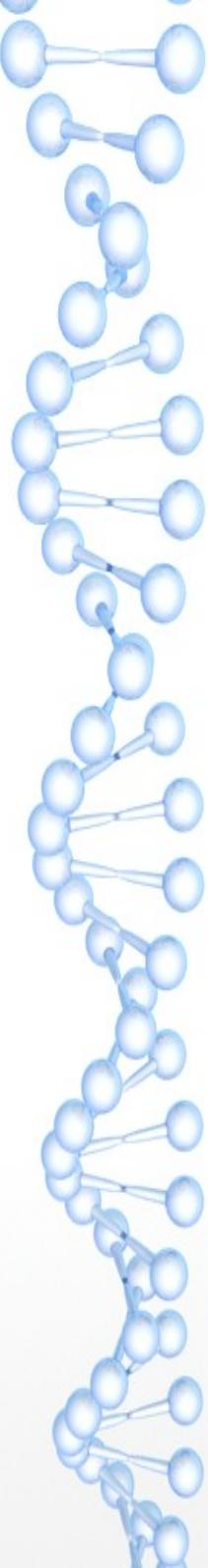


I. Structure and diversity of the genomes



I. Structure and diversity of the genomes

Definition of the term “genome”

Background question of the module

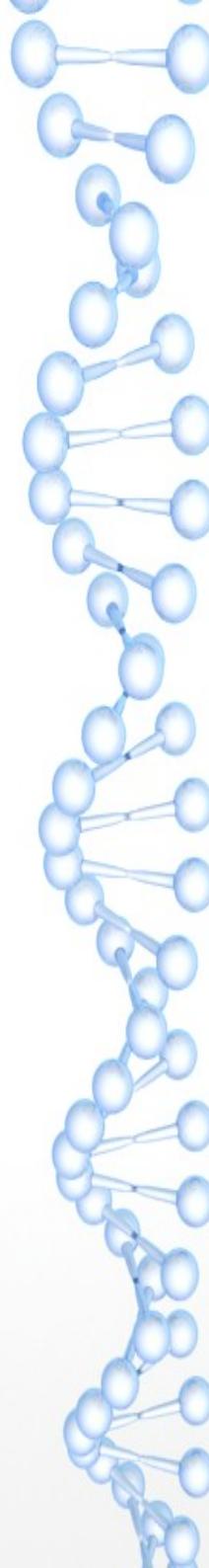
What is a genome?

Several commonly used definitions

Genome is

- the entire set of **DNA** instructions found in a **cell**.
- genetic material of haploid set of **chromosomes**
- **information repository** of organism
- all **hereditary material** of organism

Definitions overlap significantly, but not completely!



I. Structure and diversity of the genomes

Genome size

Variation of genome size

How small is the smallest genome on the Earth?

How big is the biggest one?

Your suggestions?

Variation of genome size

How small is the smallest genome on the Earth?

How big is the biggest one?

Your suggestions?

Hints:

human	3.3 Gbp
drosophila	240 Mbp
baker yeasts	12 Mbp
<i>E. coli</i>	4.2 Mbp

Gbp = giga base pairs
Mbp = mega base pairs

Biggest genome

Paris japonica

キヌガサソウ, *Kinugasasō*

150 Gbp / 152.23 pg

(Pellicer et al, 2010)



Protopterus aethiopicus

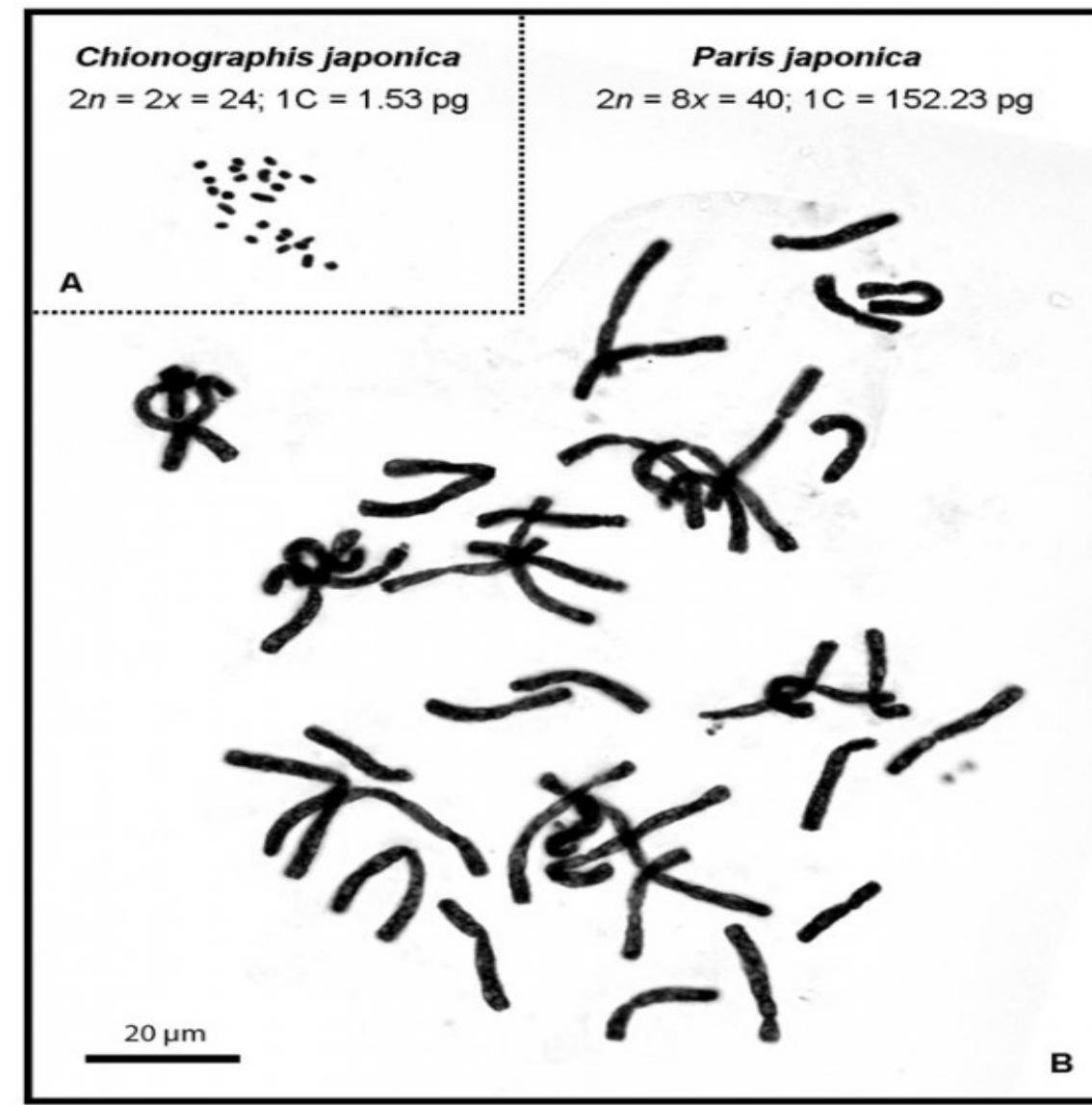
Marbled lungfish

130 Gbp / 132.83 pg

(Pedersen, 1971)



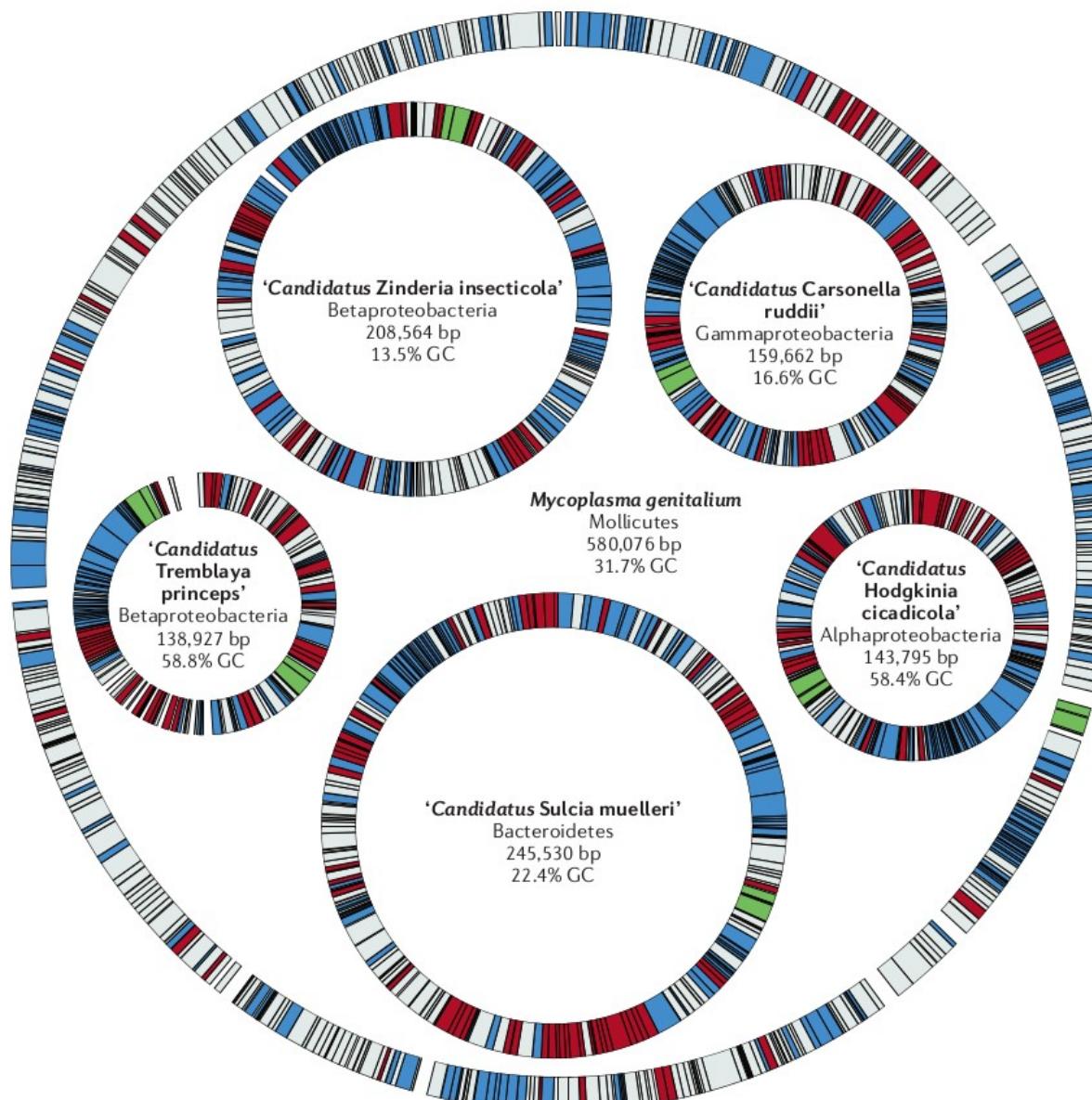
100x genome size difference in the cell



150 Gbp vs 1.5 Gbp

Pellicer et al, 2010

Small bacterial genomes

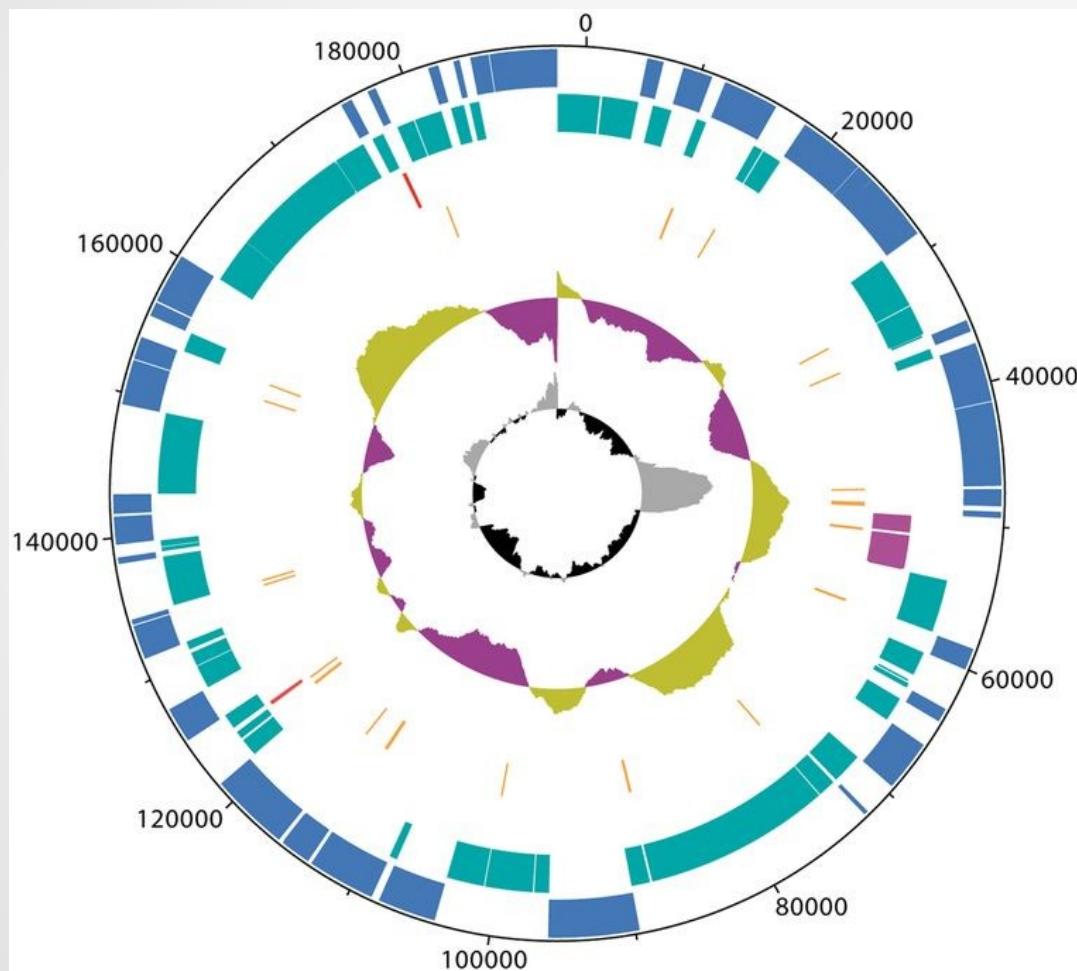


Free-living *M. genitalium*
vs
endosymbionts

580 kbp
vs
138 - 245 kbp

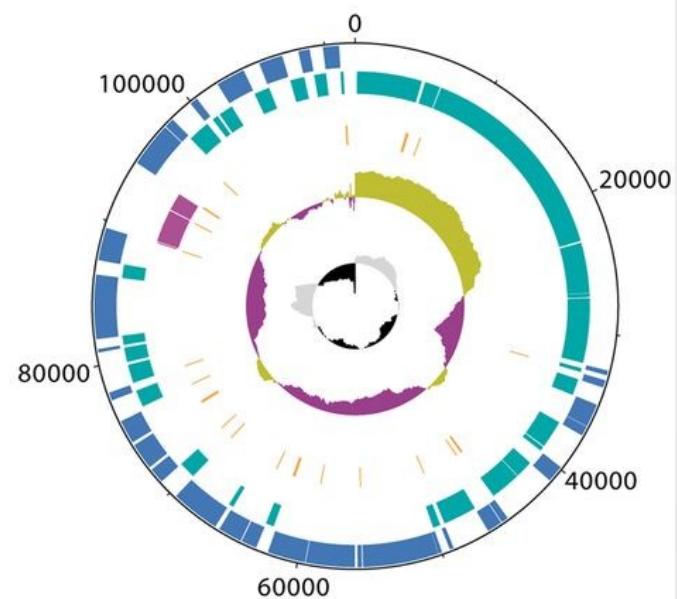
Smallest bacterial genomes

190 kbp



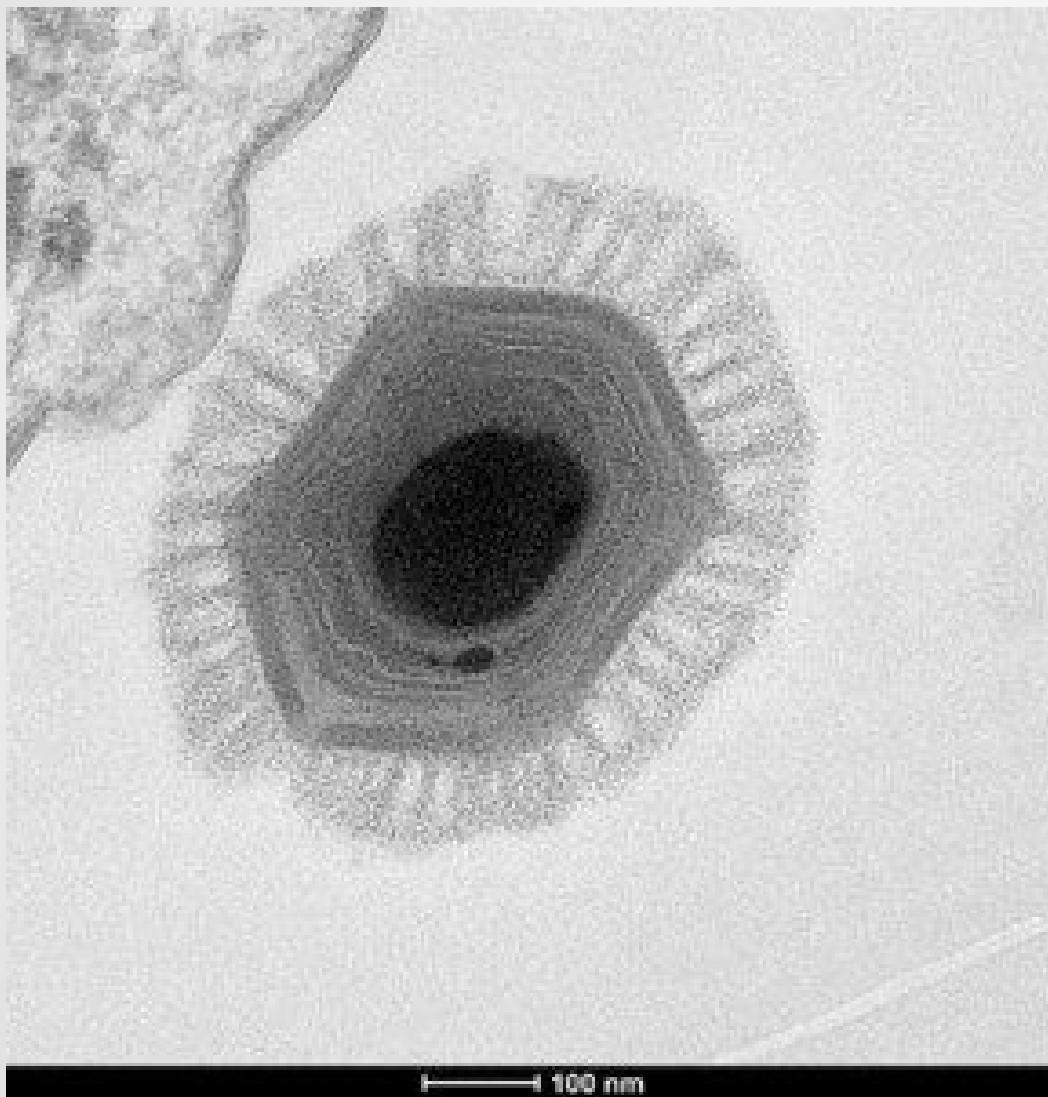
Sulcia muelleri

112 kbp



Nasuia deltocephalinicola

Biggest viral genome

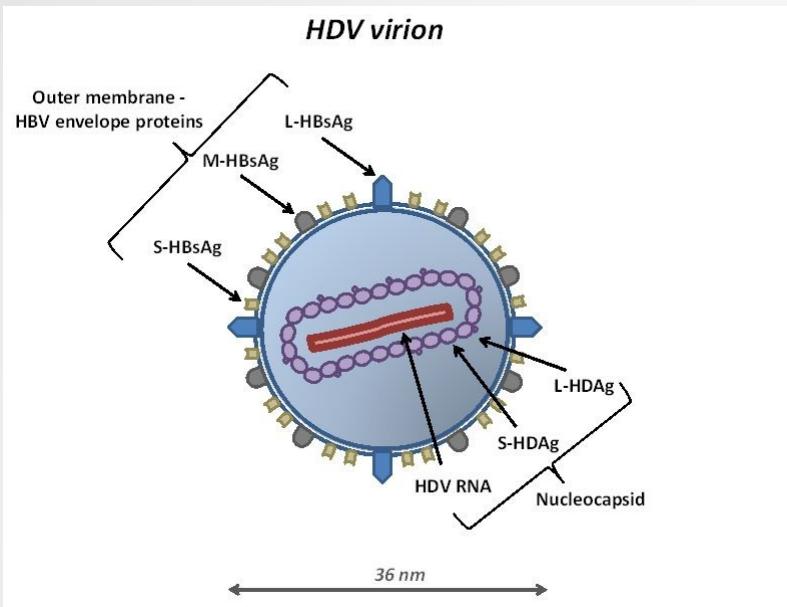


Megavirus chilensis

**1.26 Mbp
1120 genes**

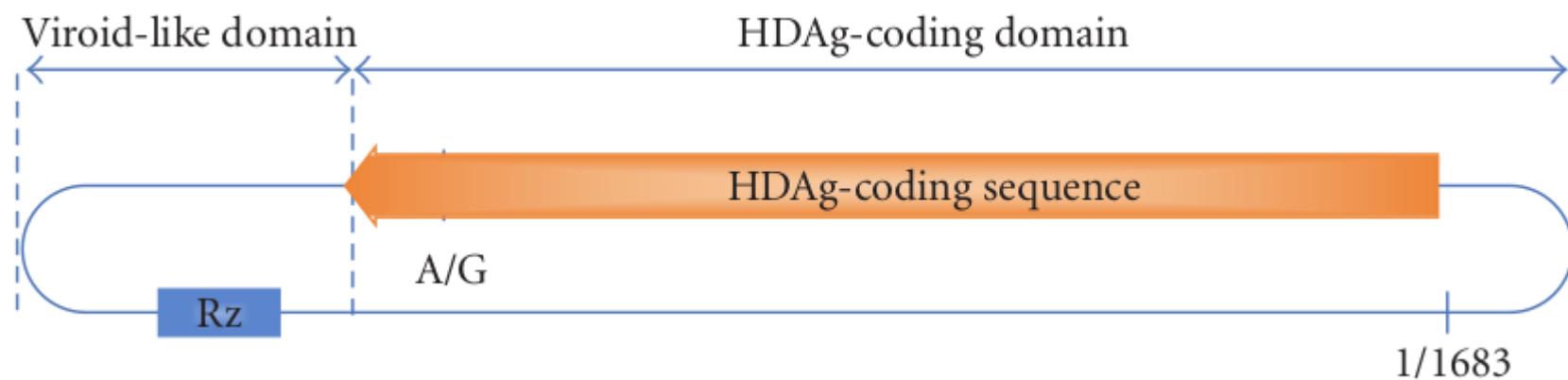
Legendre et al, 2012

Smallest viral genome



Hepatitis D virus
1683 bp

Symbiont of
Hepatitis B virus

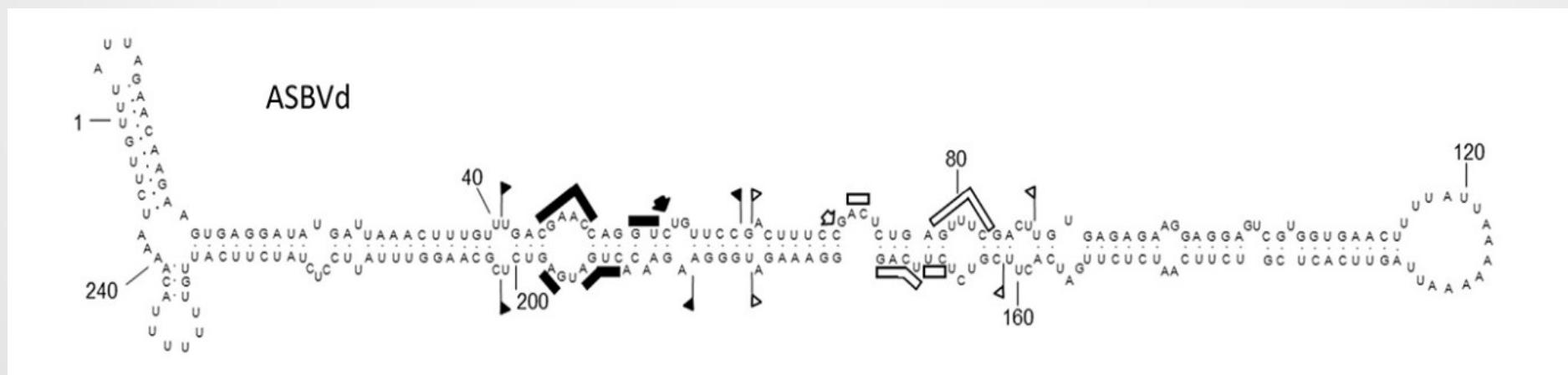


Smallest genome

Avocado sunblotch viroid

248 bp of circular RNA.

No capsid. No envelop of any type.
ONLY pure naked RNA.
Nothing more.



1 molecule = 1 organism(?)

Smallest and biggest genomes

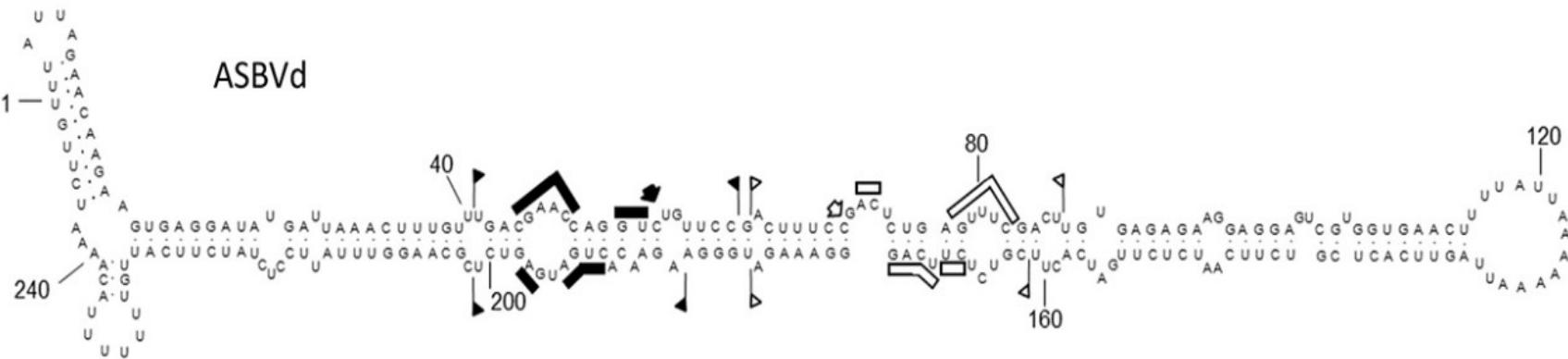


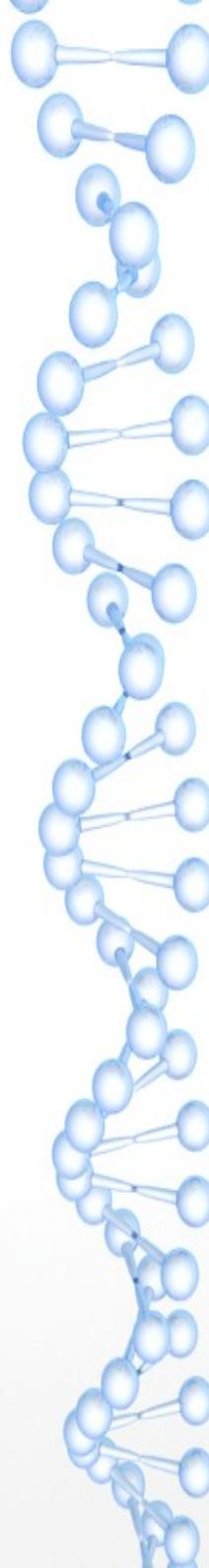
VS

Plant vs viroid
Multicellular complex organism
vs
short circular RNA only organism

150 000 000 000 bp
vs
248 bp

ASBVd





I. Structure and diversity of the genomes

Viruses and(?) viroids

Taxonomy of viruses and(?) viroids

ICTV

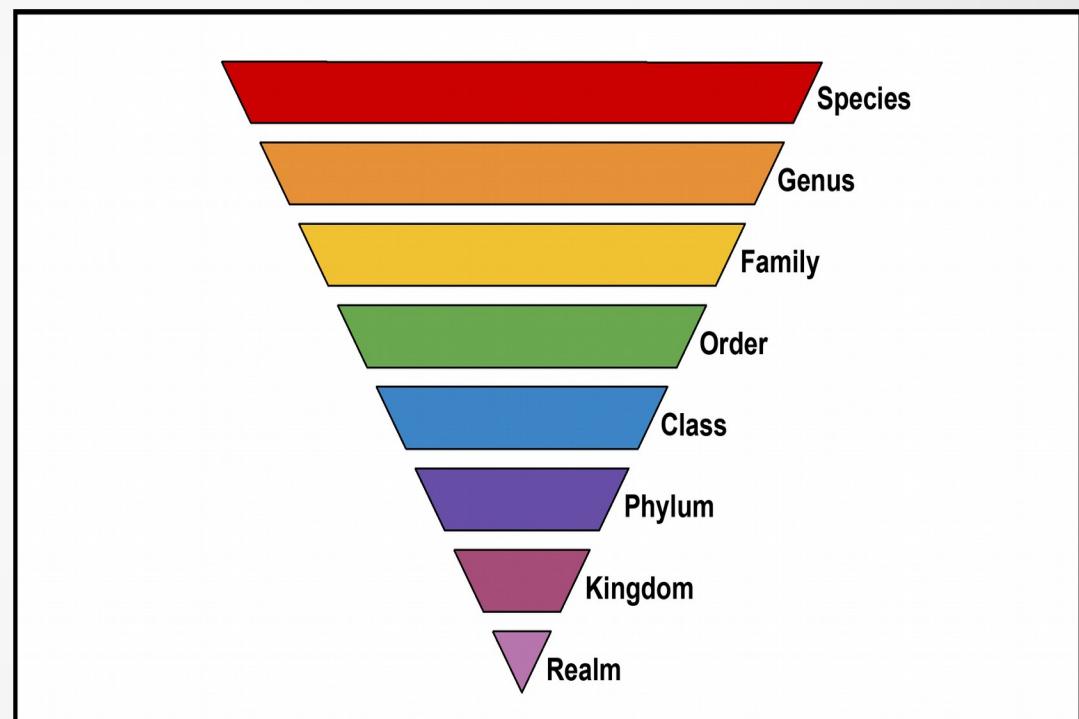
International Committee on Taxonomy of Viruses

Latest Taxonomy Release #37 (2021)

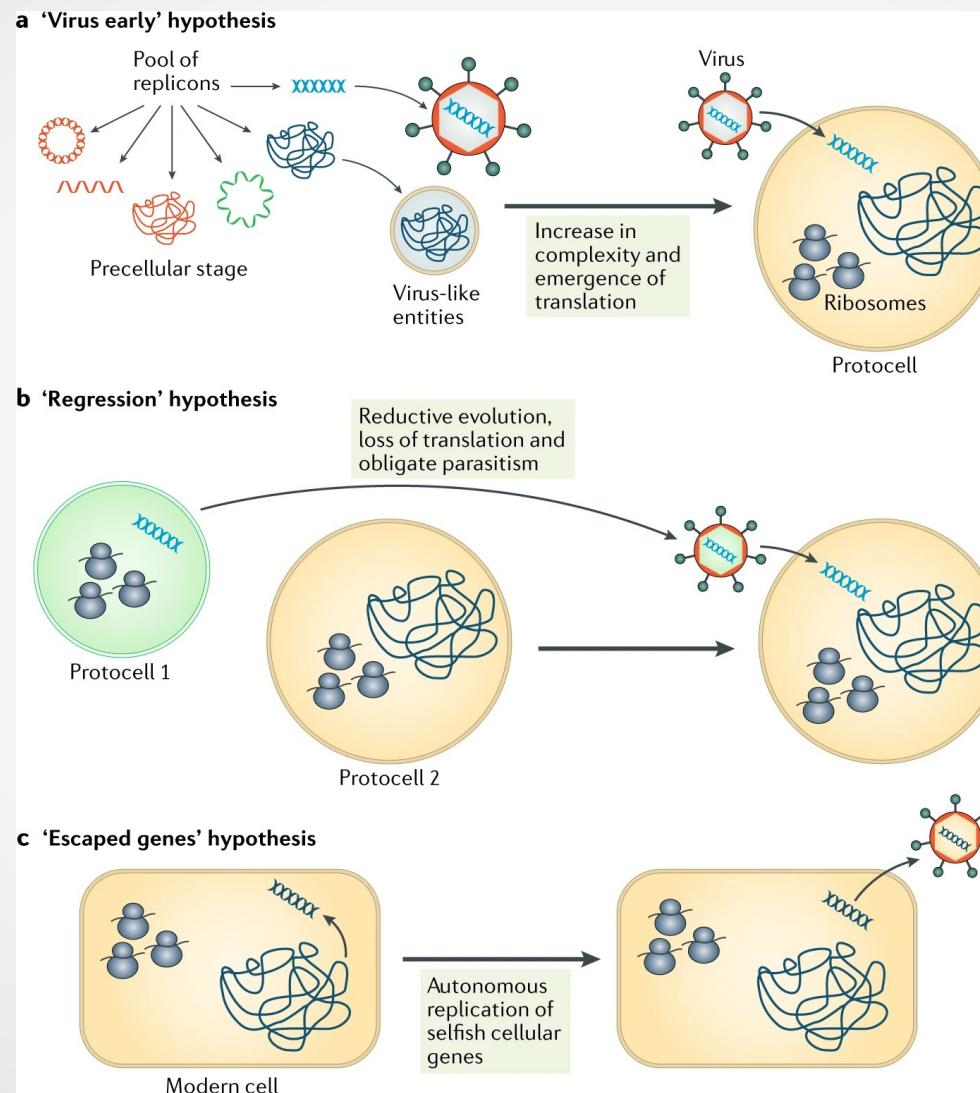
- 8 realms
- 1 class without higher classification
- 19 families without higher classification
- 2 genera without higher classification

Classification of viruses is very difficult.

Some realms are not monophyletic!

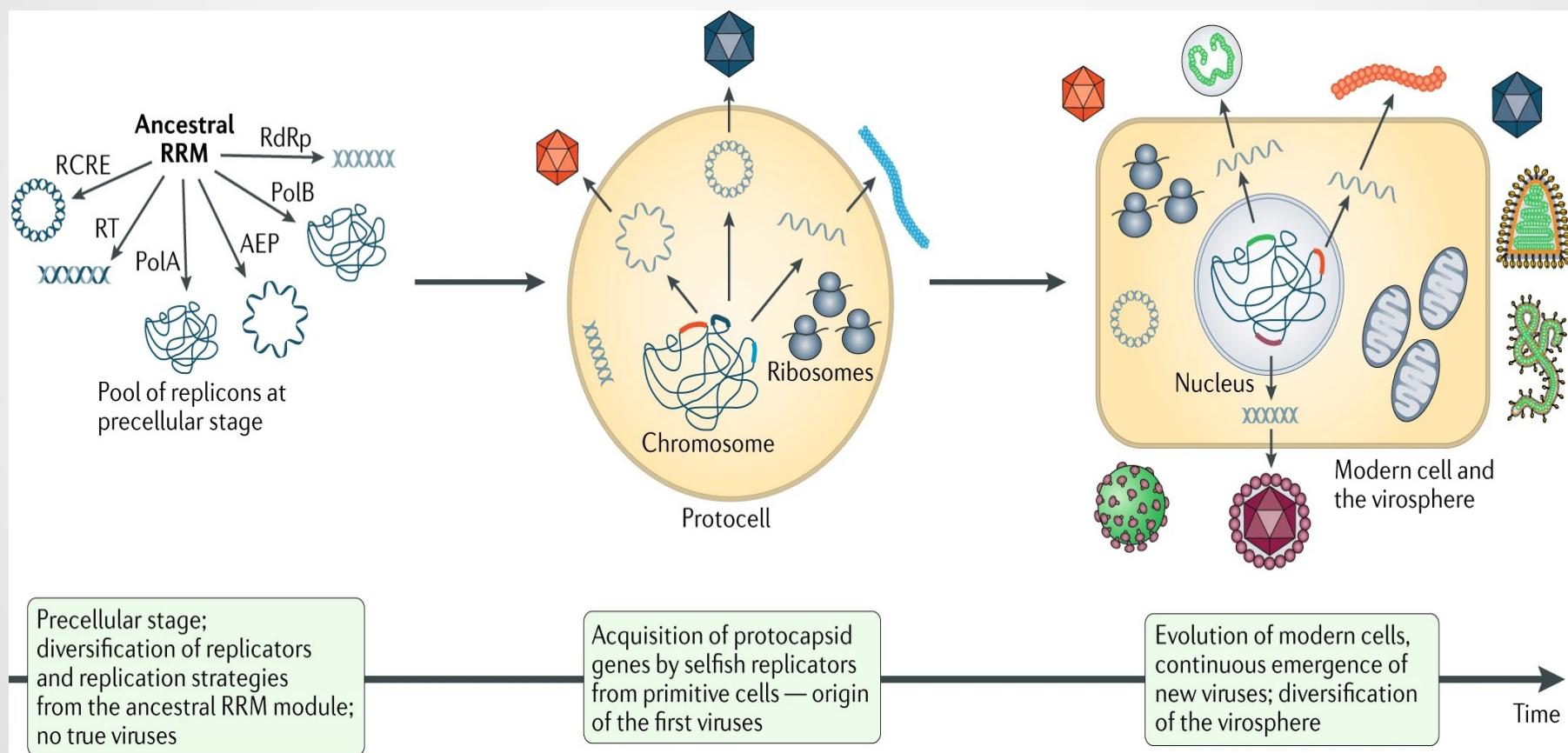


Origin of viruses: three basic scenarios



Krupovic et al, 2019

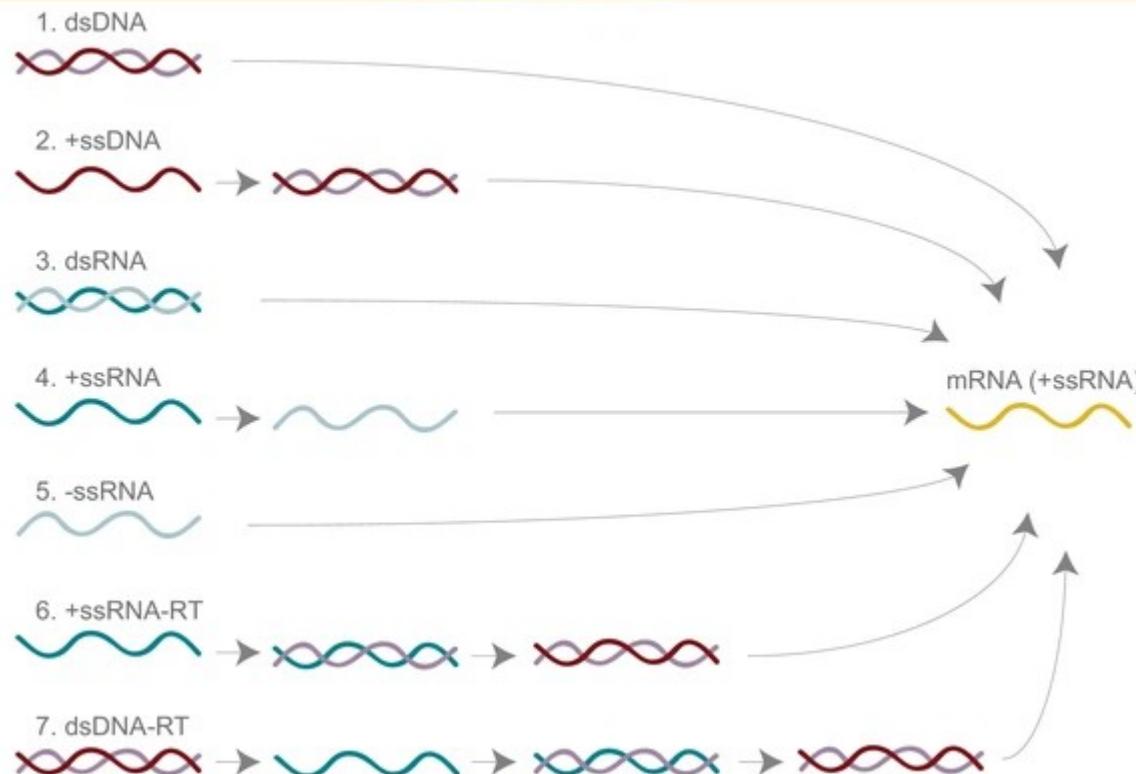
Origin of viruses: chimeric scenario



Krupovic et al, 2019

Types of viral genomes (1)

A. Baltimore Classification

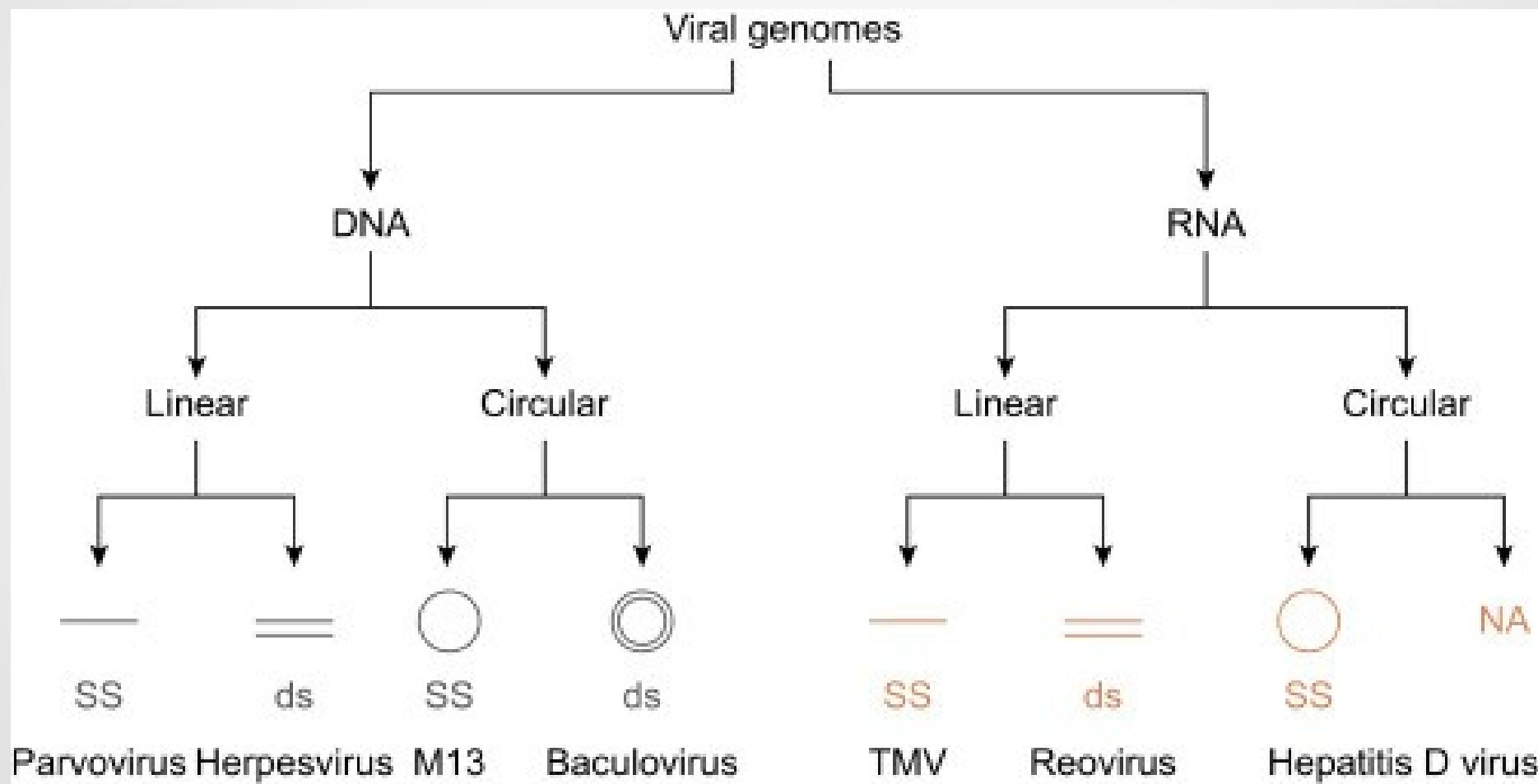


dsDNA - double stranded DNA
ssDNA - single stranded DNA
dsRNA - double stranded RNA
ssRNA - single stranded RNA

RT - reverse transcription
(RNA → DNA)

+/- - relation to expressed mRNAs

Types of viral genomes (2)



Types of viral genomes (3)

nonsegmented genomes

genome = single molecule

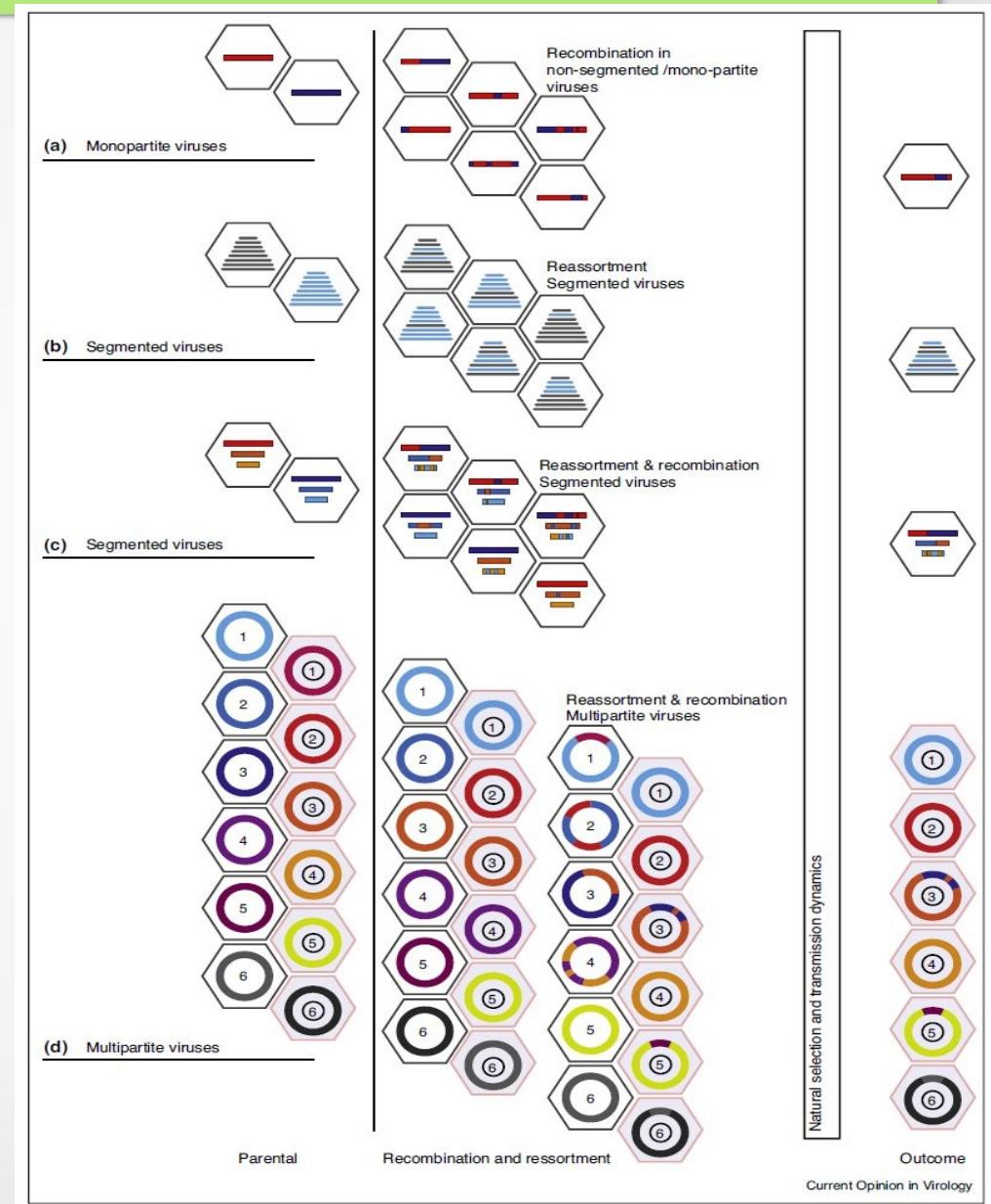
segmented genomes

genome = several molecules,
must be packed in a single virion

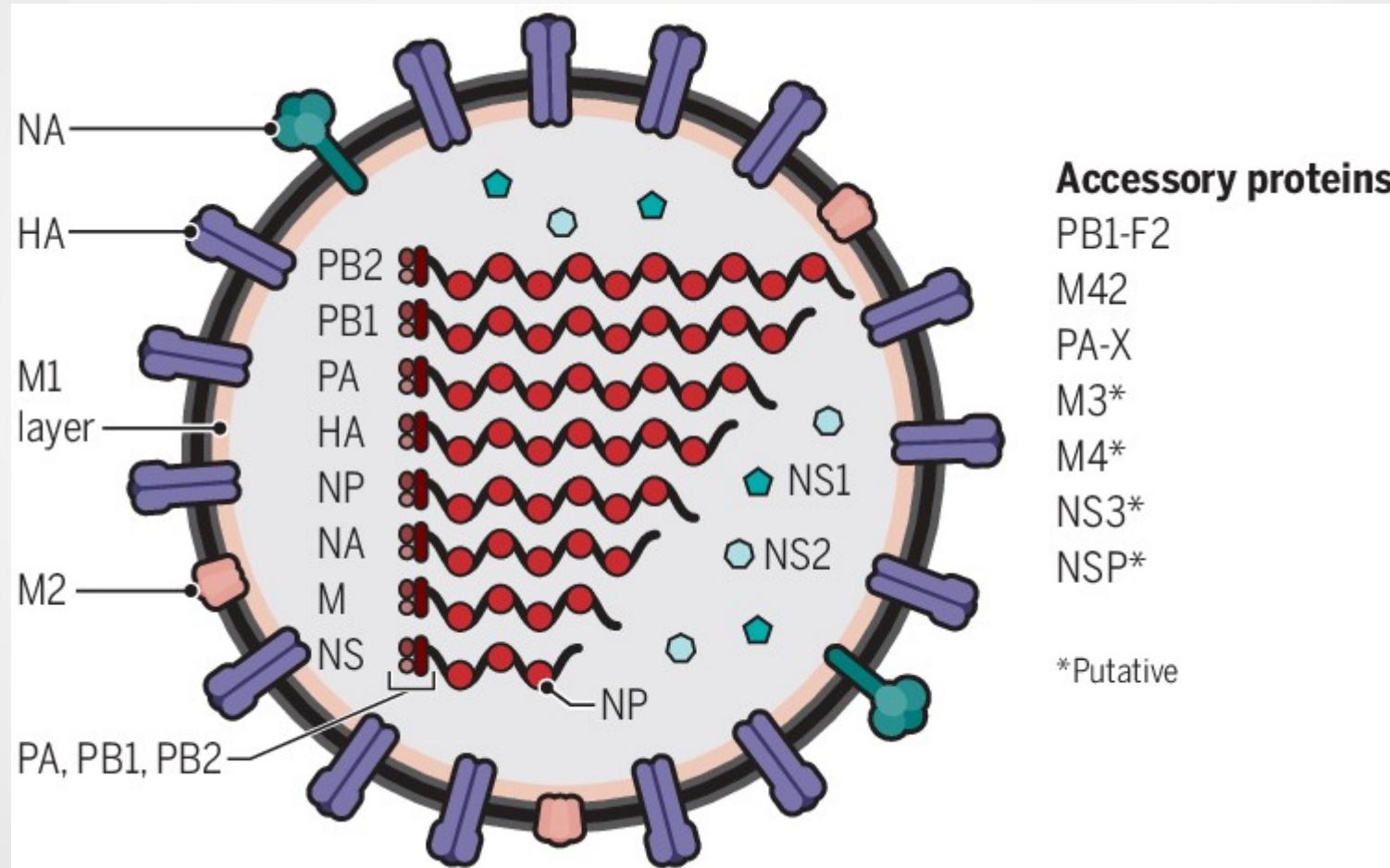
multipartite genomes

genome = several molecules, is
packed in several virions

Varsani et al, 2018



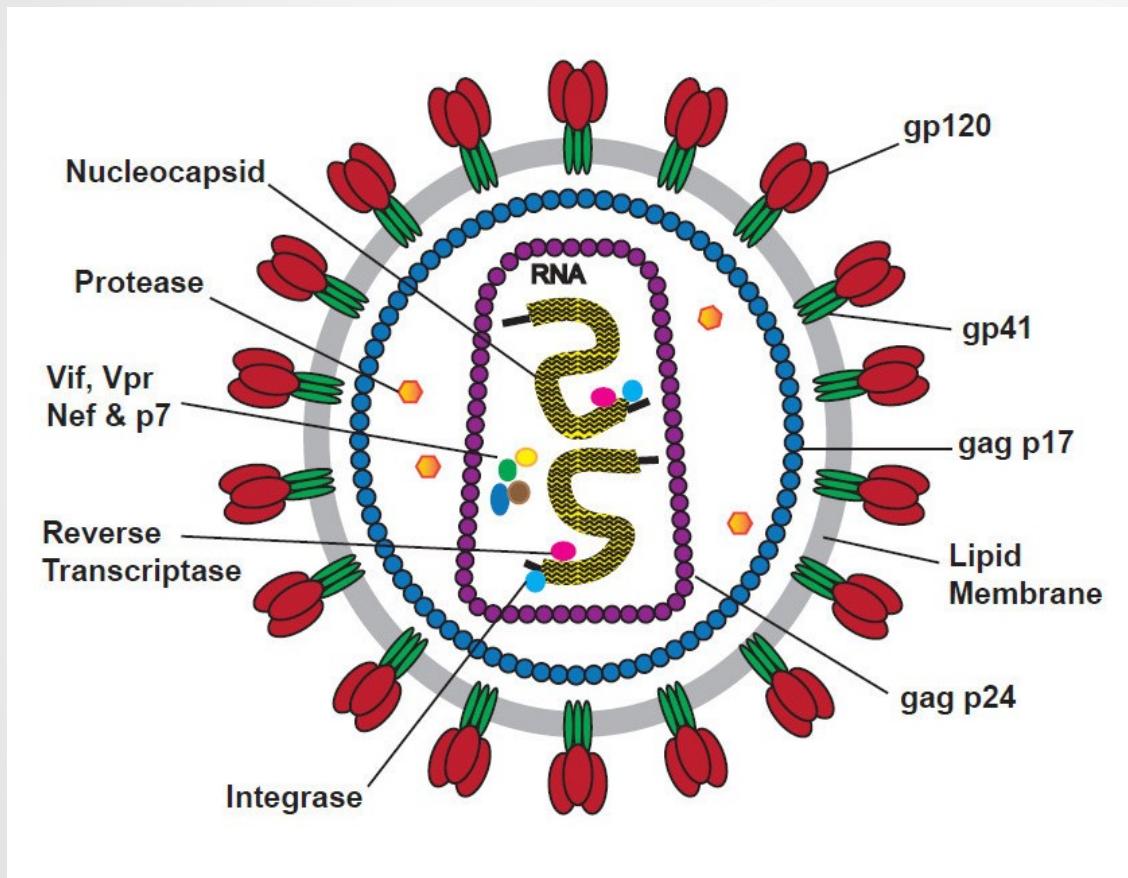
Genome of Influenza A virus



-ssRNA linear segmented genome
8 genome segments

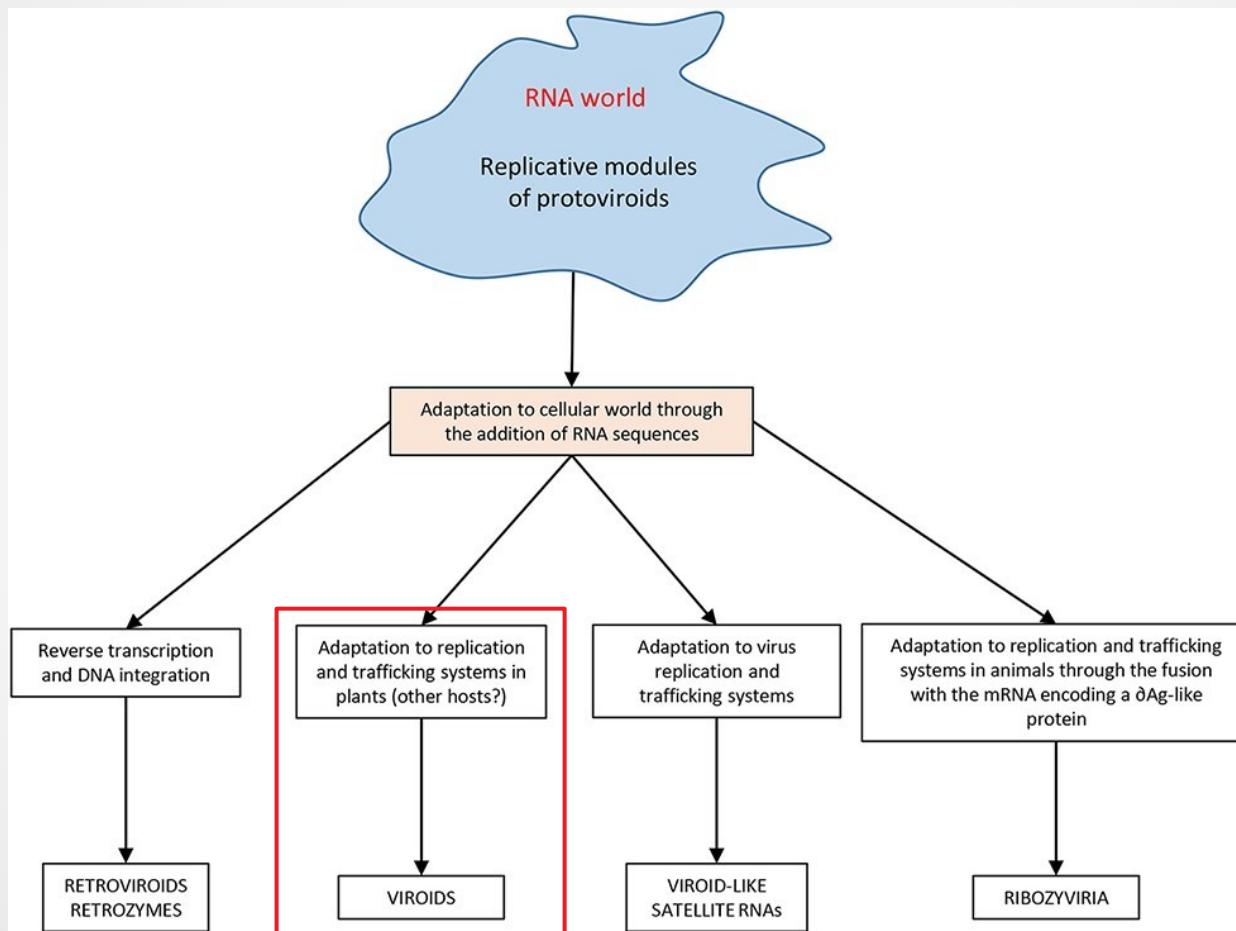
Taubenberger et al, 2019

Human immunodeficiency virus (HIV)

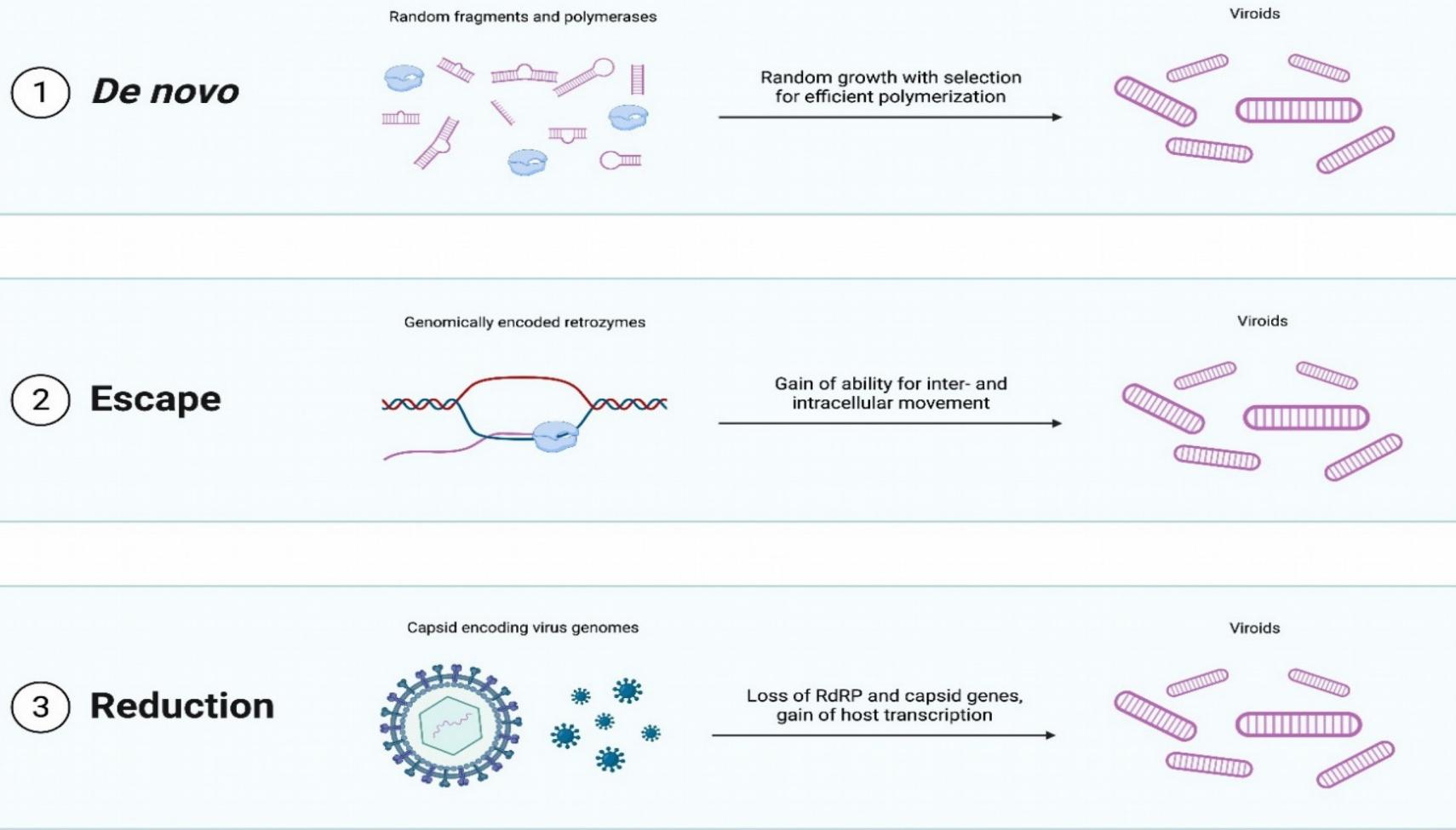


+ssRNA-RT linear
non-segmented genome
but in two copies in each virion

Origin of viroids (1): remnants of RNA world?



Origin of viroids (2): “modern” scenarios



Lee and Koonin, 2022

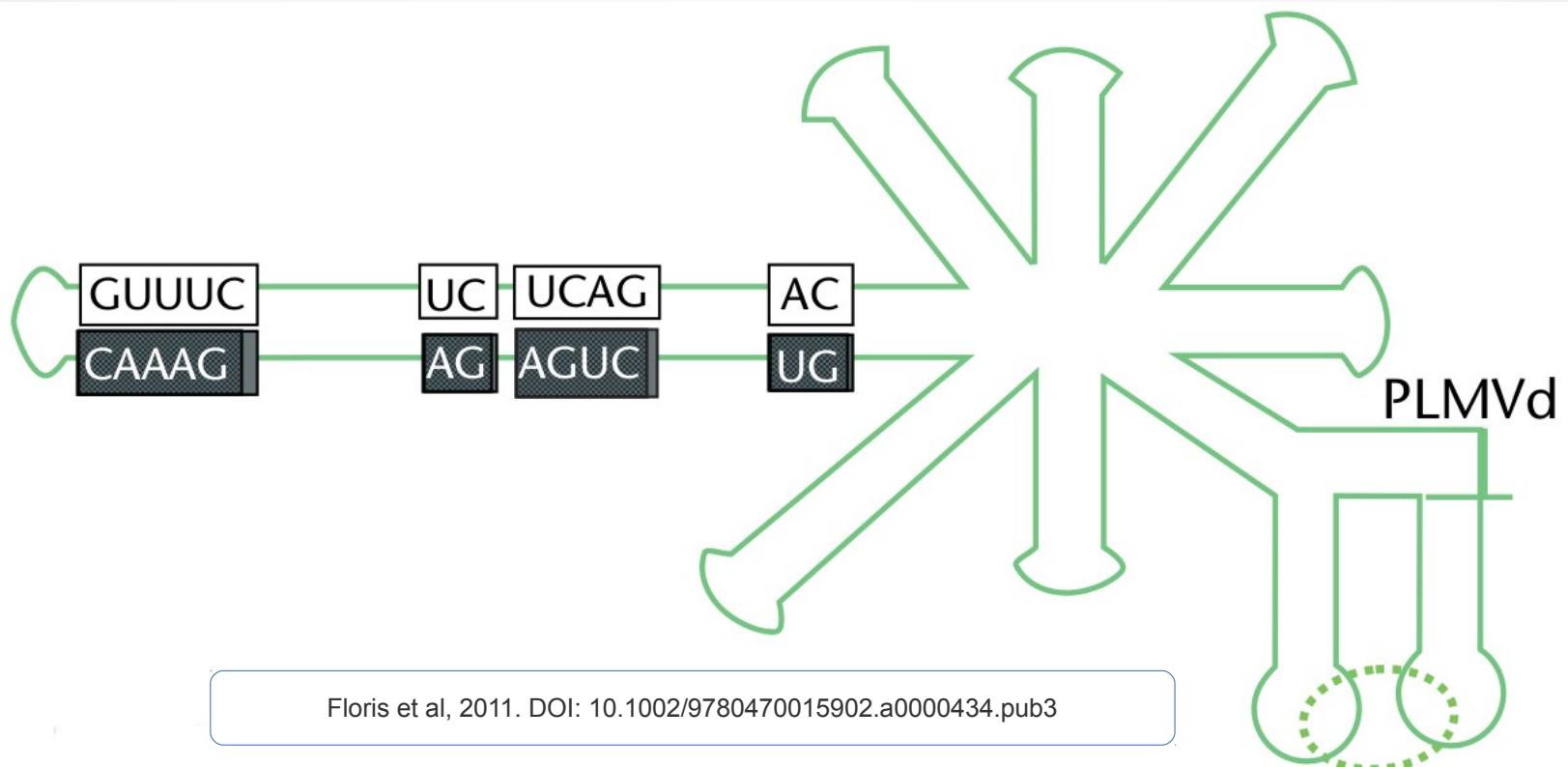
Structure of viroids

Avocado sunblotch viroid



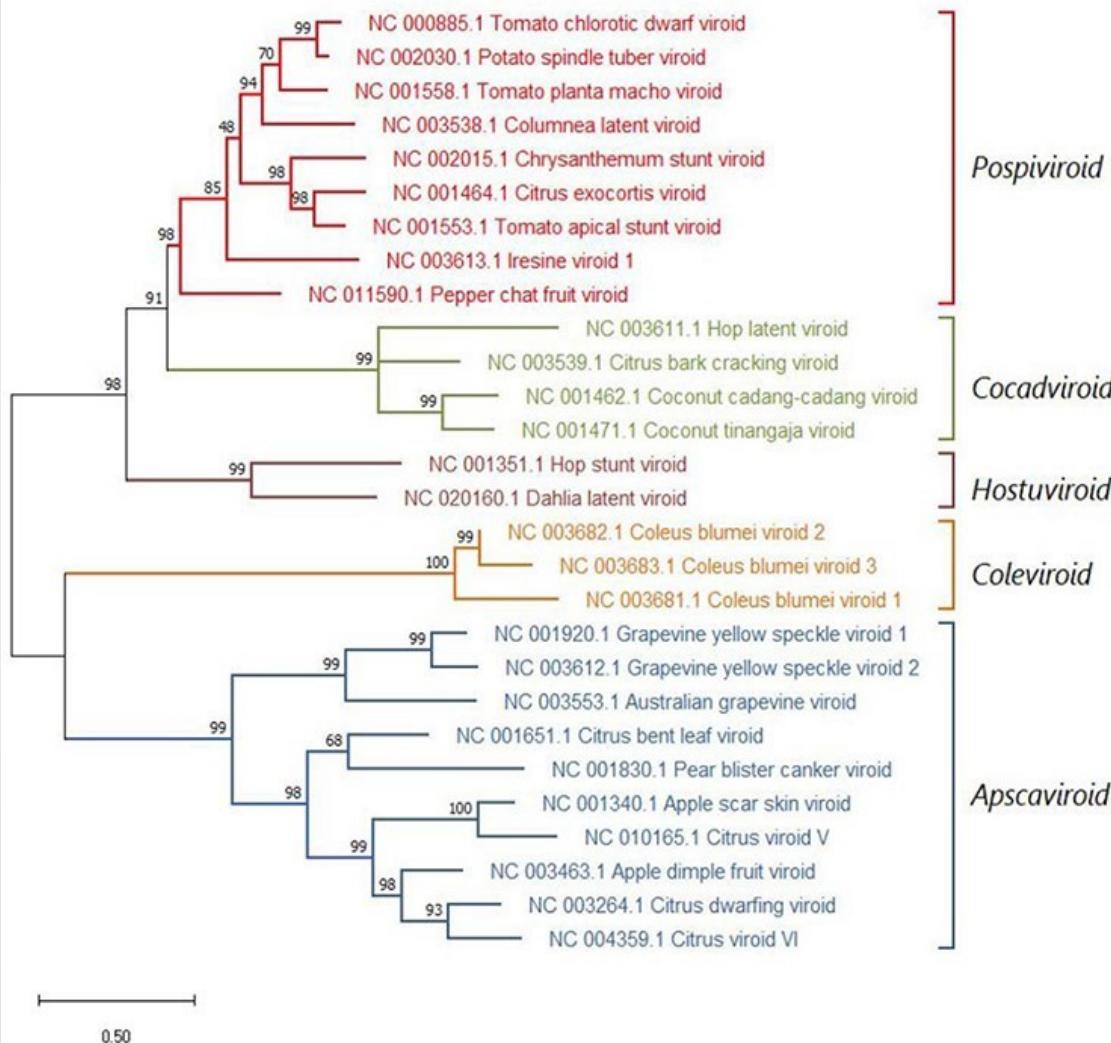
ASBVd

Peach latent mosaic viroid



PLMvd

Viroid taxonomy



species

Family Avsunviroidae

Genus <i>Avsunviroid</i>	1
Genus <i>Pelamoviroid</i>	2
Genus <i>Elaviroid</i>	1

- Use plastid RNA polymerase to replicate

Family Pospiviroidae

Genus <i>Pospiviroid</i>	10
Genus <i>Hostuviroid</i>	1
Genus <i>Cocadviroid</i>	4
Genus <i>Apscaviroid</i>	10
Genus <i>Coleviroid</i>	3

- Use nuclear RNA polymerase II to replicate

Summary for viral(-like) genomes

- Genomes are very simple but very diverse
- Genome size vary between 248 (1683) bp to 1.26 Mbp
- Genomes are based on RNA or DNA
- Single stranded or double stranded genomes
- Circular or linear genomes
- Nonsegmented, segmented and multipartite genomes
- Multipartite genome is a very interesting phenomena when organism could consists of several not connected parts during transmission between hosts
- horizontal gene transfer (HGT) is very common
- viruses could be vectors for (HGT) between very distant taxa

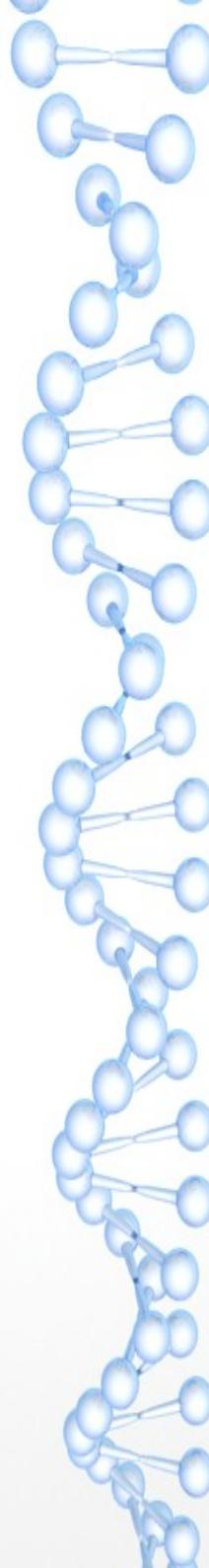
Issues with definitions of “genome” (1)

Genome is

- the entire set of ~~DNA~~ instructions found in a ~~cell~~.

RNA genomes

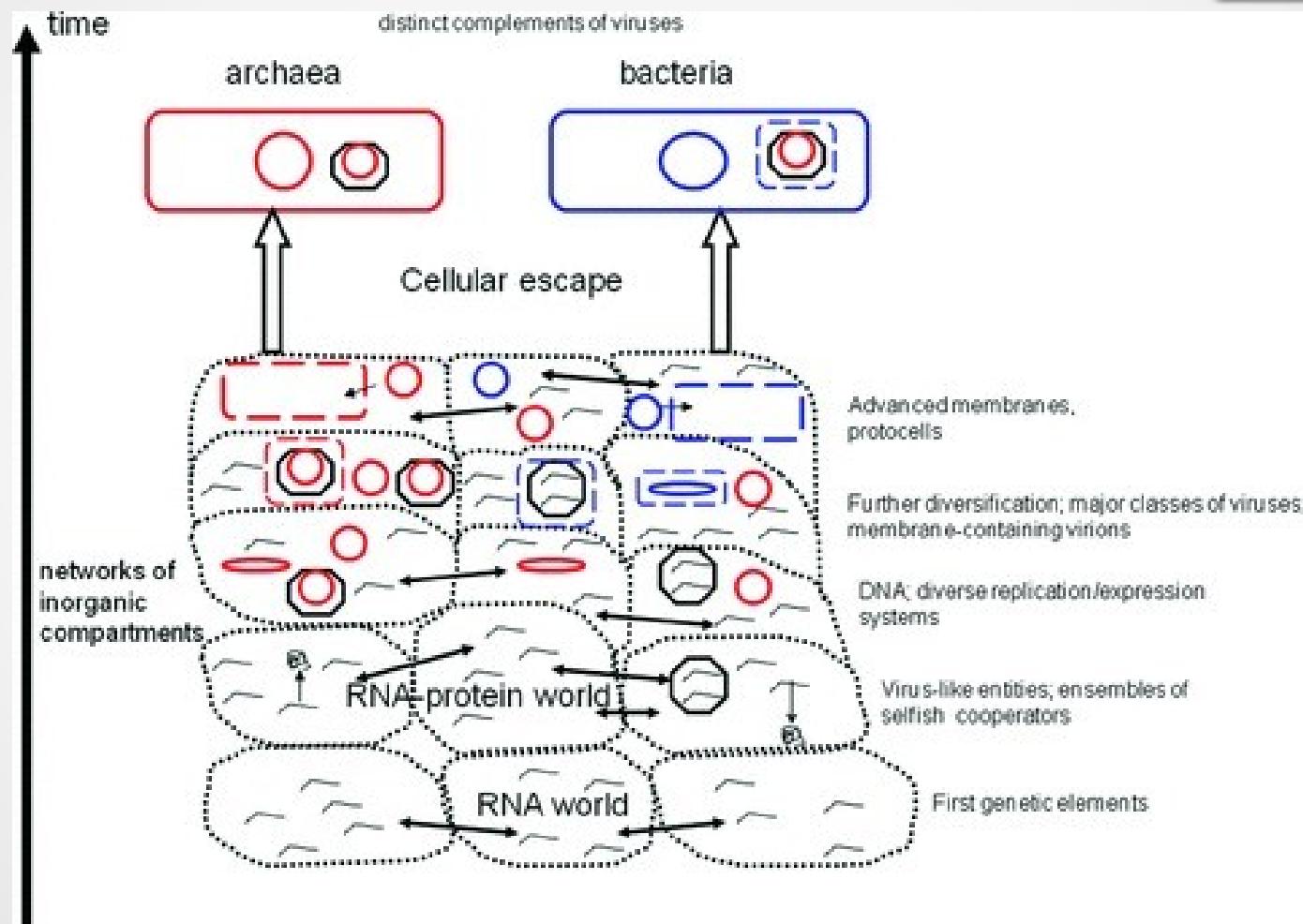
Noncellular organisms



I. Structure and diversity of the genomes

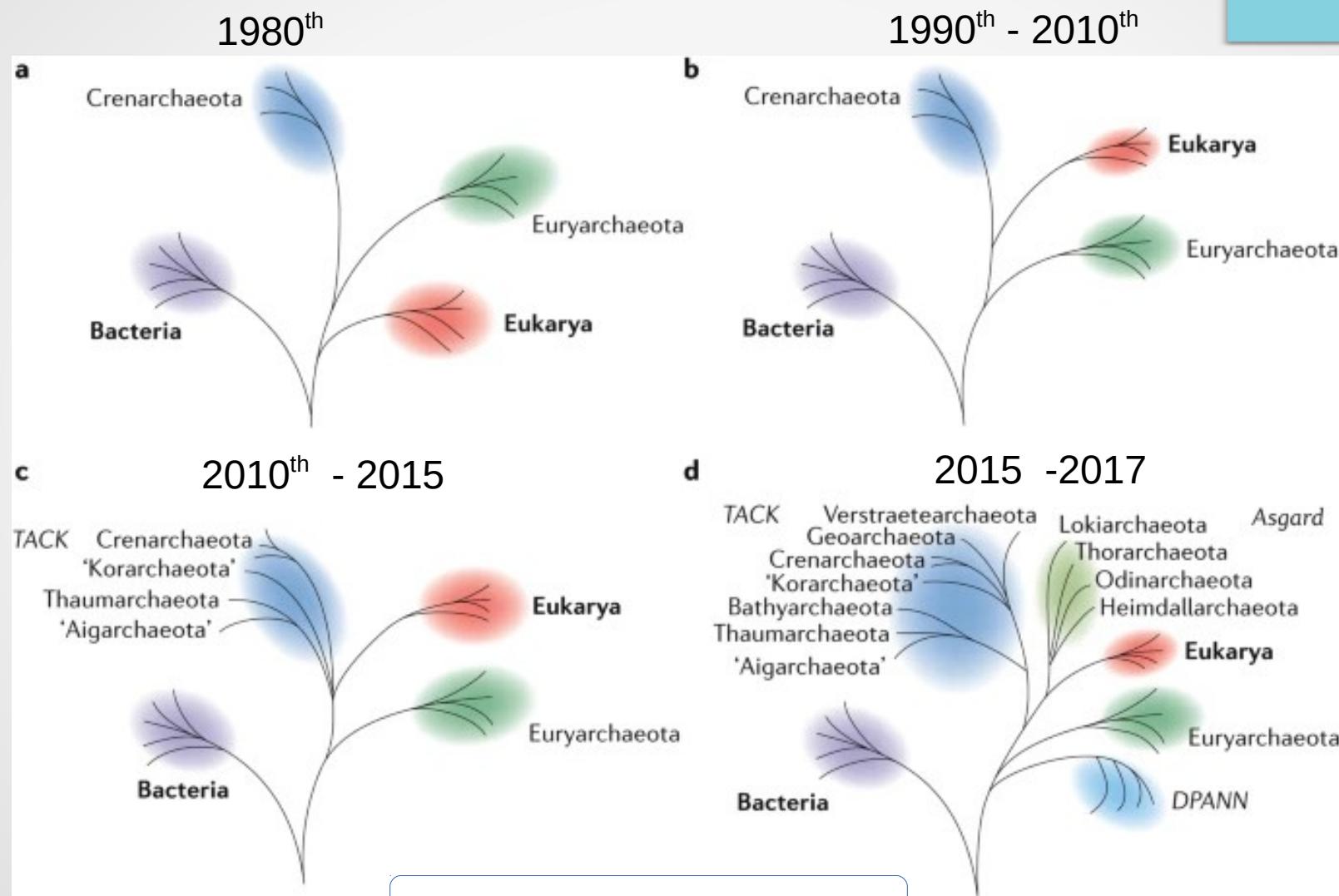
Cellular life

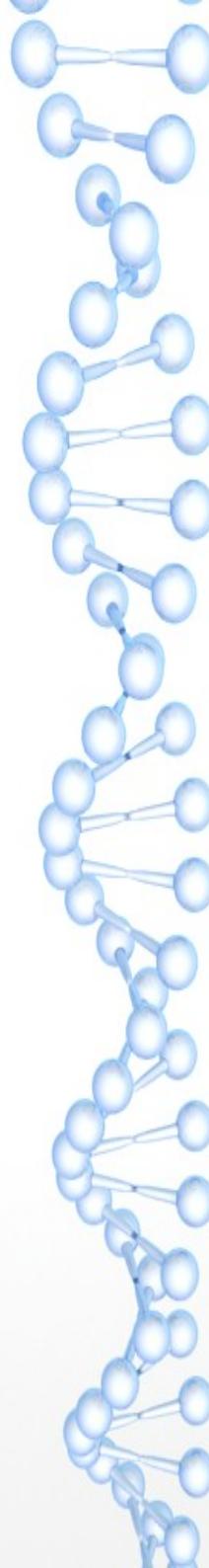
Possible origin of (cellular) life



Koonin, 2009

Evolution of tree of (cellular) life





I. Structure and diversity of the genomes

Cellular life
Prokaryotes

Structure of bacterial genome

Up to three dsDNA elements:

I. Chromosome (obligatory).

- oriC/DNA replication type
- contains core genes for cell existence

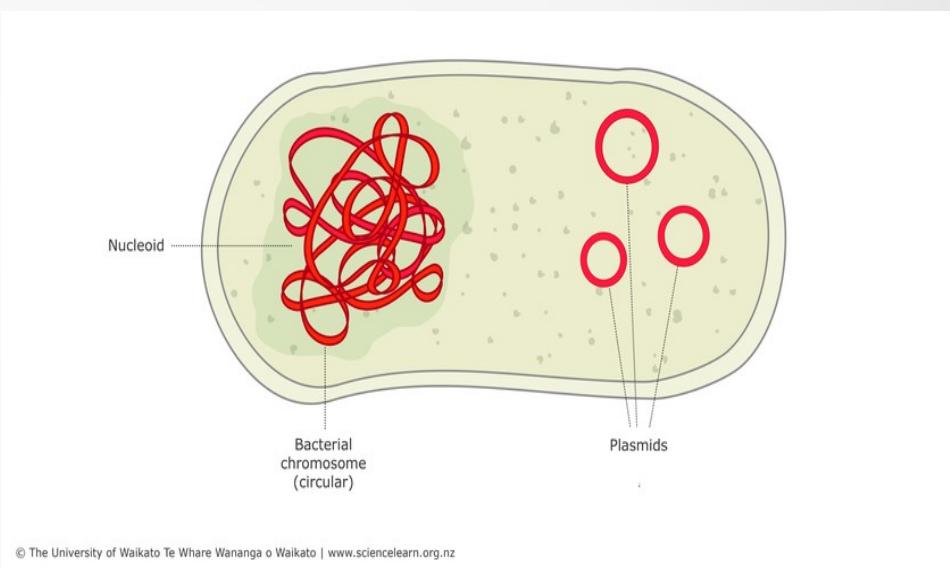
II. chromids (Harrison et al, 2010).

Obligatory for some clades

- plasmid-type replication
- could contain core genes
- comparable to chromosome size
- codon usage similar to chromosome

III. plasmids

- plasmid-type replication
- doesn't contain core genes
- could contain essential genes (resistance, conjugation)



© The University of Waikato Te Whare Wananga o Waikato | www.sciencelearn.org.nz

All element usually are circular, but sometimes could be linear

Structure of archeal genome

Up to three dsDNA elements:

I. Chromosome (obligatory).

- Orc1/Cdc6 replication type
- contains core genes for cell existence

II. mini-chromosome(s).

Obligatory for some clades

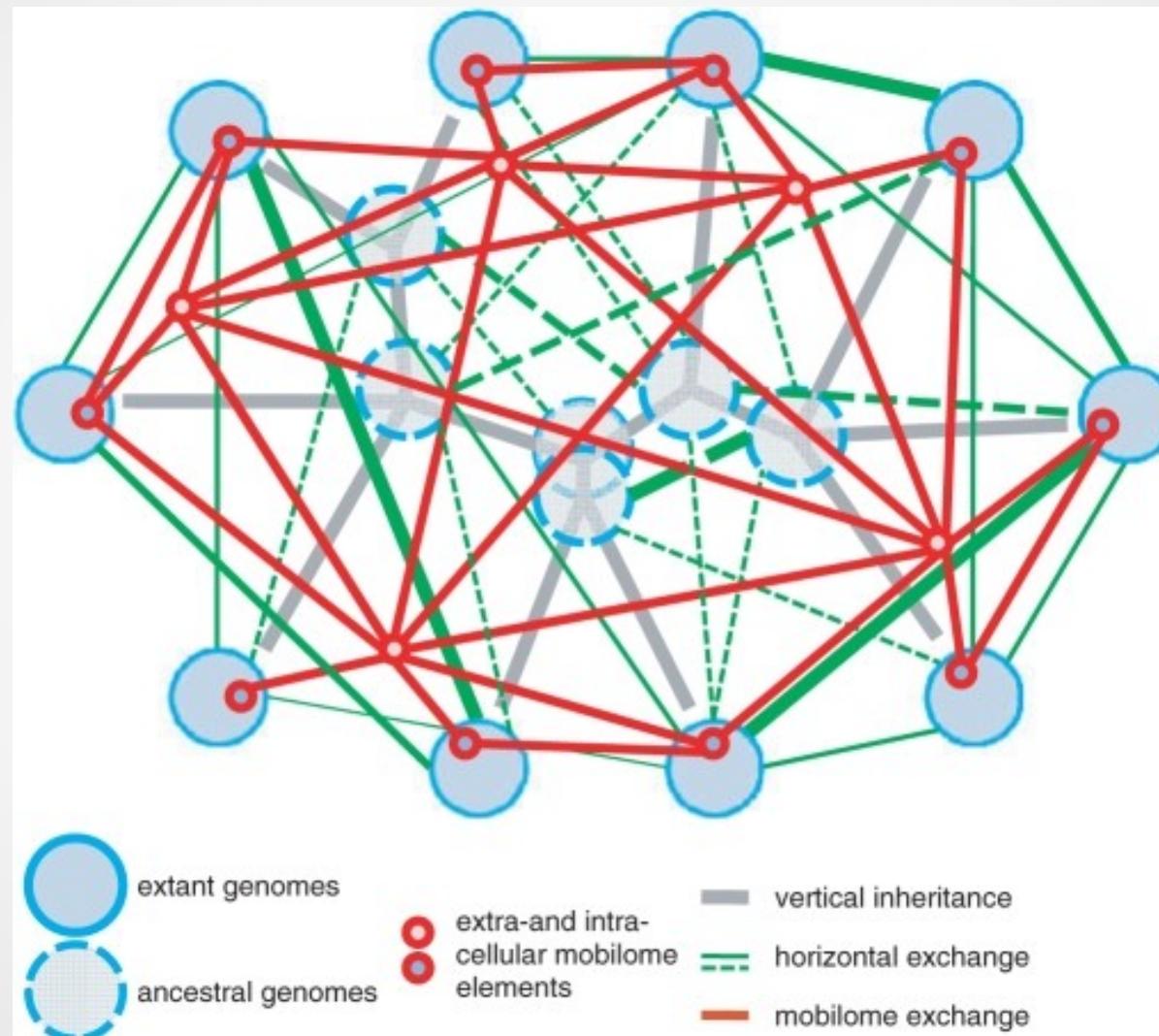
- chromosome-like replication type
- could contain core genes
- comparable to chromosome size
- codon usage similar to chromosome

III. plasmids

- plasmid-type replication
- doesn't contain core genes
- could contain essential genes (resistance, conjugation)

All element usually are circular, but sometimes could be linear

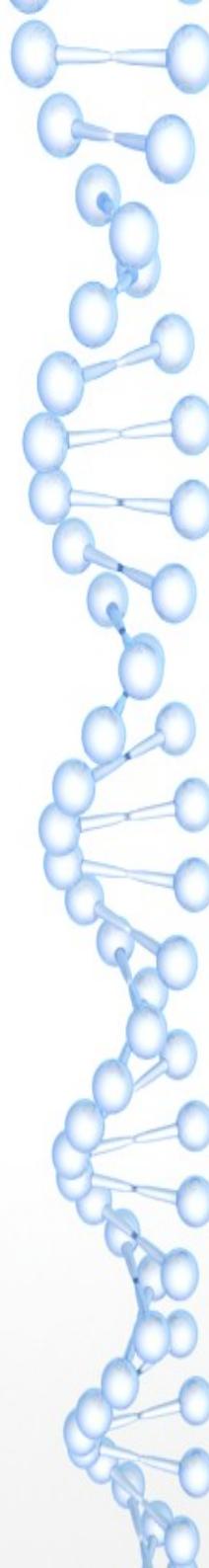
Gene flow in prokaryotic world



Koonin and Wolf, 2009

Summary for prokaryotic genomes

- Only double stranded DNA genomes
- Genome size vary between 112 kbp and 14.8 Mbp (*Sorangium cellulosum*)
- Up to three element types in the genome (chromosome + chromid/minichromosome + plasmids)
- deep difference between Bacteria and Archea in replication system, but high similarity in genome organization
- horizontal gene transfer (HGT) is very common as in viruses
- modern research supports existence of only two domains of life (Archea and Bacteria) with Eukaryotes included in Archea.



I. Structure and diversity of the genomes

Cellular life
Eukaryotes

Eukaryotic genomes

I. Basic type:

Nuclear genome

+ mitochondrial genome (many copies)

Ia. Basic type with lost mitochondria:

Nuclear genome only

One protist is known

Monocercomonoides sp

II. Plant type:

Nuclear genome

+ mitochondrial genome (many copies)

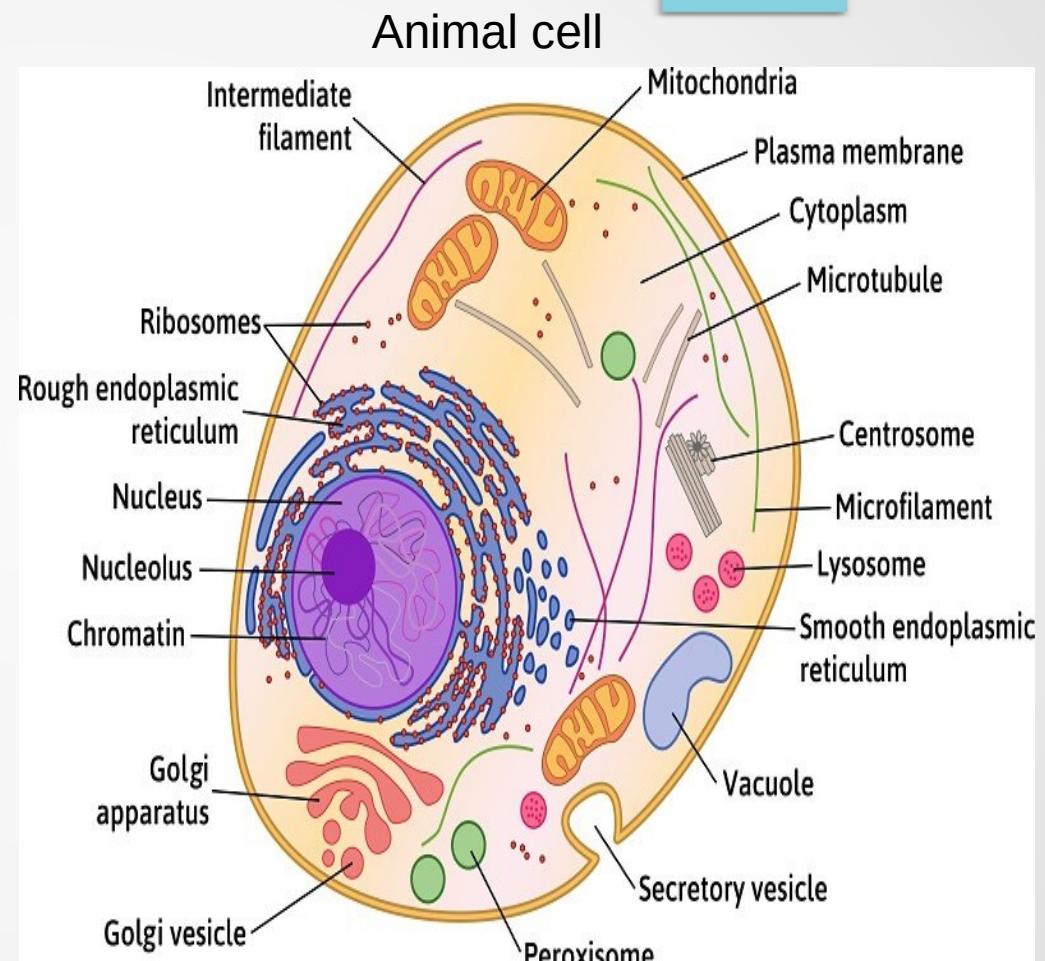
+ chloroplast genome (many copies)

IIa. Plant type with lost chloroplasts:

Nuclear genome

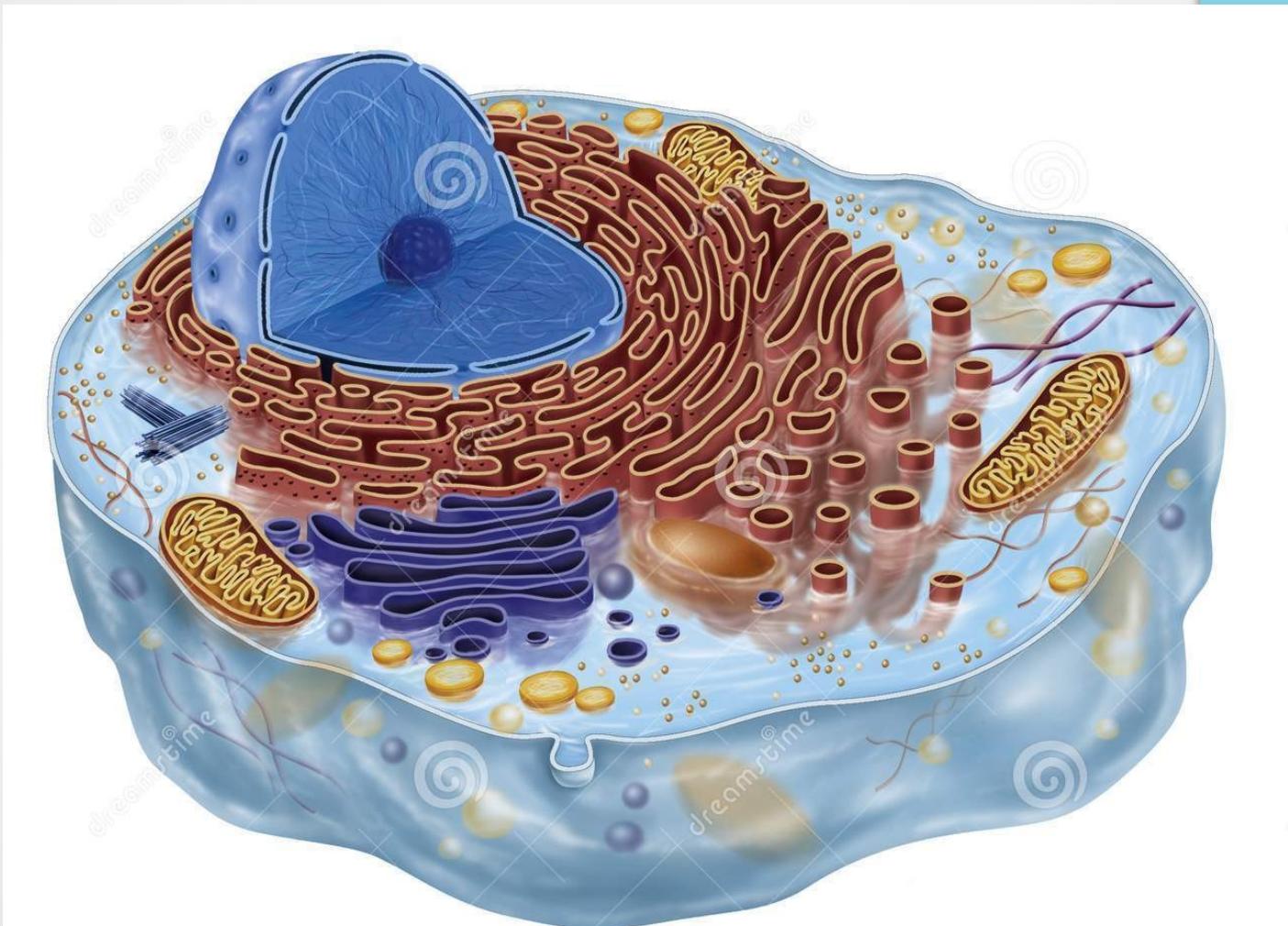
+ mitochondrial genome (many copies)

One (?) genus *Rafflesia*



<https://biologydictionary.net/animal-cell/>

Animal cell



Download from
Dreamstime.com

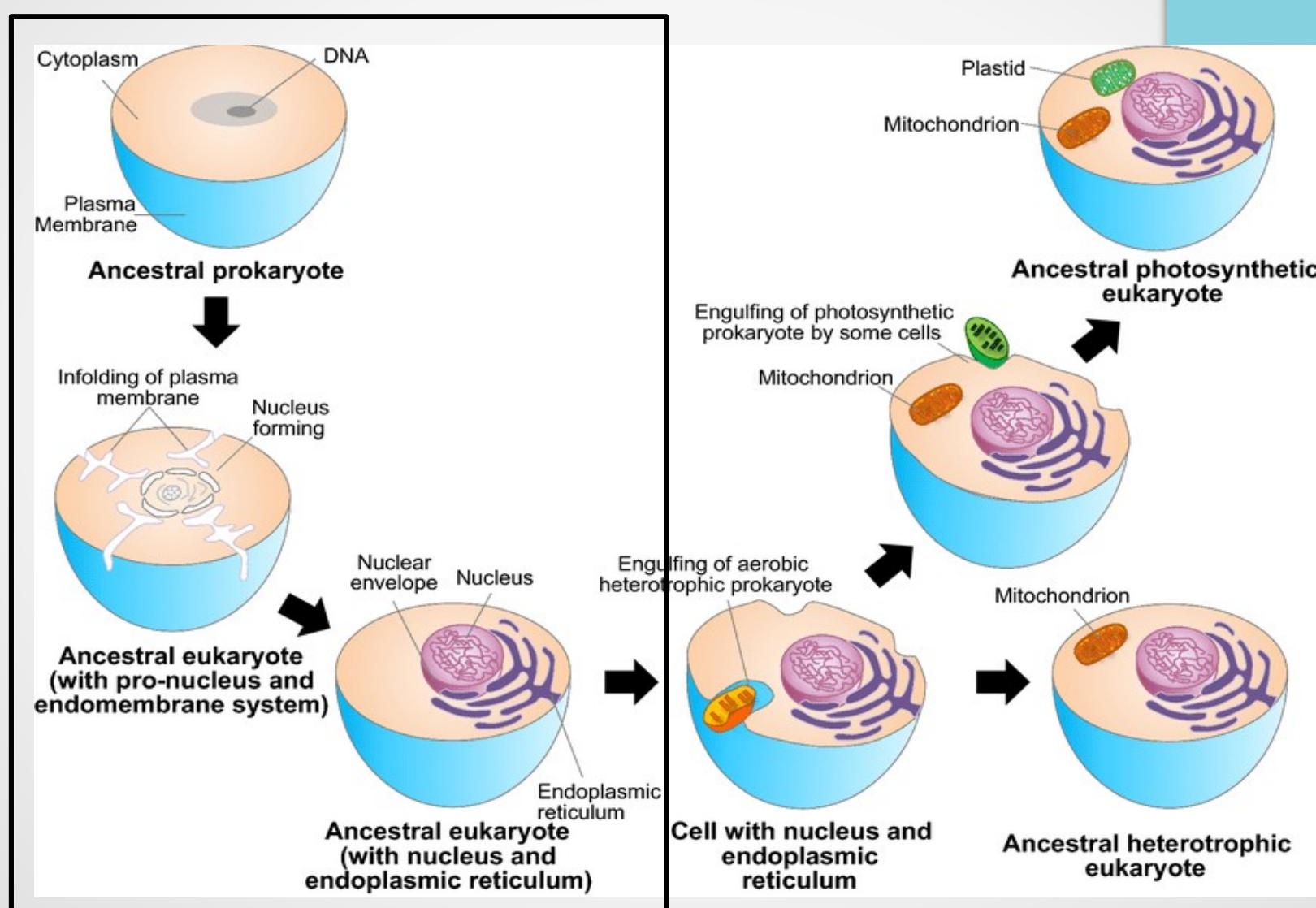
This watermarked comp image is for previewing purposes only.



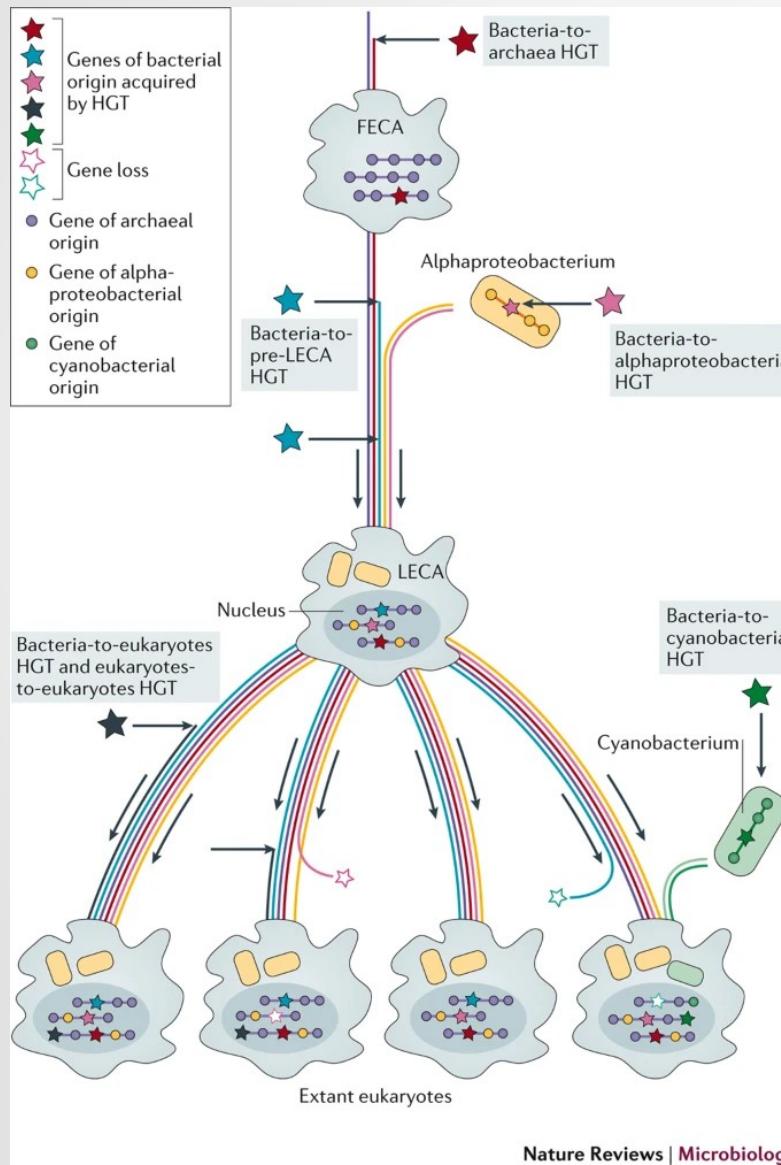
ID 37889236

© Frank Liu | Dreamstime.com

Origin of nucleus* and organells

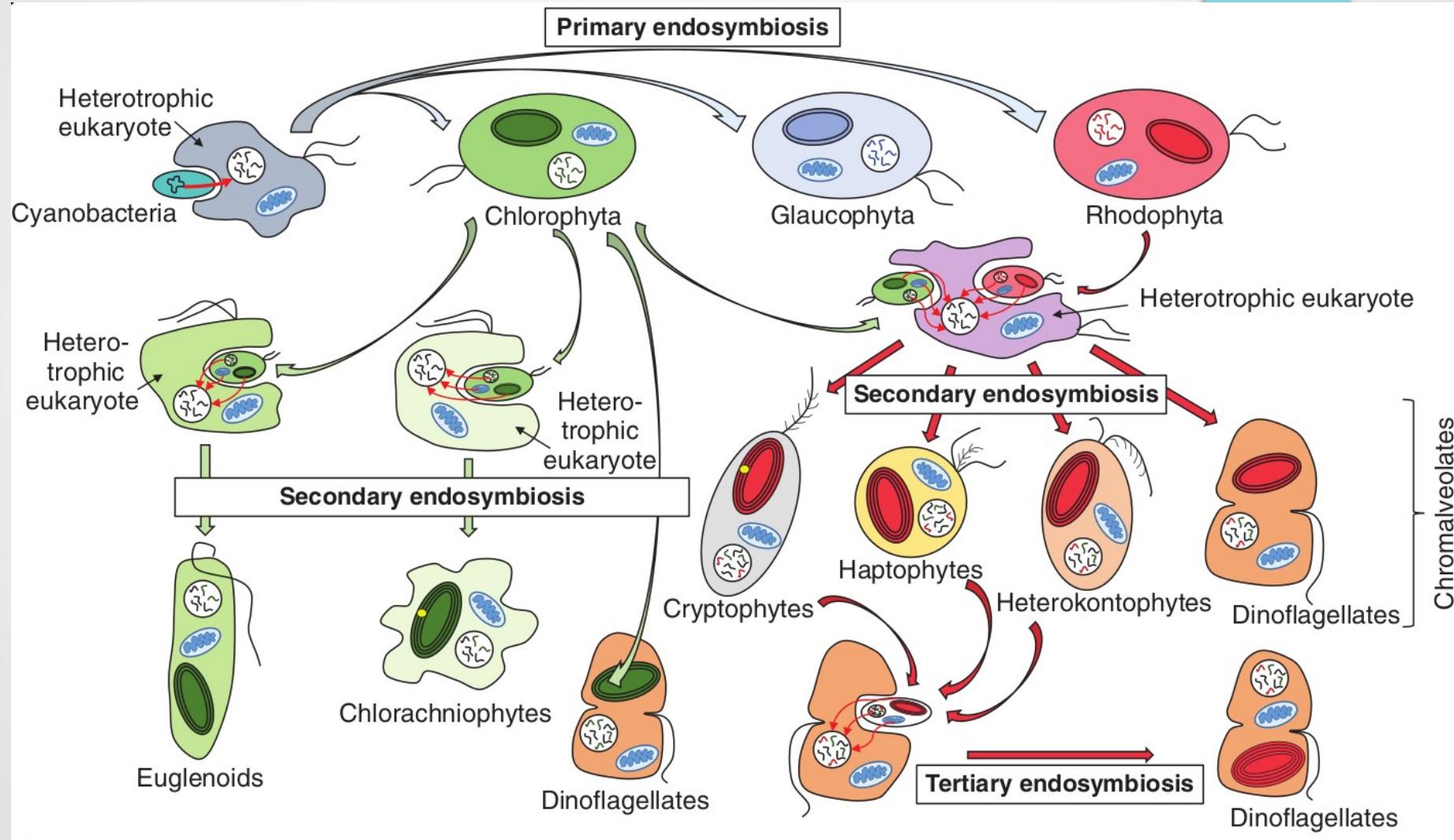


Chimeric origin of Eukaryotes



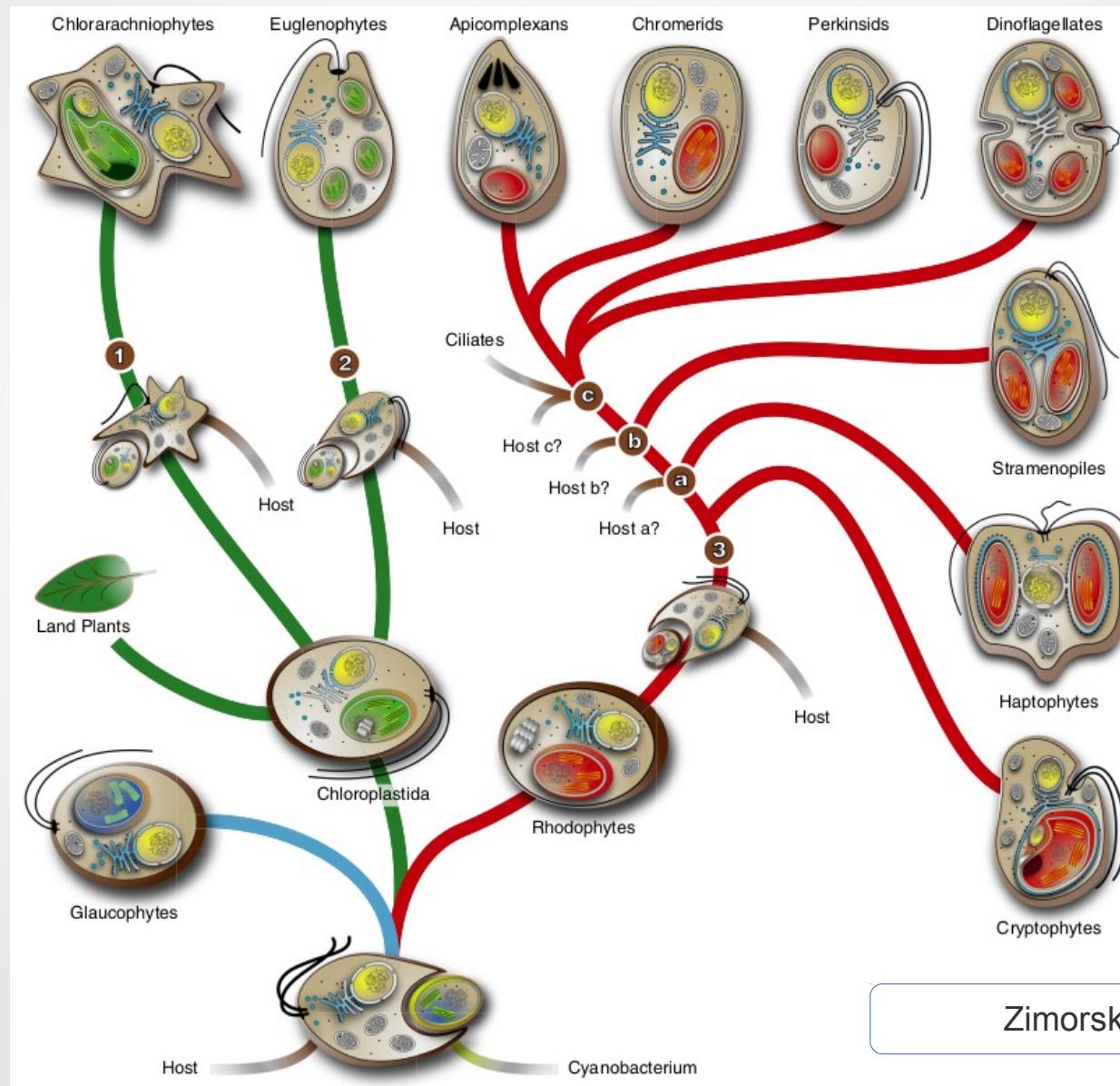
Eme et al, 2017

Origin of plastids (1)

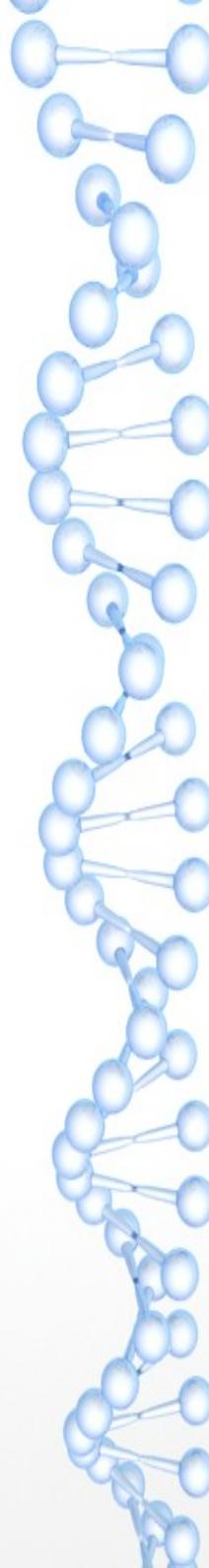


Hopes and Mock, 2015

Origin of chloroplasts (2)



Zimorski et al, 2014



I. Structure and diversity of the genomes

Genome in vivo
Karyotypes and chromosomes

Definitions

Chromosome

DNA and protein-containing structures in cells of eukaryotes, microscopically visible as a rod-shaped body during cell division metaphase

Karyotype

complete set of metaphase chromosomes in cells of an organism of a particular species.

Chromatine

DNA-protein complex. Material of chromosomes

Euchromatine

lightly packed and transcriptionally active chromatine. Enriched by genes.

Heterochromatine

highly packed chromatine. Low transcriptional activity. Few genes

Issues with definitions of “genome” (2)

- **Genome** is
genetic material of haploid set of **chromosomes**

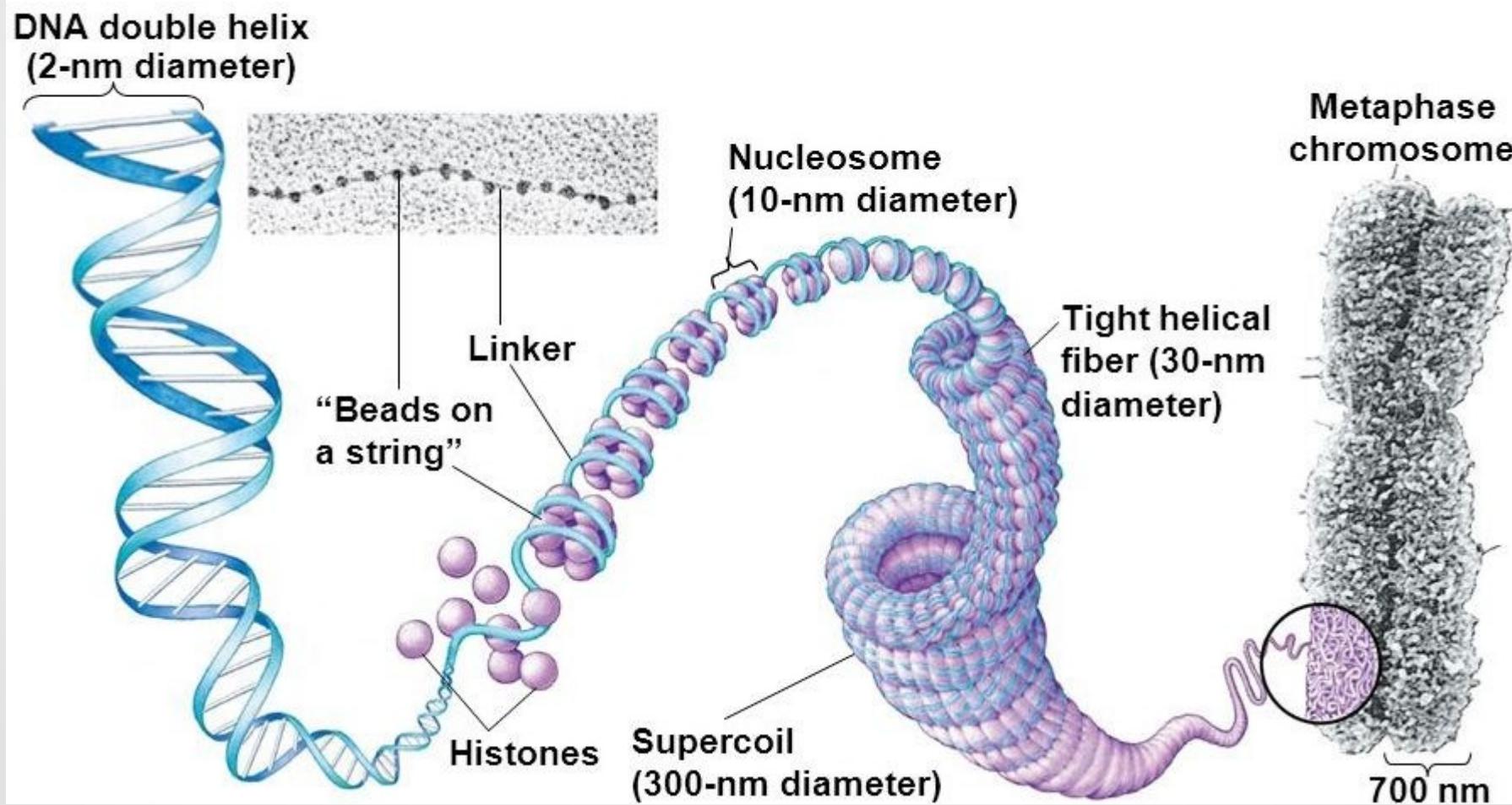
Chromosome

DNA and protein-containing structure in **cells of eukaryotes**,
microscopically visible as a **rod-shaped body during cell division**
metaphase

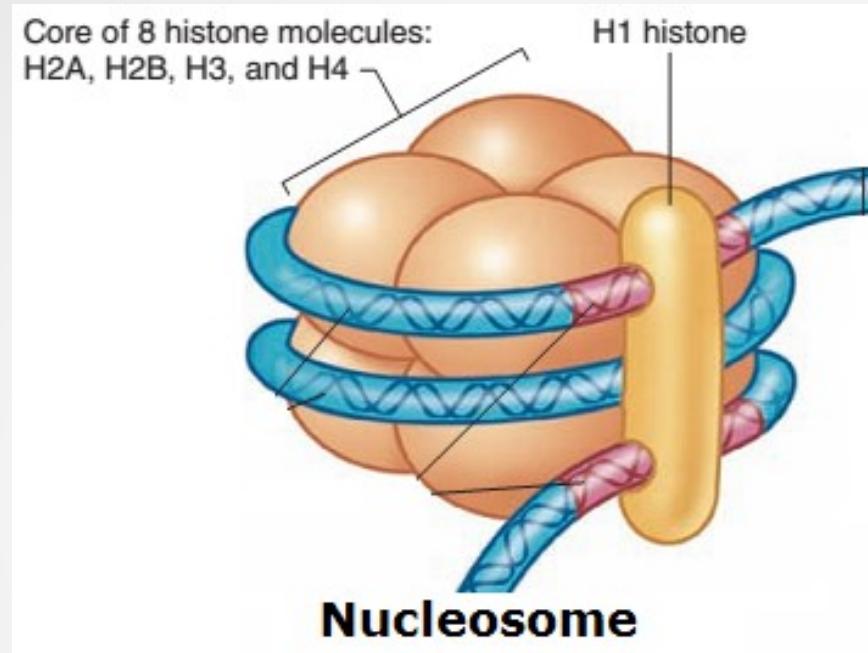
Bacteria? Viruses? mtDNA? chloroplast DNA?
Why not included?

**Definition (this) of chromosome was proposed by cytogeneticists of
“high” eukaryotes (plants and animals)**

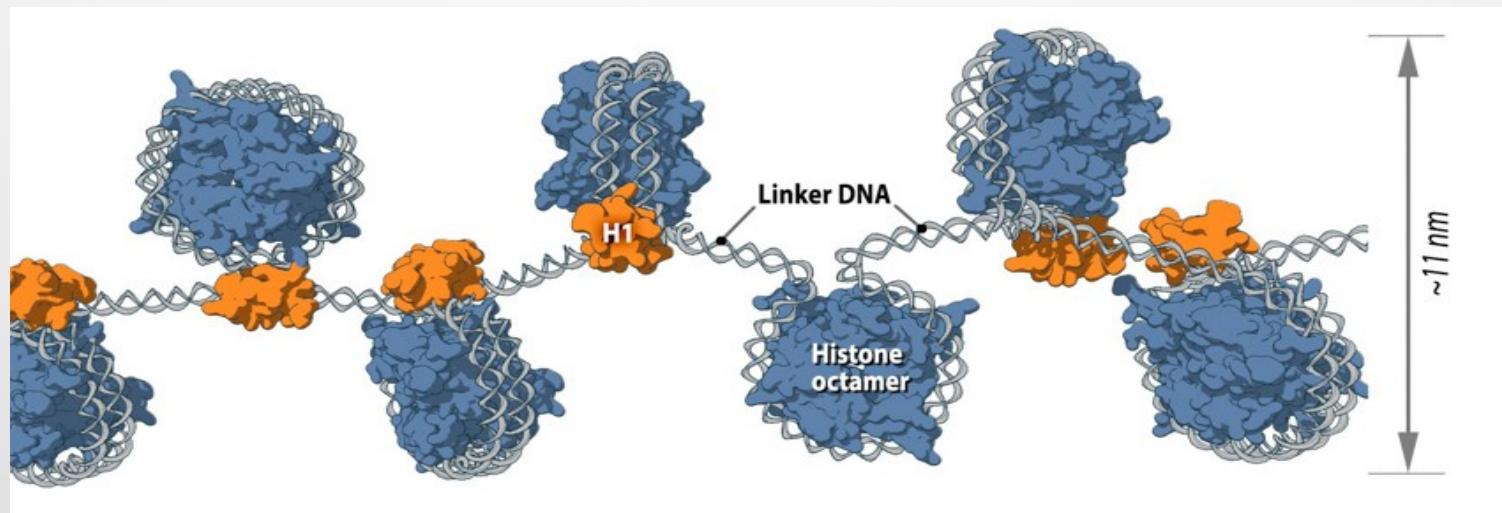
From DNA to chromosome



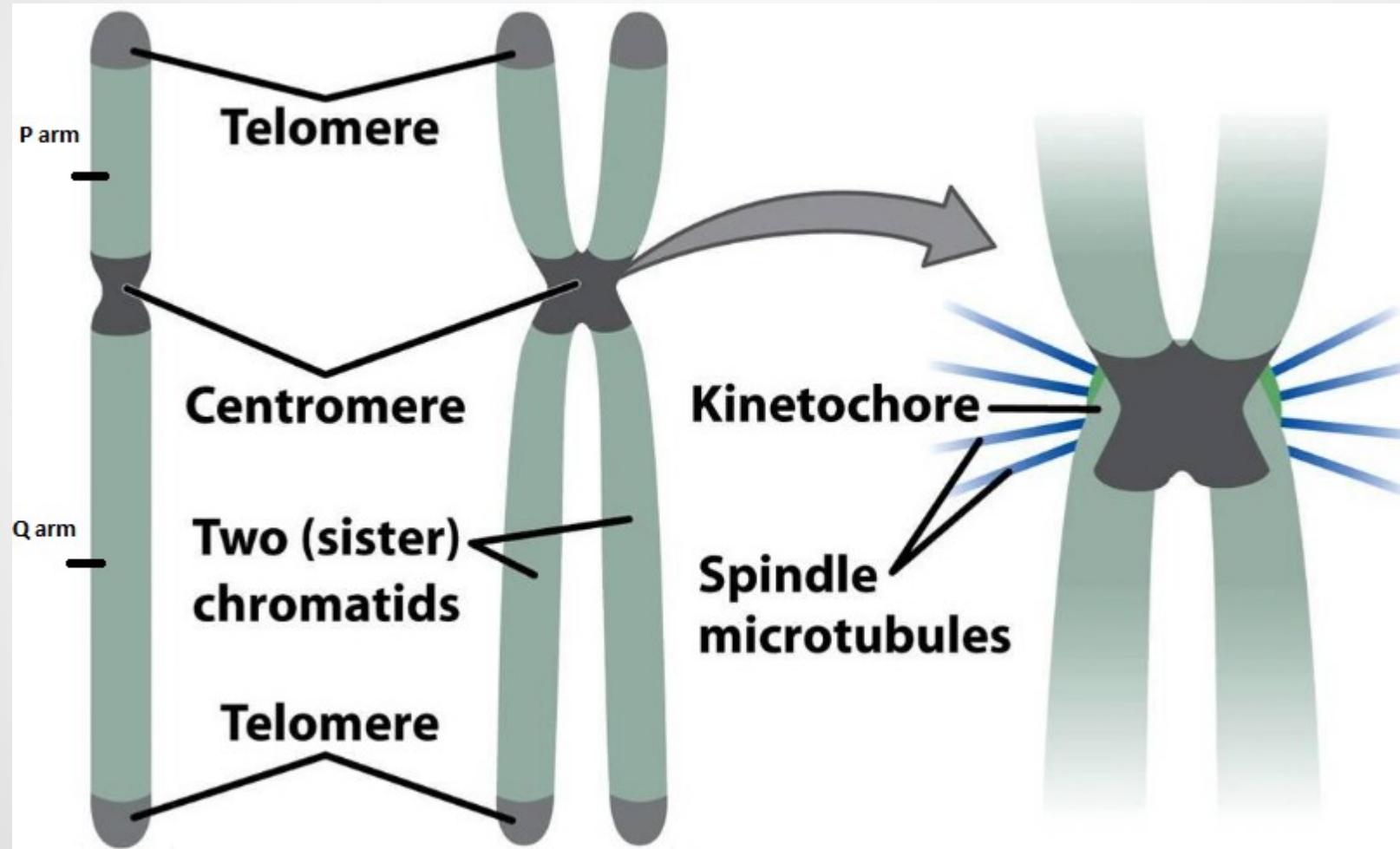
Nucleosome



147 bp per nucleosome



Chromosome structure



Chromosome types



Satellite



Metacentric

Sub-Metacentric

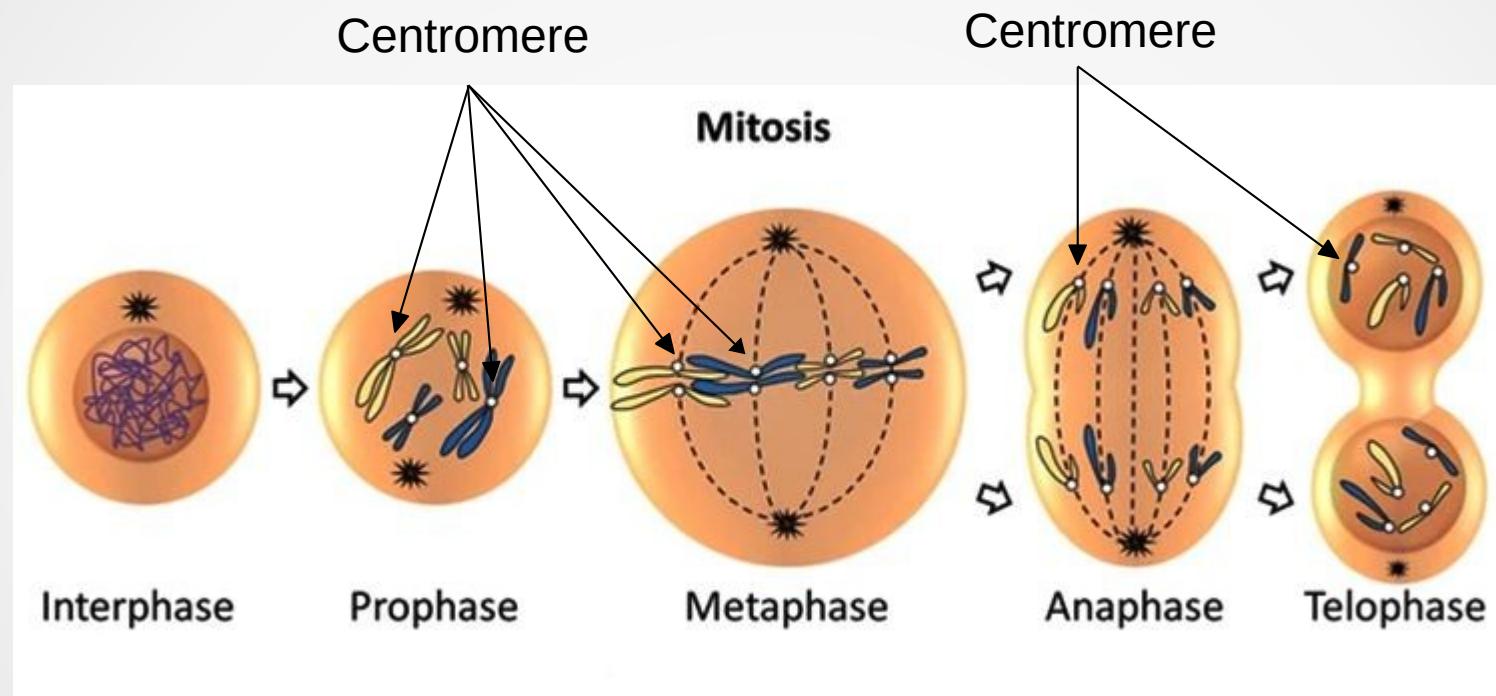


Acrocentric



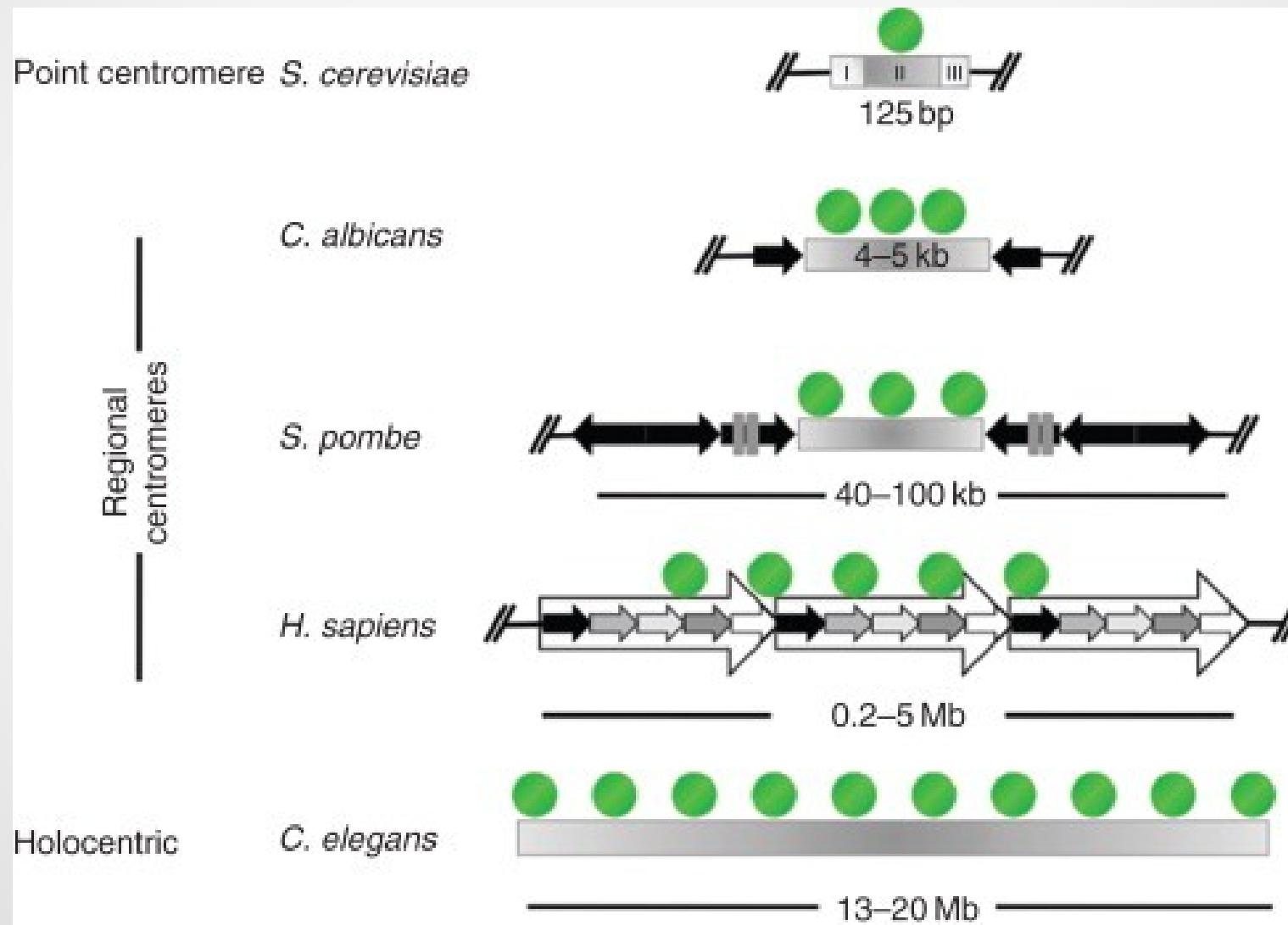
Telocentric

Role of centromere

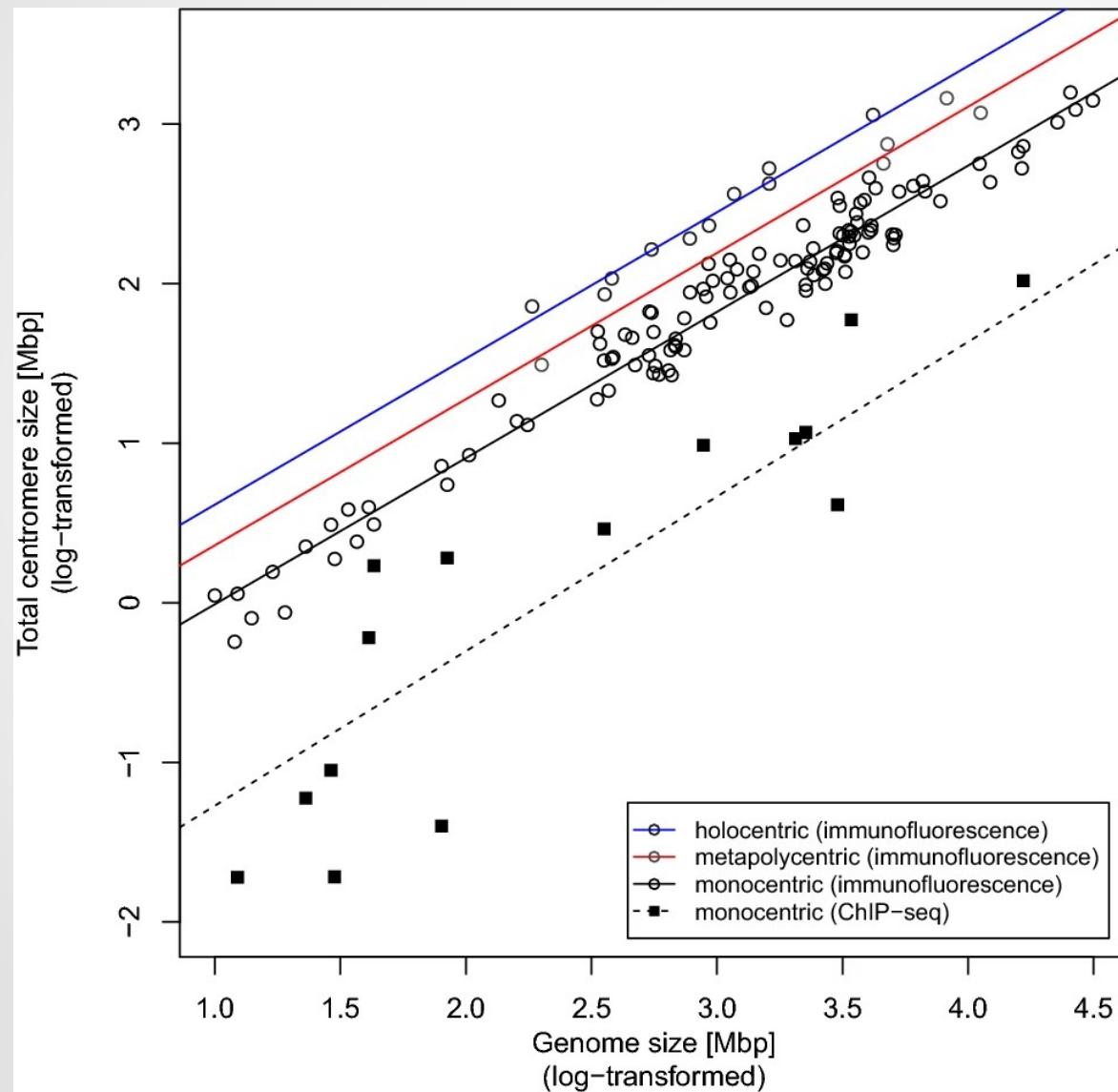


Centromeres are crucial for proper separation of chromosomes between daughter cells

Centromere structure and length



Centromere length and genome size



Plačková et al, 2021

Telomers are:

- located at the end of linear chromosomes
- highly repetitive regions
- protection system against end-replication losses
- protection system against interchromosomal fusions
- stabilizers of genome
- marker of cell age
- maintained by telomerase (or by similar mechanism) but
not in all cells

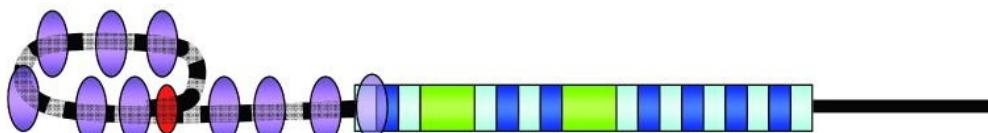
Variants of telomere structure

Drosophila



Cap	<i>HeT-A, TART, TAHRE</i> (HTT) array	Telomere associated sequence	Euchromatin
HOAP	H3Me3K9	PC	
HP1	H3Me3K4	E(Z)	
	Z4	PH	
	JIL-1	PSC	
	HP1	SCM	
	PROD	H3Me3K27	

Mammalian



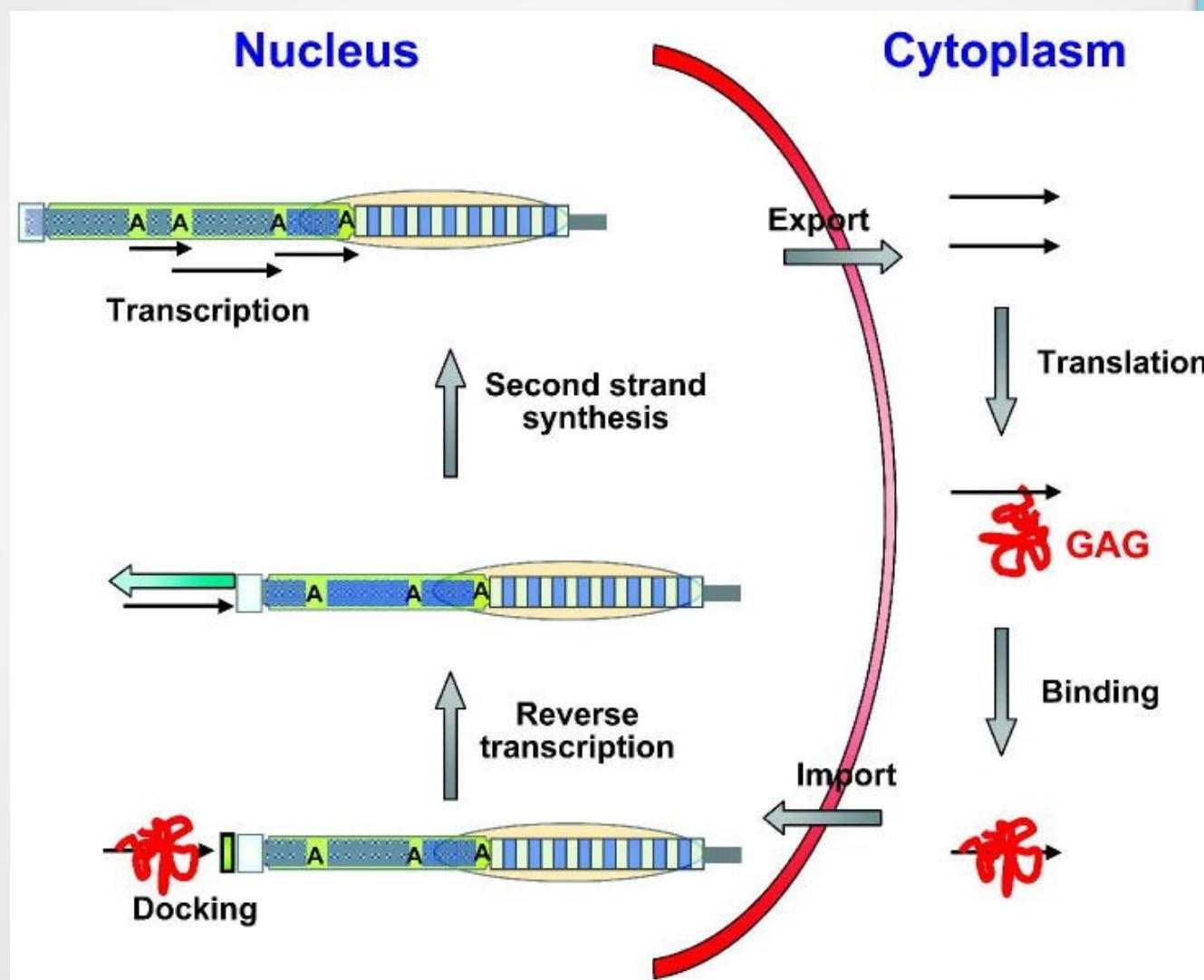
Shelterin	Telomere associated sequence	Euchromatin
TRF1		
TRF2	HP1a	HP1a
RAP1	HP1b	HP1b
TIN2	H3Me3K9	H3Me3K9
POT1	H4Me3K20	H4Me3K20
TPP1	Ku70/Ku80?	

Array of mobile elements

Tandem telomeric repeats

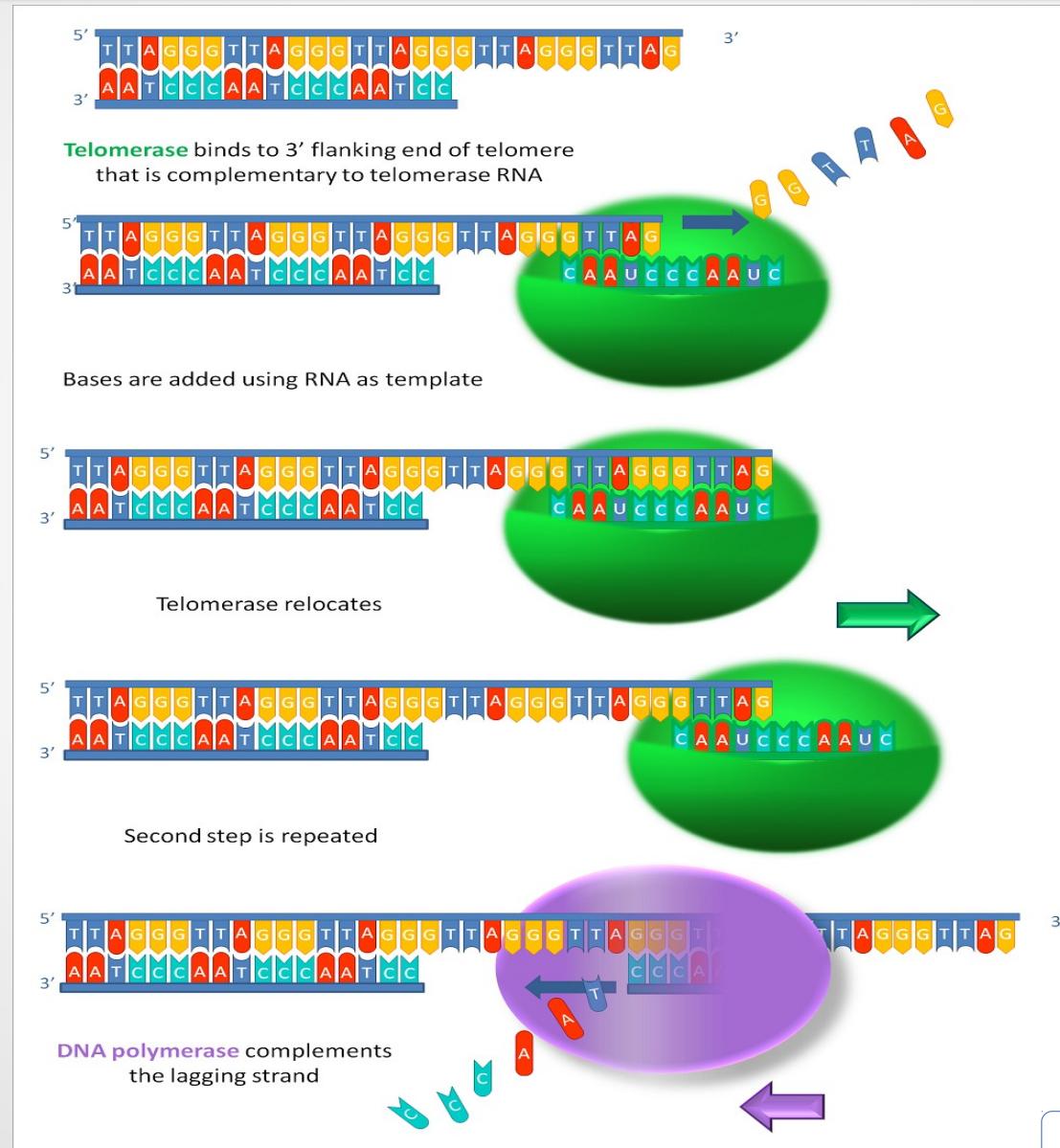
Mason et al, 2008

Elongation of Drosophila telomeres



Mason et al, 2008

Telomere elongation

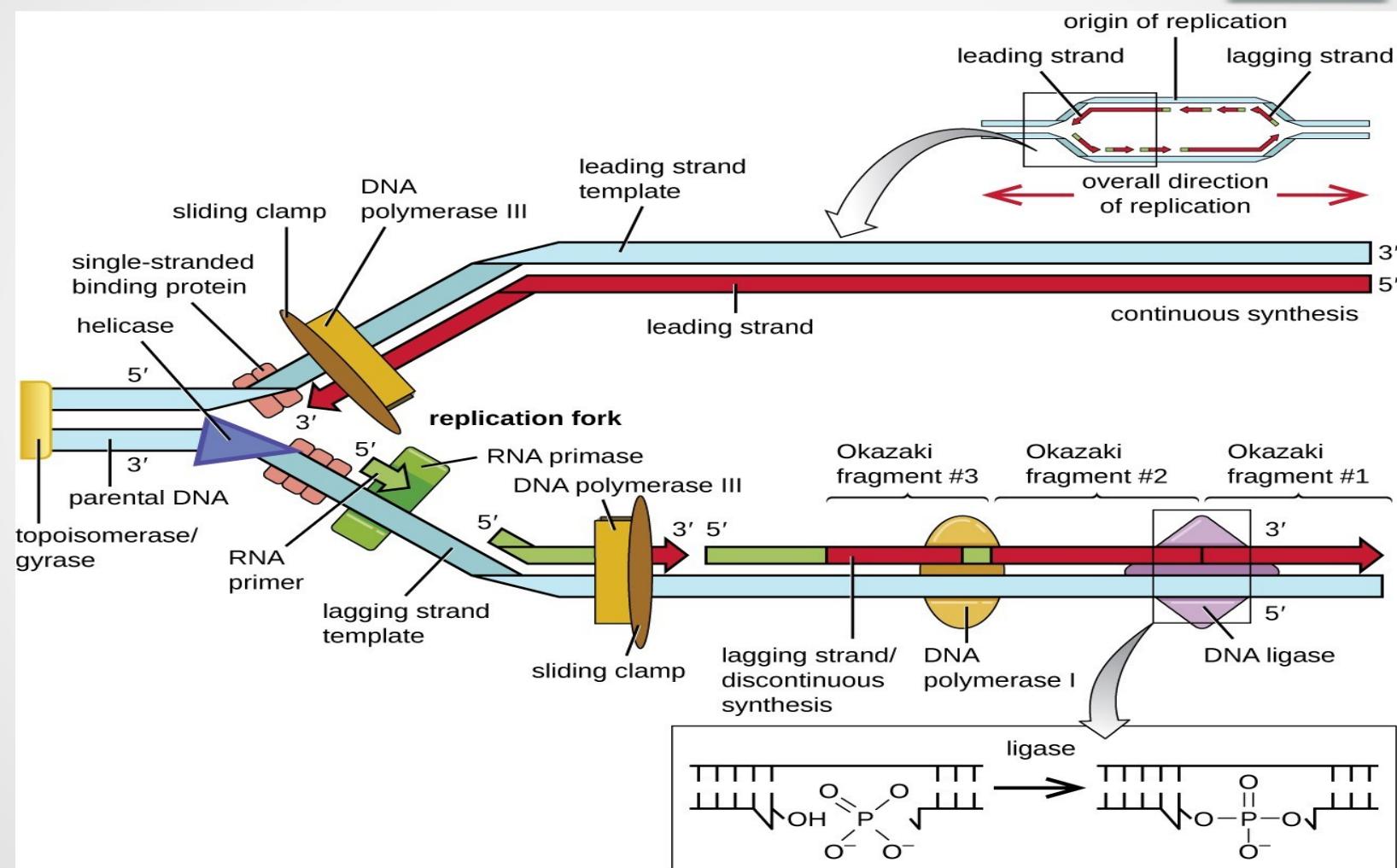


Telomerase. Wikipedia.

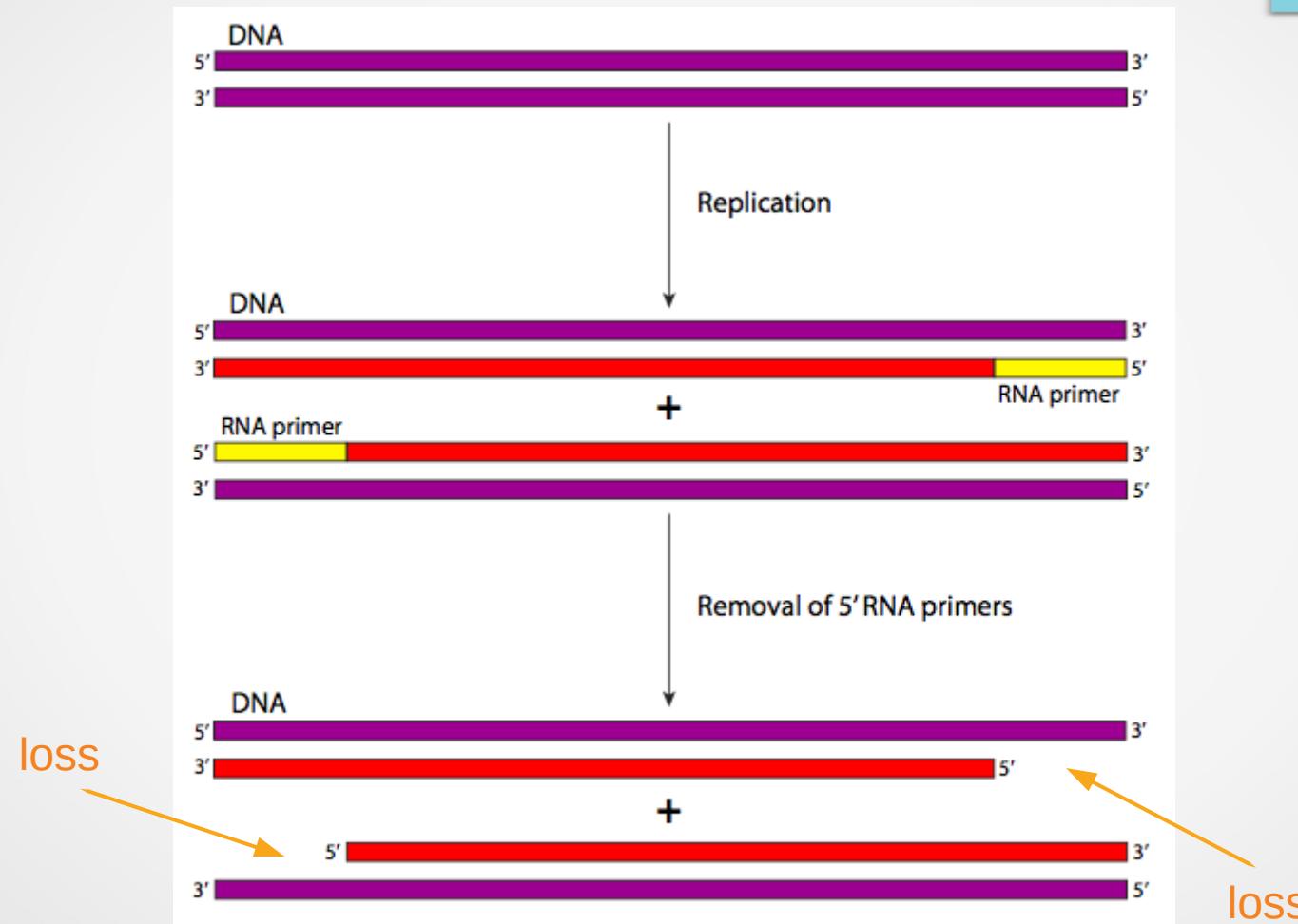
Telomeric repeats

Group	Organism	Repeat monomer
Vertebrates	Human, mouse, <i>Xenopus</i>	TTAGGG
Filamentous fungi	<i>Neurospora crassa</i>	TTAGGG
Slime moulds	<i>Physarum, Didymium</i>	TTAGGG
	<i>Dictyostelium</i>	AG(1-8)
Kinetoplastid protozoa	<i>Trypanosoma, Crithidia</i>	TTAGGG
Ciliate protozoa	<i>Tetrahymena, Glaucoma</i>	TTGGGG
	<i>Oxytricha, Stylonychia, Euplates</i>	TTTTGGGG
Apicomplexan protozoa	<i>Plasmodium</i>	TTAGGG(T/C)
Higher plants	<i>Arabidopsis thaliana</i>	TTTAGGG
Green algae	<i>Chlamydomonas</i>	TTTTAGGG

Replication



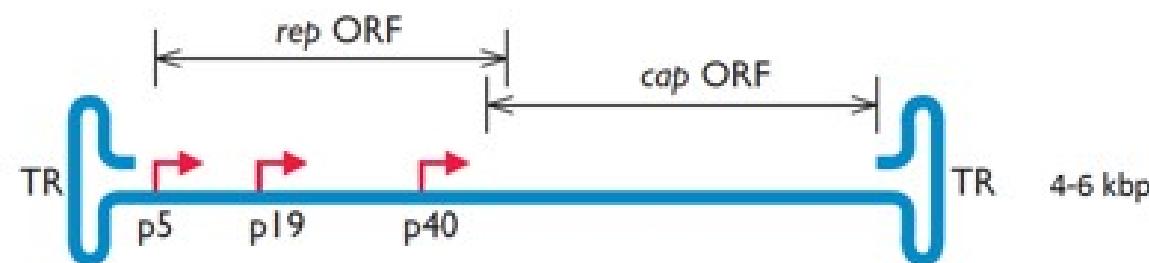
End replication problem



Specific problem of linear chromosome replication

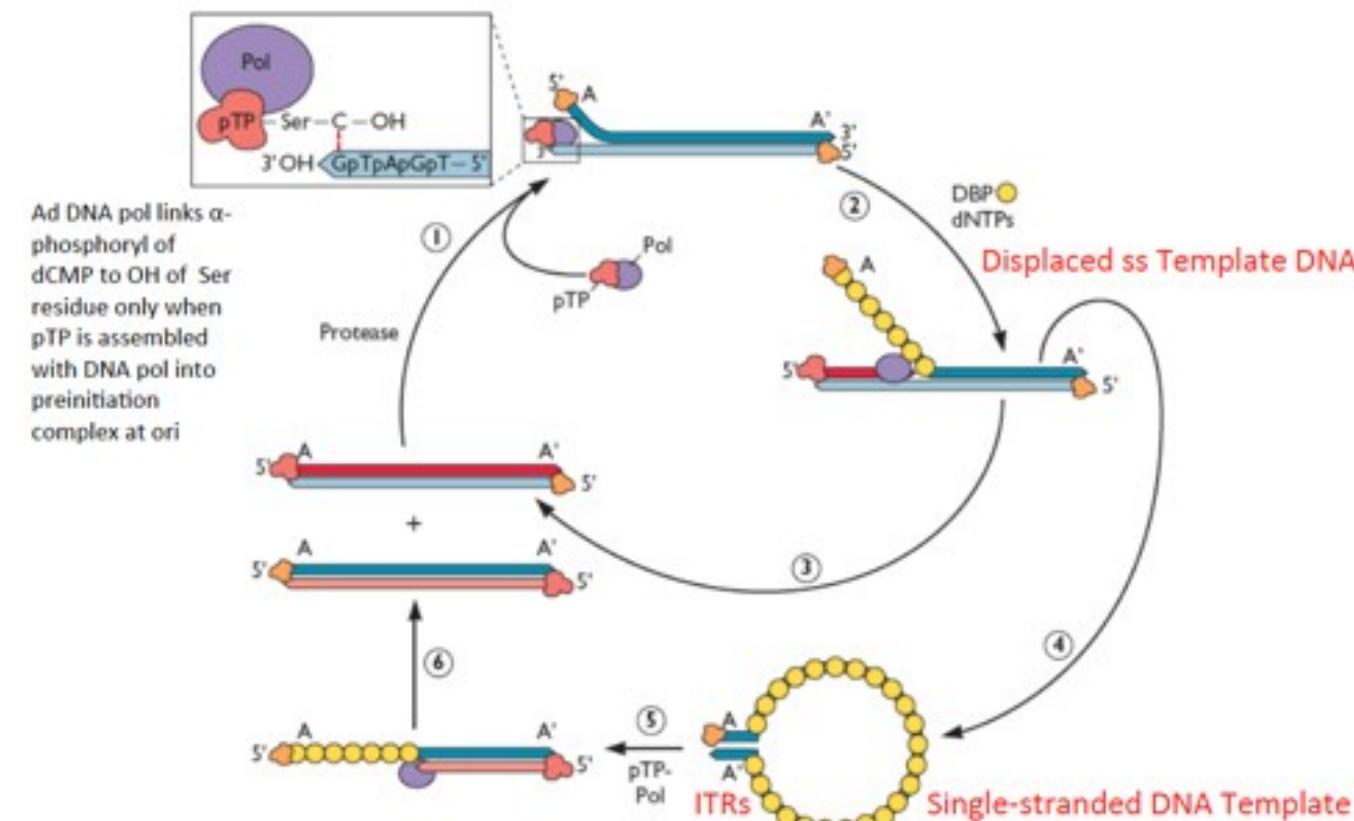
Viral solutions of end problem: parvoviruses (1)

DNA priming: Parvoviruses

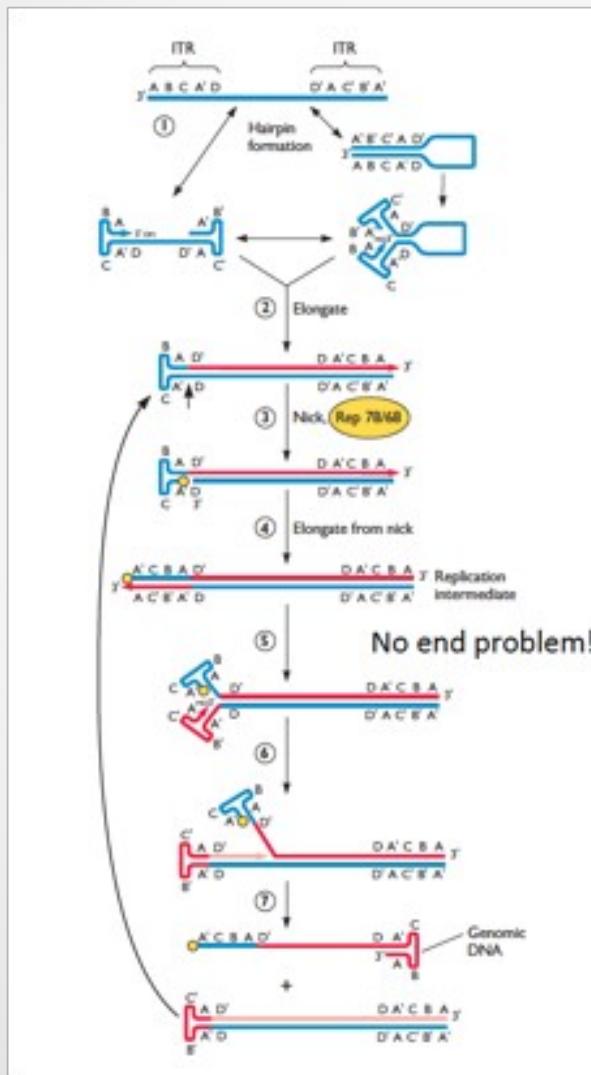


Viral solution of end problem: adenoviruses

Protein Priming

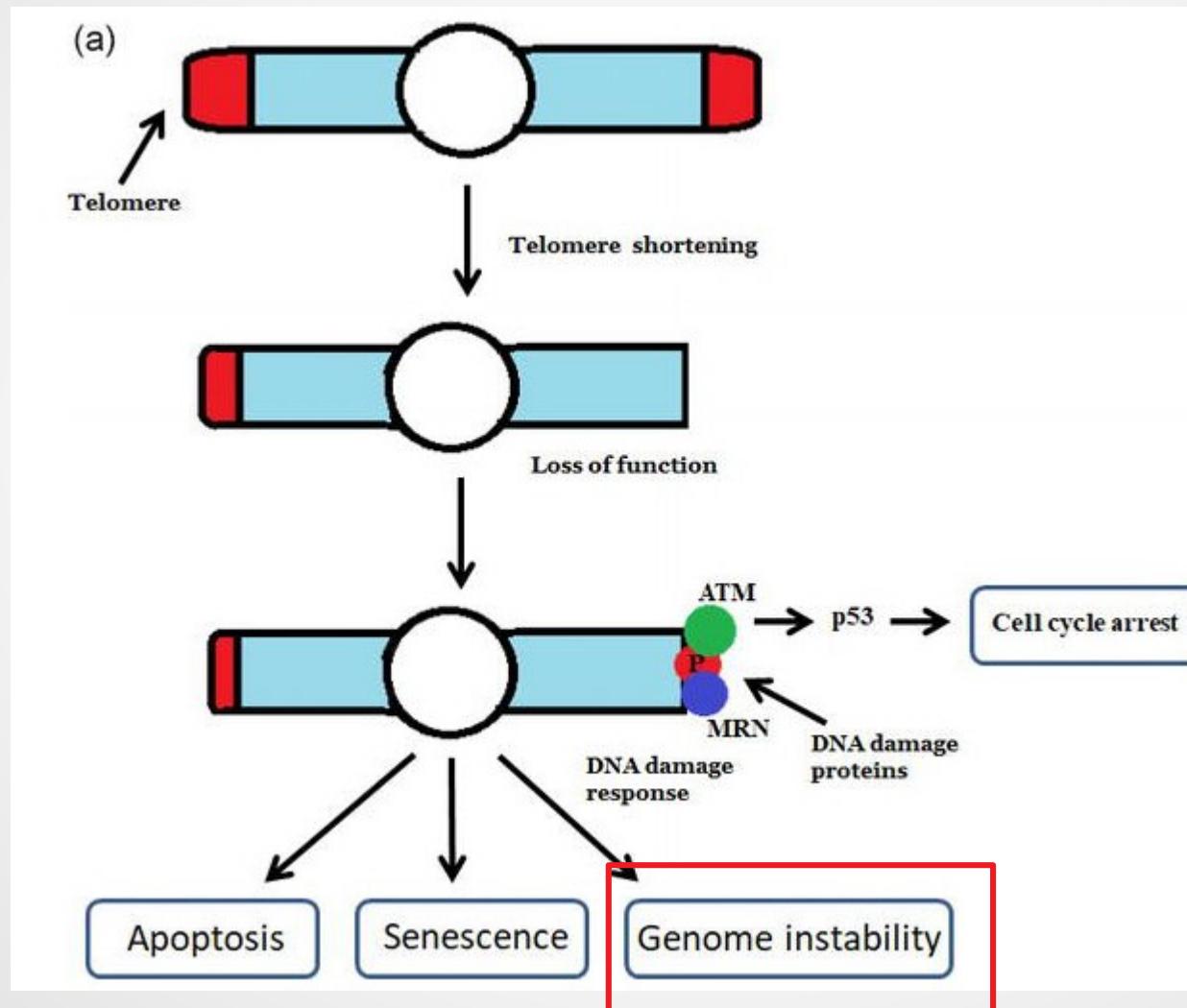


Viral solution of end problem: parvoviruses (2)



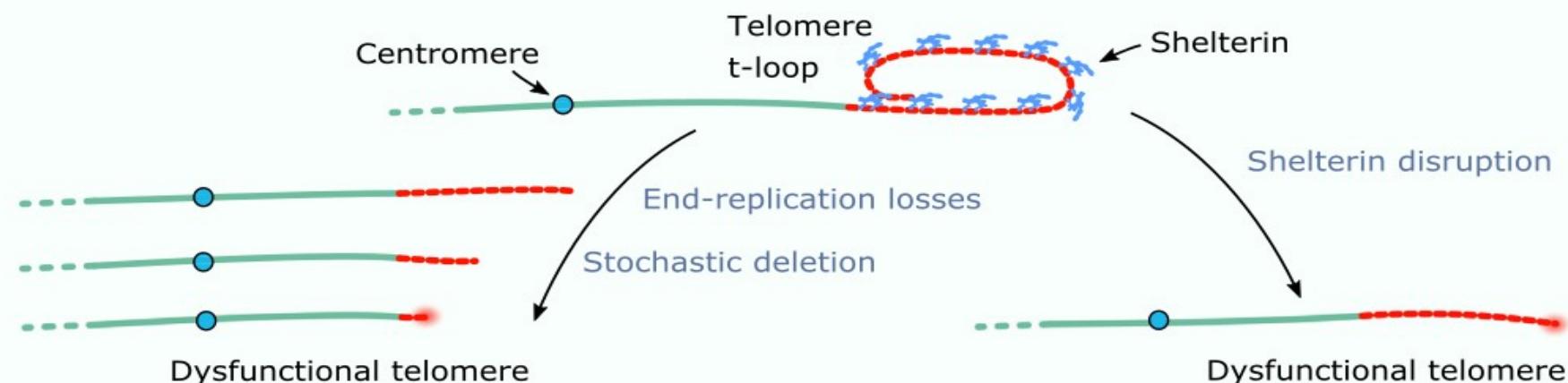
- Replication is continuous
- No pol α , uses ITR to self-prime
- Requires pol δ , Rf-C and PCNA
- Rep78/68 proteins are required for initiation and resolution: endonuclease, helicase, binds 5' terminus
- No replication fork, strand displacement

Results of telomere shortening (1)



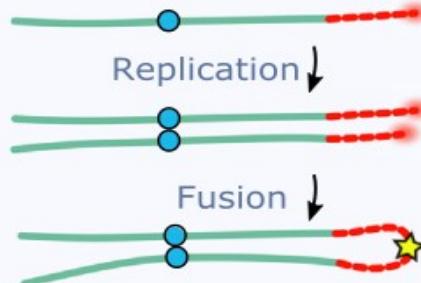
Results of telomere shortening (2)

(A) Telomere shortening and dysfunction

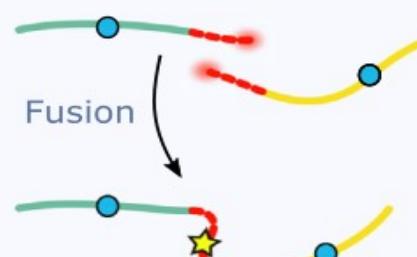


(B) Telomere fusion and dicentric chromosome formation

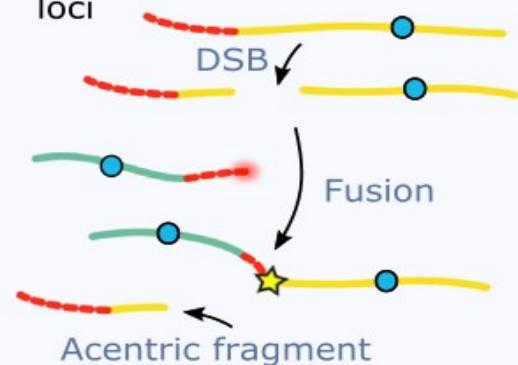
i Sister chromatid fusion



ii End-to-end fusion with other chromosome



iii Fusion with nontelomeric loci



Cleal and Baird, 2020

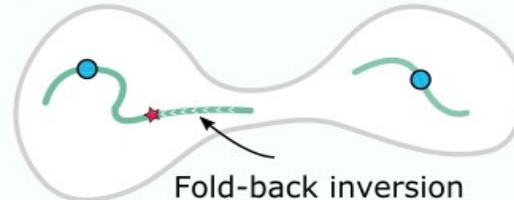
Results of telomere shortening (3)

(A) Chromatin bridge formation

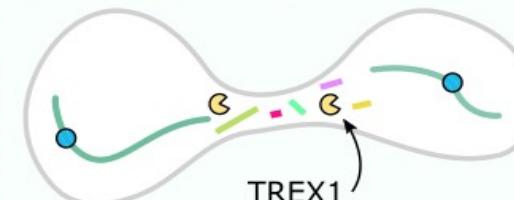


(B) Dicentric resolution

i Simple break



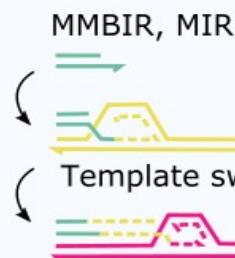
ii Complex break



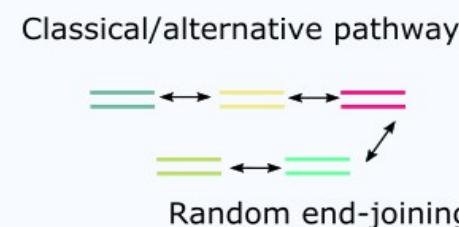
Mechanisms of fragmentation:
nuclease attack,
mechanical force,
others?

(C) Chromosome repair

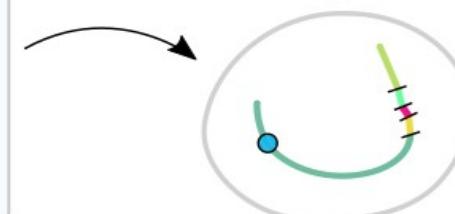
i Replicative repair



ii NHEJ



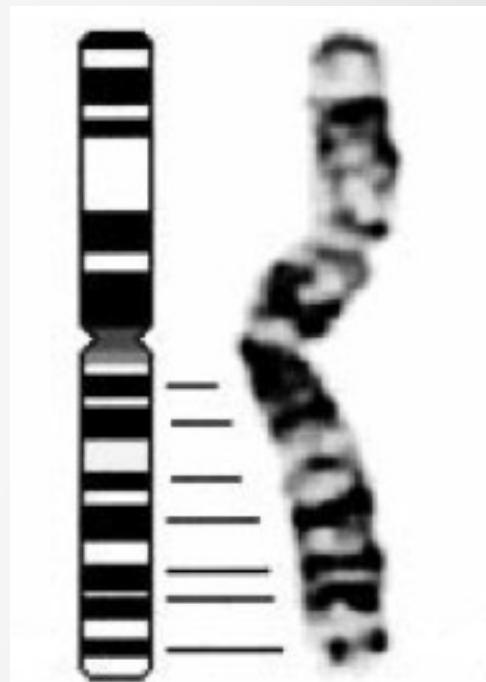
(D) Complex genome rearrangements



Can we see the chromosomes?

Yes, in microscope, but we need... to stain them first!

Staining type	Light bands	Dark bands
G-Banding	GC-rich	AT-rich
R-Banding	AT-rich	GC-rich
C-banding	euchromatine	constitutive heterochromatine



G-banded chromosome (right)
and
its ideogram (left)

Karyotype description

Karyotyping is

process of visualizing and evaluating the karyotype, i.e. all set of chromosomes from the cell

Ploidy is

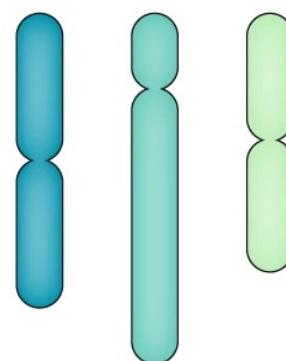
the number of complete sets of chromosomes in a cell.

Somatic number (SN) is

number of all chromosomes in the cell

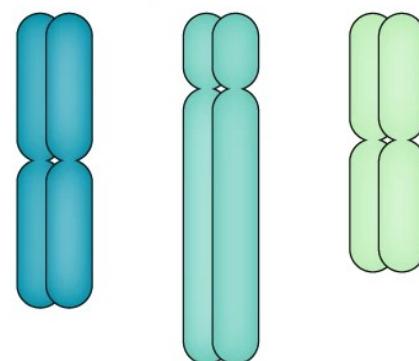
sex
chromosomes
 $2n = 32, XX$
ploidy
SN

haploid



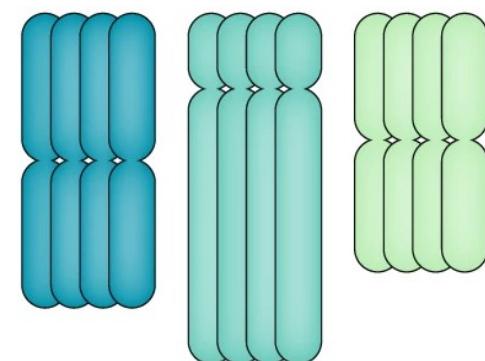
$1n = 3$

diploid



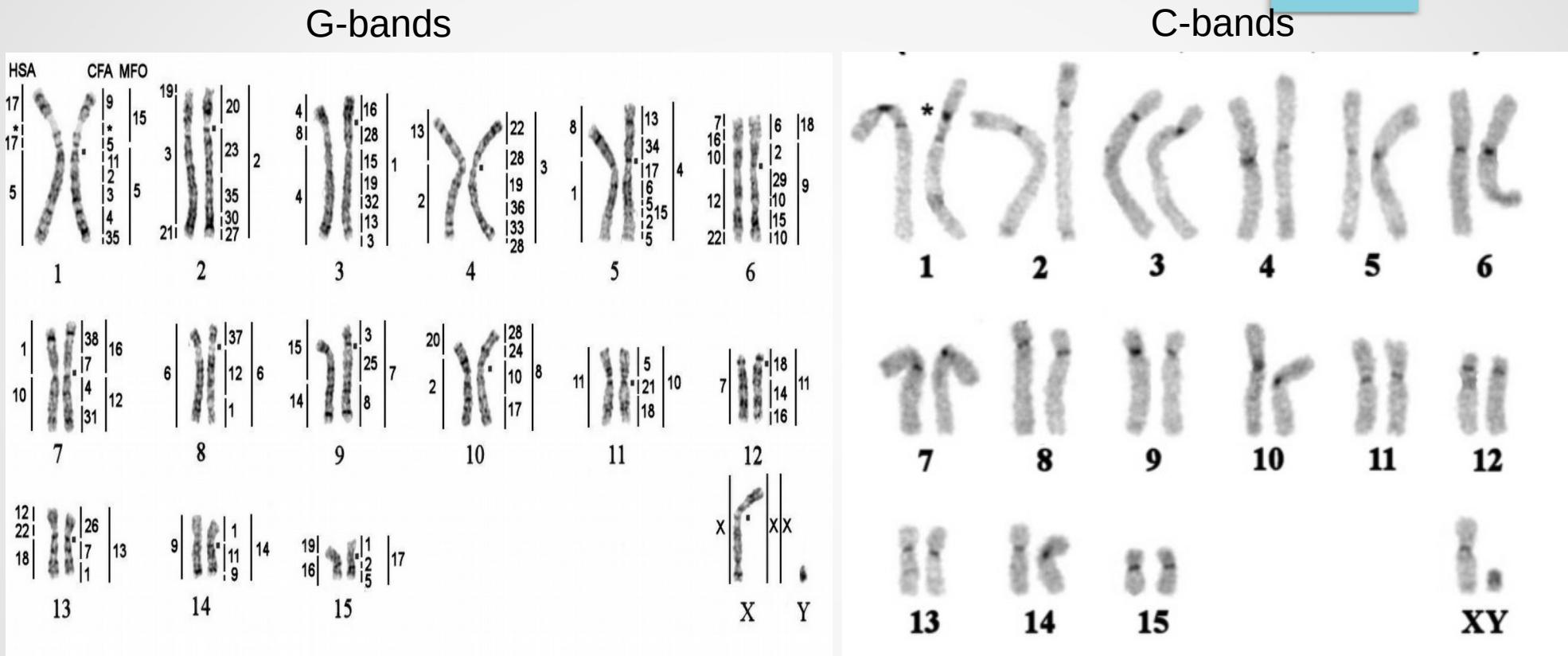
$2n = 6$

triploid



$3n = 9$

Examples of karyotypes (1)



Beklemisheva et al, 2016

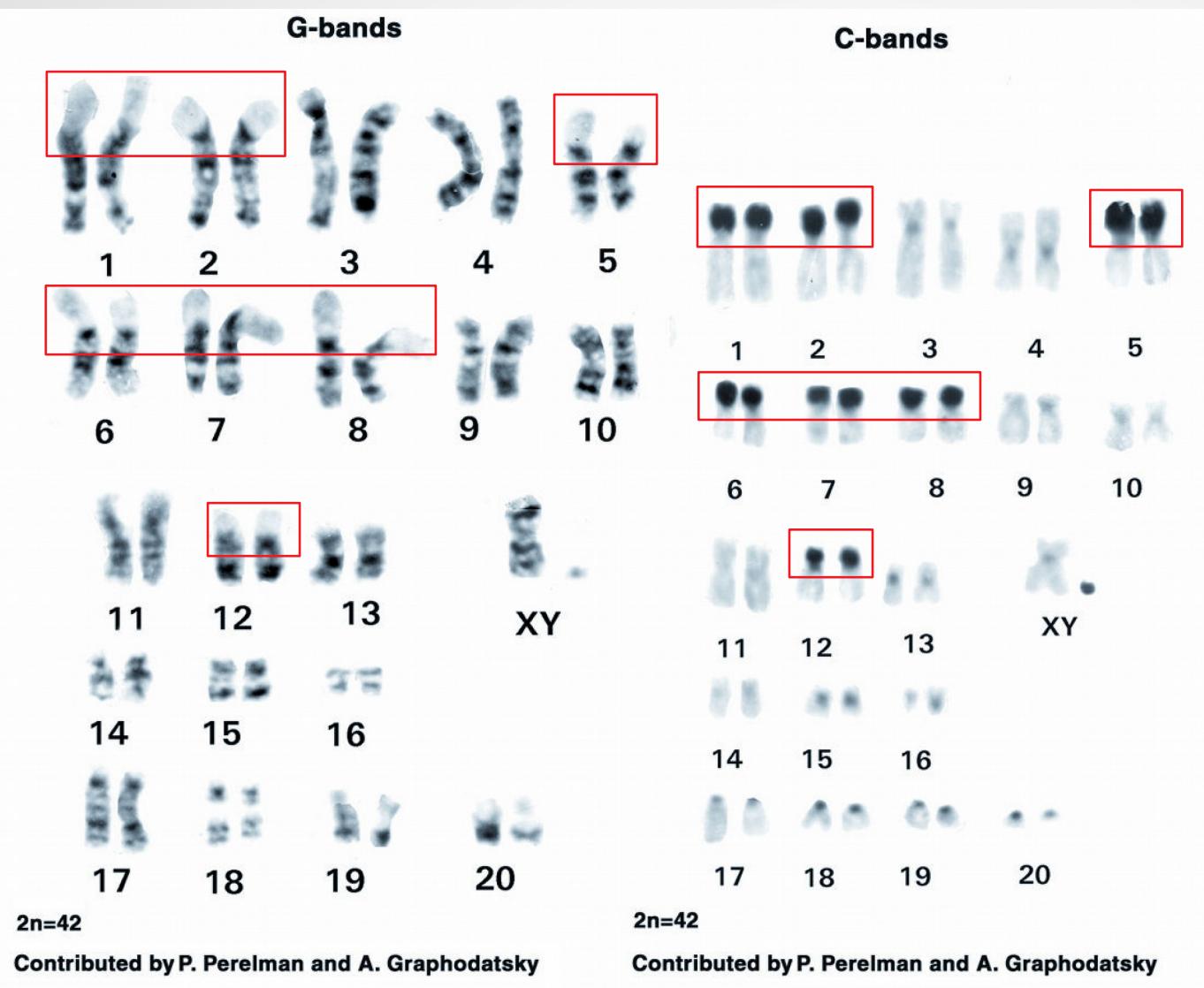
Pusa sibirica
Baikal seal



Beklemisheva et al, 2020

~ 2.2-2.3 Gbp

Examples of karyotypes (2)



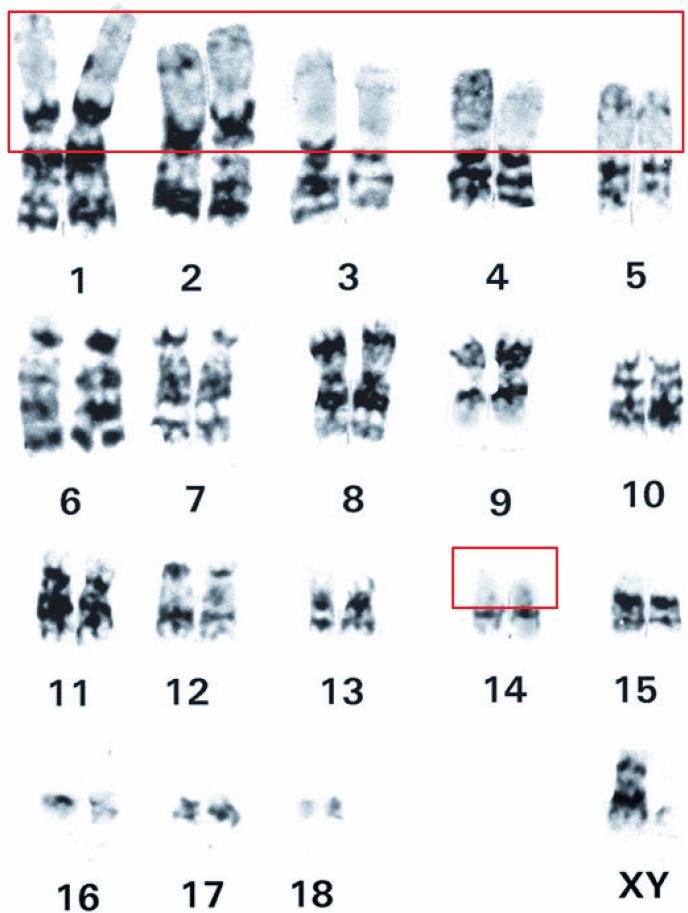
Mustela nivalis
the least weasel



~ 3.0 - 3.2 Gbp

Examples of karyotypes (3)

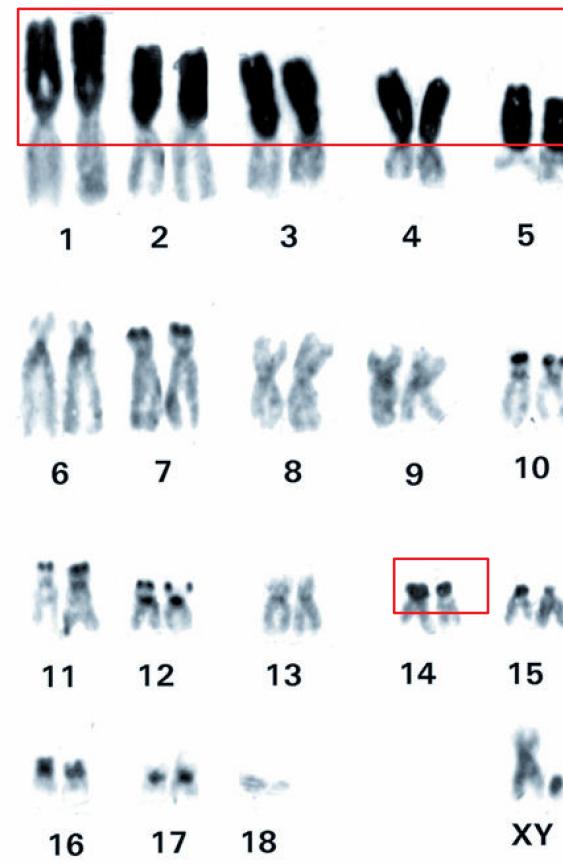
G-bands



2n=38

A. Graphodatsky (unpublished)

C-bands

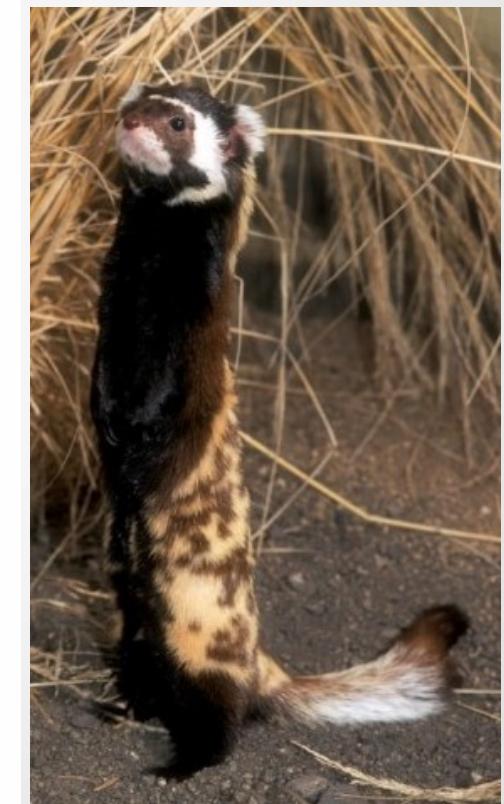


2n=38

A. Graphodatsky (unpublished)

Vormela peregusna

Marbled polecat



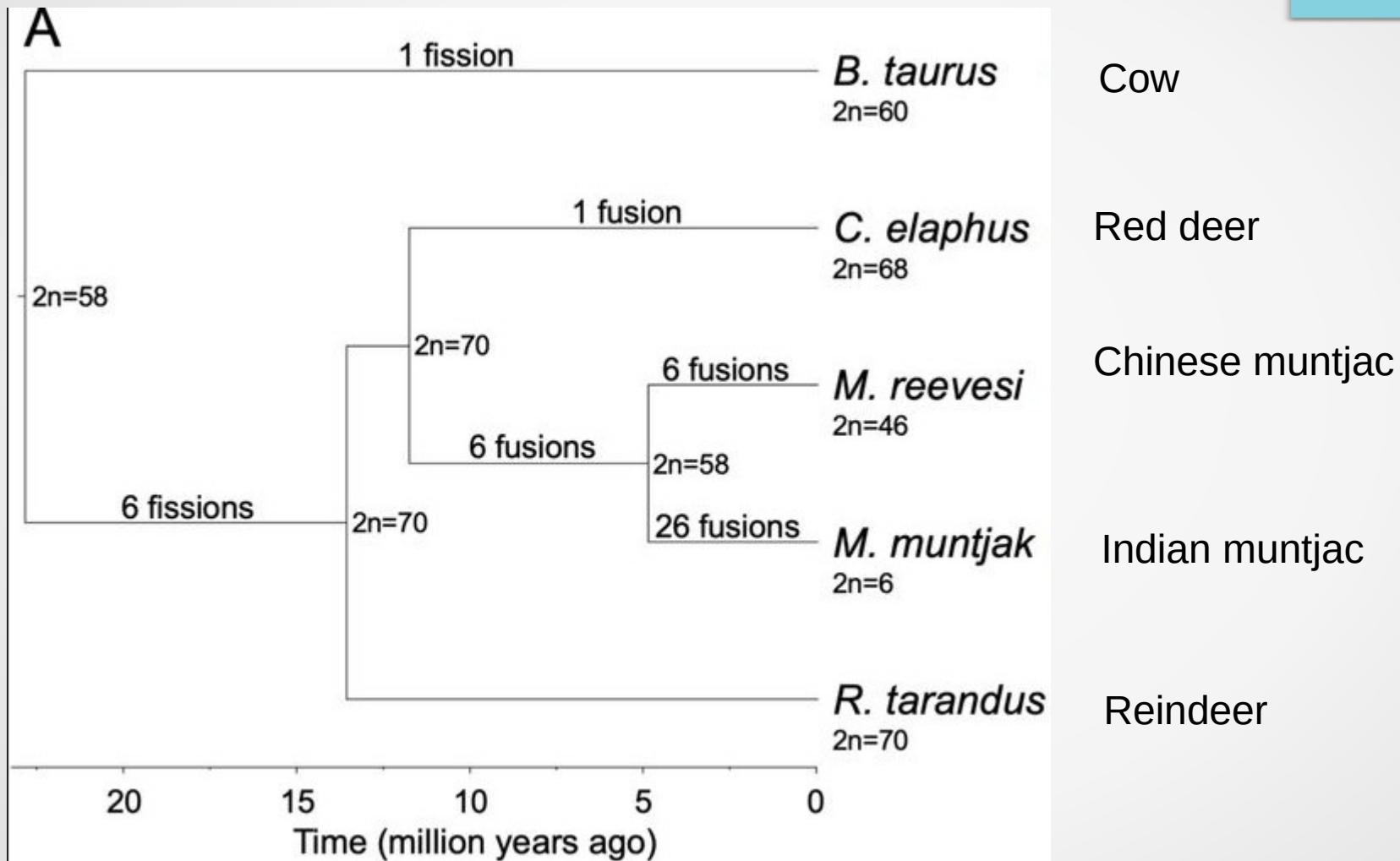
How stable are karyotypes? Muntjac case



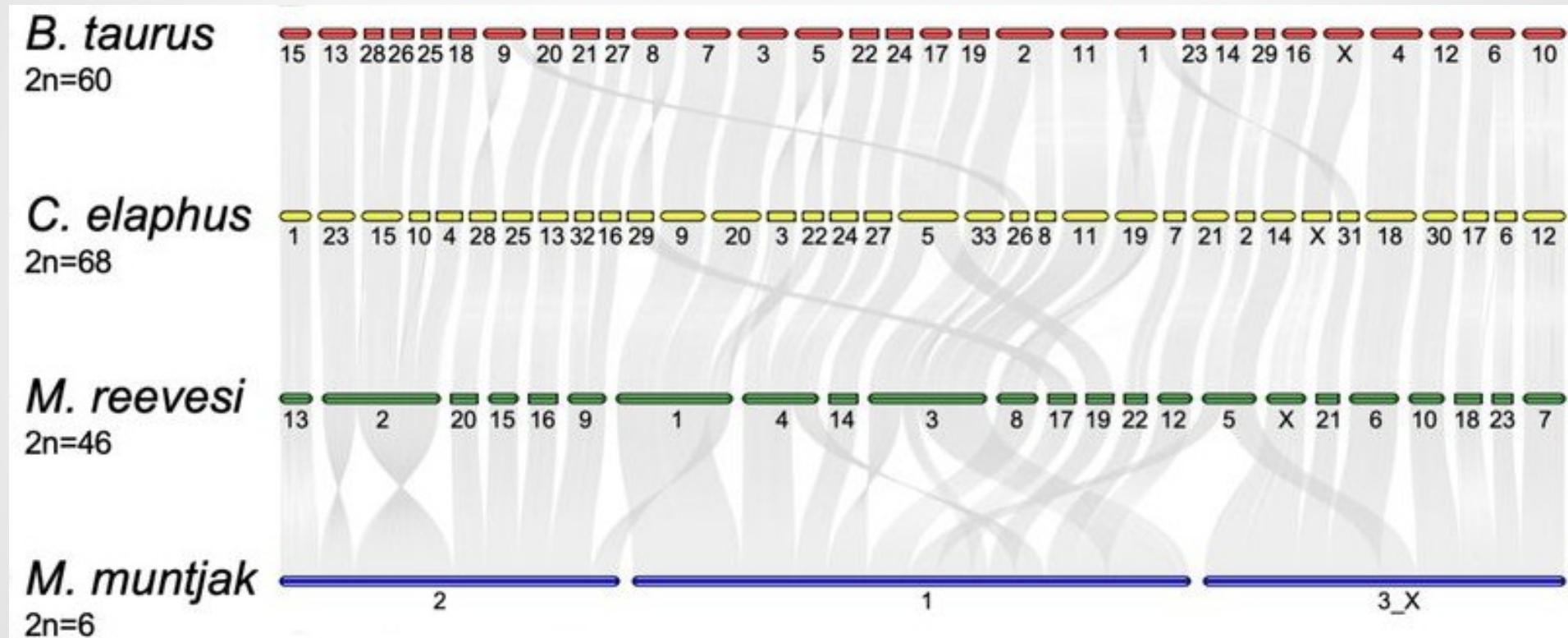
© Roland Wirth, ZGAP

Indian muntjac (*Muntiacus muntjac*)

How stable are karyotypes? Muntjac case (1)

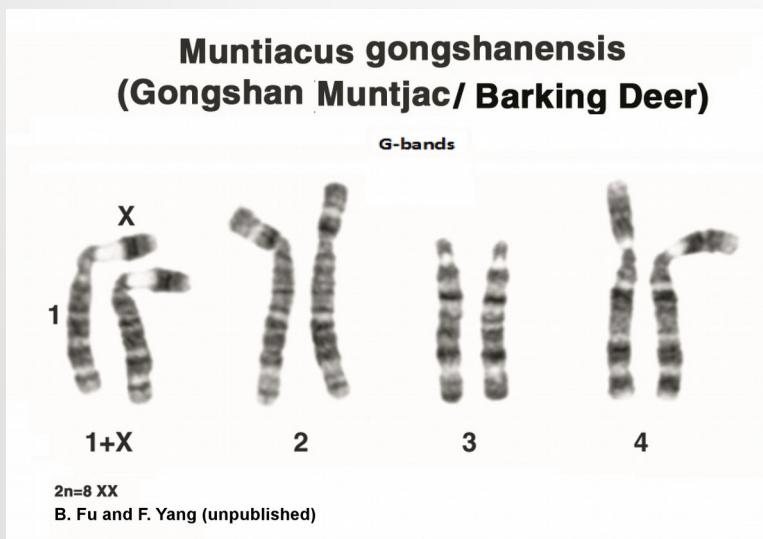
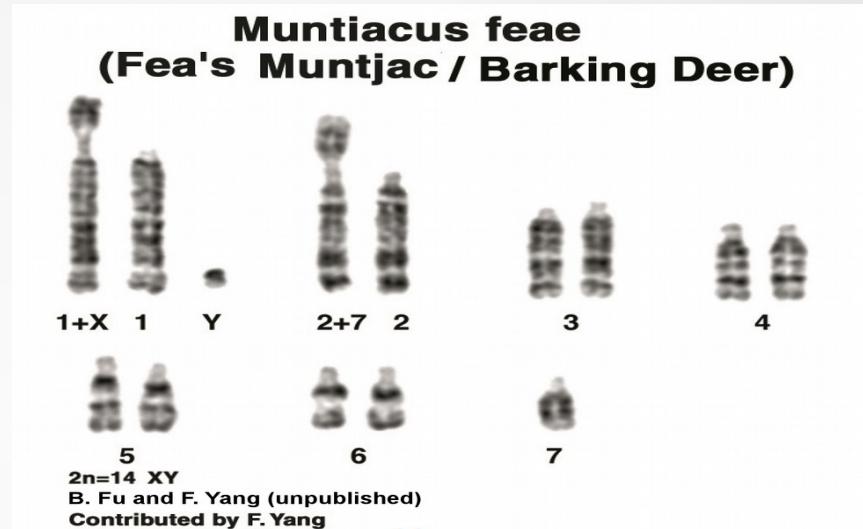
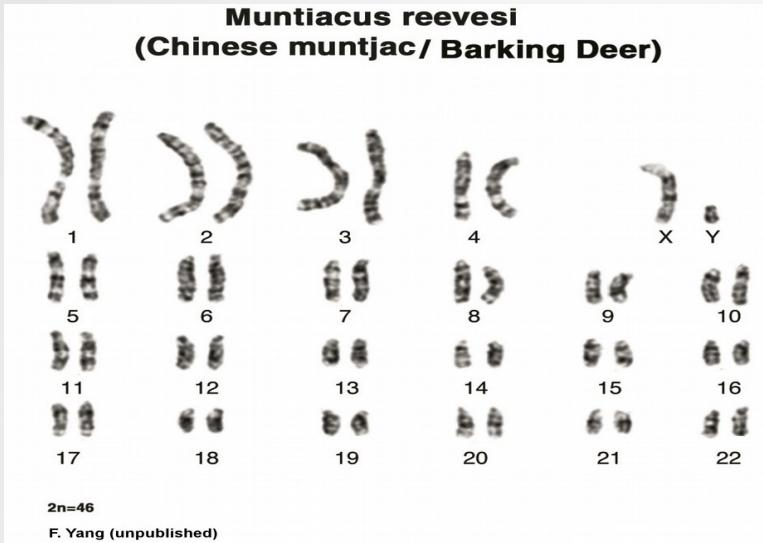


How stable are karyotypes? Muntjac case (2)



Mudd et al, 2020

How stable are karyotypes? Muntjac case (3)

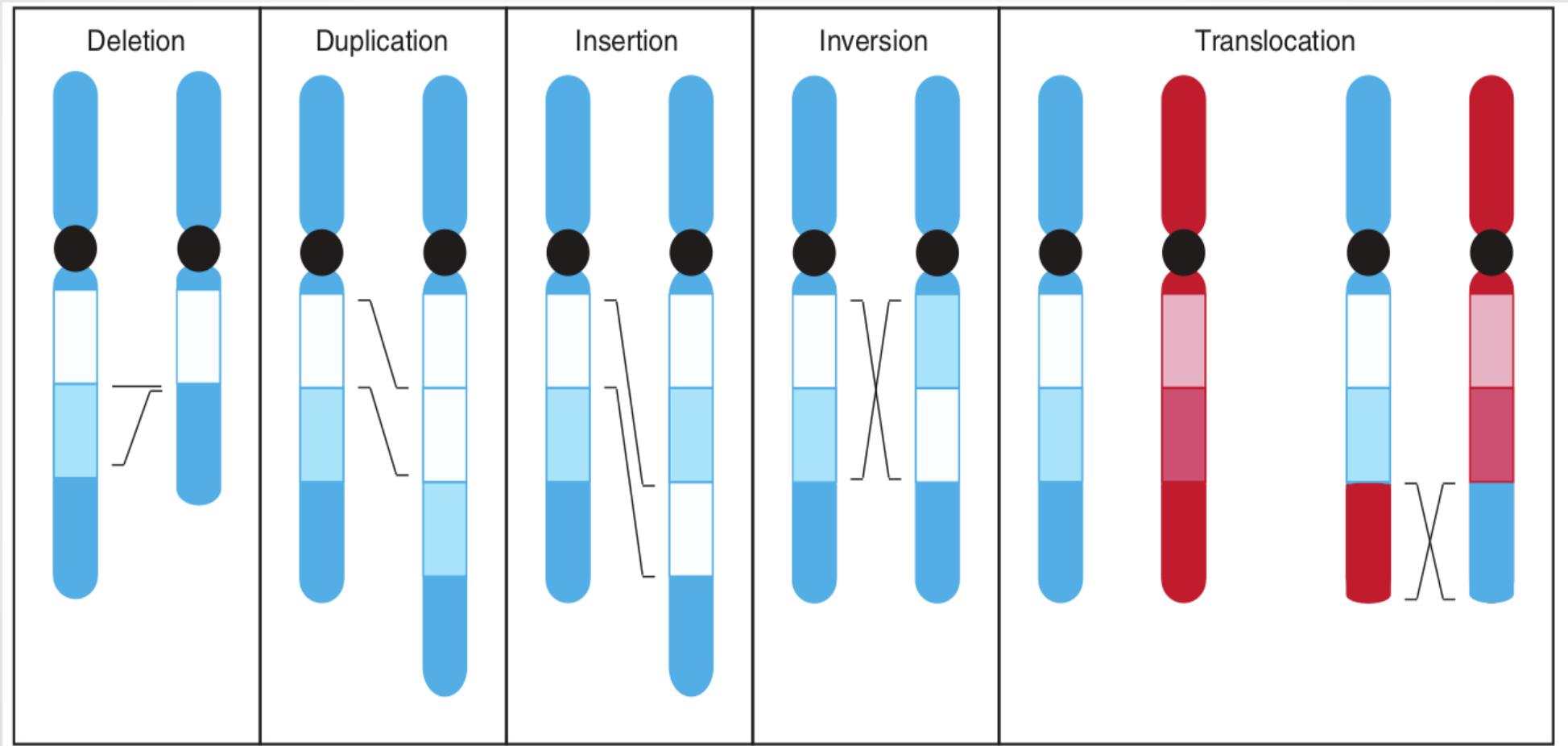


Chinese muntjak	2n=46
Fea's muntjac	2n=14
Gongshan muntjac	2n=8
Indian muntjak	2n=6

Same genus, but very different number of chromosome

Graphodatsky et al, 2020

Basic types of genome rearrangements



How stable are genomes? Case of sturgeons



© Lubomir Hlasek
www.hlasek.com
Acipenser baerii 7/92

Siberian sturgeon (*Acipenser baerii*)

Ploidy of different sturgeon species

Species	Chromosome number	Ploidy level		Reference
		“evolutionary scale”	“recent scale”	
<i>Polyodon spathula</i> (Walbaum)	120	4	2	Dingerkus and Howell, 1976
<i>Scaphirhynchus platorynchus</i> (Rafinesque)	~112	4	2	Ohno et al., 1969
<i>Acipenser sturio</i> Linnaeus	116 ± 4	4	2	Fontana and Colombo, 1974
<i>A. nudiventris</i> Lovetsky	118 ± 2	4	2	Arefjev, 1983; Sokolov and Vasil'ev, 1989a
<i>A. ruthenus</i> Linnaeus	118 ± 2	4	2	Fontana et al., 1975; Vasil'ev, 1985; Birstein and Vasil'ev, 1987
	118 ± 4	4	2	Ráb, 1986; Fontana, 1994
<i>A. stellatus</i> Pallas	118 ± 2	4	2	Vasil'ev, 1985; Birstein and Vasil'ev, 1987
<i>A. oxyrinchus</i> Mitchell	121 ± 3	4	2	Fontana et al., 2008b
<i>A. huso</i> Linnaeus	116 ± 4	4	2	Fontana and Colombo, 1974; Vasil'ev, 1985
<i>A. gueldenstaedtii</i> Brandt and Ratzeburg	250 ± 8	8	4	Vasil'ev, 1985; Vlasenko et al., 1989
<i>A. persicus</i> Borodin	~258	8	4	Nowruzfashkhami et al., 2000
<i>A. baerii</i> Brandt	249 ± 5	8	4	Vasil'ev et al., 1980; Sokolov and Vasil'ev, 1989b
<i>A. naccarii</i> Bonaparte	239 ± 7	8	4	Fontana and Colombo, 1974
<i>A. brevirostrum</i> Lesueur	~372	12	6	Kim et al., 2005
	372 ± 6	12	6	Fontana et al., 2008a
<i>A. transmontanus</i> Richardson	248 ± 8	8	4	Fontana, 1994
	~271	8	4	Van Eenennaam et al., 1998
<i>A. sinensis</i> Gray	264 ± 4	8	4	Yu et al., 1987
<i>A. fulvescens</i> Rafinesque	262 ± 6	8	4	Fonatana et al., 2004
<i>A. mikadoi</i> Hilgendorf	262 ± 4	8	4	Vasil'ev et al., 2008, 2009; this study
<i>A. medirostris</i> Ayres	249 ± 8	8	4	Van Eenennaam et al., 1999
<i>A. schrenskii</i> Brandt	238 ± 8	8	4	Song et al., 1997
	266 ± 4	8	4	This study
<i>A. dauricus</i> Georgi	268 ± 4	8	4	Vasil'ev et al., 2008, 2009; this study

Polyploidization is an increase of ploidy.

Most common cases are
 $2n \rightarrow 4n$
 $2n \rightarrow 3n$

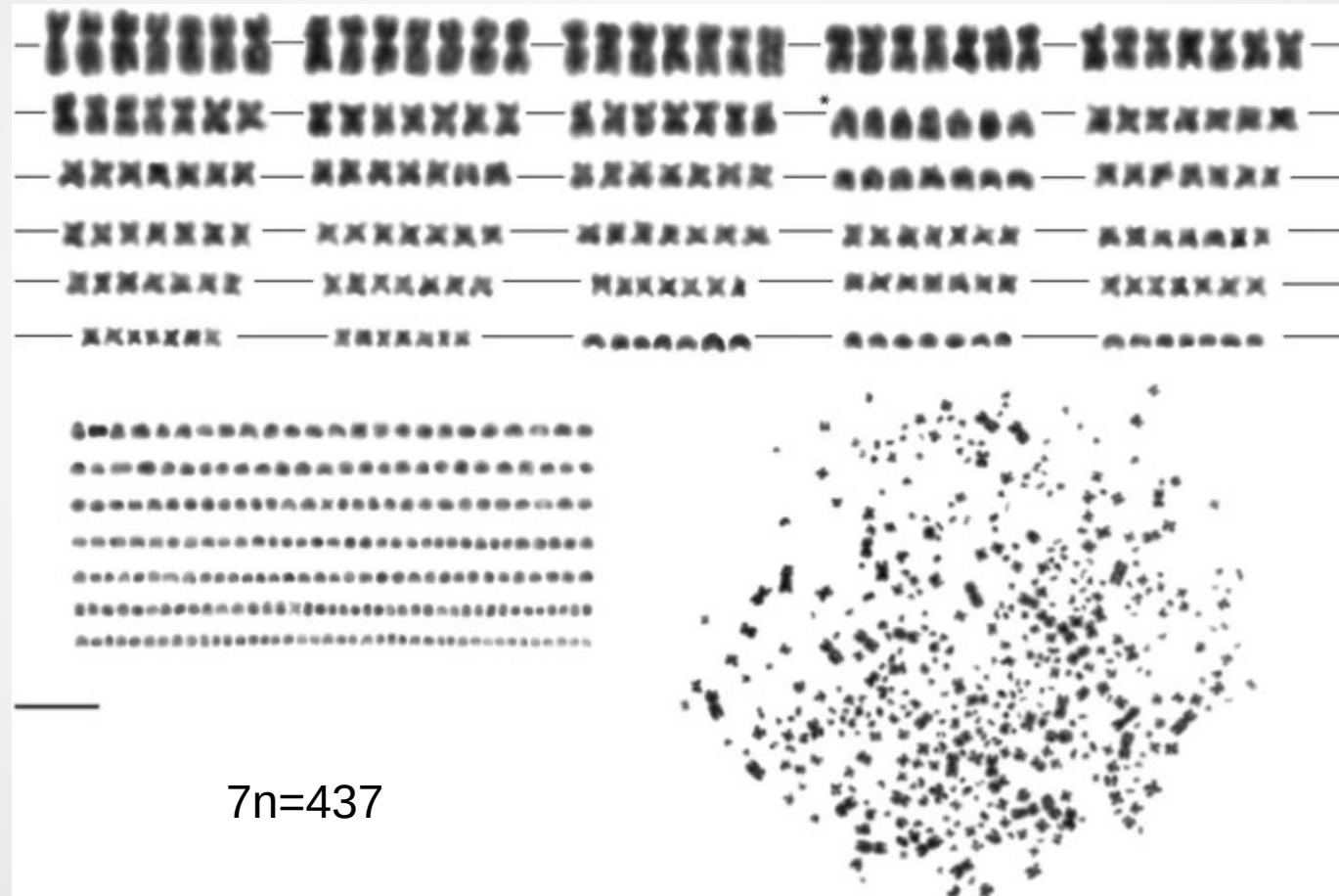
Diploidization is a conversion of polyplloid species to diploid.

Some chromosomes might be lost, the rest diverge.

Vasi'ev et al, 2009

Spontaneous polyploidization

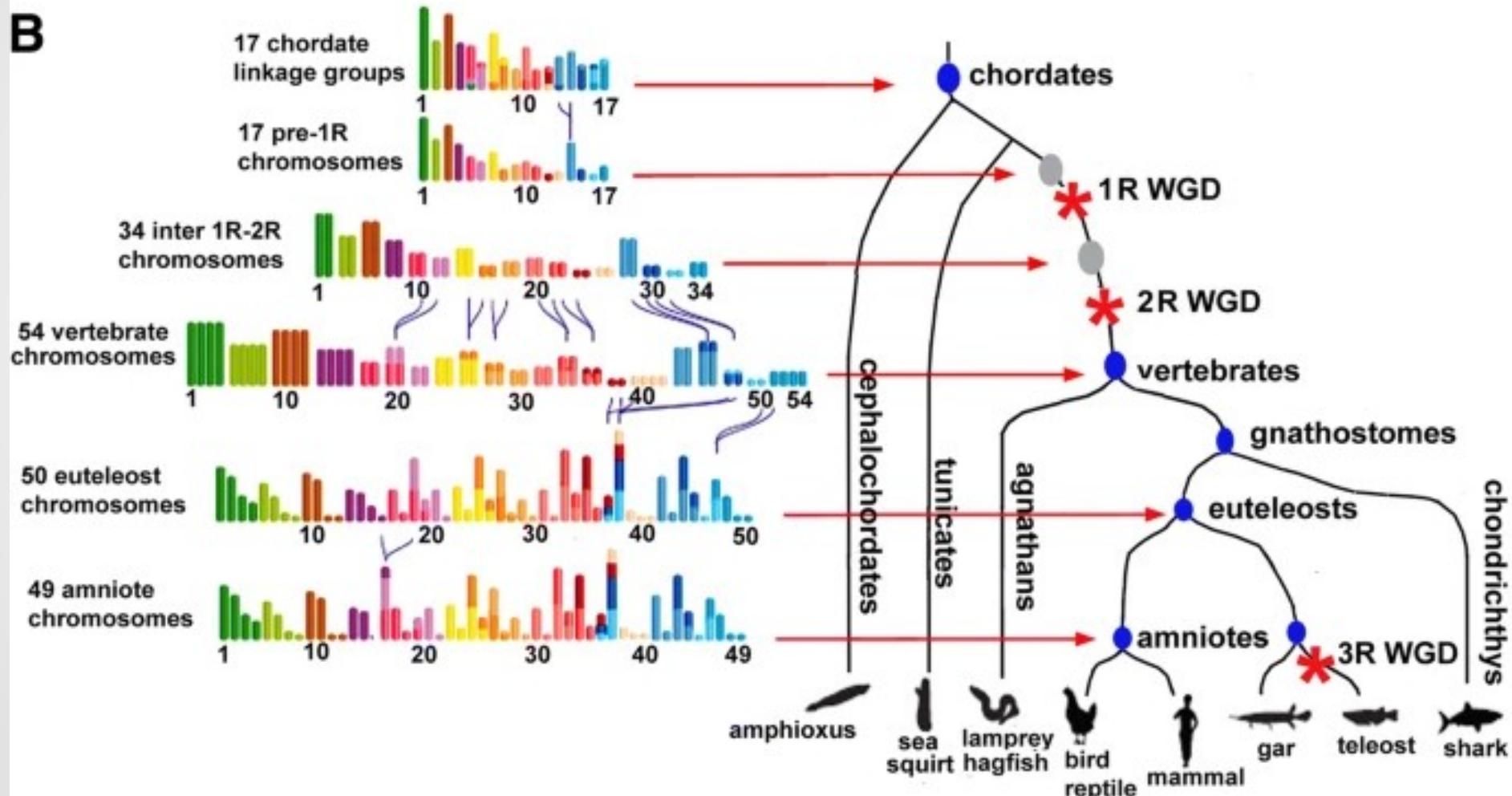
... is common in sturgeons!



Karyotype of heptaploid Siberian Sturgeon individual
Usual karyotype is 4n=250!

Whole genome duplications in evolution

B



Summary

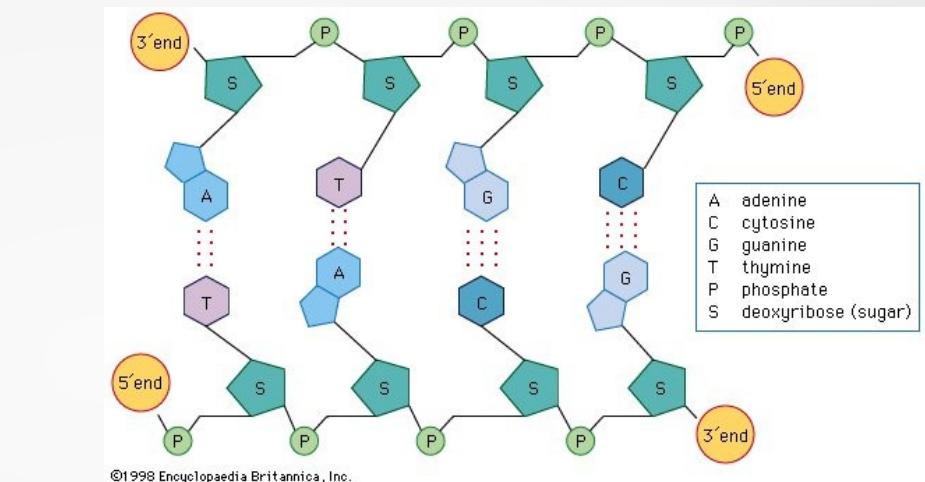
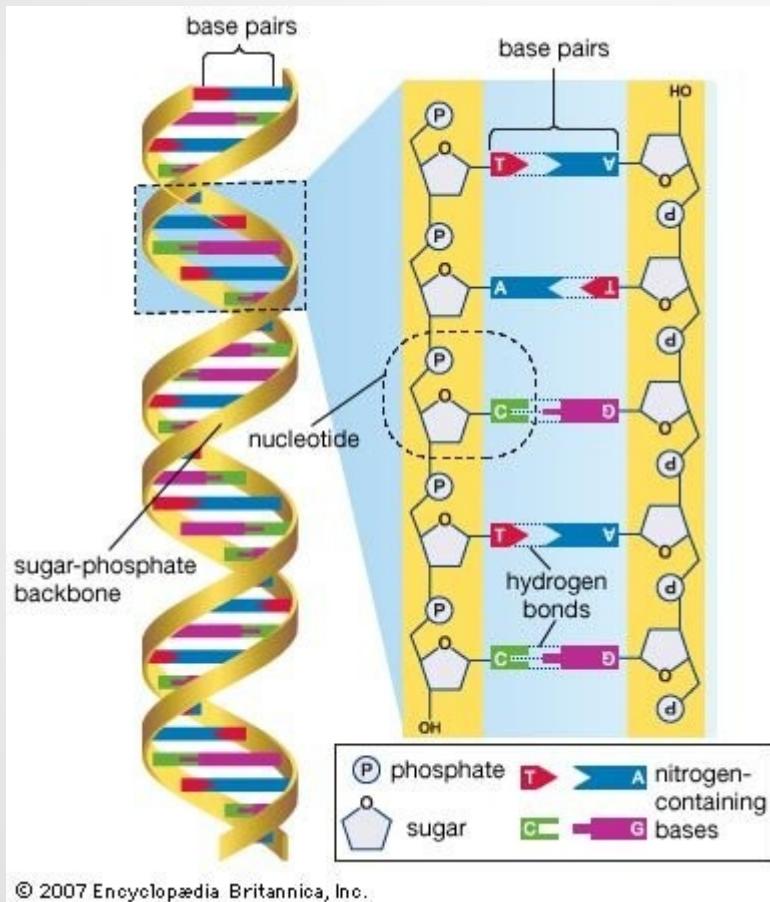
- DNA is not naked in cell. It is associated with proteins and with them forms chromosomes
- chromosomes have structural elements important for its stability and inheritance: telomeres and centromers
- linear chromosomes have a specific problem of end replication, which is solved in eukaryotes by telomeres and telomerase
- Karyotyping and differential staining provides important information about genome structure: presence or absence of heterochromatine, chromosome number, ploidy and etc
- Genome size, ploidy and number of chromosomes can vary even in closely related species



I. Structure and diversity of the genomes

Genome *in silico*
Assemblies and databases

DNA structure



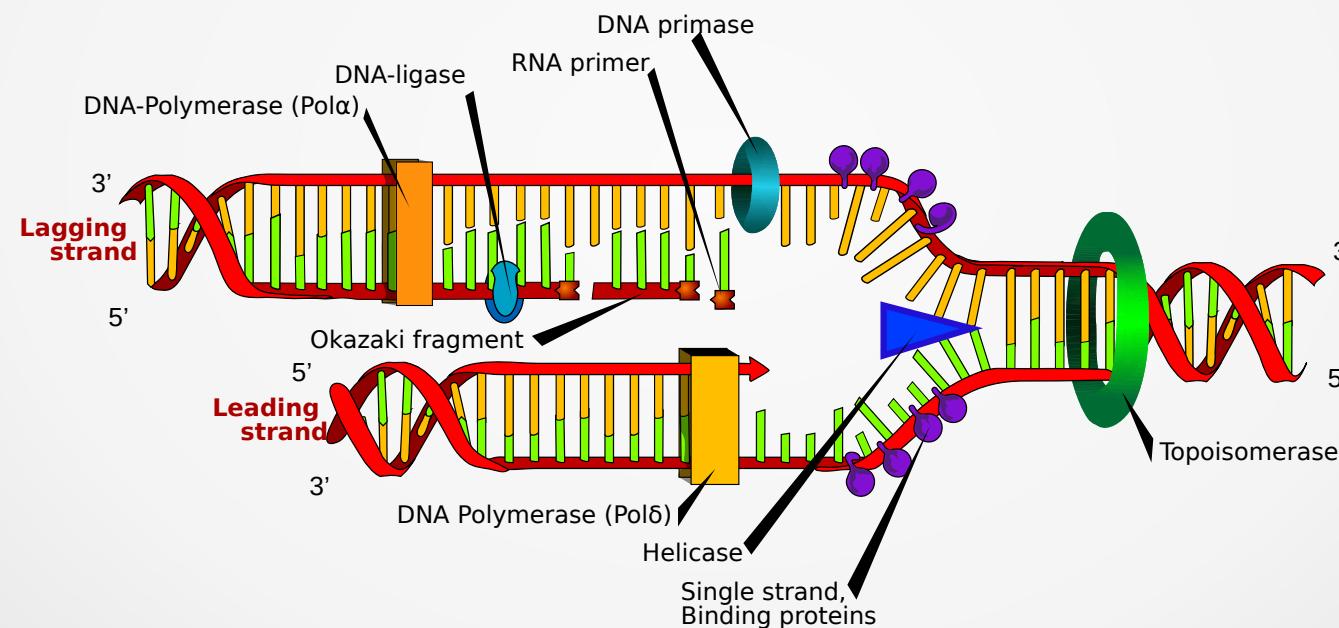
forward strand
reverse strand

5' - TACG - 3'
3' - ATGC - 5'

In databases all DNA/RNA sequences are stored as **one strand in 5' → 3' direction**

Why $5' \rightarrow 3'$ orientation?

All DNA/RNA matrix processes in cell happen in $5' \rightarrow 3'$ direction



During replication DNA polymerase synthesize complementary strand in $5' \rightarrow 3'$

Nucleotides in databases

Symbol	Meaning	Description Origin
G	G	Guanine
A	A	Adenine
T	T	Thymine
C	C	Cytosine
R	G or A	puRine
Y	T or C	pYrimidine
M	A or C	aMino
K	G or T	Ketone
S	G or C	Strong interaction
W	A or T	Weak interaction
H	A or C or T	H follows G in alphabet
B	G or T or C	B follows A in alphabet
V	G or C or A	V follows U in alphabet
D	G or A or T	D follows C in alphabet
N	G or A or T or C	aNy

unambiguous
nucleotides

ambiguous
nucleotides

Sequence formats

FASTA

GENBANK

NEXUS

Fasta format

header1

sequence1

header2

sequence2

header3

sequence3

>sequence1 Description1

gtaacatgcagtcttccggacttttctattatacgactcgtt
ctgctattgtcactaagtccctgcaagctctctccggct
ctctactcttgcaattcttttaatgttcagaacGGCCTGg
ctaacaaaaaggcatgataattgcggcccttggccgggctt
ctgggttcataggggtgtactgcctaaaggccttattgctcc
ctgtaagaaggctgtggactctcatccttaccctgttcaca
tcatagacctggccagattggtagtctgcacgtggcagcct

>sequence2 Description2

ggagacctgccattagagtctggcggtggaccggagtcgctc
cttacccctcaggcgagttgtgtgccttggggagtaaaaacggg
tttagccaaggagggatttcgatcatgtccgtacaggcc
aggatgttggcatctggtcagggtggccatctgggttgc
tgaagttcacggccttaacctgttaagataataggttaggtgaaa
>sequence3 Description3

tcacggccttaacctgttaagataataggttaggtgaaa

Extensions

- most common .fa, .fasta, .fna, .faa
- often used .cds, .pep

Multiple sequence in file

- each sequence is preceded by header
- sequence could be one line or multiline (prefered for long sequences)

Header format

- >sequence2 Description2
- sequence id description
space

Genbank format

Extensions

most common .gb, .genbank

Contains metadata

Contains annotations

Multiple sequence in file

sequence records are separated by //
followed and preceded by newline symbol

Most of fields are optional

It even could contain only annotation

```
LOCUS      MZ571765          406 bp    DNA     linear   INV 10-OCT-2021
DEFINITION Haemoproteus tinnunculi isolate R3 cytochrome b (cytb) gene,
partial cds; mitochondrial.
ACCESSION MZ571765
VERSION   MZ571765.1
KEYWORDS  .
SOURCE    mitochondrialion Haemoproteus tinnunculi
ORGANISM  Haemoproteus tinnunculi
Eukaryota; Sar; Alveolata; Apicomplexa; Aconoidasida; Haemosporida;
Haemoproteidae; Haemoproteus.
REFERENCE 1 (bases 1 to 406)
AUTHORS   Alrefaei,A.F.
TITLE     Molecular diversity of falcon blood parasites in Saudi Arabia:
          Cytochrome b lineages of the genera Haemoproteus (Haemosporida)
          from Riyadh
JOURNAL   Unpublished
REFERENCE 2 (bases 1 to 406)
AUTHORS   Alrefaei,A.F.
TITLE     Direct Submission
JOURNAL   Submitted (15-JUL-2021) Zoology Department, King Saud University,
          King Khalid Road, Riyadh, Riyadh 2455, KSA
COMMENT   ##Assembly-Data-START##
          Sequencing Technology :: Sanger dideoxy sequencing
          ##Assembly-Data-END##
FEATURES
  source      Location/Qualifiers
              1..406
              /organism="Haemoproteus tinnunculi"
              /organelle="mitochondrion"
              /mol_type="genomic DNA"
              /isolate="R3"
              /db_xref="taxon:2588639"
              <1..>406
              /gene="cytb"
  gene        <1..>406
              /gene="cytb"
              /codon_start=3
              /transl_table=4
              /product="cytochrome b"
              /protein_id="UBY00644.1"
              /translation="HILRGLNYSYYVPLSWITGLVIFLISIVTAFMGYVLPWGQMSF
              WGATVITNLLYFIPGLVSICGGYTISDPTLKRFFVLHFIFPPFIALCIVFIHFFLHL
              QGSSNPGLGYDTALKIPFYPSSLCLDVKGFFNNV"
  CDS         <1..>406
              /gene="cytb"
              /codon_start=3
              /transl_table=4
              /product="cytochrome b"
              /protein_id="UBY00644.1"
              /translation="HILRGLNYSYYVPLSWITGLVIFLISIVTAFMGYVLPWGQMSF
              WGATVITNLLYFIPGLVSICGGYTISDPTLKRFFVLHFIFPPFIALCIVFIHFFLHL
              QGSSNPGLGYDTALKIPFYPSSLCLDVKGFFNNV"
ORIGIN
              1 tacatattct aagaggattta aattttatcat atgtatattt acctttatca tggataactg
              61 gatttagttat attcctaatac tctatagtca ctgttttat gggttatgtt ctacccctgg
              121 gtcaaatggat ttctcggtt gcaacagtta ttactaaattt attatatttc ataccaggac
              181 tagtacatcg gatttggtt ggatatacta ttatgtaccc tactttaaa agattcttg
              241 tattacacctt tatattccca ttatagccc ttatgtatcg atttatacat atattctct
              301 tacatttaca aggttagctt aaccctttat gagatgtatcat agctttaaaat accctttct
              361 atccaatgtt attatgttta gatgttaaag gatftaaata tgatt
//
```

The diagram illustrates the hierarchical structure of a Genbank file. On the left, the file content is shown as a series of text lines. On the right, three blue curly braces group specific sections: the top brace groups the first few lines under 'metadata'; the middle brace groups the 'FEATURES' section and its sub-sections ('source', 'gene', 'CDS') under 'annotations'; and the bottom brace groups the 'ORIGIN' section and its sequence under 'sequence'.

Nexus format

Extensions

most common .nex, .nexus

Could contain nearly everything.

Even commands for tools

Multiple sequence in file

Most of fields are optional

Mostly used in phylogenetics
and phylogeography

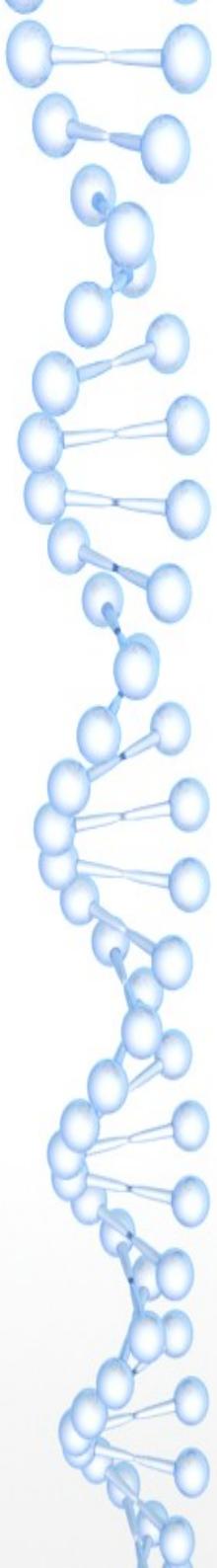
```
#NEXUS
Begin TAXA;
    Dimensions ntax=4;
    TaxLabels SpaceDog SpaceCat SpaceOrc SpaceElf
End;

Begin data;
    Dimensions nchar=15;
    Format datatype=dna missing=? gap=- matchchar=.;
    Matrix
        SpaceDog  atgcttagctagctcg
        SpaceCat  ....??.a.
        SpaceOrc  ...t.....g.
        SpaceElf  ...t.....a.
    ;
End;

BEGIN TREES;
    Tree tree1 =
        (((SpaceDog, SpaceCat), SpaceOrc, SpaceElf));
END;
```

Available genome databases

- **NCBI Genome** <https://www.ncbi.nlm.nih.gov/genome>
 - huge whale of bioinformatics
- **Ensembl** <https://www.ensembl.org>
 - just an elephant
- **GenomeArk** <https://vgp.github.io/genomeark/>
 - contains genomes assembled according VGP standards
- **DNAzoo** <https://dnazoo.org>
 - contains genomes assembled by DNAzoo team to chromosome level



I. Structure and diversity of the genomes

Genome *in silico*
Connection between assembly and karyotype

How to assign chromosome names to C-scaffolds?

```
>HiC_scaffold_1  
gtaacatgcagttctccggacttttctattataactcgtt  
ctgctattgtcactaagtccctgcaagctttctccgggtct  
ctctactcttgcaattttttatgttcagaacGGCCTGg  
ctaacaaaagccatgataattgcggccttggccgggctt  
ctgggttcataggggtgtactgcctaaaagccttattgctcc  
ctgttaagaaggctgctggactctcatccttaccctgtctcaca  
tcataagacctggccagatttgttagtctgcacgtggcagcct
```

...

...

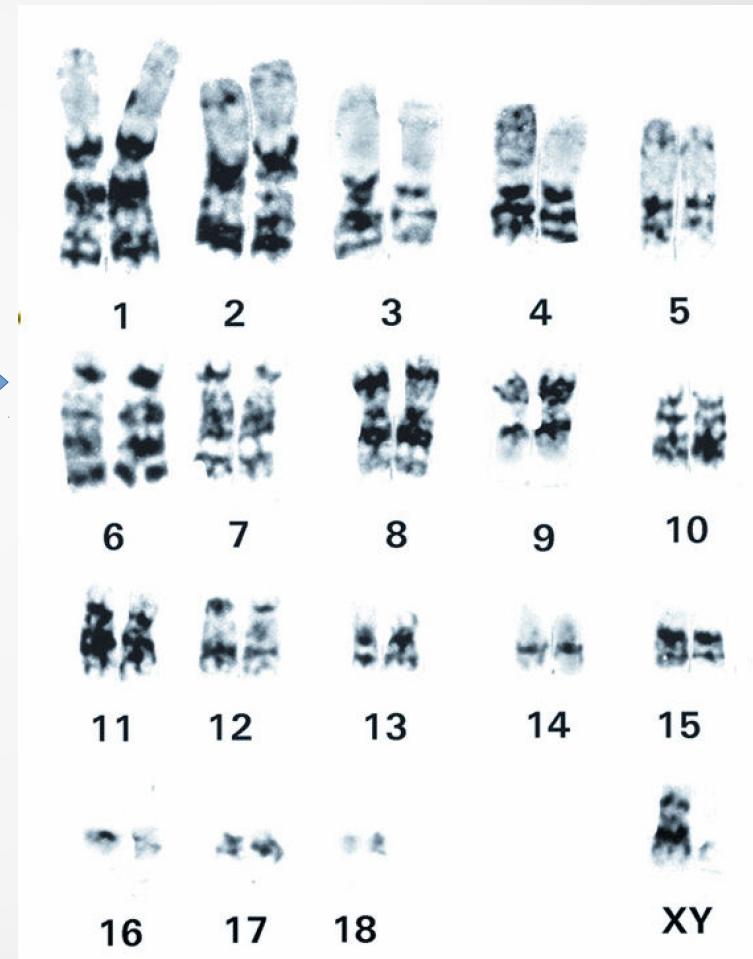
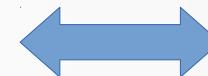
```
ttaatgttcagaacGGCCTGgctaacaaaagccatgataatt
```

```
>HiC_scaffold_2
```

```
ggagacctgccattagagtctggcggtggaccggagtcgctc  
cttaccttcaggcgagttgtgtgcttggggagtaaaaacggg  
tttagccaaggagggaattctcgatcatgtcctgacaggcc  
aggatgttggcatctggtcagggtggccatctgggttgc  
tgaagttcacggcctaacctgttaagataataggtaggtgaaa
```

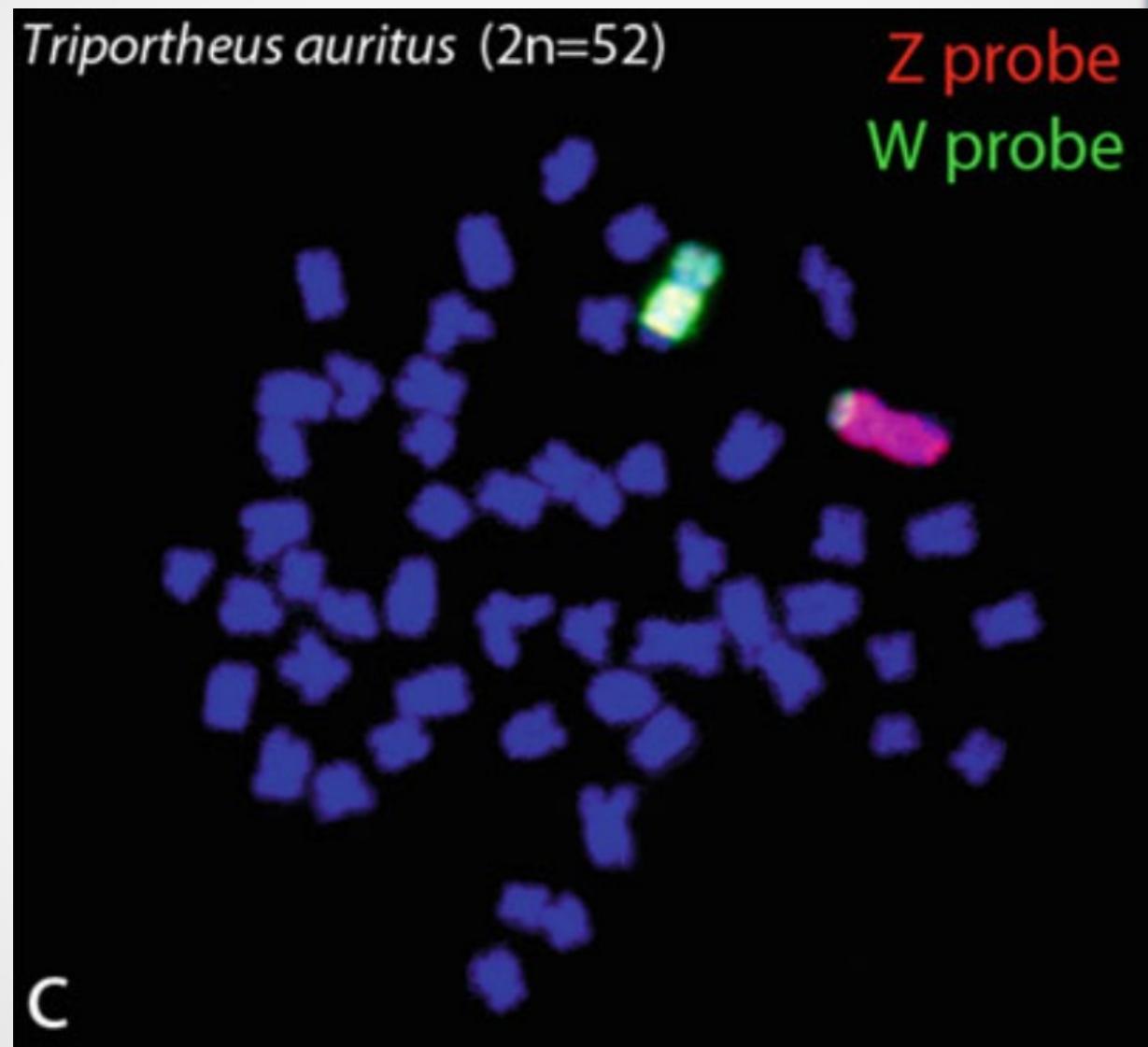
...

...

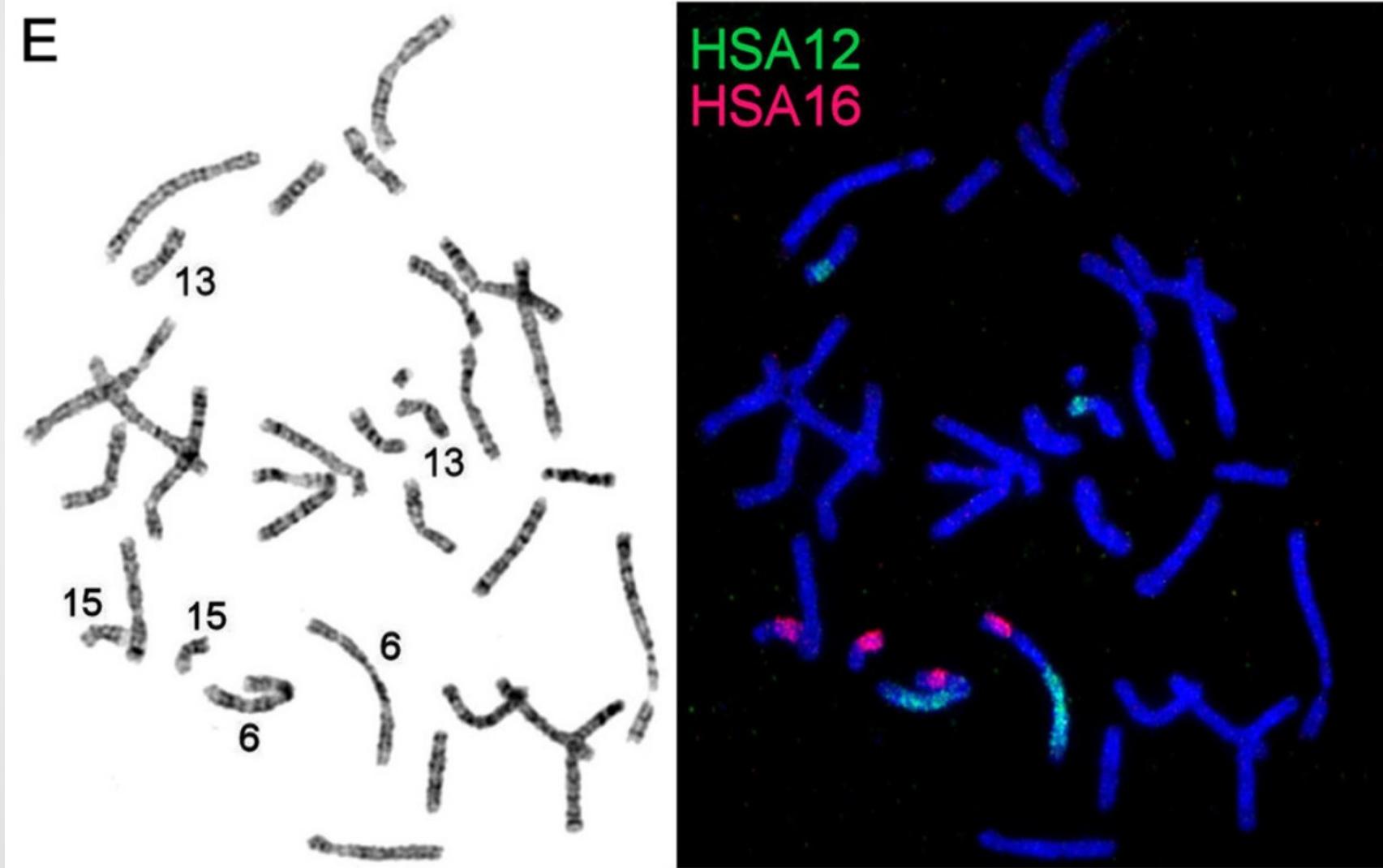


We need markers which we could detect both in the assembly and karyotype.
What can we use as marker? How many of them do we need?

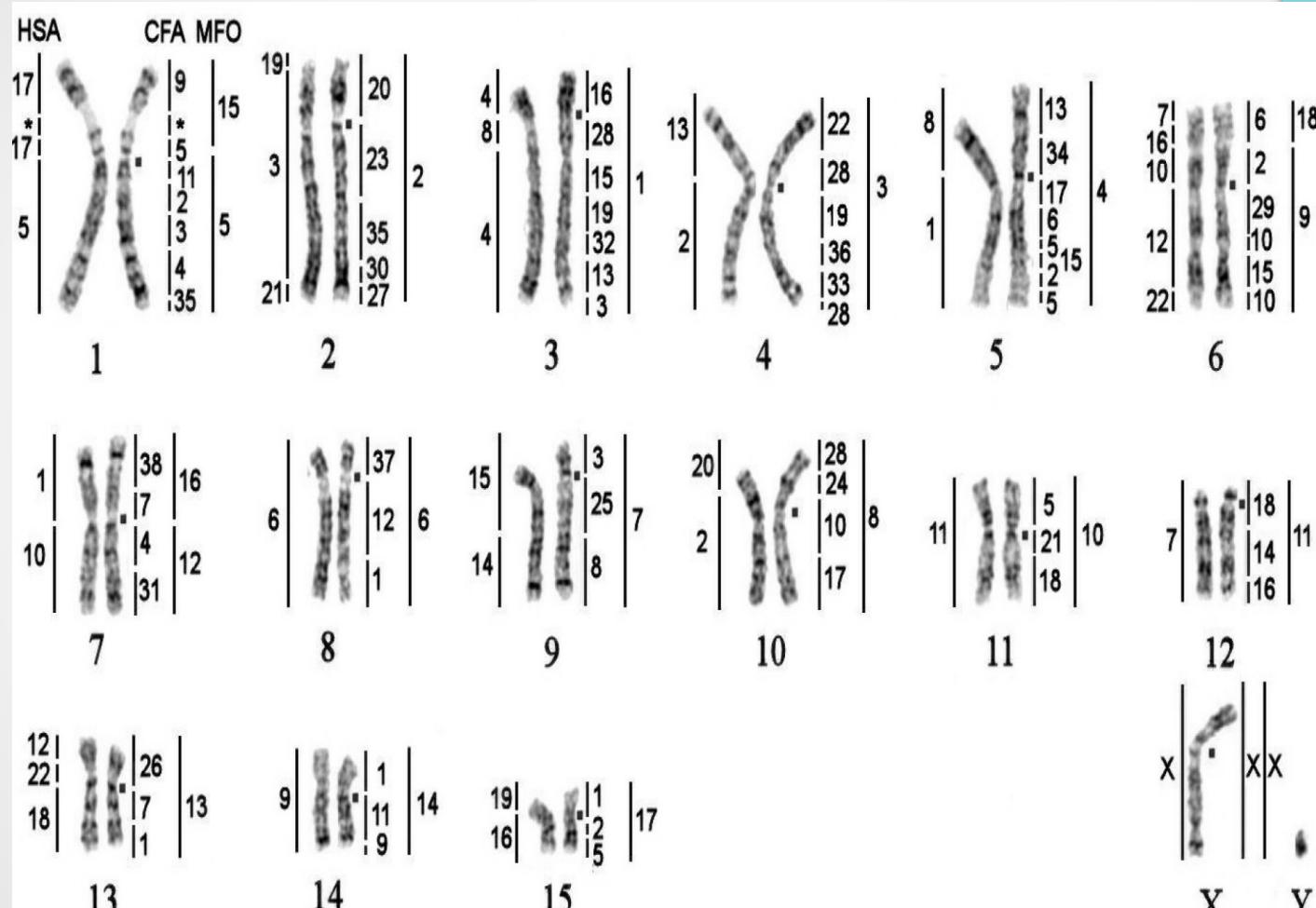
Cross-species chromosome painting (ZooFISH)



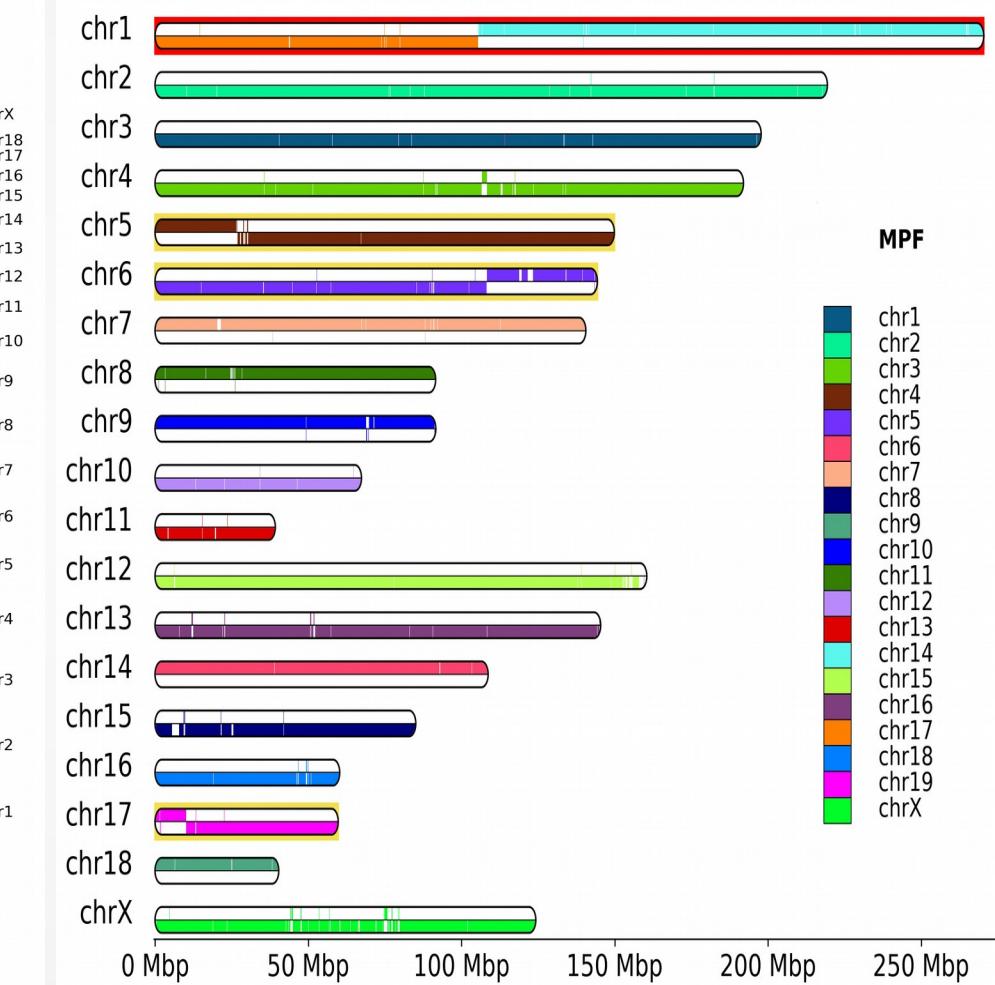
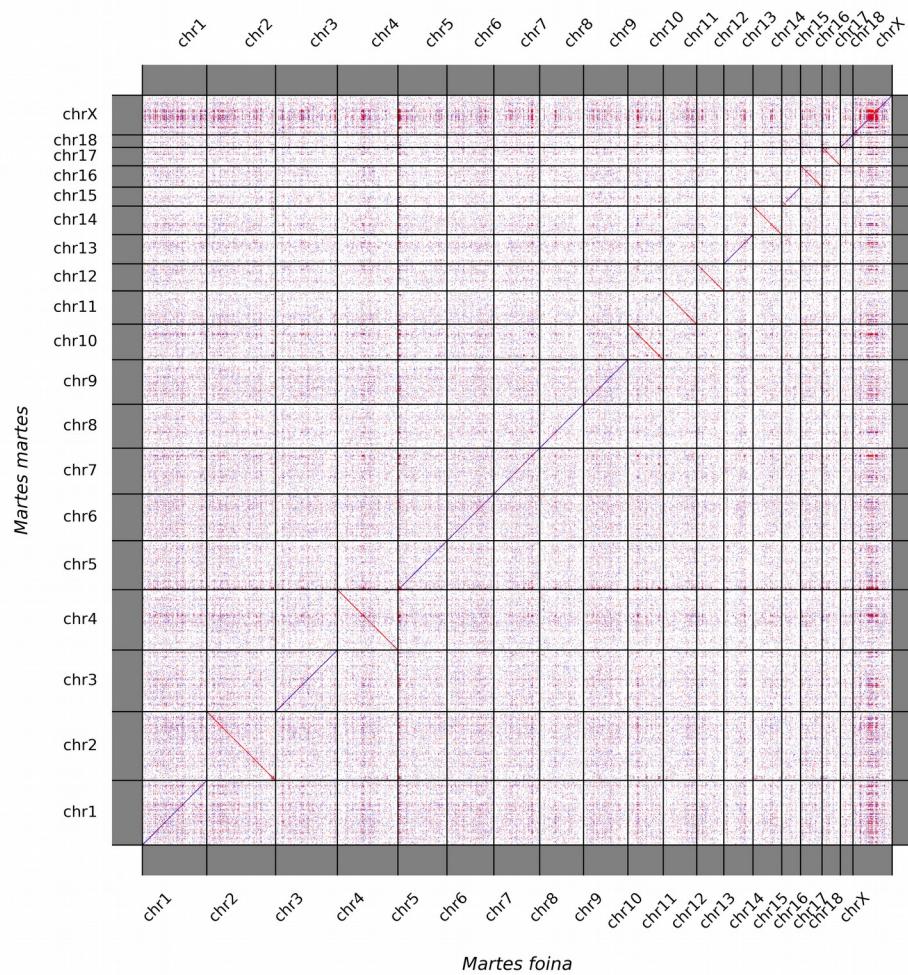
Cross-species chromosome painting (ZooFISH)



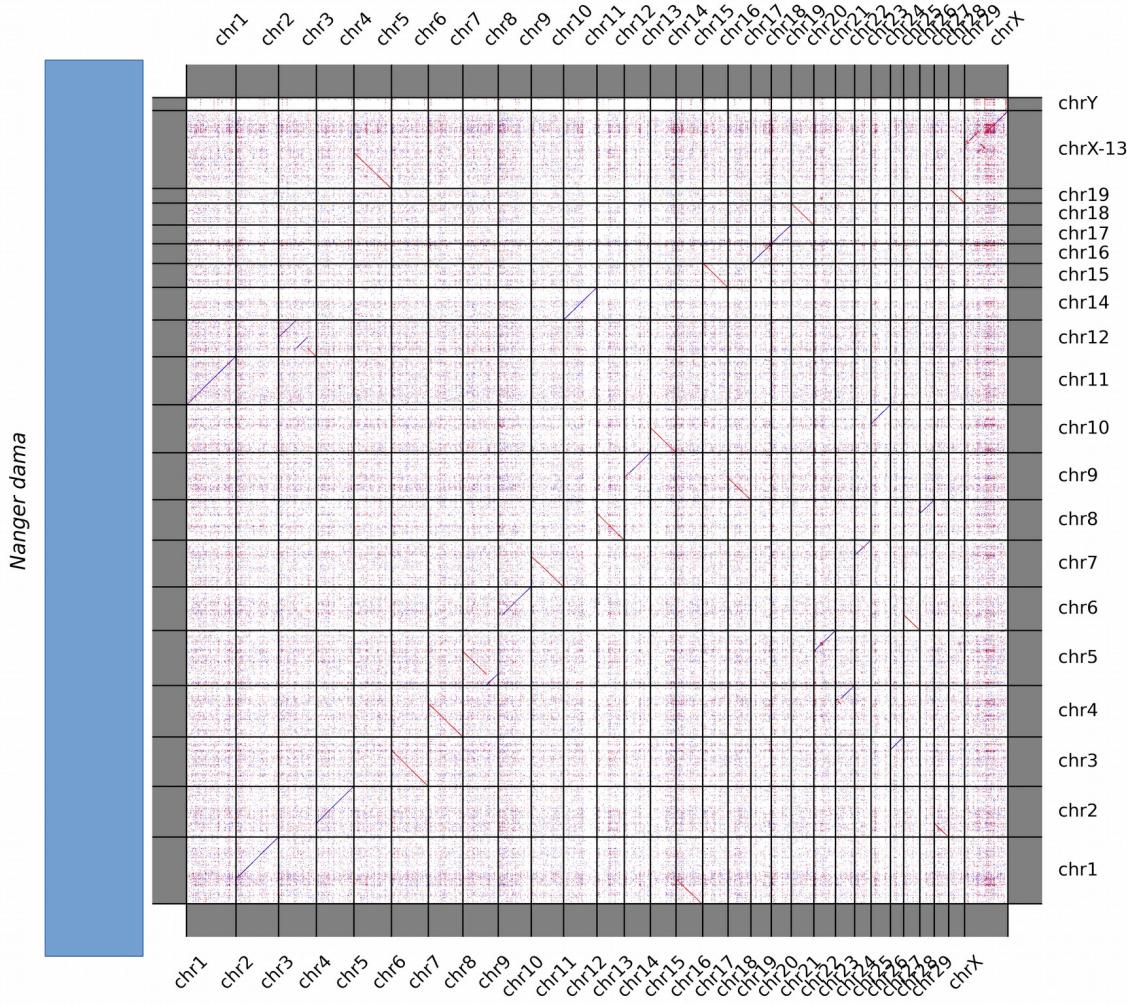
FISH maps



Whole genome alignment

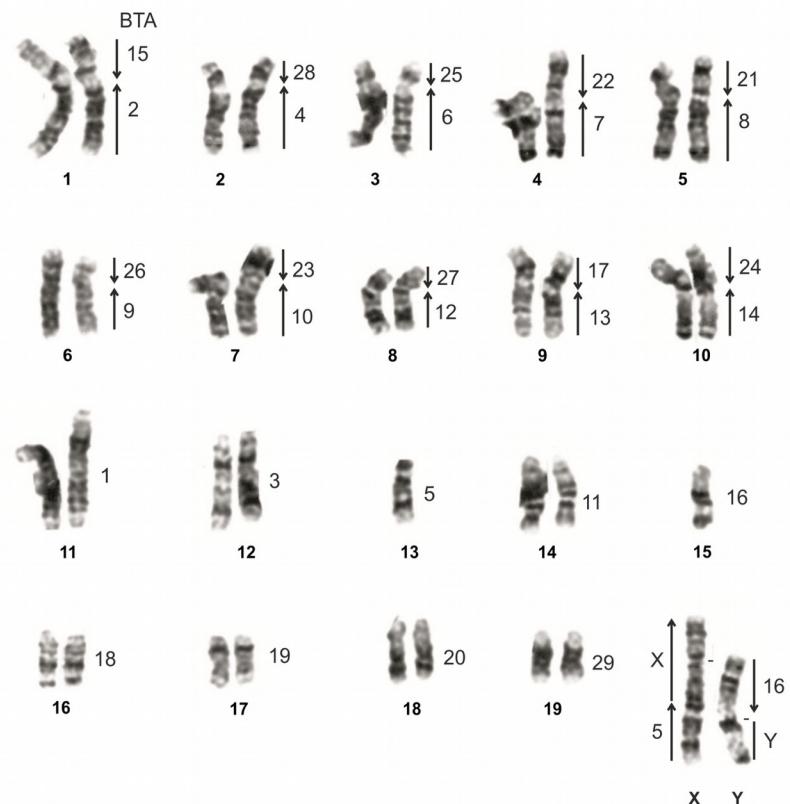


Correspondence between WGA and FISH map (1)



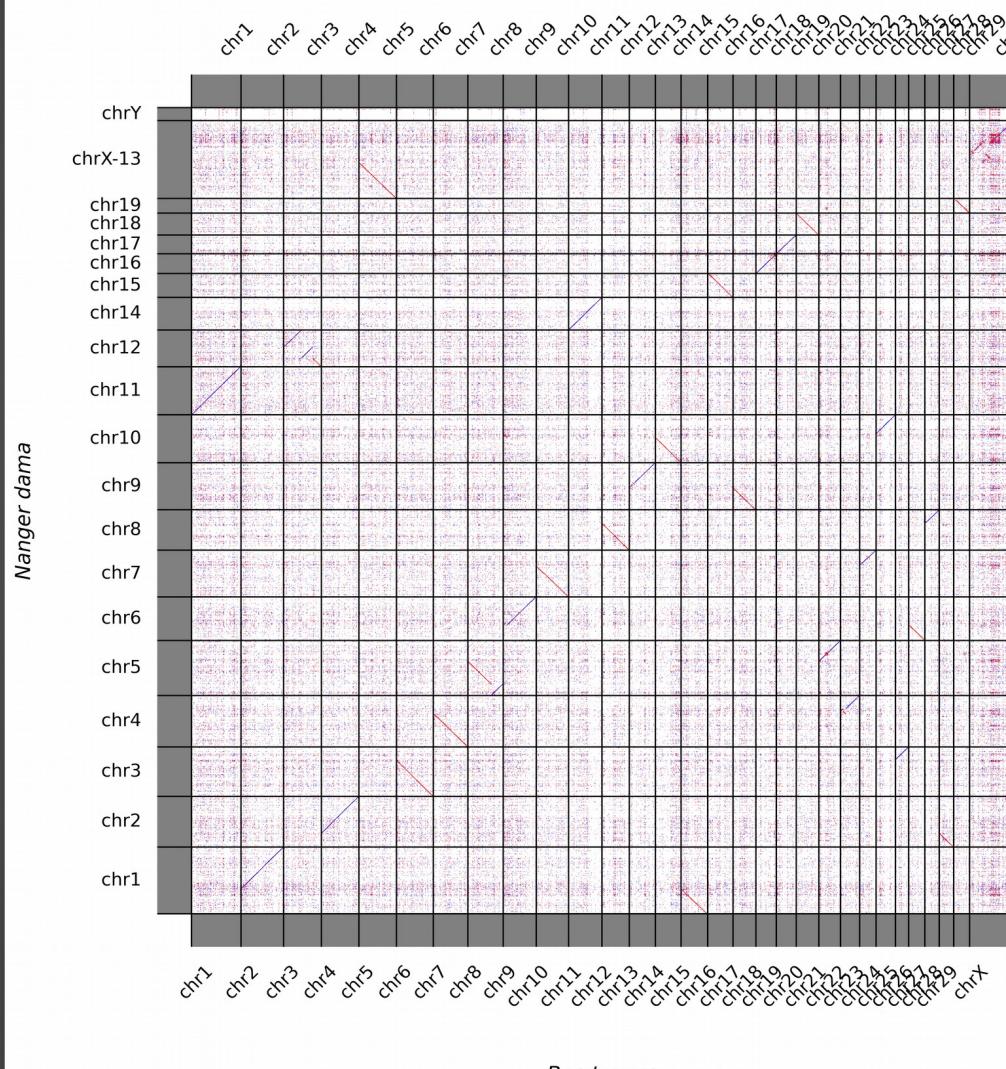
Bos taurus

***Nanger (Gazella) dama mhorr*
(Mhorr Gazelle)**

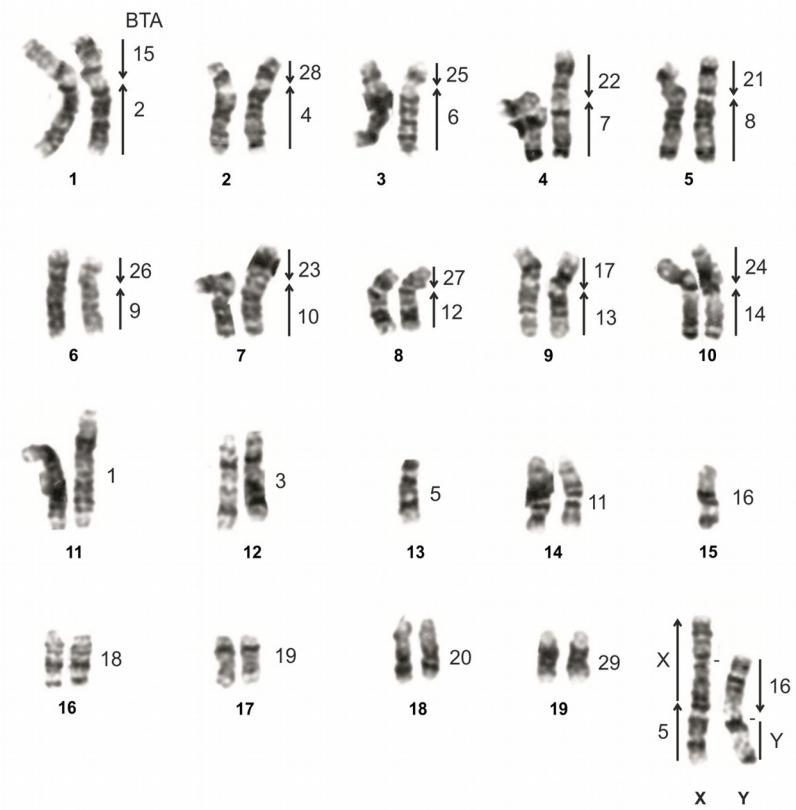


Chernohorska et al, 2012

Correspondence between WGA and FISH map (2)

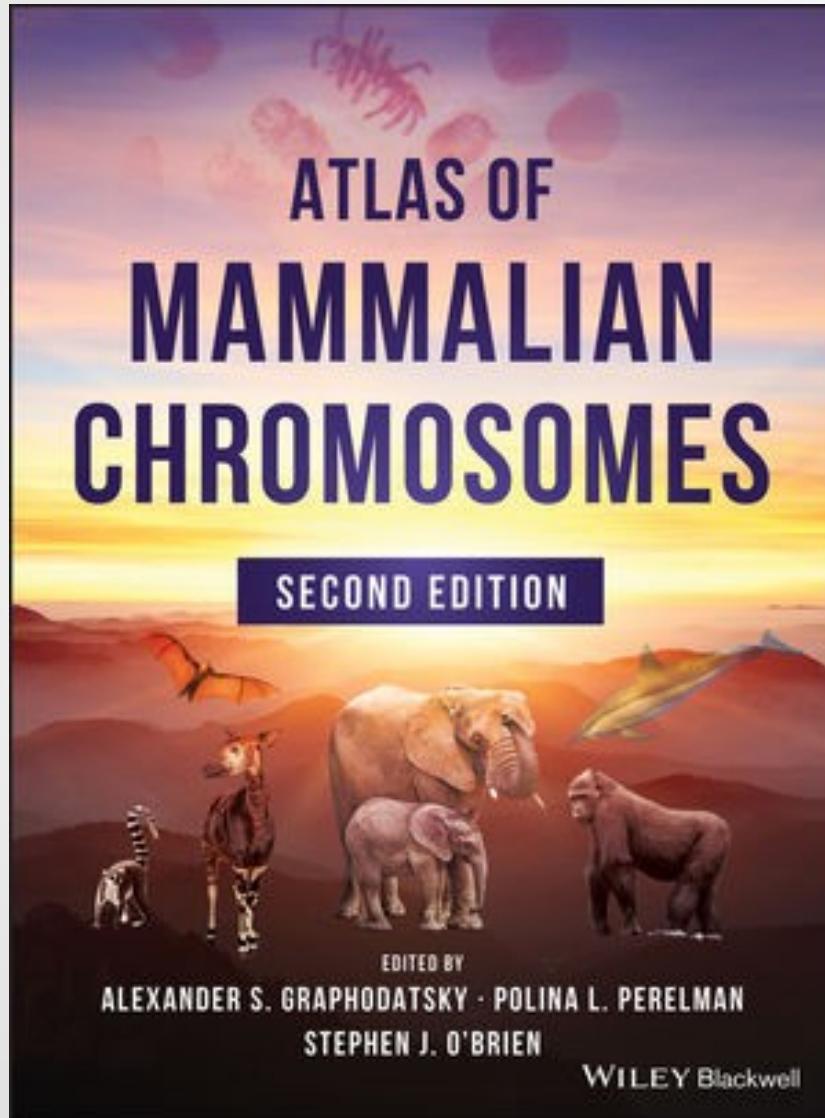


***Nanger (Gazella) dama mhorr*
(Mhorr Gazelle)**



Chernohorska et al, 2012

Atlas of mammalian chromosomes



Big book contains
hundreds karyotypes and FISH maps

Very useful for *de novo* assembly

Rules for ordering and orientation of chromosomes

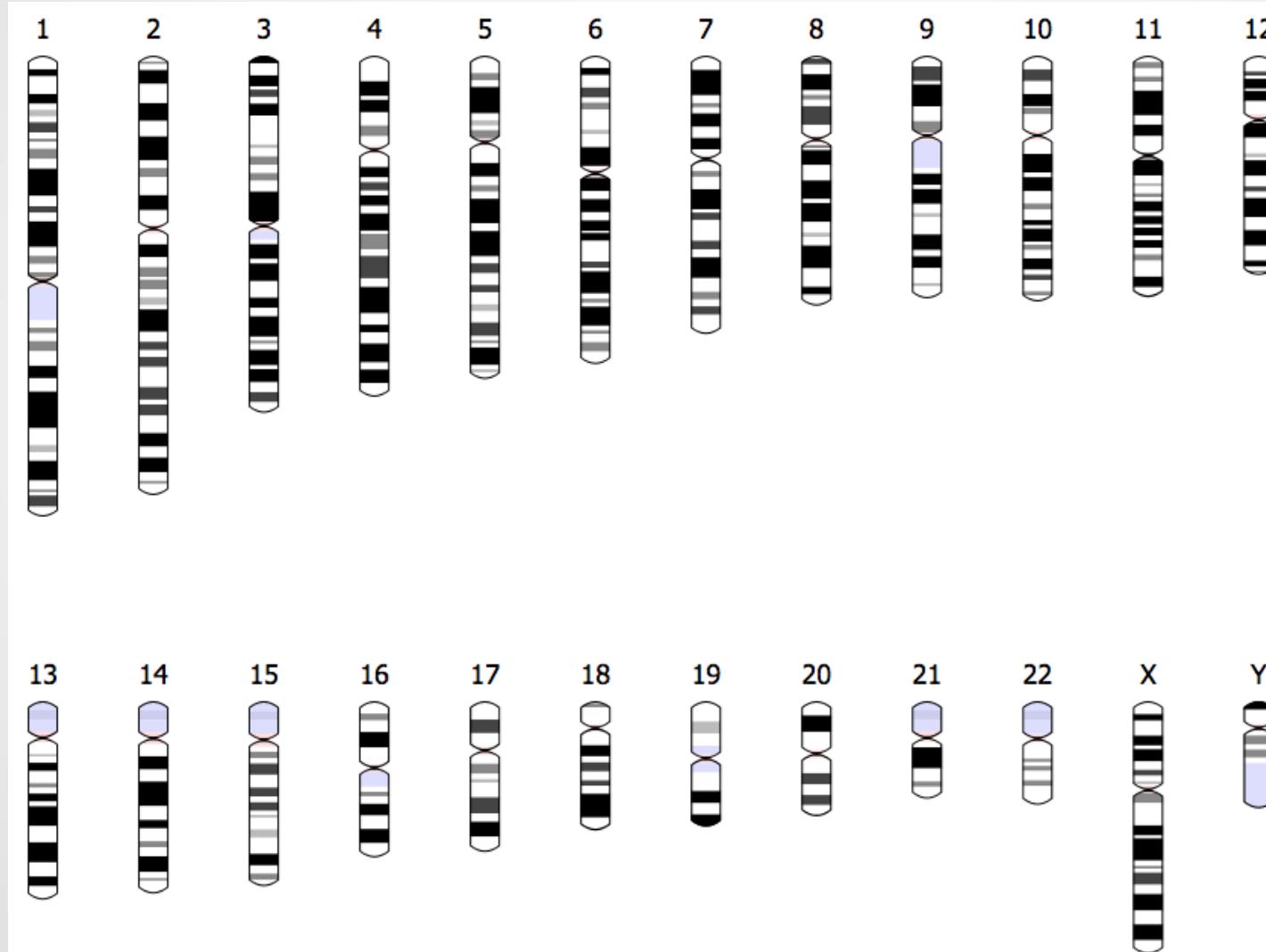
Rules for orientation and ordering of chromosomes and C-scaffolds:

- Autosomes and C-scaffolds should be ordered by length, from longest to shortest
- Each chromosome and C-scaffold should start from p-arm (short arm)

Issues:

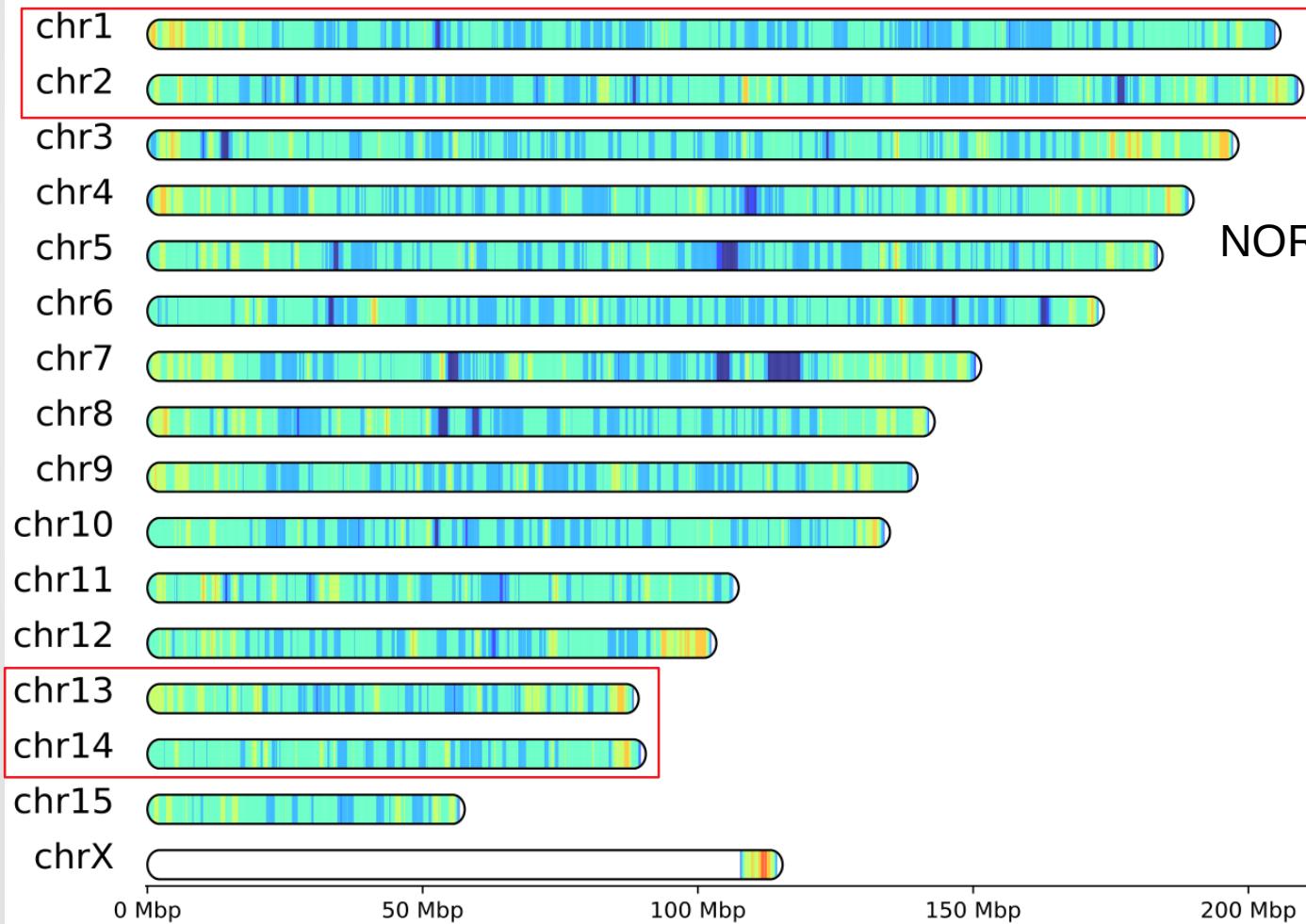
- Length of C-scaffold might significantly differ from length of corresponding chromosome
- Chromosomes might have a very similar length
- For proper orientation of C-scaffolds we need to know the location of centromere in each C-scaffold to understand where is p-arm and q-arm.
- for nearly ideal metacentric (p- and q- arm are nearly equal in length) chromosomes we might not be able to distinguish p- and q-arm
- Old nomenclature, contradicting to rules, might exist

Orientation of human chromosomes

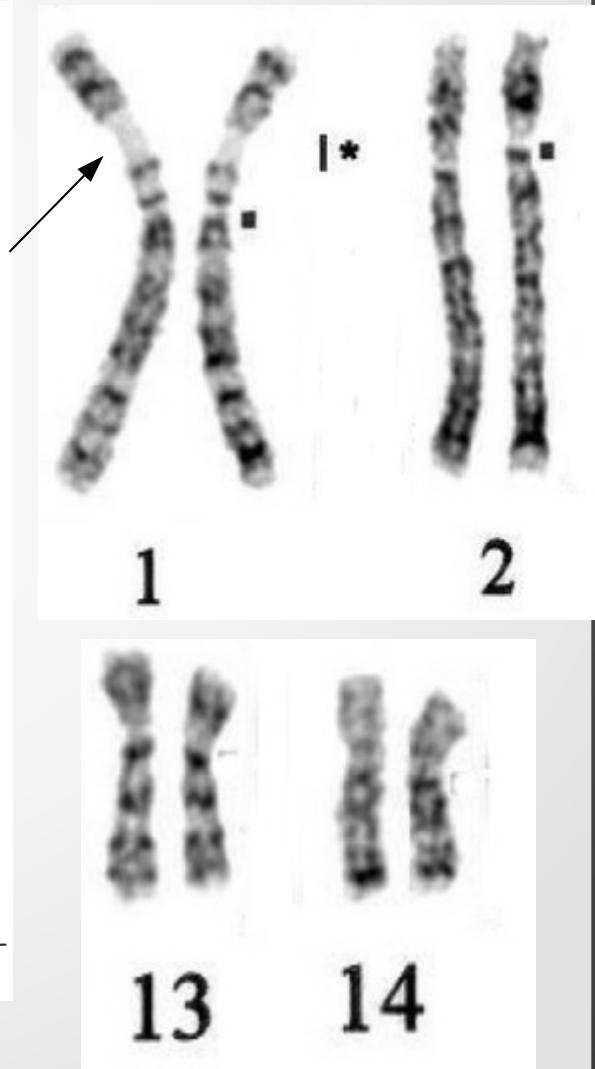


Difference between length of C-scaffolds and chromosomes: case of baikal seal

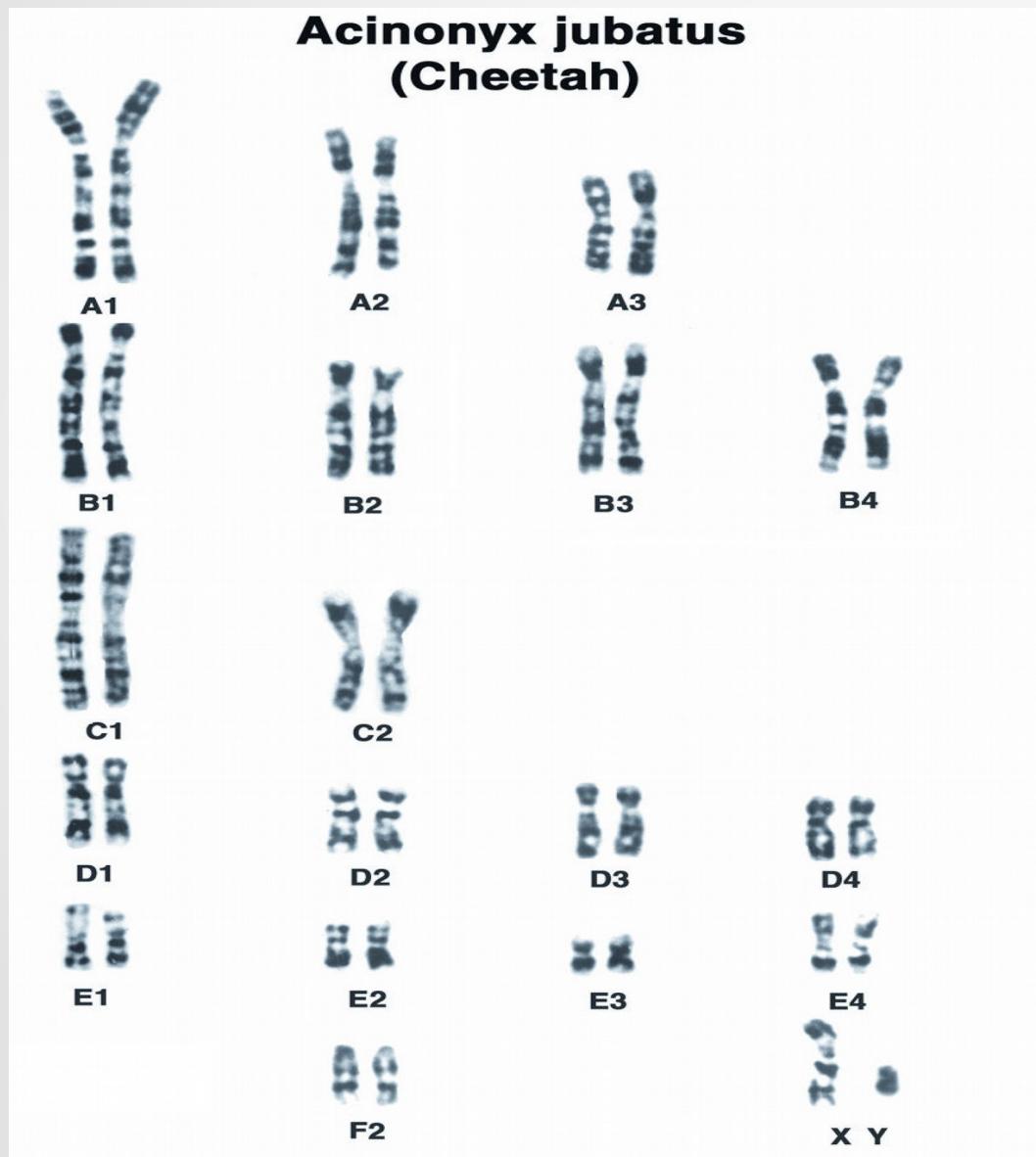
Assembly



Karyotype



Old nomenclature contradicting rules: cat case

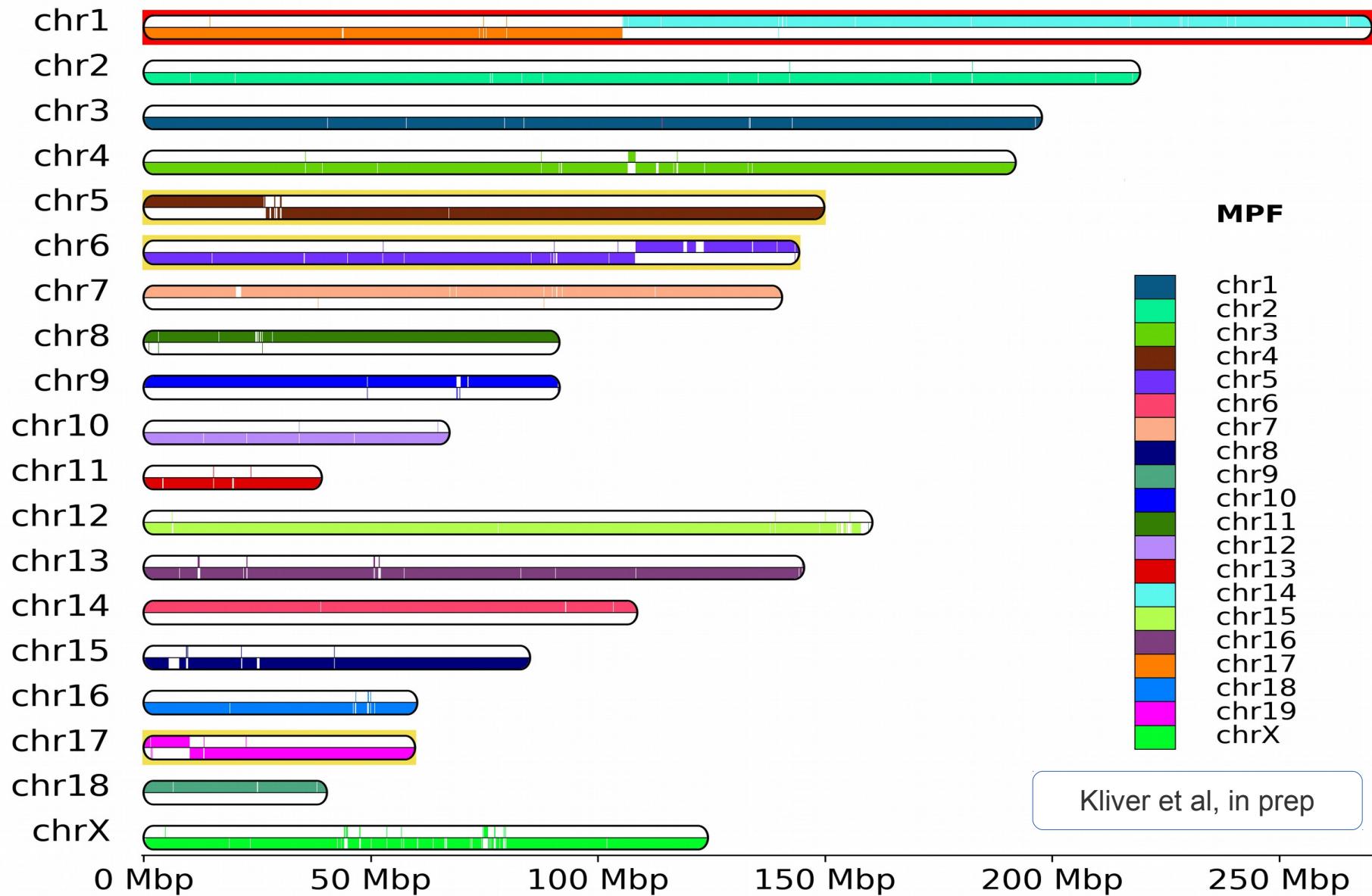


A, B, C, D, E, F groups are used for genomes of all cats.

Inherited from ... 1965

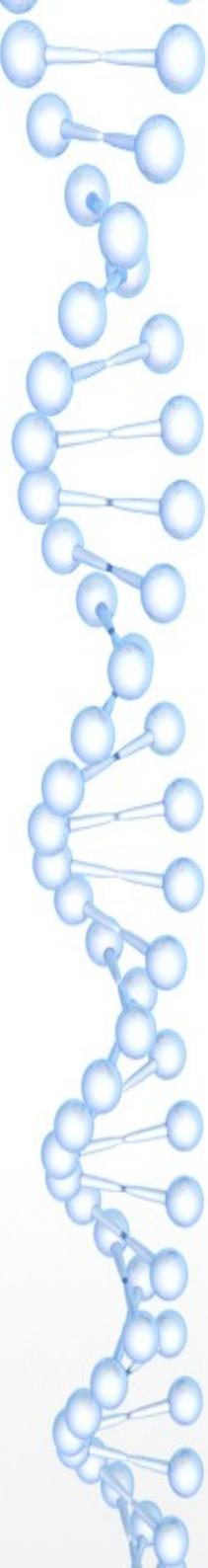
Graphodatsky et al, 2020

Example of different C-scaffold orientation in related species: black-footed ferret and domestic ferret



Summary

- Rules of chromosome nomenclature and ordering are very often impossible to follow. You should consider them as recommendation
- Always check if elder nomenclature exists for your species. Follow it if it exists.
- Do not forget that even closely related species with similar chromosome numbers and karyotypes could have different orientation of C-scaffolds(chromosomes) in the assembly. Assembler/scaffolder choose orientation of sequences randomly
- To assign chromosome names to C-scaffolds of your new chromosome-level genome assembly you need to find a reference species for which it was already done and for chromosomes of which you already have FISH map. Then just compare FISH map and WGA between your ad reference species. Sometimes you to do it multiple times if required FISH map is unavailable.
- If no cytogenetic data is available directly or indirectly, for mammals you could identify X-chromosome. And order autosomes by length.



I. Structure and diversity of the genomes

Mutations

Classification of mutations

- genomic mutations - change of chromosome number
 - whole genome duplication
 - aneuploidy (loss or gain individual chromosomes)
- chromosomal mutations
 - fusions and fissions (Robertson translocations)
 - big insertions, inversions, deletions, translocations
 - small translocations
- gene mutations
 - small insertions, inversions and deletions
 - single nucleotide substitutions



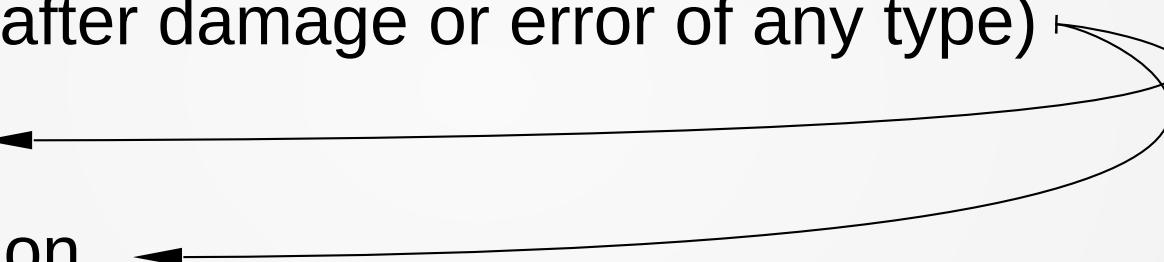
Large-scale mutations

Small-scale mutations

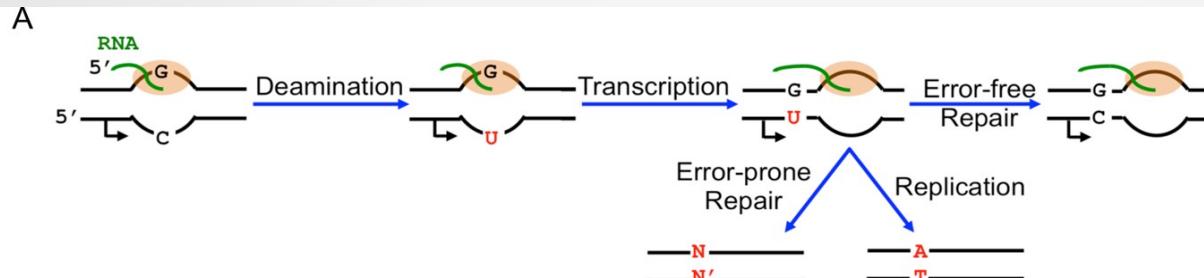
This commonly used classification is ambiguous !

Sources of mutations

Mutations are **results** of errors of:

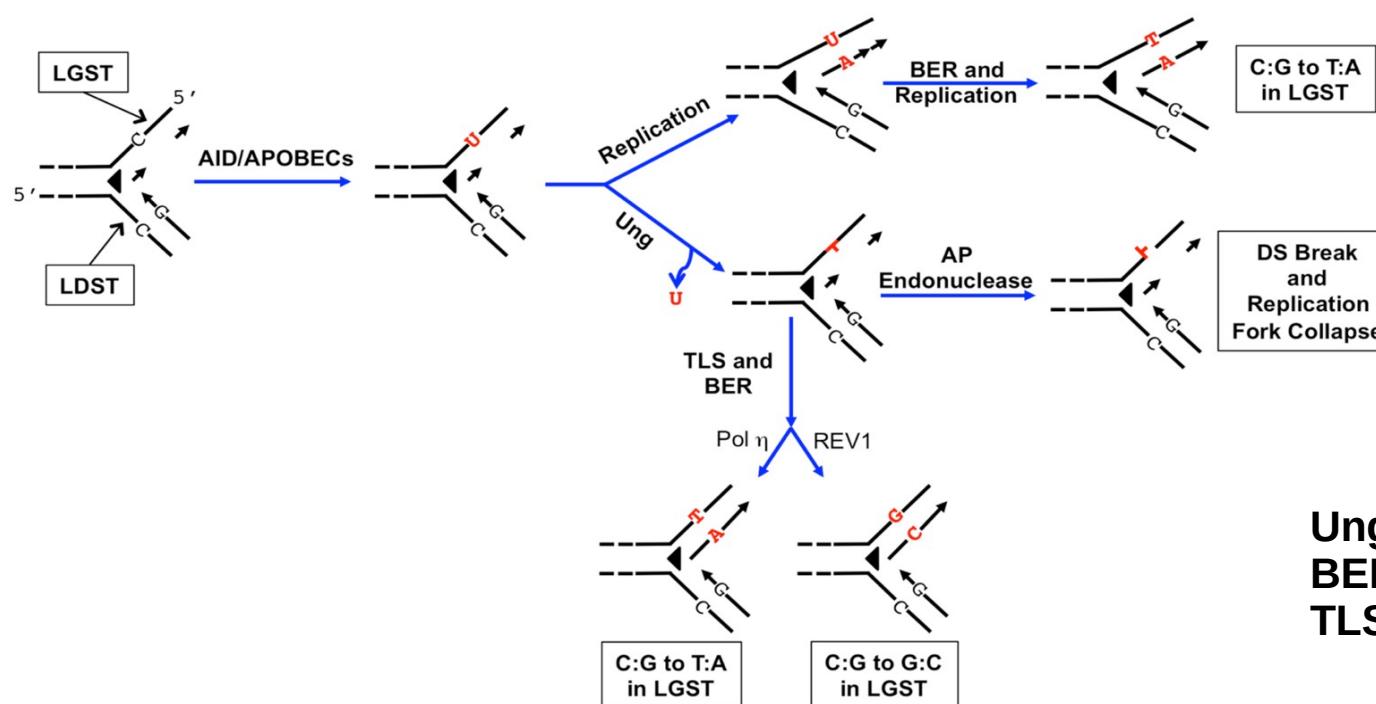
- Reparation (after damage or error of any type)
 - Replication
 - Recombination
 - Separation of chromosomes between daughter cells
- 

Difference between DNA damage and mutation: deaminase case



deaminases convert cytosine to uracile in ssDNA

B



Examples of deaminases:

AID - Activation-induced deaminase

APOBEC - apolipoprotein B mRNA-editing catalytic polypeptide-like

Ung - uracil-DNA glycosylase
BER - base-excision repair
TLS - translesion synthesis

LDST - leading strand template
LGST - lagging strand template

DNA polymerases δ and ϵ

- eukaryotic polymerases
- multi subunit protein complex
- participate in genome replication
- three types of activities:
 - 5' \rightarrow 3' polymerase
 - 5' \rightarrow 3' exonuclease
 - 3' \rightarrow 5' exonuclease
- strand elongation
- excision of obstacles
- self-proofreading

Difference between error and mutation. Precision of replication

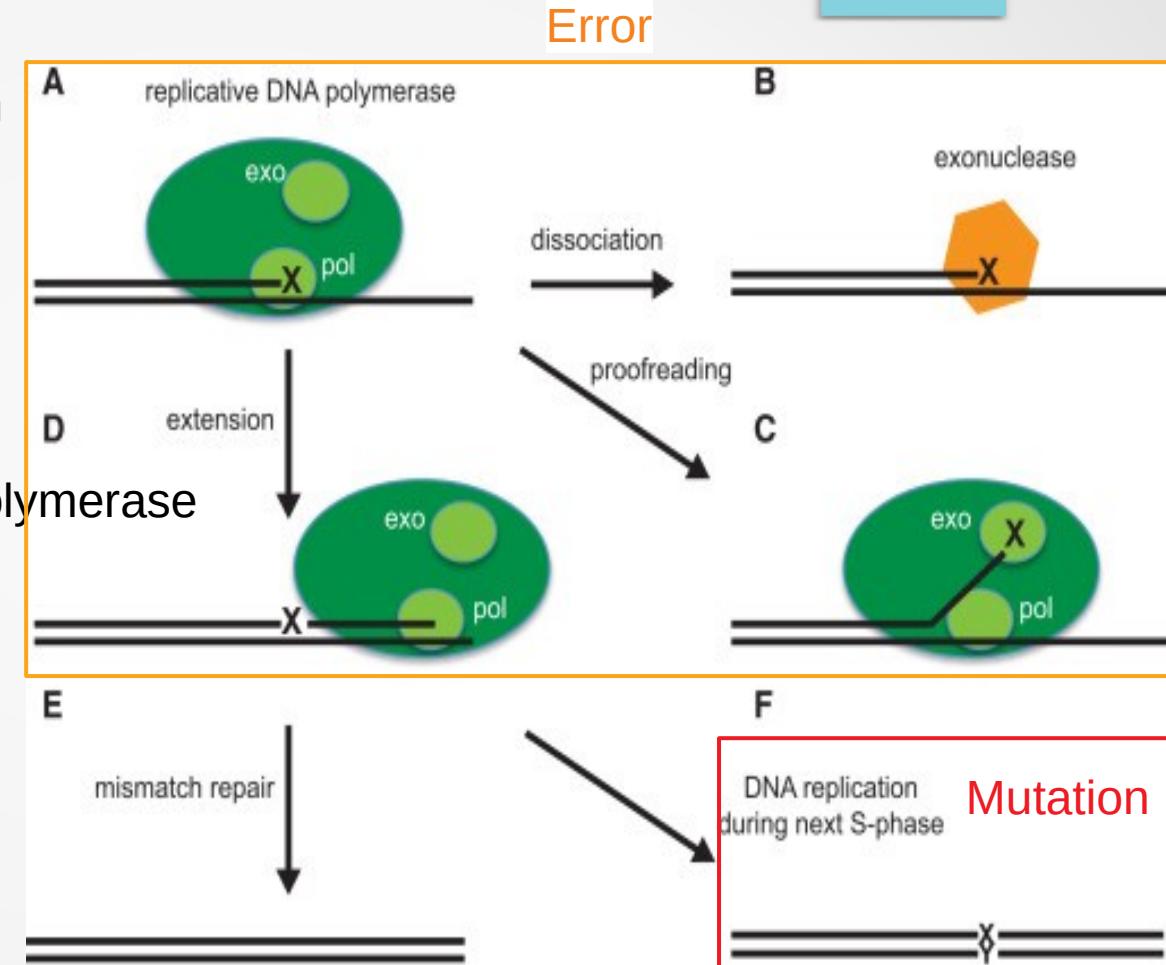
Error rate of polymerase reaction

$\sim 10^{-4} - 10^{-5}$

Error rate after proofreading by polymerase

$\sim 10^{-7} - 10^{-8}$

Ganai and Johansson et al, 2016



X - mismatching nucleotide

Error rate after final proofreading
(methylation-dependent)

$\sim 10^{-10}$

Mutation problem in big multicellular organisms

Replication error rate: 10^{-10} per base per division

Human genome size: 3.3×10^9 bp

Cells in human body:

human:

all:

$3.0 \cdot 10^{13}$ cells

without unnnuclear cells:

$0.3 \cdot 10^{13}$ cells

bacterial:

$3.8 \cdot 10^{13}$ cells

Sender et al, 2016

How to reduce a number of inherited mutations?

Mutation problem in big multicellular organisms

Replication error rate: 10^{-10} per base per division

Human genome size: 3.3×10^9 bp

Cells in human body:

human:

all:

$3.0 \cdot 10^{13}$ cells

without unnnuclear cells:

$0.3 \cdot 10^{13}$ cells

bacterial:

$3.8 \cdot 10^{13}$ cells

Sender et al, 2016

How to reduce a number of inherited mutations?

Reduce number of divisions.

How? Less divisions - less cells.

Mutation problem in big multicellular organisms

Replication error rate: 10^{-10} per base per division

Human genome size: 3.3×10^9 bp

Cells in human body:

human:

all:

$3.0 \cdot 10^{13}$ cells

without unnnuclear cells:

$0.3 \cdot 10^{13}$ cells

bacterial:

$3.8 \cdot 10^{13}$ cells

Sender et al, 2016

How to reduce a number of inherited mutations?

Reduce number of divisions

How? Less divisions - less cells

Separate germline and somatic cells as earlier as possible

Somatic and germline cells

Germline (germ) cells

gametes (eggcells, spermatozoa, etc) and their ancestor cells

Somatic cells

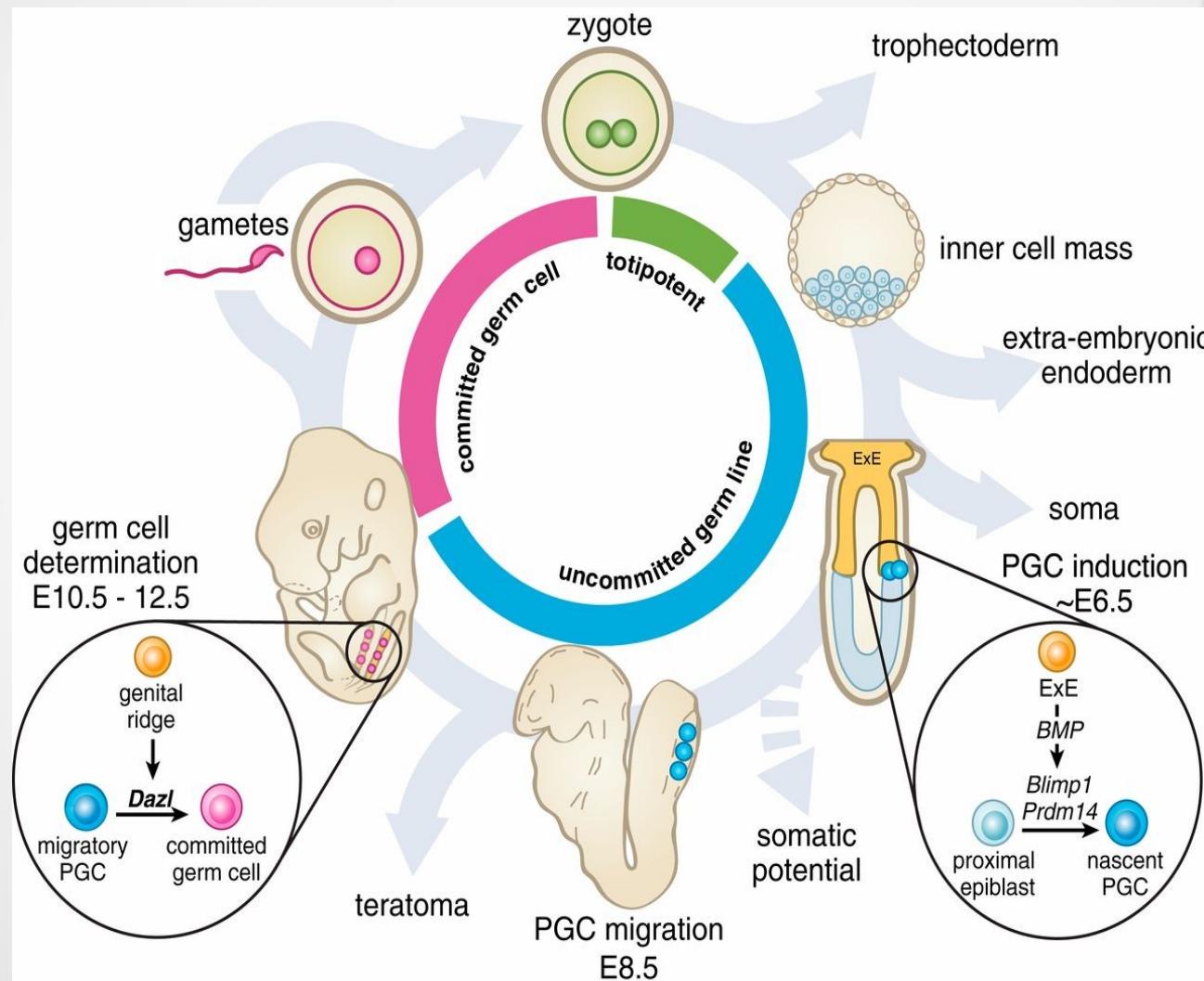
all other cells

Primordial germ cells (PGCs)

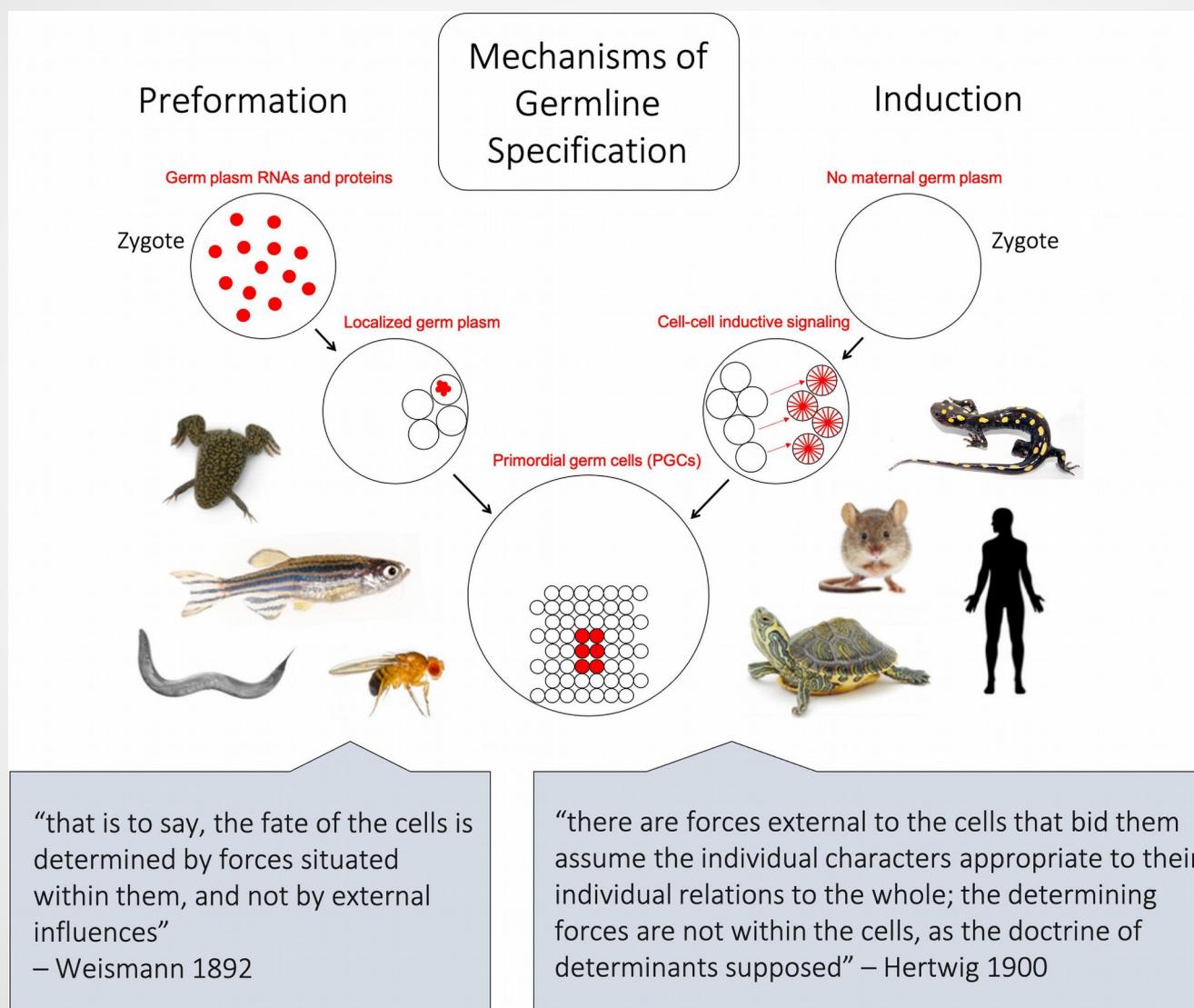
embryonic precursors of germline cells

Divergence of somatic and germline cells happens very early in embryogenesis !

Germline cells in development of mammals



Germline cells in development of animals



Variation of genome in somatic cells

Polyplloid cells in human:

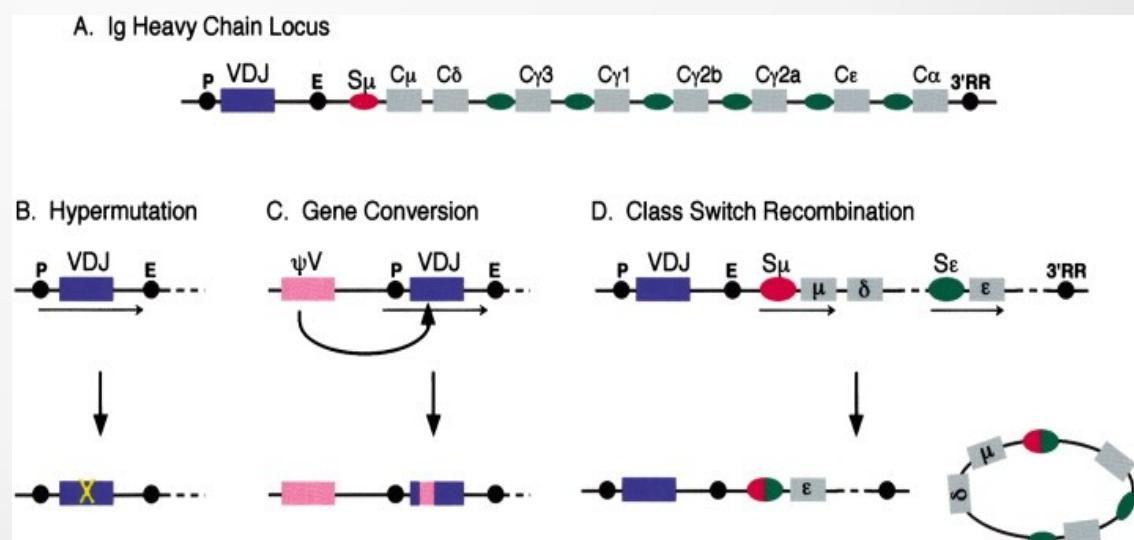
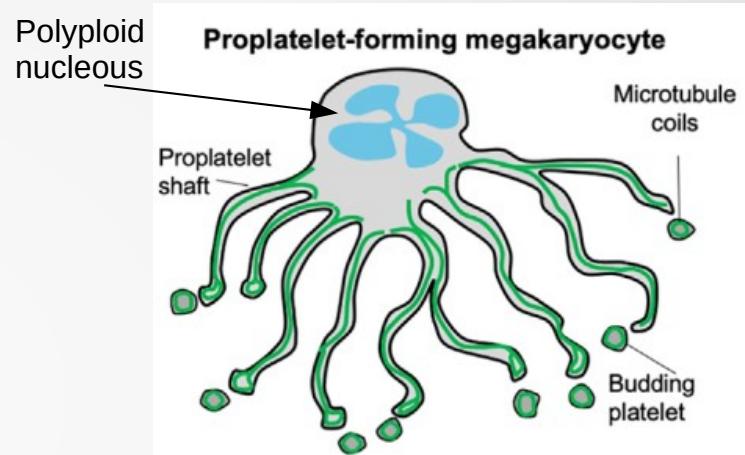
- trophoblast giant cells, critical for pregnancy.
- megakaryocytes, produce platelets (thrombocytes)
- hepatocytes, liver cells
- etc

Somatic hypermutation

loci of immunoglobulin genes

in B cells

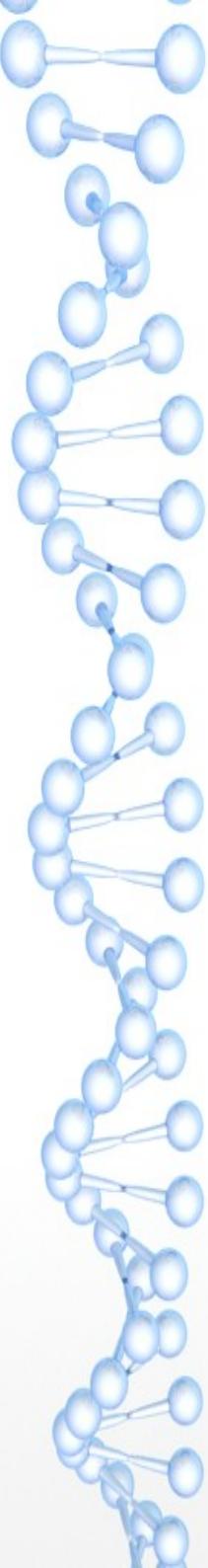
- point hypermutagenesis
- gene conversion
- deletions



Heib et al, 2021; Papavasiliou and Schatz, 2002

Summary

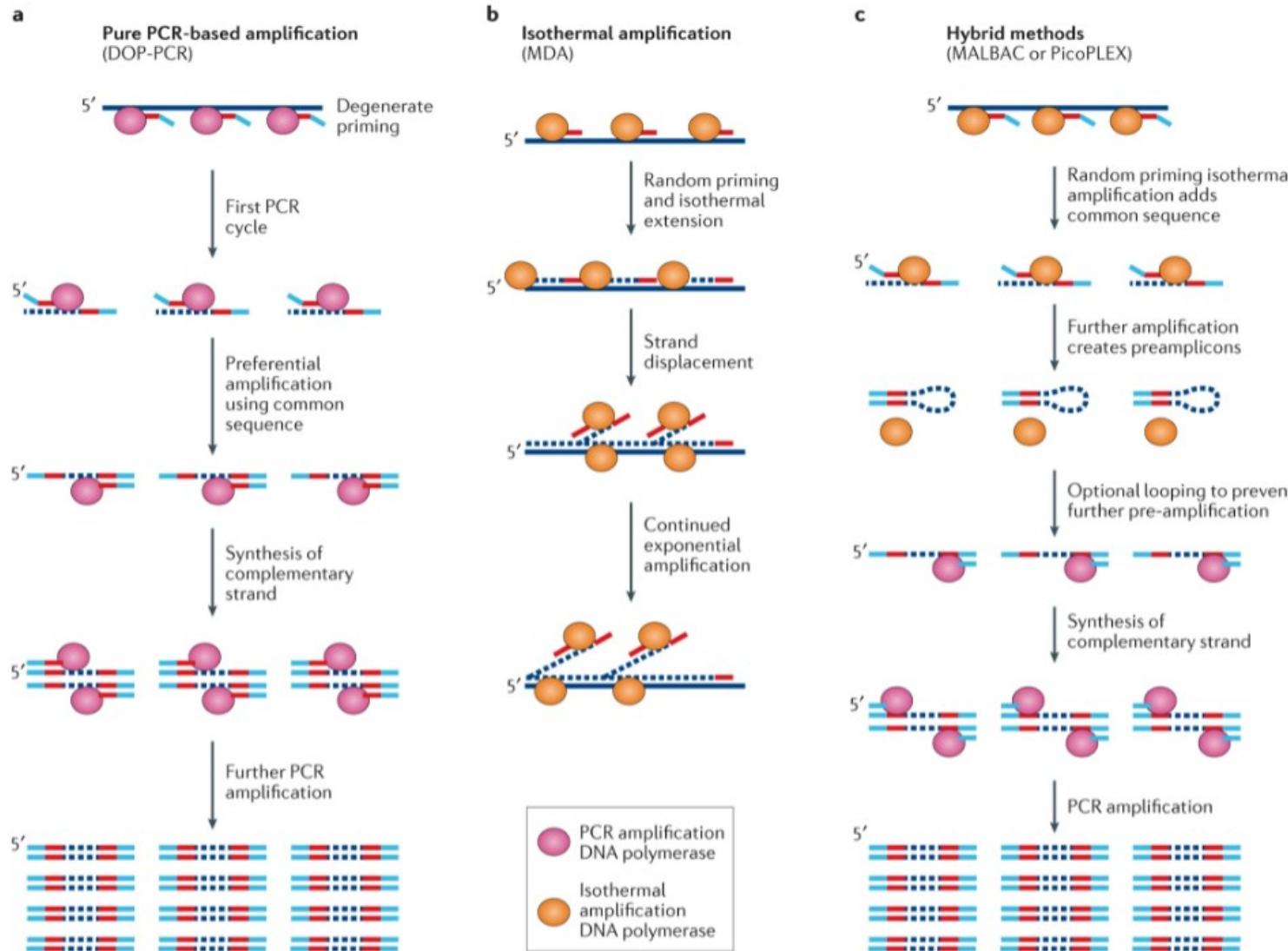
- Mutation classification is ambiguous
- There is a difference between damage or error of polymerase and mutation. Mutation is a DNA change that was fixated in next replication.
- There several proofreading systems that reduce mutation rate during replication to 10^{-10} per base per division
- methylation, and therefore epigenome, is important for error correction
- separation of somatic and germline cells is one more mechanism reducing mutation rate by lowering the number of replication cycles. It happens very earlier in embryogenesis.
- isolation of somatic cells from inheritance reduces requirements for genome stability and integrity
- Mutagenesis (even with prefix super-) might be a part of normal functioning of organism



I. Structure and diversity of the genomes

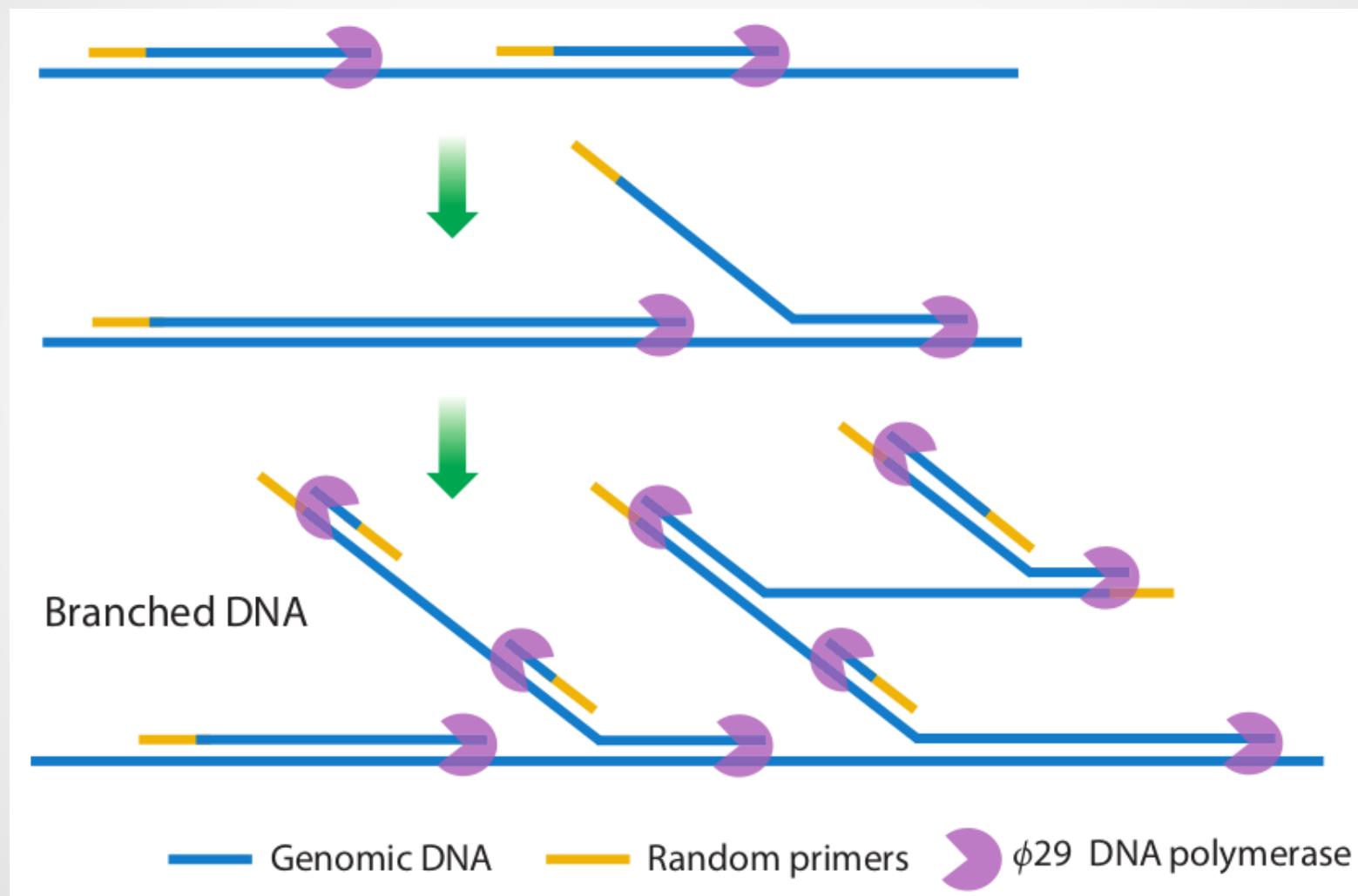
Genome “levels”

Could we assemble a genome from a single cell?



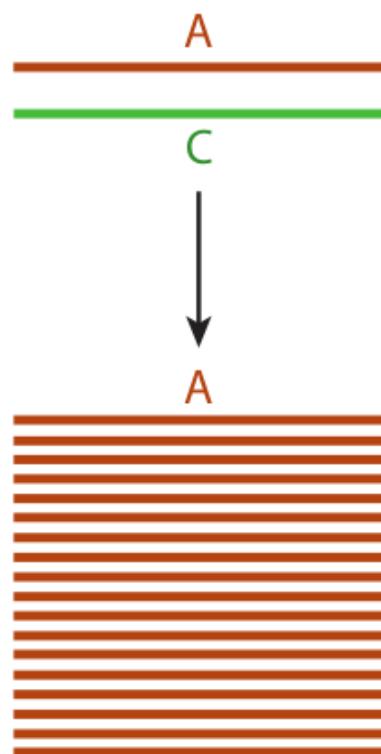
In general yes,
but we will need
to do a
**whole genome
amplification
(WGA)**
before
sequencing.

Multiple displacement amplification (MDA)

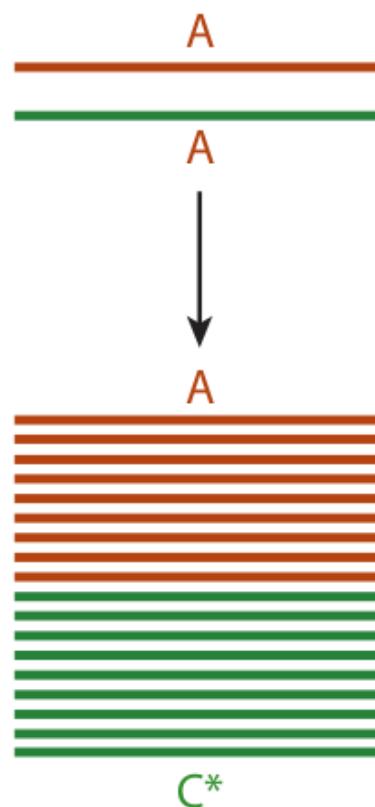


Typical errors of whole genome amplification

a Two alleles
in a single cell



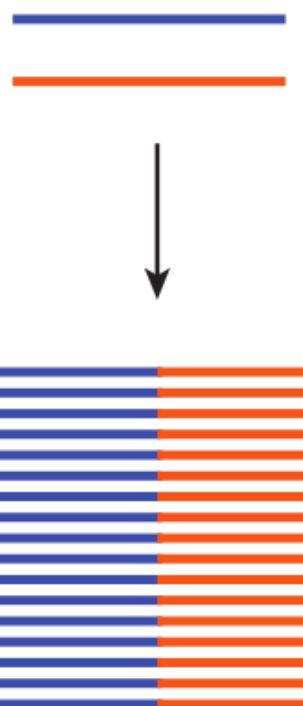
b Two alleles
in a single cell



Allele dropout

False positives

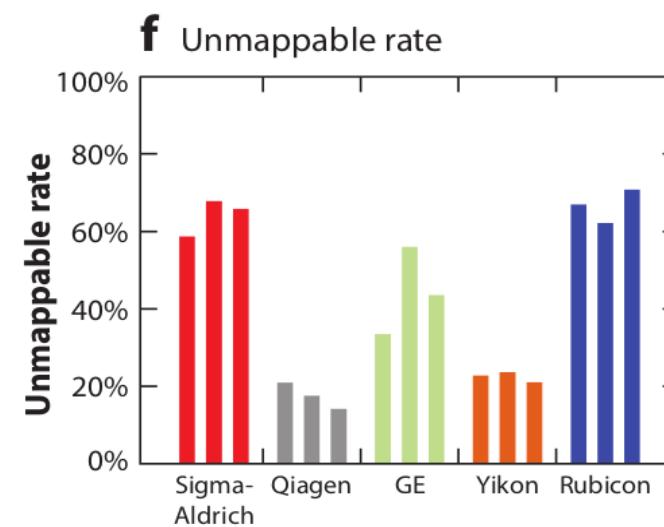
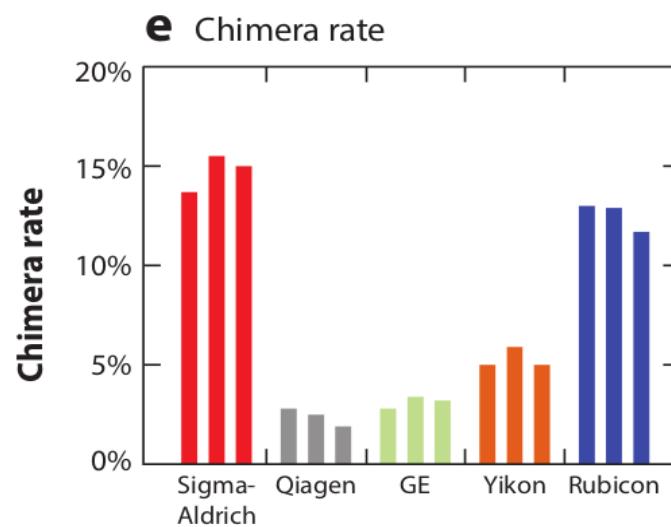
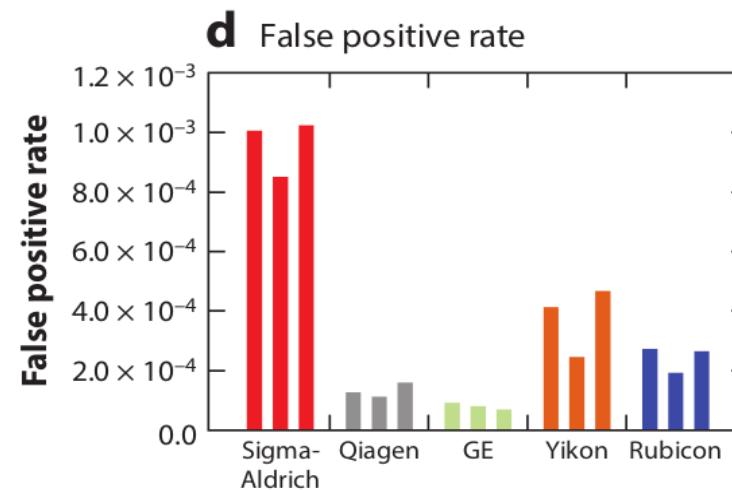
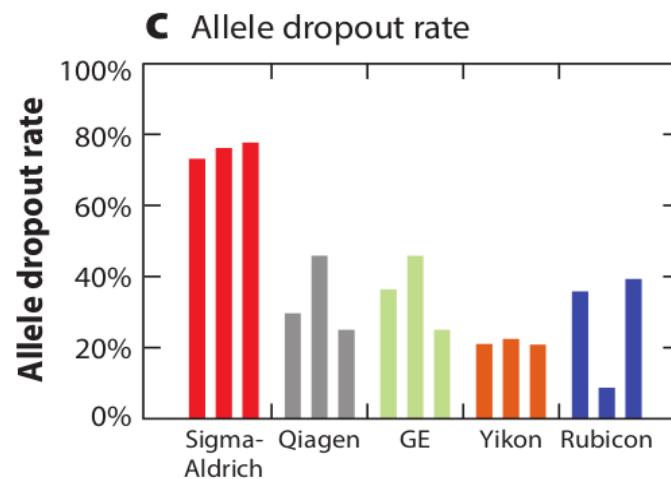
c Two separate DNA
segments in a single cell



Chimera formation

Some
regions
could
just ...
**not
amplify**

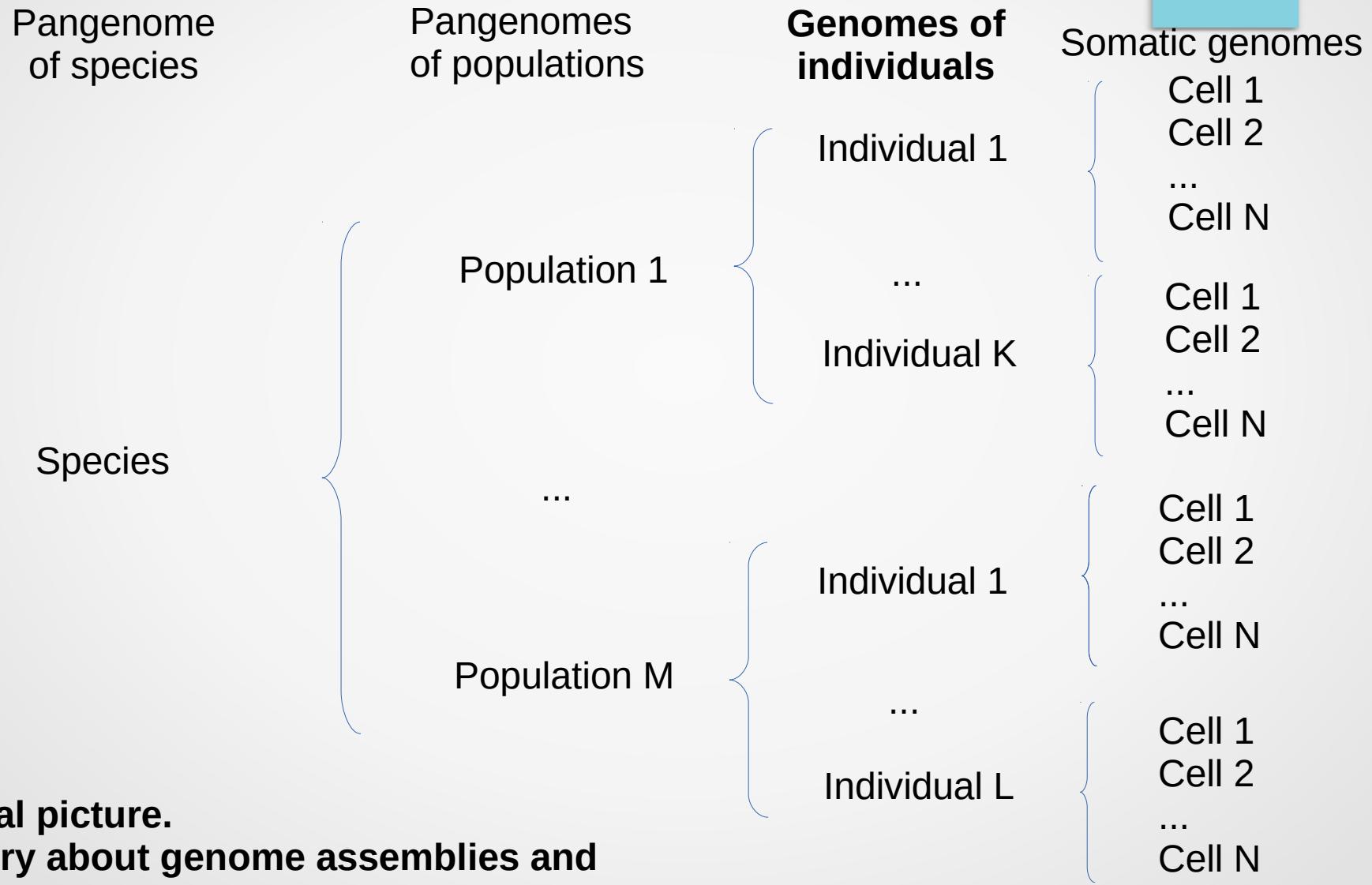
Efficiency of five commercial kits for WGA



Could we assemble a **reference genome** from a single cell?

No.

Genome and pangenome



Important features of genome assemblies

- source for DNA extraction - somatic cells:
 - blood or muscles
 - primary cell lines (usually fibroblasts)
- we try to approach original genome of zygote by simultaneous sequencing genomes of multiple somatic cells and assembling the consensus.
- in ideal case all sequencing should be done from the same individual by often it is impossible. So often, very often, genome assembly is chimera created from multiple individuals.
- Most of genome assemblies are haploid, i.e. their sequences are haploid consensus between maternal and paternal haplotypes

Human reference genome GRCh38

GRC38 (latest release GRCh38.p14) includes:

- **C-scaffolds(chromosomes) and mtDNA**
- **unlocalized scaffolds (unknown exact position, but known chromosome)**
- **unplaced scaffolds (completely unknown position)**
- alternative scaffolds
- Epstein-Barr virus sequence and additional decoy sequences

GRC38 is something intermediate between haploid assembly and graph assembly

Pangenome projects

Human pangenome

<https://humanpangenome.org>

Rice pangenome

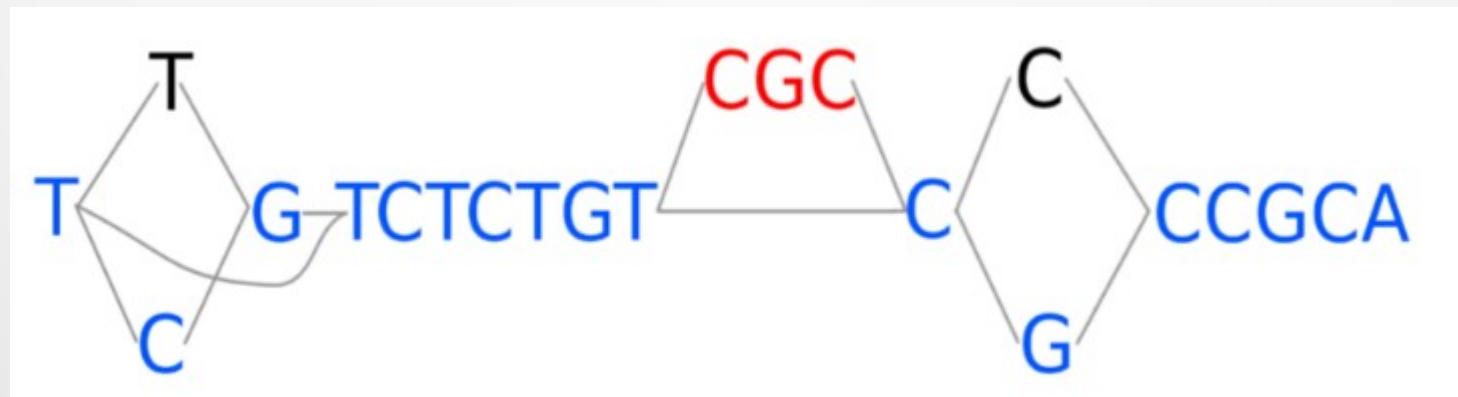
<https://cgm.sjtu.edu.cn/3kricedb/>

Maize pangenome

<https://maize-pangenome.gramene.org>

Tomato pangenome

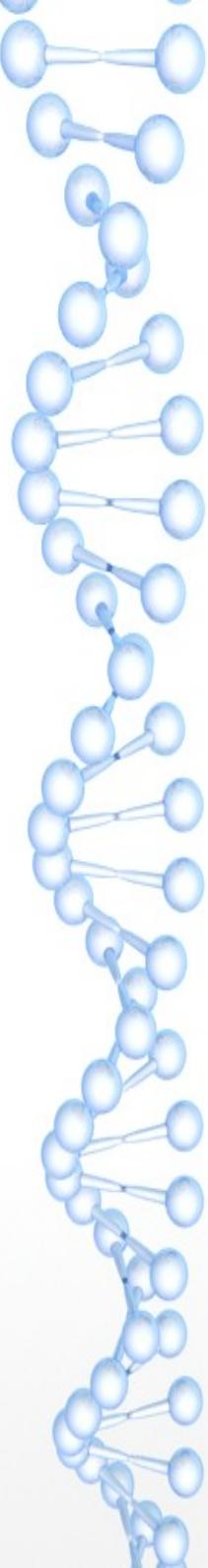
<https://solgenomics.net/projects/tgg>



scheme of graph-based references

Summary

- genome assembly of big multicellular organism nearly always is an approximation to the genome of the zygote (from which organism was developed) based on sequencing of multiple somatic cells and assembling the consensus. In some cases assembly might be even a consensus derived from multiple individuals.
- Under genome of species we usually mean a genome of a single reference individual. You should remember about this bias.
- Some assemblies like human GRCh38 could include alternative contigs to represent differences between different populations. Also additional sequences (absent in genome) might be present to improve read mapping.
- There is a trend of switching from haploid flat assemblies to diploid and graph-based assemblies



I. Structure and diversity of the genomes

Definition of the term “genome”

Several commonly used definitions

Genome is

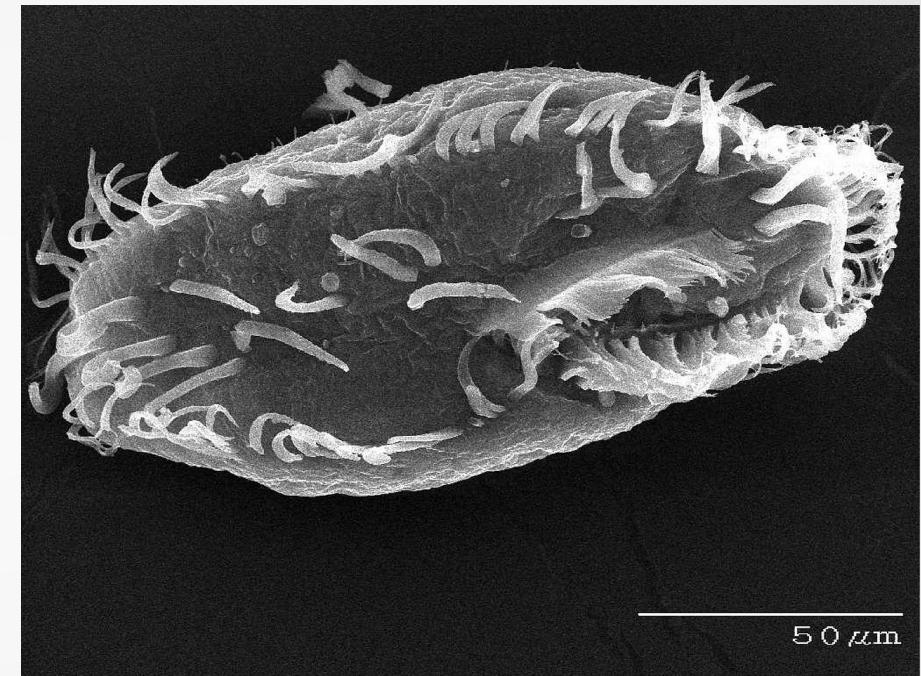
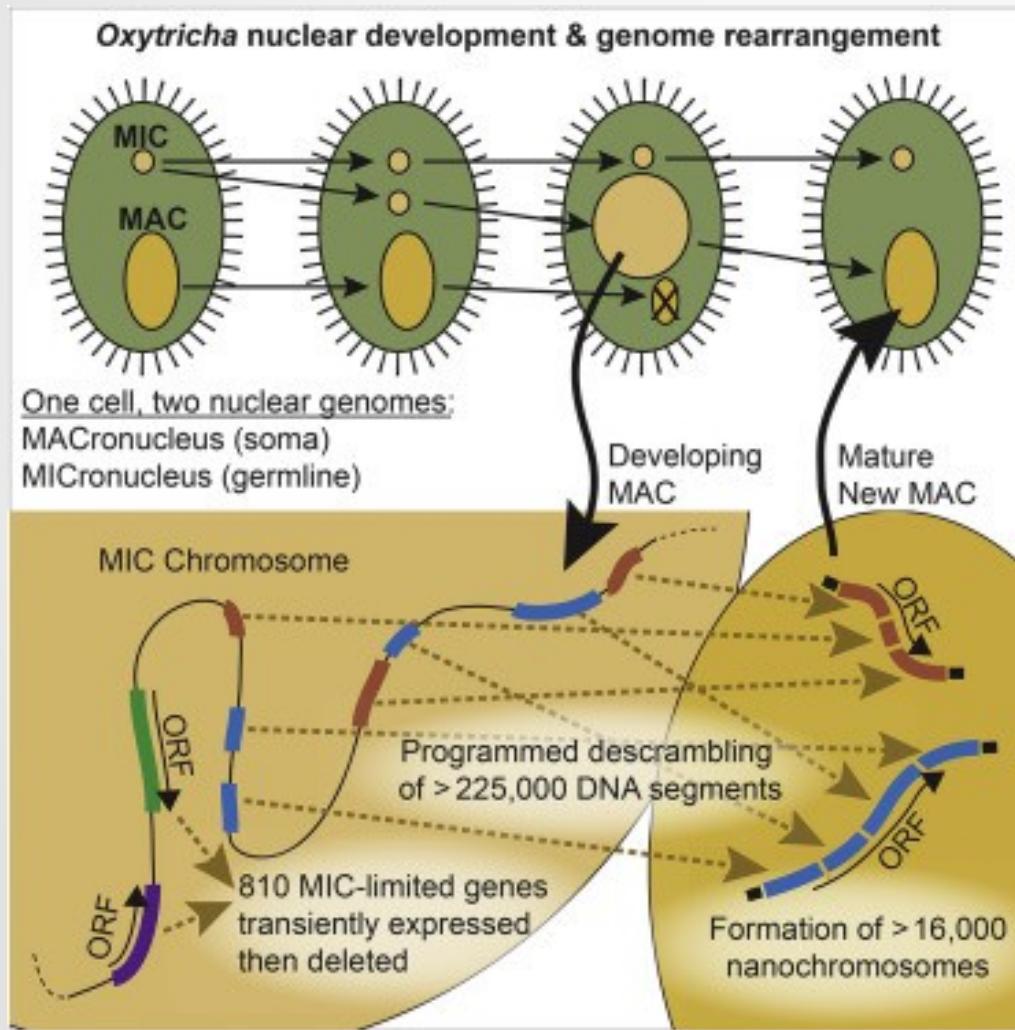
- the entire set of **DNA** instructions found in a **cell**.
- genetic material of haploid set of **chromosomes**
- **information repository** of organism
- **all hereditary material** of organism

include epigenome
(methylation and other modifications)

Definitions overlap significantly, but not completely!

RNAs and genome: case of Ciliates

Sterkiella histriomuscorum (or *Oxytricha trifallax*)



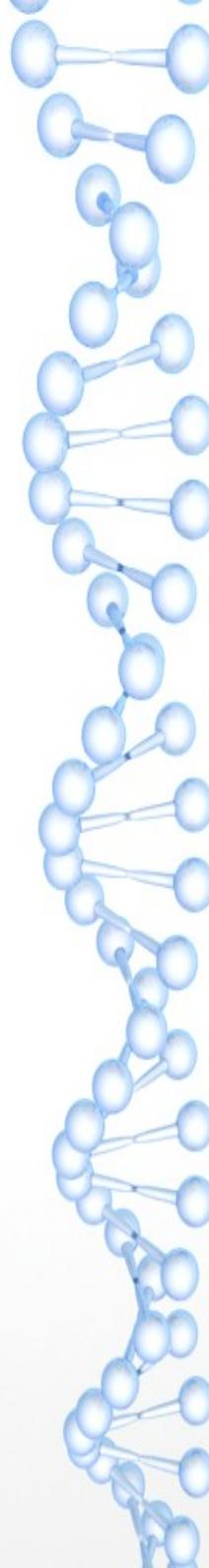
Chen et al, 2014

Information about the order of segments in scrambled genes is stored in RNA copies of nanochromosomes during at least one stage

Definitions of term genome

- **information repository** of organism
- **all hereditary material** of organism

Two definitions pretending to encompass all the variety of genomes.
But still there are issues.



End of Module I