

# **Compression-based Machine Learning**

## **for DNA analysis**

Diogo Pratas *et al.*

pratas@ua.pt

<http://sweet.ua.pt/pratas>

IEETA  
UNIVERSITY OF AVEIRO

**DATA SCIENCE  
PORTUGAL**



**DS**  
PORTUGAL

# **What is compression ?**

The process of reducing storage

# **What is relative compression ?**

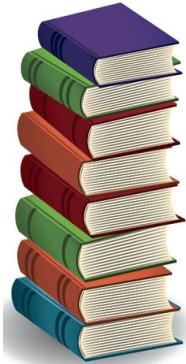
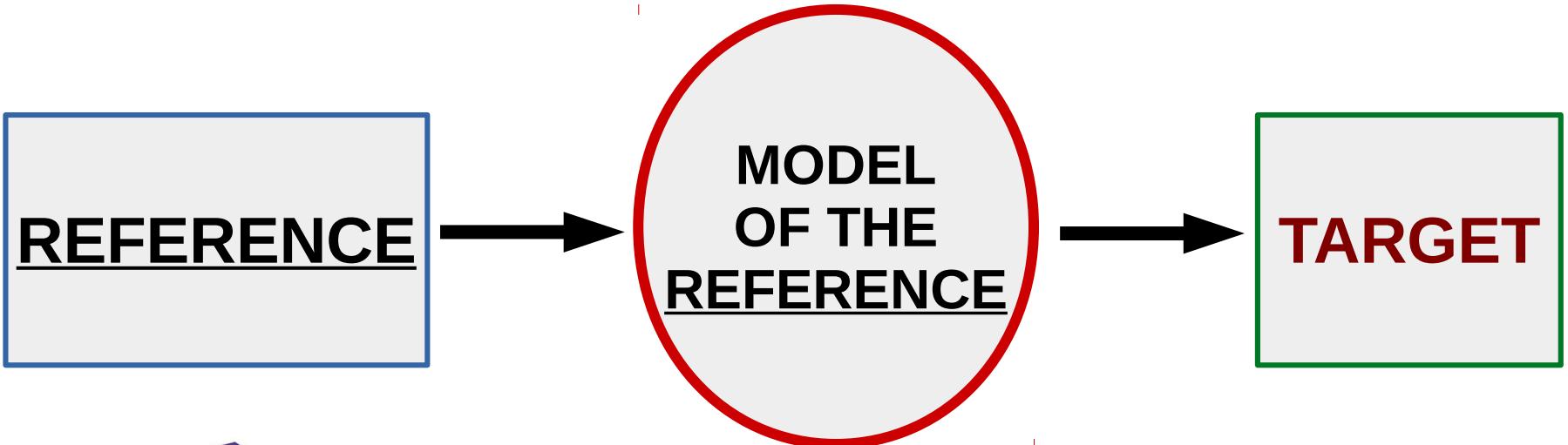
The process of reducing storage using  
exclusively a reference

**REFERENCE: THIS IS MORE THAN A WORD**



**TARGETS:**

THAT IS MORE THAN A WORD  
THESE ARE WORDS  
THOSE ARE SWORDS  
ARE YOU READING THIS ?



$$\log P(x^n|k) = \sum_{i=1}^{n-1} \log P(x_i|k, x^{i-1}) + \log P(x_n|k, x^{n-1})$$

$$\log p_{k,n} = \gamma \log p_{k,n-1} + \log P(x_n|k, x^{n-1}),$$

$$p_{k,n} = p_{k,n-1}^\gamma P(x_n|k, x^{n-1})$$

$$w_{k,n} = \frac{p_{k,n}}{\sum_{k \in \mathcal{K}} p_{k,n}}.$$

Genomic (DNA) sequences are large codified messages, from an alphabet of four symbols  $\Theta = \{A, C, G, T\}$ , describing most of the structure of all known living organisms. A huge amount of genomic data has been generated, with diverse characteristics, rendering the data deluge phenomenon a serious problem in most genomics centers. As such, most of the data are discarded (when possible), while others are compressed using general purpose algorithms, often attaining modest data reduction results.

Several specific algorithms have been proposed for the compression of genomic sequences, such as [1-7], but only some of them have been made available as usable and reliable compression tools, and those have been developed to some specific purpose or data characteristic [8-11].

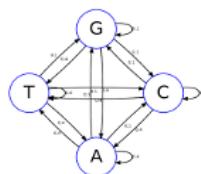
The dramatic increase of sequenced genomes, given the reduced sequencing costs, and the high redundancy characteristics, led to the development of genomic reference sequence compression. Several compressors of this type have been proposed [12-16], although most of them seem to be less efficient in handling sequences with higher rates of mutation.

## MODEL OF THE REFERENCE

**HOW WELL DOES IT LEARN  
ABOUT THE REFERENCE ?**

$$\log P(x^n|k) = \sum_{i=1}^{n-1} \log P(x_i|k, x^{i-1}) + \log P(x_n|k, x^{n-1})$$

$$\log p_{k,n} = \gamma \log p_{k,n-1} + \log P(x_n|k, x^{n-1}),$$



$$p_{k,n} = p_{k,n-1}^\gamma P(x_n|k, x^{n-1})$$

$$w_{k,n} = \frac{p_{k,n}}{\sum_{k \in \mathcal{K}} p_{k,n}}.$$

- COMPUTATIONAL TIME
- COMPUTATIONAL MEMORY
- COMPRESSION RATIO

**WE NEED EFFICIENT MODELS  
TO REPRESENT THE REFERENCE**

## **MARKOV MODEL (FINITE-CONTEXT MODEL)**

→ A model that provides **probability estimates** that depend on the recent past of the sequence

Portuguese Football League winners 2002-2016:

**9-Porto  
5-Benfica  
1-Sporting**

Who is going to win in 2017?

CONTEXT K=1

	P	B	S	TOTAL
P	6	3	0	9
B	2	3	0	5
S	1	0	0	1
TOTAL	10	6	0	

**SPPBPPPPPBPBPPPBBB?**

2002

2003

2004

2005

2006

2007

2008

2009

2010

2011

2012

2013

2014

2015

2016

2017

2017?

$$\Pr(P|B) = 2/5$$

$$\Pr(B|B) = 3/5$$

$$\Pr(S|B) = 0$$

(CONTEXT=1)

CONTEXT K=2

	P	B	S	Total
PP	3	3	0	6
PB	1	1	0	2
PS	0	0	0	0
BB	0	1	0	1
BP	1	0	0	1
BS	0	0	0	0
SS	0	0	0	0
SP	1	0	0	1
SB	0	0	0	0
Total	6	5	0	

**S P P B P P P P B P P P B B B ?**

2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017

2017?

$\Pr(P|BB)=0$

$\Pr(B|BB)=1$

$\Pr(S|BB)=0$

(CONTEXT=2)

**Which context Markov model is  
*better*? Context 1, 2, 3, 4, 5, ...?**

**Which context Markov model is  
*better*? Context 1, 2, 3, 4, 5, ...?**

→ It depends on the **nature**  
and **size** of the training data

Therefore, we consider **multiple**  
**Markov models** using different  
**contexts**

Approaches:

1. Competition
2. Cooperation

# Approaches:

- 1. Competition**
2. Cooperation

# 1. COMPETITION



What is the best model for  
the next 10 years?

...  
2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016  
SPPBPPPPBPBPPPBBB  
PAST

2017 2018 2019 2020 2021 2022 2023 2024 2025 2026  
**BSPBBBBBPPP** CONTEXT 1  
**BPPPPPPPBPP** CONTEXT 2  
**BBSBBBBBPS** CONTEXT 3  
**BBBBBBBBBBBB** CONTEXT 5  
  
**BBBBBBBBBBBB** REALITY

# 1. COMPETITION



What is the best model for  
the next 10 years?

...  
2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016  
SPPBPPPPBPBPPPBBB  
PAST

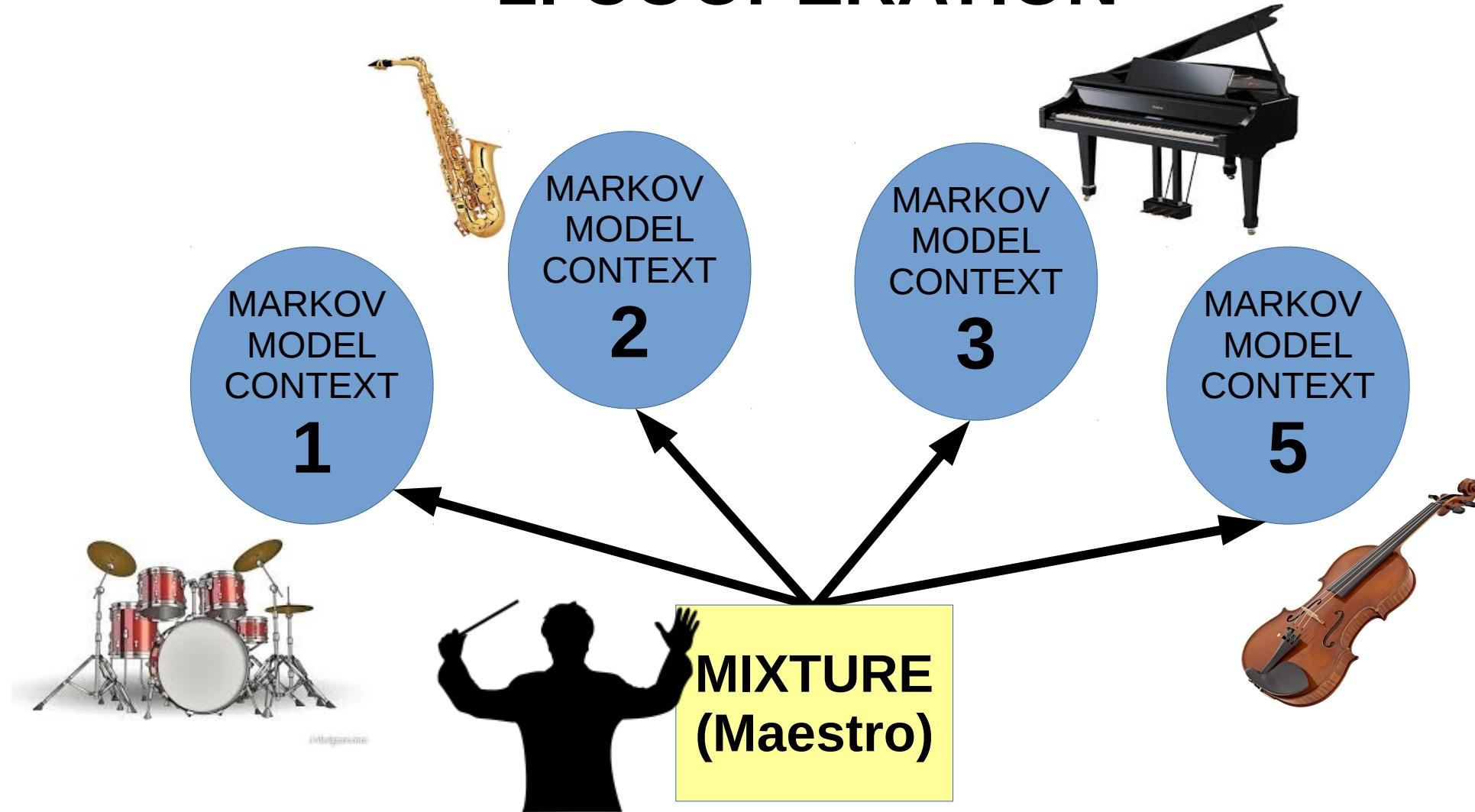
2017 2018 2019 2020 2021 2022 2023 2024 2025 2026  
**BSPBBBBBPPP** CONTEXT 1  
**BPPPPPPPBPP** CONTEXT 2  
**BBSBBBBBPS** CONTEXT 3  
**BBBBBBBBBBBB** CONTEXT 5  
BBBBBBBBBBBB REALITY



# Approaches:

1. Competition
2. Cooperation

## 2. COOPERATION



**EACH MODEL IS WEIGHTED ACCORDING  
TO ITS PERFORMANCE ALONG THE TIME**

**Generally cooperation works better  
than competition...**

**specially in DNA sequences**

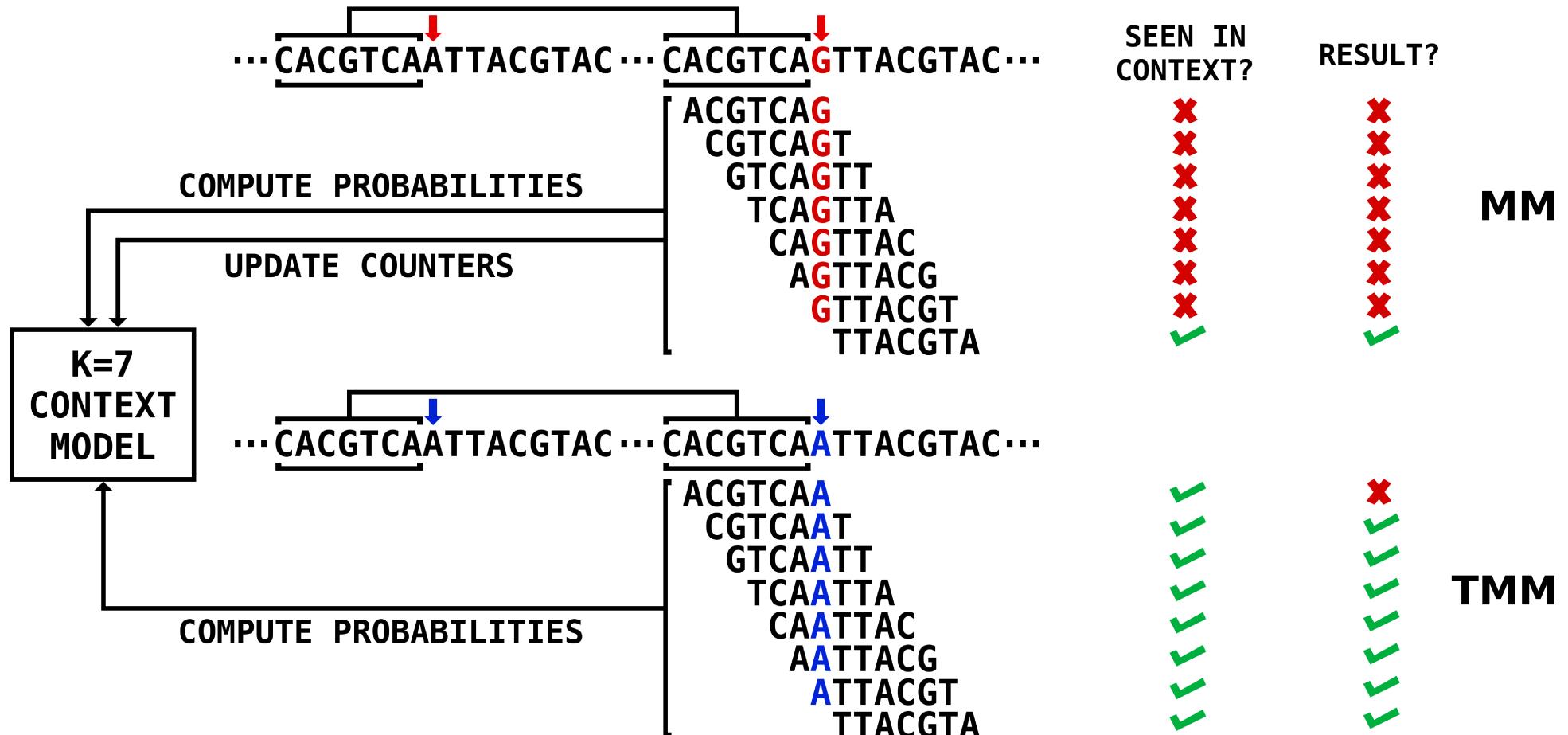
# Tolerant Markov Model (TMM)



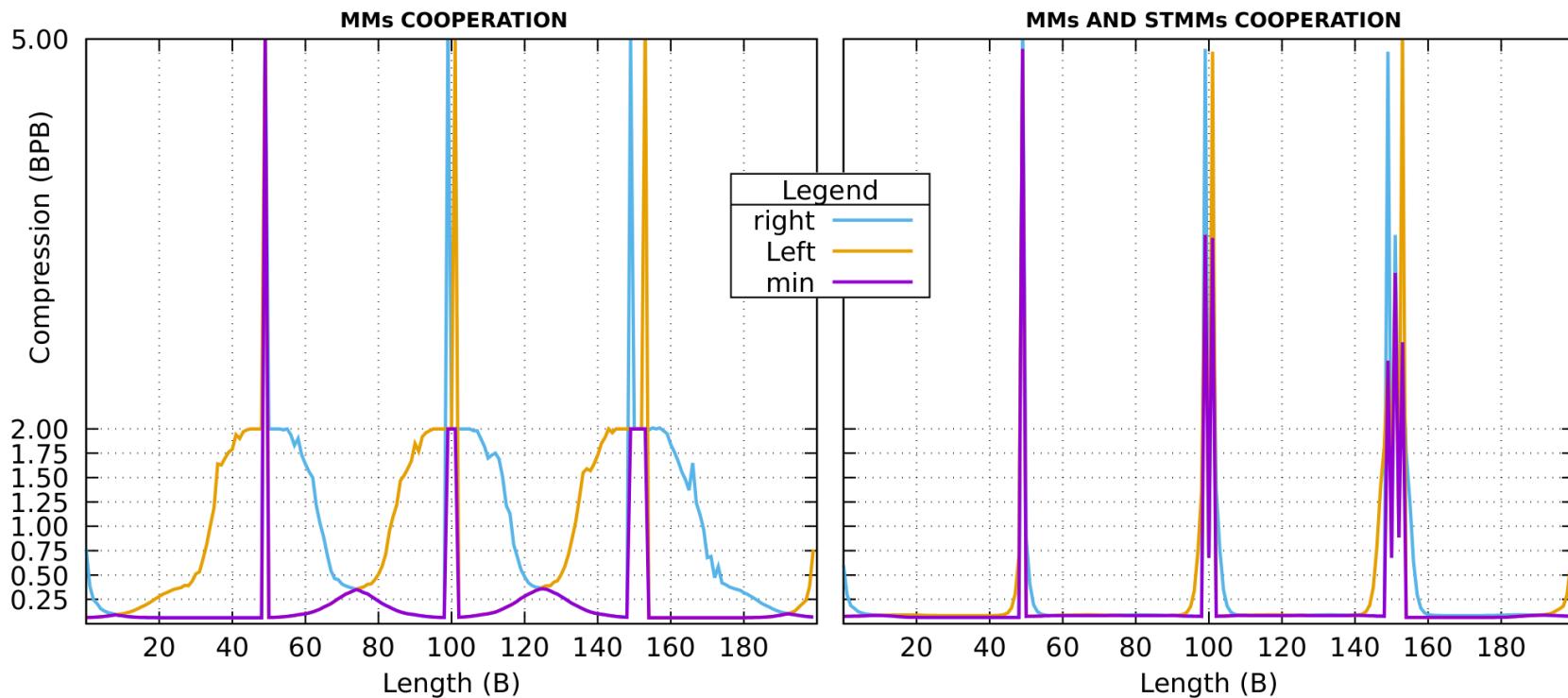
a probabilistic-algorithmic Markov model

It assigns probabilities according to a conditioning context that considers the last symbol, from the sequence to occur, as the most probable, given the occurrences stored in the memory, such as those from y, instead of the true occurring symbol.

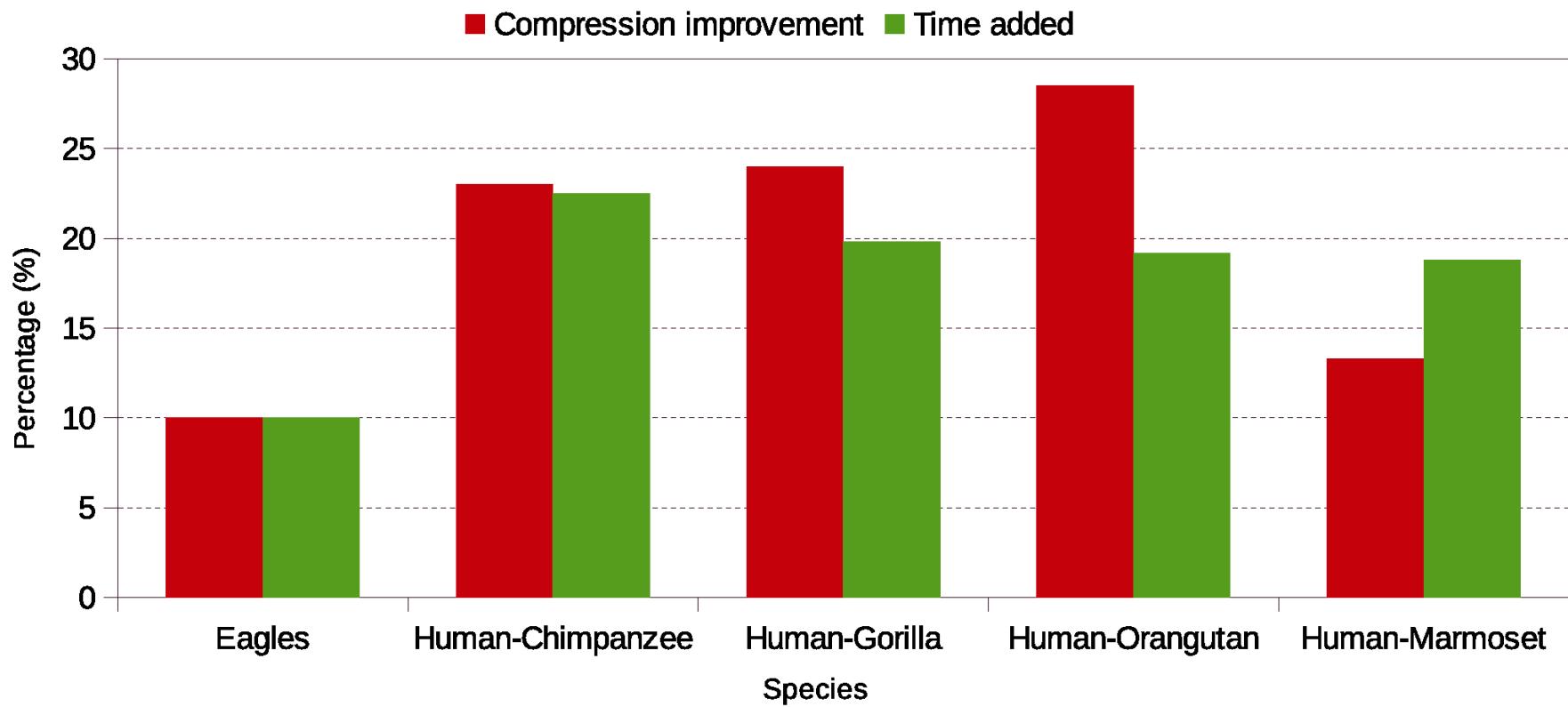
# Tolerant Markov Model (TMM)



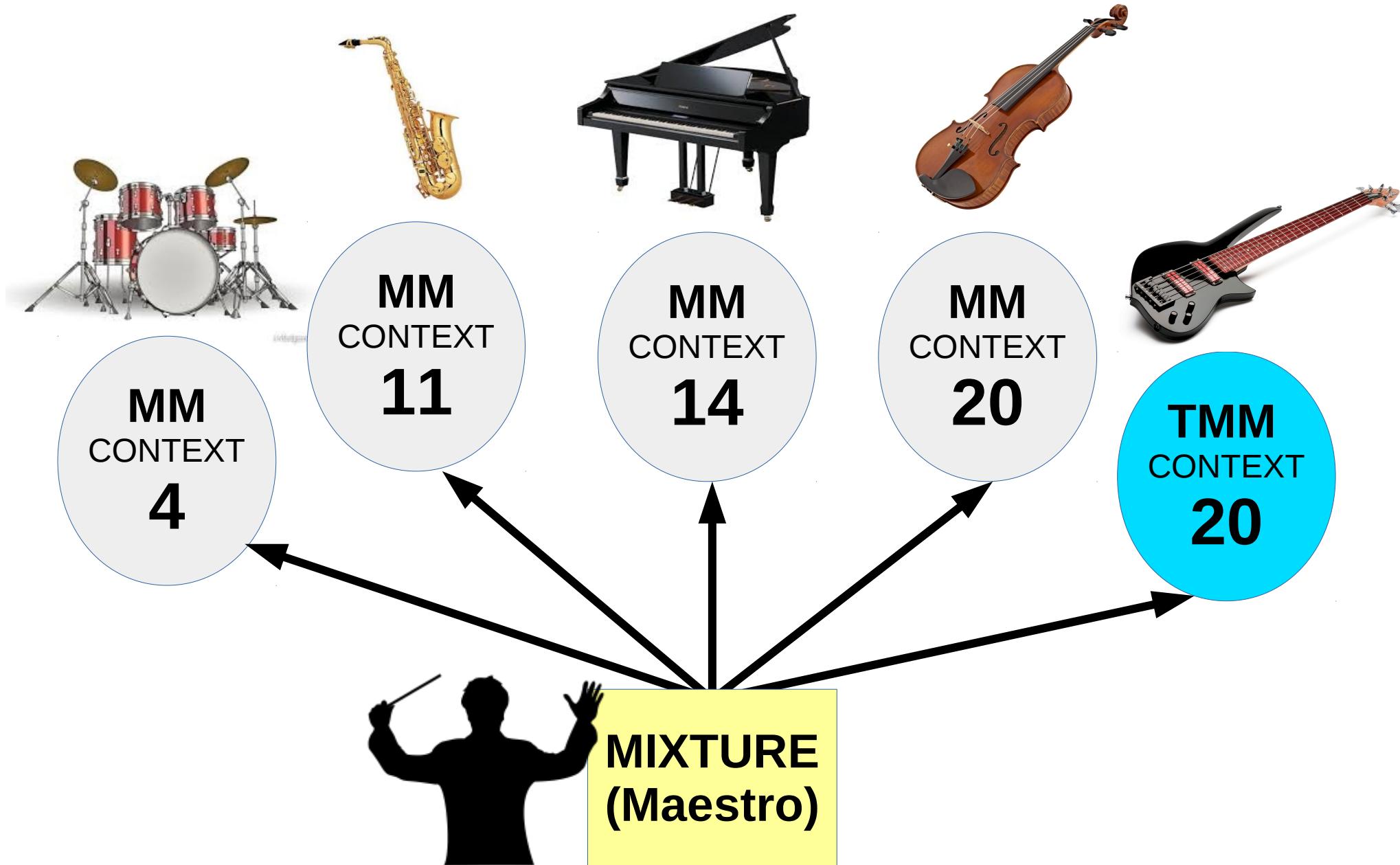
# Tolerant Markov Model (TMM)



# Tolerant Markov Model (TMM)

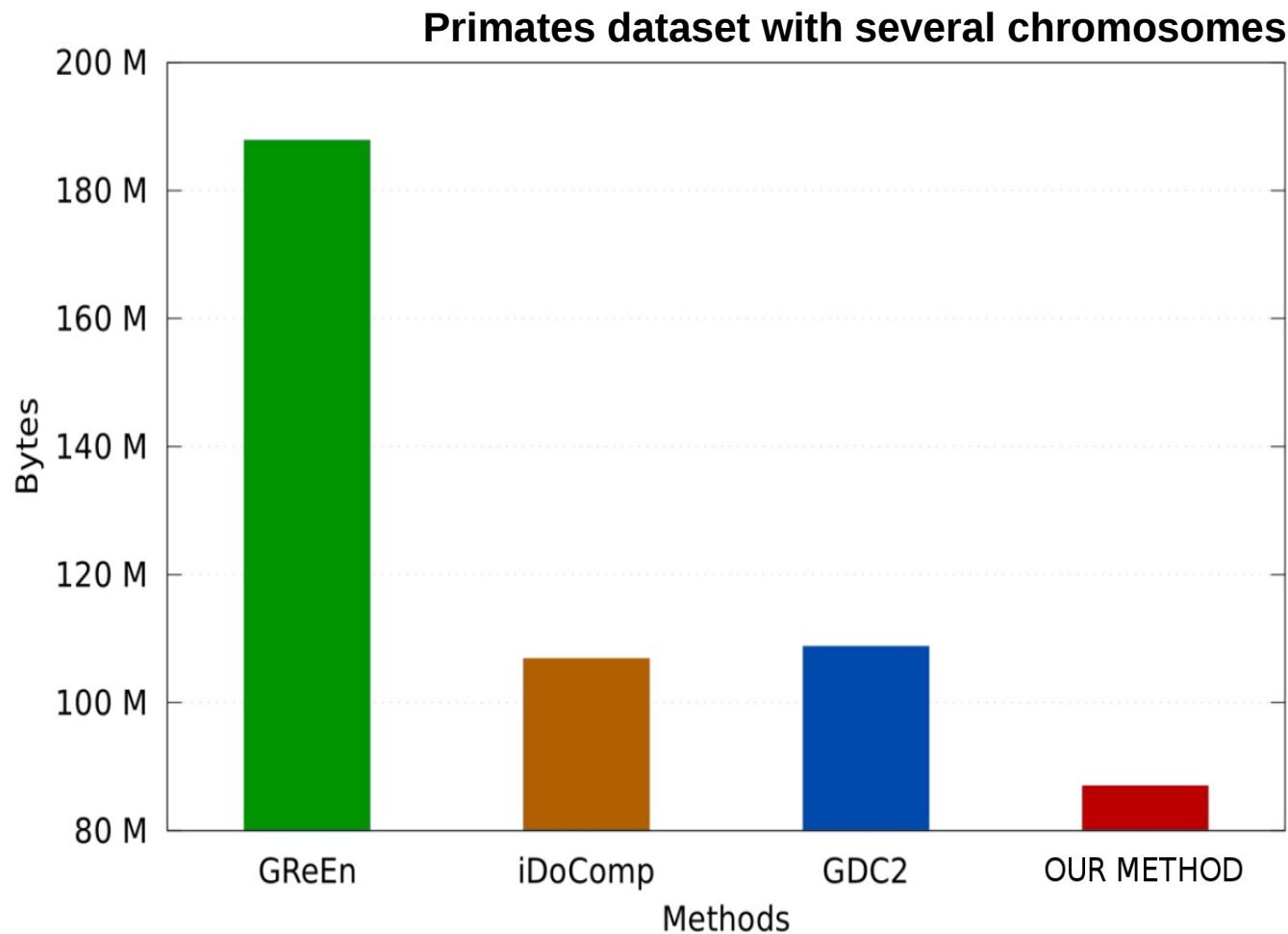


# DNA relative compressor

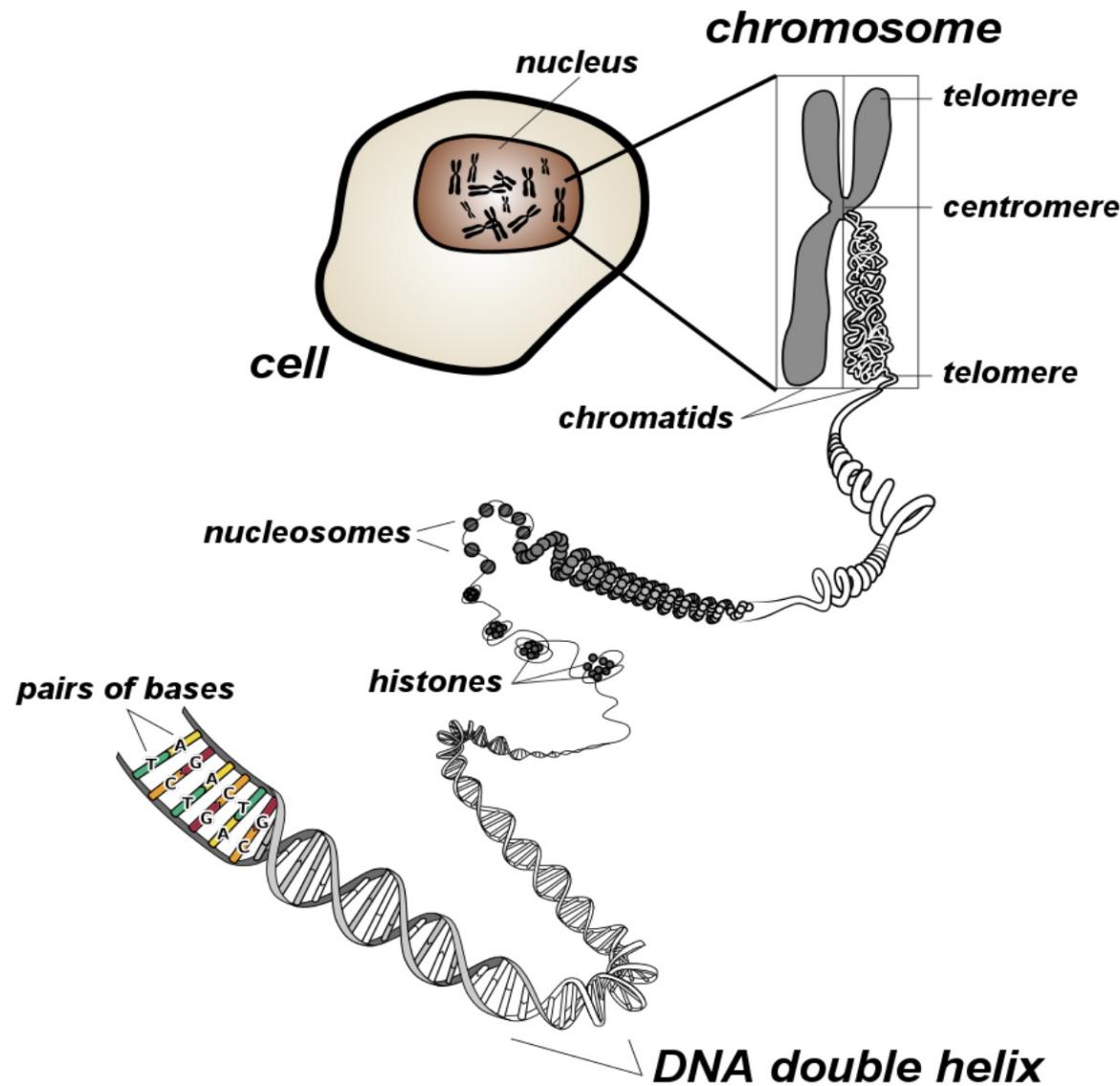


We use this compressor to measure relative similarity

# DNA relative compressor



# DNA sequences

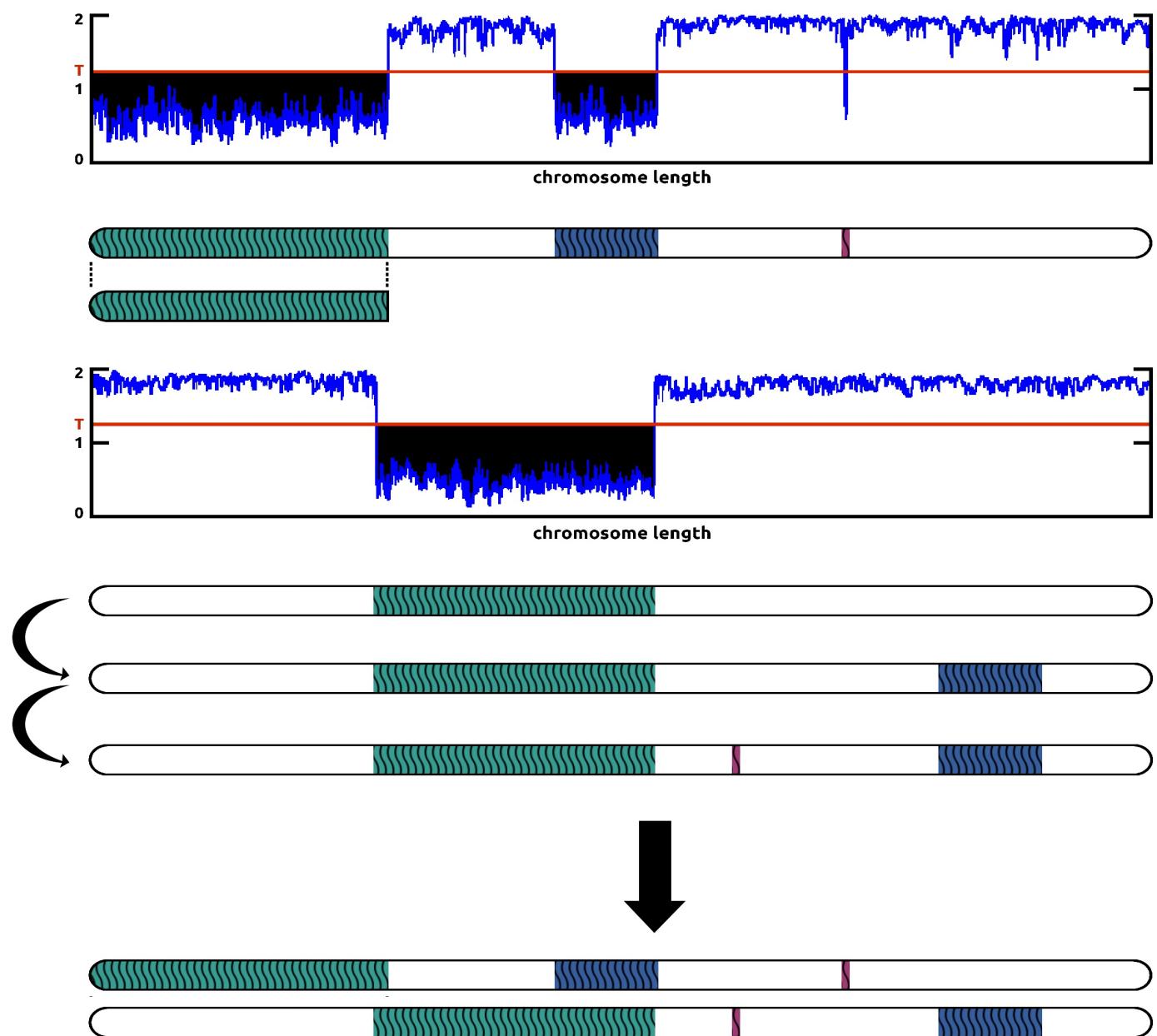


# DNA sequences

A grid of DNA sequence data showing multiple rows of nucleotide bases (A, T, C, G) in a color-coded format. The grid consists of approximately 15 rows and 40 columns. Each row represents a different DNA sequence. The colors used are: A (blue), T (red), C (green), and G (orange). The sequence "TGGAAACGGGACGCCAT" is repeated multiple times across the rows.

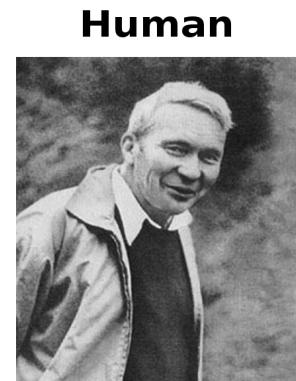
A	O	T	G	A	G	T	T	C	C	C	T	G	G	A	A	C	G	G	G	A	C	G	C	C	A	T	A	A	A	A	A	A	
T	A	O	T	G	A	G	T	T	C	C	C	T	G	G	A	A	C	G	G	G	A	C	G	C	C	A	T	A	A	A	A	A	A
C	C	G	T	C	T	G	G	T	A	G	G	A	C	A	C	C	C	A	G	C	C	C	G	G	G	A	T	A	A	A	A	A	A
T	T	C	C	G	A	G	T	T	C	C	C	T	G	G	A	A	C	G	G	G	A	C	G	C	C	A	T	A	A	A	A	A	A
C	T	T	C	C	G	A	G	T	T	C	C	C	T	G	G	A	A	C	G	G	G	A	C	G	C	C	A	T	A	A	A	A	A
T	C	C	G	A	G	T	T	C	C	C	T	G	G	A	A	C	G	G	G	A	C	G	C	C	A	T	A	A	A	A	A	A	
G	G	A	T	A	A	C	C	G	T	G	G	T	A	A	T	T	C	T	A	G	A	G	C	C	C	G	G	G	A	A	A	A	A
A	C	G	C	C	A	T	A	G	A	G	G	G	T	G	A	G	A	G	G	C	C	C	C	G	G	G	A	A	A	A	A	A	
T	T	C	C	G	A	G	T	T	C	C	C	T	G	G	A	A	C	G	G	G	A	C	G	C	C	A	T	A	A	A	A	A	
C	G	G	A	C	G	C	C	A	T	A	G	A	G	G	T	G	A	G	A	G	C	C	C	C	G	G	G	A	A	A	A	A	A
C	G	T	C	T	G	G	T	A	G	G	A	C	A	C	C	C	A	G	C	C	C	G	G	G	A	A	A	A	A	A	A	A	

# Similar regions detection



# Rearrangements detection

Chromosome  
**5**

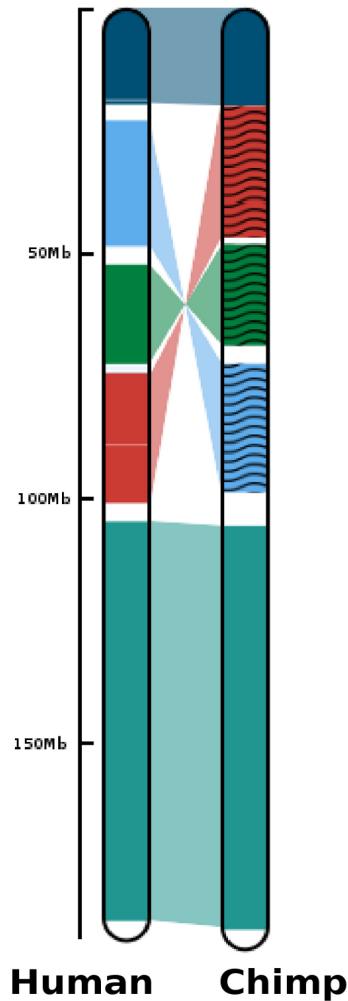


**Vs**

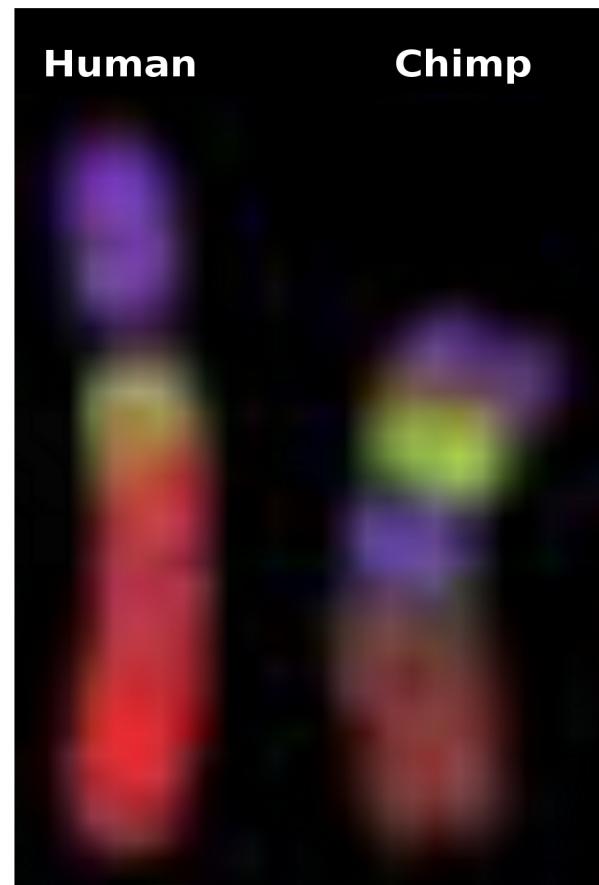


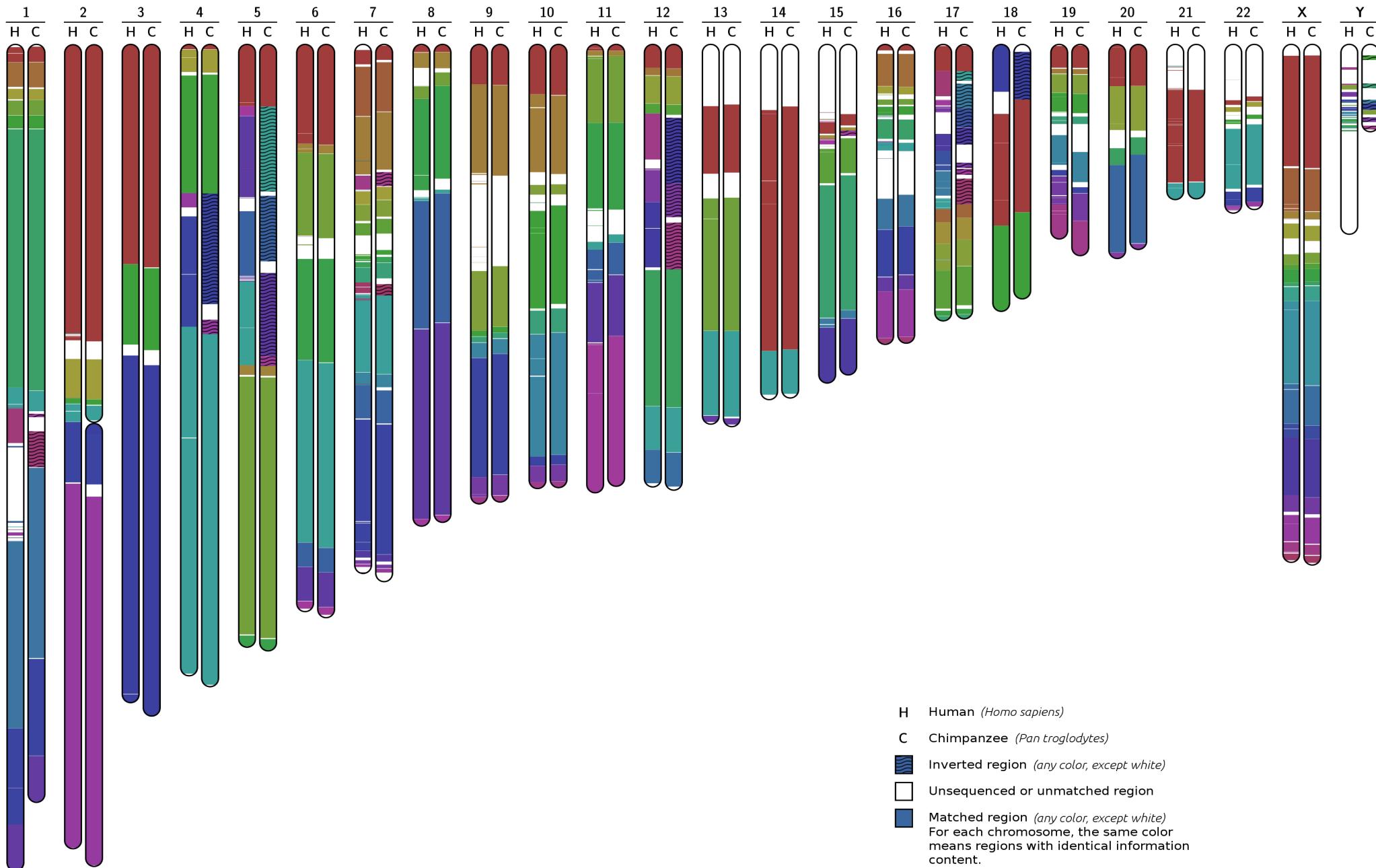
Chimp

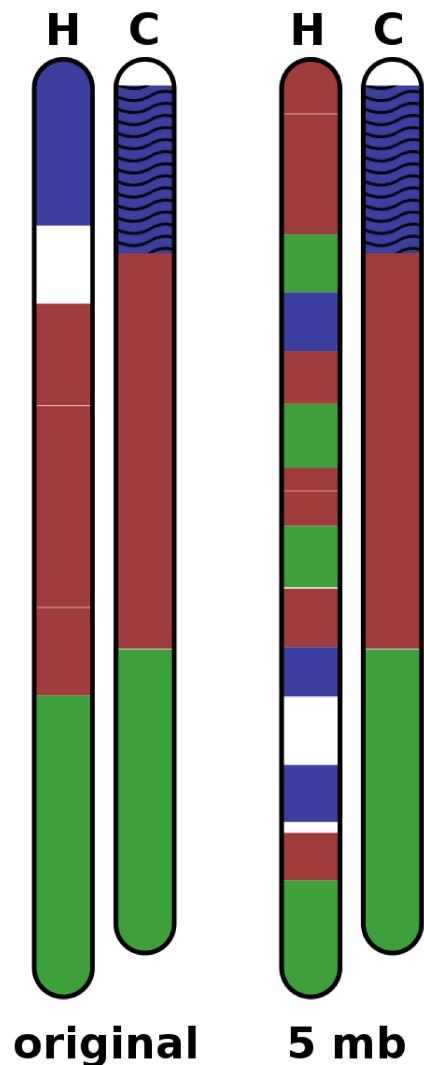
Computational approach

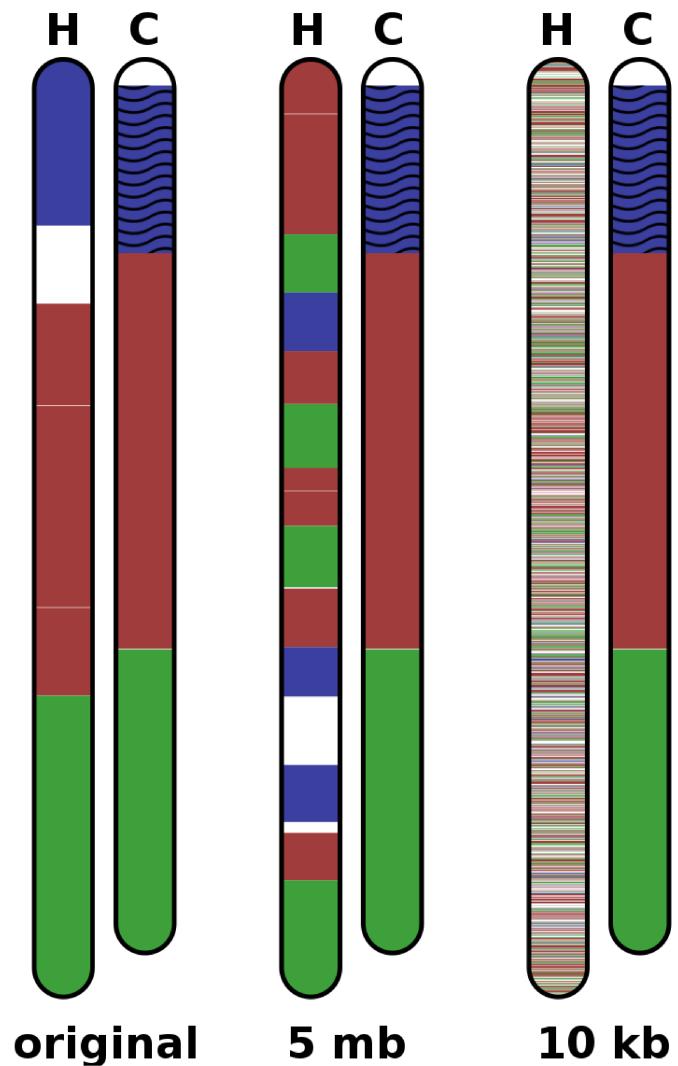


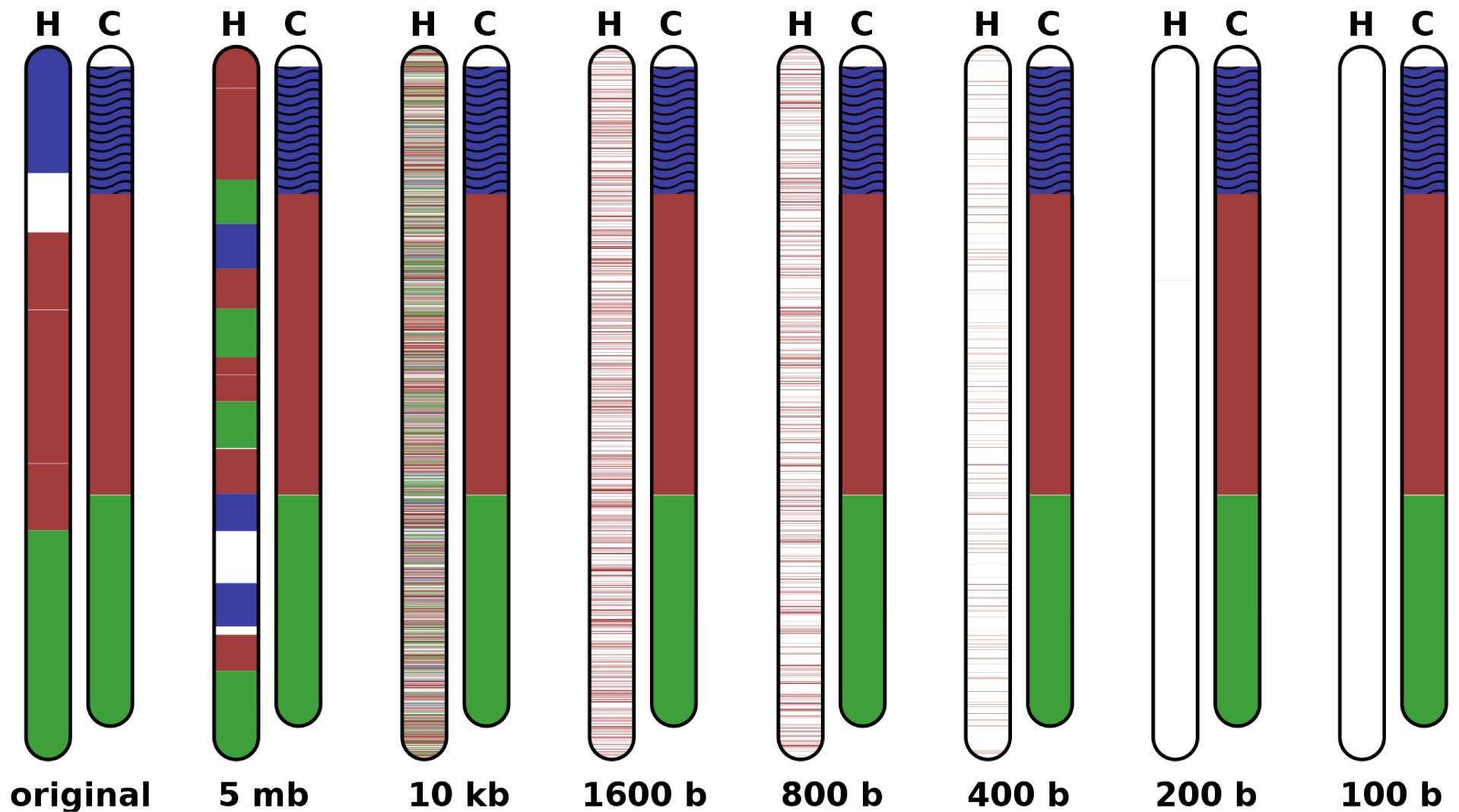
FISH (biological approach)











# DNA extraction kit



# FASTQ Reads

```
@SRR038845.3 HWI-EAS038:6:1:0:1938 length=36
CAAGCAGTTCACACCTTGGCCGACAGGCCGGTAA
+SRR038845.3 HWI-EAS038:6:1:0:1938 length=36
BA@>B=>:>>7@7@>9=BAA?;>52;>:9=8.=A
@SRR038845.41 HWI-EAS038:6:1:0:1474 length=36
CCAATGATTTCGTTTGTCAAATACGGTAA
+SRR038845.41 HWI-EAS038:6:1:0:1474 length=36
BCCBA@BB@BBBBBAB@B9B@=BABA@A:@693:@B=
@SRR038845.53 HWI-EAS038:6:1:1:360 length=36
GTTCAAAAAGAACTAAATTGTGTCAATAGAAAACTC
+SRR038845.53 HWI-EAS038:6:1:1:360 length=36
BBCBBBBB@qBAB?BBBBBCB>BBBAA>BBBAA@
```

# DNA Sequence



# Assemble reads

The local structure of turbulence in incompressible viscous fluid for very large Reynolds numbers<sup>1</sup>

By A. N. KOLMOGOROV

<sup>1</sup>) We shall denote by

$$u_j^P = u_j(x_1, x_2, x_3, t), \quad j = 1, 2, 3.$$

The components of velocity at the instant  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In compressible fluid we have to subtract the components  $u_1, u_2, u_3$  of the mean flow. Let us assume that the components  $u_1, u_2, u_3$  of the mean flow  $P$  at every point  $P = (x_1, x_2, x_3)$  of the domain  $G$  of the first-dimensional space  $(x_1, x_2, x_3)$  are random variables in the sense of the theory of probability. Then the components  $u_j^P$  of the random field  $u^P$  are random variables in the same sense. This follows from the definition of the random field  $u^P$  given by Kolmogorov (1939).

Denoting by  $J$  the mathematical expectation of the random variable we see that

$$u_j^P = J(u_j^P) \quad \text{and} \quad \langle u_k^P \rangle = 0.$$

We shall denote by  $\delta_{jk}$  every bounded subdomain of the domain  $G$ .

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  the new coordinate

$$y_{\alpha} = x_{\alpha} - (t - T)^{\mu_{\alpha}}(x_1^{\mu_1} - x_1^{\mu_2}), \quad (\alpha = 1, 2, 3), \quad (1)$$

where  $T$  is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_{\alpha}$  of any point  $P$  depend on the random variables  $x_1^{\mu_1}, x_1^{\mu_2}$  and hence are themselves random variables. The velocity components  $u_j^P$  are random variables in the same sense.

$$u_j^P = u_j(y_1^{\mu_1}, y_1^{\mu_2}, t). \quad (2)$$

Suppose that for some fixed values of  $x_1^{\mu_1}, x_1^{\mu_2}$  the points  $y_1^{\mu_1}, y_1^{\mu_2}, t = 1, 2, \dots$  having the uniform distribution in the domain  $\delta_{jk}$  are statistically independent. Then we may take a 3-dimensional distribution law of probabilities  $F_{jk}$  for the quantities

$$u_j^P = u_j(y_1^{\mu_1}), \quad x = 1, 2, \dots, \quad k = 1, 2, \dots, n,$$

where  $y_1^{\mu_1} = x_1^{\mu_1} - (t - T)^{\mu_1}(x_1^{\mu_1} - x_1^{\mu_2})$ .

Physically speaking, the distribution law  $F_{jk}$  depends on the parameters  $x_1^{\mu_1}, x_1^{\mu_2}, t$ .

**Definition 1.** The solution  $u^P$  is called *locally homogeneous* in the domain  $G$ , if for every point  $P$  of  $G$  the distribution law  $F_{jk}$  is independent from  $x_1^{\mu_1}, x_1^{\mu_2}$  as long as all points  $y_1^{\mu_1}$  are chosen in  $\delta_{jk}$ .

<sup>1</sup>) Published in Russian in *Zhurn. vopr. stat. i teor. mekhan.* No. 1, Paper received 26 December 1948. The translation is based on the Russian edition which contains the details of this article.

<sup>2</sup>) *Prob. teor. i mat. statist.* 1946, No. 3-4.

<sup>3</sup>) *Prob. teor. i mat. statist.* 1946, No. 3-4.

# The local structure of turbulence in incompressible viscous fluid for large Reynolds numbers†

By A. N. OGOROV

§1. We shall denote by

$$u_\alpha(P) = u_\alpha(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_\alpha(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are *random variables* in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$\bar{u}_\alpha^2 \text{ and } \overline{(du_\alpha/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain  $G$ .

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} y_\alpha &= x_\alpha - x_\alpha^{(0)} - u_\alpha(P^{(0)}) (t - t^{(0)}), \\ s &= t - t^{(0)}, \end{aligned} \right\}$$

where

$$P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t)$$

is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_\alpha$  of any point  $P$  depend on the random variables  $u_\alpha(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_\alpha(P) = u_\alpha(P) - u_\alpha(P^{(0)}). \quad (2)$$

Suppose that for some fixed values of  $u_\alpha(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_\alpha^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_n$  for the quantities

$$w_\alpha^{(k)} = w_\alpha(P^{(k)}), \quad \alpha = 1, 2, 3; \quad k = 1, 2, \dots, n$$

where

$$u_\alpha^{(0)} = u_\alpha(P^{(0)})$$

are given.

Generally speaking, the distribution law  $F_n$  depends on the parameters  $x_\alpha^{(0)}, t^{(0)}, u_\alpha^{(0)}$ ,  $y_\alpha^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_\alpha^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_n$  is independent from  $x_\alpha^{(0)}, t^{(0)}$  and  $u_\alpha^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

† First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.

# The local structure of turbulence in incompressible viscous fluid for very large Reynolds numbers†

By A. N. KOLMOGOROV

§1. We shall denote by

$$u_\alpha(P) = u_\alpha(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_\alpha(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are *random variables* in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$\bar{u}_\alpha^2 \text{ and } \overline{(du_\alpha/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} y_\alpha &= x_\alpha - x_\alpha^{(0)} - u_\alpha(P^{(0)}) (t - t^{(0)}), \\ s &= t - t^{(0)}, \end{aligned} \right\} \quad (1)$$

where

$$P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t)$$

is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_\alpha$  of any point  $P$  depend on the random variables  $u_\alpha(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_\alpha(P) = u_\alpha(P) - u_\alpha(P^{(0)}).$$

Suppose that for some fixed values of  $u_\alpha(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_\alpha^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_n$  for the quantities

$$w_\alpha^{(k)} = w_\alpha(P^{(k)}), \quad \alpha = 1, 2, 3; \quad k = 1, 2, \dots, n,$$

where

$$u_\alpha^{(0)} =$$

are given.

Generally speaking, the distribution law  $F_n$  depends on the parameters  $x_\alpha^{(0)}, t^{(0)}, u_\alpha^{(0)}$ ,  $y_\alpha^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_\alpha^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_n$  is independent from  $x_\alpha^{(0)}, t^{(0)}$  and  $u_\alpha^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

† First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.

$$u_a(P) = u_a(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

city at the moment  $t$  at the point with rectangular cartesian coordinates  $(x_1, x_2, x_3)$ . In considering the turbulence it is natural to assume that the velocity  $u_a(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are *random variables* in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Mathematical expectation of the random variable  $A$  we suppose

$$w_a^2 \text{ and } \overline{(du_a/dx_\beta)^2}$$

in every bounded subdomain of the domain  $G$ . Introducing new coordinates

$$\begin{aligned} x_a - x_a^{(0)} - u_a(P^{(0)}) (t - t^{(0)}), \\ s = t - t^{(0)}, \end{aligned} \quad \left. \begin{aligned} P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t) \end{aligned} \right\}$$

from the domain  $G$ . Observe that the coordinates  $y_a$  of any random variables  $u_a(P^{(0)})$  and hence are themselves random variables in the new coordinates are

$$w_a(P) = u_a(P) - u_a(P^{(0)}). \quad (2)$$

for fixed values of  $u_a(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_a^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_a$  for the

accessible members†

Definition 1. The probability that the random variables  $y_a^{(k)}$  and  $s^{(k)}$  are given is called the probability of the event

$$P^{(k)}, \quad \alpha = 1, 2, 3; \quad k = 1, \dots, n, \quad u_a^{(0)} = u_a(P^{(0)})$$

$$(y_a^{(k)}, s^{(k)})$$

the distribution law of the random variables  $y_a^{(k)}$  and  $s^{(k)}$  is called the distribution law of the random variables  $y_a^{(k)}$  and  $s^{(k)}$ .

The parameters  $x_a^{(0)}, t^{(0)}, u_a^{(0)}$  are called the initial parameters.

in incompressible fluid for very large Reynolds numbers

By A. N.

## structure of turbulence for very large Reynolds numbers

### The local structure of viscous fluid

### The local structure of turbulence in incompressible viscous fluid for very large Reynolds numbers†

By A. N. KOLMOGOROV

§1. We shall denote by

$$u_a(P) = u_a(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_a(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are *random variables* in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$w_a^2 \text{ and } \overline{(du_a/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain  $G$ . Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\begin{aligned} x_a - x_a^{(0)} - u_a(P^{(0)}) (t - t^{(0)}), \\ s = t - t^{(0)}, \end{aligned} \quad \left. \begin{aligned} P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t) \end{aligned} \right\} \quad (1)$$

where  $w_a^2 = w_a^2(P^{(0)})$  and  $P^{(0)}$  is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_a$  of any point  $P$  depend on the random variables  $u_a(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_a(P) = u_a(P) - u_a(P^{(0)}). \quad (2)$$

Suppose that for some fixed values of  $u_a(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_a^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_a$  for the quantities

$$w_a^{(k)} = w_a(P^{(k)}), \quad \alpha = 1, 2, 3; \quad k = 1, 2, \dots, n,$$

where

$$u_a^{(0)} =$$

are given.

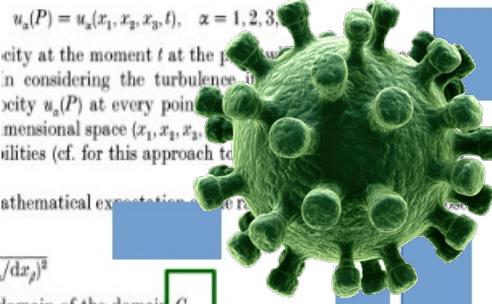
Generally speaking, the distribution law  $F_a$  depends on the parameters  $x_a^{(0)}, t^{(0)}, u_a^{(0)}$ ,  $y_a^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_a^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_a$  is independent from  $x_a^{(0)}, t^{(0)}$  and  $u_a^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.

Proc. R. Soc. Lond. A (1991) 434, 9–13

Printed in Great Britain



in every bounded subdomain of the domain  $G$ , in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} x_a - x_a^{(0)} - u_a(P^{(0)}) (t - t^{(0)}), \\ s = t - t^{(0)}, \end{aligned} \right\}$$

$$P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t)$$

from the domain  $G$ . Observe that the coordinates  $y_a$  of any random variables  $u_a(P^{(0)})$  and hence are themselves random variables in the new coordinates are

$$w_a(P) = u_a(P) - u_a(P^{(0)}). \quad (2)$$

at fixed values of  $u_a(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_a^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_a$  for the

accessible members†

Definition 1. The probability that the point  $P$  is situated in the domain  $G$  is given by the formula

$$F_a(P^{(0)}), \quad \alpha = 1, 2, 3; \quad k = 1, \dots, n,$$

$$u_a^{(0)} = u_a(P^{(0)})$$

$$l_a w_a = \frac{d}{dt}$$

The parameters  $x_a^{(0)}, t^{(0)}, u_a^{(0)}$

are called the initial parameters of the process. The initial conditions for the velocity components  $w_a$  are given by the formula

Proc. R. Soc. Lond. A (1991) 434, 9–13  
Printed in Great Britain

## The local structure of turbulence in incompressible viscous fluid for very large Reynolds numbers†

By A. N. KOLMOGOROV

§1. We shall denote by

$$u_a(P) = u_a(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_a(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are random variables in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$u_a^2 \text{ and } \overline{(du_a/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain  $G$ . Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} y_a = x_a - x_a^{(0)} - u_a(P^{(0)}) (t - t^{(0)}), \\ s = t - t^{(0)}, \end{aligned} \right\} \quad (1)$$

where  $P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t)$

is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_a$  of any point  $P$  depend on the random variables  $u_a(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_a(P) = u_a(P) - u_a(P^{(0)}). \quad (2)$$

Suppose that for some fixed values of  $u_a(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_a^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_a$  for the quantities

$$w_a^{(k)} = w_a(P^{(k)}), \quad \alpha = 1, 2, 3; \quad k = 1, 2, \dots, n,$$

where

$$u_a^{(k)} =$$

are given.

Generally speaking, the distribution law  $F_a$  depends on the parameters  $x_a^{(0)}, t^{(0)}, u_a^{(0)}$ ,  $y_a^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_a^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_a$  is independent from  $x_a^{(0)}, t^{(0)}$  and  $u_a^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

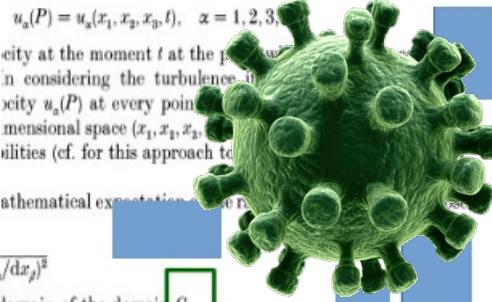
† First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.

Proc. R. Soc. Lond. A (1991) 434, 9–13

Printed in Great Britain

# **Ancient DNA**

**Much more complex...**



in every bounded subdomain of the domain  $G$ , in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} x_a - x_a^{(0)} - u_a(P^{(0)}) (t - t^{(0)}), \\ s = t - t^{(0)}, \end{aligned} \right\}$$

$$P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t)$$

from the domain  $G$ . Observe that the coordinates  $y_a$  of any random variables  $u_a(P^{(0)})$  and hence are themselves random variables in the new coordinates are

$$w_a(P) = u_a(P) - u_a(P^{(0)}). \quad (2)$$

at fixed values of  $u_a(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_a^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_a$  for the

accessible members†

Definition 1. The probability that the point  $P$  is situated in the domain  $G$  is given by the formula

$$F_a(P^{(0)}), \quad \alpha = 1, 2, 3; \quad k = 1, \dots, n,$$

$$u_a^{(0)} = u_a(P^{(0)})$$

where  $x_a^{(0)}, t^{(0)}$  are the parameters of the point  $P^{(0)}$ ;  $y_a^{(k)}$  are the coordinates of the point  $P^{(k)}$  in the new coordinates;  $w_a^{(k)}$  are the coordinates of the point  $P^{(k)}$  in the new coordinates;  $u_a^{(k)}$  are the components of velocity at the moment  $t^{(0)}$  at the point  $P^{(0)}$ .

Proc. R. Soc. Lond. A (1991) 434, 9–13  
Printed in Great Britain

## The local structure of turbulent viscous fluid for very large Reynolds numbers†

By A. N. KOLMOGOROV

§1. We shall denote by

$$u_a(P) = u_a(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_a(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are random variables in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$u_a^2 \text{ and } \overline{(du_a/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} y_a = x_a - x_a^{(0)} - u_a(P^{(0)}) (t - t^{(0)}), \\ s = t - t^{(0)}, \end{aligned} \right\} \quad (1)$$

where

$$P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t)$$

is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_a$  of any point  $P$  depend on the random variables  $u_a(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_a(P) = u_a(P) - u_a(P^{(0)}). \quad (2)$$

Suppose that for some fixed values of  $u_a(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_a^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_a$  for the quantities

$$w_a^{(k)} = w_a(P^{(k)}), \quad \alpha = 1, 2, 3; \quad k = 1, 2, \dots, n,$$

where

$$u_a^{(k)} =$$

are given.

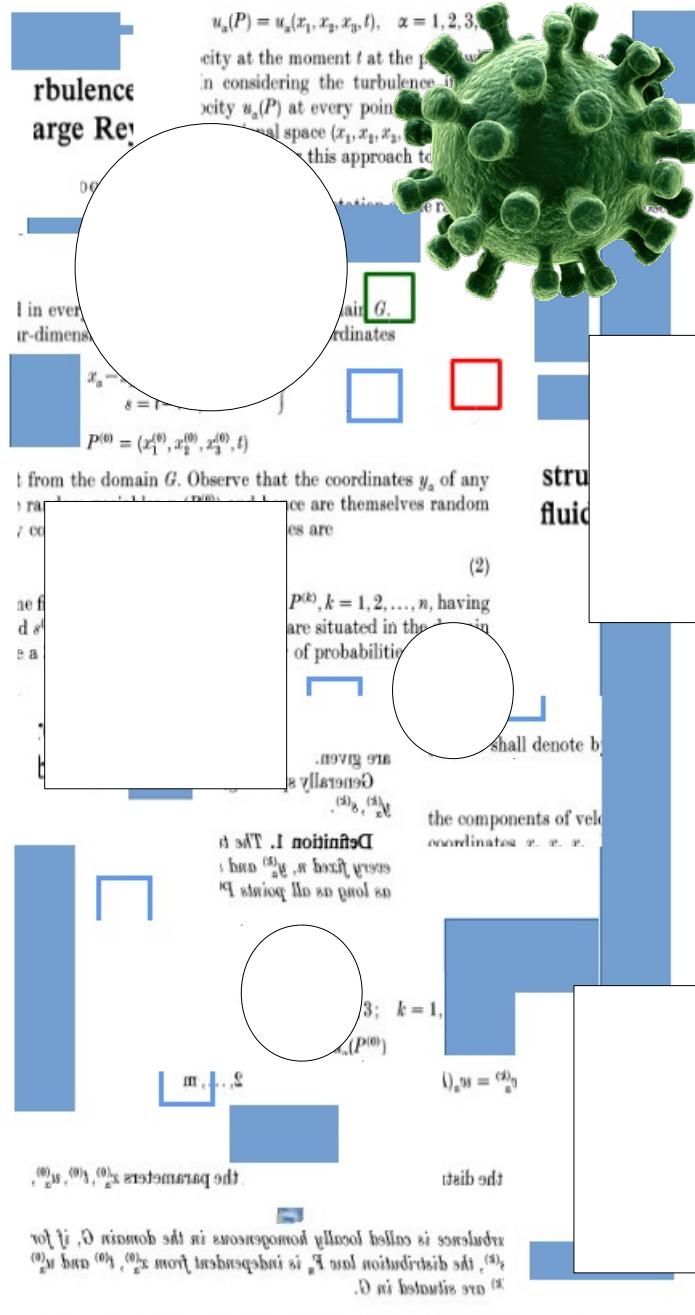
Generally speaking, the distribution law  $F_a$  depends on the parameters  $x_a^{(0)}, t^{(0)}, u_a^{(0)}$ ,  $y_a^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_a^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_a$  is independent from  $x_a^{(0)}, t^{(0)}$  and  $u_a^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

† First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.

Proc. R. Soc. Lond. A (1991) 434, 9–13

Printed in Great Britain



## The local viscous

from the domain  $G$ . Observe that the coordinates  $y_z$  of any point  $P$  of the velocity field are themselves random variables.

are situated in the domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are random variables (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

and bounded

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

where

is a certain fixed point

point  $P$  depend on the variables. The velocity

Suppose that for some fixed values of  $u_a(P^{(0)})$  the coordinates  $y_z^{(0)}$  and  $s^{(0)}$  in the coordinate system (1) are situated in the domain  $G$ . Then we may define quantities

published in Russian translation by V. Levin, re

Proc. R. Soc. Lond. A (1991) 434, 9–13  
Printed in Great Britain

## The local structure of turbulence in incompressible viscous fluid for very large Reynolds numbers†

By A. N. KOLMOGOROV

§1. We shall denote by

$$u_a(P) = u_a(x_1, x_2, x_3, t), \quad a = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_a(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are random variables in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$\bar{u}_a^2 \quad \text{and} \quad \overline{(du_a/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} y_a &= x_a - x_a^{(0)} - u_a(P^{(0)})(t - t^{(0)}), \\ s &= t - t^{(0)}, \end{aligned} \right\} \quad (1)$$

where

$$P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t)$$

is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_a$  of any point  $P$  depend on the random variables  $u_a(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_a(P) = u_a(P) - u_a(P^{(0)}). \quad (2)$$

Suppose that for some fixed values of  $u_a(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_z^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_a$  for the quantities

$$w_a^{(k)} = w_a(P^{(k)}), \quad a = 1, 2, 3; \quad k = 1, 2, \dots, n,$$

where

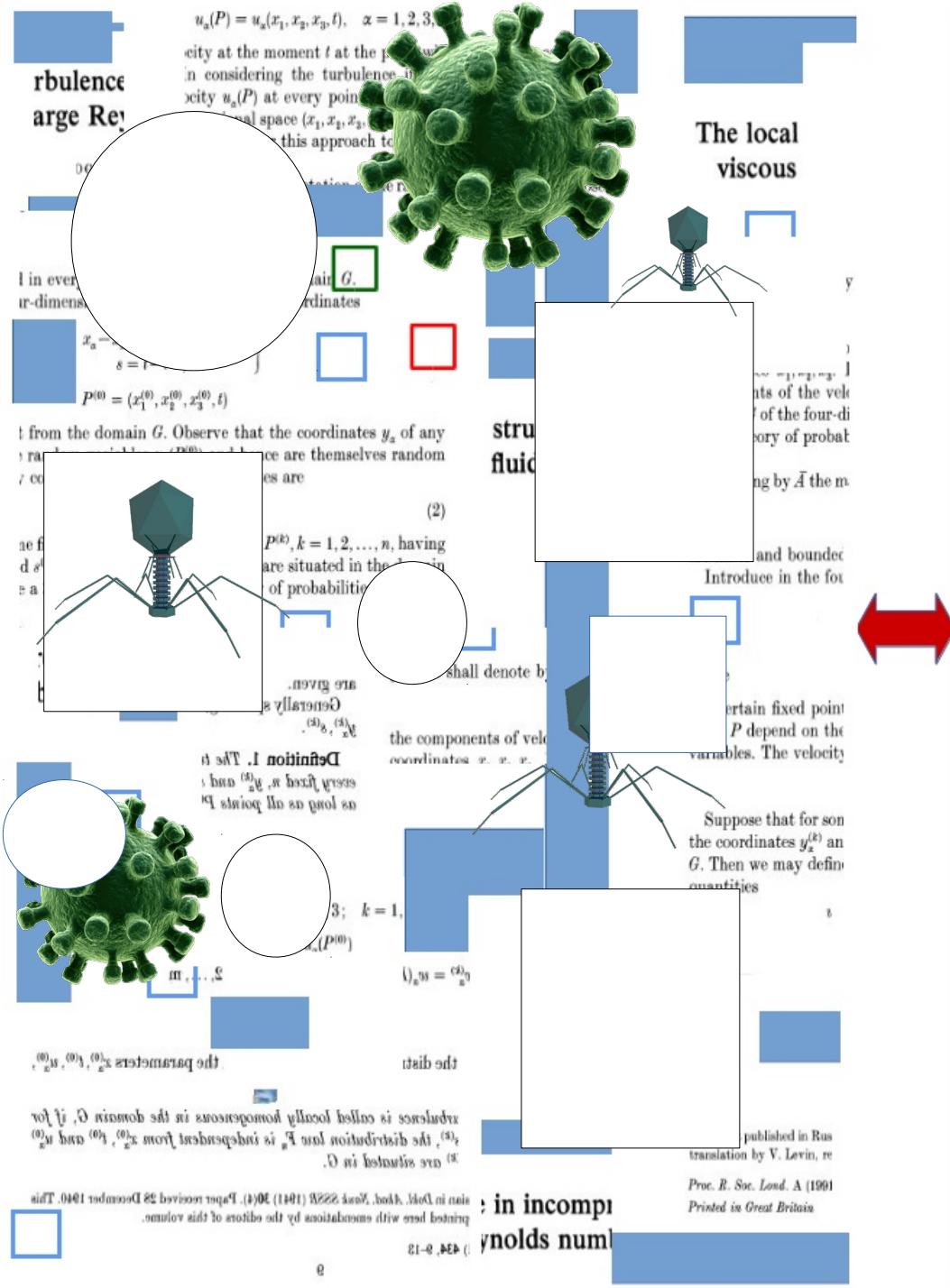
$$u_a^{(k)} =$$

are given.

Generally speaking, the distribution law  $F_a$  depends on the parameters  $x_z^{(0)}, t^{(0)}, u_z^{(0)}$ ,  $y_z^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_z^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_n$  is independent from  $x_z^{(0)}, t^{(0)}$  and  $u_z^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

† First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.



## The local structure of turbulence in incompressible viscous fluid for very large Reynolds numbers†

By A. N. KOLMOGOROV

§1. We shall denote by

$$u_a(P) = u_a(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_a(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are random variables in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$\bar{u}_a^2 \text{ and } \overline{(du_a/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} y_a &= x_a - x_a^{(0)} - u_a(P^{(0)}) (t - t^{(0)}), \\ s &= t - t^{(0)}, \end{aligned} \right\} \quad (1)$$

where

$$P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t)$$

is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_a$  of any point  $P$  depend on the random variables  $u_a(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_a(P) = u_a(P) - u_a(P^{(0)}). \quad (2)$$

Suppose that for some fixed values of  $u_a(P^{(0)})$  the points  $P^{(k)}, k = 1, 2, \dots, n$ , having the coordinates  $y_a^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_n$  for the quantities

$$w_a^{(k)} = w_a(P^{(k)}), \quad \alpha = 1, 2, 3; \quad k = 1, 2, \dots, n,$$

where

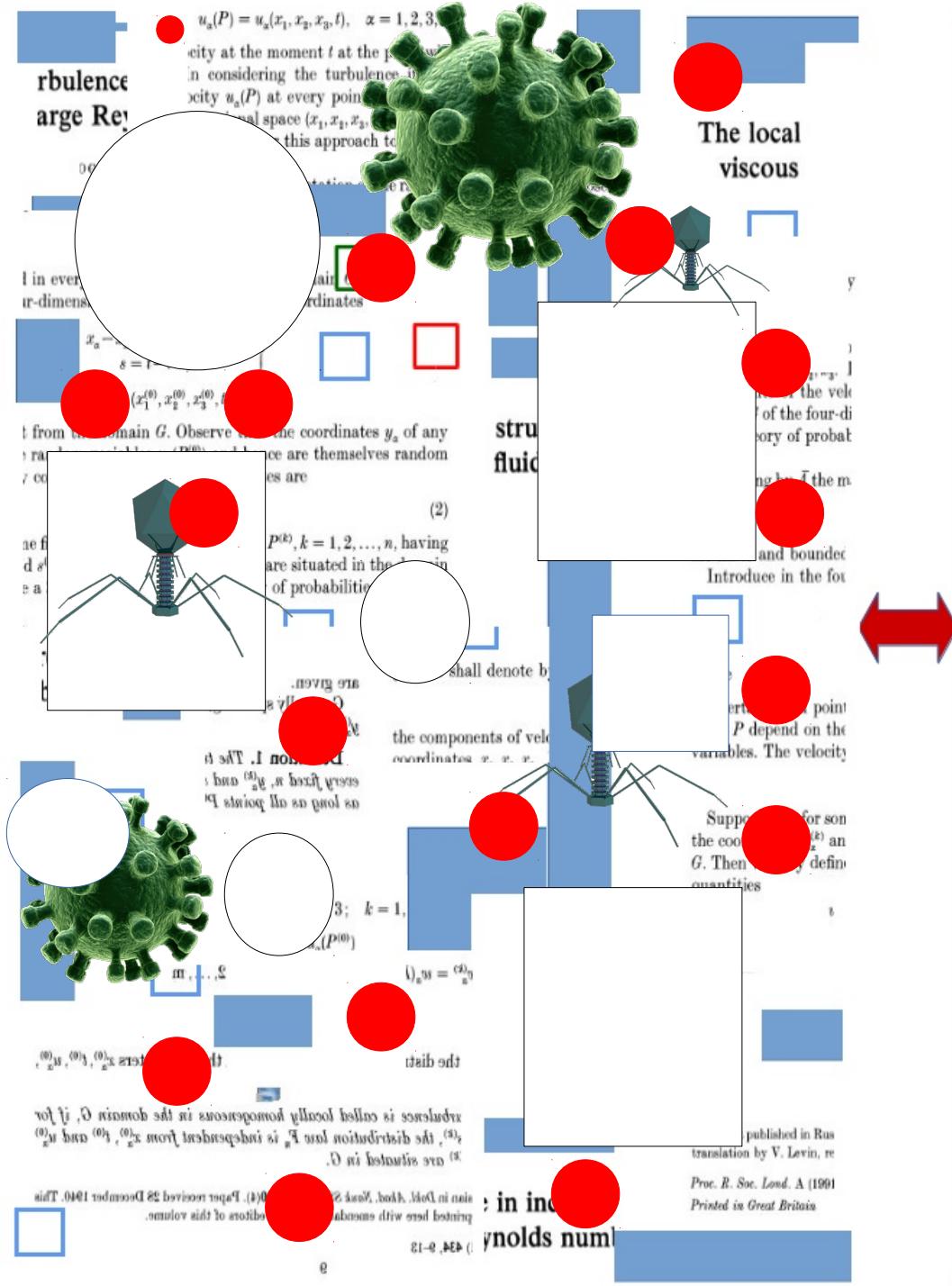
$$u_a^{(k)} =$$

are given.

Generally speaking, the distribution law  $F_n$  depends on the parameters  $x_a^{(0)}, t^{(0)}, u_a^{(0)}$ ,  $y_a^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_a^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_n$  is independent from  $x_a^{(0)}, t^{(0)}$  and  $u_a^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

† First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.



## The local structure of turbulence in incompressible viscous fluid for very large Reynolds numbers†

By A. N. KOLMOGOROV

§1. We shall denote by

$$u_\alpha(P) = u_\alpha(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_\alpha(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are random variables in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$\bar{u}_\alpha^2 \text{ and } \overline{(du_\alpha/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} y_\alpha &= x_\alpha - x_\alpha^{(0)} - u_\alpha(P^{(0)})t - t^{(0)}, \\ s &= t - t^{(0)}, \end{aligned} \right\} \quad (1)$$

where  $P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t^{(0)})$

is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_\alpha$  of any point  $P$  depend on the random variables  $u_\alpha(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_\alpha(P) = u_\alpha(P) - u_\alpha(P^{(0)}). \quad (2)$$

Suppose that for some fixed values of  $u_\alpha(P^{(0)})$  the points  $P^{(k)}, k = 1, 2, \dots, n$ , having the coordinates  $y_\alpha^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_n$  for the quantities

$$w_\alpha^{(k)} = w_\alpha(P^{(k)}), \quad \alpha = 1, 2, 3; \quad k = 1, 2, \dots, n,$$

where  $u_\alpha^{(k)} = u_\alpha(P^{(k)})$

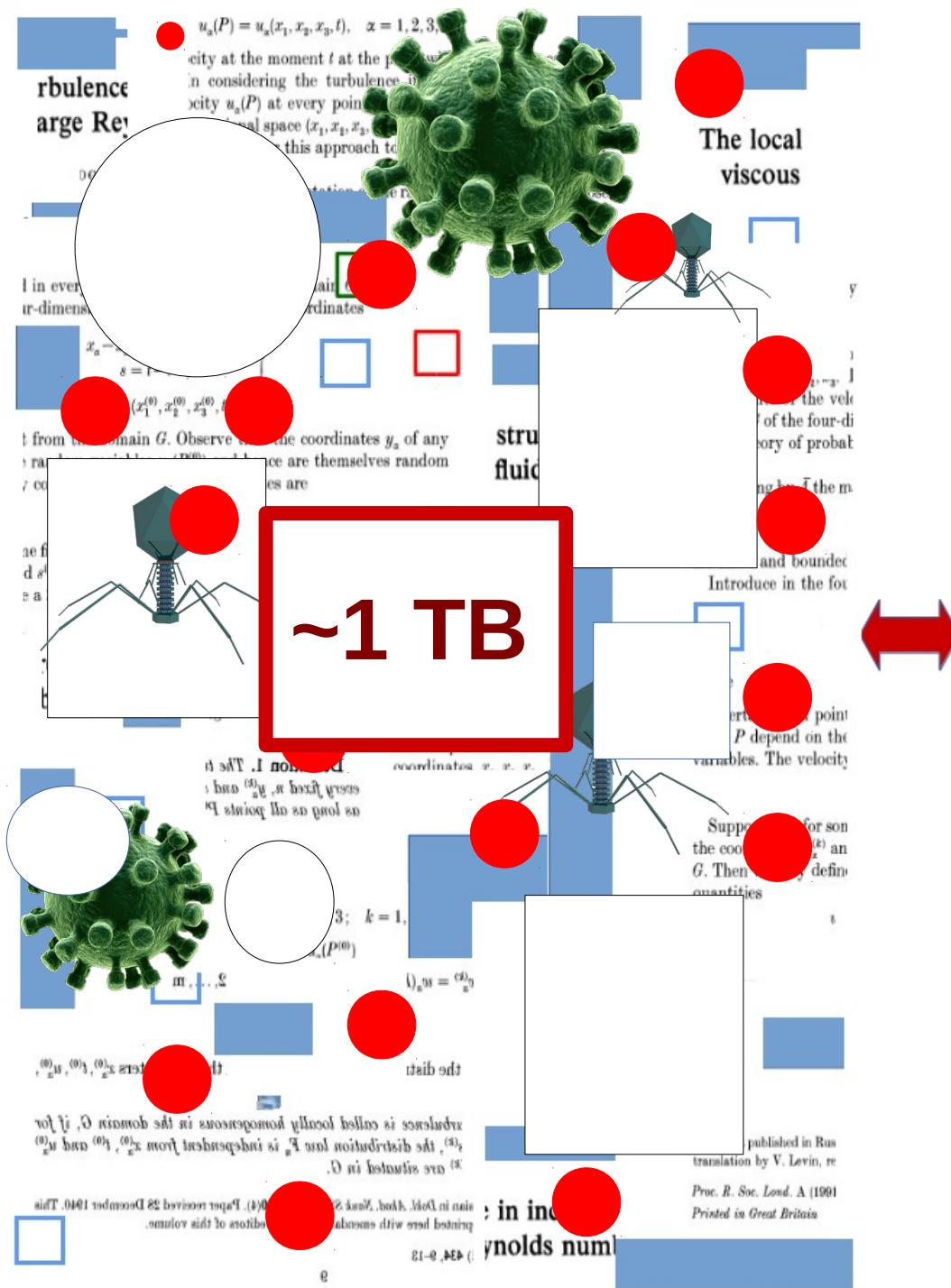
$$y_\alpha^{(k)}, s^{(k)}$$

are given.

Generally speaking, the distribution law  $F_n$  depends on the parameters  $x_\alpha^{(0)}, t^{(0)}, u_\alpha^{(0)}, y_\alpha^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_\alpha^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_n$  is independent from  $x_\alpha^{(0)}, t^{(0)}$  and  $u_\alpha^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

† First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.



## The local structure of turbulence in incompressible viscous fluid for very large Reynolds numbers†

By A. N. KOLMOGOROV

§1. We shall denote by

$$u_\alpha(P) = u_\alpha(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_\alpha(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are random variables in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$u_\alpha^2 \text{ and } \overline{(du_\alpha/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} y_\alpha &= x_\alpha - x_\alpha^{(0)} - u_\alpha(P^{(0)})(t - t^{(0)}), \\ s &= t - t^{(0)}, \end{aligned} \right\} \quad (1)$$

where

$$P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t)$$

is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_\alpha$  of any point  $P$  depend on the random variables  $u_\alpha(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_\alpha(P) = u_\alpha(P) - u_\alpha(P^{(0)}). \quad (2)$$

Suppose that for some fixed values of  $u_\alpha(P^{(0)})$  the points  $P^{(k)}$ ,  $k = 1, 2, \dots, n$ , having the coordinates  $y_\alpha^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_n$  for the quantities

$$w_\alpha^{(k)} = w_\alpha(P^{(k)}), \quad \alpha = 1, 2, 3; \quad k = 1, 2, \dots, n,$$

where

$$u_\alpha^{(k)} =$$

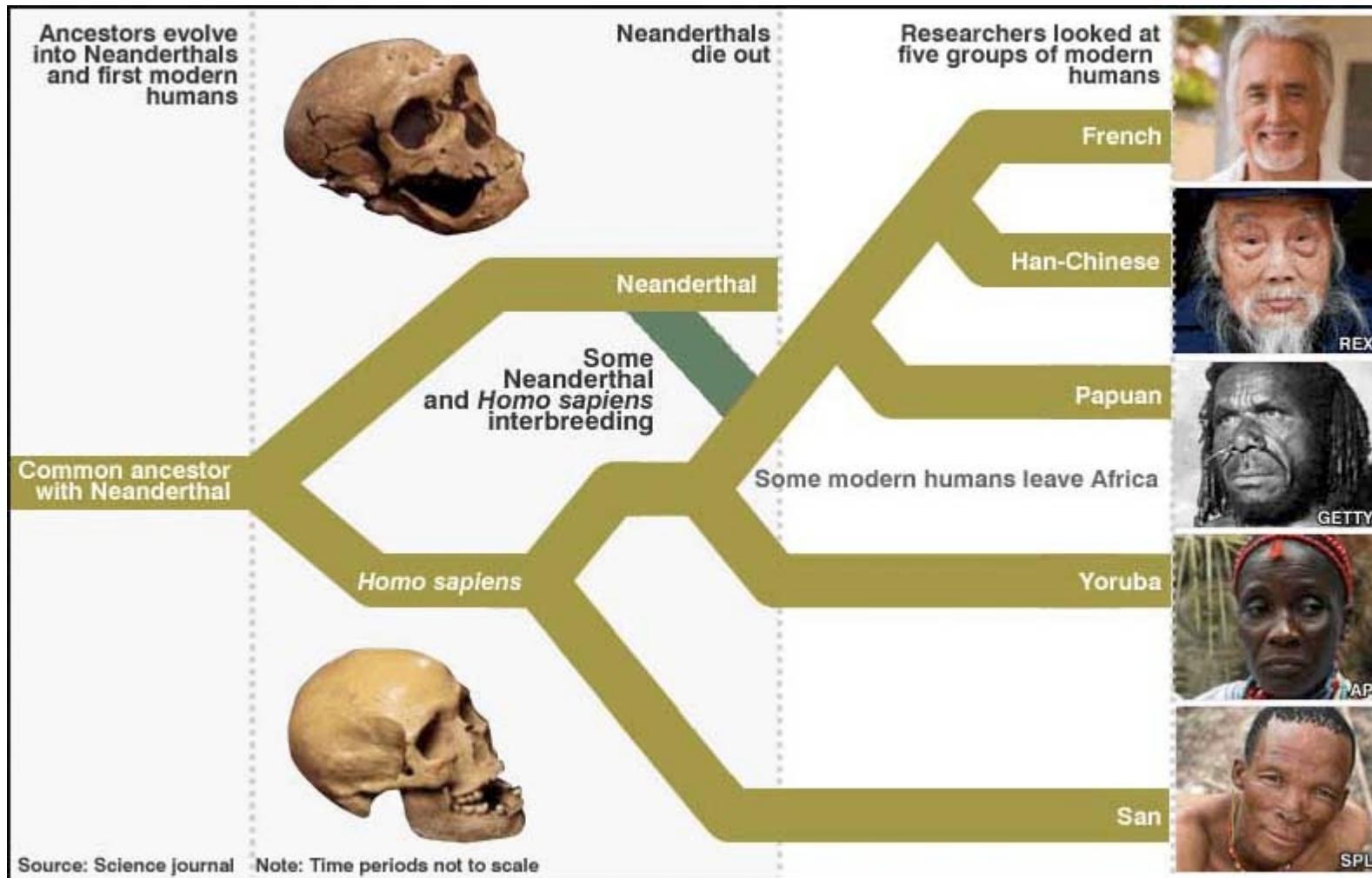
are given.

Generally speaking, the distribution law  $F_n$  depends on the parameters  $x_\alpha^{(0)}, t^{(0)}, u_\alpha^{(0)}, y_\alpha^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_\alpha^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_n$  is independent from  $x_\alpha^{(0)}, t^{(0)}$  and  $u_\alpha^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

† First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.

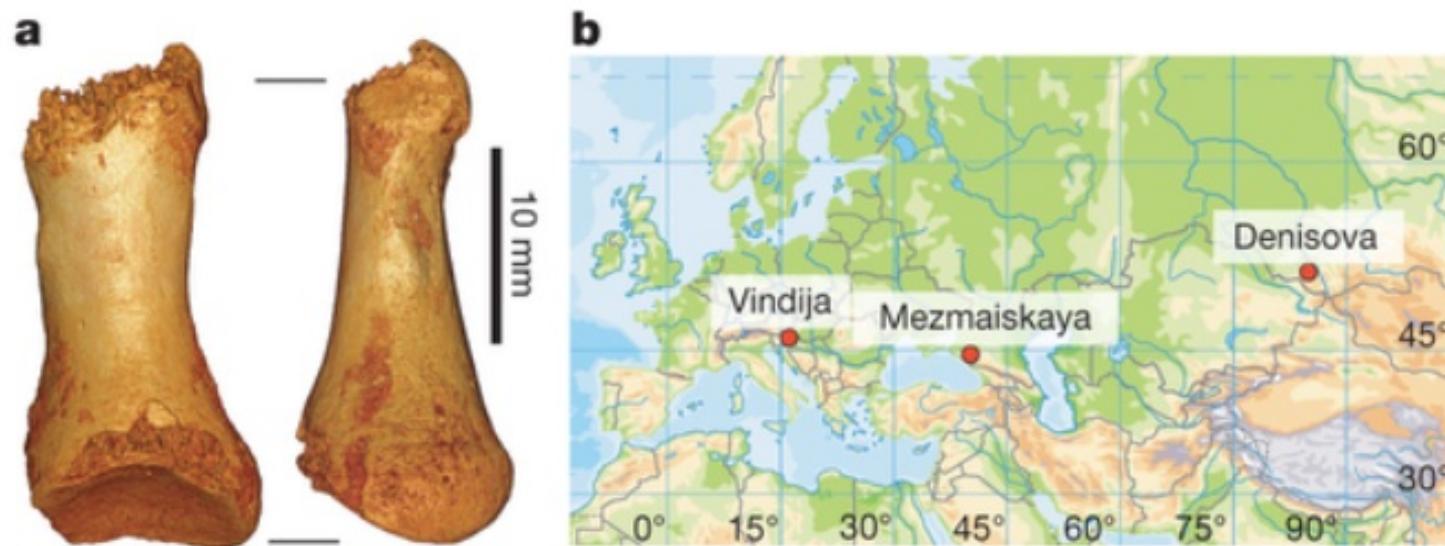
# Timeline



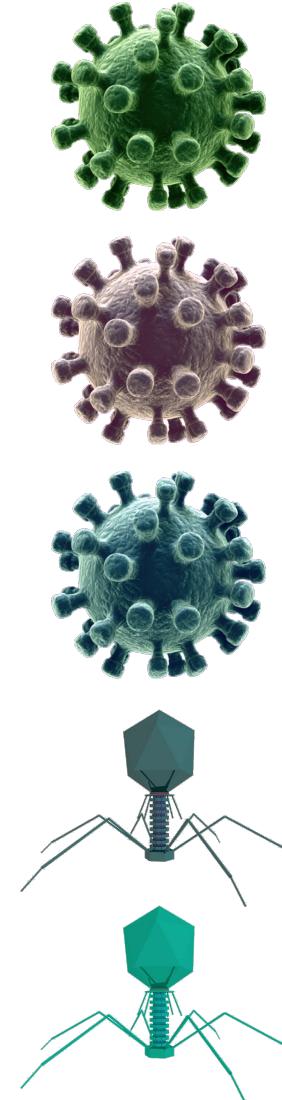
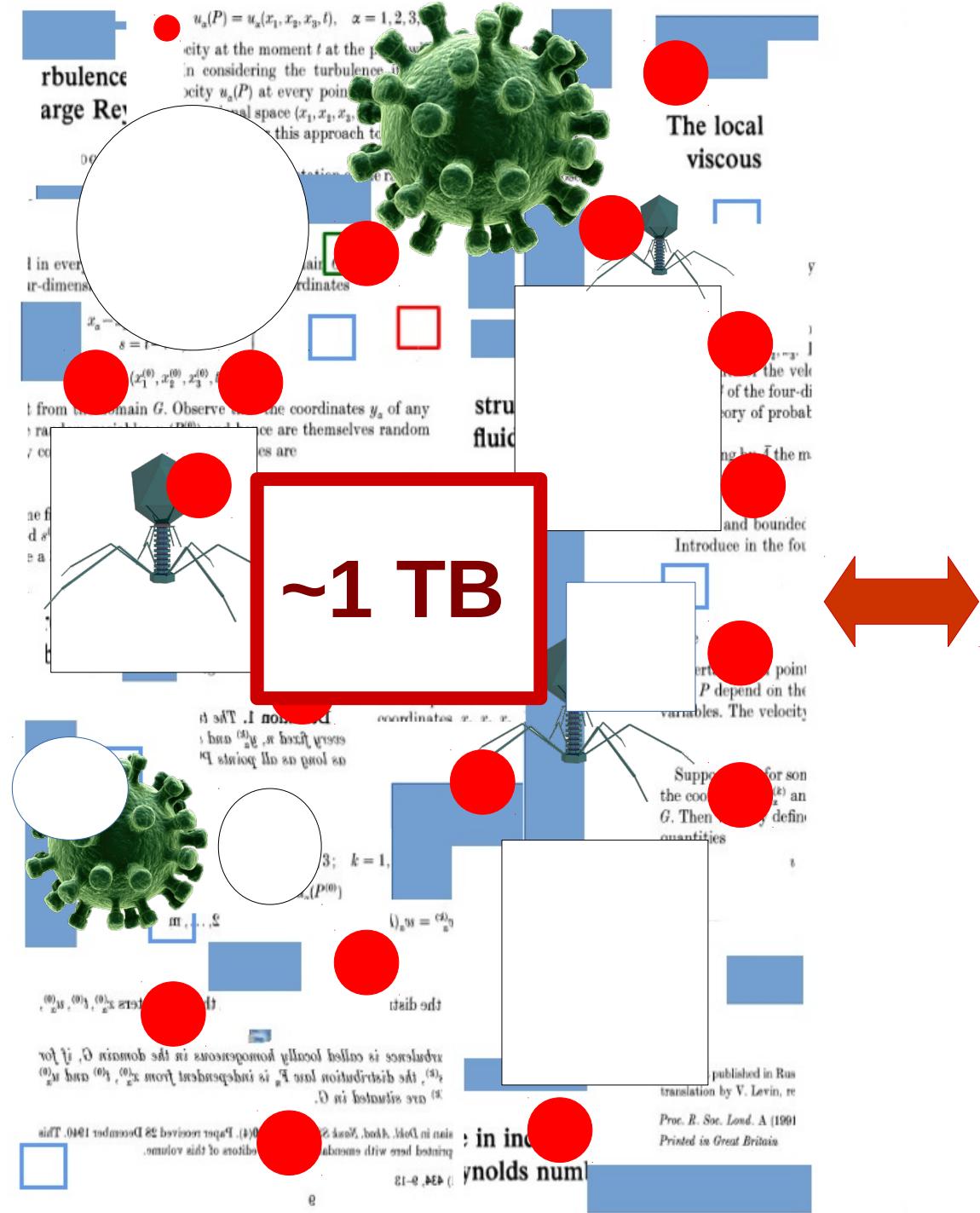
Neanderthals approximate period: ~400,000 to ~40,000 years ago; Unknown reason of extinction.

# Neanderthal (~50,000 years ago)

Figure 1: Toe phalanx and location of Neanderthal samples for which genome-wide data are available.

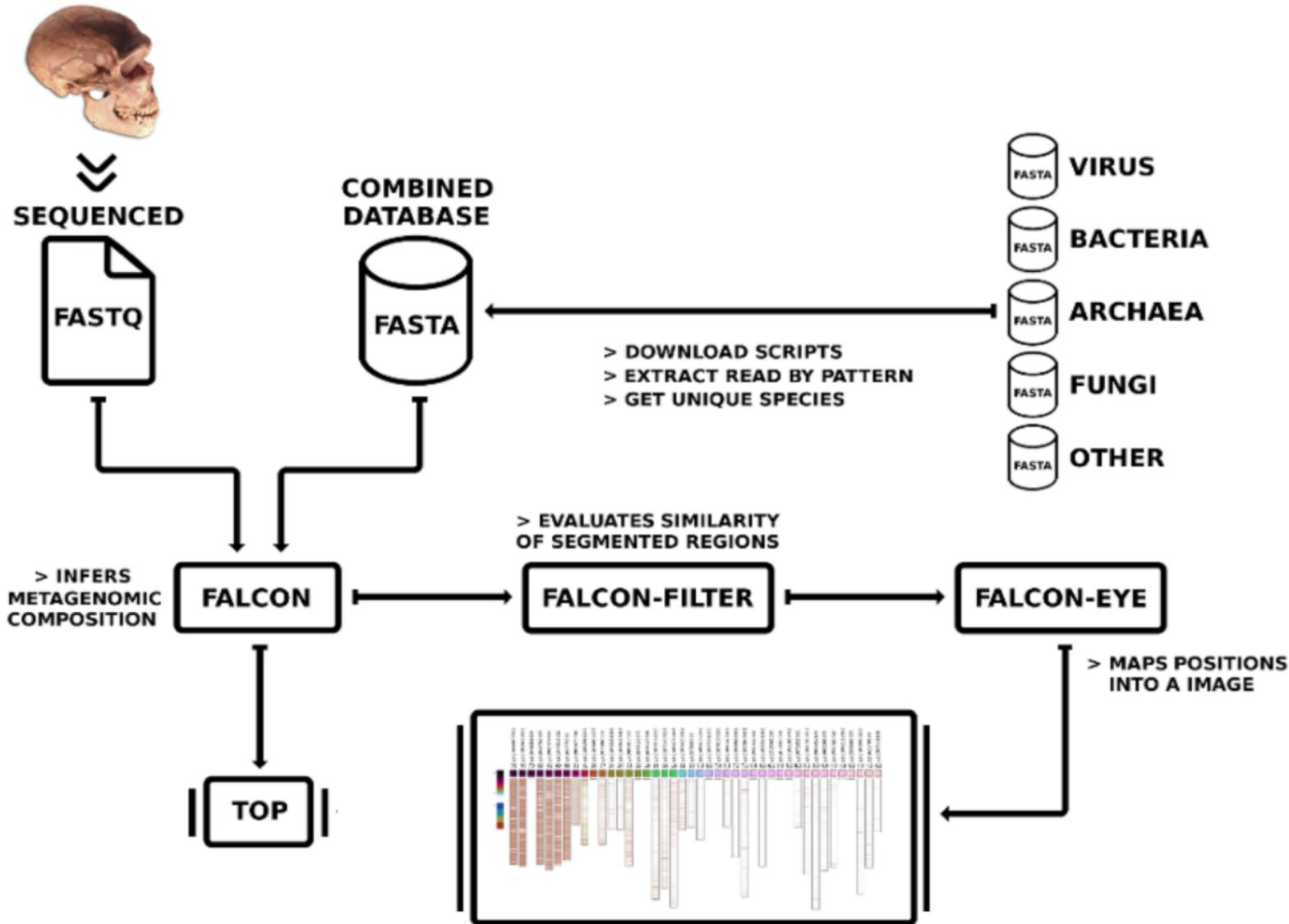


**a**, The toe phalanx found in the east gallery of Denisova Cave in 2010. Dorsal view (left image), left view (right image). Total length of the bone is 26mm. **b**, Map of Eurasia showing the location of Vindija Cave, Mezmaiskaya Cave and D..

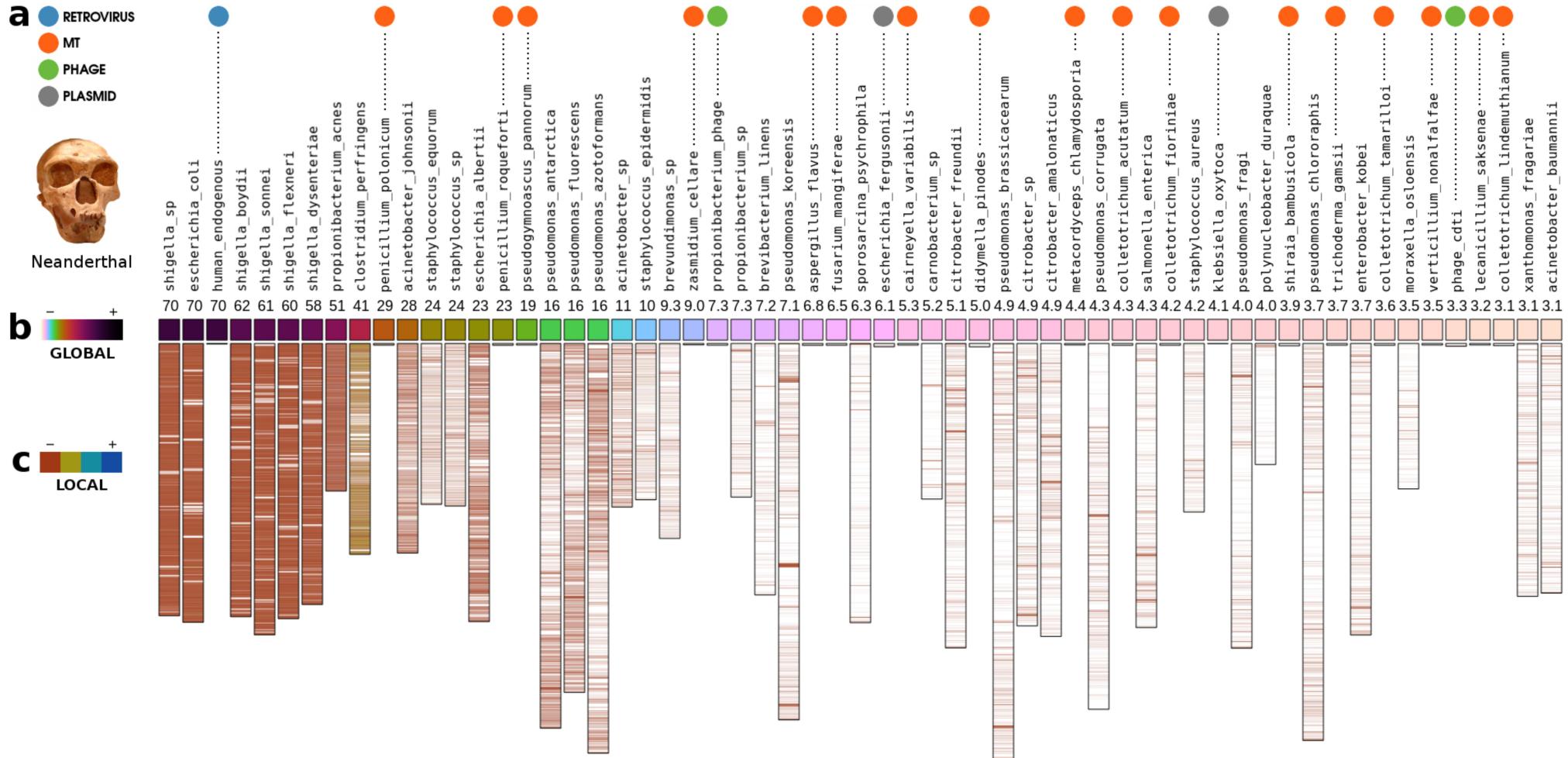


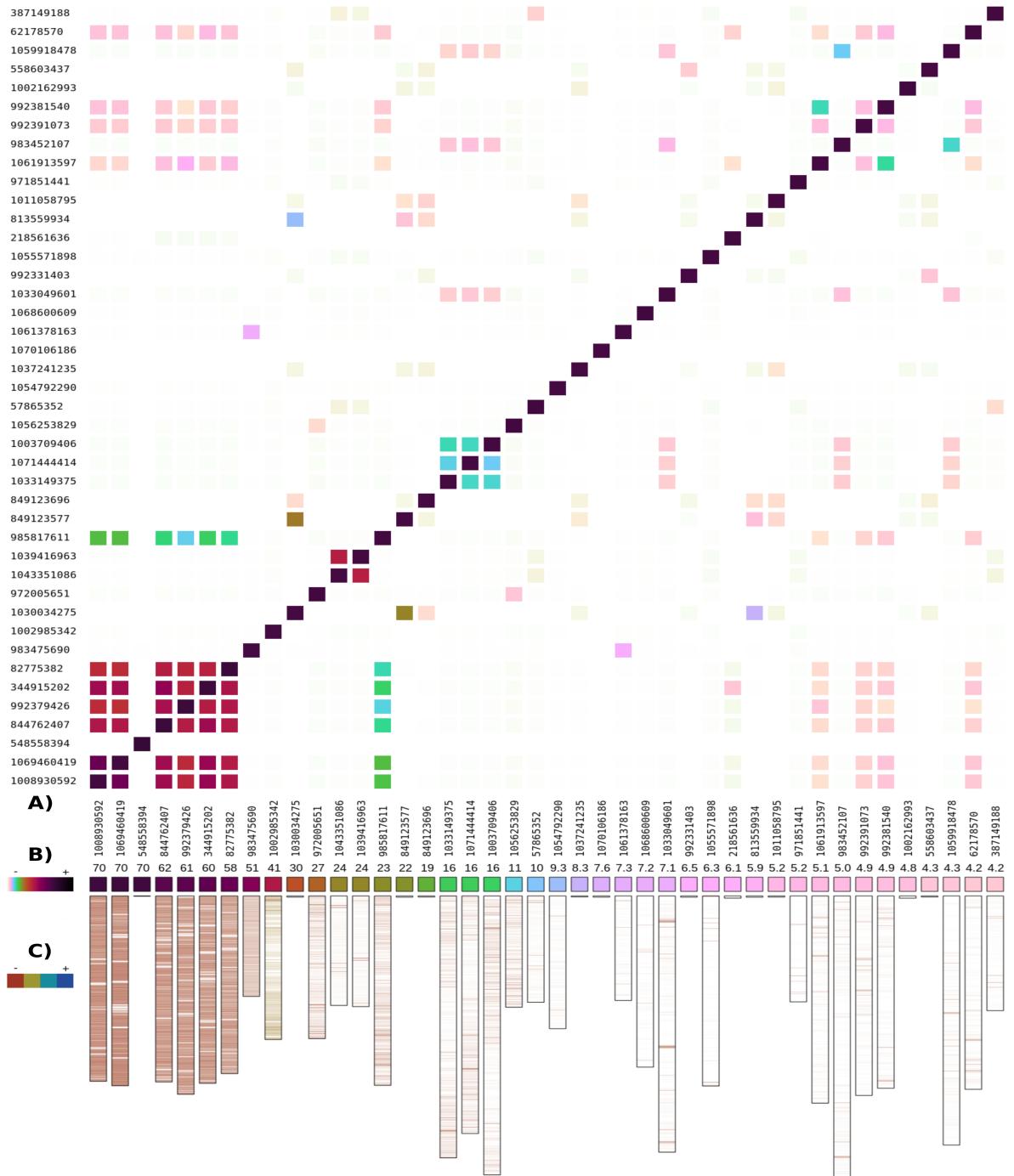
# All known microorganisms

# Metagenomic sample composition pipeline

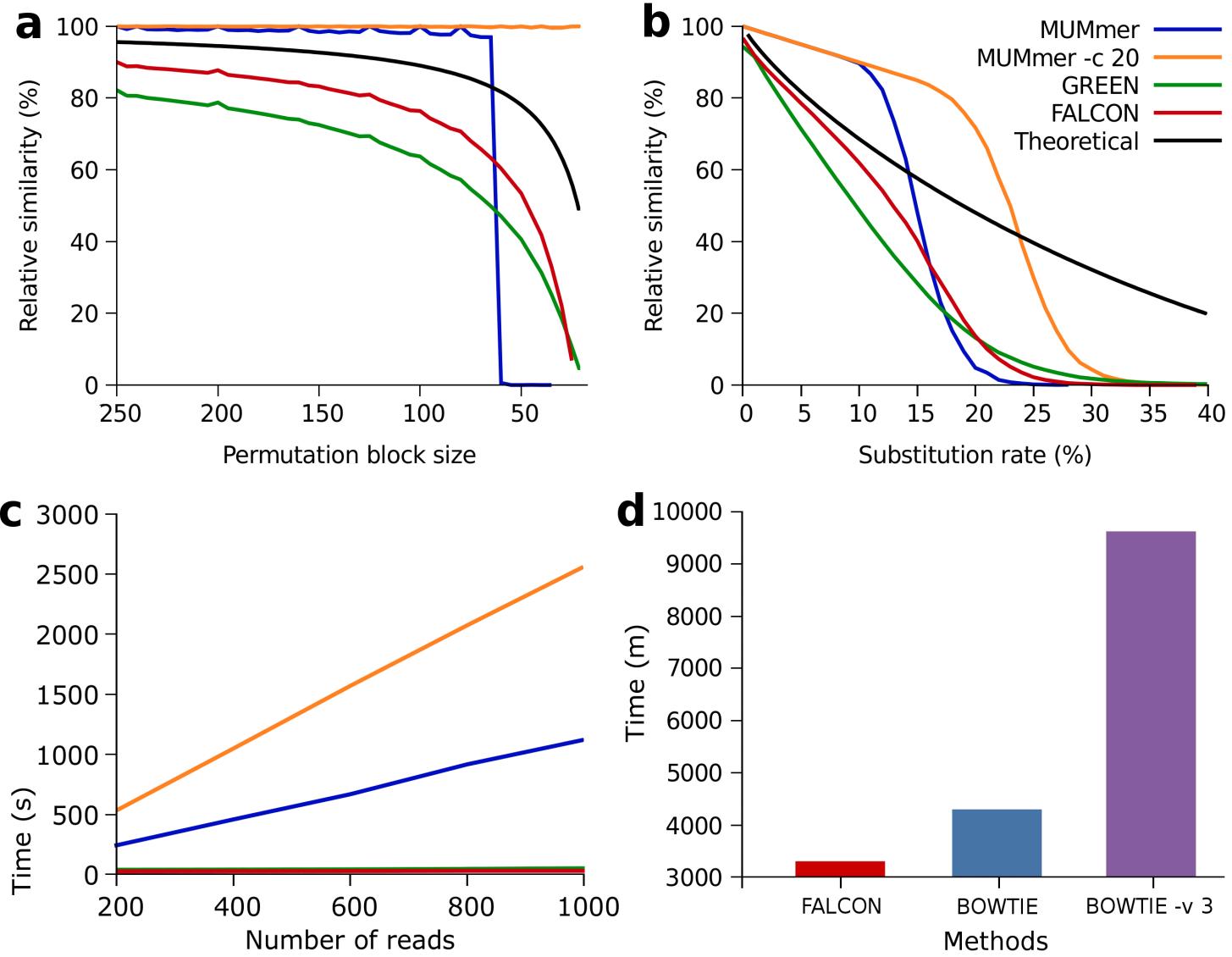


# Neanderthal metagenomic sample composition

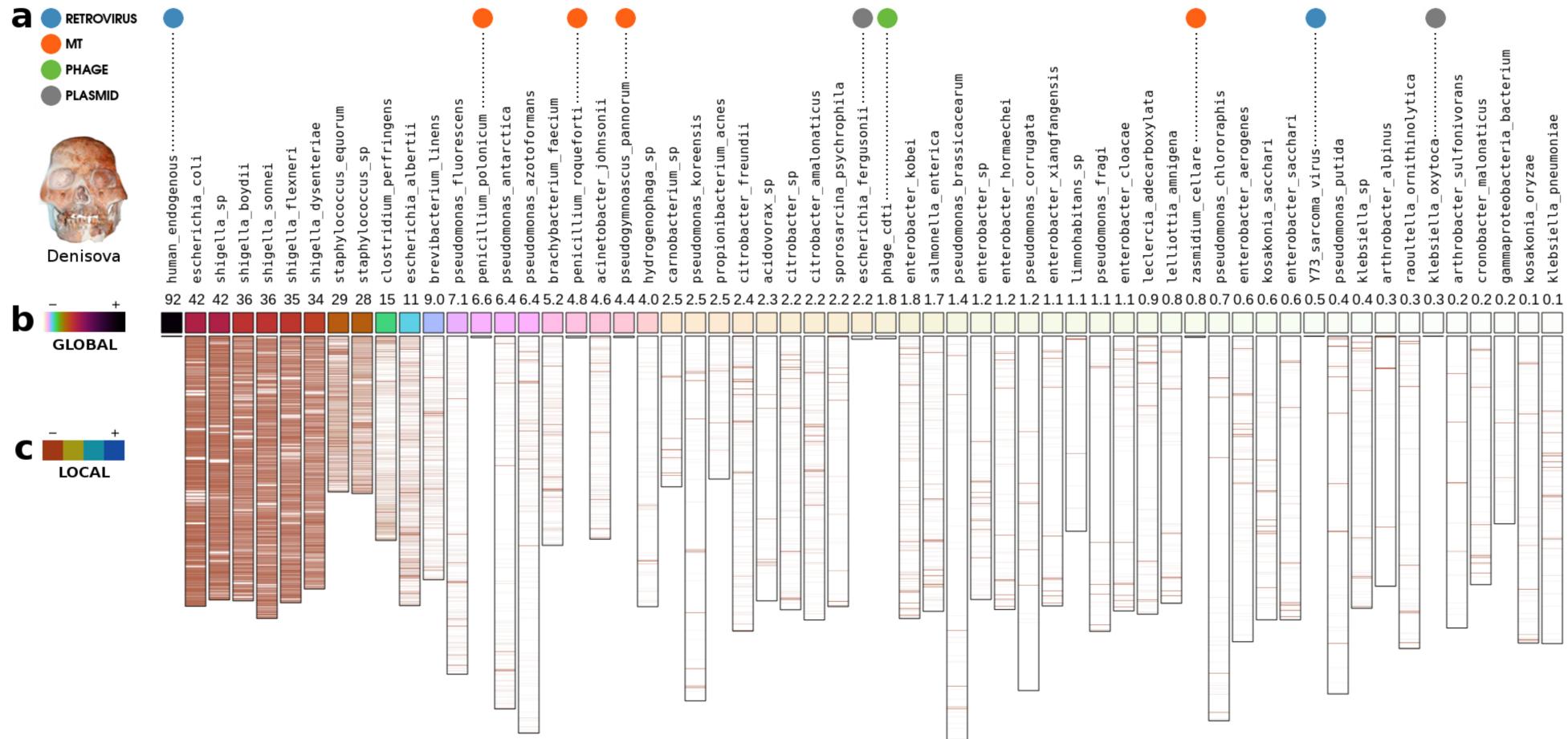




D. Pratas, R. Silva, A. J. Pinho, P. J. S. G. Ferreira. Unsupervised measurement of relative similarity between non-assembled genomic sequences. (submitted).

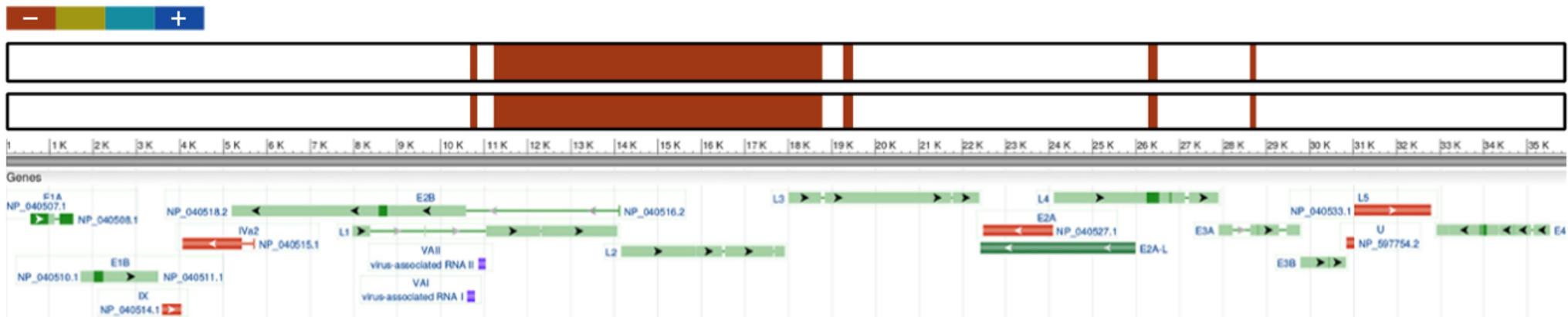


# Denisova metagenomic sample composition



# Denisova metagenomic sample composition

Local similarities of Human Adenovirus with an ancient virus:

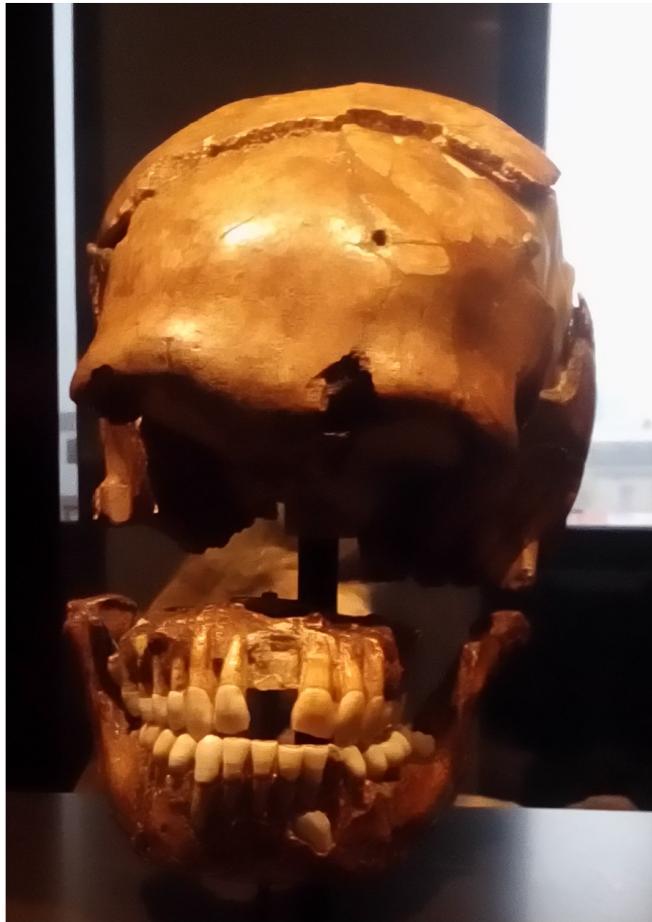


Associated with respiratory infections

**Ancient genomes hide a key for efficient therapeutics**

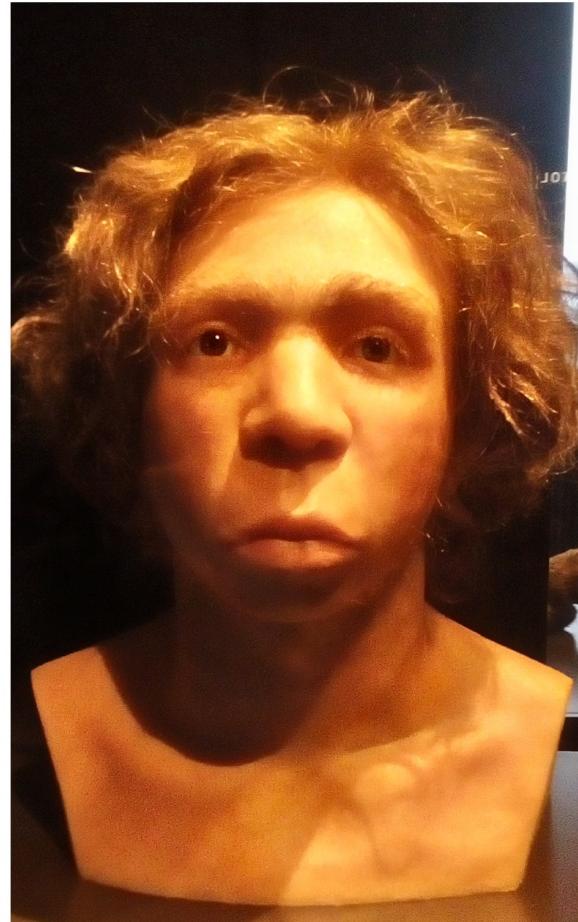
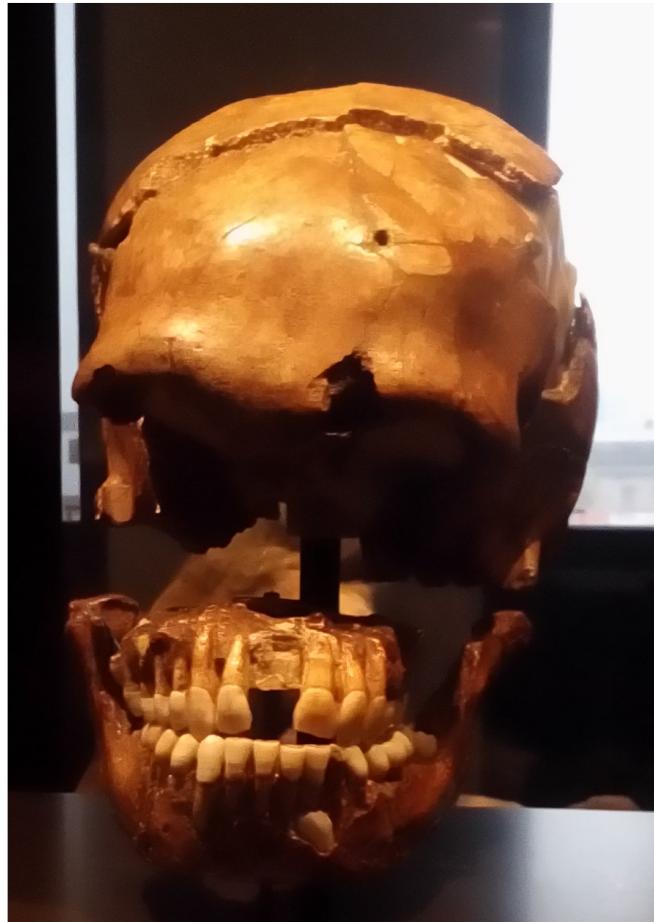
**What makes us human?**

# Differences between Humans and Neanderthals

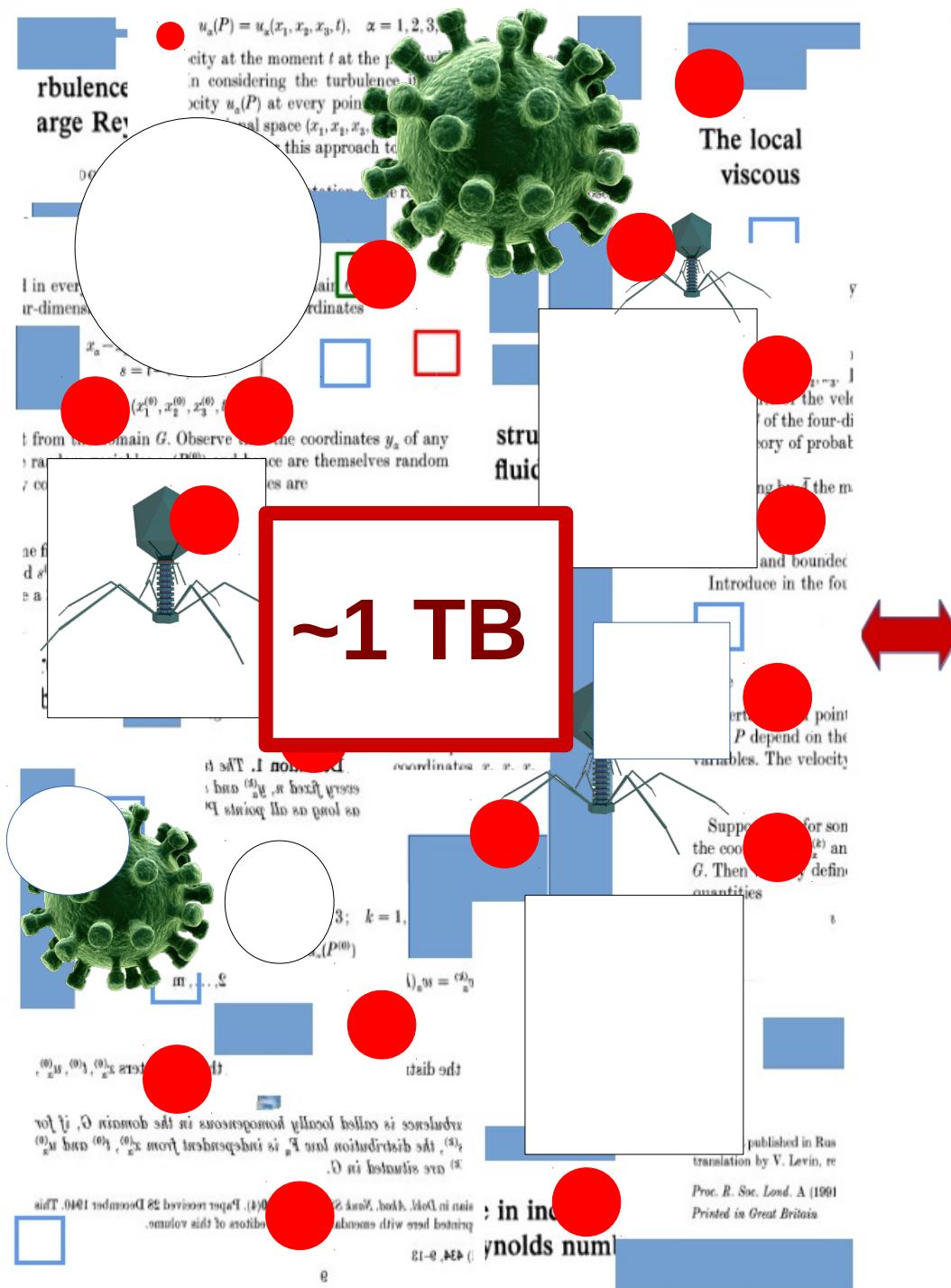


45,000 years old Neanderthal in the Neues museum at Berlin. Photo and edition by Pratas *et al*, May, 2017.

# Differences between Humans and Neanderthals



45,000 years old Neanderthal in the Neues museum at Berlin. Photo and edition by Pratas et al, May, 2017.



## The local structure of turbulence in incompressible viscous fluid for very large Reynolds numbers†

By A. N. KOLMOGOROV

§1. We shall denote by

$$u_\alpha(P) = u_\alpha(x_1, x_2, x_3, t), \quad \alpha = 1, 2, 3,$$

the components of velocity at the moment  $t$  at the point with rectangular cartesian coordinates  $x_1, x_2, x_3$ . In considering the turbulence it is natural to assume the components of the velocity  $u_\alpha(P)$  at every point  $P = (x_1, x_2, x_3, t)$  of the considered domain  $G$  of the four-dimensional space  $(x_1, x_2, x_3, t)$  are random variables in the sense of the theory of probabilities (cf. for this approach to the problem Millionshtchikov (1939)).

Denoting by  $\bar{A}$  the mathematical expectation of the random variable  $A$  we suppose that

$$\overline{u_\alpha^2} \quad \text{and} \quad \overline{(du_\alpha/dx_\beta)^2}$$

are finite and bounded in every bounded subdomain of the domain

Introduce in the four-dimensional space  $(x_1, x_2, x_3, t)$  new coordinates

$$\left. \begin{aligned} y_\alpha &= x_\alpha - x_\alpha^{(0)} - u_\alpha(P^{(0)})t - t^{(0)}, \\ s &= t - t^{(0)}, \end{aligned} \right\} \quad (1)$$

where

$$P^{(0)} = (x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, t^{(0)})$$

is a certain fixed point from the domain  $G$ . Observe that the coordinates  $y_\alpha$  of any point  $P$  depend on the random variables  $u_\alpha(P^{(0)})$  and hence are themselves random variables. The velocity components in the new coordinates are

$$w_\alpha(P) = u_\alpha(P) - u_\alpha(P^{(0)}). \quad (2)$$

Suppose that for some fixed values of  $u_\alpha(P^{(0)})$  the points  $P^{(k)}, k = 1, 2, \dots, n$ , having the coordinates  $y_\alpha^{(k)}$  and  $s^{(k)}$  in the coordinate system (1), are situated in the domain  $G$ . Then we may define a  $3n$ -dimensional distribution law of probabilities  $F_n$  for the quantities

$$w_\alpha^{(k)} = w_\alpha(P^{(k)}), \quad \alpha = 1, 2, 3; \quad k = 1, 2, \dots, n,$$

where

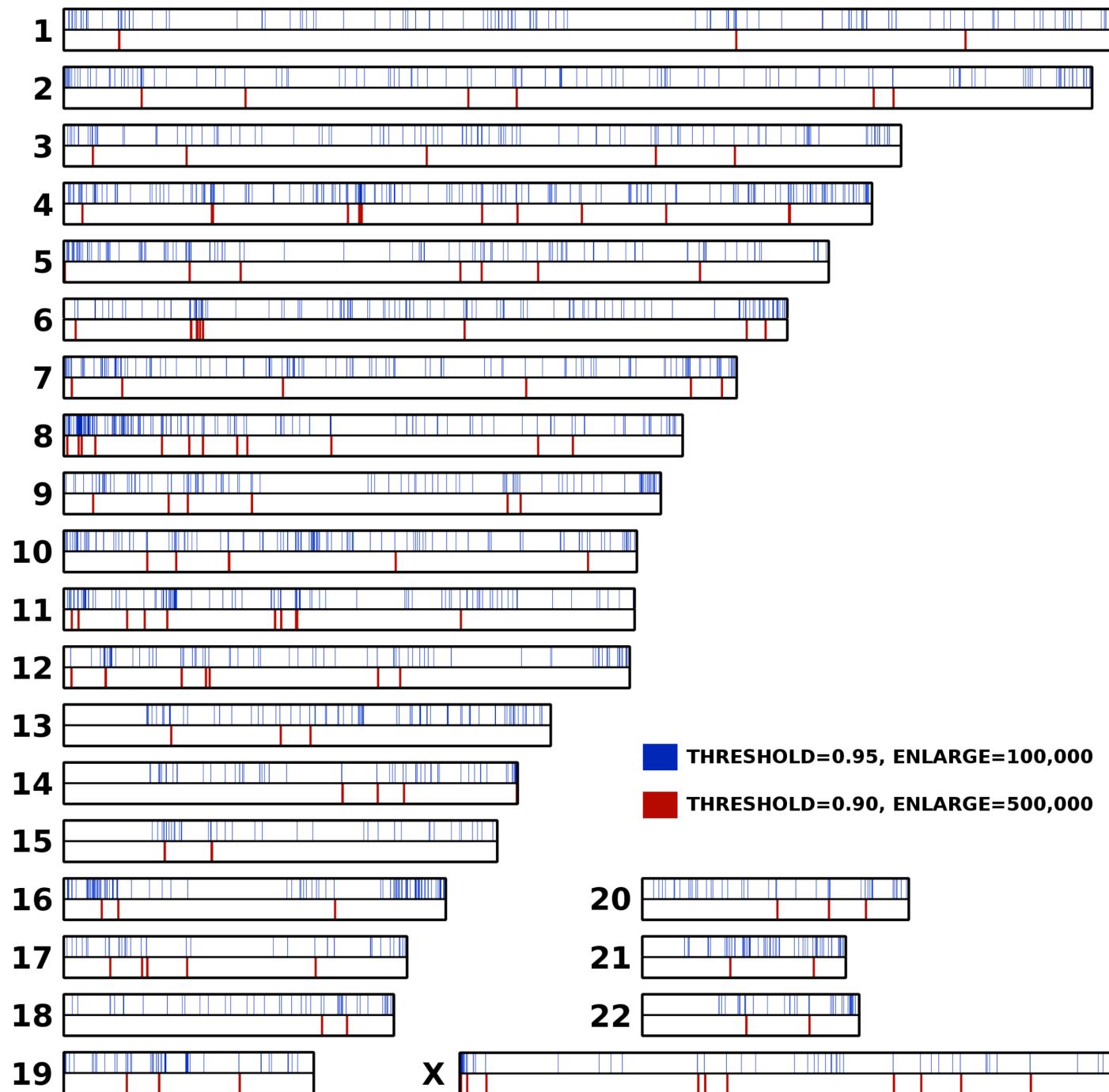
$$u_\alpha^{(k)} =$$

are given.

Generally speaking, the distribution law  $F_n$  depends on the parameters  $x_\alpha^{(0)}, t^{(0)}, u_\alpha^{(0)}, y_\alpha^{(k)}, s^{(k)}$ .

**Definition 1.** The turbulence is called locally homogeneous in the domain  $G$ , if for every fixed  $n$ ,  $y_\alpha^{(k)}$  and  $s^{(k)}$ , the distribution law  $F_n$  is independent from  $x_\alpha^{(0)}, t^{(0)}$  and  $u_\alpha^{(0)}$  as long as all points  $P^{(k)}$  are situated in  $G$ .

† First published in Russian in *Dokl. Akad. Nauk SSSR* (1941) 30(4). Paper received 28 December 1940. This translation by V. Levin, reprinted here with emendations by the editors of this volume.



# Modern-human genes associated with:



brain (**neurotransmitters** and **synapses**)

# Modern-human genes associated with:



brain (**neurotransmitters** and **synapses**)



hearing

# Modern-human genes associated with:



brain (**neurotransmitters** and **synapses**)



hearing



blood

# Modern-human genes associated with:



**brain (neurotransmitters and synapses)**



**hearing**



**blood**



**fertility**

# Modern-human genes associated with:



**brain (neurotransmitters and synapses)**



**hearing**



**blood**



**fertility**



**immune system**

# Modern-human genes associated with:



**brain (neurotransmitters and synapses)**



**hearing**



**blood**



**fertility**



**immune system**



**and others with unknown function...**

## Next challenges:

- Visualization of distinct DNA regions of the modern human relatively to a **Denisova** genome
- Classification of **ancient** or **contamination** microorganisms present in a **Neanderthal** and **Denisova** genomes
- Reconstruction of ancient microorganisms
- Study the co-evolution between host and pathogen

**Thanks for the attention!**

**Work supported by FCT funds**

**[pratas@ua.pt](mailto:pratas@ua.pt)**

