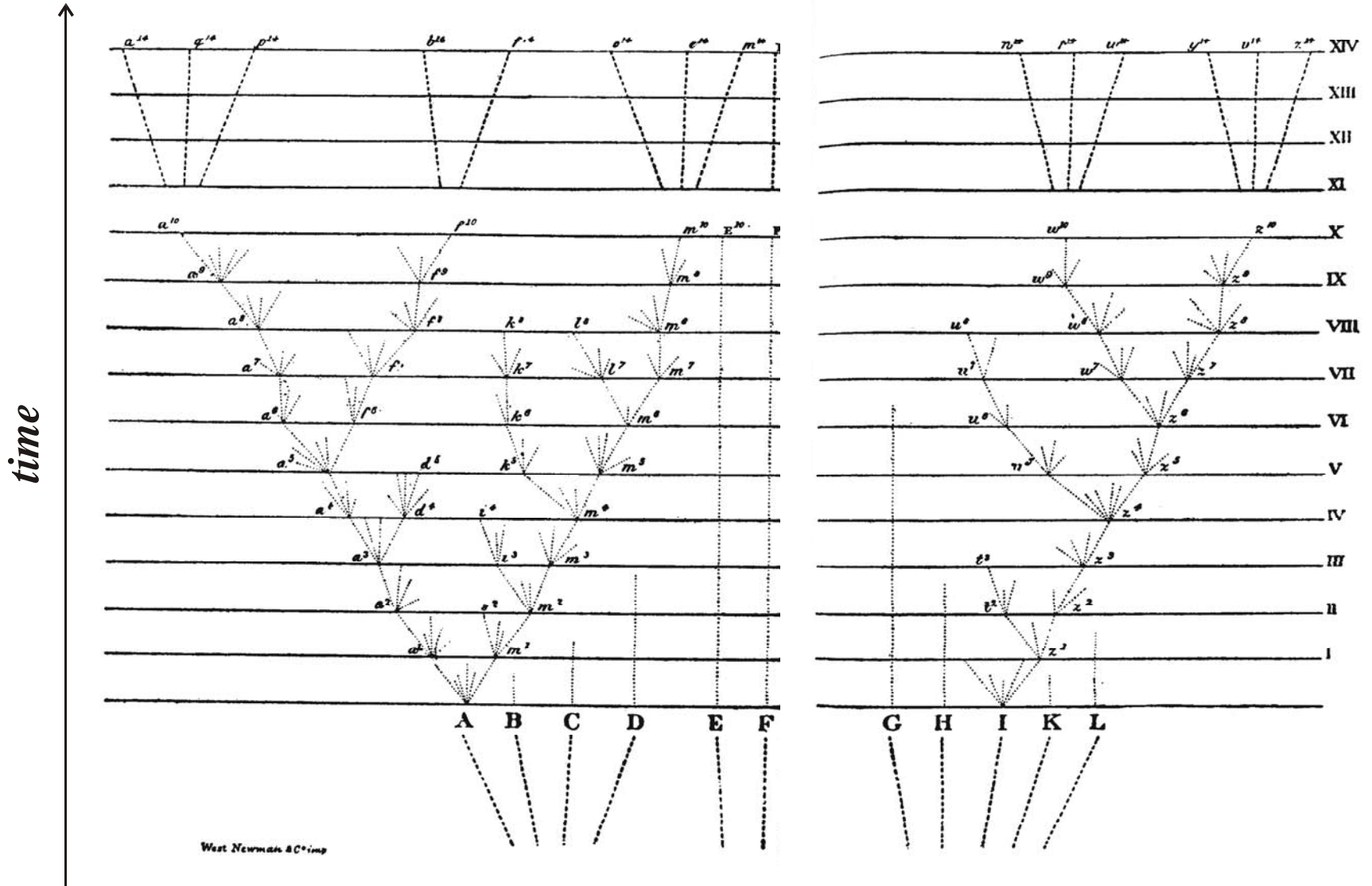
Three necessary conditions for Darwinian evolution are:

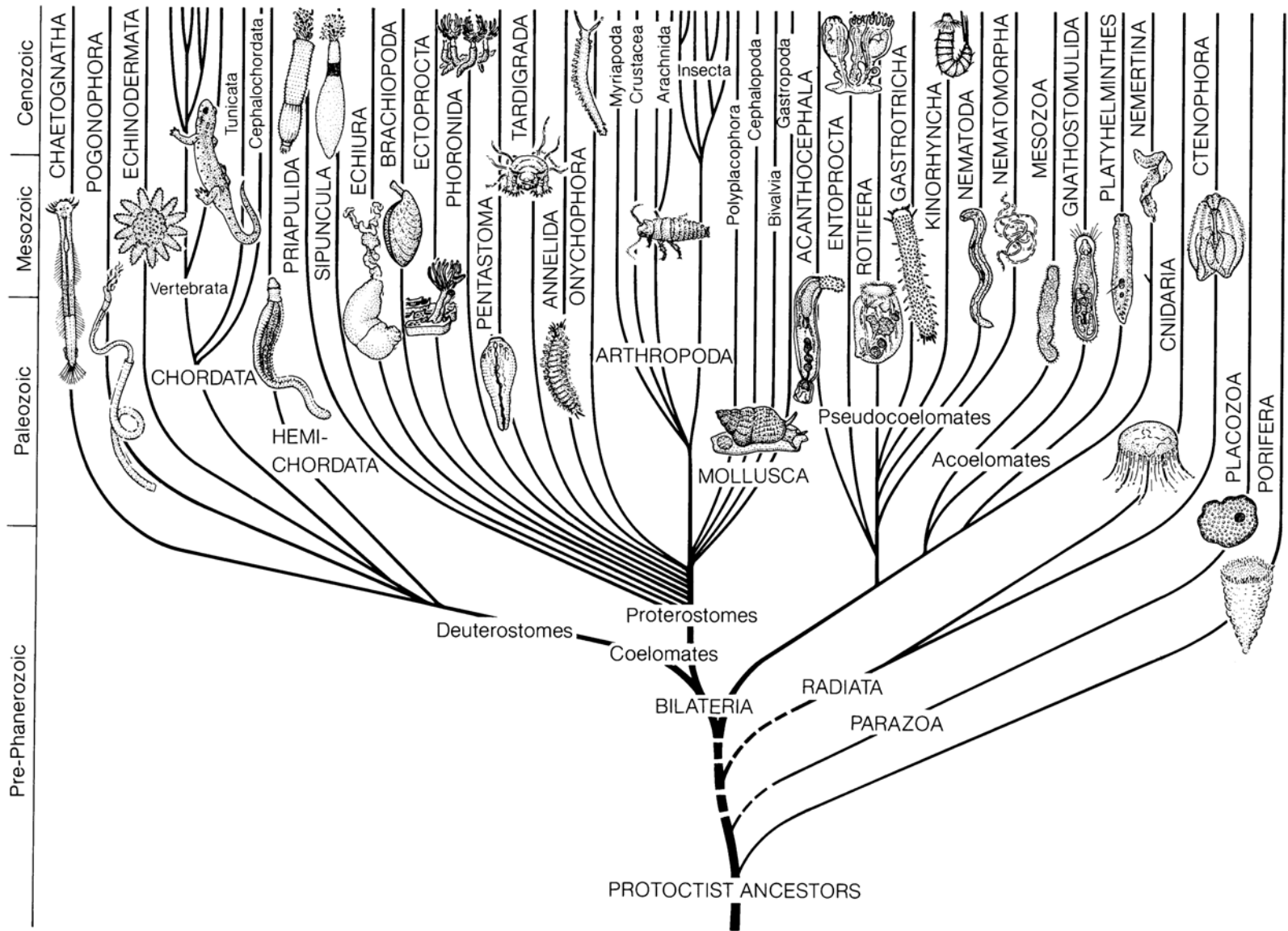1. **Multiplication,**

2. **Variation**, and

3. **Selection**.

Biologists distinguish the **genotype** – the genetic information – and the **phenotype** – the organisms and all its properties. The **genotype** is unfolded in development and yields the **phenotype**.

**Variation** operates on the **genotype** – through mutation and recombination – whereas the **phenotype** is the target of **selection**.

One important property of the Darwinian mechanism is that **variations** in the form of mutation or recombination events occur **uncorrelated** to their **effects** on the **selection** of the **phenotype**.

Charles Darwin, *The Origin of Species*, 6th edition.
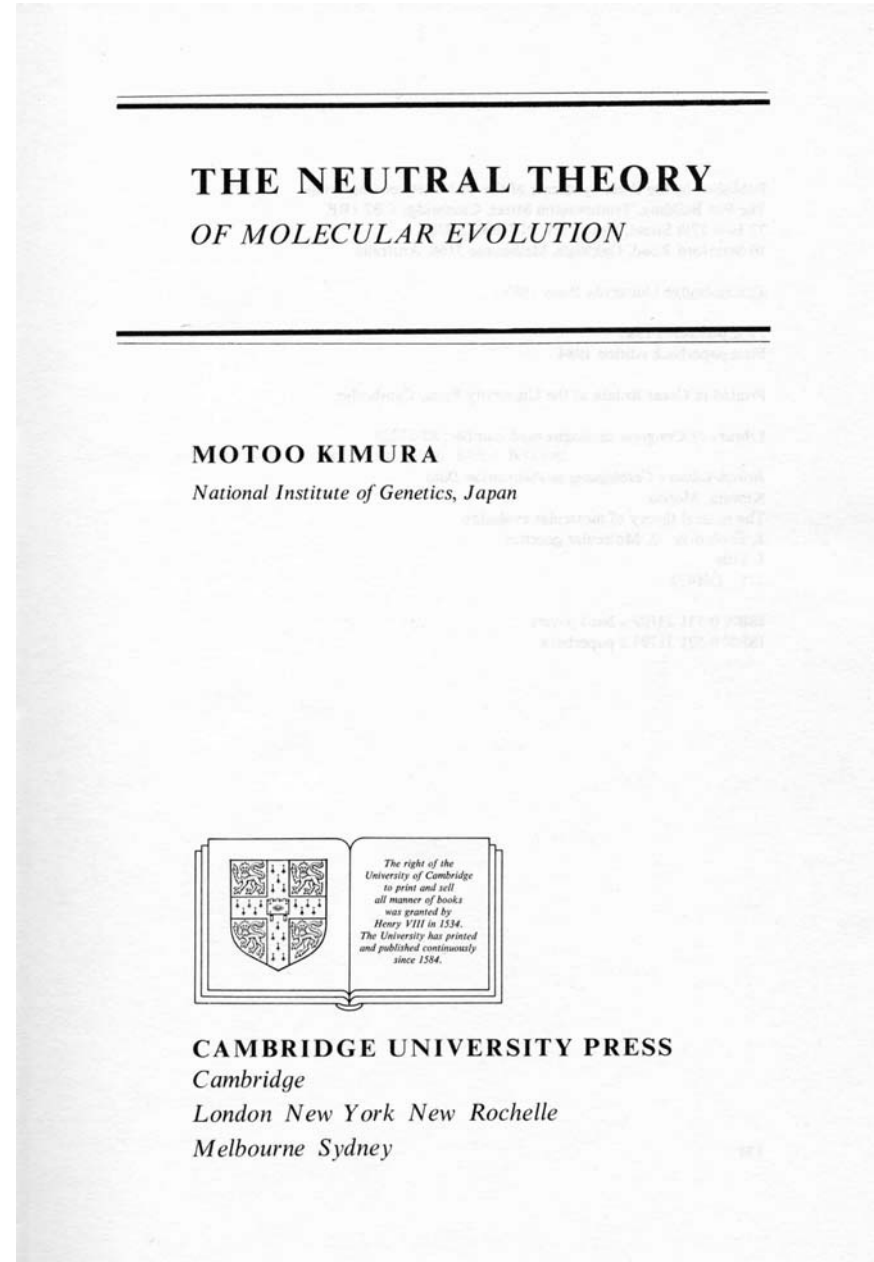Everyman's Library, Vol.811, Dent London, pp.121-122.

Modern phylogenetic tree: Lynn Margulis, Karlene V. Schwartz. *Five Kingdoms. An Illustrated Guide to the Phyla of Life on Earth*. W.H. Freeman, San Francisco, 1982.

MOTOO KIMURA



**THE NEUTRAL THEORY**
*OF MOLECULAR EVOLUTION*

**MOTOO KIMURA**
National Institute of Genetics, Japan

*The right of the
University of Cambridge
to print and sell
all manner of books
was granted by
Henry VIII in 1534.
The University has printed
and published continuously
since 1584.*

**CAMBRIDGE UNIVERSITY PRESS**
*Cambridge*
*London New York New Rochelle*
*Melbourne Sydney*

Motoo Kimuras population genetics of neutral evolution.

Evolutionary rate at the molecular level. *Nature* **217**: 624-626, 1955.

*The Neutral Theory of Molecular Evolution*. Cambridge University Press. Cambridge, UK, 1983.
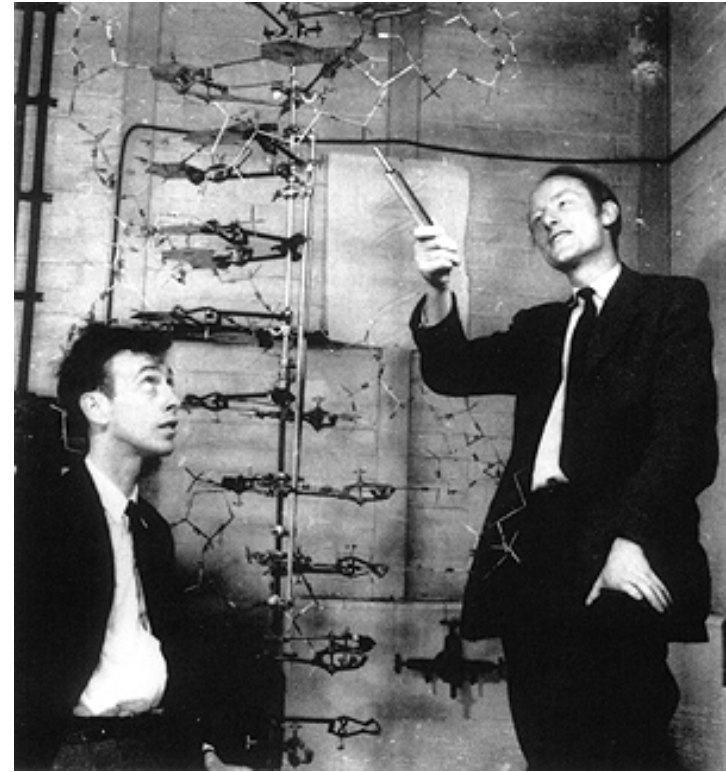
The molecular clock of evolution



Fig. 4.2. Percentage amino acid differences when the α hemoglobin chains are compared among eight vertebrates together with their phylogenetic relationship and the times of divergence.

Motoo Kimura. *The Neutral Theory of Molecular Evolution*. Cambridge University Press. Cambridge, UK, 1983.

The three-dimensional structure of a
short double helical stack of B-DNA



James D. Watson, 1928-, and Francis H.C. Crick, 1916-2004

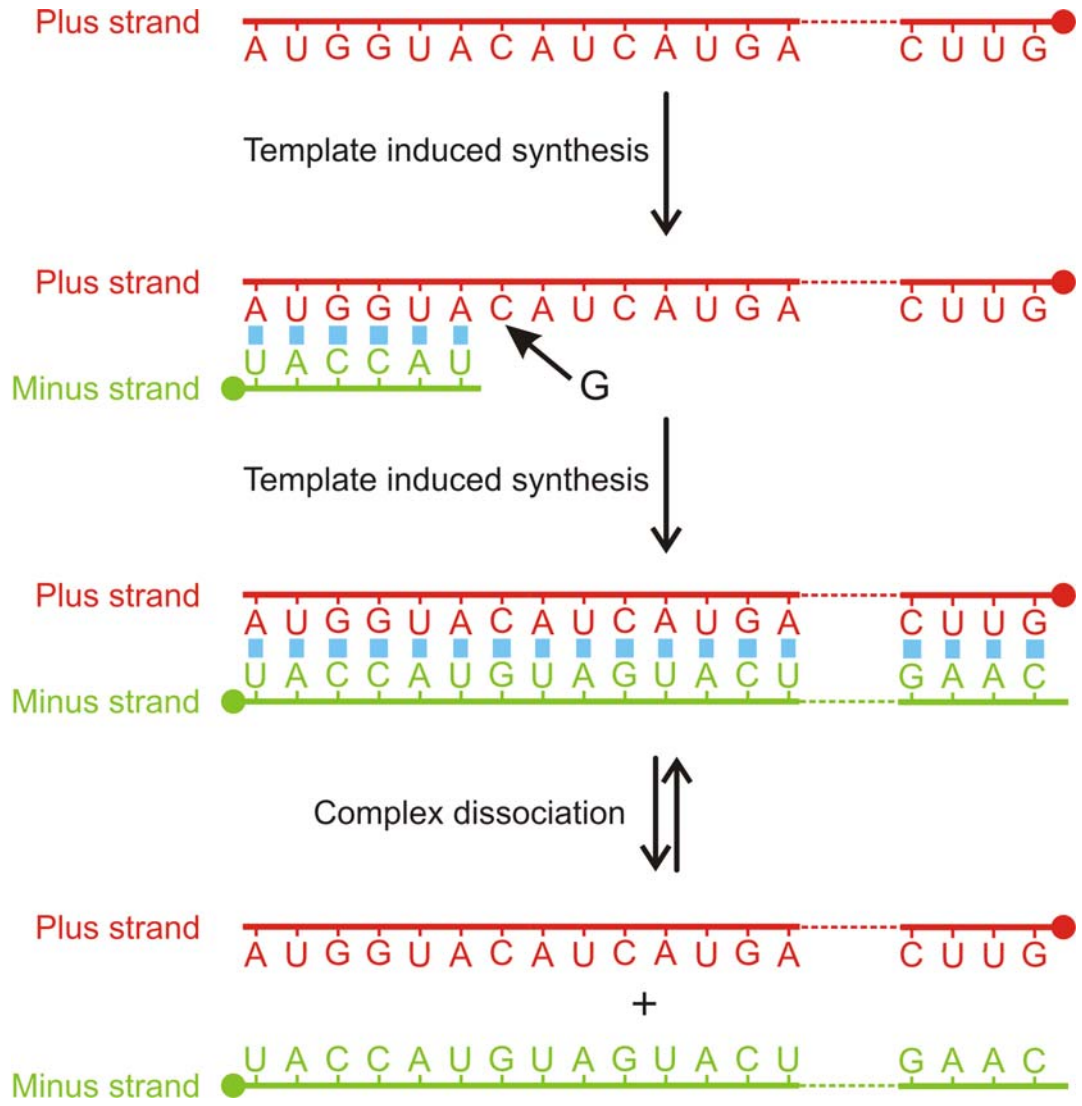Nobel prize 1962

**1953 – 2003  fifty years double helix**

The geometry of the double helix is compatible
only with the base pairs:

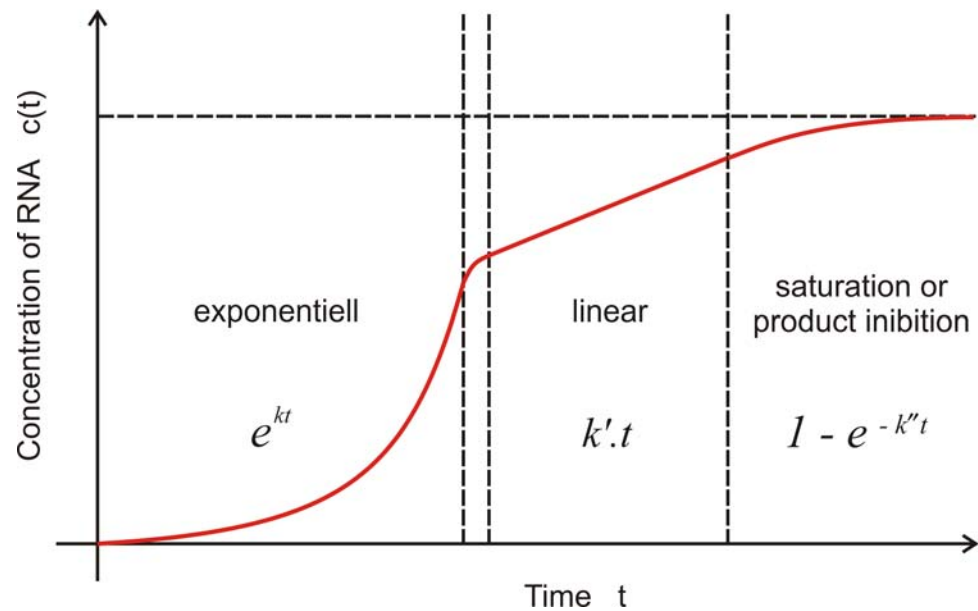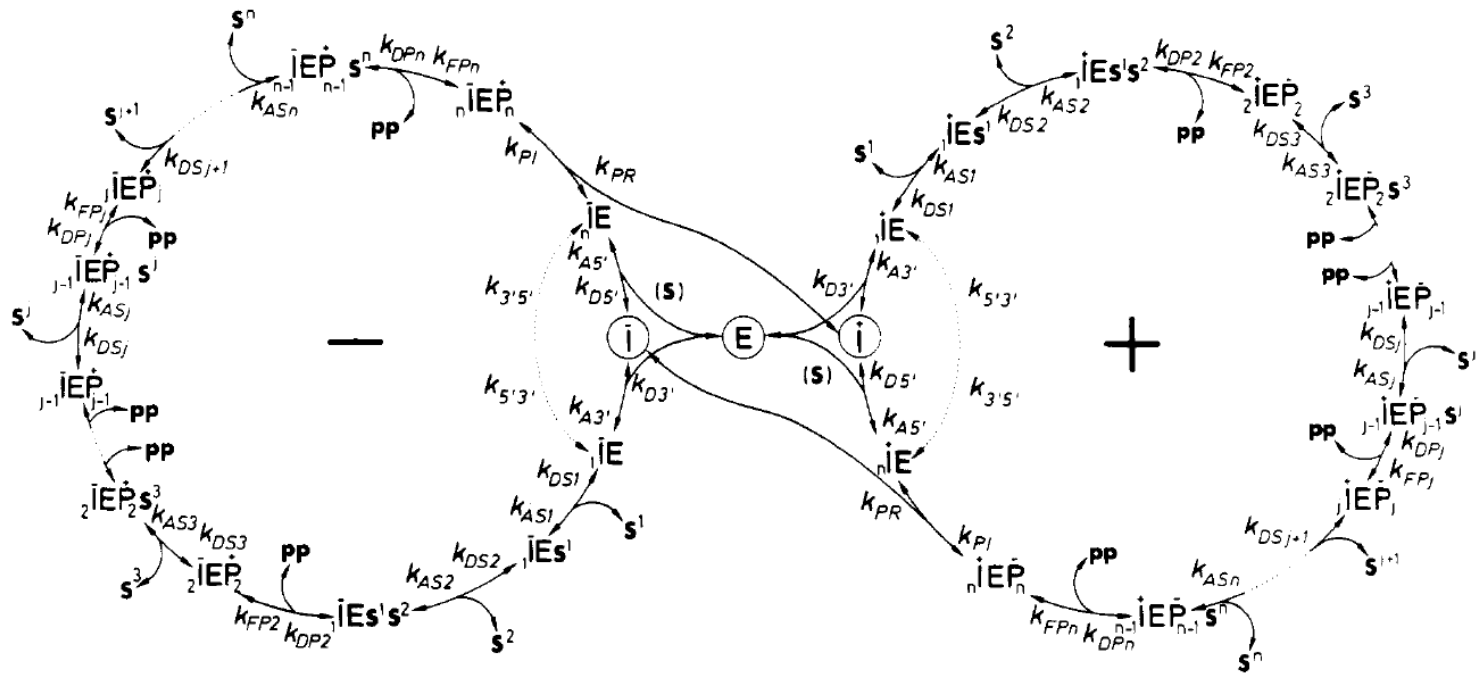**AT**, **TA**, **CG**, and **GC**

,Replication fork' in DNA replication

The mechanism of DNA replication is ,semi-conservative'

Plus strand

AUGGUACAUCAUGA CUUG

Template induced synthesis

Plus strand

AUGGUACAUCAUGA CUUG
UACCAU
G

Template induced synthesis

Plus strand

AUGGUACAUCAUGA CUUG
UACCAUGUAGUACU GAAC

Minus strand

Complex dissociation

Plus strand

AUGGUACAUCAUGA CUUG

+

Minus strand

UACCAUGUAGUACU GAAC

Complementary replication is the simplest copying mechanism of RNA.
Complementarity is determined by Watson-Crick base pairs:

**G≡C** and **A=U**
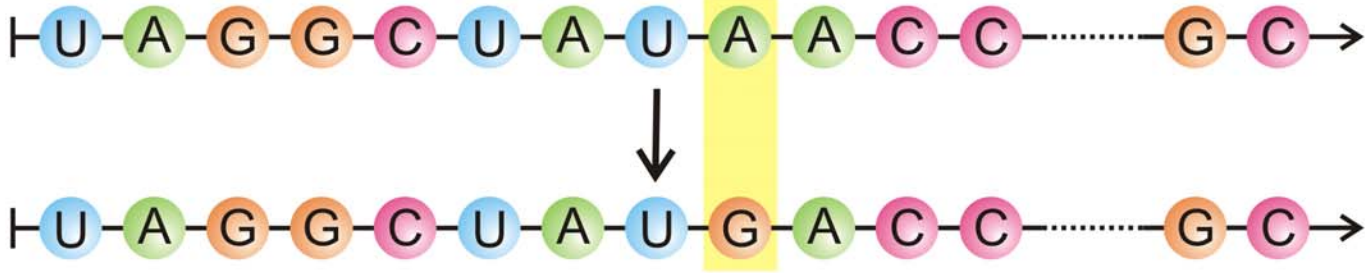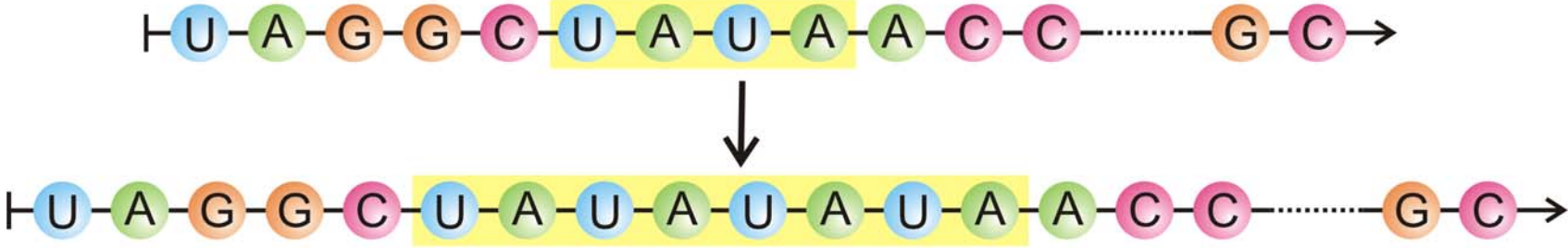
Kinetics of RNA replication

C.K. Biebricher, M. Eigen, W.C. Gardiner, Jr.
*Biochemistry* **22**:2544-2559, 1983

Concentration of RNA c(t)

exponentiell
$e^{kt}$

linear
$k'.t$

saturation or
product inibition
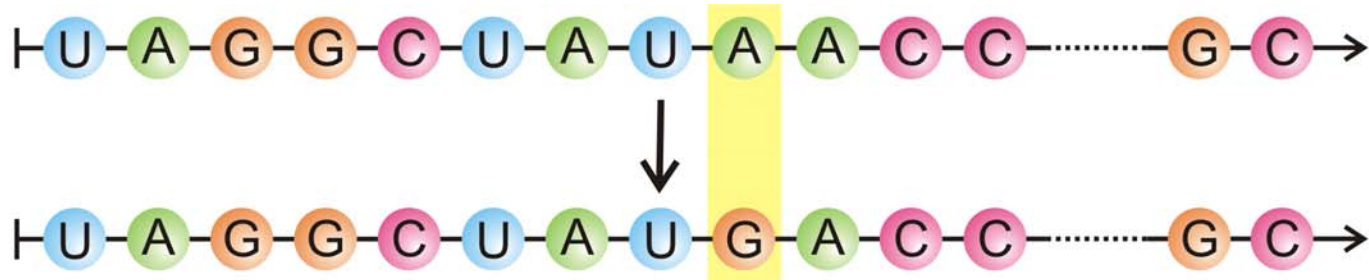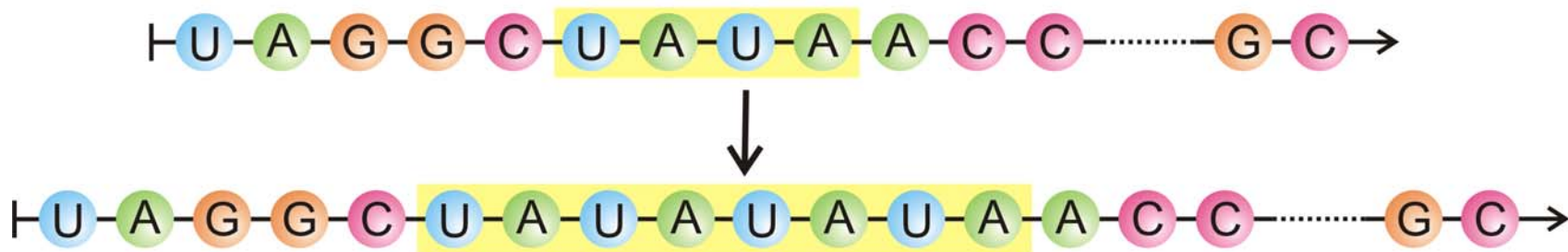$1 - e^{-k''t}$

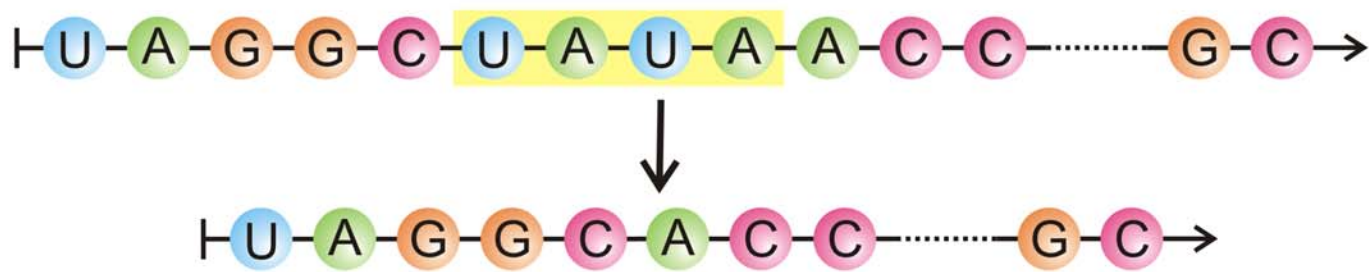Time t

Punktmutation

Punktmutation

Insertion

Punktmutation

Insertion

Deletion

# Evolution of RNA molecules based on Qβ phage

D.R.Mills, R.L.Peterson, S.Spiegelman, *An extracellular Darwinian experiment with a self-duplicating nucleic acid molecule*. Proc.Natl.Acad.Sci.USA **58** (1967), 217-224

S.Spiegelman, *An approach to the experimental analysis of precellular evolution*. Quart.Rev.Biophys. **4** (1971), 213-253
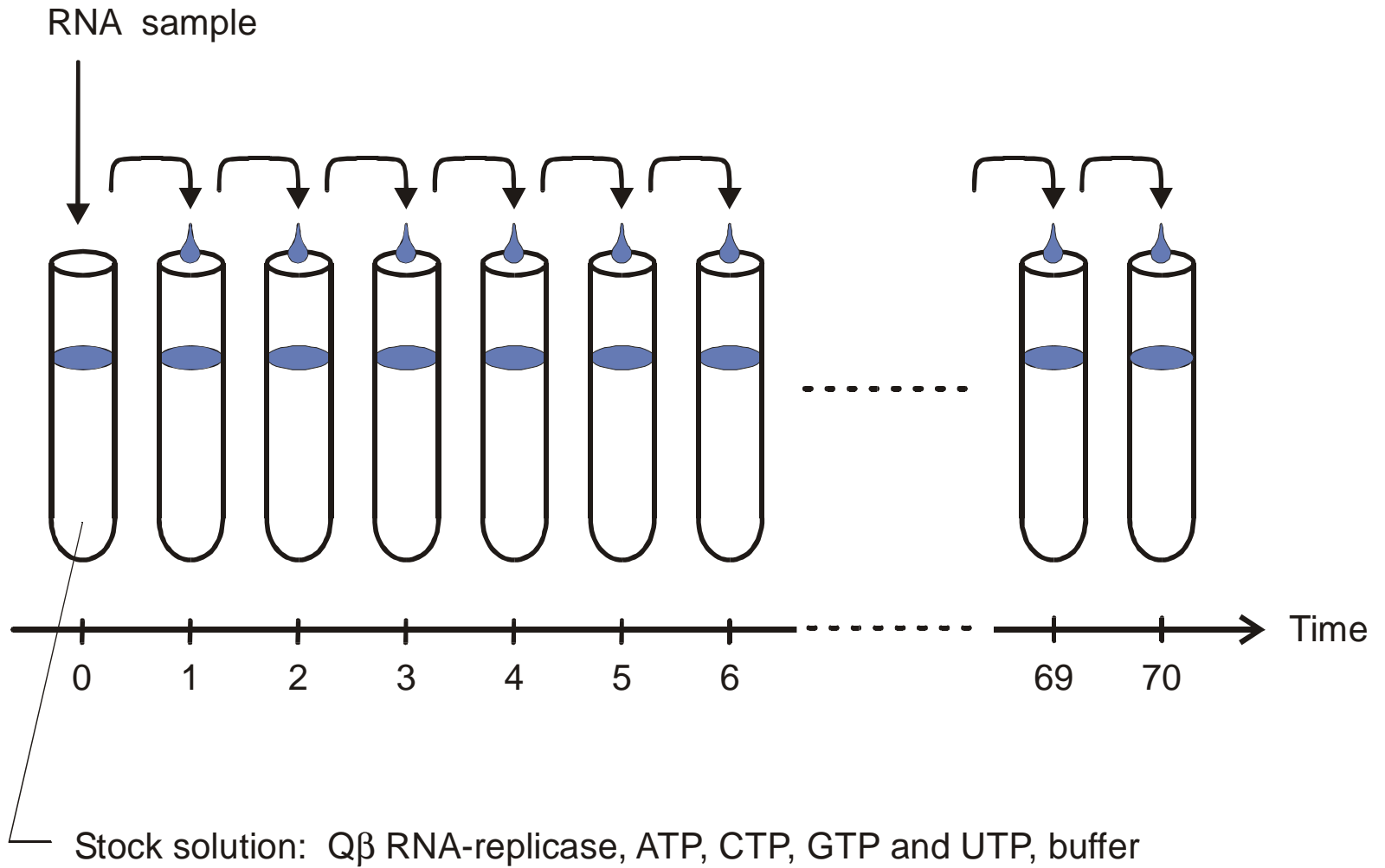
C.K.Biebricher, *Darwinian selection of self-replicating RNA molecules*. Evolutionary Biology **16** (1983), 1-52

G.Bauer, H.Otten, J.S.McCaskill, *Travelling waves of in vitro evolving RNA.* *Proc.Natl.Acad.Sci.USA* **86** (1989), 7937-7941

C.K.Biebricher, W.C.Gardiner, *Molecular evolution of RNA* **in vitro**. Biophysical Chemistry **66** (1997), 179-192

G.Strunk, T.Ederhof, *Machines for automated evolution experiments* **in vitro** *based on the serial transfer concept*. Biophysical Chemistry **66** (1997), 193-202

F.Öhlenschlager, M.Eigen, *30 years later – A new approach to Sol Spiegelman's and Leslie Orgel's* **in vitro** *evolutionary studies*. Orig.Life Evol.Biosph. **27** (1997), 437-457
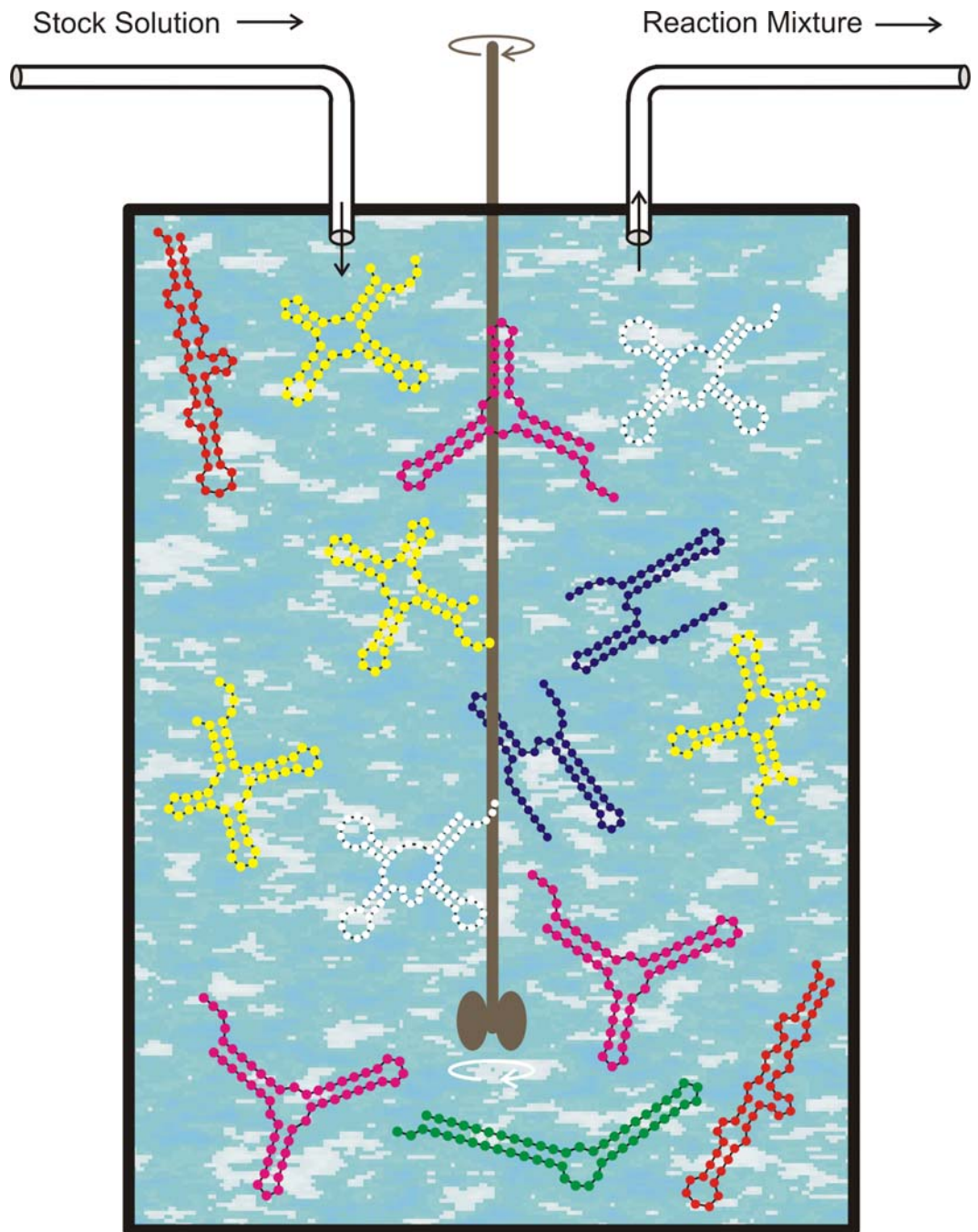
Application of serial transfer to RNA evolution in the test tube

**Stock solution**:

activated monomers, **ATP, CTP, GTP, UTP (TTP);**
a replicase, an enzyme that performs complemantary replication;
buffer solution

The flowreactor is a device for **studies** of evolution *in vitro* and *in silico.*

Stock Solution ⟶

Reaction Mixture ⟶

# Evolutionary design of RNA molecules

A.D. Ellington, J.W. Szostak, **In vitro** *selection of RNA molecules that bind specific ligands*. Nature **346** (1990), 818-822
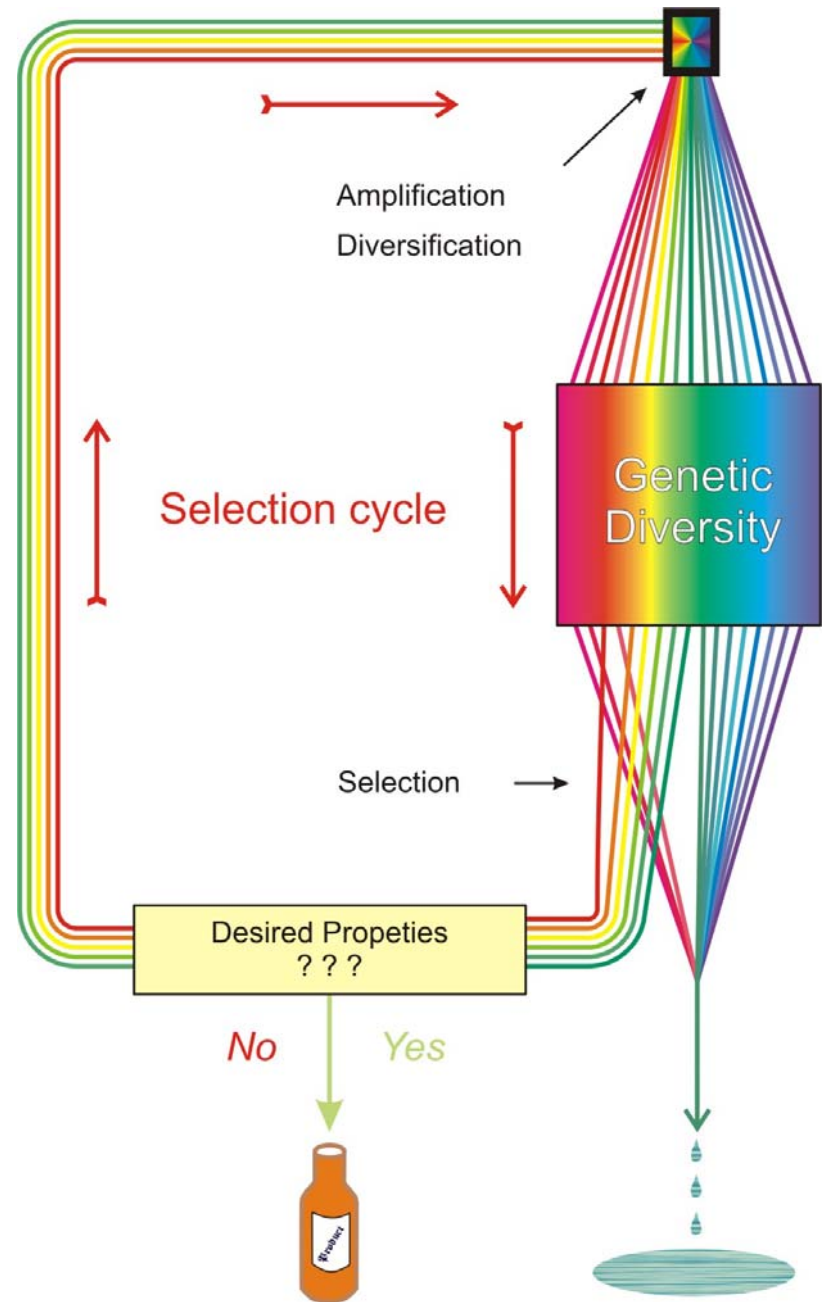
C. Tuerk, L. Gold, **SELEX** - *Systematic evolution of ligands by exponential enrichment: RNA ligands to bacteriophage* **T4** *DNA polymerase*. Science **249** (1990), 505-510

D.P. Bartel, J.W. Szostak, *Isolation of new ribozymes from a large pool of random sequences*. Science **261** (1993), 1411-1418
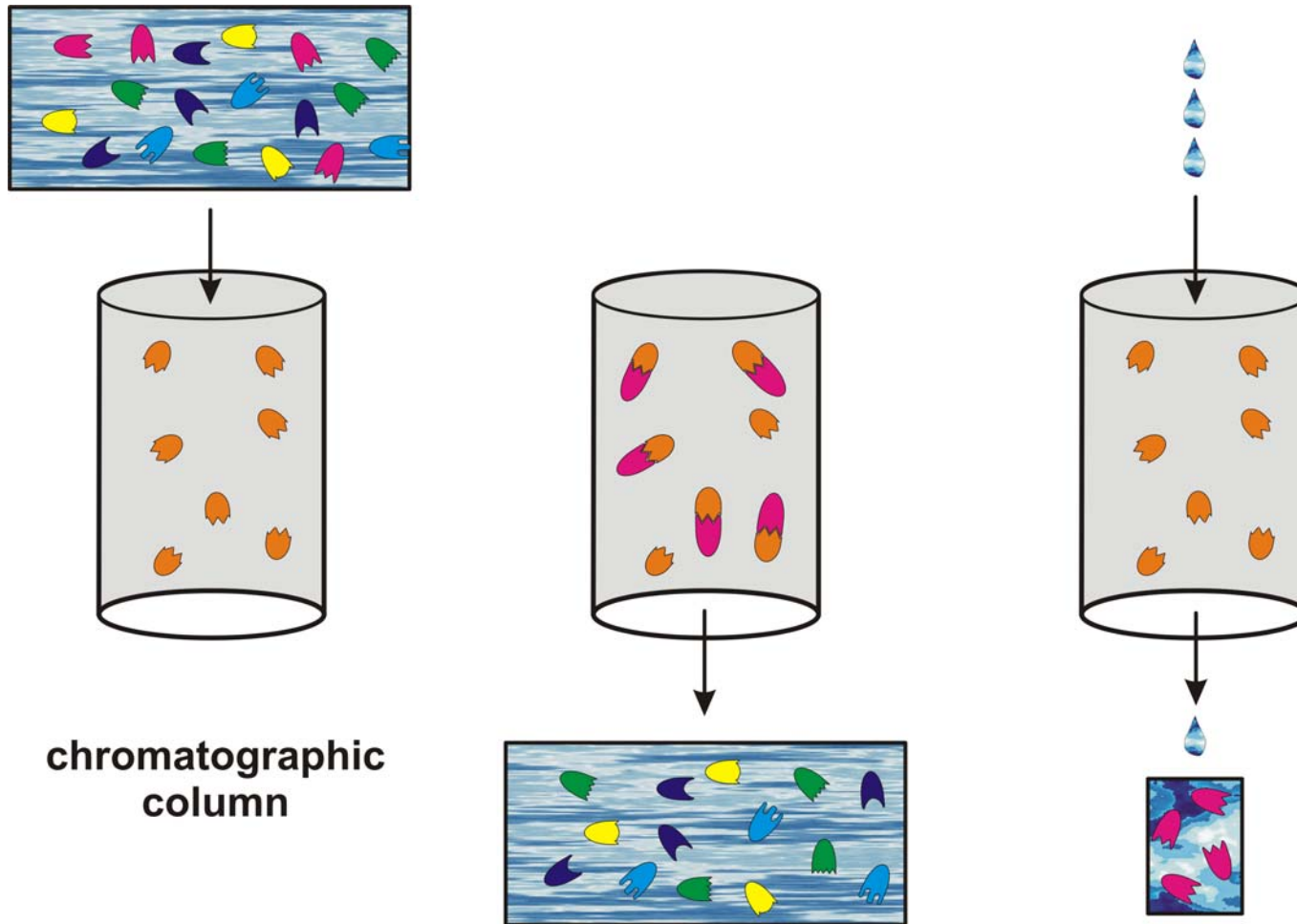
R.D. Jenison, S.C. Gill, A. Pardi, B. Poliski, *High-resolution molecular discrimination by RNA*. Science **263** (1994), 1425-1429

Y. Wang, R.R. Rando, *Specific binding of aminoglycoside antibiotics to RNA*. Chemistry & Biology **2** (1995), 281-290
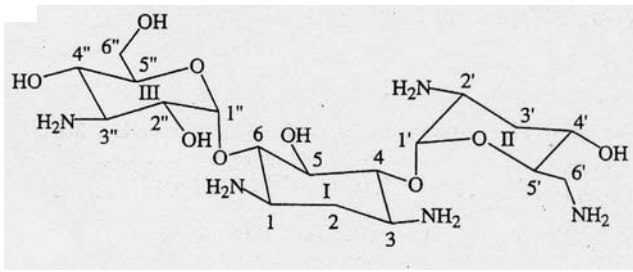
L. Jiang, A. K. Suri, R. Fiala, D. J. Patel, *Saccharide-RNA recognition in an aminoglycoside antibiotic-RNA aptamer complex*. Chemistry & Biology **4** (1997), 35-50
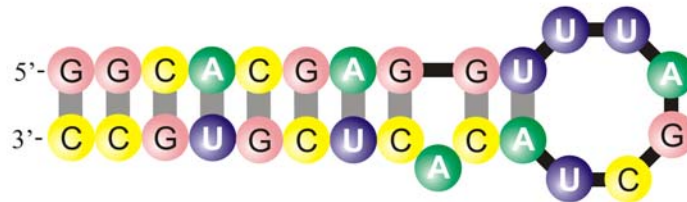
Amplification
Diversification

Selection cycle

Genetic Diversity

Selection

Desired Propeties
? ? ?

No        Yes

An example of 'artificial selection' with RNA molecules or 'breeding' of biomolecules

chromatographic
column

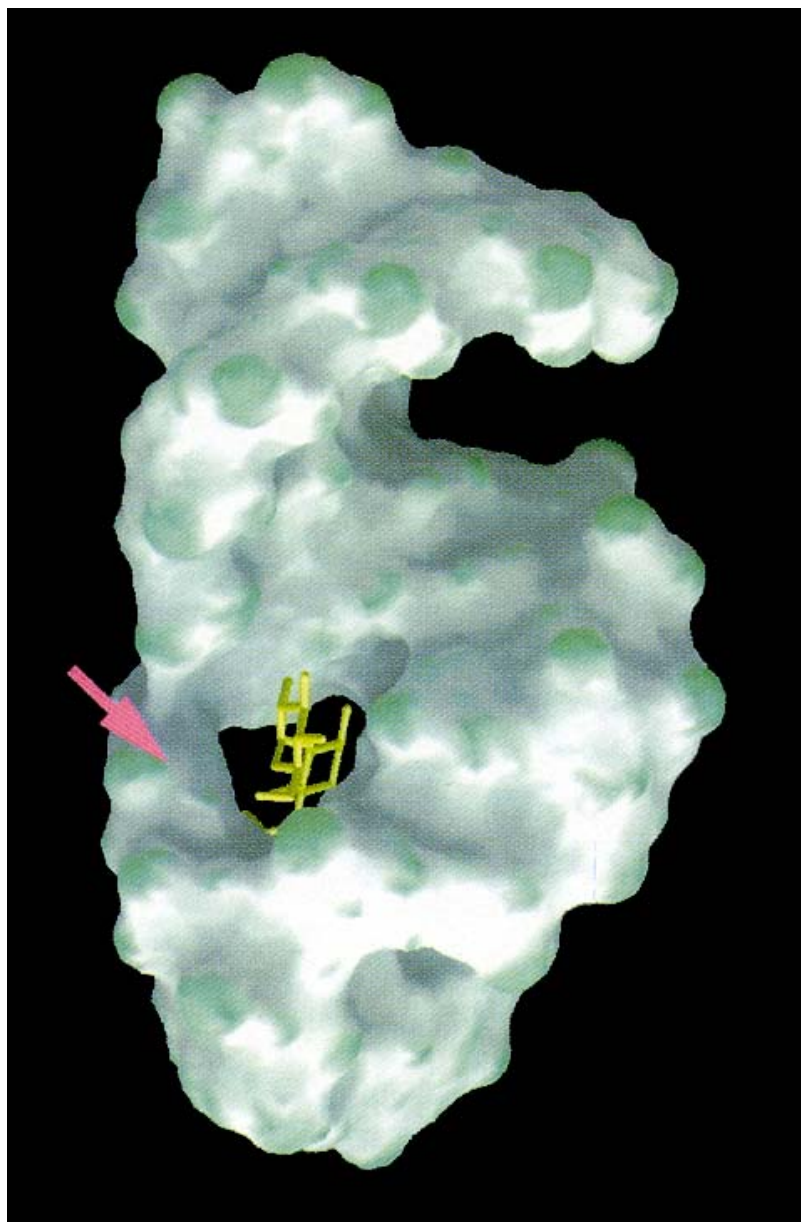Die SELEX-Technik zur evolutionären Erzeugung von stark bindenden Molekülen
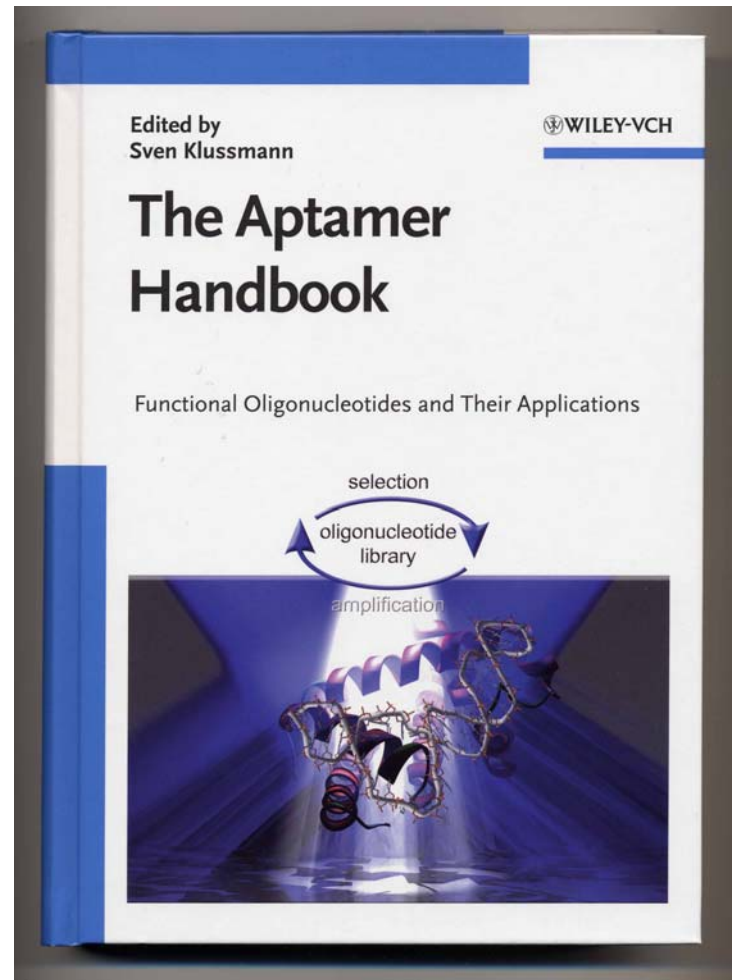
tobramycin

RNA aptamer, $n = 27$
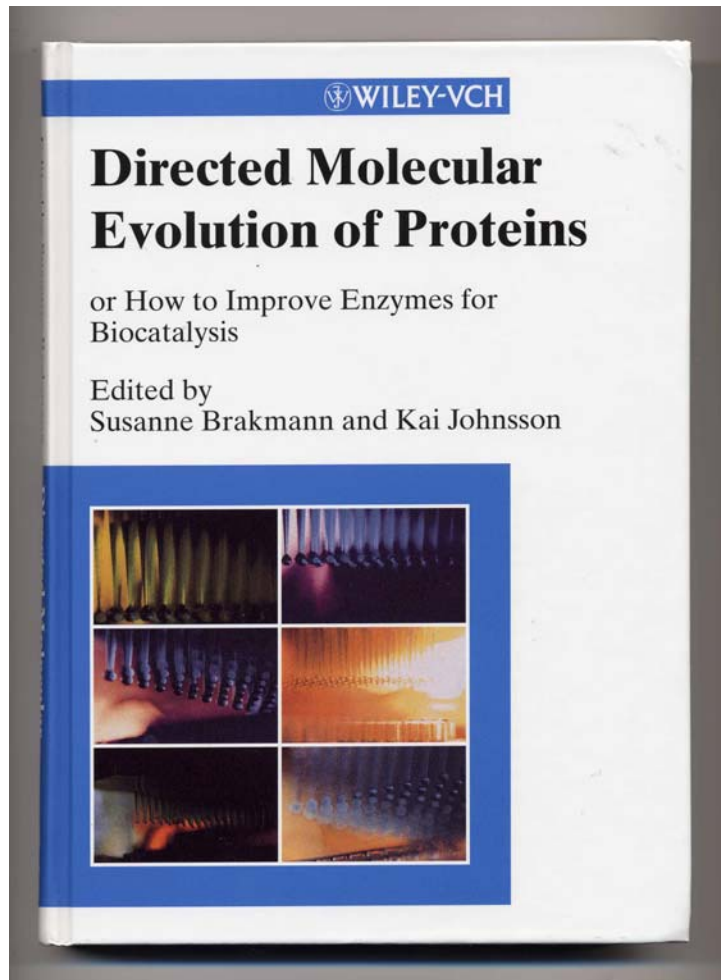
Formation of secondary structure of the tobramycin binding RNA aptamer with  $K_D = 9$ nM

L. Jiang, A. K. Suri, R. Fiala, D. J. Patel, *Saccharide-RNA recognition in an aminoglycoside antibiotic-RNA aptamer complex.*  Chemistry & Biology **4**:35-50 (1997)

The three-dimensional structure of the tobramycin aptamer complex

L. Jiang, A. K. Suri, R. Fiala, D. J. Patel, Chemistry & Biology **4**:35-50 (1997)

Application of molecular evolution to problems in biotechnology

**Artificial evolution in biotechnology and pharmacology**

G.F. Joyce. 2004. Directed evolution of nucleic acid enzymes. *Annu.Rev.Biochem.* **73**:791-836.

C. Jäckel, P. Kast, and D. Hilvert. 2008. Protein design by directed evolution. *Annu.Rev.Biophys.* **37**:153-173.

S.J. Wrenn and P.B. Harbury. 2007. Chemical evolution as a tool for molecular discovery. *Annu.Rev.Biochem.* **76**:331-349.

**Results from evolution experiments**:

• Replication of RNA molecules *in vitro* gives rise to exponential growth under suitable conditions.

•Evolutionary optimization does not require cells and occurs as well in cell-free molecular systems.

•*In vitro* evolution allows for production of molecules for predefined purposes and gave rise to a branch of biotechnology.

1. Charles Darwins pathbreaking thoughts

2. Evolution without cellular life

3. **Chemical kinetics of molecular evolution**

4. Neutrality in replication

5. Modeling optimization of molecules

6. Complexity of biology

DIE NATURWISSENSCHAFTEN

58. Jahrgang, 1971                                          Heft 10 Oktober

Selforganization of Matter
and the Evolution of Biological Macromolecules

Manfred Eigen*

Max-Planck-Institut für Biophysikalische Chemie,
Karl-Friedrich-Bonhoeffer-Institut, Göttingen-Nikolausberg

1971

Die Naturwissenschaften    64. Jahrgang    Heft 11    November 1977

The Hypercycle

A Principle of Natural Self-Organization

Part A: Emergence of the Hypercycle

Manfred Eigen
Max-Planck-Institut für biophysikalische Chemie, D-3400 Göttingen

Peter Schuster
Institut für theoretische Chemie und Strahlenchemie der Universität, A-1090 Wien

1977

Molecular Quasi-Species†

Manfred Eigen,* John McCaskill,
Max Planck Institut für biophysikalische Chemie, Am Fassberg, D 3400 Göttingen-Nikolausberg, BRD

and Peter Schuster*
Institut für theoretische Chemie und Strahlenchemie, der Universität Wien, Währinger Strasse 17,
A-1090 Wien, Austria (Received: June 9, 1988)

1988

Chemical kinetics of molecular evolution

$$(A) + I_1 \xrightarrow{f_1} I_2 + I_1$$

$$(A) + I_2 \xrightarrow{f_2} I_1 + I_2$$

$$\frac{dx_1}{dt} = f_2\, x_2 \quad \text{and} \quad \frac{dx_2}{dt} = f_1\, x_1$$

$$x_1 = \sqrt{f_2}\,\xi_1 \,, \quad x_2 = \sqrt{f_1}\,\xi_2 \,, \quad \zeta = \xi_1 + \xi_2 \,, \quad \eta = \xi_1 - \xi_2 \,, \quad f = \sqrt{f_1 f_2}$$

$$\eta(t) = \eta(0)\, e^{-ft}$$

$$\zeta(t) = \zeta(0)\, e^{ft}$$

Complementary replication as the simplest molecular mechanism of reproduction

Chemical kinetics of replication and mutation as parallel reactions

$$\frac{dc_i}{dt} = \sum_{j=1}^{N} Q_{ij}\, f_j\, c_j\,; \quad i = 1, 2, \ldots, N$$

$$\frac{d\mathbf{c}}{dt} = \mathrm{W} \cdot \mathbf{c}\,; \quad \sum_{1=1}^{N} c_i(t) = c(t)\,; \quad \mathrm{W} = \{W_{ij} \doteq Q_{ij}\, f_j\}$$

Normalization

$$x_i = c_i / c\,; \quad \sum_{i_1}^{n} x_i = 1$$

$$\frac{d\mathbf{x}}{dt} = \mathrm{W} \cdot \mathbf{x} - \bar{f}\,\mathbf{x} = (\mathrm{G} \cdot \mathrm{F} - \bar{f}\,\mathbb{E}) \cdot \mathbf{x}\,; \quad \bar{f} = \sum_{i=1}^{N} x_i\, f_i$$

# Decomposition of matrix W

$$
W = \begin{pmatrix} w_{11} & w_{12} & \ldots & w_{1n} \\ w_{21} & w_{22} & \ldots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \ldots & w_{nn} \end{pmatrix} = Q \cdot F \ \ \text{with}
$$

$$
Q = \begin{pmatrix} Q_{11} & Q_{12} & \ldots & Q_{1n} \\ Q_{21} & Q_{22} & \ldots & Q_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ Q_{n1} & Q_{n2} & \ldots & Q_{nn} \end{pmatrix} \ \ \text{and} \ \ F = \begin{pmatrix} f_1 & 0 & \ldots & 0 \\ 0 & f_2 & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & f_n \end{pmatrix}
$$

Stationary population or quasispecies as a function of the mutation or error rate *p*

Fitness landscapes showing error thresholds

Error threshold: Individual sequences

$n = 10$, $\sigma = 2$ and $d = 0, 1.0, 1.85$,

$s = 491$

The error threshold in replication

Preface

# Antiviral strategy on the horizon

Error catastrophe had its conceptual origins in the middle of the XXth century, when the consequences of mutations on enzymes involved in protein synthesis, as a theory of aging. In those times biological processes were generally perceived differently from today. Infectious diseases were regarded as a fleeting nuisance which would be eliminated through the use of antibiotics and antiviral agents. Microbial variation, although known in some cases, was not thought to be a significant problem for disease control. Variation in differentiated organisms was seen as resulting essentially from exchanges of genetic material associated with sexual reproduction. The problem was to unveil the mechanisms of inheritance, expression of genetic information and metabolism. Few saw that genetic change is occurring at present in all organisms, and still fewer recognized Darwinian principles as essential to the biology of pathogenic viruses and cells. Population geneticists rarely used bacteria or viruses as experimental systems to define concepts in biological evolution. The extent of genetic polymorphism among individuals of the same biological species came as a surprise when the first results on comparison of electrophoretic mobility of enzymes were obtained. With the advent of in vitro DNA recombination, and rapid nucleic acid sequencing techniques, molecular analyses of genomes reinforced the conclusion of extreme inter-individual genetic variation within the same species. Now, due largely to spectacular progress in comparative genomics, we see cellular DNAs, both prokaryotic and eukaryotic, as highly dynamic. Most cellular processes, including such essential information-bearing and transferring events as genome replication, transcription and translation, are increasingly perceived as inherently inaccurate. Viruses, and in particular RNA viruses, are among the most extreme examples of exploitation of replication inaccuracy for survival.

Error catastrophe, or the loss of meaningful genetic information through excess genetic variation, was formulated in quantitative terms as a consequence of quasispecies theory, which was first developed to explain self-organization and adaptability of primitive replicons in early stages of life. Recently, a conceptual extension of error catastrophe that could be defined as "induced genetic deterioration" has emerged as a possible antiviral strategy. This is the topic of the current special issue of *Virus Research*.

Few would nowadays doubt that one of the major obstacles for the control of viral disease is short-term adaptability of viral pathogens. Adaptability of viruses follows the same Darwinian principles that have shaped biological evolution over eons, that is, repeated rounds of reproduction with genetic variation, competition and selection, often perturbed by random events such as statistical fluctuations in population size. However, with viruses the consequences of the operation of these very same Darwinian principles are felt within very short times. Short-term evolution (within hours and days) can be also observed with some cellular pathogens, with subsets of normal cells, and cancer cells. The nature of RNA viral pathogens begs for alternative antiviral strategies, and forcing the virus to cross the critical error threshold for maintenance of genetic information is one of them.

The contributions to this volume have been chosen to reflect different lines of evidence (both theoretical and experimental) on which antiviral designs based on genetic deterioration inflicted upon viruses are being constructed. Theoretical studies have explored the copying fidelity conditions that must be fulfilled by any information-bearing replication system for the essential genetic information to be transmitted to progeny. Closely related to the theoretical developments have been numerous experimental studies on quasispecies dynamics and their multiple biological manifestations. The latter can be summarized by saying that RNA viruses, by virtue of existing as mutant spectra rather than defined genetic entities, remarkably expand their potential to overcome selective pressures intended to limit their replication. Indeed, the use of antiviral inhibitors in clinical practice and the design of vaccines for a number of major RNA virus-associated diseases, are currently presided by a sense of uncertainty. Another line of growing research is the enzymology of copying fidelity by viral replicases, aimed at understanding the molecular basis of mutagenic activities. Error catastrophe as a potential new antiviral strategy received an important impulse by the observation that ribavirin (a licensed antiviral nucleoside analogue) may be exerting, in some systems, its antiviral activity through enhanced mutagenesis. This has encouraged investigations on new mutagenic base analogues, some of them used in anticancer chemotherapy. Some chapters summarize these important biochemical studies on cell entry pathways and metabolism of mutagenic agents, that may find new applications as antiviral agents.

This volume intends to be basically a progress report, an introduction to a new avenue of research, and a realistic appraisal of the many issues that remain to be investigated. In this respect, I can envisage (not without many uncertainties) at least three lines of needed research: (i) One on further understanding of quasispecies dynamics in infected individuals to learn more on how to apply combinations of virus-specific mutagens and inhibitors in an effective way, finding synergistic combinations and avoiding antagonistic ones as well as severe clinical side effects. (ii) Another on a deeper understanding of the metabolism of mutagenic agents, in particular base and nucleoside analogues. This includes identification of the transporters that carry them into cells, an understanding of their metabolic processing, intracellular stability and alterations of nucleotide pools, among other issues. (iii) Still another line of needed research is the development of new mutagenic agents specific for viruses, showing no (or limited) toxicity for cells. Some advances may come from links with anticancer research, but others should result from the designs of new molecules, based on the structures of viral polymerases. I really hope that the reader finds this issue not only to be an interesting and useful review of the current situation in the field, but also a stimulating exposure to the major problems to be faced.

The idea to prepare this special issue came as a kind invitation of Ulrich Desselberger, former Editor of *Virus Research*, and then taken enthusiastically by Luis Enjuanes, recently appointed as Editor of *Virus Research*. I take this opportunity to thank Ulrich, Luis and the Editor-in-Chief of *Virus Research*, Brian Mahy, for their continued interest and support to the research on virus evolution over the years.

My thanks go also to the 19 authors who despite their busy schedules have taken time to prepare excellent manuscripts, to Elsevier staff for their prompt responses to my requests, and, last but not least, to Ms. Lucía Horrillo from Centro de Biología Molecular "Severo Ochoa" for her patient dealing with the correspondence with authors and the final organization of the issue.

Esteban Domingo
*Universidad Autónoma de Madrid*
*Centro de Biología Molecular "Severo Ochoa"*
*Consejo Superior de Investigaciones Científicas*
*Cantoblanco and Valdeolmos*
*Madrid, Spain*
Tel.: + 34 91 497 84858/9; fax: +34 91 497 4799
*E-mail address:* edomingo@cbm.uam.es
Available online 8 December 2004

Molecular evolution of viruses

## Results from kinetic theory of molecular evolution:

•Replicating ensembles of molecules form stationary populations called **quasispecies**, which represent the genetic reservoir of asexually reproducing species.

• For stable inheritance of genetic information mutation rates must not exceed a precisely defined and computable **error-threshold**.

•The error-threshold can be exploited for the development of novel antiviral strategies.

A fitness landscape including neutrality

Fig. 3.1. Behavior of mutant genes following their appearance in a finite population. Courses of change in the frequencies of mutants destined to fixation are depicted by thick paths. $N_e$ stands for the effective population size and $v$ is the mutation rate.

Motoo Kimura

Is the Kimura scenario correct for frequent mutations?

# STATIONARY MUTANT DISTRIBUTIONS AND EVOLUTIONARY OPTIMIZATION

■ PETER SCHUSTER and JÖRG SWETINA
Institut für theoretische Chemie
und Strahlenchemie der Universität Wien,
Währingerstraße 17,
A 1090 Wien,
Austria

Molecular evolution is modelled by erroneous replication of binary sequences. We show how the selection of two species of equal or almost equal selective value is influenced by its nearest neighbours in sequence space. In the case of perfect neutrality and sufficiently small error rates we find that the Hamming distance between the species determines selection. As the error rate increases the fitness parameters of neighbouring species become more and more important. In the case of almost neutral sequences we observe a critical replication accuracy at which a drastic change in the "quasispecies", in the stationary mutant distribution occurs. Thus, in frequently mutating populations fitness turns out to be an ensemble property rather than an attribute of the individual.

In addition we investigate the time dependence of the mean excess production as a function of initial conditions. Although it is optimized under most conditions, cases can be found which are characterized by decrease or non-monotonous change in mean excess productions.

*1. Introduction.* Recent data from populations of RNA viruses provided direct evidence for vast sequence heterogeneity (Domingo *et al.*, 1987). The origin of this diversity is not yet completely known. It may be caused by the low replication accuracy of the polymerizing enzyme, commonly a virus specific, RNA dependent RNA synthetase, or it may be the result of a high degree of selective neutrality of polynucleotide sequences. Eventually, both factors contribute to the heterogeneity observed. Indeed, mutations occur much more frequently than previously assumed in microbiology. They are by no means rare events and hence, neither the methods of conventional population genetics (Ewens, 1979) nor the neutral theory (Kimura, 1983) can be applied to these virus populations. Selectively neutral variants may be close with respect to Hamming distance and then the commonly made assumption that the mutation backflow from the mutants to the wilde type is negligible does not apply.

A kinetic theory of polynucleotide evolution which was developed during the past 15 years (Eigen, 1971; 1985; Eigen and Schuster, 1979; Eigen *et al.*, 1987; Schuster, 1986); Schuster and Sigmund, 1985) treats correct replication and mutation as parallel reactions within one and the same reaction network

Neutral network

$\lambda = 0.01$, $s = 367$

$d_H = 1$

$\lim_{p \to 0} x_1(p) = x_2(p) = 0.5$

Neutral network

$\lambda = 0.01$, $s = 877$

$d_H = 2$

$\lim_{p \to 0} x_1(p) = a$

$\lim_{p \to 0} x_2(p) = 1 - a$

$d_H \quad 3$

random fixation in the sense of
Motoo Kimura

Pairs of genotypes in neutral replication networks

Neutral network

$\lambda = 0.01$, s = 367

Neutral network: Individual sequences

n = 10, $\sigma$ = 1.1, d = 1.0

Consensus sequence of a quasispecies of two strongly coupled sequences of Hamming distance $d_H(X_i, X_j) = 1$.

Neutral network

$\lambda = 0.01, \ s = 877$

Neutral network: Individual sequences

$n = 10, \ \sigma = 1.1, \ d = 1.0$

Consensus sequence of a quasispecies of two strongly coupled sequences of
Hamming distance $d_H(X_i, X_j) = 2$.

Relative concentration $\bar{x}(p)$

Error rate $p \rightarrow$

N = 7

Neutral network

$\lambda = 0.10, \quad s = 229$

Neutral networks with increasing $\lambda$: $\lambda = 0.10$, s = 229

Relative concentration $\bar{x}(p)$

N = 7

Error rate $p \rightarrow$

184   600
248   760   728
504   729

Neutral network

$\lambda = 0.10, \ s = 229$

Perturbation matrix W

$$W = \begin{pmatrix} f & 0 & \varepsilon & 0 & 0 & 0 & 0 \\ 0 & f & \varepsilon & 0 & 0 & 0 & 0 \\ \varepsilon & \varepsilon & f & \varepsilon & 0 & 0 & 0 \\ 0 & 0 & \varepsilon & f & \varepsilon & 0 & 0 \\ 0 & 0 & 0 & \varepsilon & f & \varepsilon & \varepsilon \\ 0 & 0 & 0 & 0 & \varepsilon & f & 0 \\ 0 & 0 & 0 & 0 & \varepsilon & 0 & f \end{pmatrix}$$

Eigenvalues of W

$$\lambda_0 = f + 2\varepsilon,$$

$$\lambda_1 = f + \sqrt{2}\varepsilon,$$
$$\lambda_{2,3,4} = f,$$
$$\lambda_5 = f - \sqrt{2}\varepsilon,$$
$$\lambda_6 = f - 2\varepsilon.$$

Largest eigenvector of W

$$\xi_0 = (0.1, 0.1, 0.2, 0.2, 0.2, 0.1, 0.1).$$

Neutral networks with increasing $\lambda$: $\lambda = 0.10, \ s = 229$

N = 24

Neutral network

λ = 0.15, s = 229

Neutral networks with increasing λ: λ = 0.15, s = 229

N = 70

Neutral network

λ = 0.20, s = 229

Neutral networks with increasing λ:  λ = 0.20, s = 229

5'-end **GCGGAU**UUA**GCUC**AGUUGGGA**GAGC**G**CCAGA**CUGAAGA**UCUGG**AGGUC**CUGUG**UUCGAUC**CACAG**A**AUUCGC**ACCA 3'-end

5' - end

**N₁**

**N₂**

Na ⊕

**N₃**

Na ⊕

Definition of RNA structure

3'-end

5'-end

70

60

10

50

20

30   40

Sequence: GCGGAUUUAGCUCAGDDGGGAGAGCMCCAGACUGAAYAUCUGGAGMUCCUGUGTPCGAUCCACAGAAUUCGCACCA

$N = 4^n$

5'-End    3'-End

Secondary structure

$N_S < 3^n$

Symbolic notation: 5'-End ((((((···((((········)))))·(((((·······)))))····(((((·······)))))·)))))))···· 3'-End

Criterion:  Minimum free energy (mfe)

Rules:  _ ( _ ) _  ∈ {**AU**,**CG**,**GC**,**GU**,**UA**,**UG**}

A symbolic notation of RNA secondary structure that is equivalent to the conventional graphs

*Neutral network*

Sequence space          Structure space

many genotypes     ⇒     one phenotype

**Evolution *in silico***

W. Fontana, P. Schuster,
*Science* **280** (1998), 1451-1455

random individuals. The primer pair used for genomic DNA amplification is 5'-TCTCCCTGGATTCT-CATTTA-3' (forward) and 5'-TCTTTGTCTTCTGT-TCCACC-3' (reverse). Reactions were performed in 25 μl using 1 unit of Taq DNA polymerase with each primer at 0.4 μM; 200 μM each dATP, dTTP, dGTP, and dCTP; and PCR buffer [10 mM tris-HCl (pH 8.3), 50 mM KCl₂,1.5 mM MgCl₂] in a cycle condition of 94°C for 1 min and then 35 cycles of 94°C for 30 s, 55°C for 30 s, and 72°C for 30 s followed by 72°C for 6 min. PCR products were purified (Qiagen), digested with Xmn I, and separated in a 2% agarose gel.

32. A nonsense mutation may affect mRNA stability and result in degradation of the transcript [L. Maquat, *Am. J. Hum. Genet.* **59**, 279 (1996)].

33. Data not shown; a dot blot with poly (A)⁺ RNA from 50 human tissues (The Human RNA Master Blot, 7770-1, Clontech Laboratories) was hybridized with a probe from exons 29 to 47 of *MYO15* using the same condition as Northern blot analysis (13).

34. Smith–Magenis syndrome (SMS) is due to deletions of 17p11.2 of various sizes, the smallest of which includes *MYO15* and perhaps 20 other genes [(6); K-S Chen, L. Potocki, J. R. Lupski, *MRDD Res. Rev.* **2**, 122 (1996)]. *MYO15* expression is easily detected in the pituitary gland (data not shown). Haploinsufficiency for *MYO15* may explain a portion of the SMS

35. R. A. Fridell, data not shown.

36. K. B. Avraham et al., *Nature Genet.* **11**, 369 (1995); X-Z. Liu et al., ibid. **17**, 268 (1997); F. Gibson et al., *Nature* **374**, 62 (1995); D. Weil et al., ibid., p. 60.

37. RNA was extracted from cochlea (membranous labyrinths) obtained from human fetuses at 18 to 22 weeks of development in accordance with guidelines established by the Human Research Committee at the Brigham and Women's Hospital. Only samples without evidence of degradation were pooled for poly (A)⁺ selection over oligo(dT) columns. First-strand cDNA was prepared using an Advantage RT-for-PCR kit (Clontech Laboratories). A portion of the first-strand cDNA (4%) was amplified by PCR with Advantage cDNA polymerase mix (Clontech Laboratories) using human *MYO15*-specific oligonucleotide primers (forward, 5'-GCATGACCTGCCGGCTAAT-GGG-3'; reverse, 5'-CTCACGGCTTCTGCATGGT-GCTCGGCTGGC-3'). Cycling conditions were 40 s at 94°C; 40 s at 66°C (3 cycles), 60°C (5 cycles), and 55°C (29 cycles); and 45 s at 68°C. PCR products were visualized by ethidium bromide staining after fractionation in a 1% agarose gel. A 688-bp PCR

38. We are grateful to the people of Bengkala, Bali, and the two families in India. We thank J. R. Lupski and K.-S. Chen for providing the human chromosome 17 cosmid library. For technical and computational assistance, we thank N. Dietrich, M. Fergusson, A. Gupta, E. Sorbello, R. Torkzadeh, C. Varner, M. Walker, G. Bouffard, and S. Beckstrom-Sternberg (National Institutes of Health Intramural Sequencing Center). We thank J. T. Hinnant, I. N. Arhya, and S. Winata for assistance in Bali, and T. Barber, S. Sullivan, E. Green, D. Drayna, and J. Battey for helpful comments on this manuscript. Supported by the National Institute on Deafness and Other Communication Disorders (NIDCD) (Z01 DC 00035-01 and Z01 DC 00038-01 to T.B.F. and E.R.W. and R01 DC 03402 to C.C.M.), the National Institute of Child Health and Human Development (R01 HD30428 to S.A.C.) and a National Science Foundation Graduate Research Fellowship to F.J.P. This paper is dedicated to J. B. Snow Jr. on his retirement as the Director of the NIDCD.

product is expected from amplification of the human *MYO15* cDNA. Amplification of human genomic DNA with this primer pair would result in a 2903-bp fragment.

9 March 1998; accepted 17 April 1998

## Continuity in Evolution: On the Nature of Transitions

Walter Fontana and Peter Schuster

To distinguish continuous from discontinuous evolutionary change, a relation of nearness between phenotypes is needed. Such a relation is based on the probability of one phenotype being accessible from another through changes in the genotype. This nearness relation is exemplified by calculating the shape neighborhood of a transfer RNA secondary structure and provides a characterization of discontinuous shape transformations in RNA. The simulation of replicating and mutating RNA populations under selection shows that sudden adaptive progress coincides mostly, but not always, with discontinuous shape transformations. The nature of these transformations illuminates the key role of neutral genetic drift in their realization.

A much-debated issue in evolutionary biology concerns the extent to which the history of life has proceeded gradually or has been punctuated by discontinuous transitions at the level of phenotypes (1). Our goal is to make the notion of a discontinuous transition more precise and to understand how it arises in a model of evolutionary adaptation.

We focus on the narrow domain of RNA secondary structure, which is currently the simplest computationally tractable, yet realistic phenotype (2). This choice enables the definition and exploration of concepts that may prove useful in a wider context. RNA secondary structures represent a coarse level of analysis compared with the three-dimensional structure at atomic resolution. Yet, secondary structures are empir-

ically well defined and obtain their biophysical and biochemical importance from being a scaffold for the tertiary structure. For the sake of brevity, we shall refer to secondary structures as "shapes." RNA combines in a single molecule both genotype (replicatable sequence) and phenotype (selectable shape), making it ideally suited for in vitro evolution experiments (3, 4).

To generate evolutionary histories, we used a stochastic continuous time model of an RNA population replicating and mutating in a capacity-constrained flow reactor under selection (5, 6). In the laboratory, a goal might be to find an RNA aptamer binding specifically to a molecule (4). Although in the experiment the evolutionary end product was unknown, we thought of its shape as being specified implicitly by the imposed selection criterion. Because our intent is to study evolutionary histories rather than end products, we defined a target shape in advance and assumed the replication rate of a sequence to be a function of

the similarity between its shape and the target. An actual situation may involve more than one best shape, but this does not affect our conclusions.

An instance representing in its qualitative features all the simulations we performed is shown in Fig. 1A. Starting with identical sequences folding into a random shape, the simulation was stopped when the population became dominated by the target, here a canonical tRNA shape. The black curve traces the average distance to the target (inversely related to fitness) in the population against time. Aside from a short initial phase, the entire history is dominated by steps, that is, flat periods of no apparent adaptive progress, interrupted by sudden approaches toward the target structure (7). However, the dominant shapes in the population not only change at these marked events but undergo several fitness-neutral transformations during the periods of no apparent progress. Although discontinuities in the fitness trace are evident, it is entirely unclear when and on the basis of what the series of successive phenotypes itself can be called continuous or discontinuous.

A set of entities is organized into a (topological) space by assigning to each entity a system of neighborhoods. In the present case, there are two kinds of entities: sequences and shapes, which are related by a thermodynamic folding procedure. The set of possible sequences (of fixed length) is naturally organized into a space because point mutations induce a canonical neighborhood. The neighborhood of a sequence consists of all its one-error mutants. The problem is how to organize the set of possible shapes into a space. The issue arises because, in contrast to sequences, there are

3'-end

5'-end

Structure of
randomly chosen
initial sequence

3'-end

5'-end

70

60

10

50

20

30  40

Phenylalanyl-tRNA as
target structure

Stock Solution $\longrightarrow$

Reaction Mixture $\longrightarrow$

**Replication rate constant (Fitness)**:

$$f_k = \gamma \, / \, [\alpha + \Delta d_S^{(k)}]$$

$$\Delta d_S^{(k)} = d_H(S_k, S_\tau)$$

**Selection pressure**:

The population size,

$N = \#$ RNA moleucles,

is determined by the flux:

$$N(t) \approx \overline{N} \pm \sqrt{\overline{N}}$$

**Mutation rate**:

$p = 0.001$ / Nucleotide $\times$ Replication

The flow reactor as a device for studying the evolution of molecules *in vitro* and *in silico*.

*In silico* optimization in the flow reactor: Evolutionary Trajectory

**28 neutral point mutations** during a long quasi-stationary epoch

Average structure distance to target Δd$_S$ — Number of relay step

08
10
12
14

Evolutionary trajectory

Time (arbitrary units)

| | |
|---|---|
| entry | GGUAUGGGCGUUGAAUAGUAGGGUUUAAACCAAUCGGCCAACGAUCUCGUGUGCGCAUUUCAUAUCCCGUACAGAA |
| 8 | .(((((((((((........(((....))).......)))))).....(((((......))))))))))).... |
| exit | GGUAUGGGCGUUGAAUAAUAGGGUUUAAACCAAUCGGCCAACGAUCUCGUGUGCGCAUUUCAUAUCCCAUACAGAA |
| entry | GGUAUGGGCGUUGAAUAAUAGGGUUUAAACCAAUCGGCCAACGAUCUCGUGUGCGCAUUUCAUAUACCAUACAGAA |
| 9 | .((((((.(((((........(((....))).......)))))).....(((((......))))).))))).... |
| exit | UGGAUGGACGUUGAAUAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACACCGUCCCAAG |
| entry | UGGAUGGACGUUGAAUAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACAGCGUCCCAAG |
| 10 | .(((((..(((((........(((....))).......)))))).....(((((......)))))..))))).... |
| exit | UGGAUGGACGUUGAAUAACAAGGUAUCGACCAAACAACCAACGAGUAAGUGUGUACGCCCCACACAGCGUCCCAAG |

**Transition inducing point mutations** change the molecular structure

**Neutral point mutations** leave the molecular structure unchanged

Neutral genotype evolution during phenotypic stasis

Evolutionary trajectory

Spreading of the population
on neutral networks

Drift of the population center
in sequence space

start of optimization

start of optimization

end of optimization

end of optimization

Cost function

Genotype space

start of optimization

start of optimization

Cost function

Genotype space

target

A sketch of optimization on neutral networks

## Neutrality in molecular structures and its role in evolution:

• Neutrality is an essential feature in biopolymer structures at the resolution that is relevant for function.

• Neutrality manifests itself in the search for minimum free energy structures.

• Diversity in function despite neutrality in structures results from differences in suboptimal conformations and folding kinetics.

• Neutrality is indispensible for optimization and adaptation.

A model genome with 12 genes

Regulatory gene
Structural gene
Regulatory protein or RNA
Enzyme
Metabolite

Sketch of a genetic and metabolic network

# Dynamic patterns of gene regulation I: Simple two-gene systems

Stefanie Widder[a], Josef Schicho[b], Peter Schuster[a,c,*]

[a]*Institut für Theoretische Chemie der Universität Wien, Währingerstraße 17, A-1090 Wien, Austria*
[b]*RICAM—Johann Radon Institute for Computational and Applied Mathematics of the Austrian Academy of Sciences, Altenbergerstraße 69, A-4040 Linz, Austria*
[c]*Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA*

## Abstract

Regulation of gene activities is studied by means of computer assisted mathematical analysis of ordinary differential equations (ODEs) derived from binding equilibria and chemical reaction kinetics. Here, we present results on cross-regulation of two genes through activator and/or repressor binding. Arbitrary (differentiable) binding function can be used but systematic investigations are presented for gene–regulator complexes with integer valued Hill coefficients up to $n = 4$. The dynamics of gene regulation is derived from bifurcation patterns of the underlying systems of kinetic ODEs. In particular, we present analytical expressions for the parameter values at which one-dimensional (transcritical, saddle-node or pitchfork) and/or two-dimensional (Hopf) bifurcations occur. A classification of regulatory states is introduced, which makes use of the sign of a 'regulatory determinant' $D$ (being the determinant of the block in the Jacobian matrix that contains the derivatives of the regulator binding functions): (i) systems with $D < 0$, observed, for example, if both proteins are activators or repressors, to give rise to one-dimensional bifurcations only and lead to bistability for $n \geqslant 2$ and (ii) systems with $D > 0$, found for combinations of activation and repression, sustain a Hopf bifurcation and undamped oscillations for $n > 2$. The influence of basal transcription activity on the bifurcation patterns is described. Binding of multiple subunits can lead to richer dynamics than pure activation or repression states if intermediates between the unbound state and the fully saturated DNA initiate transcription. Then, the regulatory determinant $D$ can adopt both signs, plus and minus.
© 2007 Elsevier Ltd. All rights reserved.

*Keywords:* Basal transcription; Bifurcation analysis; Cooperative binding; Gene regulation; Hill coefficient; Hopf bifurcation

## 1. Introduction

Theoretical work on gene regulation goes back to the 1960s (Monod et al., 1963) soon after the first repressor protein had been discovered (Jacob and Monod, 1961). A little later the first paper on oscillatory states in gene regulation was published (Goodwin, 1965). The interest in gene regulation and its mathematical analysis never ceased (Tiwari et al., 1974; Tyson and Othmer, 1978; Smith, 1987) and saw a great variety of different attempts to design models of genetic regulatory networks that can be used in systems biology for computer simulation of *gen*(etic and met)*abolic* networks.[1] Most models in the literature aim at a minimalist dynamic description which, nevertheless, tries to account for the basic regulatory functions of large networks in the cell in order to provide a better understanding of cellular dynamics. A classic in general regulatory dynamics is the monograph by Thomas and D'Ari (1990). The currently used mathematical methods comprise application of Boolean logic (Thomas and Kaufman, 2001b; Savageau, 2001; Albert and Othmer, 2003), stochastic processes (Hume, 2000) and deterministic dynamic models, examples are Cherry and Adler (2000), Bindschadler and Sneyd (2001) and Kobayashi et al. (2003) and the recent elegant analysis of bistability (Craciun et al.,

*Corresponding author. Institut für Theoretische Chemie der Universität Wien, Währingerstraße 17, A-1090 Wien, Austria.
Tel.: +43 1 4277 527 43; fax: +43 1 4277 527 93.
*E-mail address:* pks@tbi.univie.ac.at (P. Schuster).

[1]Discussion and analysis of combined genetic and metabolic networks has become so frequent and intense that we suggest to use a separate term, *genabolic networks*, for this class of complex dynamical systems.

|   | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | **Biochemical Pathways** | | | | | | | | | | | |
| 2 | | | | | | | | | | | | |
| 3 | | | | | | | | | | | | |
| 4 | | | | | | | | | | | | |
| 5 | | | | | | | | | | | | |
| 6 | | | | | | | | | | | | |
| 7 | | | | | | | | | | | | |
| 8 | | | | | | | | | | | | |
| 9 | | | | | | | | | | | | |
| 10 | | | | | | | | | | | | |

The reaction network of cellular metabolism published by Boehringer-Ingelheim.

The citric acid or Krebs cycle (enlarged from previous slide).

**E. coli**: Genome length $4 \times 10^6$ nucleotides

Number of cell types 1

Number of genes 4 460



**Man**: Genome length $3 \times 10^9$ nucleotides

Number of cell types 200

Number of genes $\approx 30\ 000$



Complexity in biology

Wolfgang Wieser. 1998. ‚*Die Erfindung der Individualität*' oder ‚*Die zwei Gesichter der Evolution*'. Spektrum Akademischer Verlag, Heidelberg 1998

The difficulty to define the notion of „gene".

Helen Pearson,
*Nature* **441**: 399-401, 2006

---

# WHAT IS A GENE?

The idea of genes as beads on a DNA string is fast fading. Protein-coding sequences have no clear beginning or end and RNA is a key part of the information package, reports **Helen Pearson**.

'Gene' is not a typical four-letter word. It is not offensive. It is never bleeped out of TV shows. And where the meaning of most four-letter words is all too clear, that of gene is not. The more expert scientists become in molecular genetics, the less easy it is to be sure about what, if anything, a gene actually is.

Rick Young, a geneticist at the Whitehead Institute in Cambridge, Massachusetts, says that when he first started teaching as a young professor two decades ago, it took him about two hours to teach fresh-faced undergraduates what a gene was and the nuts and bolts of how it worked. Today, he and his colleagues need three months of lectures to convey the concept of the gene, and that's not because the students are any less bright. "It takes a whole semester to teach this stuff to talented graduates," Young says. "It used to be we could give a one-off definition and now it's much more complicated."

In classical genetics, a gene was an abstract concept — a unit of inheritance that ferried a characteristic from parent to child. As biochemistry came into its own, those characteristics were associated with enzymes or proteins, one for each gene. And with the advent of molecular biology, genes became real, physical things — sequences of DNA which when converted into strands of so-called messenger RNA could be used as the basis for building their associated protein piece by piece. The great coiled DNA molecules of the chromosomes were seen as long strings on which gene sequences sat like discrete beads.
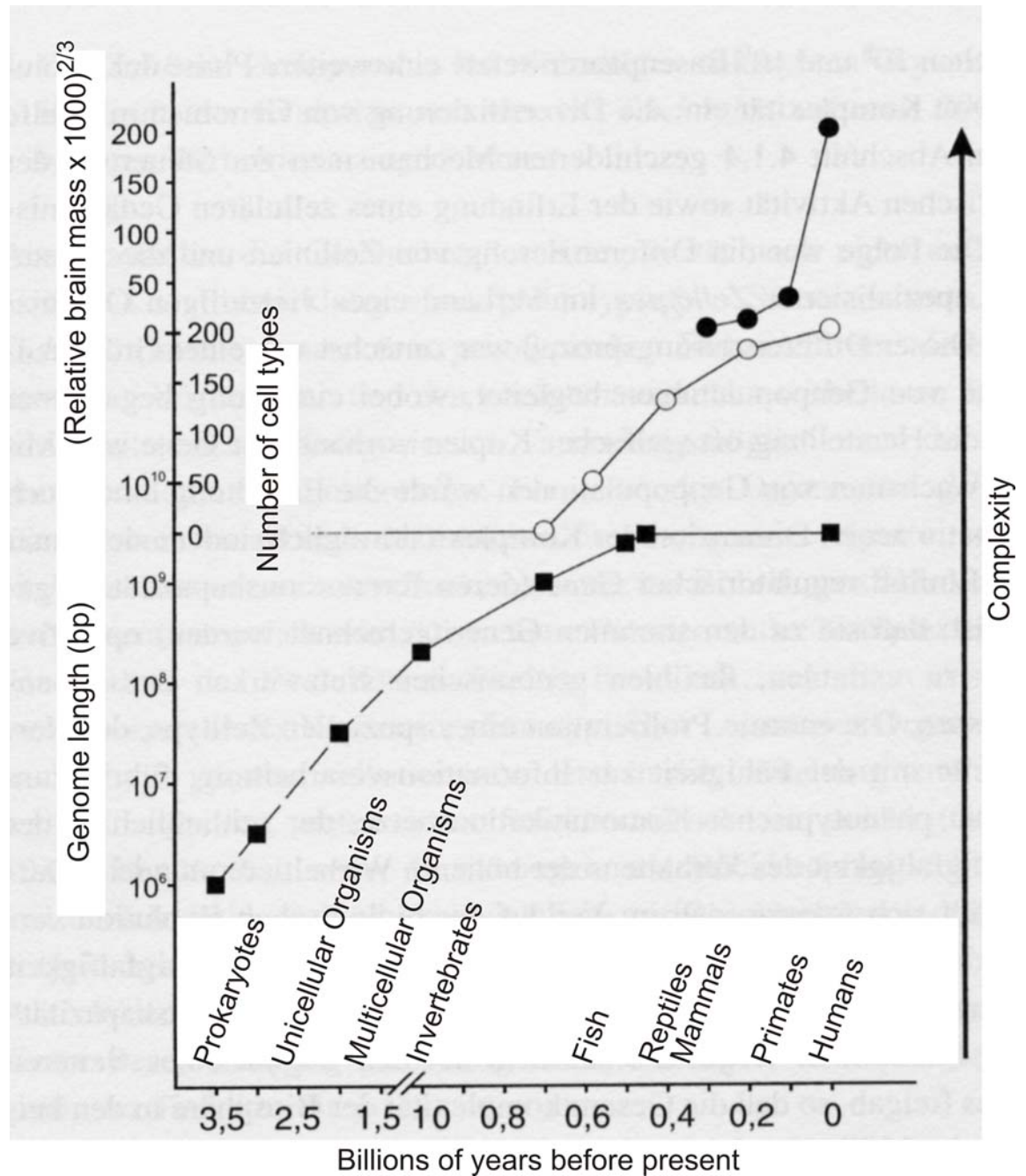
This picture is still the working model for many scientists. But those at the forefront of genetic research see it as increasingly old-fashioned — a crude approximation that, at best, hides fascinating new complexities and, at worst, blinds its users to useful new paths of enquiry.

Information, it seems, is parceled out along chromosomes in a much more complex way than was originally supposed. RNA molecules are not just passive conduits through which the gene's message flows into the world but active regulators of cellular processes. In some cases, RNA may even pass information across generations — normally the sole preserve of DNA.

An eye-opening study last year raised the possibility that plants sometimes rewrite their DNA on the basis of RNA messages inherited from generations past[1]. A study on page 469 of this issue suggests that a comparable phenomenon might occur in mice, and by implication in other mammals[2]. If this type of phenomenon is indeed widespread, it "would have huge implications," says evolutionary geneticist Laurence Hurst at the University of Bath, UK.

"All of that information seriously challenges our conventional definition of a gene," says molecular biologist Bing Ren at the University of California, San Diego. And the information challenge is about to get even tougher. Later this year, a glut of data will be released from the international Encyclopedia of DNA Elements (ENCODE) project. The pilot phase of ENCODE involves scrutinizing roughly 1% of the human genome in unprecedented detail; the aim is to find all the sequences that serve a useful purpose and explain what that purpose is. "When we started the ENCODE project I had a different view of what a gene was," says contributing researcher Roderic Guigo at the Center for Genomic Regulation in Barcelona. "The degree of complexity we've seen was not anticipated."

"We've come to the realization that the genome is full of overlapping transcripts." — Phillip Kapranov

### Under fire

The first of the complexities to challenge molecular biology's paradigm of a single DNA sequence encoding a single protein was alternative splicing, discovered in viruses in 1977 (see 'Hard to track', overleaf). Most of the DNA sequences describing proteins in humans have a modular arrangement in which exons, which carry the instructions for making proteins, are interspersed with non-coding introns. In alternative splicing, the cell snips out introns and sews together the exons in various different orders, creating messages that can code for different proteins. Over the years geneticists have also documented overlapping genes, genes within genes and countless other weird arrangements (see 'Muddling over genes', overleaf).

Alternative splicing, however, did not in itself require a drastic reappraisal of the notion of a gene; it just showed that some DNA sequences could describe more than one protein. Today's assault on the gene concept is more far reaching, fuelled largely by studies that show the pre-viously unimagined scope of RNA.

The one gene, one protein idea is coming under particular assault from researchers who are comprehensively extracting and analysing the RNA messages, or transcripts, manufactured by genomes, including the human and mouse genome. Researchers led by Thomas Gingeras at the company Affymetrix in Santa Clara, California, for example, recently studied all the transcripts from ten chromosomes across eight human cell lines and worked out precisely where on the chromosomes each of the transcripts came from[3].

The picture these studies paint is one of mind-boggling complexity. Instead of discrete genes dutifully mass-producing identical RNA transcripts, a teeming mass of transcription converts many segments of the genome into multiple RNA ribbons of differing lengths. These ribbons can be generated from both strands of DNA, rather than from just one as was conventionally thought. Some of these transcripts come from regions of DNA previously identified as holding protein-coding genes. But many do not. "It's somewhat revolutionary," says Gingeras's colleague Phillip Kapranov. "We've come to the realization that the genome is full of overlapping transcripts."

Other studies, one by Guigo's team[4], and one by geneticist Rotem Sorek[5], now at Tel Aviv University, Israel, and his colleagues, have hinted at the reasons behind the mass of transcription. The two teams investigated occasional reports that transcription can start at a DNA sequence associated with one protein and run straight through into the gene for a completely different protein, producing a fused transcript. By delving into databases of human RNA transcripts, Guigo's team estimate that 4–5% of the DNA in regions conventionally recognized as genes is transcribed in this way. Producing fused transcripts could be one way for a cell to generate a greater variety of proteins from a limited number of exons, the researchers say.

Many scientists are now starting to think that the descriptions of proteins encoded in DNA know no borders — that each sequence reaches into the next and beyond. This idea will be one of the central points to emerge from the ENCODE project when its results are published later this year.
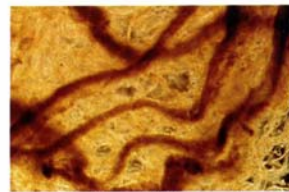
Kapranov and others say that they have documented many examples of transcripts in which protein-coding exons from one part of the genome combine with exons from another

Spools of DNA (above) still harbour surprises, with one protein-coding gene often overlapping the next.

ENCODE stands for
**ENC**yclopedia **O**f **D**NA **E**lements.

**ENCODE** Project Consortium.
Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project.
*Nature* **447**:799-816, 2007

# Acknowledgement of support

**Universität Wien**

# Coworkers

**Peter Stadler**, **Bärbel M. Stadler**, Universität Leipzig, GE

**Paul E. Phillipson**, University of Colorado at Boulder, CO

**Heinz Engl, Philipp Kügler**, **James Lu**, **Stefan Müller**, RICAM Linz, AT

**Jord Nagel**, **Kees Pleij**, Universiteit Leiden, NL

**Walter Fontana**, Harvard Medical School, MA

**Christian Reidys**, **Christian Forst**, Los Alamos National Laboratory, NM

**Ulrike Göbel, Walter Grüner**, **Stefan Kopp**, **Jaqueline Weber,** Institut für
Molekulare Biotechnologie, Jena, GE

**Ivo L.Hofacker**, **Christoph Flamm**, **Andreas Svrček-Seiler**, Universität Wien, AT

**Kurt Grünberger, Michael Kospach** , **Andreas Wernitznig**, **Stefanie Widder,**
**Stefan Wuchty**, Universität Wien, AT

**Jan Cupal**, **Stefan Bernhart, Lukas Endler, Ulrike Langhammer**, **Rainer Machne,**
**Ulrike Mückstein, Hakim Tafer, Thomas Taylor,** Universität Wien, AT

**Universität Wien**

Web-Page for further information:

http://www.tbi.univie.ac.at/~pks