

The theory of evolution in the light of 21st century's science

Peter Schuster

Institut für Theoretische Chemie, Universität Wien, Österreich

und

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



Conference on Evolutionism and Religion

Florence, 19.-21.11.2009

Web-Page for further information:

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



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.**

1. Darwin's natural selection
2. The tree of life
3. From evolution *in vitro* to biotechnology
4. Genotypes with multiple functions
5. How complex is biology?

1. **Darwin's natural selection**
2. The tree of life
3. From evolution *in vitro* to biotechnology
4. Genotypes with multiple functions
5. How complex is biology?

Genotype, Genome

Collection of genes

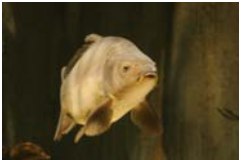
Developmental program

Highly specific environmental conditions

Unfolding of the genotype

Phenotype

Evolution explains the origin of species and their interactions



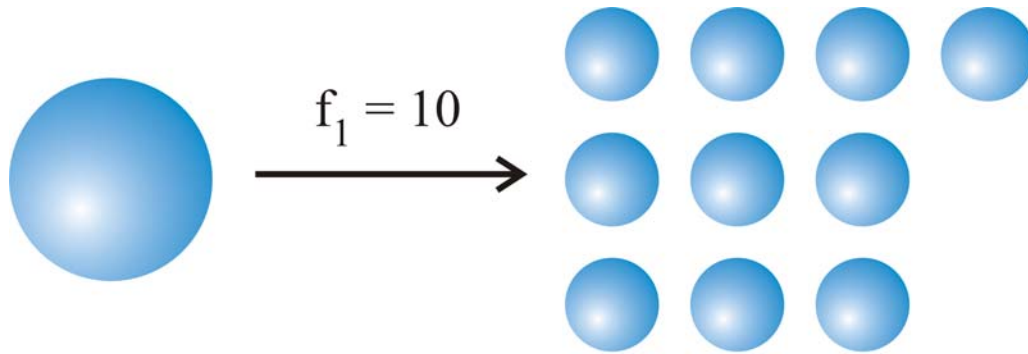


Three necessary conditions for Darwinian evolution are:

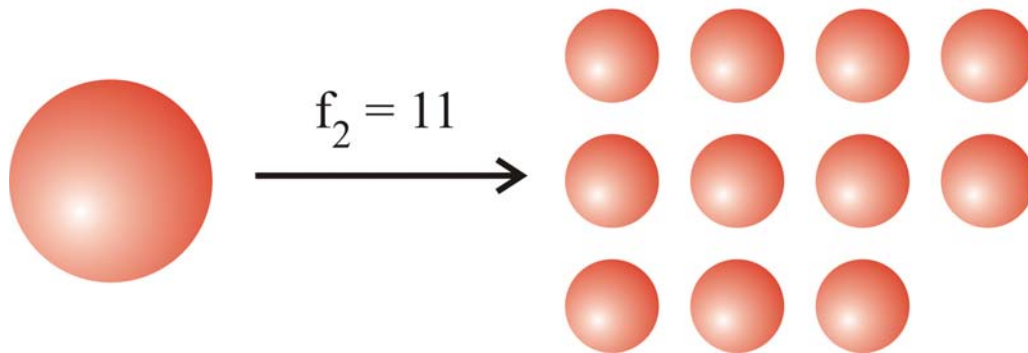
1. **Multiplication,**
2. **Variation,** and
3. **Selection.**

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

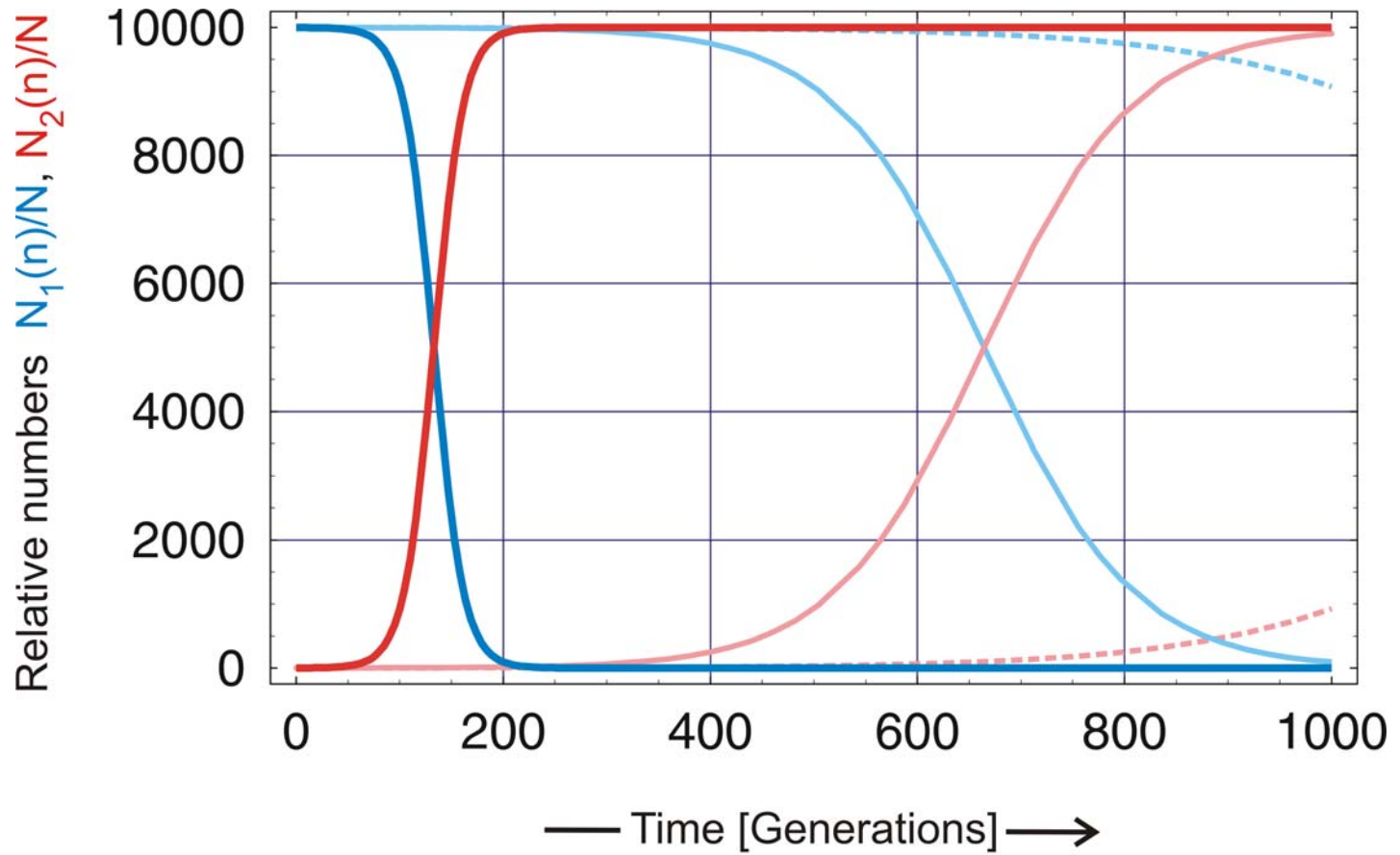
One important property of the Darwinian scenario is that **variations** in the form of mutations or recombination events occur **uncorrelated** with their **effects on the selection process**.



$$s = \frac{f_2 - f_1}{f_1} = 0.1$$



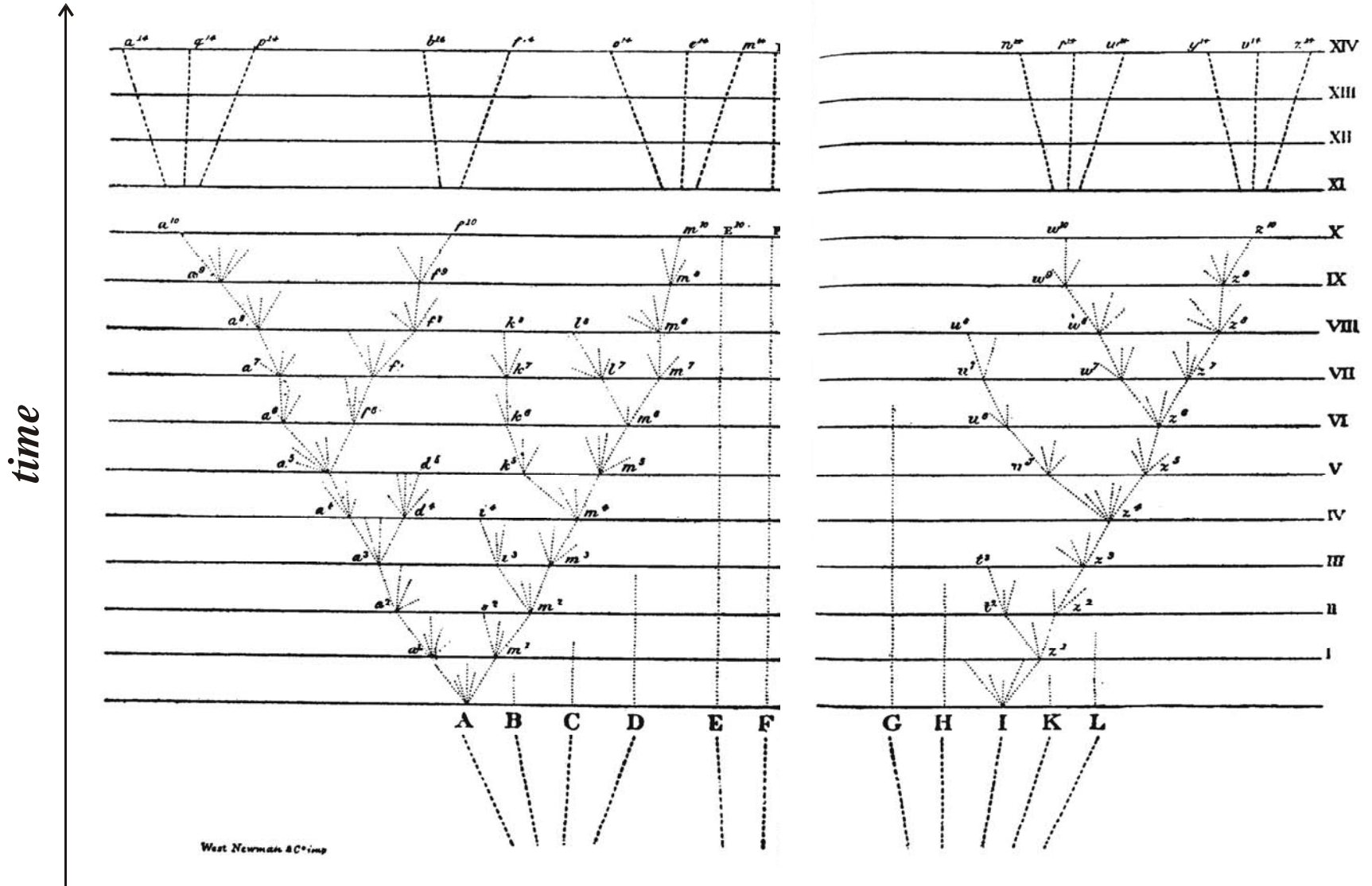
Two variants with a mean progeny of ten or eleven descendants



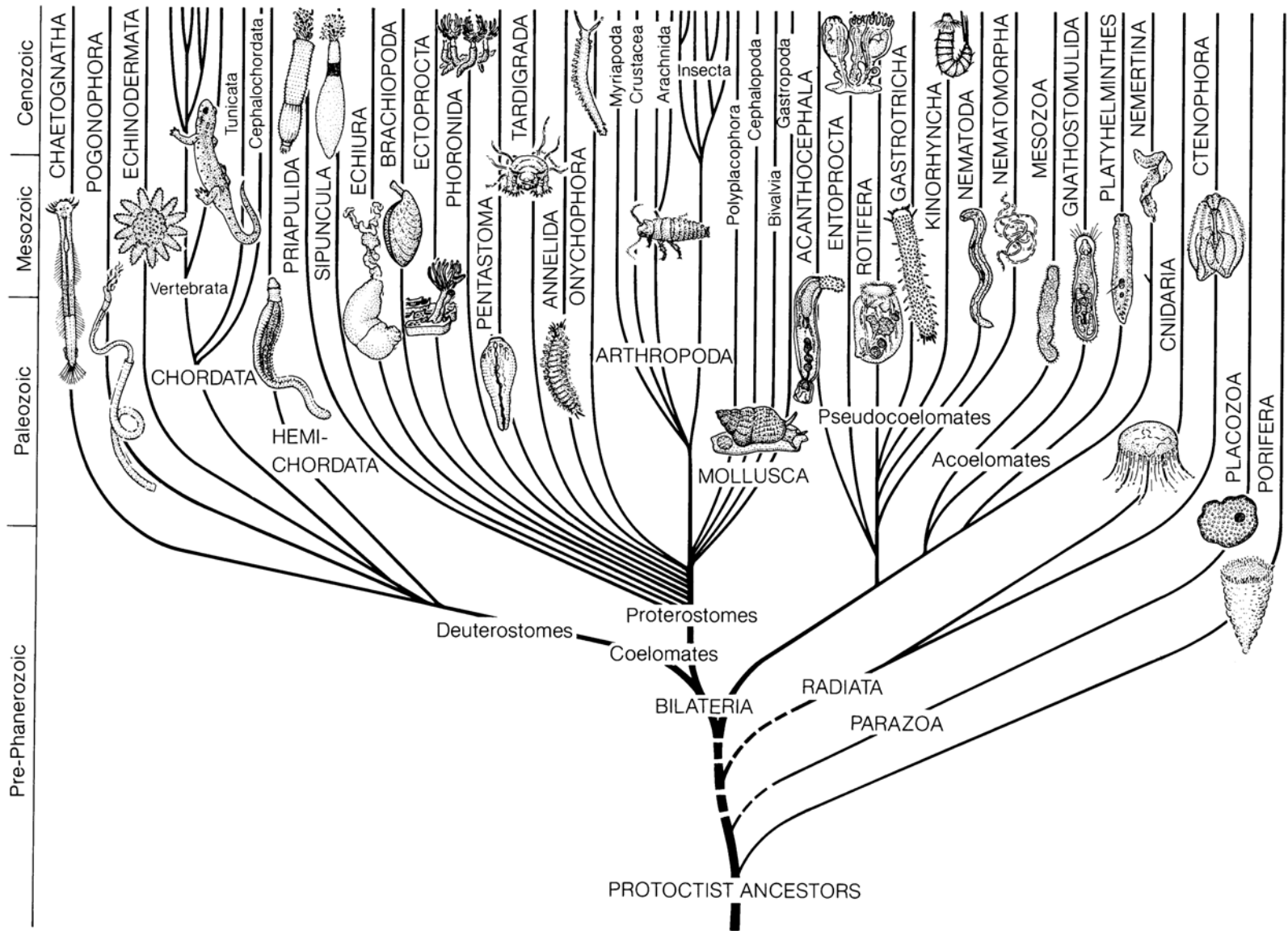
$$N_1(0) = 9999, N_2(0) = 1; \quad s = 0.1, 0.02, 0.01$$

Selection of advantageous mutants in populations of $N = 10\,000$ individuals

1. Darwin's natural selection
- 2. The tree of life**
3. From evolution *in vitro* to biotechnology
4. Genotypes with multiple functions
5. How complex is biology?



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



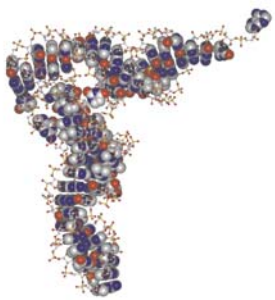
Modern phylogenetic tree: Lynn Margulis, Karlene V. Schwartz. *Five Kingdoms. An Illustrated Guide to the Phyla of Life on Earth.* W.H. Freeman, San Francisco, 1982.

Genotype, Genome

CGGGATTAGCTCAGTTGGGAGAGCGCCAGACTGAAGATCTGGAGGTCCTGTGTTTCGATCCACAGAATTTCGCACCA

Quantitative biology

'the new biology is the chemistry of living matter'



evolution of RNA molecules, ribozymes and splicing, the idea of an RNA world, selection of RNA molecules, RNA editing, the ribosome is a ribozyme, small RNAs and RNA switches.

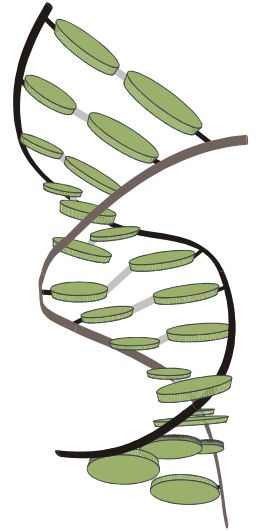
The exciting RNA story

Biochemistry
molecular biology
structural biology
molecular evolution
molecular genetics
systems biology
bioinformatics
epigenetics

Unfolding of the genotype

Phenotype

Highly specific environmental conditions



John Kendrew



Manfred Eigen



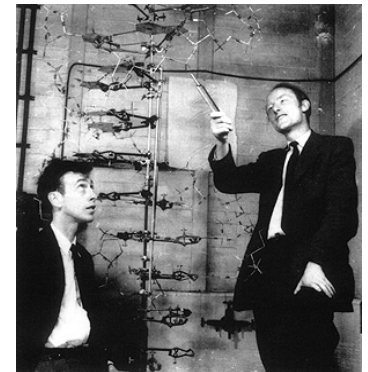
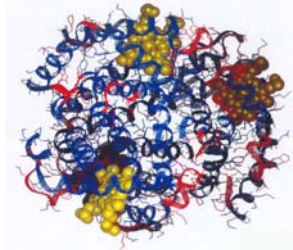
Molecular evolution
Linus Pauling and
Emile Zuckerkandl



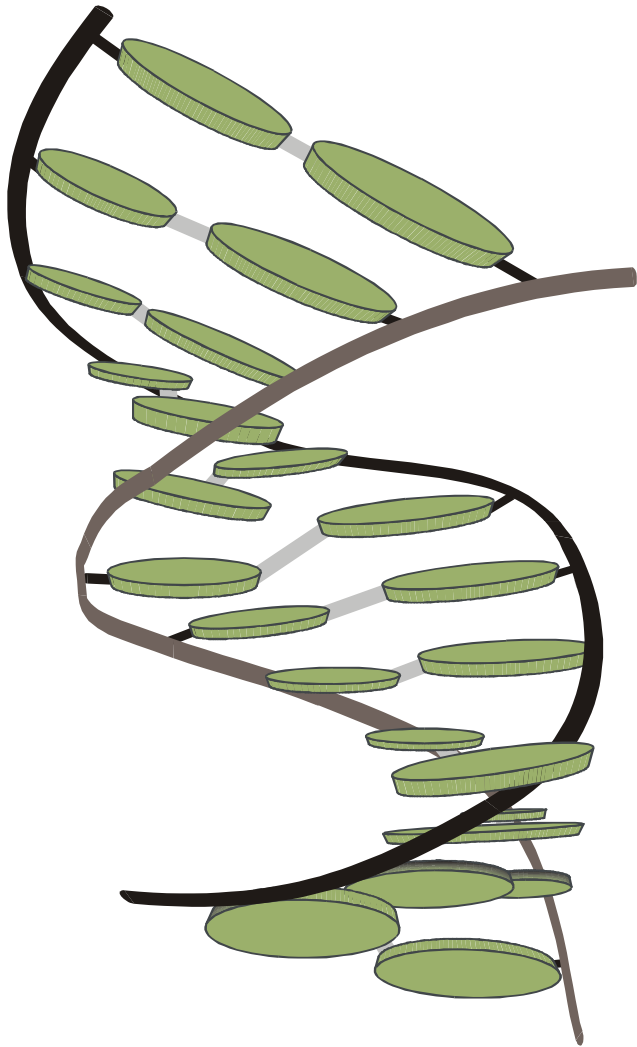
Hemoglobin sequence
Gerhard Braunitzer



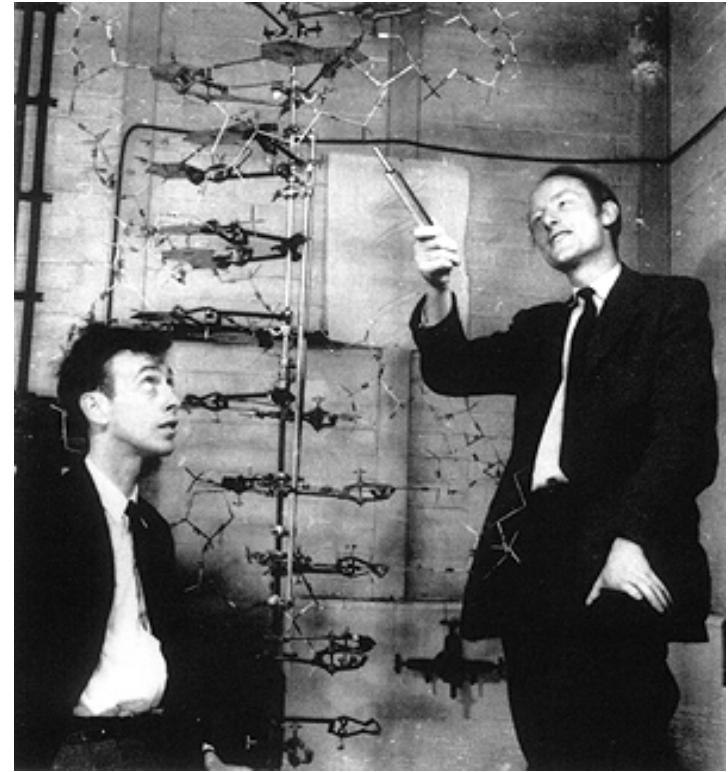
Max Perutz



James D. Watson und
Francis H.C. Crick



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

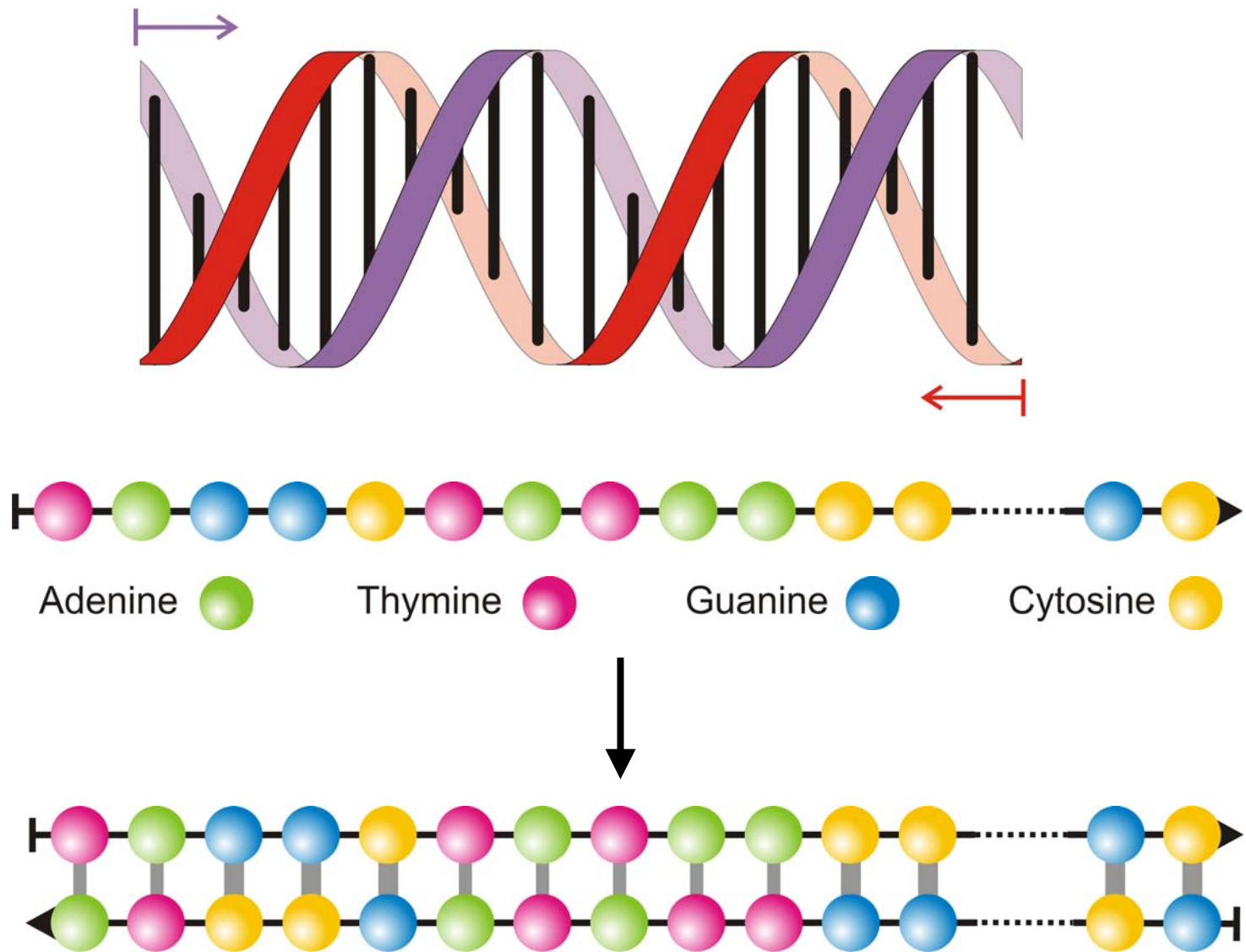


James D. Watson, 1928-, and Francis H.C. Crick, 1916-2004

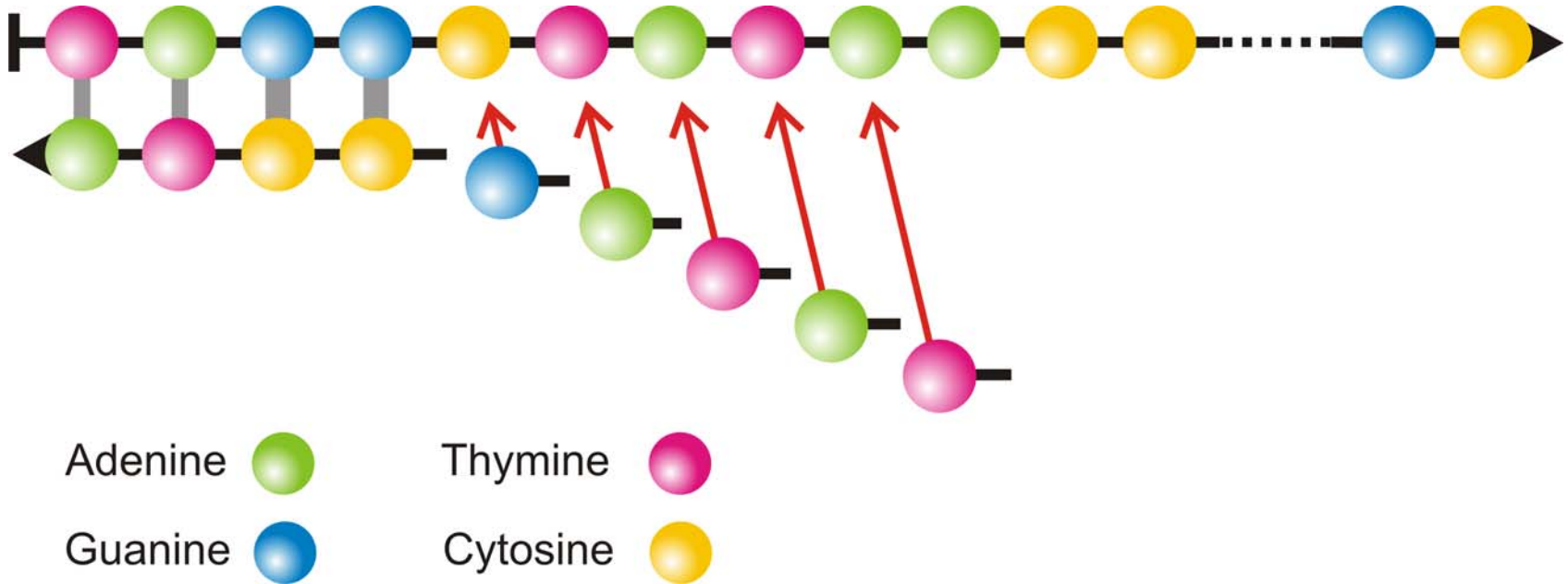
Nobel prize 1962

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

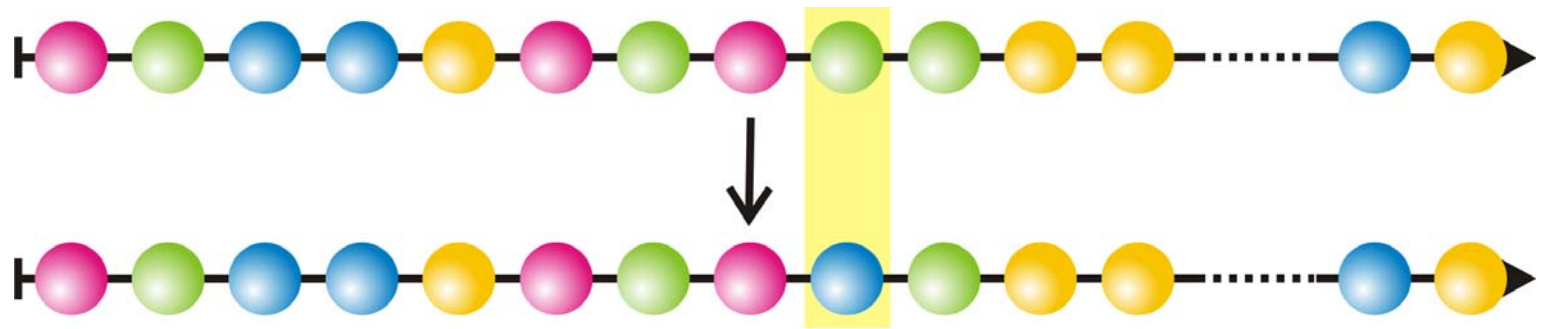
AT, TA, CG, and GC



The structure of DNA suggests a mechanism for reproduction



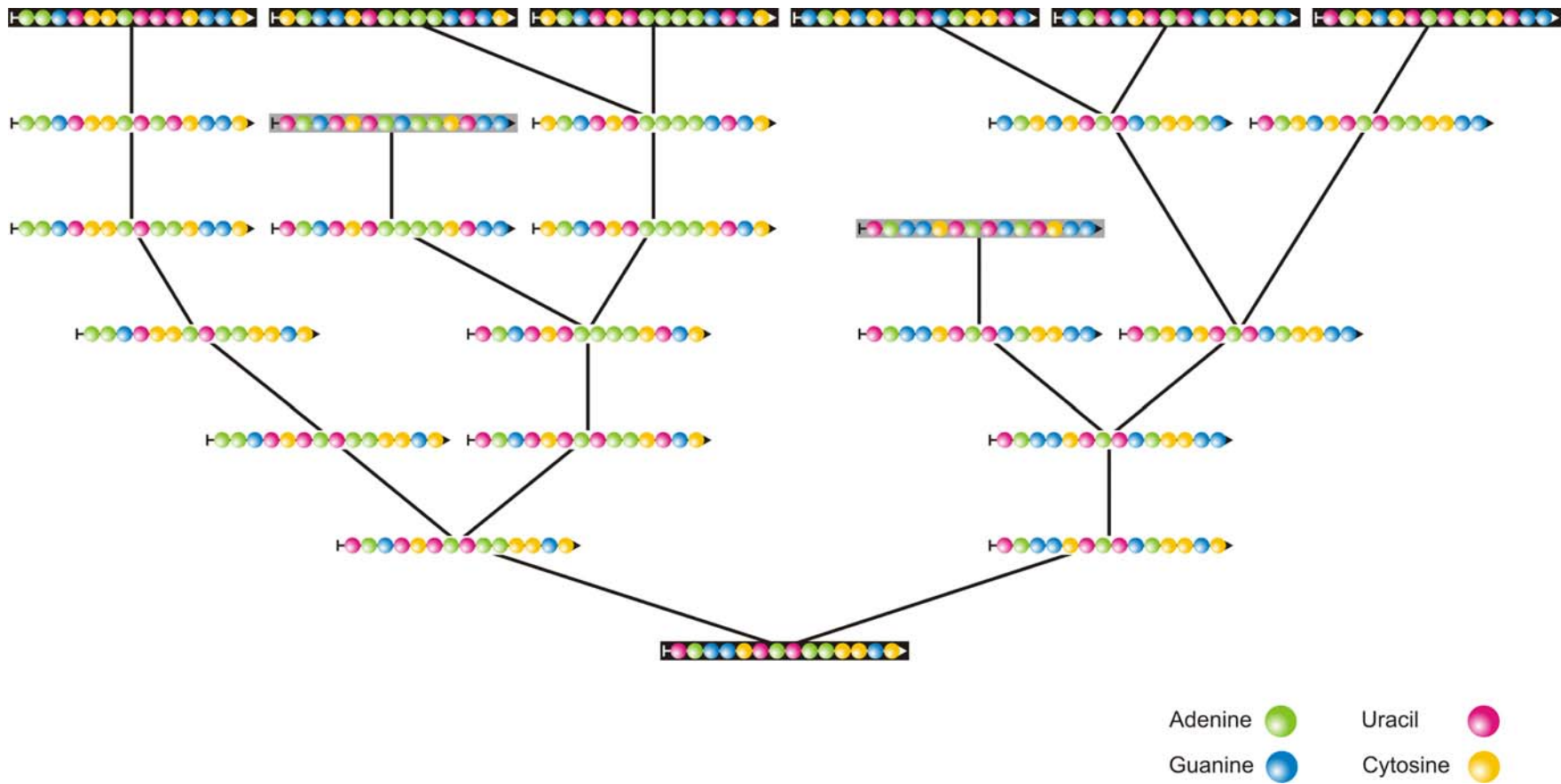
The logics of DNA replication



Adenine  Thymine  Guanine  Cytosine 

point mutation

The molecular mechanism of mutation



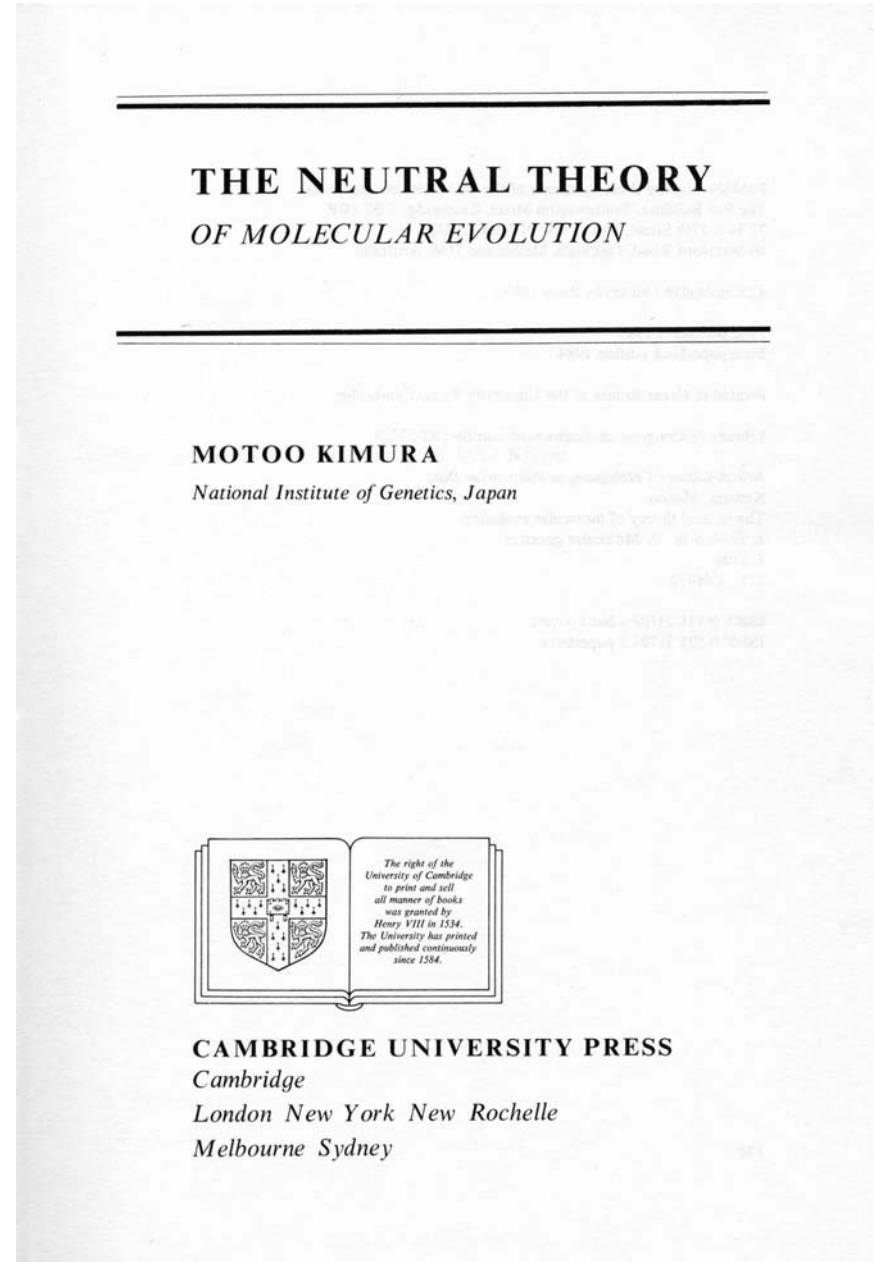
Molecular phylogeny



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.



What is neutrality ?

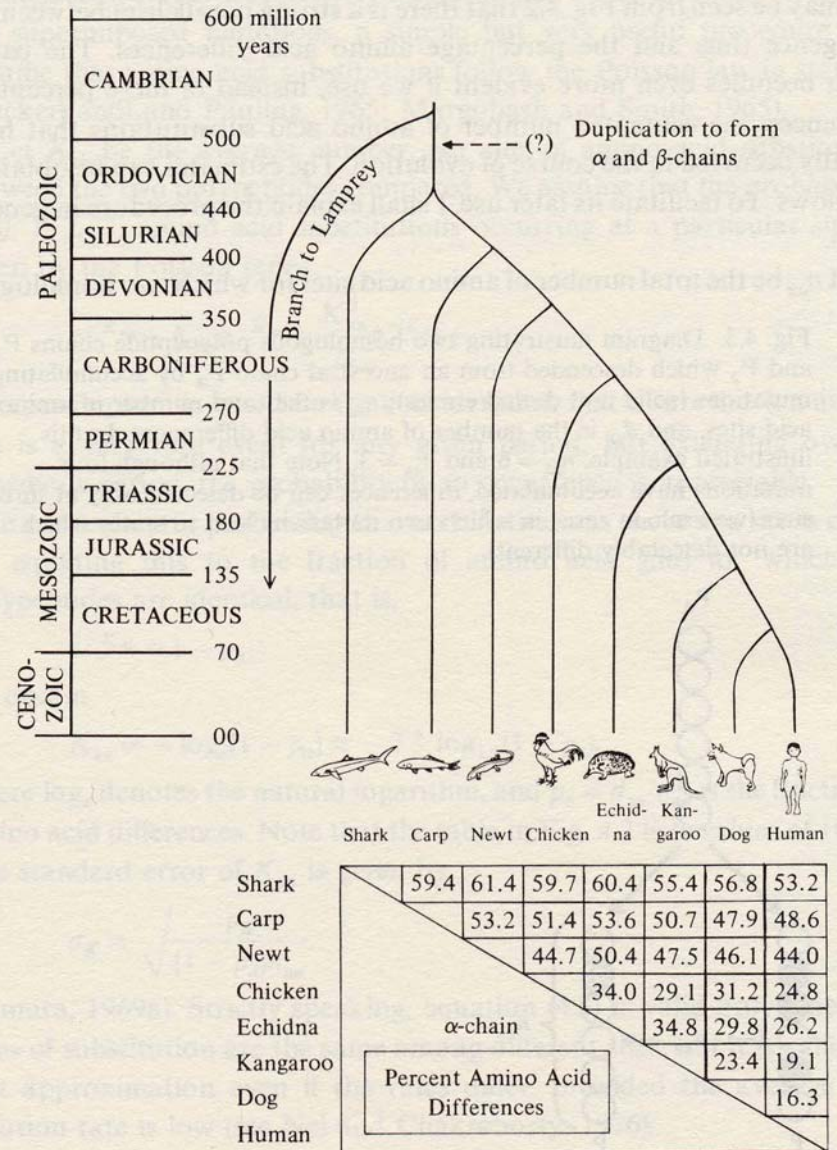
Selective neutrality =
= several genotypes having the **same fitness**.

Several genotypes \Rightarrow one phenotype

The molecular clock of evolution

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

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.



Results from molecular evolution:

- The molecular machineries of all present day cells are very similar and provide a strong hint that all life on Earth descended from one common ancestor (called „last universal common ancestor“, **LUCA**).
- Comparison of DNA sequences from present day organisms allows for a reconstruction of phylogenetic trees, which are (almost) identical with those derived from morphological comparison of species and the paleontologic record of fossils.

1. Darwin's natural selection
2. The tree of life
- 3. From evolution *in vitro* to biotechnology**
4. Genotypes with multiple functions
5. How complex is biology?



Three necessary conditions for Darwinian evolution are:

1. **Multiplication,**
2. **Variation,** and
3. **Selection.**

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

One important property of the Darwinian scenario is that **variations** in the form of mutations or recombination events occur **uncorrelated** with their **effects on the selection process**.

All conditions can be fulfilled not only by cellular organisms but also by **nucleic acid molecules** in suitable **cell-free experimental assays**.

Evolution of RNA molecules based on Q β 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.Öhlenschläger, 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

Molecular Evolution

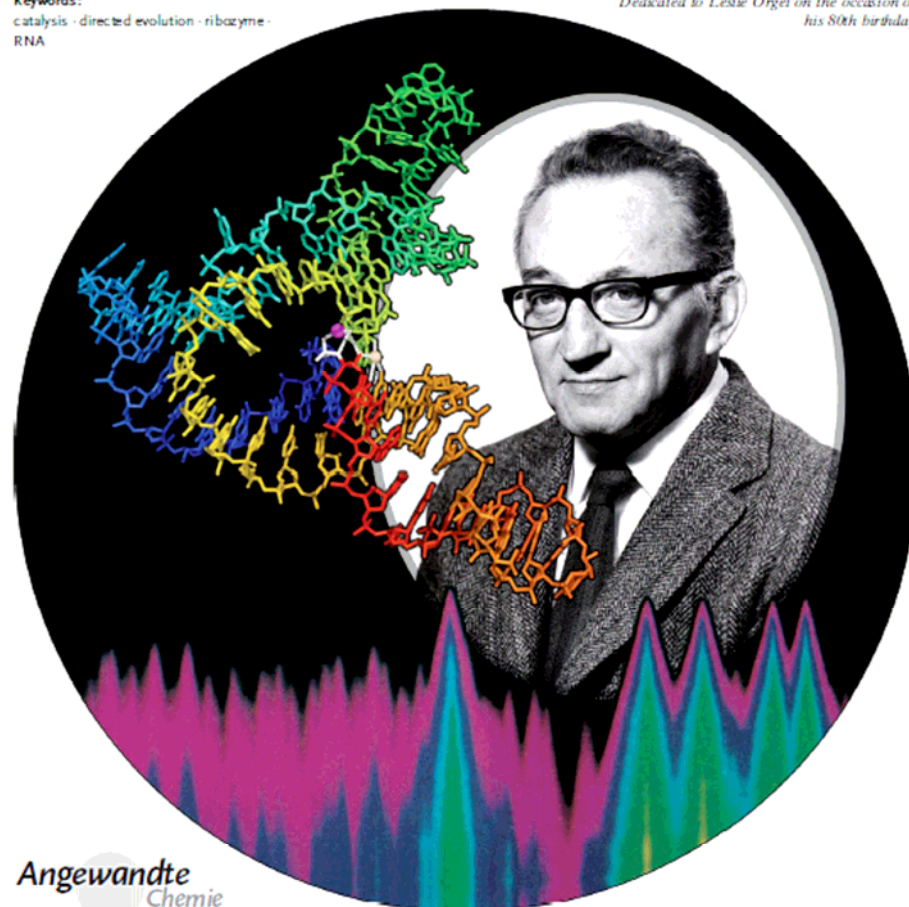
DOI: 10.1002/anie.200701369

Forty Years of In Vitro Evolution**

Gerald F. Joyce*

Keywords:
catalysis · directed evolution · ribozyme ·
RNA

*Dedicated to Leslie Orgel on the occasion of
his 80th birthday*



Evolution in the test tube:

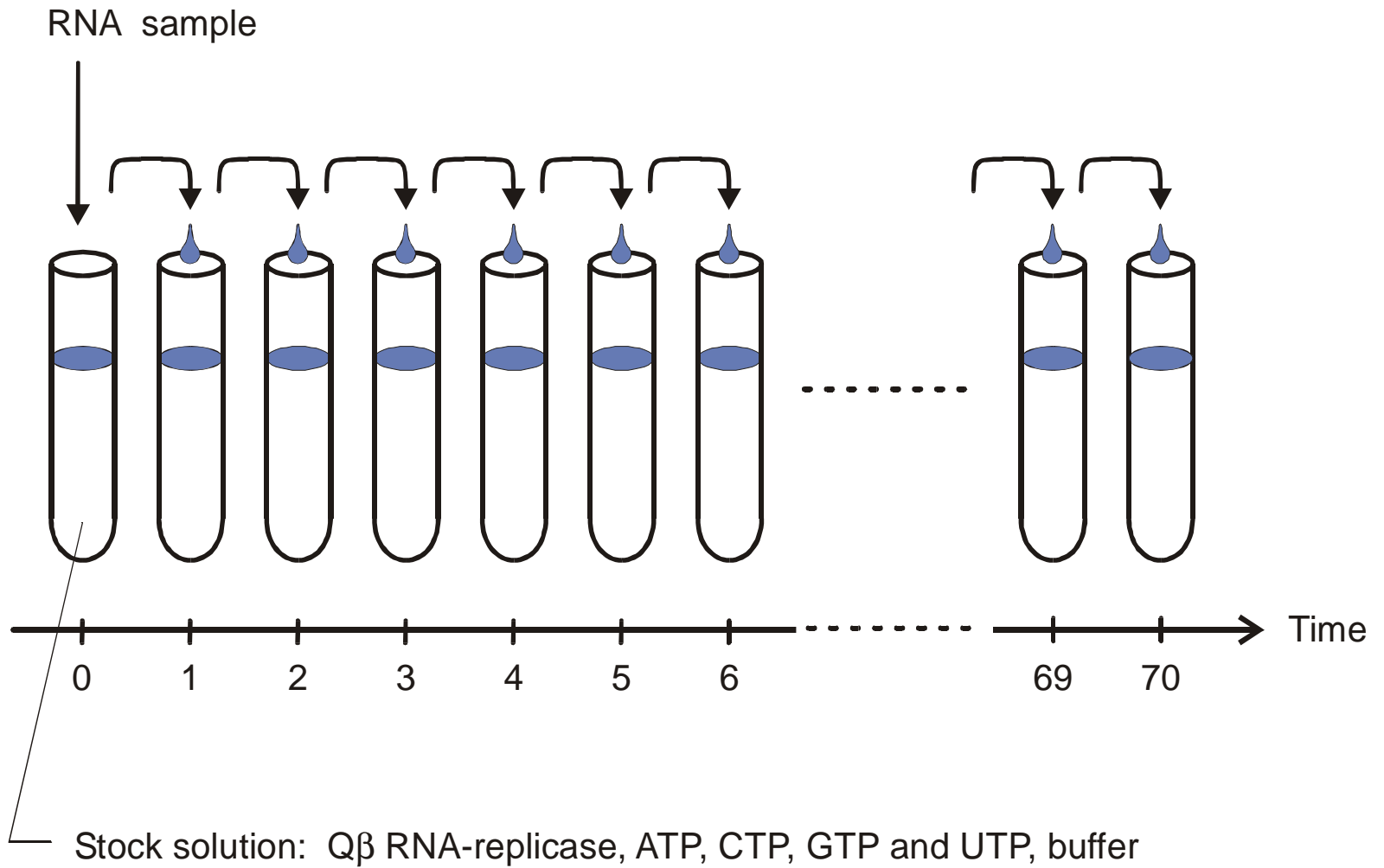
G.F. Joyce, *Angew.Chem.Int.Ed.*
46 (2007), 6420-6436

Angewandte
Chemie

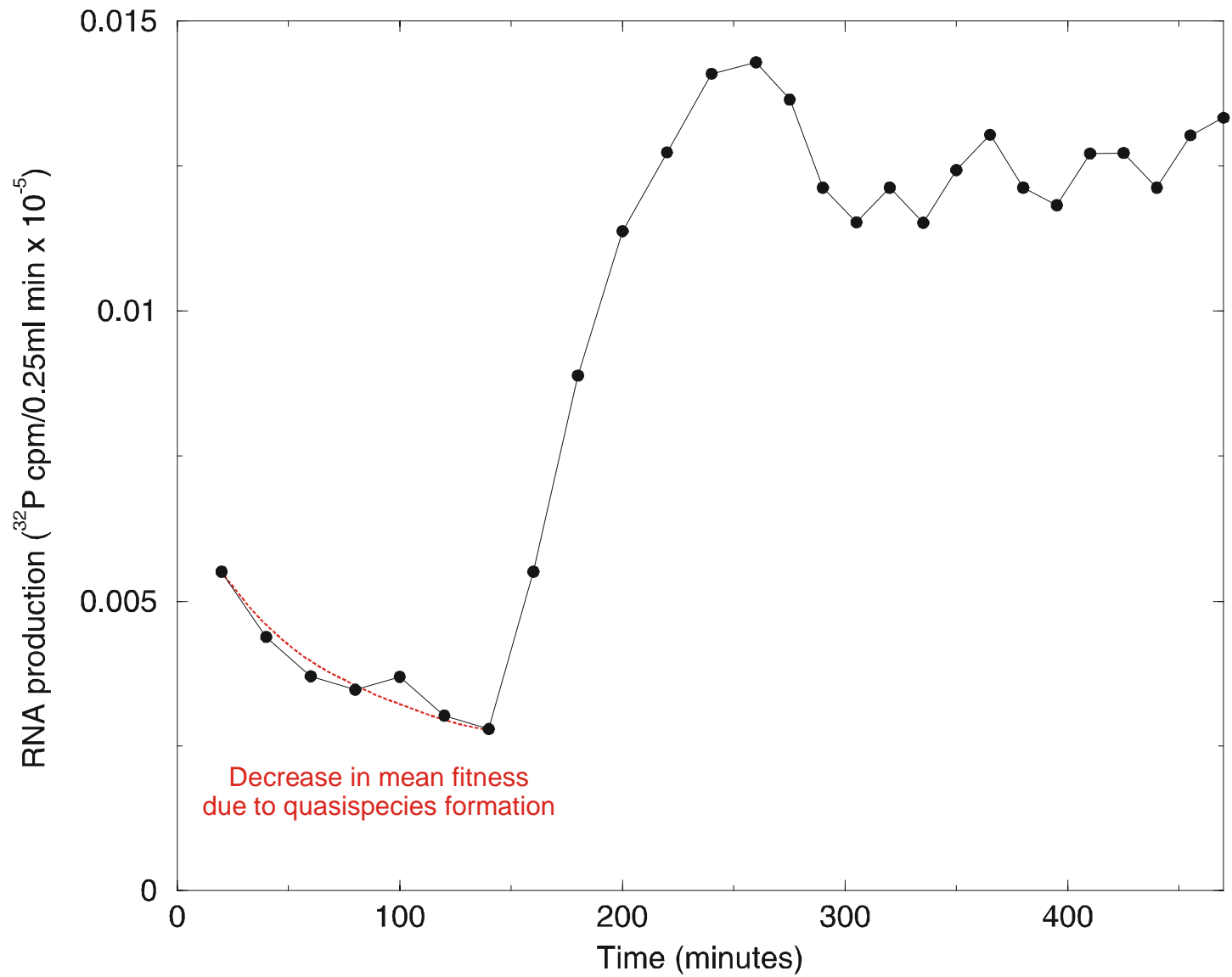
6420 www.angewandte.org

© 2007 Wiley-VCH Verlag GmbH & Co. KGaA, Weinheim

Angew. Chem. Int. Ed. 2007, 46, 6420-6436



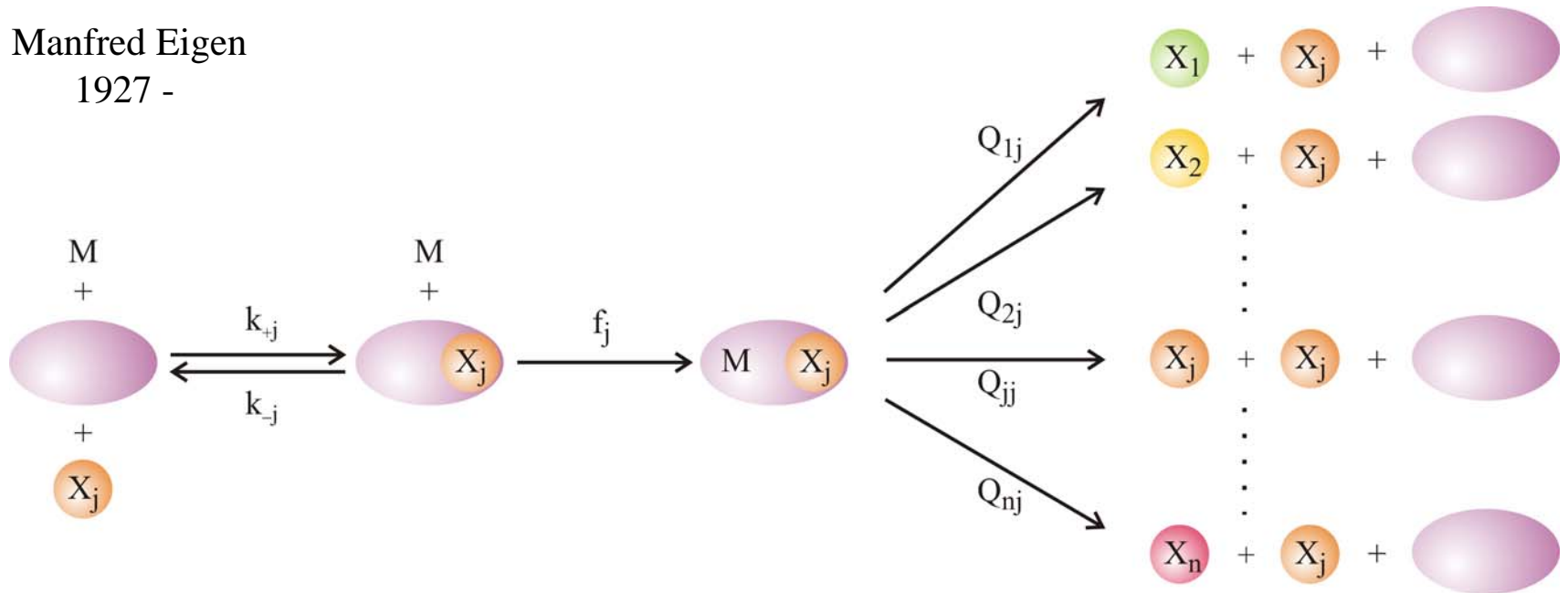
Application of serial transfer technique to evolution of RNA in the test tube



The increase in RNA production rate during a serial transfer experiment



Manfred Eigen
1927 -



Mutation and (correct) replication as parallel chemical reactions

M. Eigen. 1971. *Naturwissenschaften* 58:465,

M. Eigen & P. Schuster. 1977. *Naturwissenschaften* 64:541, 65:7 und 65:341

Selforganization of Matter and the Evolution of Biological Macromolecules

MANFRED EIGEN*

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

I. Introduction
1.1. Cause and Effect
1.2. Penetration of Self-Organization
1.2.1. Evolution Must Start from Random Events
1.2.2. Instructive Requires Information
1.2.3. Information Obligates or Gains Value by Selection
1.2.4. Selection Occurs with Special Instances under Special Conditions
II. Phenomenological Theory of Selection
II.1. The Concept "Information"
II.2. Phenomenological Equations
II.3. Selection Criteria
II.4. Selection Equilibrium
II.5. Quality Factor and Error Distribution
II.6. Kinetics of Selection
III. Stochastic Approach to Selection
III.1. Limitations of a Deterministic Theory of Selection
III.2. Fluctuations around Equilibrium States
III.3. Fluctuations in the Steady State
III.4. Stochastic Models in Markov Chains
III.5. Quantitative Discussion of Three Prototypes of Selection
IV. Self-Organization Based on Complementary Interactions: Nucleic Acids
IV.1. True Self-Organization
IV.2. Complementary Interaction and Selection
IV.3. Complementary Base Recognition (Experimental Data)
IV.3.1. Single Pair Formation
IV.3.2. Cooperative Interactions in Oligo- and Polynucleotides
IV.3.3. Conclusions about Recognition

I. Introduction
1.1. "Cause and Effect"

which even in its simplest forms always appears to be associated with complex macroscopic (i.e. multimolecular) systems, such as the living cell. As a consequence of the exciting discoveries of "molecular biology", a common version of the above question is: Which came first, the protein or the nucleic acid?—a modern variant of the old "chicken-and-egg" problem. The term "first" is usually meant to define a causal rather than a temporal relationship, and the words "protein" and "nucleic acid" may be substituted by "function" and "information". The question in this form, when applied to the interplay of nucleic acids and proteins as presently encountered in the living cell, leads to absurdum, because "function"

* Fully presented at the "Robbins Lectures" at Pomona College, California, in spring 1970.

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

I. Introduction
1.1. Cause and Effect
1.2. Penetration of Self-Organization
1.2.1. Evolution Must Start from Random Events
1.2.2. Instructive Requires Information
1.2.3. Information Obligates or Gains Value by Selection
1.2.4. Selection Occurs with Special Instances under Special Conditions
II. Phenomenological Theory of Selection
II.1. The Concept "Information"
II.2. Phenomenological Equations
II.3. Selection Criteria
II.4. Selection Equilibrium
II.5. Quality Factor and Error Distribution
II.6. Kinetics of Selection
III. Stochastic Approach to Selection
III.1. Limitations of a Deterministic Theory of Selection
III.2. Fluctuations around Equilibrium States
III.3. Fluctuations in the Steady State
III.4. Stochastic Models in Markov Chains
III.5. Quantitative Discussion of Three Prototypes of Selection
IV. Self-Organization Based on Complementary Interactions: Nucleic Acids
IV.1. True Self-Organization
IV.2. Complementary Interaction and Selection
IV.3. Complementary Base Recognition (Experimental Data)
IV.3.1. Single Pair Formation
IV.3.2. Cooperative Interactions in Oligo- and Polynucleotides
IV.3.3. Conclusions about Recognition

This paper is the first part of a trilogy, which comprises a detailed study of a special type of functional organization and demonstrates its relevance with respect to the origin and evolution of life. Self-replicating macromolecules, such as RNA or DNA in a suitable environment exhibit a behavior, which we may call Darwinian and which can be formally represented by the concept of the quasi-species. A quasi-species is defined as a given distribution of macromolecular species with closely interrelated sequences, dominated by one or several (hypothesized) master copies. External constraints enforce the selection of the best adapted distribution, outcompetitively referred to as the wild-type. Most important for Darwinian behavior are the criteria for internal stability of the quasi-species. If these criteria are violated, the information stored in the nucleotide sequence of the master copy will disseminate irreversibly leading to an error catastrophe. As a consequence, selection and evolution of RNA or DNA molecules is limited with respect to the amount of information that can be stored in a single replicative unit. An analysis of experimental data regarding RNA and DNA replication at various levels of organization reveals that a sufficient amount of information for the build up of a translation machinery can be gained only via integration of several different replicative units (replicative cycles) through reciprocal linkages. A stable functional organization then will cause the system to a new level of organization and thereby enlarge its information capacity correspondingly. The hypercycle appears to be such a form of organization.

Preview on Part C: The Abiotic Hypercycle

The mathematical analysis of dynamical systems using methods of differential topology yields the result that there is only one type of mechanism which fulfills the following requirements: The information stored in each single replicative unit (or replicative cycle) must be maintained, i.e., the respective master copies must cooperate favorably with their error distributions despite their competitive behavior; these units must establish a cooperation head, the cycle as a whole must consist to emerge already with any other single entity or isolated ensemble which does not contribute to its sustained function. These requirements are crucial for a selection of the best adapted functionally linked ensemble and its evolutionary optimization. Only hypercyclic organizations are able to fulfill these requirements. Non-specific linkages among the autonomous reproduction cycles, such as those of branched, one-line networks are devoid of such properties. The mathematical methods used for proving these assertions are fixed-point, Lyapunov and trajectory analysis in high-dimensional phase spaces, spanned by the concentration coordinates of the cooperating partners. The self-organizing properties of hypercycles are elucidated, using analytical as well as numerical techniques.

I. The Paradigm of Unity and Diversity in Evolution

Why do millions of species, plants and animals, exist, while there is only one basic molecular machinery of the cell: one universal genetic code and unique chemicalities of the macromolecules? The generalists of our day would not hesitate to give an immediate answer to the first part of this question: Diversity of species is the outcome of the tremendous branching process of evolution with its myriads of single steps of reproduction and mutation. It in-

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)

The molecular quasi-species model describes the physicochemical organization of monomers into an ensemble of heteropolymers with combinatorial complexity by ongoing template polymerization. Polynucleotides belong to the simplest class of such molecules. The quasi-species line represents the stationary distribution of macromolecular sequences maintained by chemical reaction effecting error-prone replication and by transport processes. It is obtained deterministically, by mass-action kinetics, as the dominant eigenvector of a square matrix, W, which is derived directly from chemical rate coefficients, but it also exhibits stochastic features, being composed of a significant fraction of unique individual macromolecular sequences. The quasi-species model demonstrates how macromolecular information originates through specific non-equilibrium autocatalytic reactions and thus forms a bridge between reaction kinetics and molecular evolution. Selection and evolutionary optimization appear as new features in physical chemistry. Concentration bias in the production of mutants is a new concept in population genetics, relevant to frequently mating populations, which is shown to greatly enhance the optimization process. The present theory relates to naturally replicating ensembles, but this restriction is not essential. A sharp transition is exhibited between a drifting population of essentially random macromolecular sequences and a localized population of close relatives. This transition at a threshold error rate was found to depend on sequence lengths, distributions of selective values, and population sizes. It has been determined generally for complex landscapes and for special cases, and, it was shown to persist generally in the presence of nearly neutral mutants. Replication dynamics has much in common with the equilibrium statistics of complex spin systems: the error threshold is equivalent to a magnetic order-disorder transition. A rational function of the replication accuracy plays the role of temperature. Experimental data obtained from test-tube evolution of polynucleotides and from studies of natural virus populations support the quasi-species model. The error threshold seems to set a limit to the genome lengths of several classes of RNA viruses. In addition, the results are relevant even in eucaryotes where they contribute to the exon-intron debate.

1. Molecular Selection

Our knowledge of physical and chemical systems is, in a final analysis, based on models derived from repeatable experiments. While none of the classic and rather besieged list of properties rounded up to support the intuition of a distinction between the living and nonliving—metabolism, self-reproduction, irritability, and adaptability, for example—intrinsically limit the application of the scientific method, a determining role by unique or individual entities comes into conflict with the requirement of repeatability. Combinatorial variety, such as that in heteropolymers based on even very small numbers of different bases, even just two, readily provides numbers of different entities so enormous that neither consecutive nor parallel physical realization is possible. The physical chemistry of finite systems of such macromolecules must deal with both known regularities and the advent of unique copolymeric sequences. Normally this would present no difficulty in a statistical mechanical analysis of typical behavior, where rare events play no significant role, but with autocatalytic polymerization processes even unique single molecules may be singled out to determine the fate of the entire system. Potentially creative, self-organizing around unique events, the dynamics of the simplest living chemical system is invested with regularities that both allow and limit efficient adaptation. The quasi-species model is a study of these regularities. The fundamental regularity in living organisms that has invited explanation is adaptation. Why are organisms so well fitted to their environments? At a more chemical level, why are enzymes

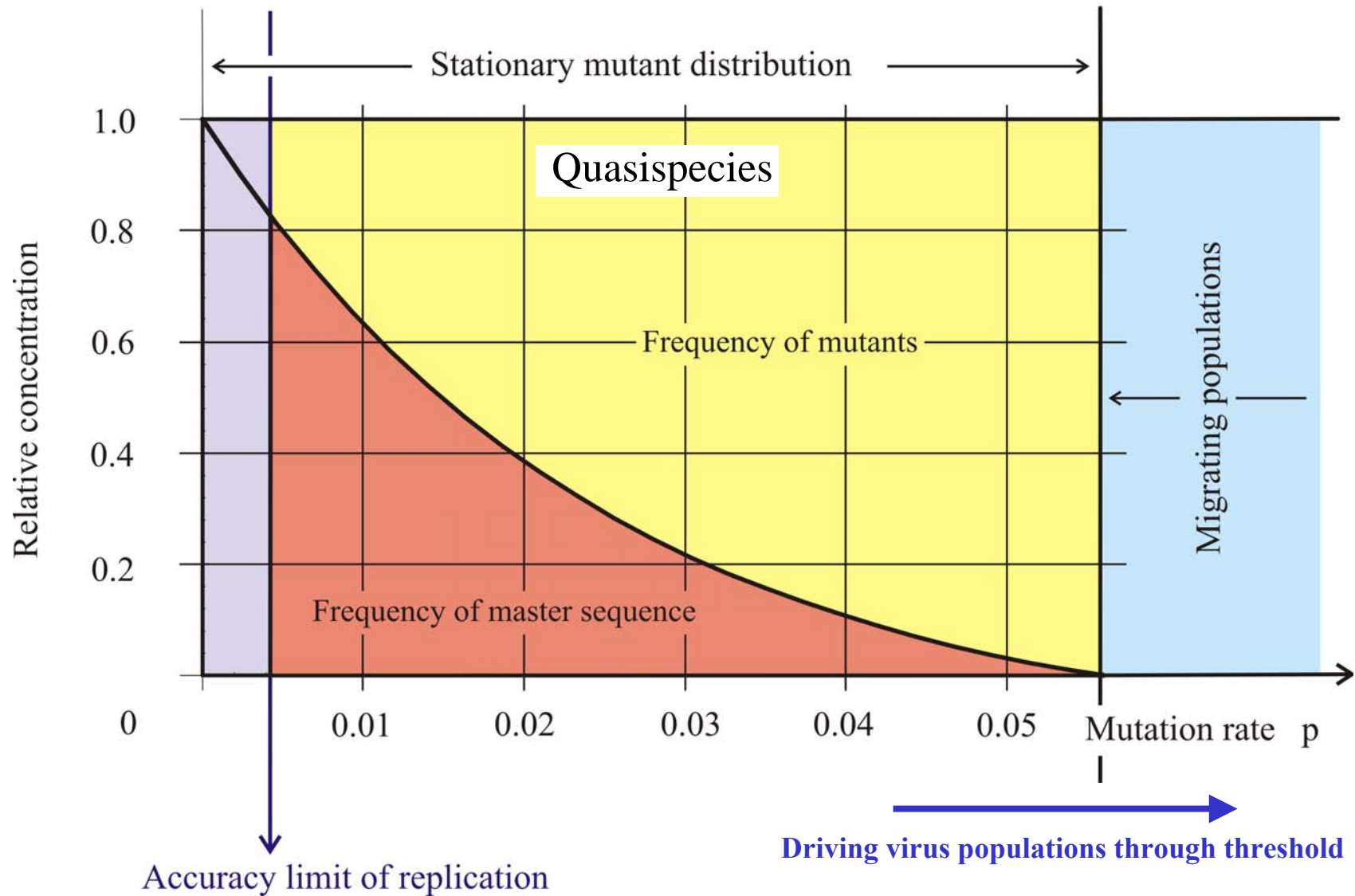
* This is an abridged account of the quasi-species theory that has been submitted in comprehensive form to Advances in Chemical Physics. (1) Eigen, M.; McCaskill, J.S.; Schuster, P. Adv. Chem. Phys., in press.

1971

1977

1988

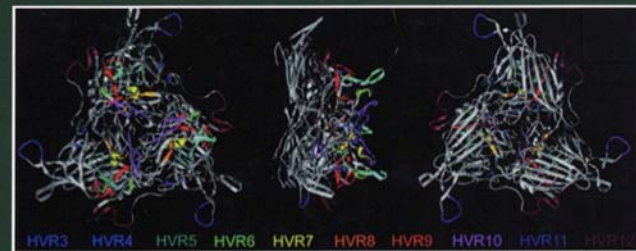
Chemical kinetics of molecular evolution



The error threshold in replication

SECOND EDITION

ORIGIN AND EVOLUTION OF VIRUSES



Edited by
ESTEBAN DOMINGO
COLIN R. PARRISH
JOHN J. HOLLAND



Molecular evolution of viruses

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

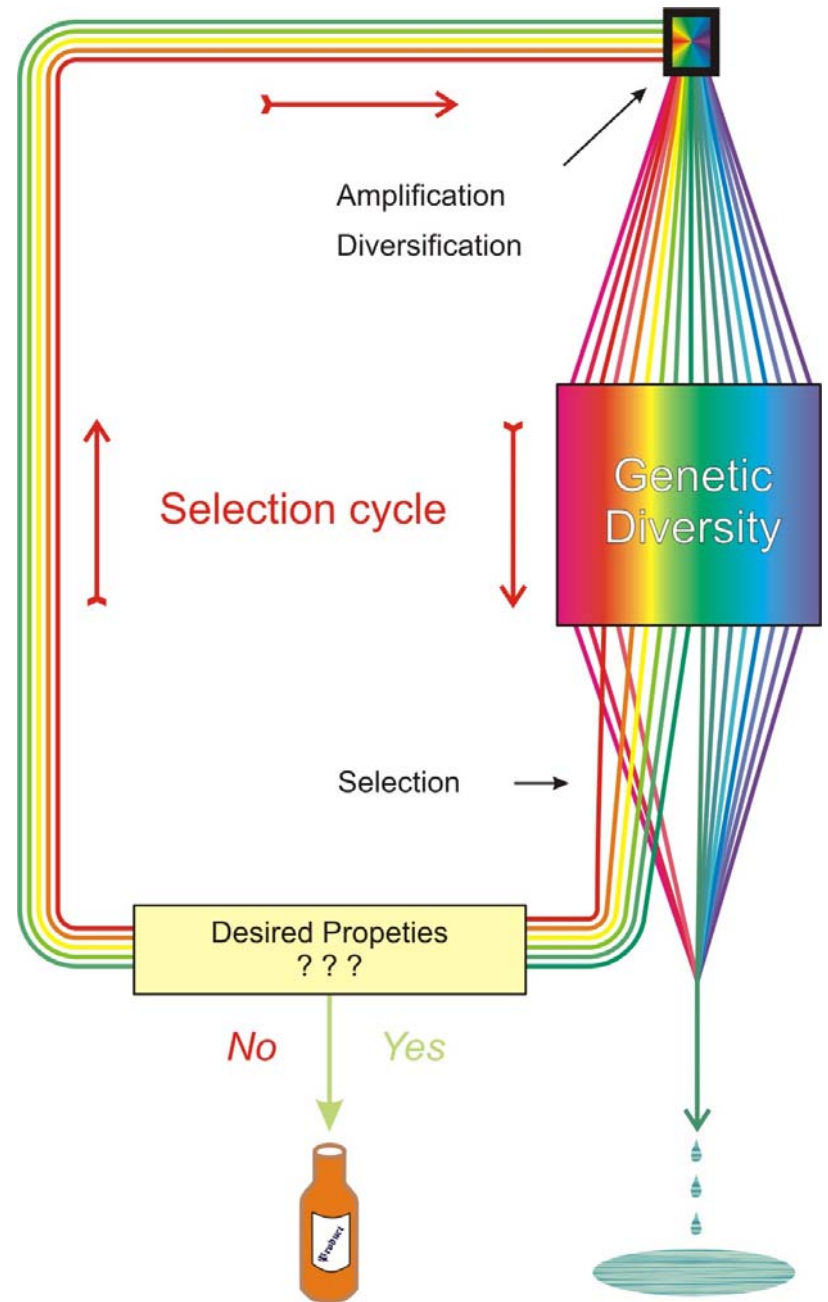
C. Tuerk, L. Gold, *SELEX - Systematic evolution of ligands by exponential enrichment: RNA ligands to bacteriophage T4 DNA polymerase*. Science **249** (1990), 505-510

D.P. Bartel, J.W. Szostak, *Isolation of new ribozymes from a large pool of random sequences*. Science **261** (1993), 1411-1418

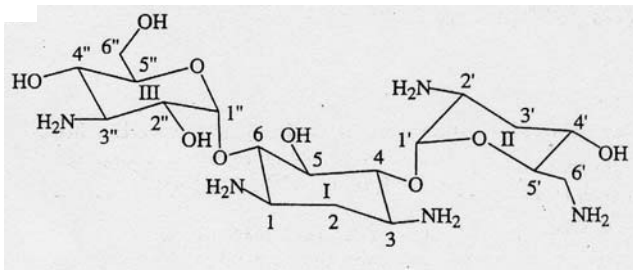
R.D. Jenison, S.C. Gill, A. Pardi, B. Poliski, *High-resolution molecular discrimination by RNA*. Science **263** (1994), 1425-1429

Y. Wang, R.R. Rando, *Specific binding of aminoglycoside antibiotics to RNA*. Chemistry & Biology **2** (1995), 281-290

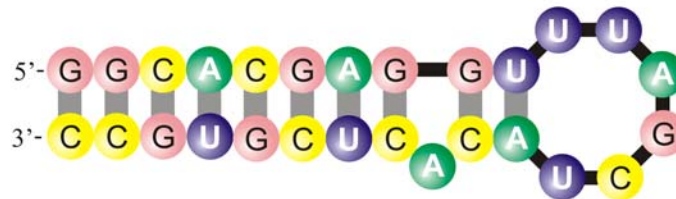
L. Jiang, A. K. Suri, R. Fiala, D. J. Patel, *Saccharide-RNA recognition in an aminoglycoside antibiotic-RNA aptamer complex*. Chemistry & Biology **4** (1997), 35-50



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



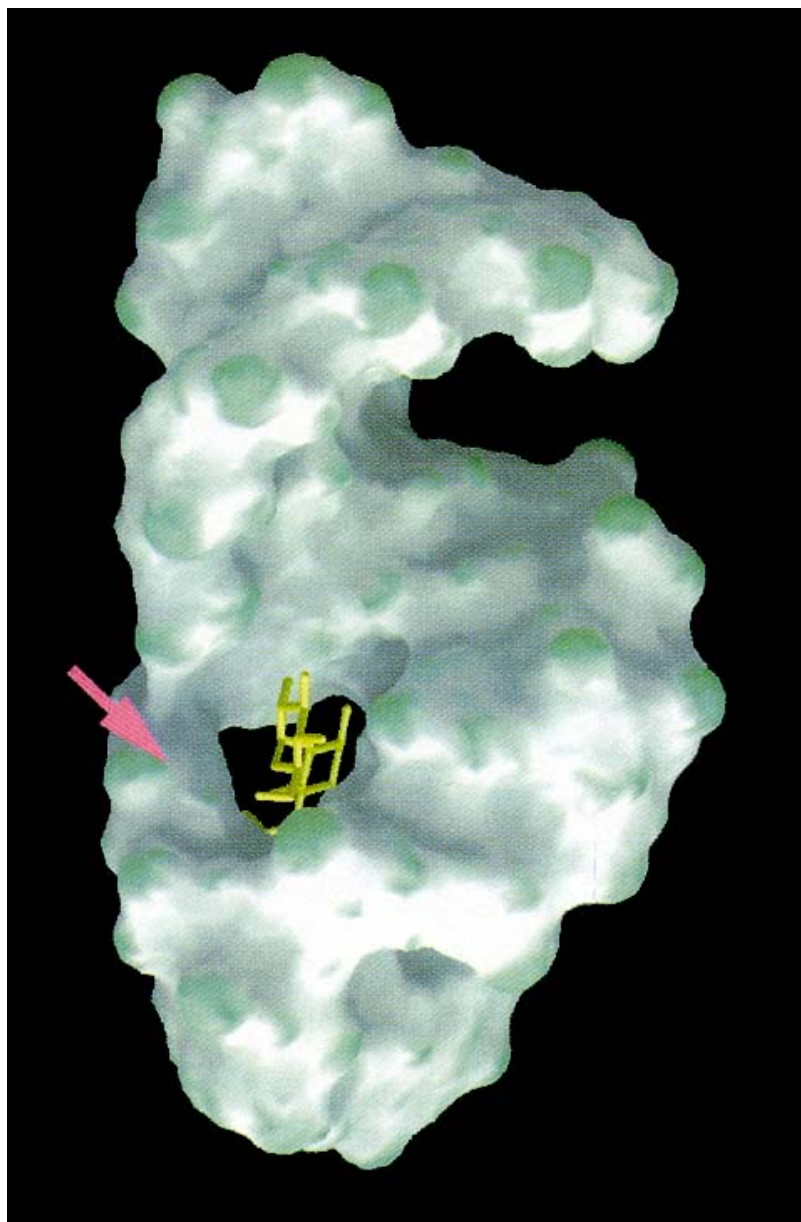
tobramycin



RNA aptamer, n = 27

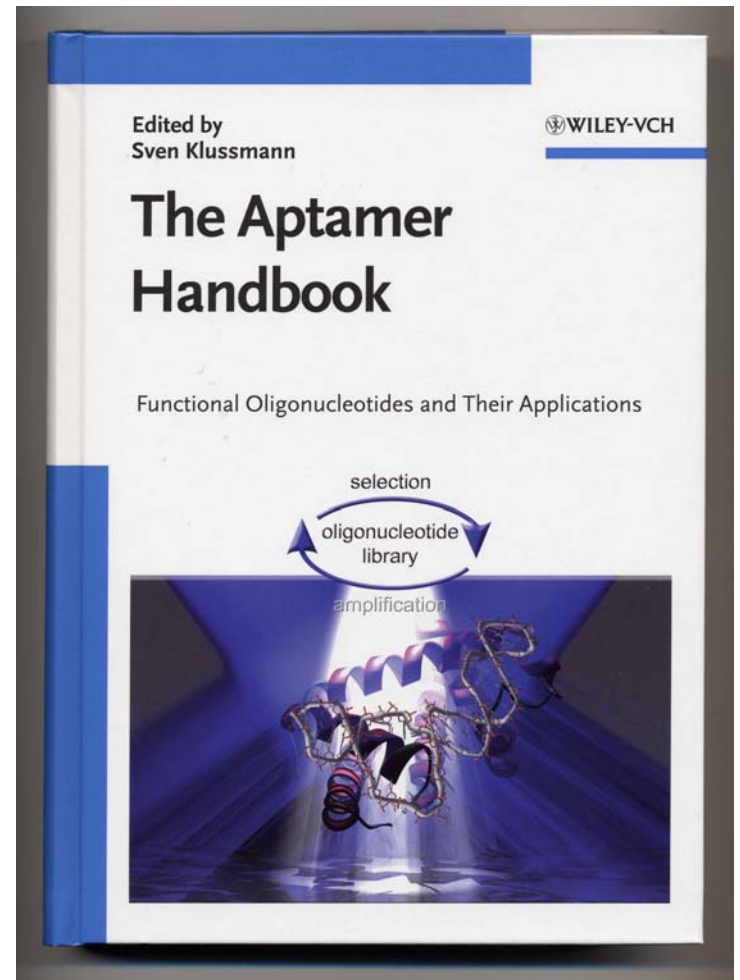
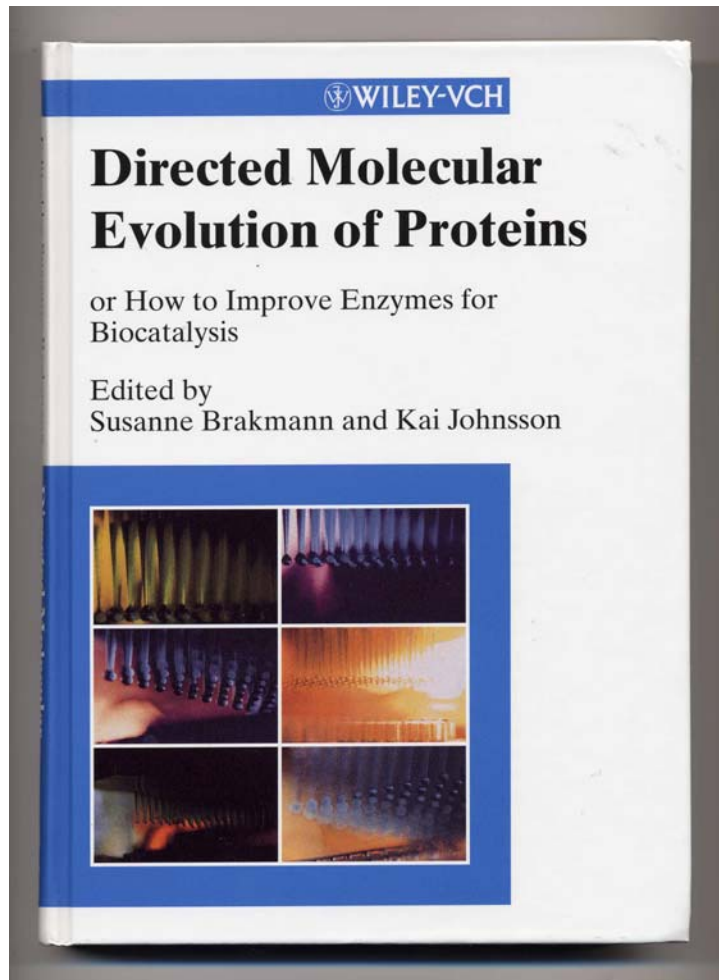
Formation of secondary structure of the tobramycin binding RNA aptamer with $K_D = 9 \text{ nM}$

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 laboratory experiments in molecular evolution:

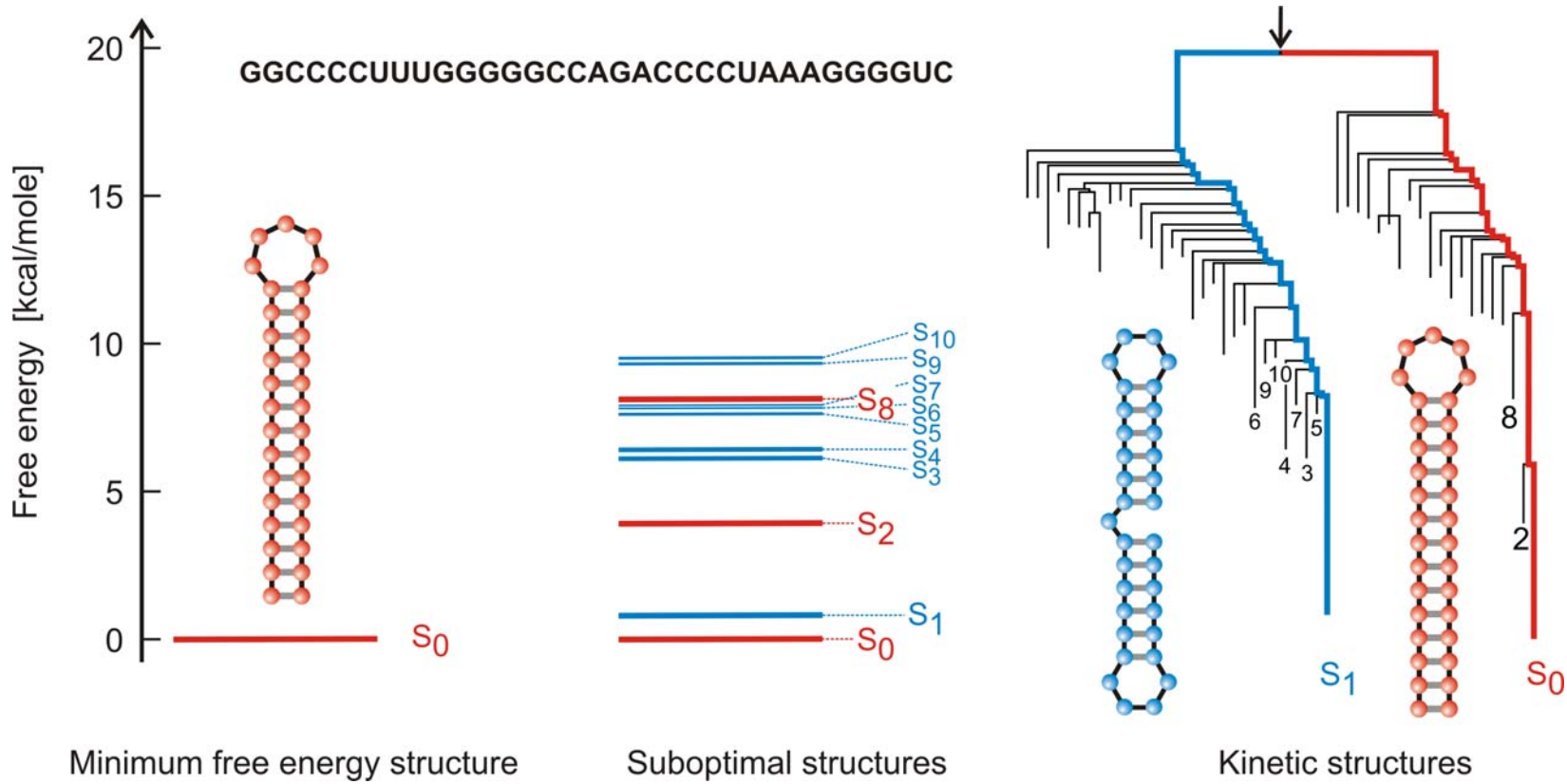
- Evolutionary optimization does not require cells and occurs in molecular systems too.
- *In vitro* evolution allows for production of molecules for predefined purposes and gave rise to a branch of biotechnology.
- Direct evidence that neutrality is a major factor for the success of evolution.

1. Darwin's natural selection
2. The tree of life
3. From evolution *in vitro* to biotechnology
4. **Genotypes with multiple functions**
5. How complex is biology?

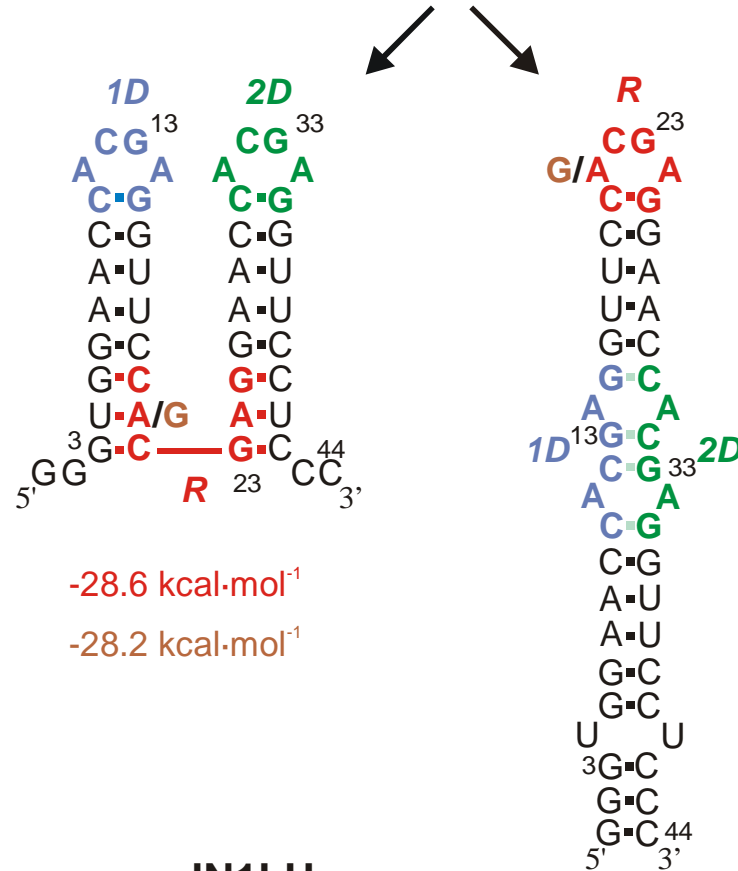
What is conformational multiplicity ?

Conformational multiplicity =
= several structures formed by one sequence.

One genotype \Rightarrow several phenotypes



Extension of the notion of structure



-28.6 kcal·mol⁻¹

-28.2 kcal·mol⁻¹

-28.6 kcal·mol⁻¹

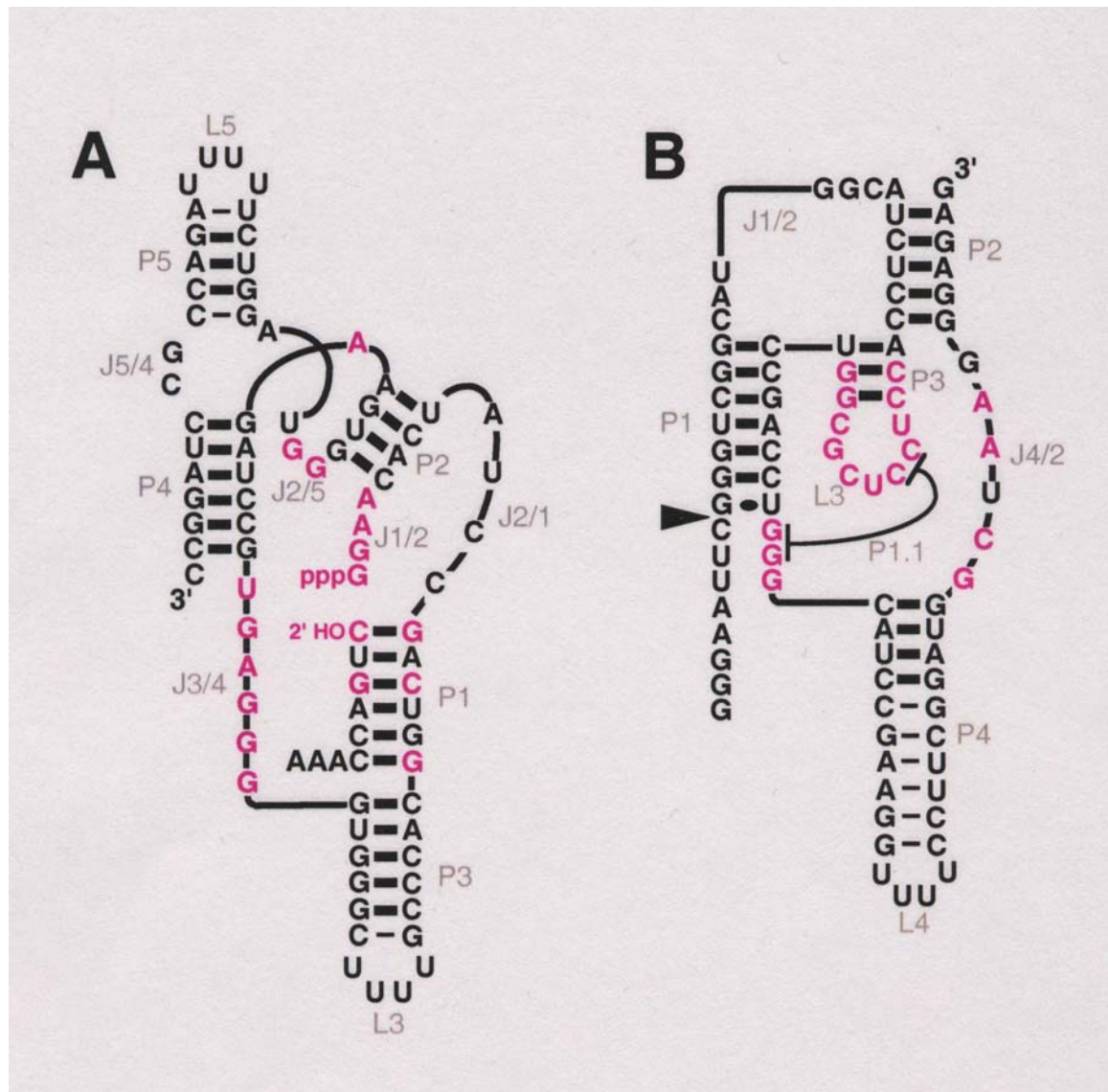
-31.8 kcal·mol⁻¹

An experimental RNA switch

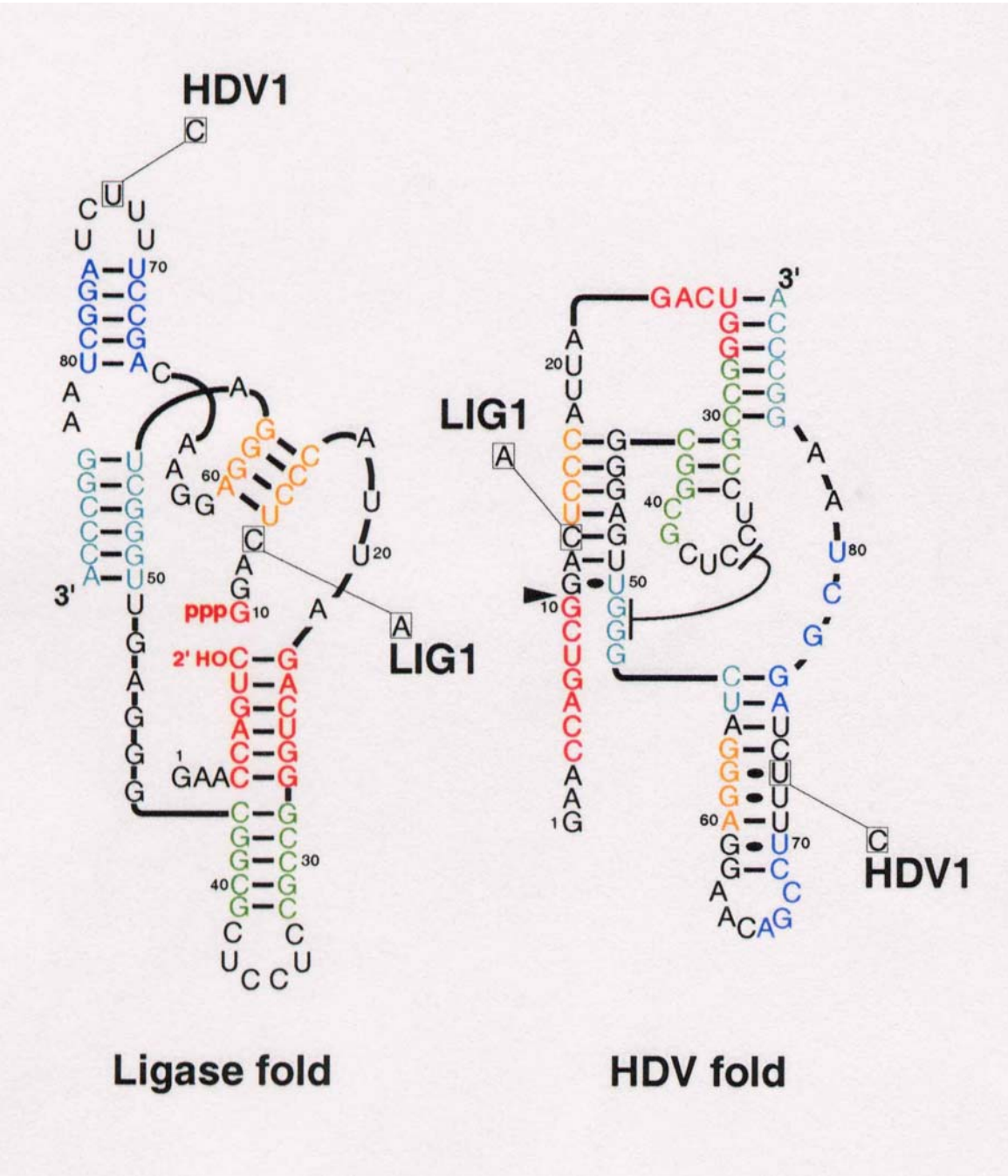
JN1LH

J.H.A. Nagel, C. Flamm, I.L. Hofacker, K. Franke,
M.H. de Smit, P. Schuster, and C.W.A. Pleij.

Structural parameters affecting the kinetic competition of RNA hairpin formation. *Nucleic Acids Res.* **34**:3568-3576 (2006)

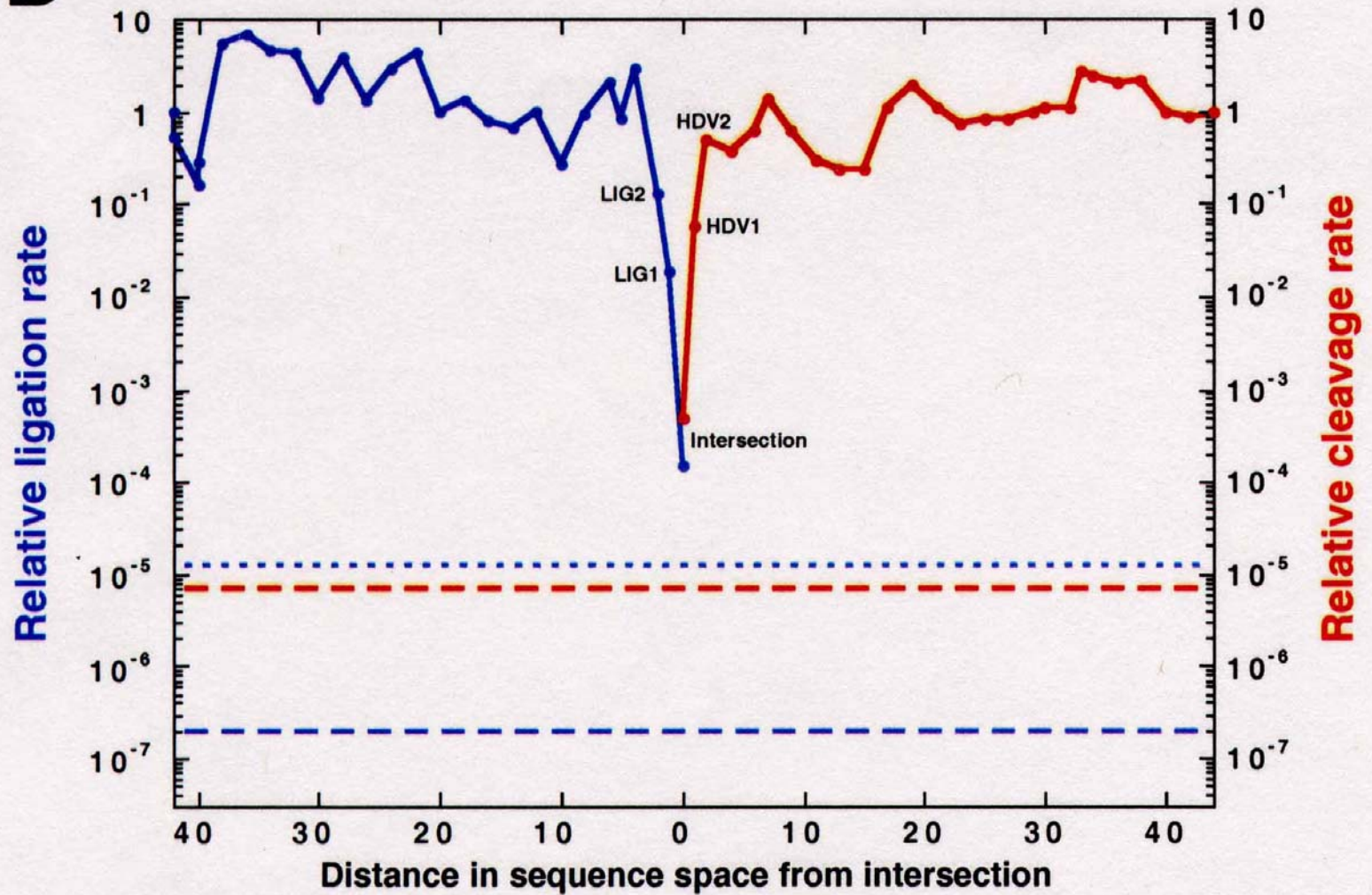


Two ribozymes of chain lengths $n = 88$ nucleotides: An artificial ligase (**A**) and a natural cleavage ribozyme of hepatitis- δ -virus (**B**)

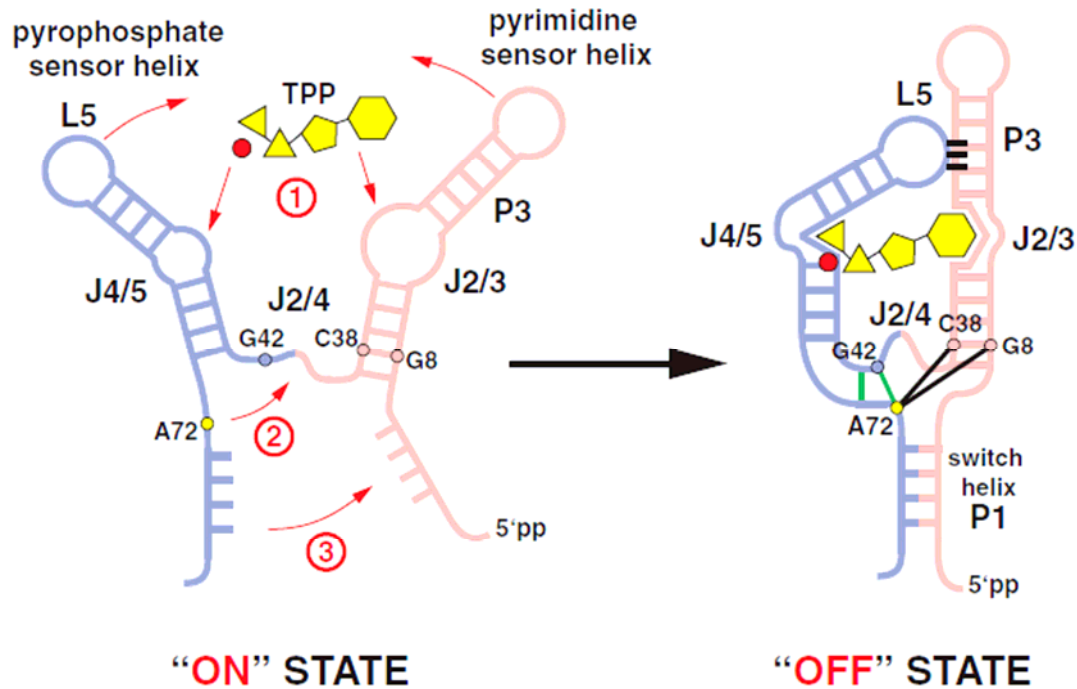


The sequence at the *intersection*:

An RNA molecules which is 88 nucleotides long and can form both structures

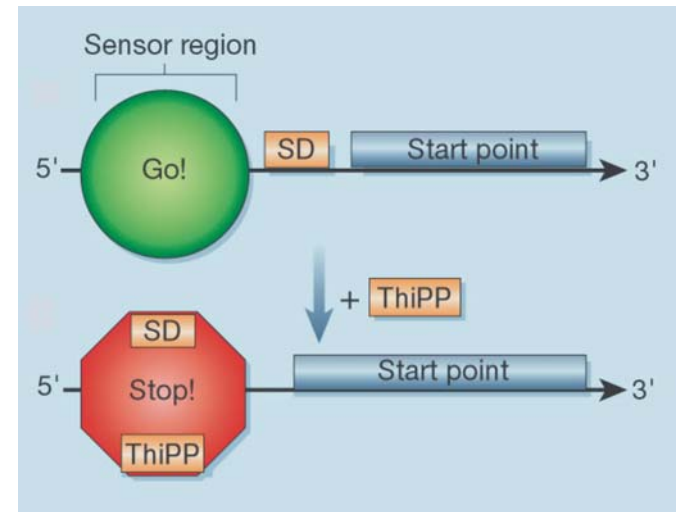
B

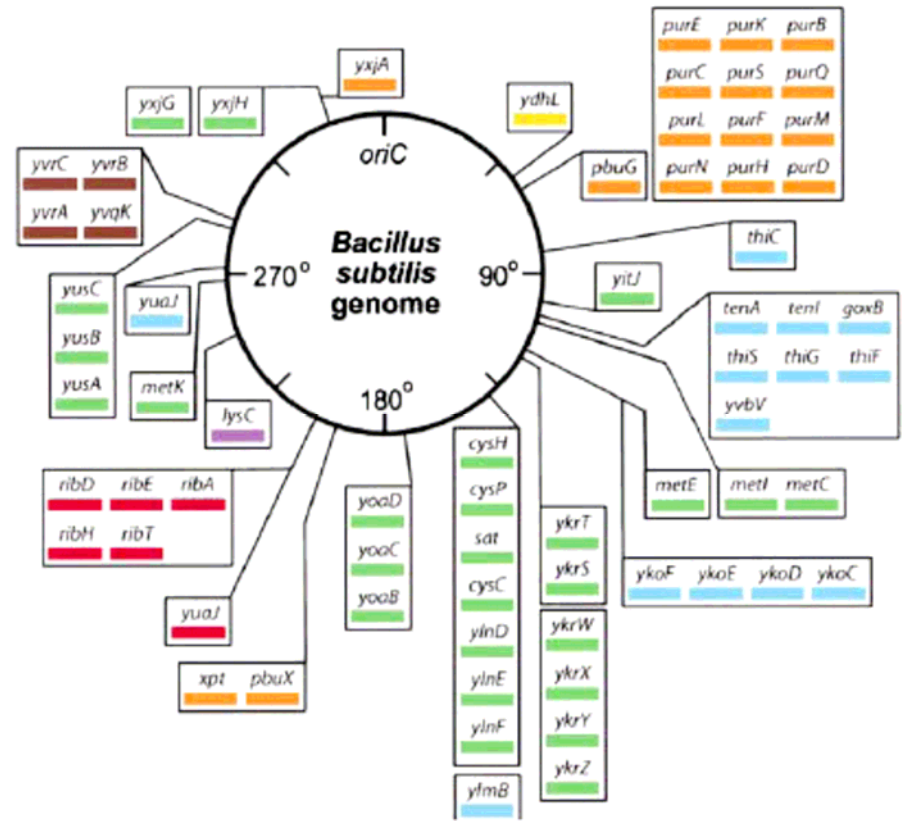
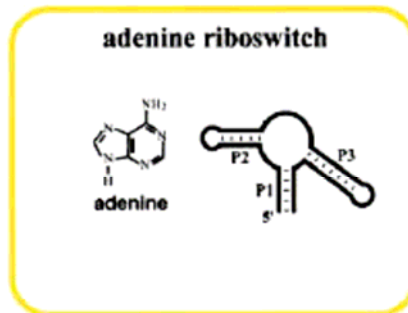
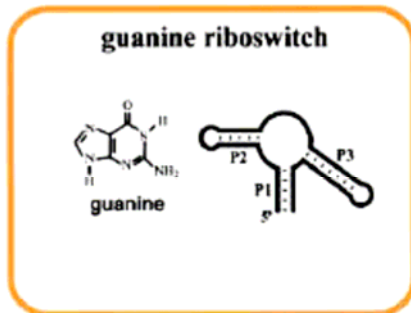
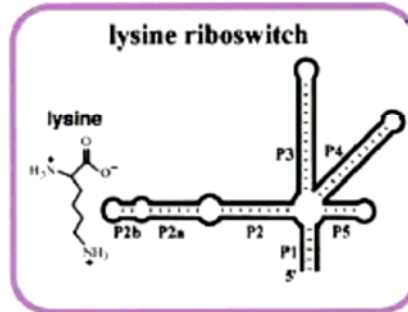
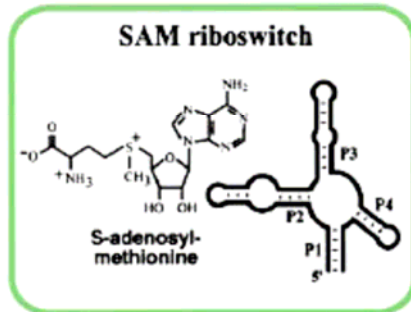
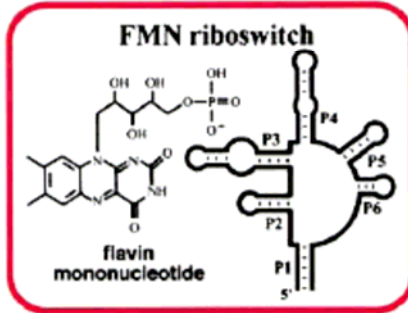
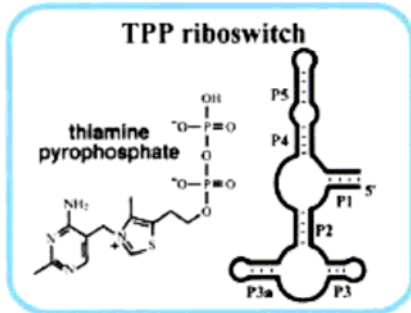
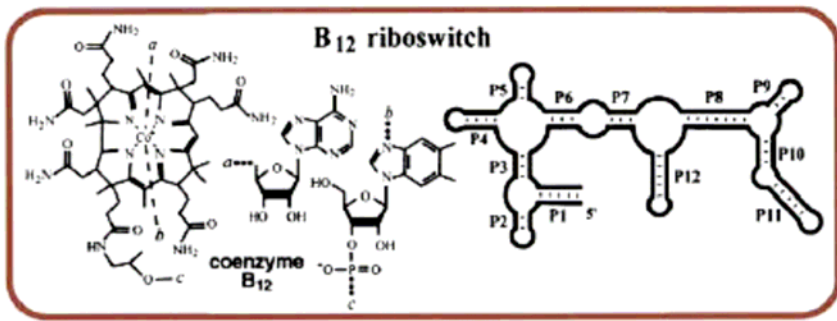
Two neutral walks through sequence space with conservation of structure and catalytic activity



The thiamine-pyrophosphate riboswitch

S. Thore, M. Leibundgut, N. Ban.
Science **312**:1208-1211, 2006.

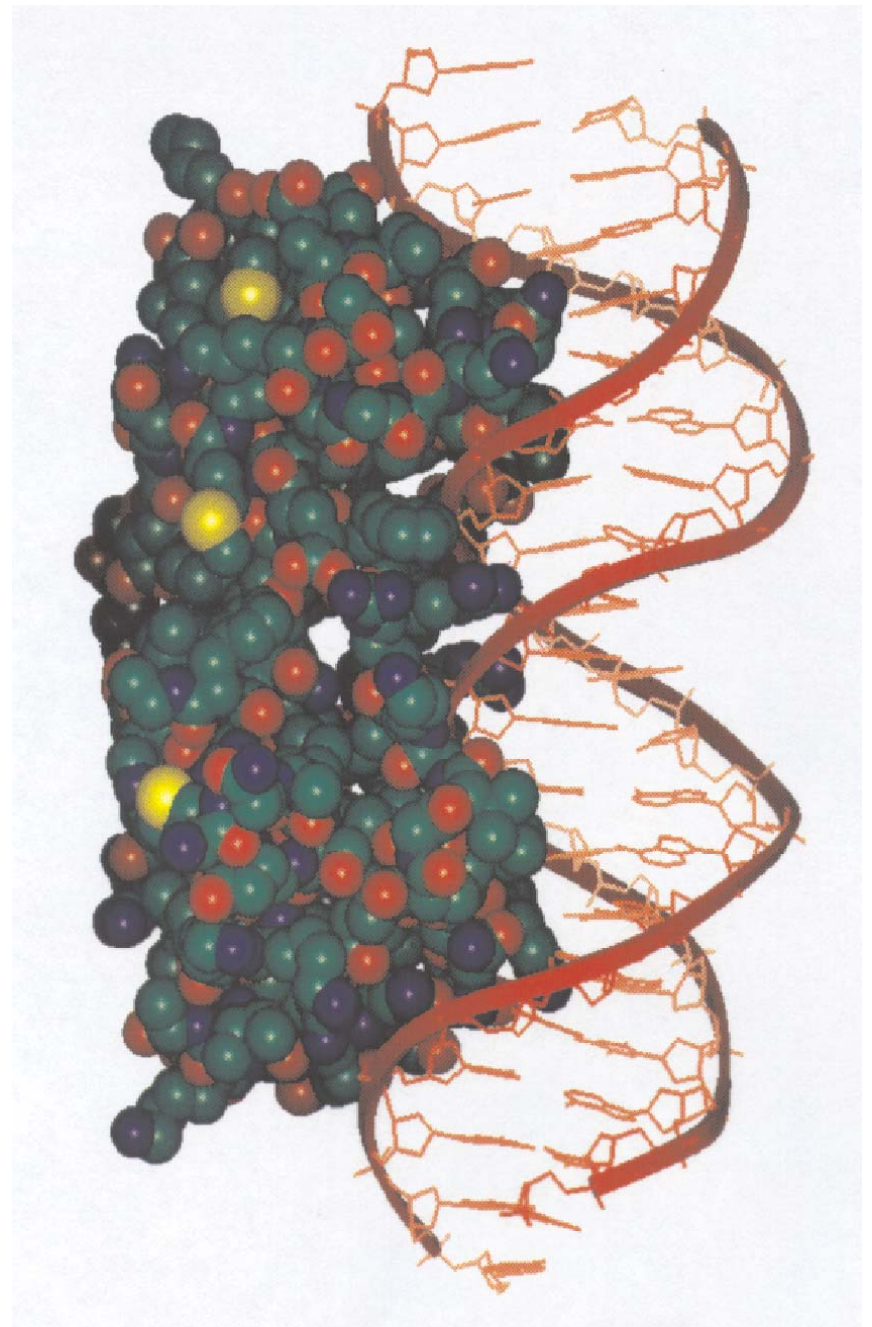




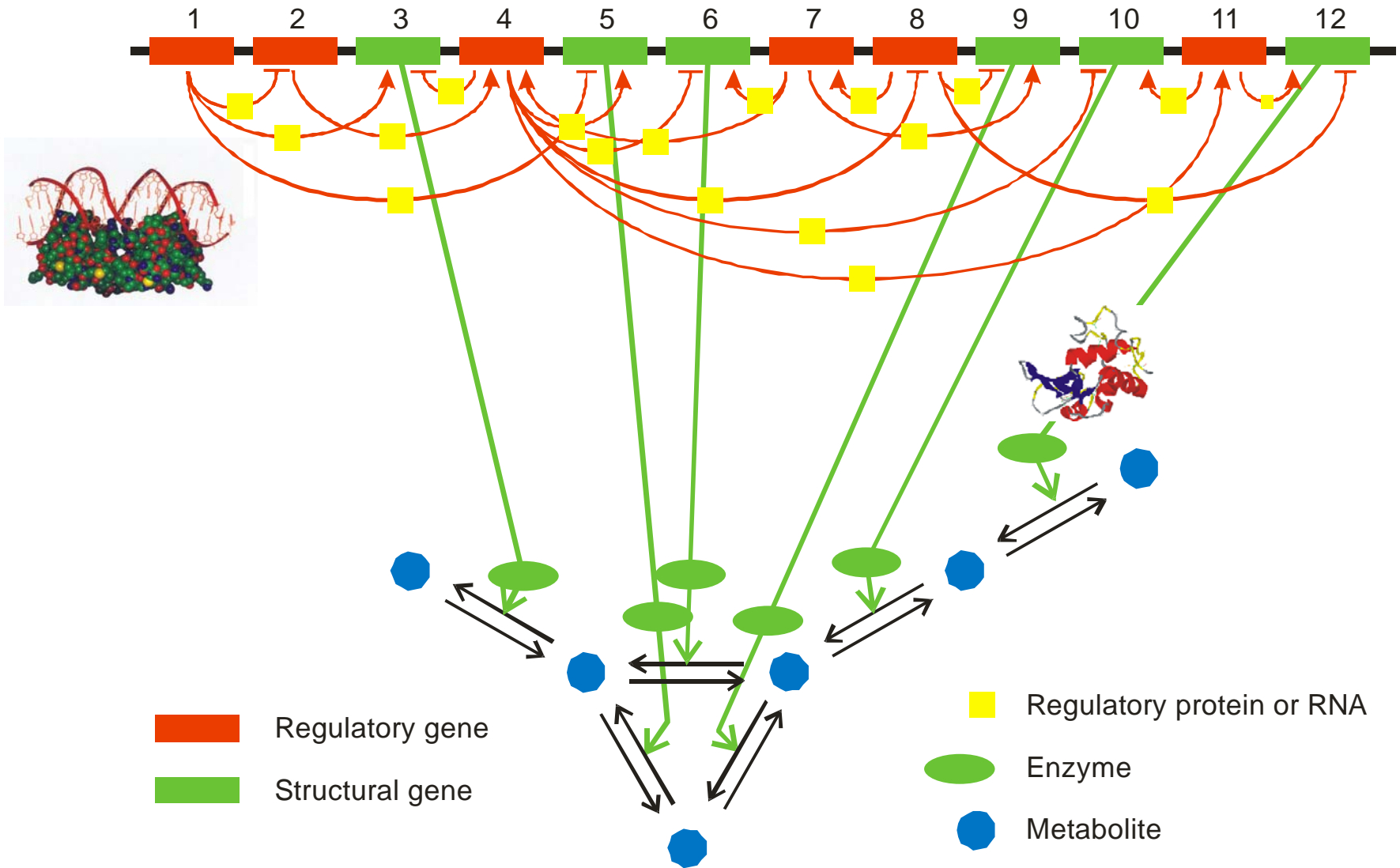
M. Mandal, B. Boese, J.E. Barrick, W.C. Winkler, R.R. Breaker. Cell 113:577-586 (2003)

1. Darwin's natural selection
2. The tree of life
3. From evolution *in vitro* to biotechnology
4. Genotypes with multiple functions
5. **How complex is biology?**

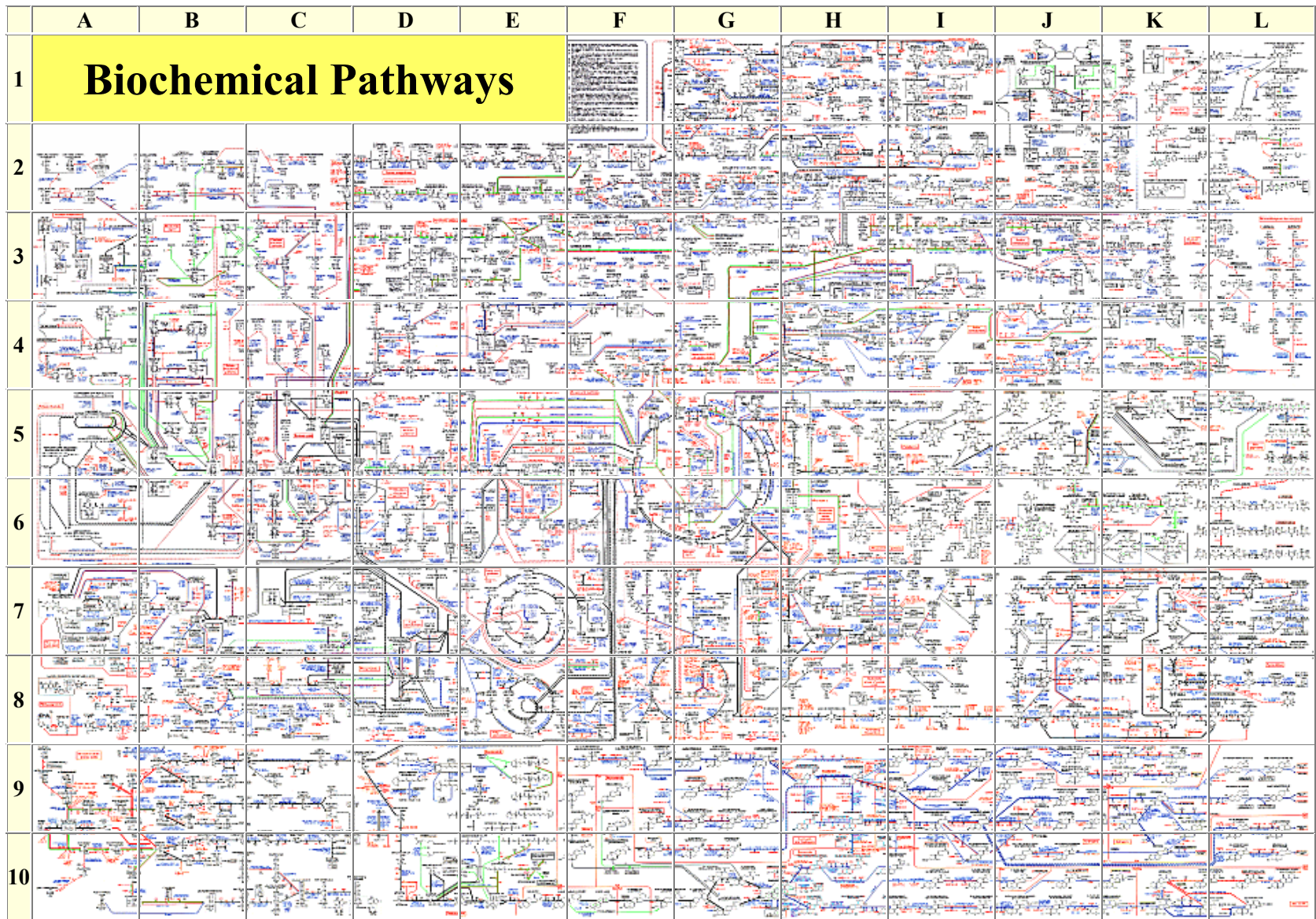
Three-dimensional structure of the complex between the regulatory protein **cro-repressor** and the binding site on λ -phage **B-DNA**



A model genome with 12 genes

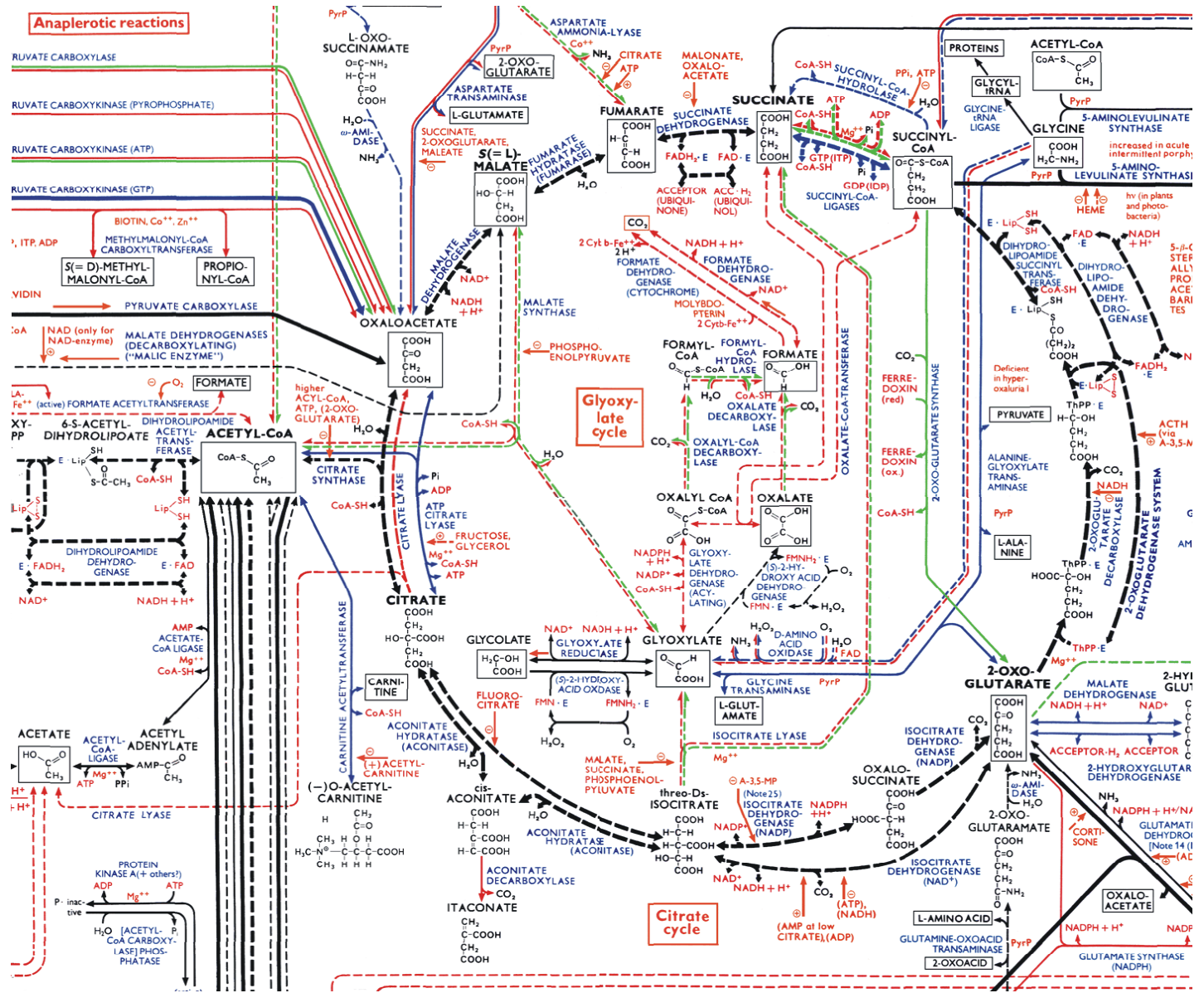


Sketch of a genetic and metabolic network



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

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

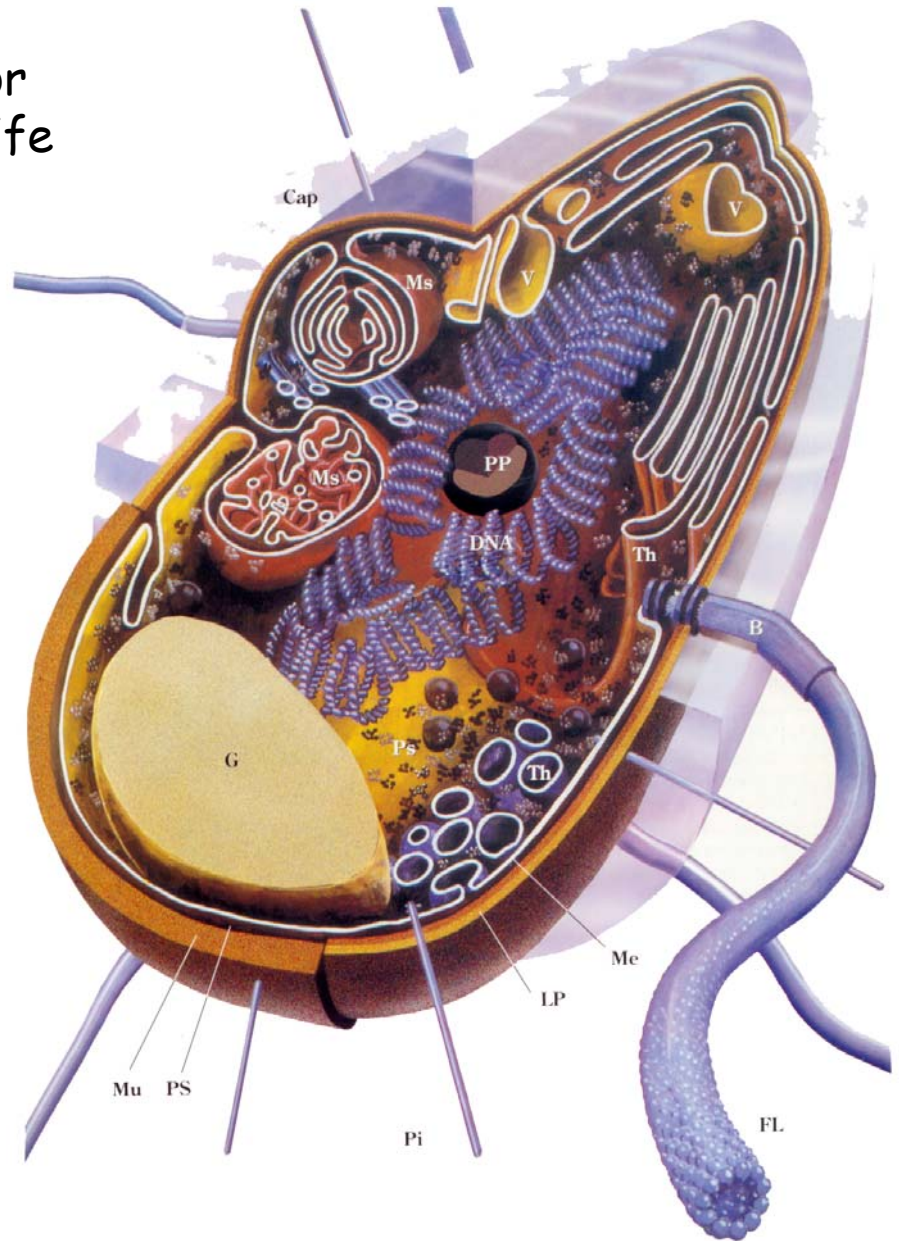


The bacterial cell as an example for the simplest form of autonomous life

Escherichia coli genome:

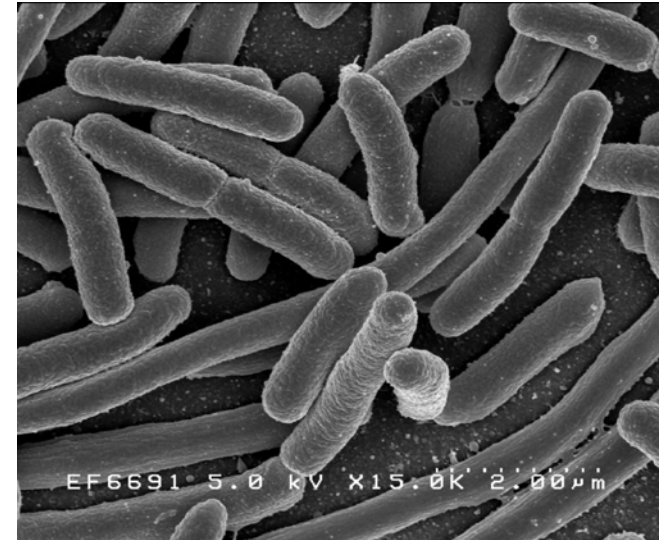
4 million nucleotides

4460 genes



The structure of the bacterium *Escherichia coli*

E. coli: Genome length 4×10^6 nucleotides
Number of cell types 1
Number of genes 4 460

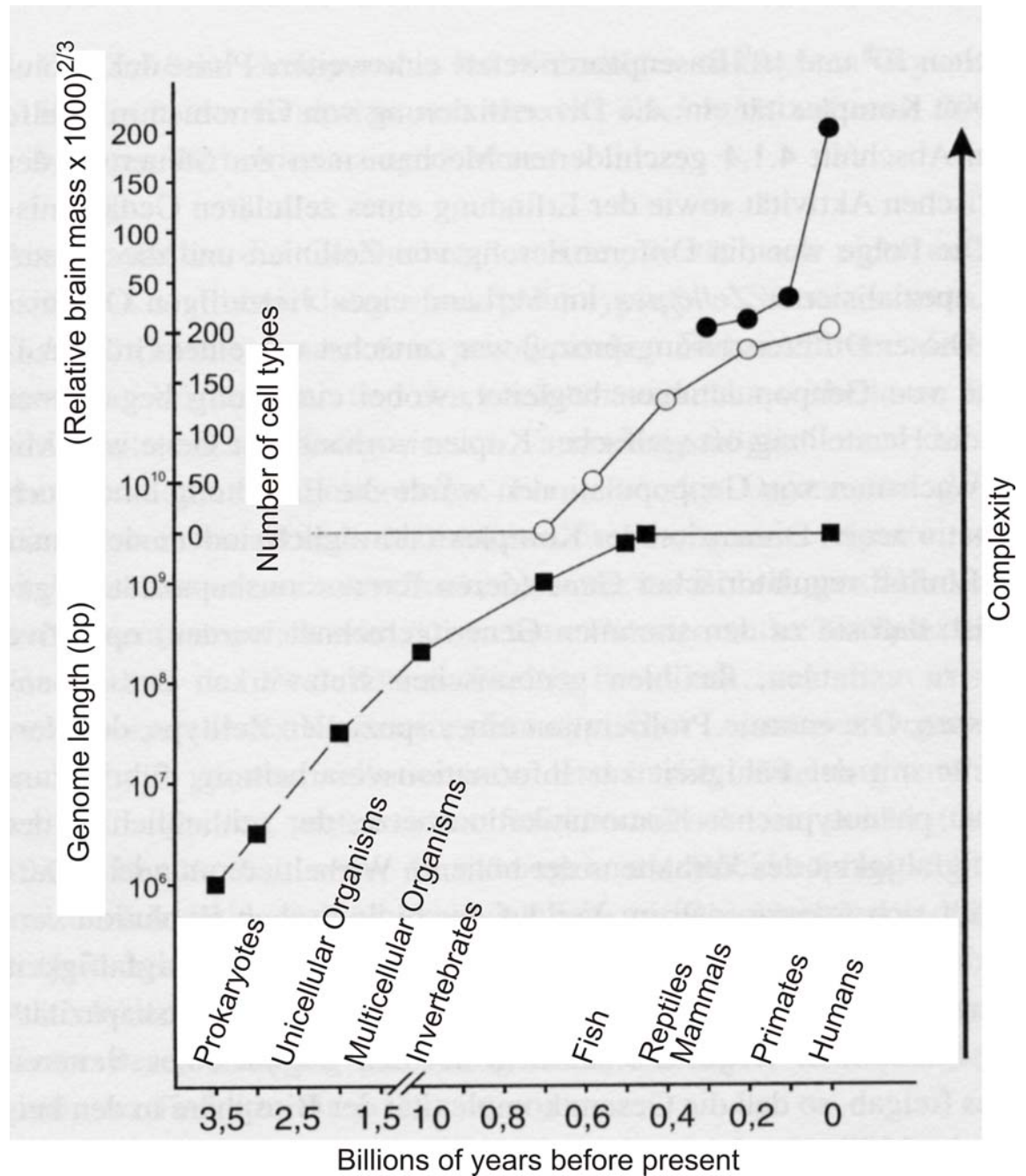


Man: Genome length 3×10^9 nucleotides
Number of cell types 200
Number of genes $\approx 30\,000$



Complexity in biology

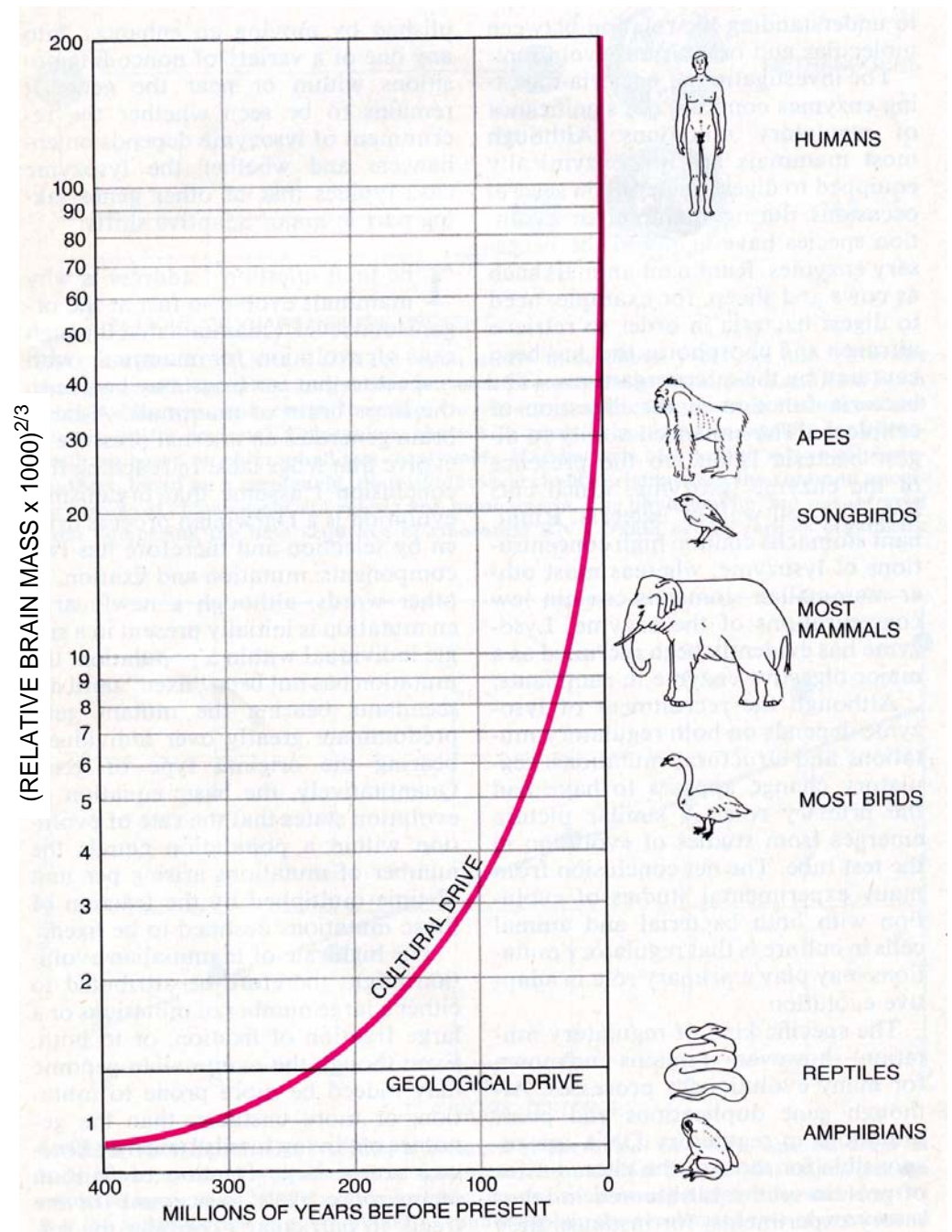
Wolfgang Wieser. 1998. *„Die Erfindung der Individualität“* oder *„Die zwei Gesichter der Evolution“*. Spektrum Akademischer Verlag, Heidelberg 1998

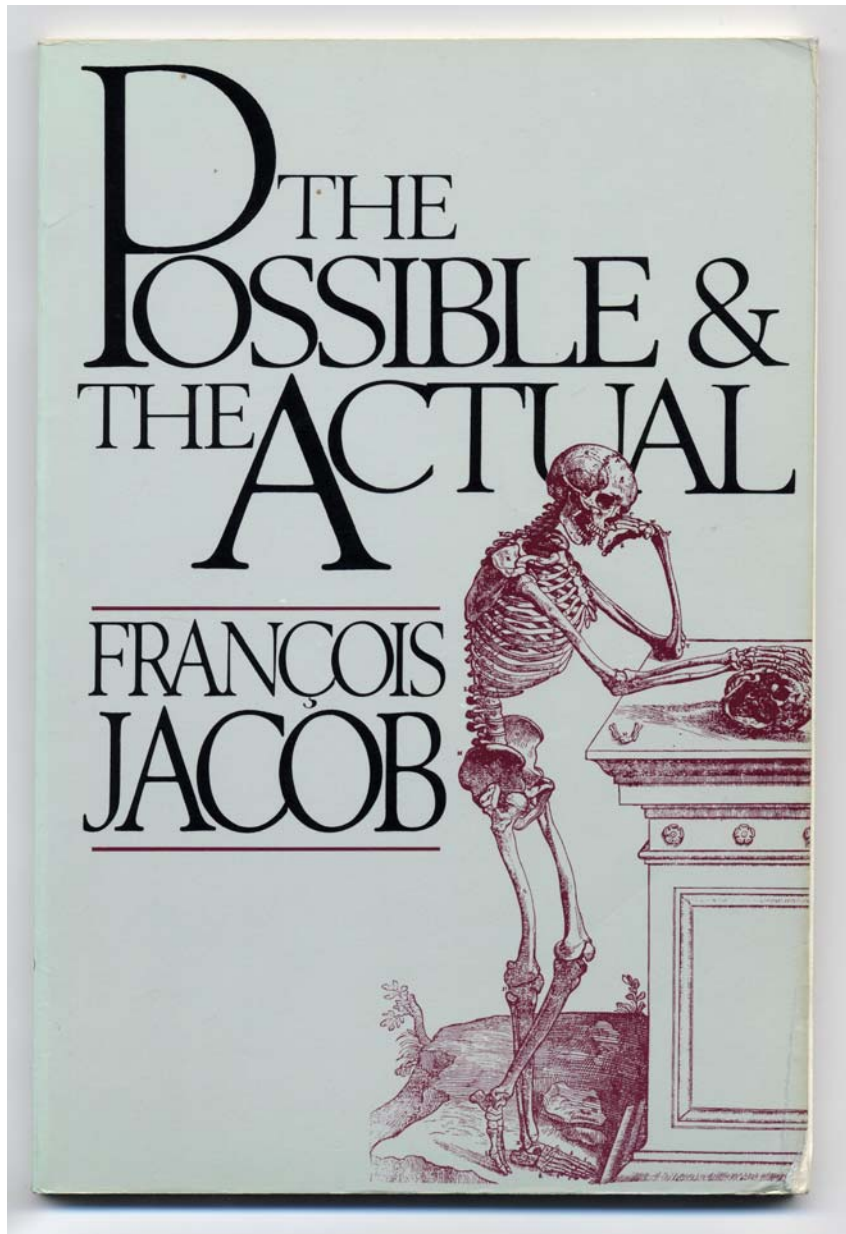




BRITISH TIT

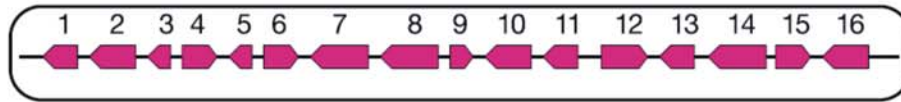
Alan C. Wilson.1985. The molecular basis of evolution.
Scientific American **253**(4):148-157.





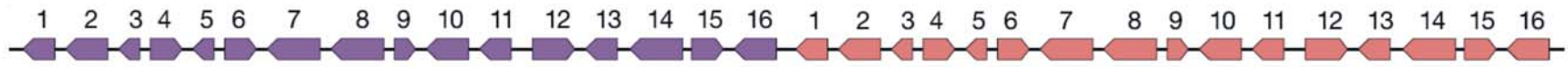
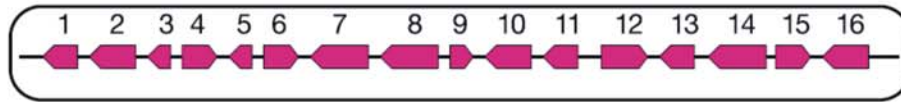
Evolution does not design with
the eyes of an engineer,
evolution works like a tinkerer.

François Jacob. *The Possible and the Actual*.
Pantheon Books, New York, 1982, and
Evolutionary tinkering. *Science* **196** (1977),
1161-1166.



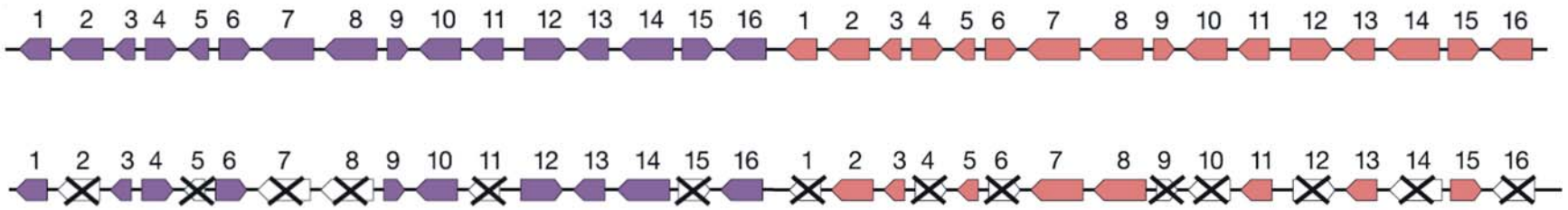
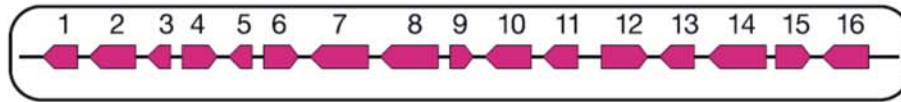
A model for the genome duplication in yeast 100 million years ago

Manolis Kellis, Bruce W. Birren, and Eric S. Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624, 2004



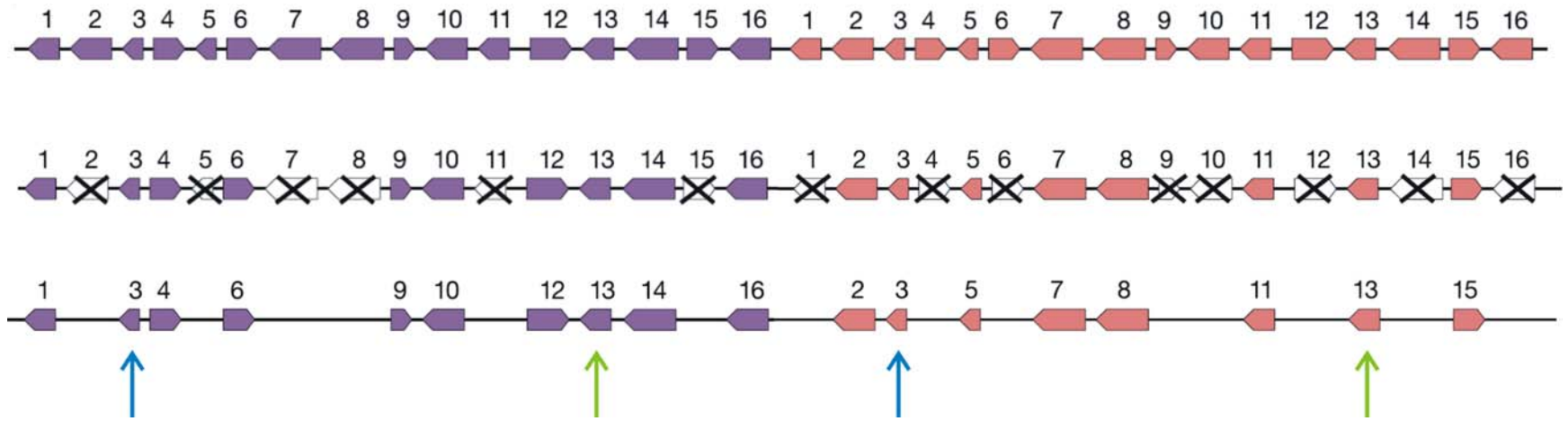
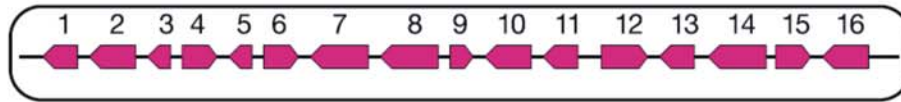
A model for the genome duplication in yeast 100 million years ago

Manolis Kellis, Bruce W. Birren, and Eric S. Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624, 2004



A model for the genome duplication in yeast 100 million years ago

Manolis Kellis, Bruce W. Birren, and Eric S. Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624, 2004



A model for the genome duplication in yeast 100 million years ago

Manolis Kellis, Bruce W. Birren, and Eric S. Lander. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* **428**: 617-624, 2004

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 past¹. A study on page 469 of this issue suggests that a comparable phenomenon might occur in mice, and by implication in other mammals². 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."

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 from³.

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 team⁴, and one by geneticist Rotem Sorek⁵, 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

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

— Phillip Kapranov

The difficulty to define the notion of „gene“.

Helen Pearson,
Nature 441: 399-401, 2006

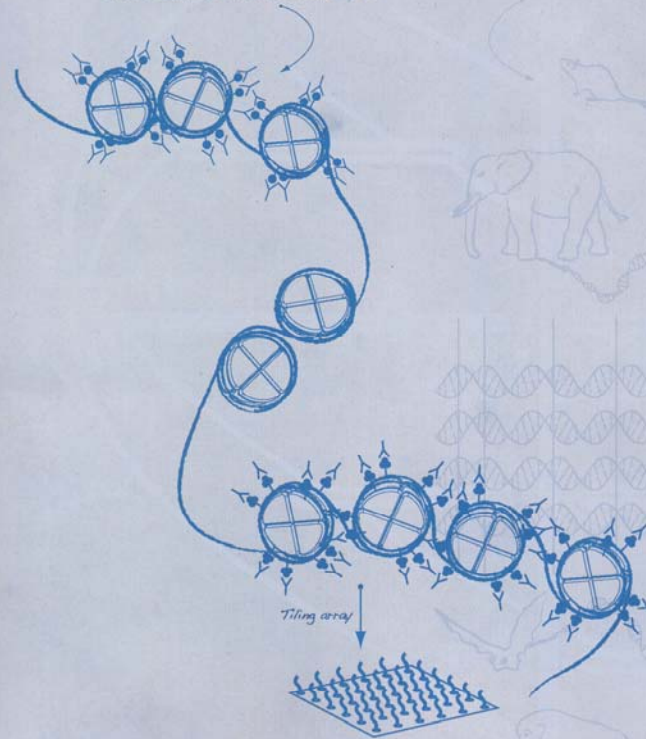


Spools of DNA (above) still harbour surprises, with one protein-coding gene often overlapping the next.

nature

Hi-stone-modification chromatin IP

Comparative syntenic alignment



**MARS'S
ANCIENT OCEAN**
Polar wander
solves an enigma

**THE DEPTHS OF
DISGUST**
Understanding the
ugliest emotion

MENTORING
How to be top

NATUREJOBS
Contract
research

DECODING THE BLUEPRINT

The ENCODE pilot maps
human genome function



ENCODE stands for
ENCyclopedia **Of** **DNA** **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

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

Martin Nowak, Harvard University, MA

Christian Reidys, Nankai University, Tien Tsin, China

Christian Forst, Los Alamos National Laboratory, NM

Thomas Wiehe, 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, Jan Cupal, Stefan Bernhart, Lukas Ender, Ulrike Langhammer,
Rainer Machne, Ulrike Mückstein, Erich Bornberg-Bauer,
Universität Wien, AT



Universität Wien

Acknowledgement of support

Fonds zur Förderung der wissenschaftlichen Forschung (FWF)
Projects No. 09942, 10578, 11065, 13093
13887, and 14898

Wiener Wissenschafts-, Forschungs- und Technologiefonds (WWTF)
Project No. Mat05

Jubiläumsfonds der Österreichischen Nationalbank
Project No. Nat-7813

European Commission: Contracts No. 98-0189, 12835 (NEST)

Austrian Genome Research Program – GEN-AU: Bioinformatics
Network (BIN)

Österreichische Akademie der Wissenschaften

Siemens AG, Austria

Universität Wien and the Santa Fe Institute



Universität Wien

Thank you for your attention!

Web-Page for further information:

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

