

# **Human viewed as a machine**

**or Biology from the point of view of a  
Mathematician**

**Alvis Brazma**

Bioinformatics Summer School, BDS 3, Uzhhorod, July 2023



# The world before 1991



Riga





With two of my fellow undergraduate students in the Carpathian mountains not too far from here in ~1978

# ChatGPT, continue:



continue 16, 32, 64, 128, 256, ...



512, 1024, 2048, 4096, 8192, 16384, 32768, 65536, 131072, 262144, ...



continue 16, 1, 6, 32, 3, 2, 64, 6, 4, 128, 1, 2, 8, 256, 2, 5, 6, ...



This sequence doesn't seem to have an obvious pattern, but here are the next few terms:

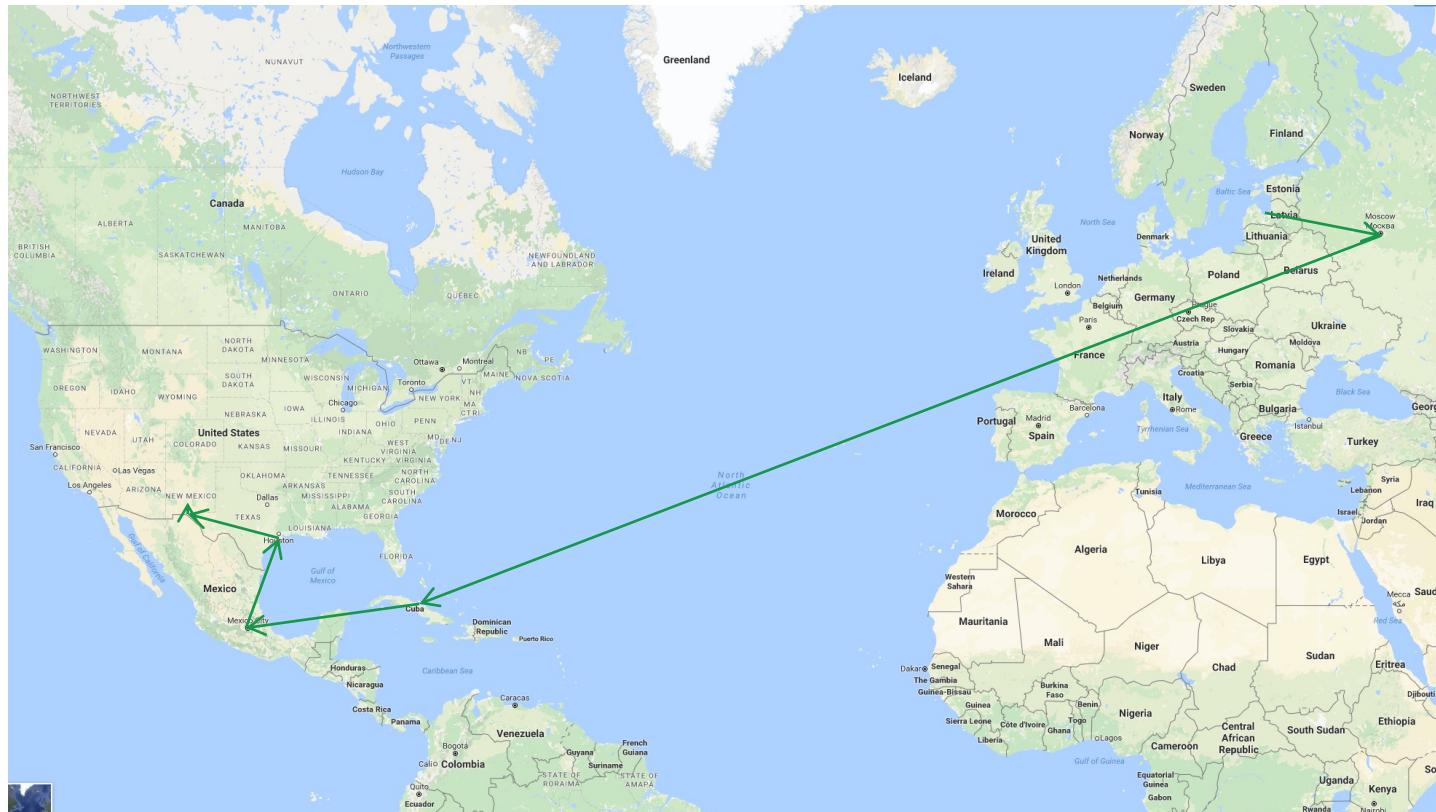
512, 4, 8, 1024, 9, 10, 2048, 12, 16, 4096, 1, 18, 32, 8192, 20, 25, 18, 16384, 27, 4, 65536, 30, 32, 512, ...

Without more information, it's difficult to determine a clear rule or formula that generates this sequence. It's possible that the sequence is arbitrary or follows a complex pattern that is not immediately apparent.

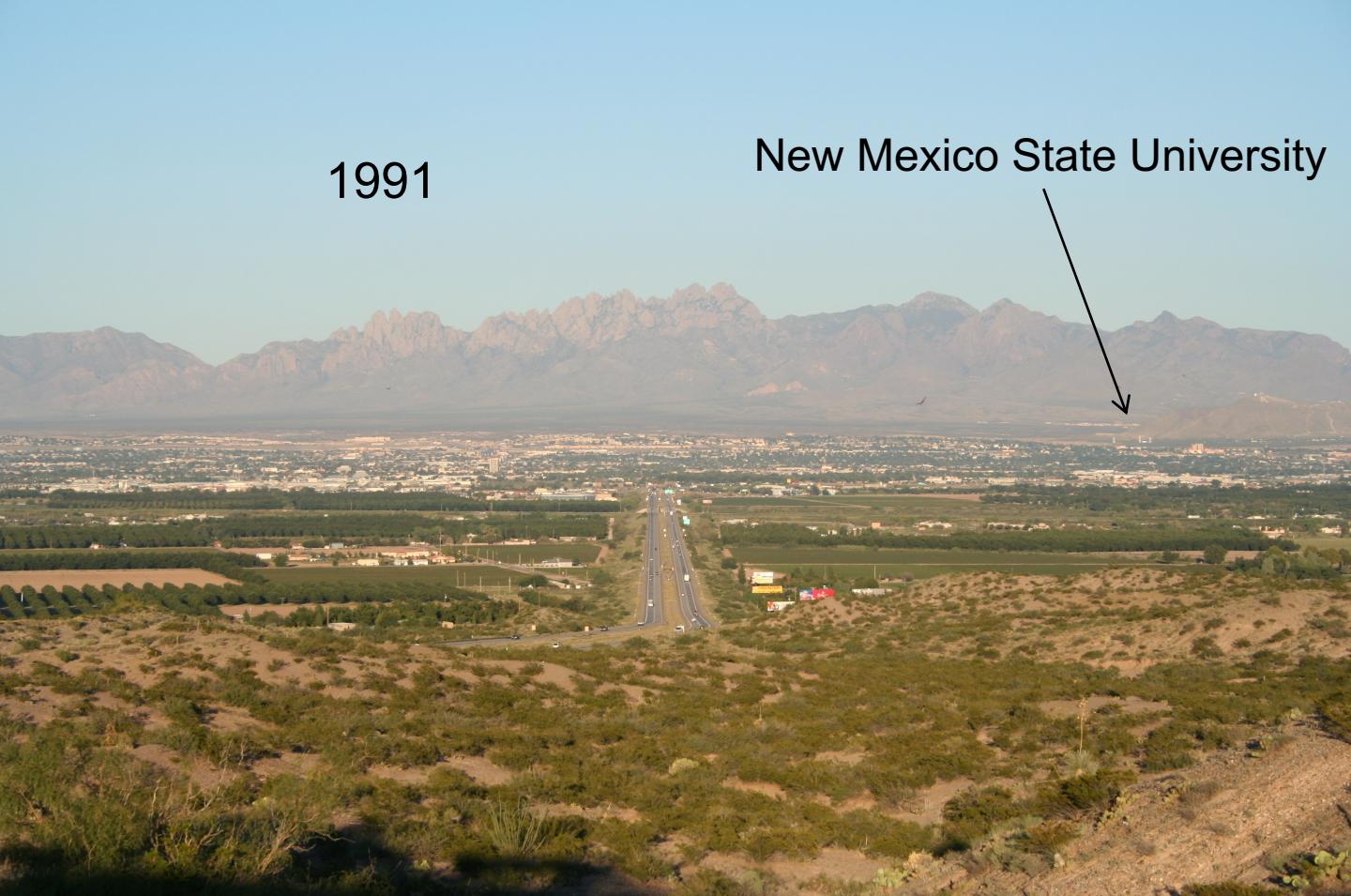
**My PhD was roughly about this**



9 November 1989



My itinerary on 14-16 January 1990



1991

New Mexico State University

# The case for the human genome

A meeting next month in London should give a clear view of where the Human Genome Project is now, and where it is likely to be heading in the years that follow. The project (which is several) continues to gather momentum. The strategic questions have been settled, but there will be endless arguments on tactics. And ethical issues, while insubstantial, must be handled delicately.

FAMILIARITY exorcizes fear. That seems to be one of the unremarkable lessons learned so far from the great Human Genome Project or, more accurately, from the dozens of national projects that have sprung up in the past few years in different places.

Fears that research in biology would in future be conducted on automated assembly lines, or that research laboratories would be robbed by the demands of the sequencing plants of a whole generation of skilled technicians, have not materialized and are unlikely to do so, perhaps because none of the projects has yet entered what might be called its production phase.

Yet on the first issue, it is curious that Professor Walter Gilbert's impassioned statement of the reasons why biologists must learn

plete sequence began.

Fears that human genome projects everywhere would rob other worthwhile projects of funds have similarly not yet bitten into the pattern of research supported by the grant-making agencies. Moreover, when people are talking realistically of automatic sequencing at a cost of no more than \$0.50 a base (with overheads included), it seems unlikely that a ten-year project to provide the complete sequence of 23 stretches of DNA, with 3,000 million bases altogether, should much exceed the annual cost of the present development phase.

The ethical objections to the sequencing of the human genome are necessarily more subtle. They range from the assertion that it would be improper that knowledge won in a

than any now available to the places in the human genome at which genetic abnormalities occasion genetic disease or a high frequency of somatic genetic change (as in cancer) — that is one of the important objectives of the exercise, after all. But the objection, if valid, also applies to what is already possible by the use of genetic markers for identifying particular variants of a genetic polymorphism.

In present practice in, for example, the insurance of persons, the trouble (and cost) that companies take to discriminate between their would-be customers is measured by the risks they are asked to undertake; it remains to be seen whether present ambitions to automate the sequencing process will be followed by a machine that can print out a

# European Bioinformatics Institute

- World leading source of public genomics and biomolecular data
- We accept, process, archive and make available to the life sciences community genomics and other life sciences data
- We are part of the European Molecular Biology Laboratory (EMBL), Europe's flagship laboratory for the life sciences.



# EMBL member states

## Member states (29)

Austria 1974	Belgium 1990
Denmark 1974	Portugal 1998
France 1974	Ireland 2003
Germany 1974	Iceland 2005
Israel 1974	Croatia 2006
Italy 1974	Luxembourg 2007
Netherlands 1974	Czech Republic 2014
Sweden 1974	Malta 2016
Switzerland 1974	Hungary 2017
United Kingdom 1974	Slovakia 2018
Finland 1984	Montenegro 2018
Greece 1984	Poland 2019
Norway 1985	Lithuania 2019
Spain 1986	Estonia 2023
	Latvia 2023

## Associate member states

Australia 2008

## Prospect member states

Serbia



# European Molecular Biology Laboratory - EMBL

Europe's centre  
of excellence in *life*  
*sciences* research,  
services and training

Founded in 1974 by 10 states  
as an intergovernmental  
organisation to promote the  
molecular life sciences in Europe  
and beyond

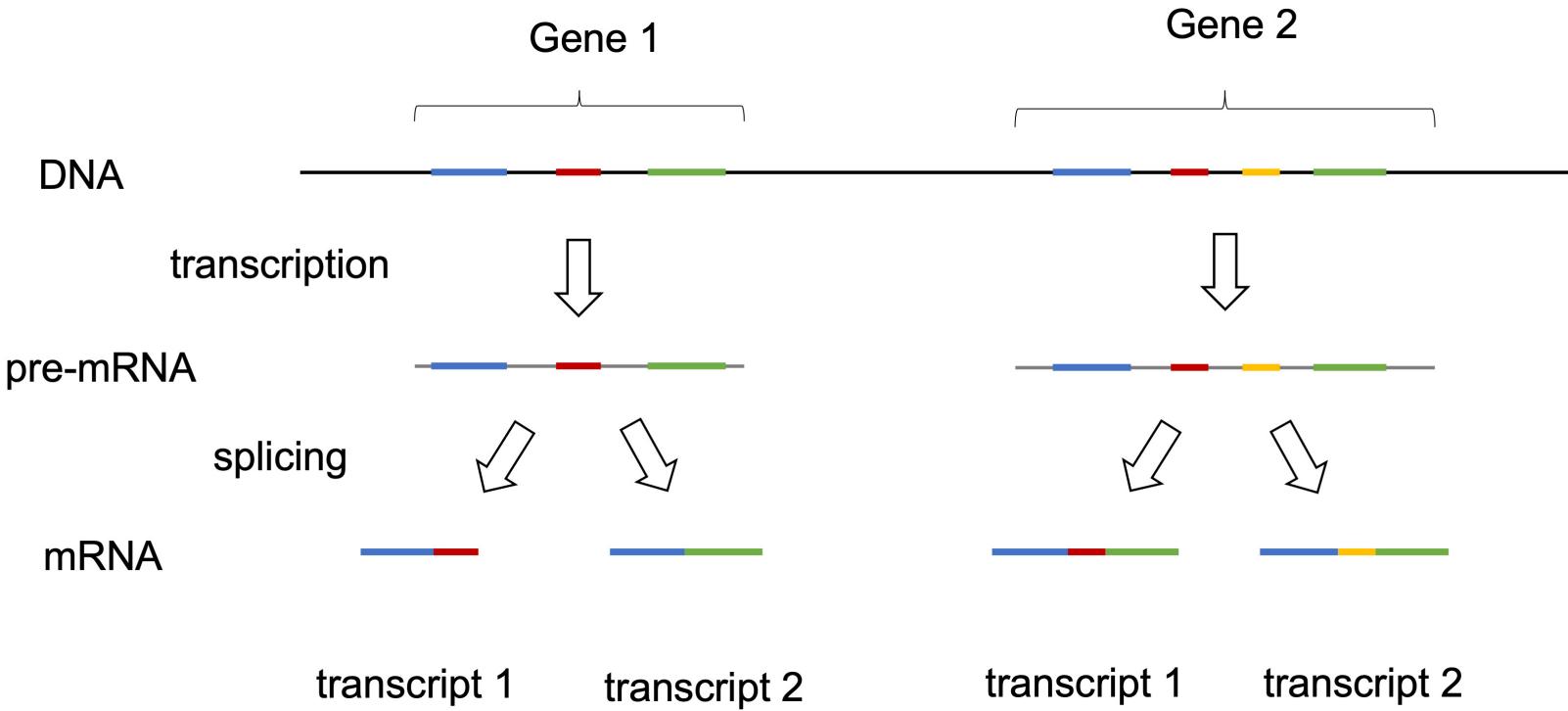
"I believe that international activity is very important in building world peace." Sir John Kendrew, EMBL's 1<sup>st</sup> DG



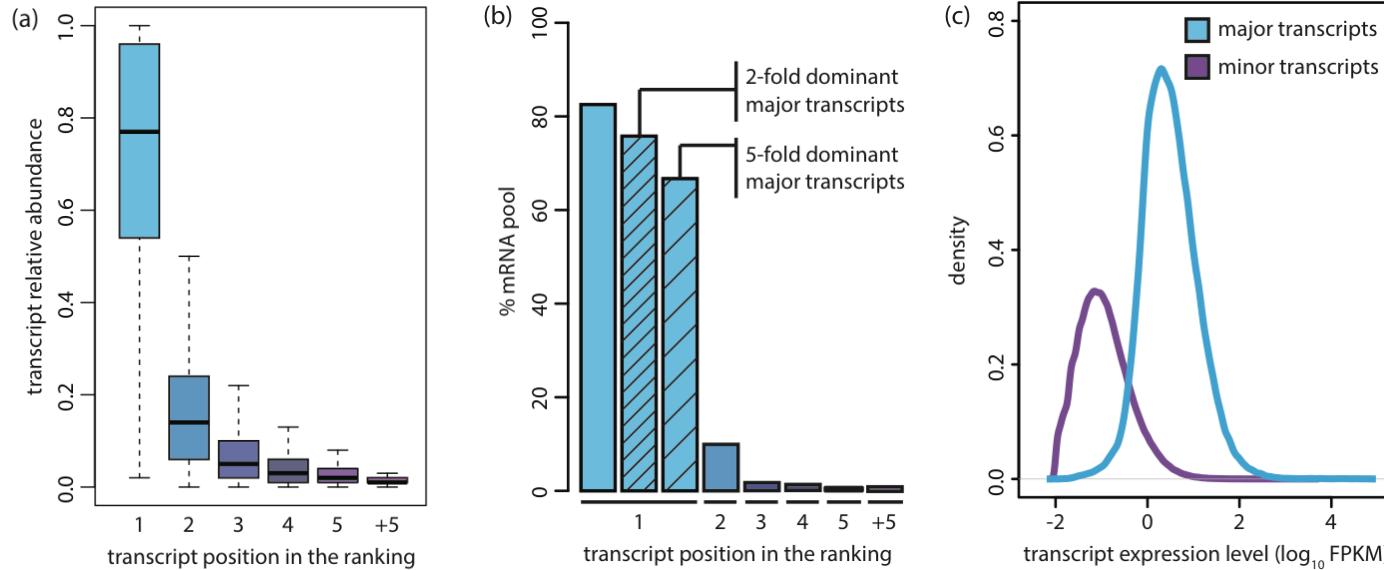
First Director General  
Sir John Kendrew  
Nobel Prize in 1962

# What's the difference between theoretical and experimental biologists?

- Experimental biologists observe what cannot be explained
- Theoretical biologists explain what cannot be observed
- Two perspectives on science
  - Mathematicians – looking for rules from which everything else follows as much as possible; exceptions can be taken care of later
  - Biologists – outliers (exceptions) are what is interesting, the rules, if there are any, usually are obvious and not interesting

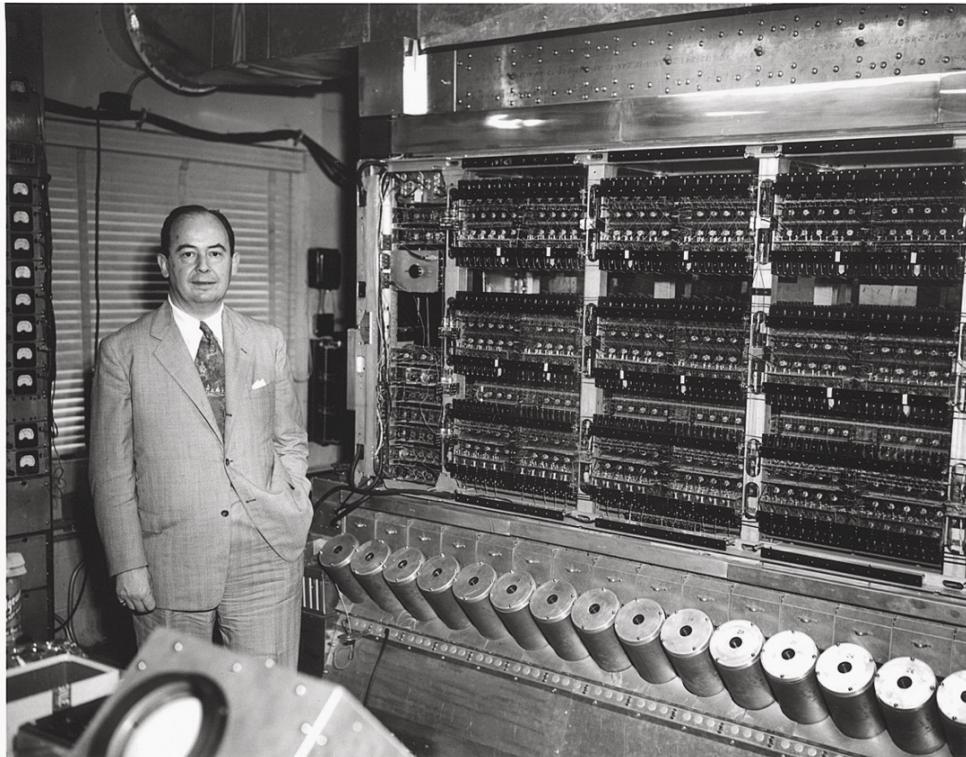


# An example of a rule (with many important exceptions) – most genes in most tissues express only one splice-form



Mar Gonzàlez-  
Porta

# Can a machine reproduce? What are the minimum needed for this?



John von Neumann (1903 – 1957)

Ability to reproduce is often seen as one of the dividing lines between living and nonliving

*Her Majesty pointed to a clock and said, 'See to it that it produces offspring'.*

From alleged conversation between Queen Christina of Sweden and René Descartes in 1649

# Von Neuman's mechanical replicator

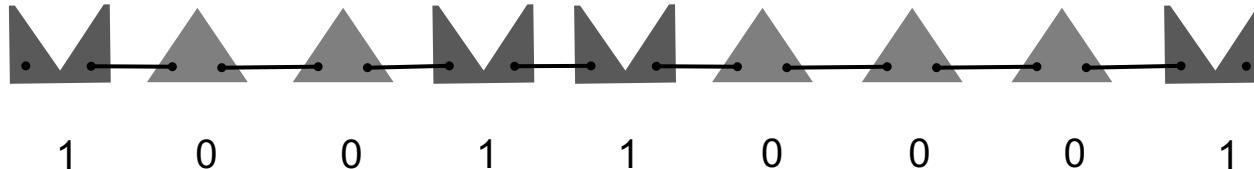
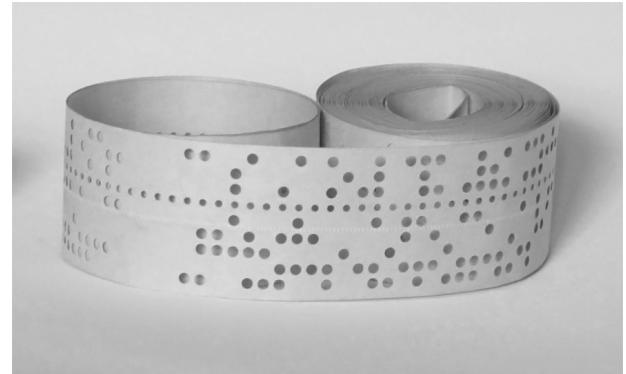
- *Draw up a list of unambiguously defined elementary parts. Imagine that there is a practically unlimited supply of these parts floating around in a large container. One can imagine an automaton functioning in the following manner: It also is floating around in this medium; its essential activity is to pick up parts and put them together (...).* Von Neumann and Burks (1966) from von Neumann's lectures in 1948
- The task of this automaton is to assemble a clone of itself

# Degenerative trend

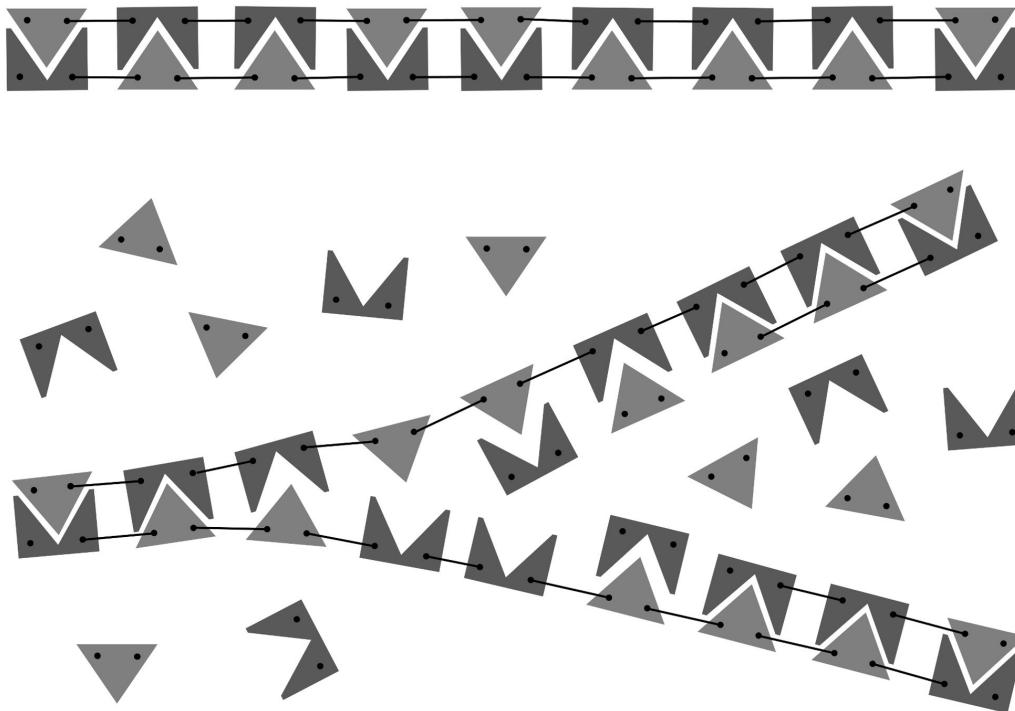
- It does appear that for an *automaton* to make a *thing*, the *automaton* needs to contain the full description of that *thing*
- If this description is a part of the *automaton*, then the automaton becomes more complex than the *thing* it can build
- Doesn't it follow that an *automaton* can only make *things* that are simpler than itself? If so, this leads to a *degenerative trend*. Isn't this a paradox?
- The solution to this paradox is in the decoupling the description (of the *thing* to build) from the *automaton* building it

# The Universal Constructor

- *Universal Constructor*—an abstraction of a **general-purpose robot**, which given appropriate instructions and necessary materials could build whatever it has been instructed to build.
- Instructions to the universal constructor can be encoded on a **tape**
- Replicating 1- or 2-dimensional objects are much simpler than replicating a complex 3-dimensional object



# Exploiting complementarity



# How to clone oneself

- $X$  – an arbitrary *thing* made of von Neumann's elementary parts; this *thing* can also be some sort of *aparatus*
- $U$  – Universal Constructor
- $\text{Code}(X)$  – a tape containing the description of  $X$  that given to  $U$  produces  $X$
- $R$  – a simple device that can replicate the tape  $\text{Code}(X)$
- $C$  – a *controlling device* that given  $\text{Code}(X)$  first runs  $R$  on  $\text{Code}(X)$  producing another copy of  $\text{Code}(X)$ , then runs  $U$  on  $\text{Code}(X)$  producing  $X$  and then ‘ties’ them together

# How to clone oneself

- The device  $U+R+C+\text{code}(X)$  produces  $X+\text{code}(X)$
- Thus, starting with  $U+R+C+\text{code}(X)$  after some time we get

$U+R+C+\text{code}(X), X+\text{code}(X)$

- But  $X$  is an arbitrary object, it can be some sort of *apparatus*. Take

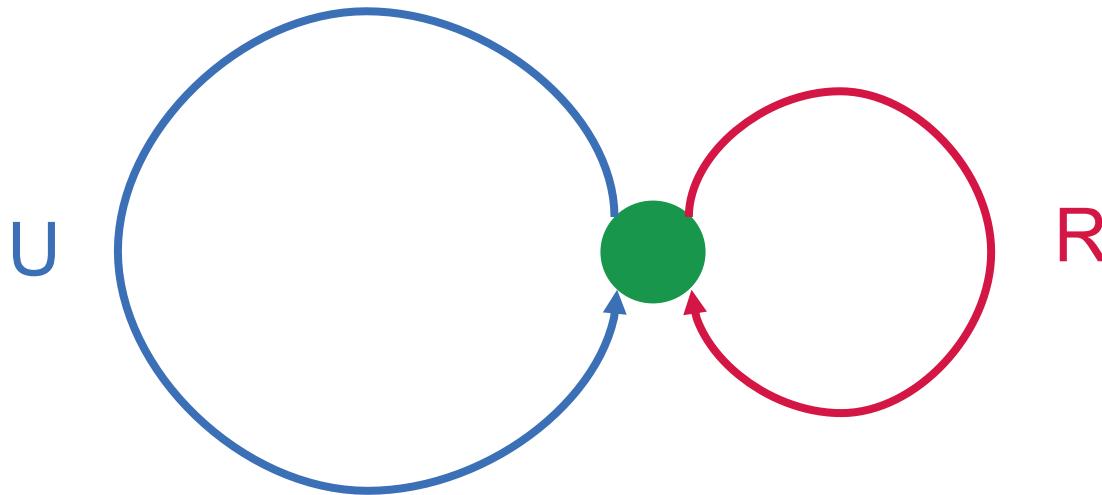
$X = U+R+C+\text{code}(U+R+C)$

- Then starting from  $U+R+C+\text{code}(U+R+C)$  we get

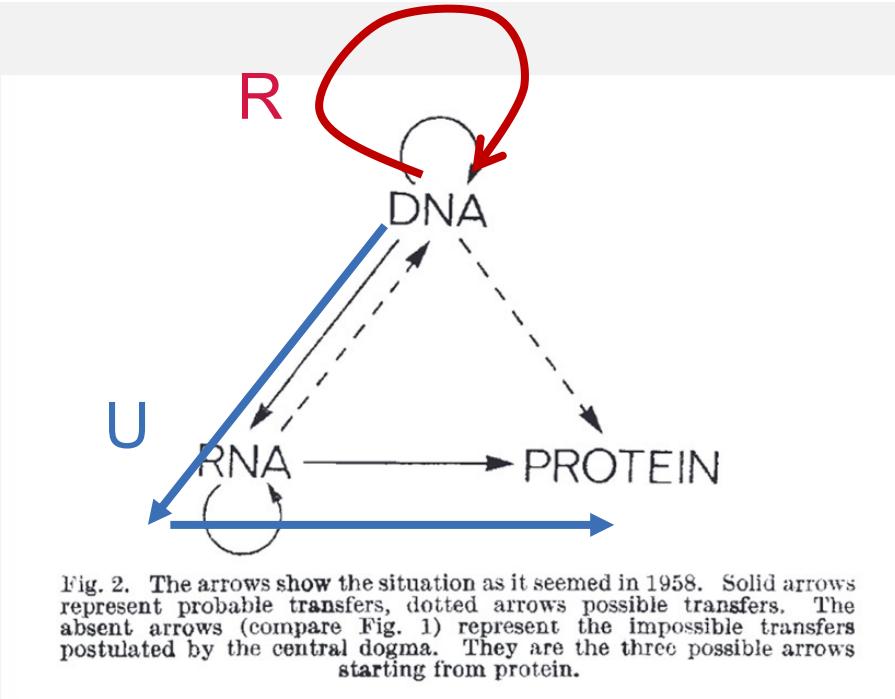
$U+R+C+\text{code}(U+R+C), U+R+C+\text{code}(U+R+C)$

**The two processes U and R have to be separated**

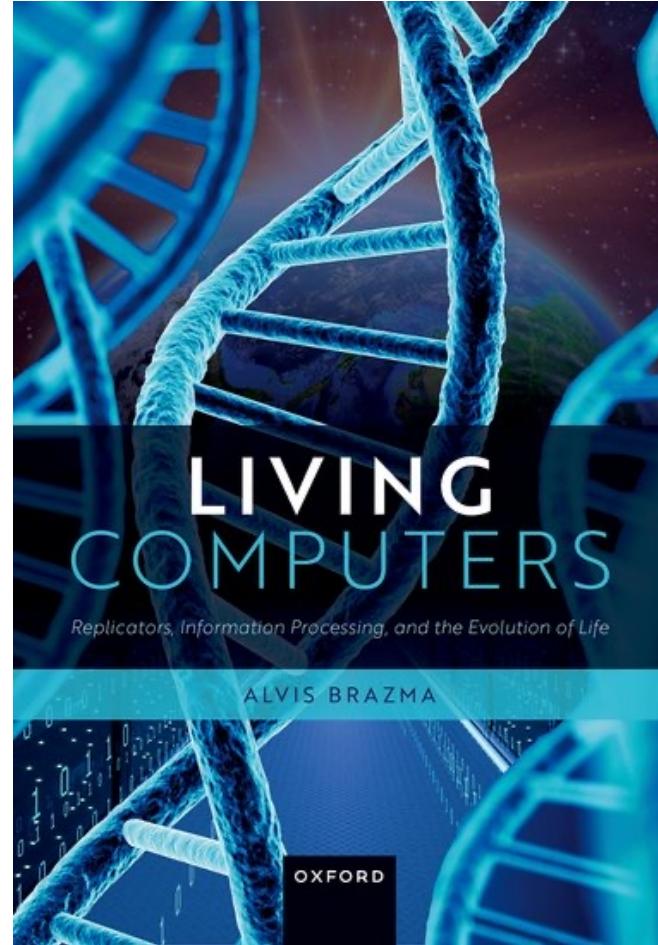
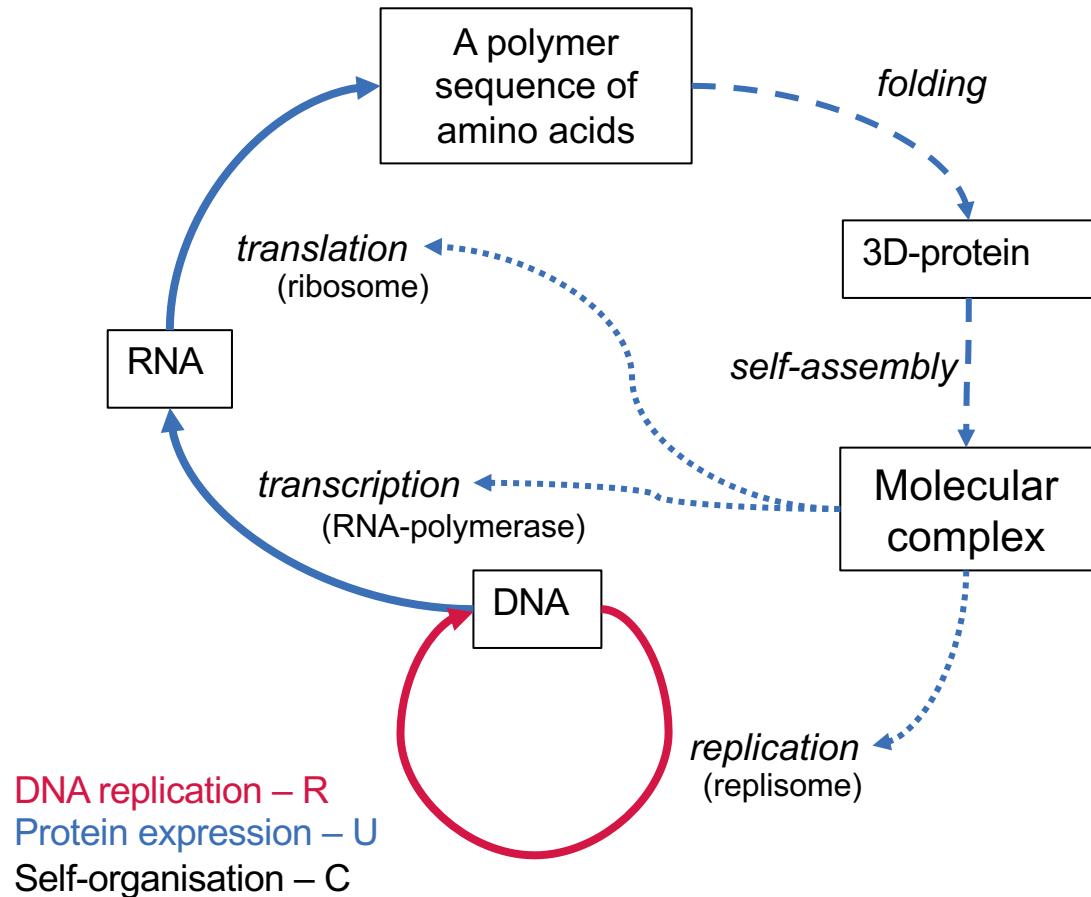
To beat the *degenerative trend* processes **U** and **R** have to be separated



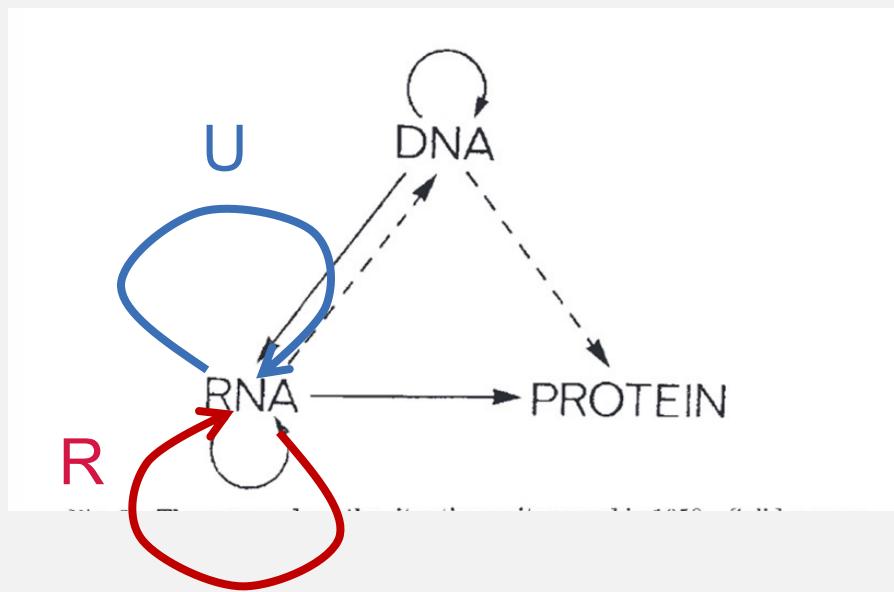
# Crick's central Dogma



# Crick's central dogma and von Neuman's replicator



# RNA life



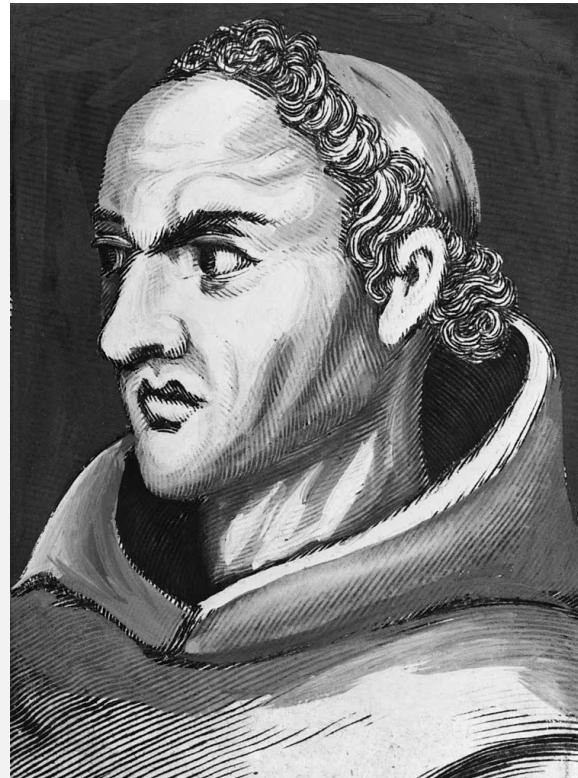
Is this how life started?

# Minimum complexity threshold

- There is minimum complexity necessary for information carrying self-replication.
- Below this complexity we have only degenerative trend possible
- Above this the system can build a machine as complex as itself and (through evolution) even a more complex system
- Given that Universal Constructor is quite complex, this complexity threshold is quite high
- Life cannot be too simple!
- Biologists will say "we knew it all along". For a mathematician this is a revelation ;-)

# Occam's Razor

- “Entities should not be multiplied without a necessity”
- The modern formulation – *if several theories explain the observations equally well, choose the simplest one*
- In my own formulation – *explain as much as possible with as little as possible*



William of Ockham (1287 – 1347)

# **Rissanen's minimum description length (MDL) principle:**

the best model of data is their shortest description (in bits)

This is a practical implementation of Occam's Razor

# Applying Rissanen's MDL principle

```
ATTGATGAGAGTTA  
TAGCAGGATGAGTAGCAG  
TTAGCAGGATGAGCTGCAT  
GAGCTGATGAGTTACA  
CAGTCTGATGAGCGTATA  
TATTGTGAAATATTTATTG  
CCATTGAAACGCAGCATGAT  
CTAGTTGTGAAAGCGAGCTGATG  
ATCTCGTGAAGGTATTCAG  
TGATGAAACGAAATGAAA
```

# Applying Rissanen's MDL principle

ATT**GATGAGAGTTA**  
TAGCAG**GATGAGTAGCAG**  
TTAGCAG**GATGAGCTGCAT**  
GAGCT**GATGAGTTACA**  
CAGTCT**GATGAGGCGTATA**  
TATTG**TGAAA**TATTTATTG  
CCAT**TGAAA**CGCAGCATGAT  
CTAGTTG**TGAAA**GCGAGCTGATG  
ATCTCG**TGAAA**GGTATTCAG  
**TGATGAAACGAAATGAAA**

# Applying Rissanen's MDL principle

ATT**GATGAG**AGTTA  
TAGCAG**GATGAG**TAGCAG  
TTAGCAG**GATGAG**CTGCAT  
GAGCT**GATGAG**TTACA  
CAGTCT**GATGAG**GCGTATA  
TATTG**TGAAA**TATTTATTG  
CCAT**TGAAA**CGCAGCATGAT  
CTAGTTG**TGAAA**GCGAGCTGATG  
ATCTCG**TGAAA**GGTATTCAG  
TGA**TGAAA**CGAAATGAAA

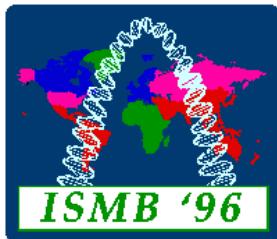
ATT**X1**AGTTA  
TAGCAG**X1**TAGCAG  
TTAGCAG**X1**CTGCAT  
GAGCT**X1**TTACA  
CAGTCT**X1**GCGTATA  
TATTG**X2**TATTTATTG  
CCAT**X2**CGCAGCATGAT  
CTAGTTG**X2**GCGAGCTGATG  
ATCTCG**X2**GGTATTCAG  
**TGAX2**CGAAATGAAA

**X1=GATGAG**

**X2=TGAAA**

189 letters

171 characters



From: ISMB-96 Proceedings. Copyright © 1996, AAAI (www.aaai.org). All rights reserved.

## Discovering Patterns and Subfamilies in Biosequences

**Alvis Brāzma\***

abrazma@cclu.lv

Institute of Mathematics and Computer Science  
University of Latvia  
29 Rainis Bulevard  
LV-1459 Riga, Latvia

**Esko Ukkonen**

Esko.Ukkonen@cs.Helsinki.FI  
Department of Computer Science  
P.O.Box 26 (Teollisuuskatu 23)  
FIN-00014 University of Helsinki  
Finland

### Abstract

We consider the problem of automatic discovery of patterns and the corresponding subfamilies in a set of biosequences. The sequences are unaligned and may contain noise of unknown level. The patterns are of the type used in PROSITE database. In our approach we discover patterns and the respective subfamilies simultaneously. We develop a theoretically substantiated significance measure for a set of such patterns and an algorithm approximating the best pattern set and the subfamilies. The approach is based on the minimum

**Inge Jonassen**

inge@ii.uib.no

Department of Informatics  
University of Bergen, HIB  
N5020 Bergen, Norway

**Jaak Vilo**

Jaak.Vilo@cs.Helsinki.FI  
Department of Computer Science  
P.O.Box 26 (Teollisuuskatu 23)  
FIN-00014 University of Helsinki  
Finland

2. as grouping or clustering a set of biosequences into subsets so that each subset shares a distinct common pattern and the noise is sorted out.

A subproblem of this is finding a pattern common to a set of sequences in the presence of an unknown level of noise.

In particular we consider sequences representing proteins and patterns of the type used in PROSITE database (Bairoch 1992), but the method can be applied

$$F'(\pi, l) = l \cdot (I(\pi) + 2 \log p + f(\pi)) - M_1(\pi),$$

where

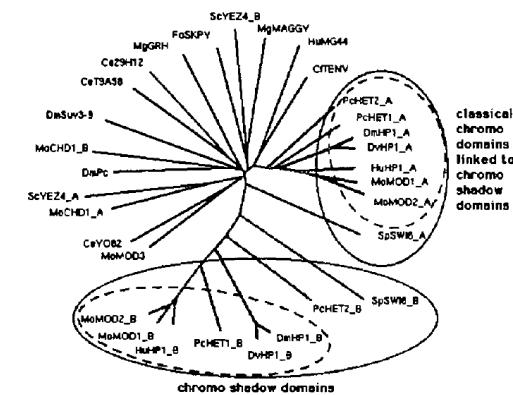
$$I(\pi) = I_1(\pi) + I_2(\pi),$$

$$I_1(\pi) = \sum_{j \in \hat{\Phi}(\pi)} I'(\pi(j)),$$

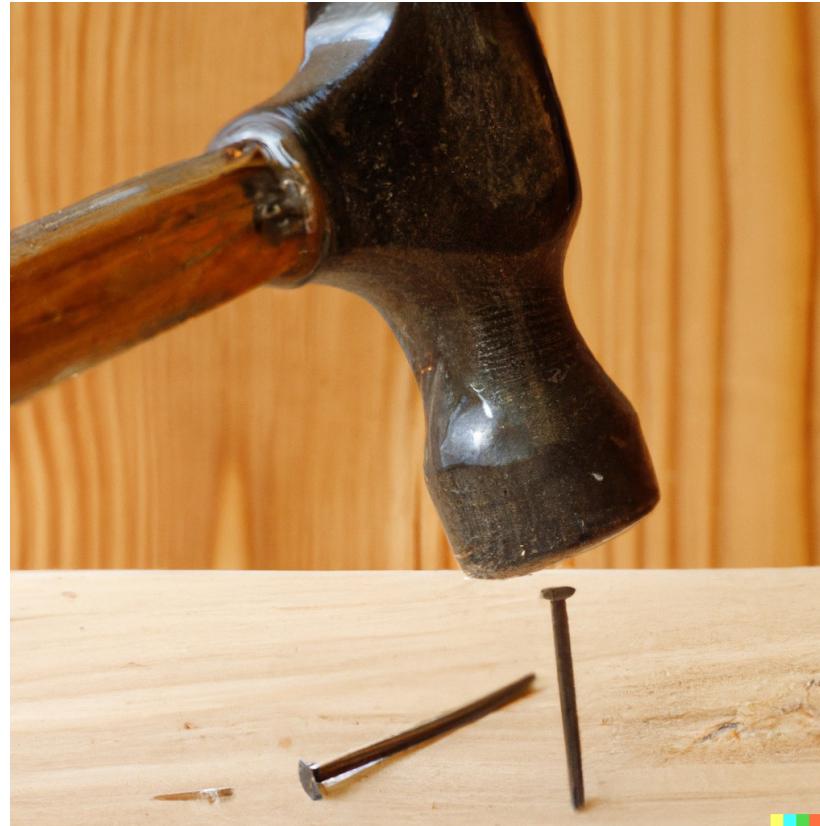
$$I'(b_h) = - \sum_{i=1}^m P_2(a_i) \log P_2(a_i) + \\ + \sum_{a_i \in K_h} P_2(a_i|b_h) \log P_2(a_i|b_h),$$

$$I_2(\pi) = - \sum_{x(t_i, v_i)} \log(v_i - t_i + 1)$$

NJ-tree (ClustalW) of chromo domains



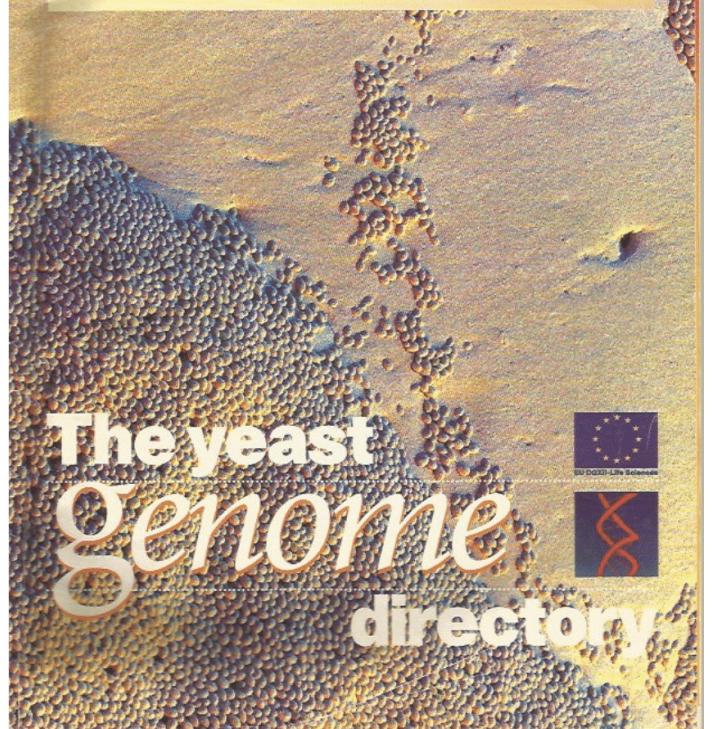
Converting from a  
mathematician to  
a biologist



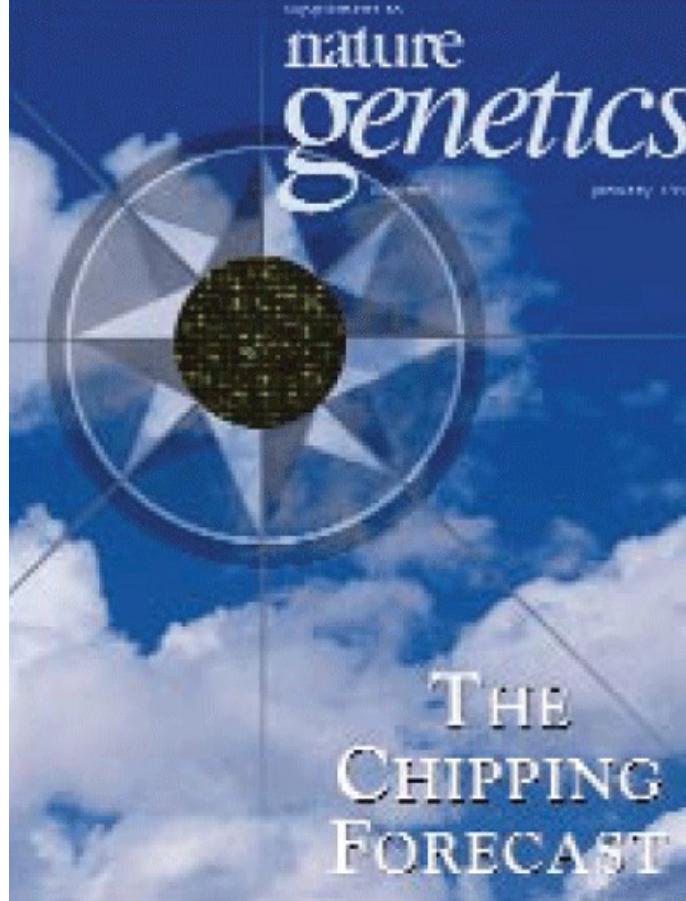
29 May 1997

International weekly journal of science

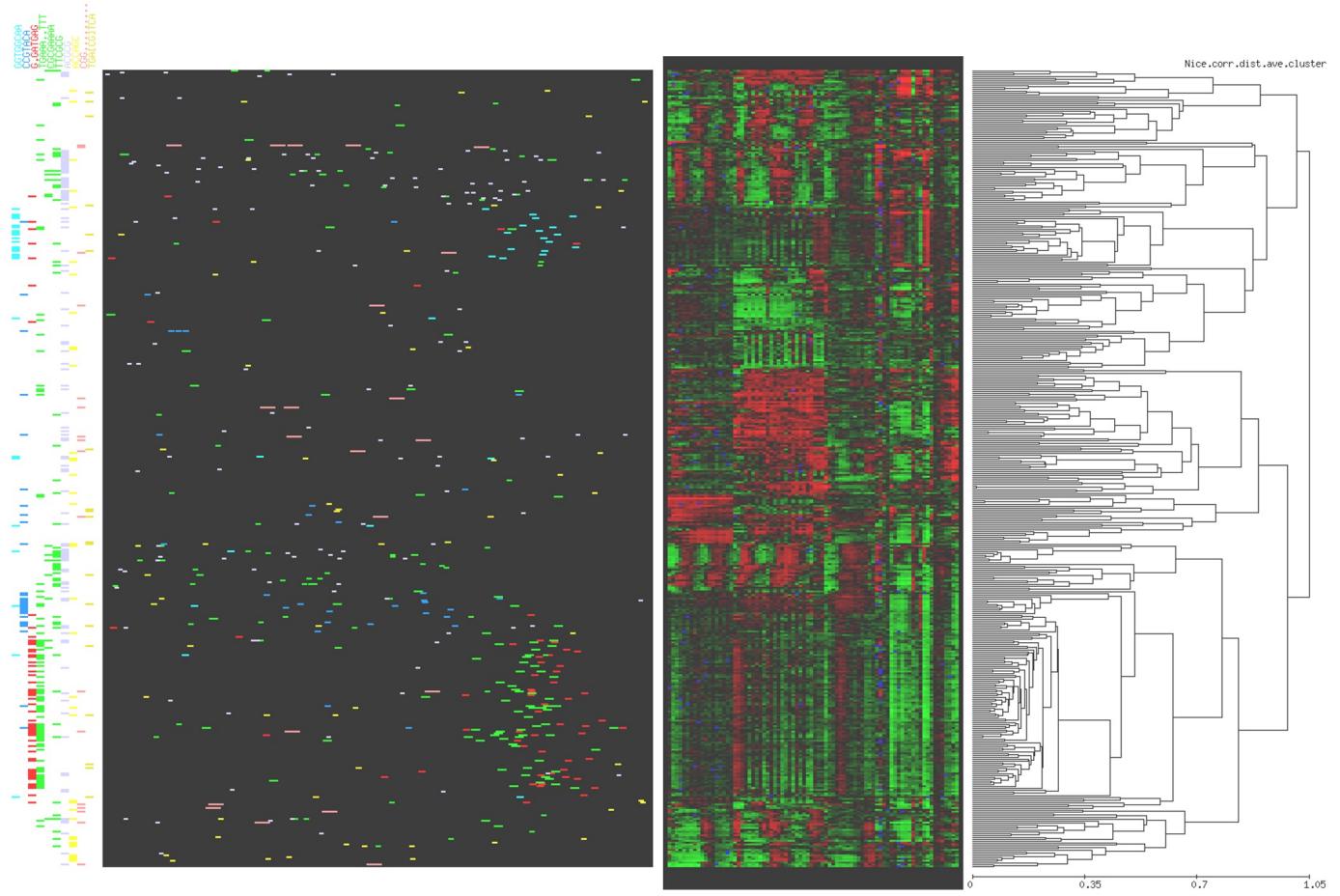
# nature



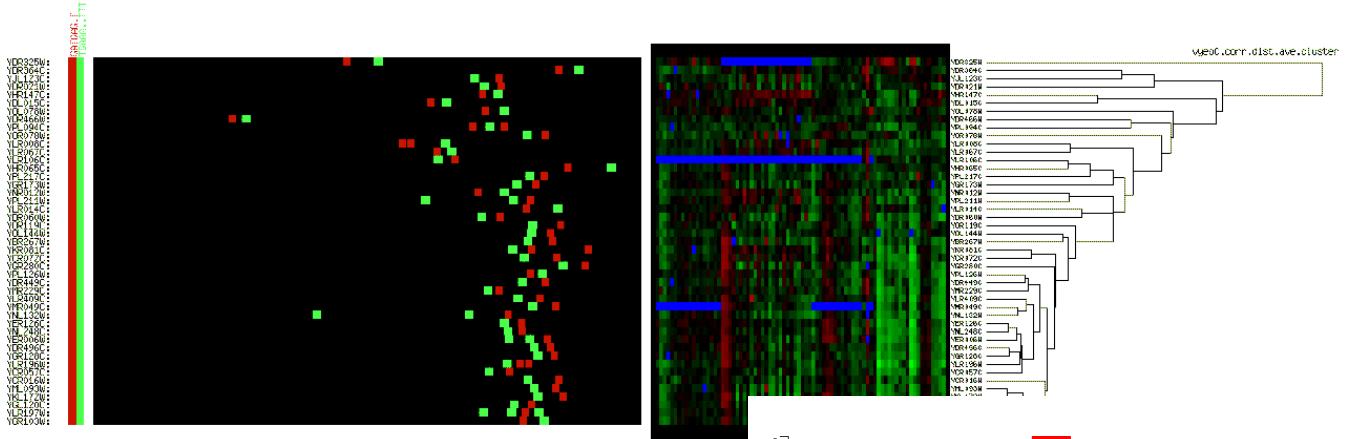
1997 – the first Eukaryotic genome ever sequenced



Microarrays – the beginning of functional genomics

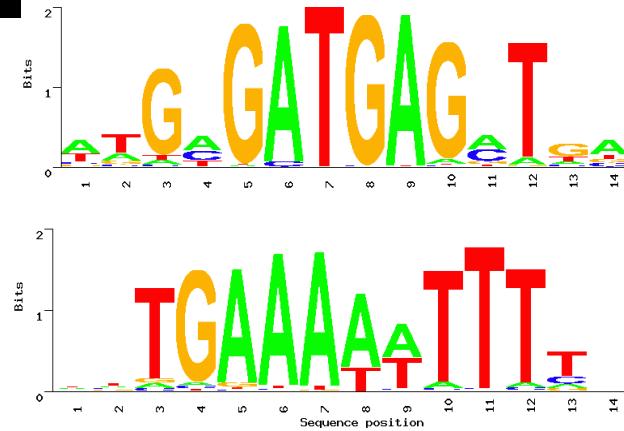


# The hypothesis – genes that are similarly expressed should be similarly regulated



GATGAG . T  
TGAAA . . TTT

Brazma A, Jonassen I, Vilo J, Ukkonen E. Predicting gene regulatory elements *in silico* on a genomic scale. *Genome Research* (1998)



**If you have a hammer but you see a screw, put the hammer aside and look for a screwdriver.**

(But do not forget where you have put your hammer, it may become handy one day)

# Lessons learned

- Trying to *explain as much as possible with as little as possible is* relevant to biology
- If you have a hammer but you see a screw, put the hammer aside and look for a screwdriver. But do not forget where you have put your hammer, it may become handy one day
- It's good to be in the field early!!!

## ChatGPT, write one paragraph about **the future of bioinformatics** in the light of AI

The future of bioinformatics is highly intertwined with the use of artificial intelligence (AI). The ability of AI to learn and analyze vast amounts of biological data presents a significant opportunity to advance our understanding of complex biological systems. With the use of machine learning algorithms and other AI techniques, bioinformatics researchers can make sense of large datasets, identify novel patterns and correlations, and develop predictive models for a wide range of applications in healthcare and life sciences. Furthermore, AI can also enable the integration of diverse data sources, such as genomics, transcriptomics, proteomics, and metabolomics, allowing researchers to gain a more comprehensive understanding of biological systems. Overall, the integration of AI in bioinformatics is poised to revolutionize the field, opening up new avenues for discovery and enabling researchers to develop more personalized and effective treatments for various diseases.





continue 16, 32, 64, 128, 256, ...



512, 1024, 2048, 4096, 8192, 16384, 32768, 65536, 131072, 262144, ...



continue 16, 1, 6, 32, 3, 2, 64, 6, 4, 128, 1, 2, 8, 256, 2, 5, 6, ...



This sequence doesn't seem to have an obvious pattern, but here are the next few terms:



512, 4, 8, 1024, 9, 10, 2048, 12, 16, 4096, 1, 18, 32, 8192, 20, 25, 18, 16384, 27, 4, 65536, 30, 32, 512, ...

Without more information, it's difficult to determine a clear rule or formula that generates this sequence. It's possible that the sequence is arbitrary or follows a complex pattern that is not immediately apparent.

To learn more what ChatGPT 4 can and cannot do, read my blog on LinkedIn