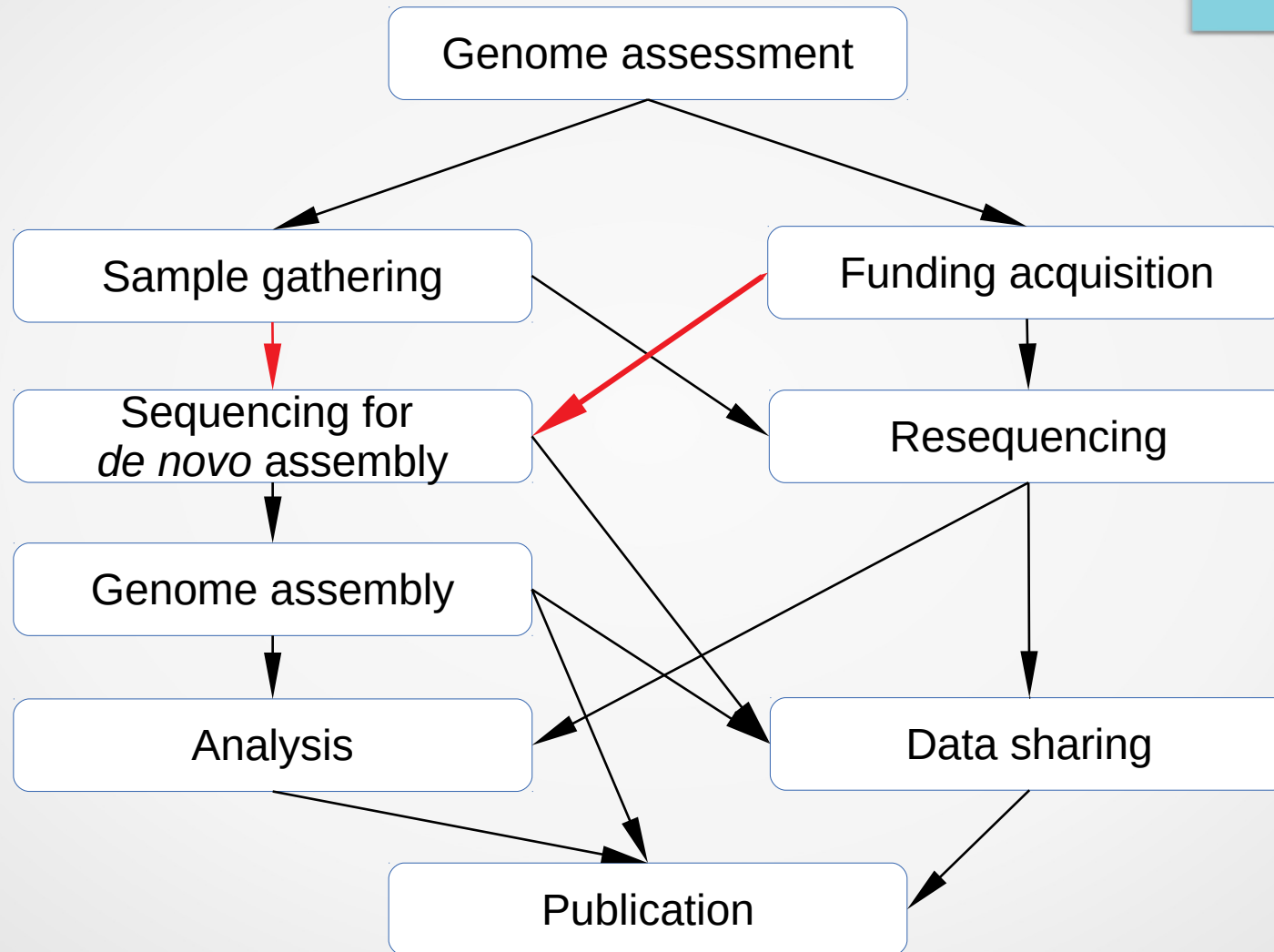
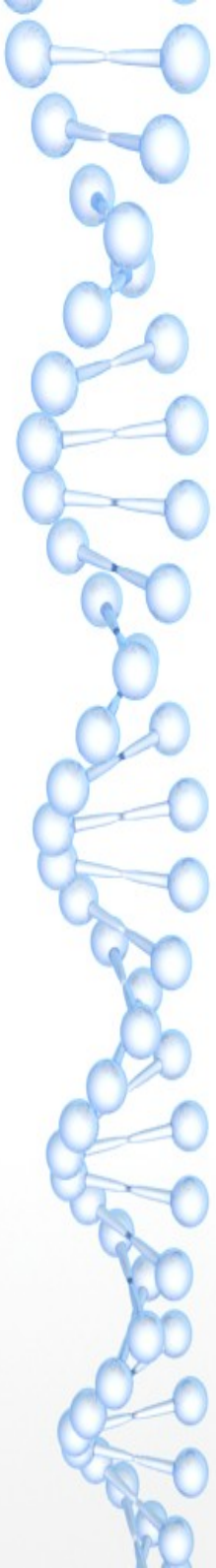


## III. Genome Projects

# Stages of the project





# III. Genome Projects

## Genome assessment

# Questions you need to answer before starting project

- I. How big is the genome?
- II. Does it have heterochromatin?
- III. Is it diploid or polyploid?
- IV. Is it highly heterozygous or not?
- V. Is your species hybrid or not?
- VI. Are the genome rearrangements widespread in your species?

# Genome size of your species

## Ways to assess

### I. Check literature and databases

for animals	<a href="https://www.genomesize.com">https://www.genomesize.com</a>
for plants	<a href="https://cvalues.science.kew.org">https://cvalues.science.kew.org</a>
for fungi	<a href="http://www.zbi.ee/fungal-genomesize/">http://www.zbi.ee/fungal-genomesize/</a>

### II. Estimate using flow cytometry

- require cell line of your species
- require reference (at least two) species with known genome size
- require cytometer and cytogenetisist

### III. Estimate from reads using k-mer distribution

- require preliminary sequencing

C-value (haploid genome size) might be in picograms and Mbp.

**1 pg = 978 Mbp**

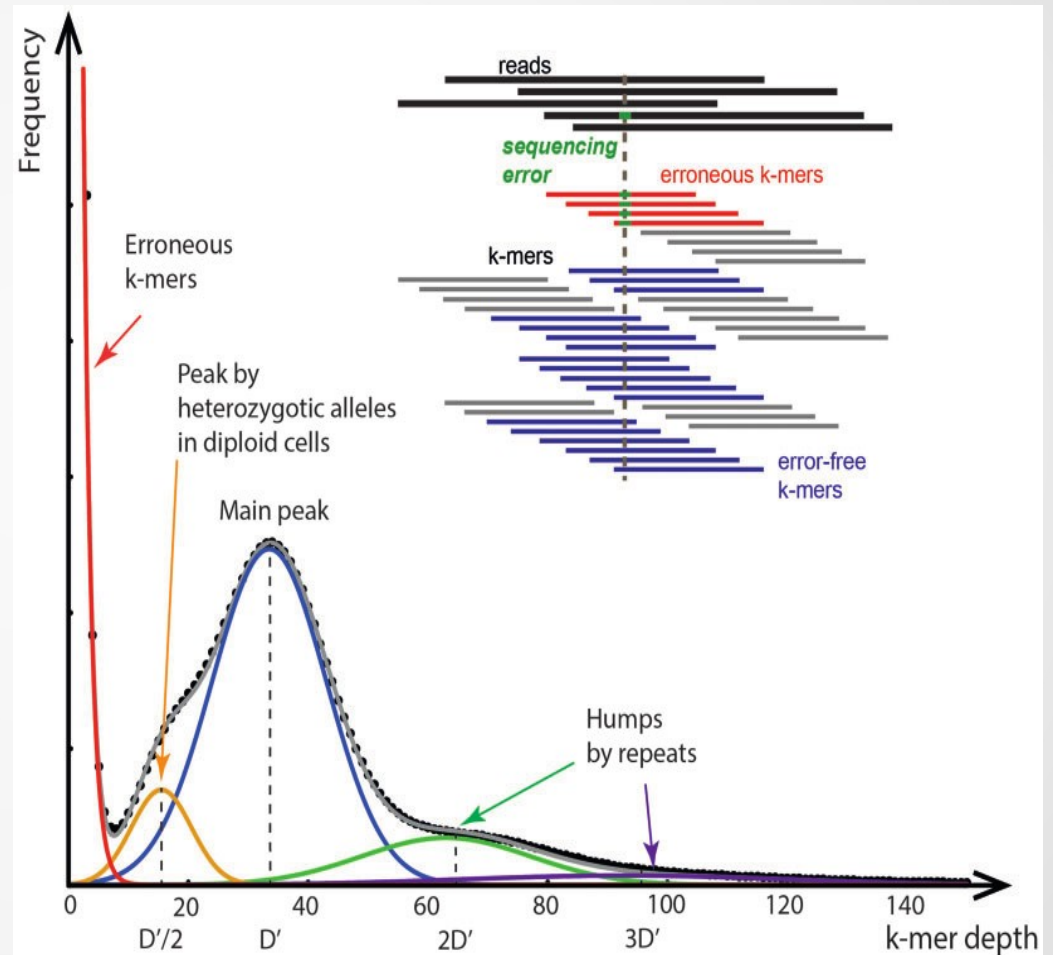
# Genome size estimation from reads (1)

## K-mer based approach:

- I. Count k-mers and create database
- II. Count histogram
- III. Assess genome size from histogram

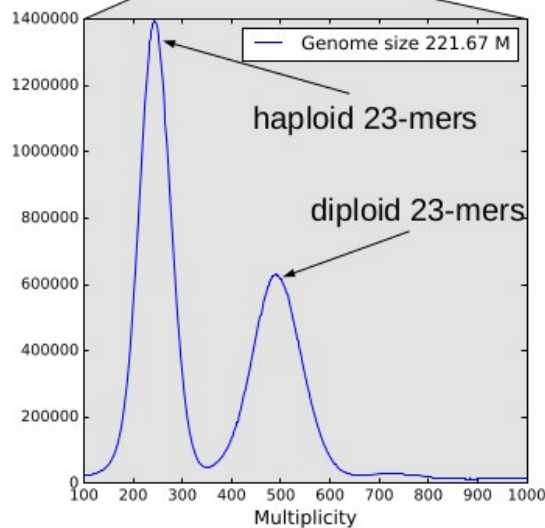
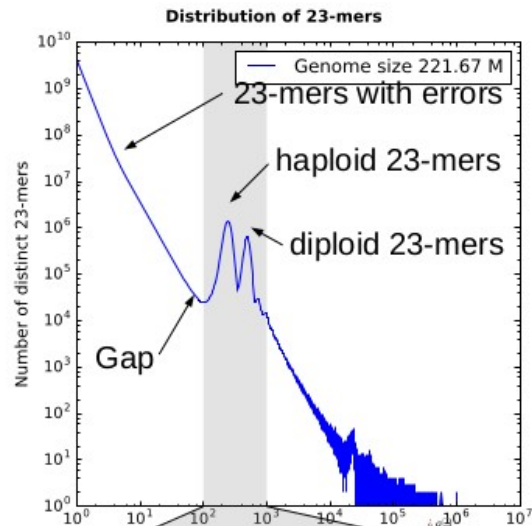
## Tools:

- I. Jellyfish 2
- II. Jellyfish 2
- III. Genomescope 2, KrATER, etc



Sohn and Nam, 2016

# Genome size estimation from reads (2)



**23-mer distribution for  
PE reads of hybrid plant *B. divaricarpa***

## Genome size estimation

$$\text{Genome size} = \frac{\sum_{i=g} N_i * m_i}{C}$$

*g* – *k* – mer multiplicity at gap between peak of *k* – mers with errors and corresponding to unique part of genome

*m<sub>i</sub>* – multiplicity of distinct *k* – mers

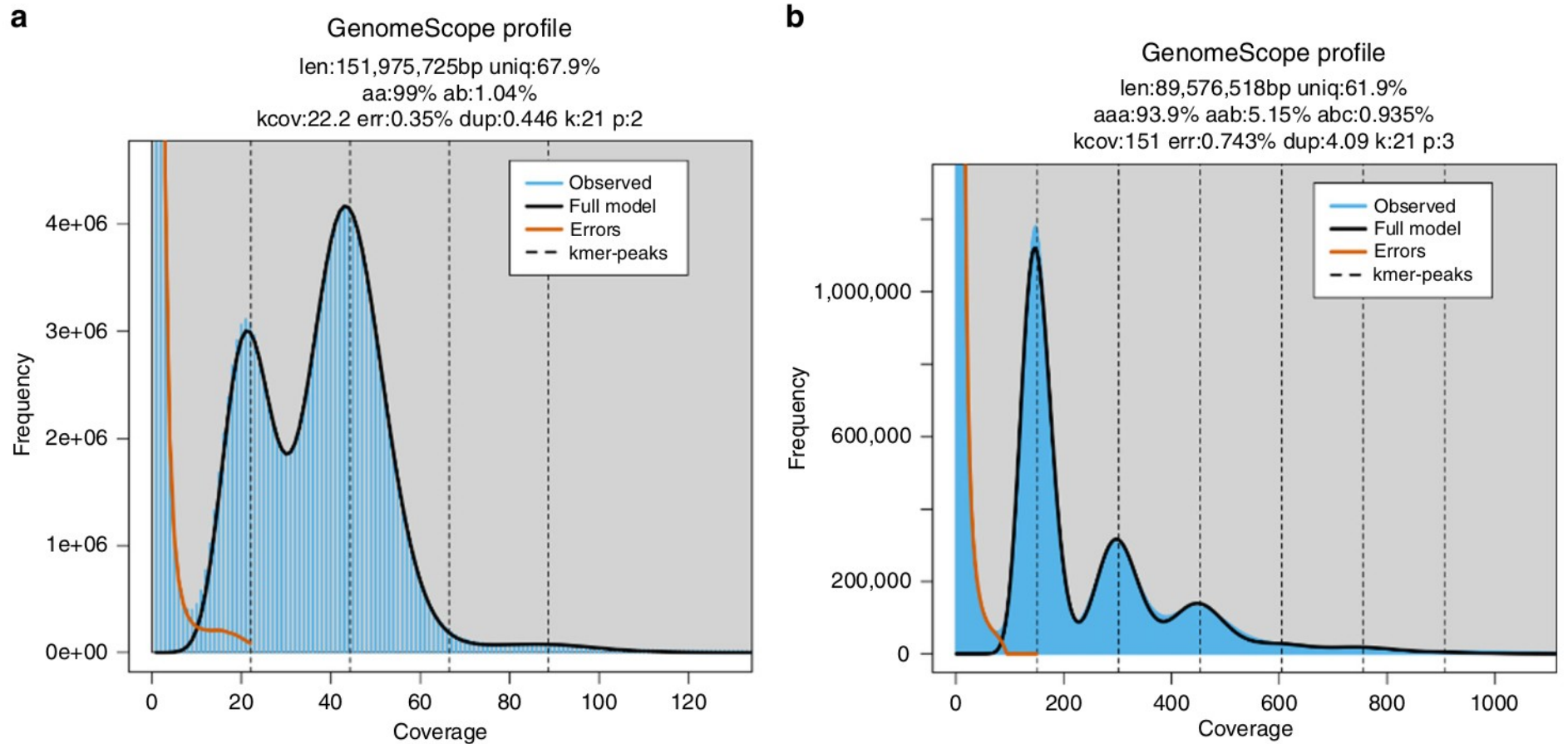
*N<sub>i</sub>* – number of distinct *k* – mers with *m<sub>i</sub>* multiplicity

*C* – sample coverage estimated by mode of multiplicity of diploid *k* – mers

## Naive approach:

- direct assessment from histogram
- works always if there distinguishable peaks
- huge error

# Genome size estimation from reads (3)



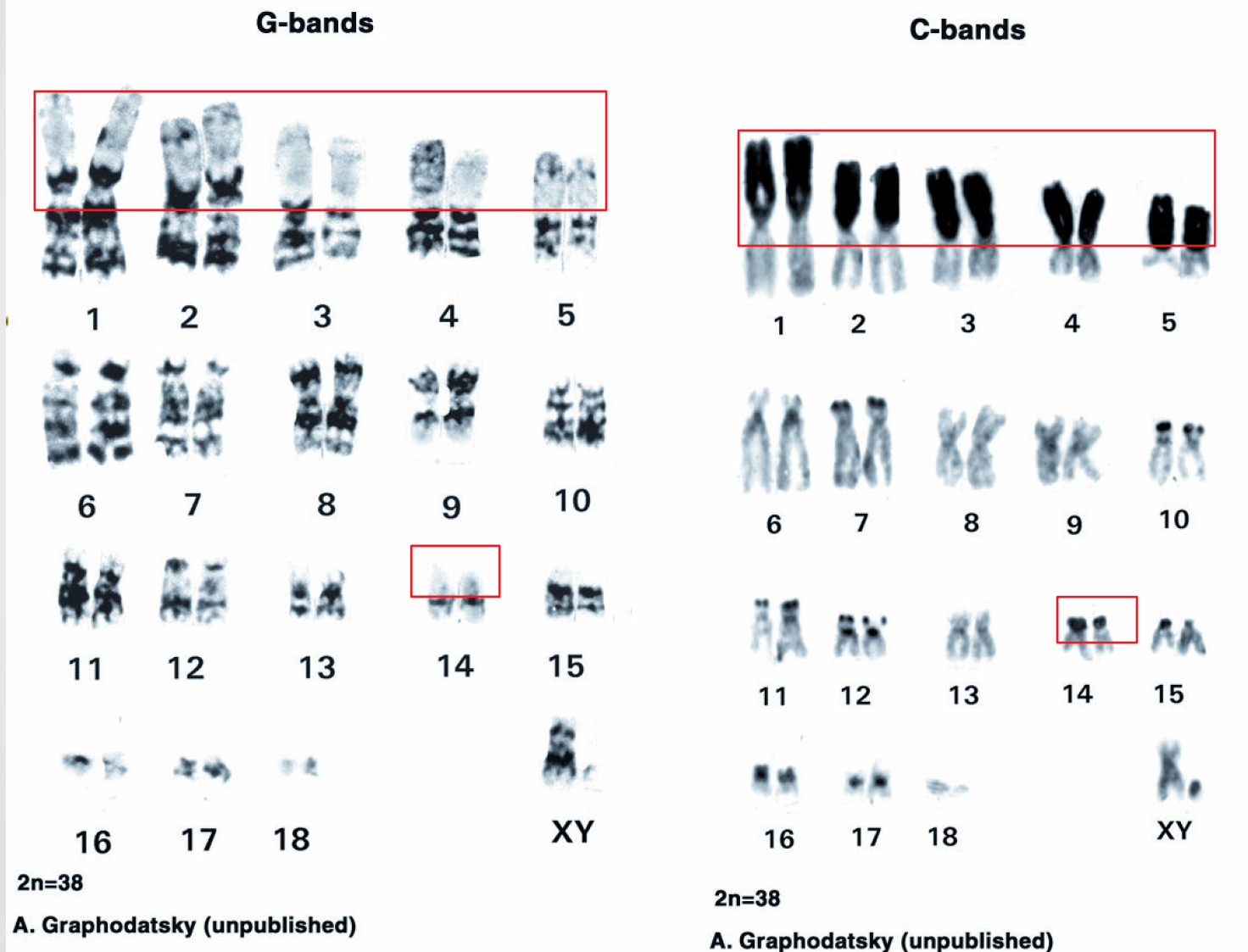
## Modeling approach:

- assessment of genome size from histogram by fitting it with model (sum of several negative binominal distributions)
- works not always
- more precise

Ranallo-Benavidez et al, 2020

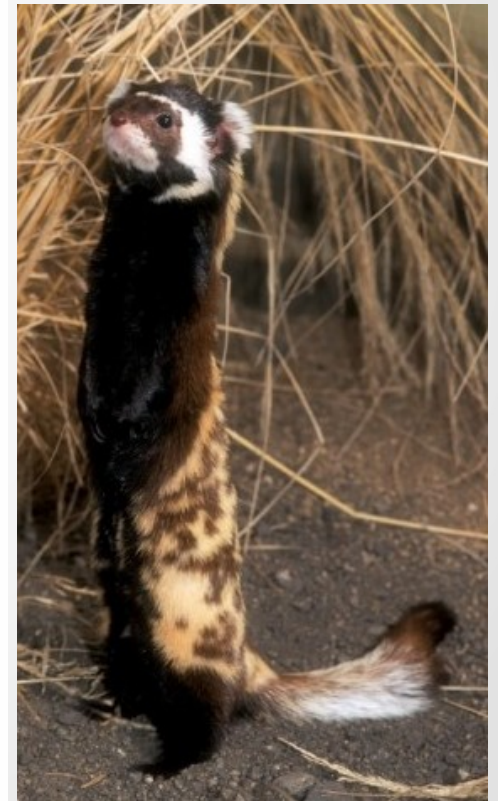


# Heterochromatine and repeat content



*Vormela peregusna*

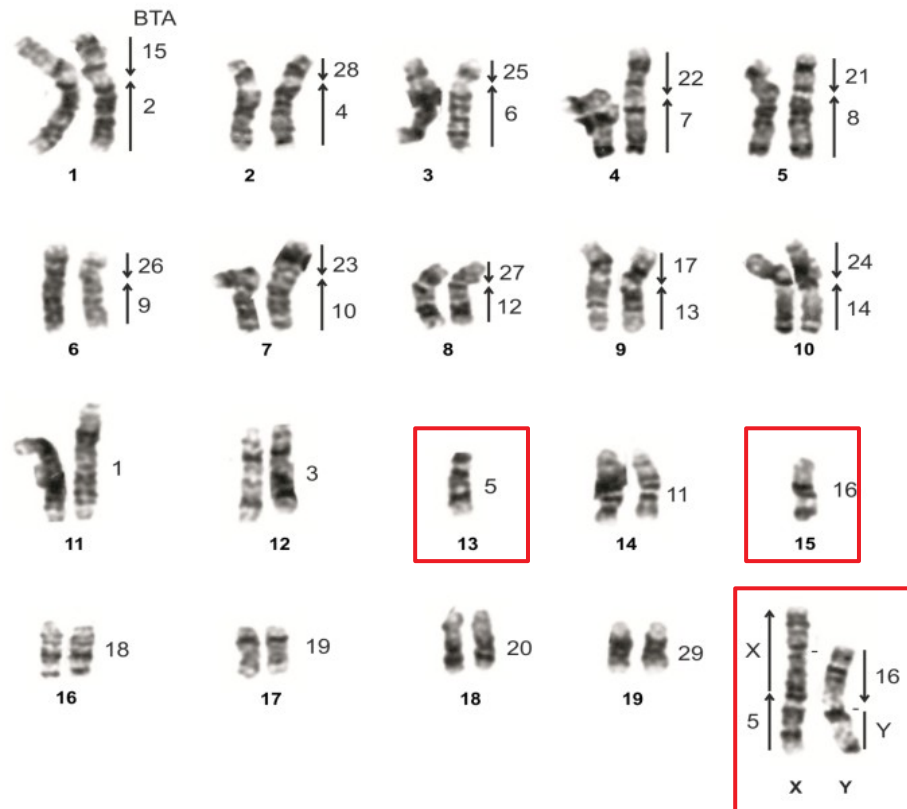
Marbled polecat



~ 6.4 Gbp (!)

# Widespread chromosomal rearrangements

## *Nanger (Gazella) dama mhorror* (Mhorror Gazelle)

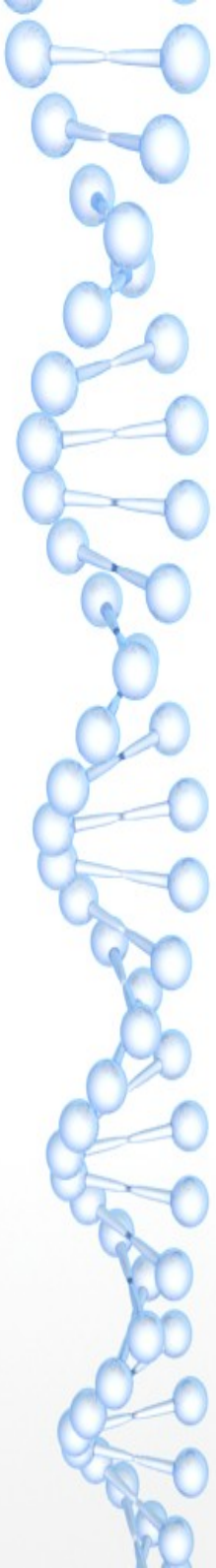


$2n=38-40$

Center for Reproduction of Endangered Species (CRES)  
Zoological Society of San Diego

Probes: *Bos taurus* (BTA)  
Cernohorska et al. (2012)

Graphodatsky et al, 2020



# III. Genome Projects

## Samples

# Fragment size of extracted DNA

## HMW DNA

high molecular weight DNA

~ 50 - 300 kbp fragments

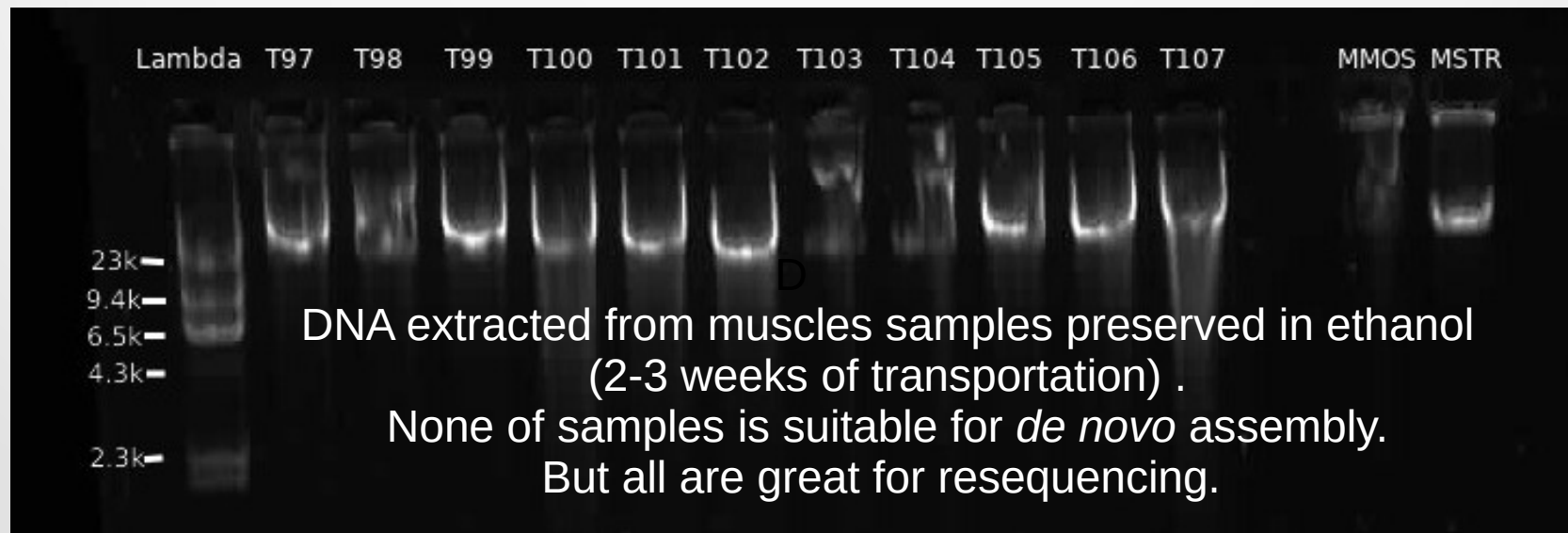
## UHMW DNA

ultrahigh molecular weight DNA

~300+ kbp fragments

**Modern approaches for *de novo* assembly require at least HMW DNA.**

**Not all sample types and DNA extraction methods could produce HMW DNA!**



# DNA sources

		Fragments, bp
<ul style="list-style-type: none"><li>• Cell lines</li><li>• Fresh tissue samples<ul style="list-style-type: none"><li>• blood</li><li>• biopsy</li><li>• necropsy</li></ul></li></ul>	Suitable for <i>de novo</i> assembly	up to 500000
		up to 300000
<ul style="list-style-type: none"><li>• Preserved tissue samples</li><li>• Secretions<ul style="list-style-type: none"><li>• saliva</li></ul></li></ul>		up to 50000 - 100000 ~100 - 10000
<ul style="list-style-type: none"><li>• Museum samples<ul style="list-style-type: none"><li>• skins</li><li>• bones</li></ul></li></ul>		~50-200
<ul style="list-style-type: none"><li>• Ancient samples<ul style="list-style-type: none"><li>• bones</li></ul></li></ul>		~25-200

# Primary cell lines

## Advantages:

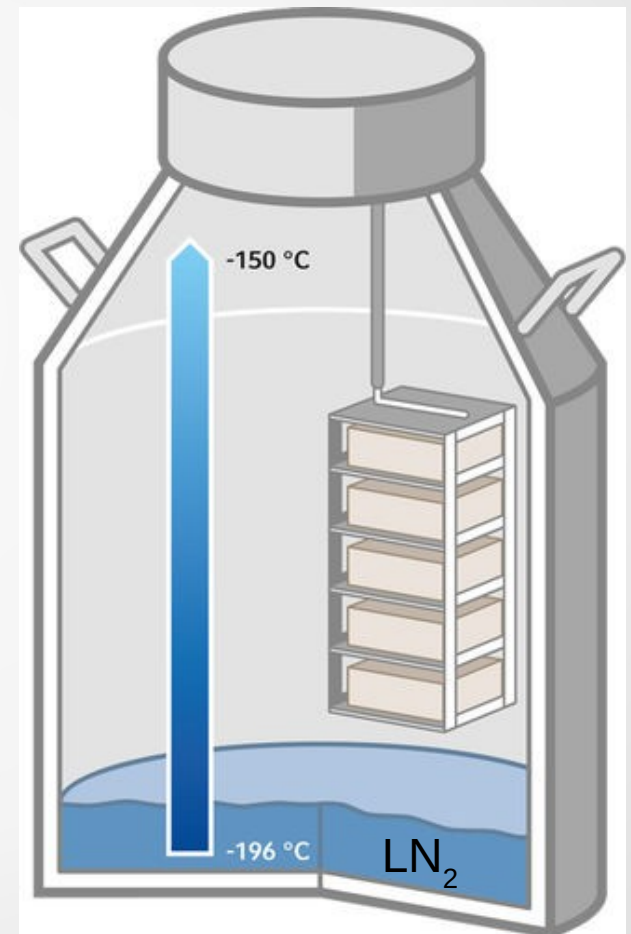
- Available living cells
- Allows cytogenetic experiments
- Source for DNA of excellent quality
- Could be stored for decades

## Disadvantages:

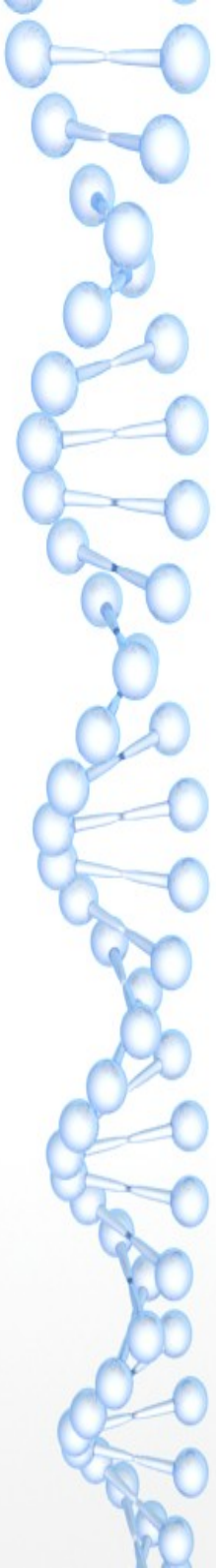
- Require cytogenetic lab and cytogenetisists
- Strict requirements for samples to establish cell line (cell must be alive)
- Storage in liquid nitrogen (LN<sub>2</sub>)
- Will stop growing after specific number of divisions

**Immortalized cell lines are not suitable for *de novo* genome assembly of new species!**

cryogenic storage dewar







# III. Genome Projects

## Genome assembly

# Major definitions related to genome assemblies

## Read

small fragment generated by sequencing

## Contig

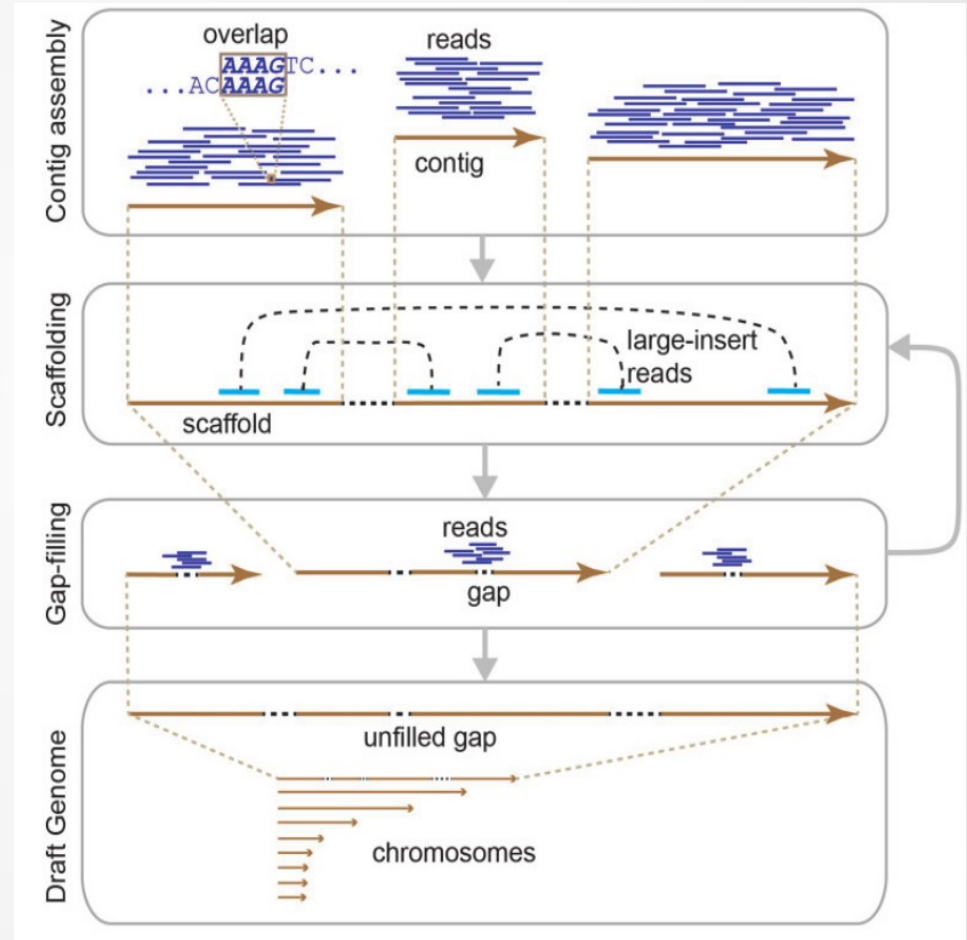
ungapped sequence assembled from reads

## Scaffold

sequence with gaps, generated from contigs (set of oriented contigs)

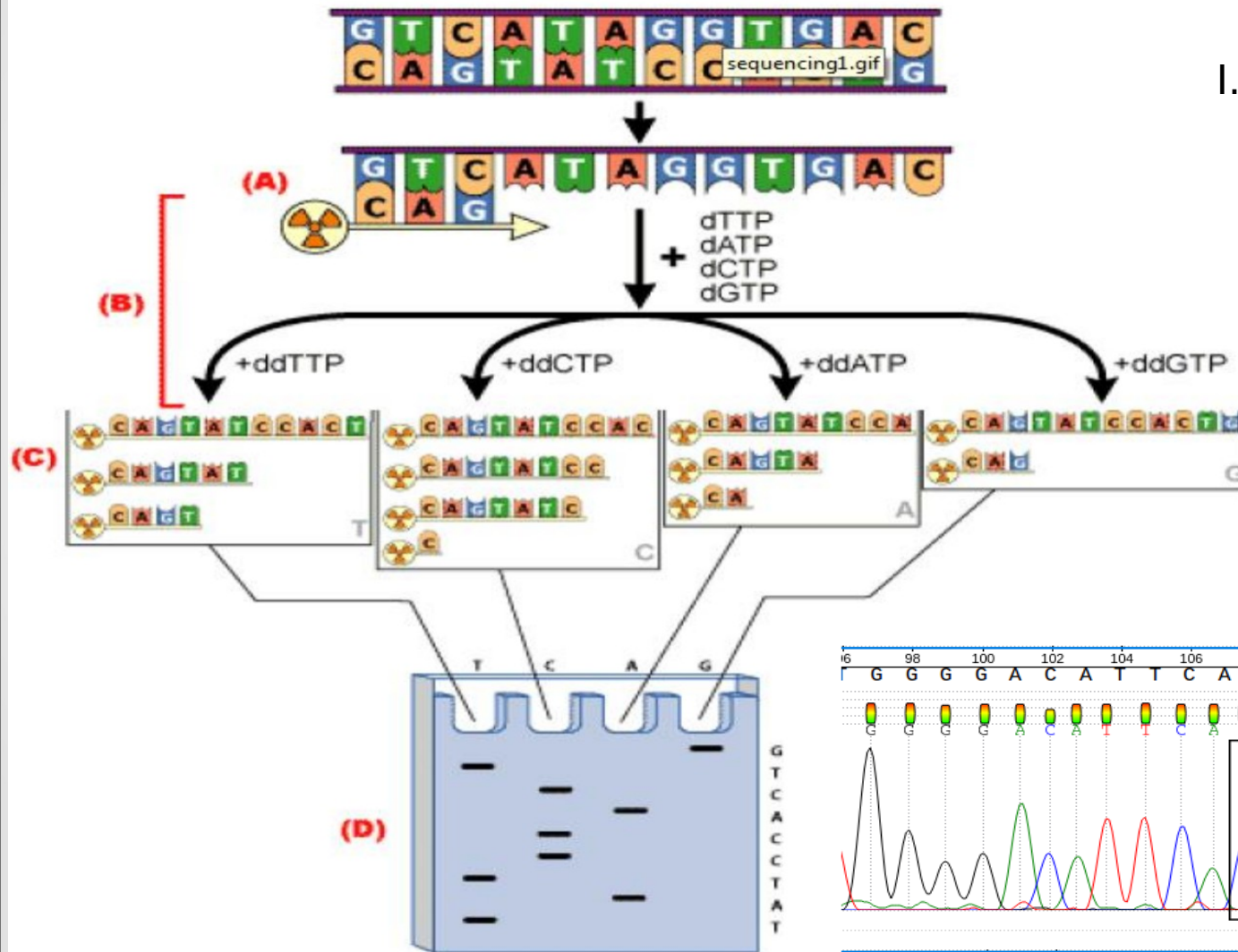
## C-scaffold

chromosome (or close to) scaffold, representing whole chromosome or its significant part





# Sequencing technologies: first generation

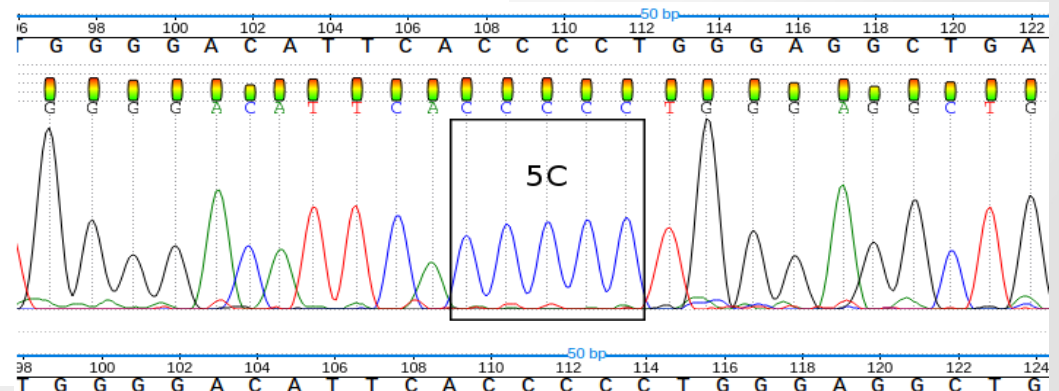


I. Maxam-Gilbert sequencing

- R.I.P

II. Sanger sequencing

- still alive
- used in validation and preliminary analysis
- reads up to 800 bp



# Sequencing technologies: second (next) generation

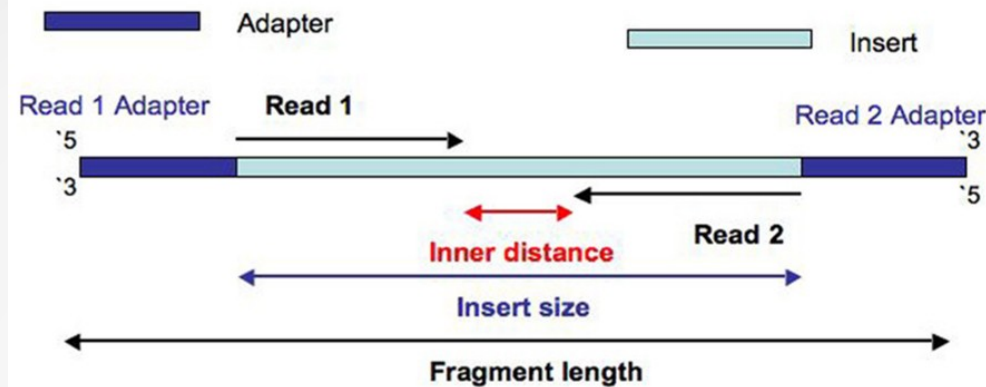
Platform	Max read length, bp	
• SOLID • R.I.P.	35	
• Roche 454 (pyrosequencing) • R.I.P.	400	Issues with sequencing of homopolymers
• IonTorrent • R.I.P. for <i>de novo</i> assembly	400	
• Illumina	250	
• MGISEQ	300	

**All NGS platforms are based on sequencing-by-synthesis (SBS)  
and can't sequence a single molecule!**

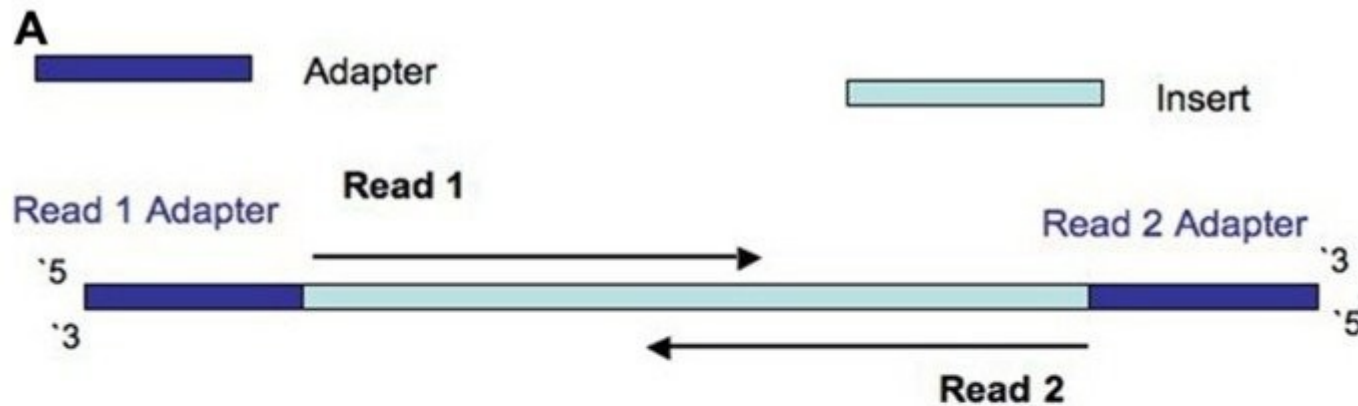
Major differences between platforms are related to what is detected during synthesis and how amplification is performed.

# Read length and insert size

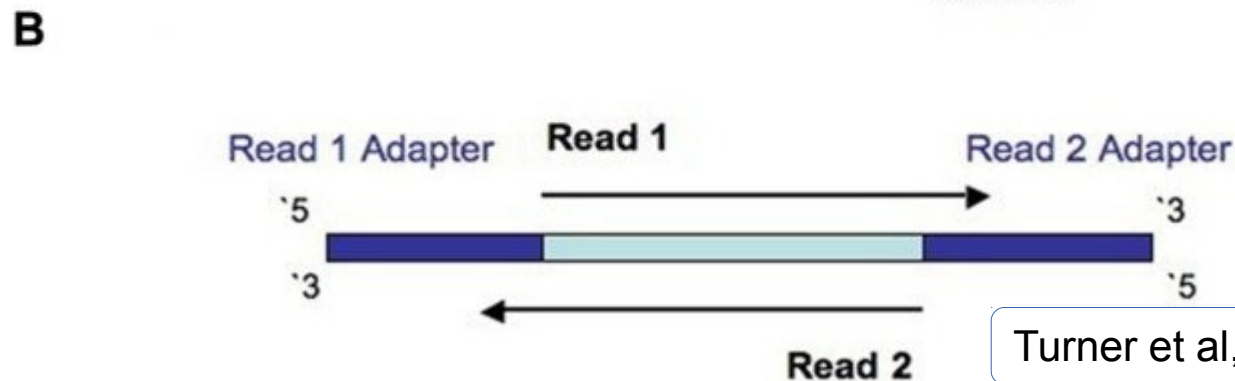
IS - insert size  
RL - read length



$$IS > 2 * RL$$



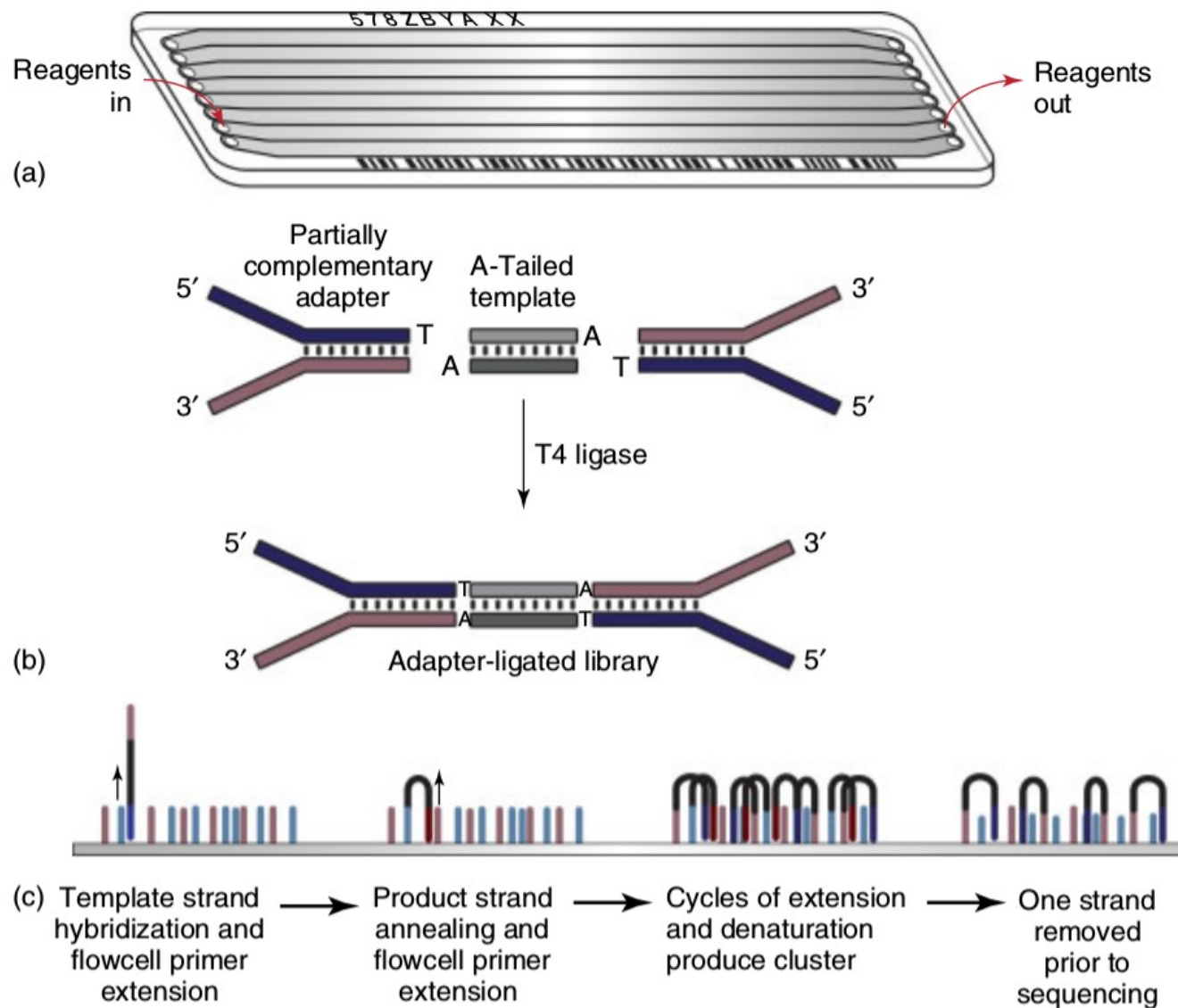
$$RL < IS < 2 * RL$$



$$IS < RL$$

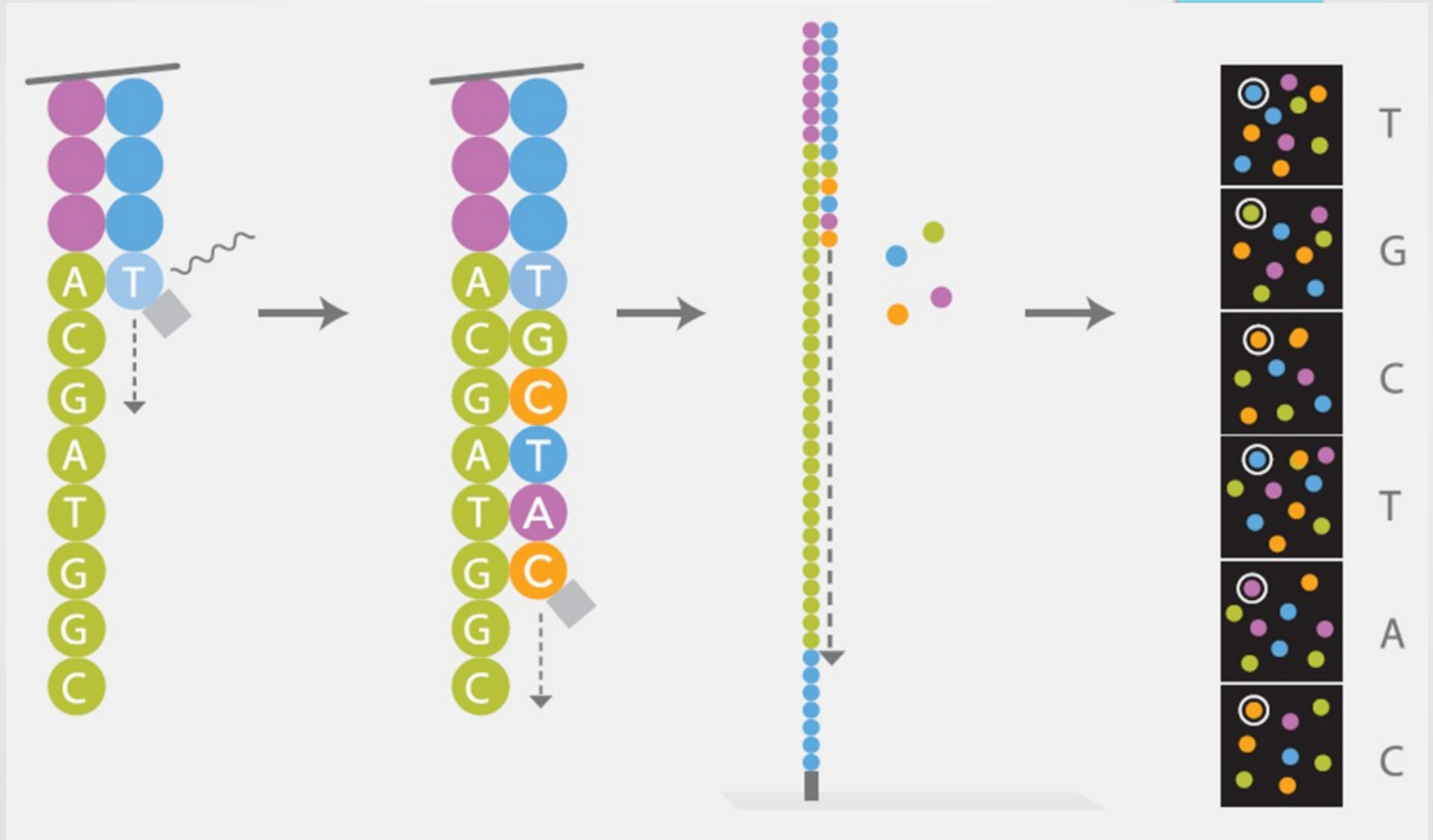
Turner et al, 2014

# Illumina platform



Turner et al, 2014

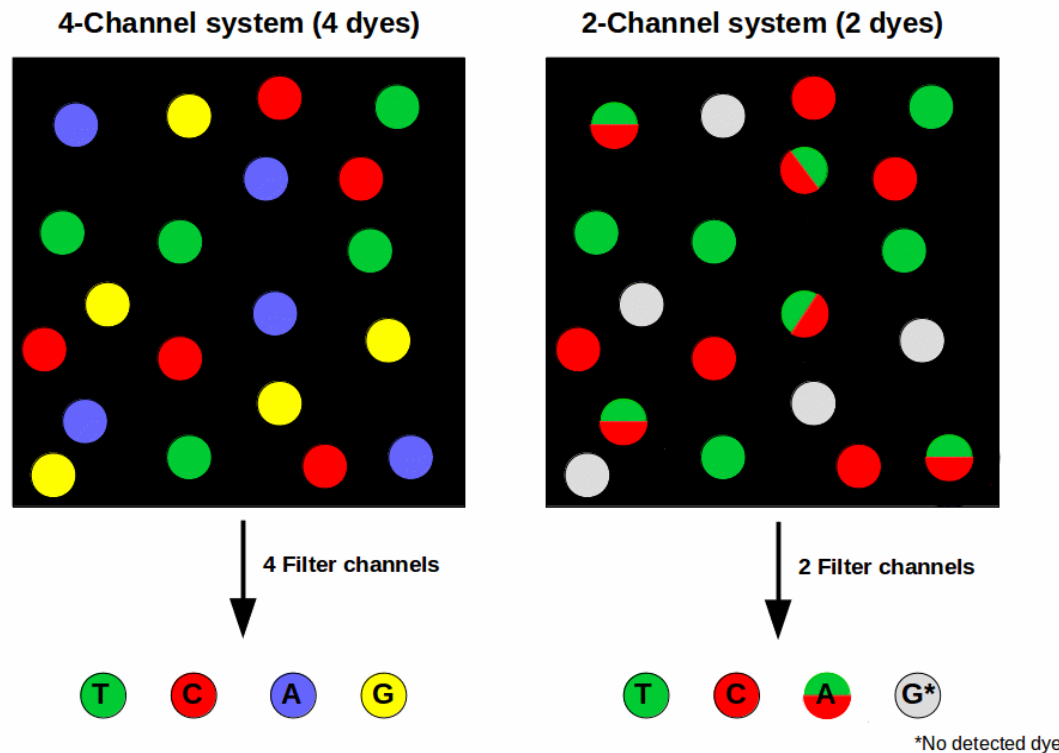
# Sequencing-by-synthesis on Illumina platform



Ansorge et al, 2009

# Illumina platform is heterogenous

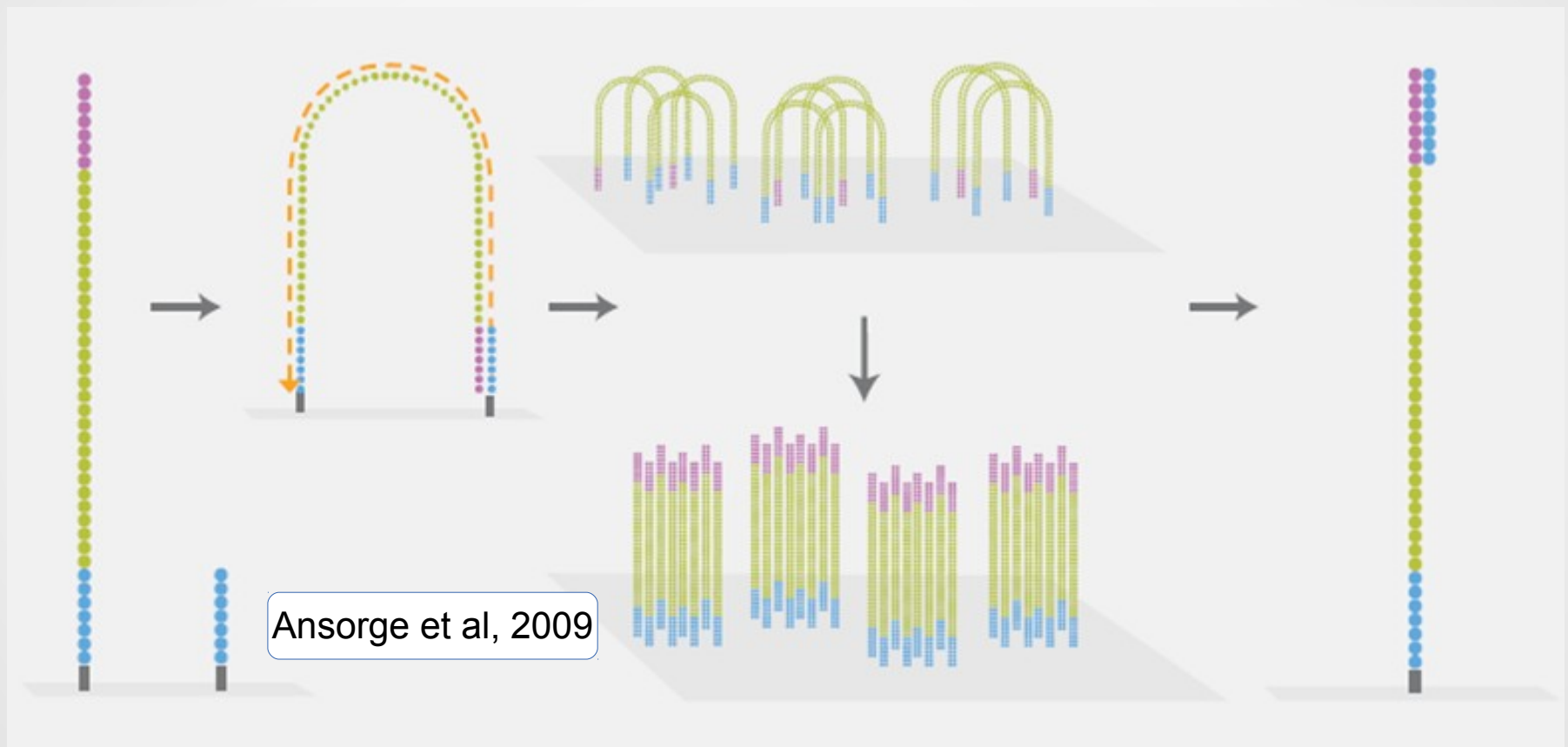
GAII  
HiSeq2000  
HiSeq2500  
HiSeq3000  
HiSeq4000



NovaSeq  
NextSeq



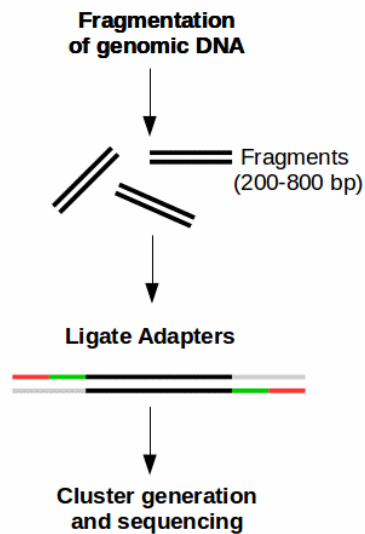
# Bridge amplification on illumina platform



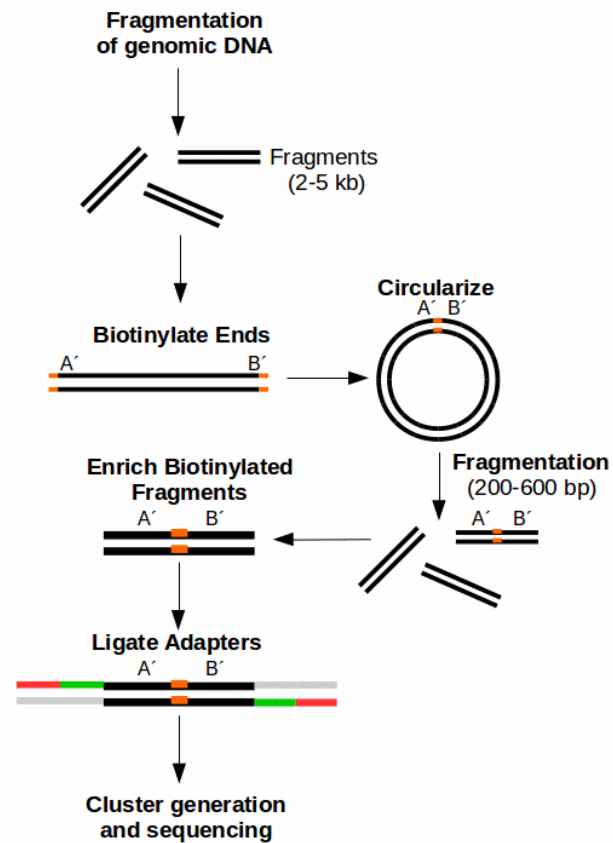
Bridge amplification doesn't work for fragments longer ~ 1500 bp.  
It is a maximal threshold for insert size (IS), and it is difficult to achieve.  
Commonly use IS are 250, **350** and 550

# Special types of libraries: Mate pairs

## Paired-End Sequencing (Short-insert paired-end reads)

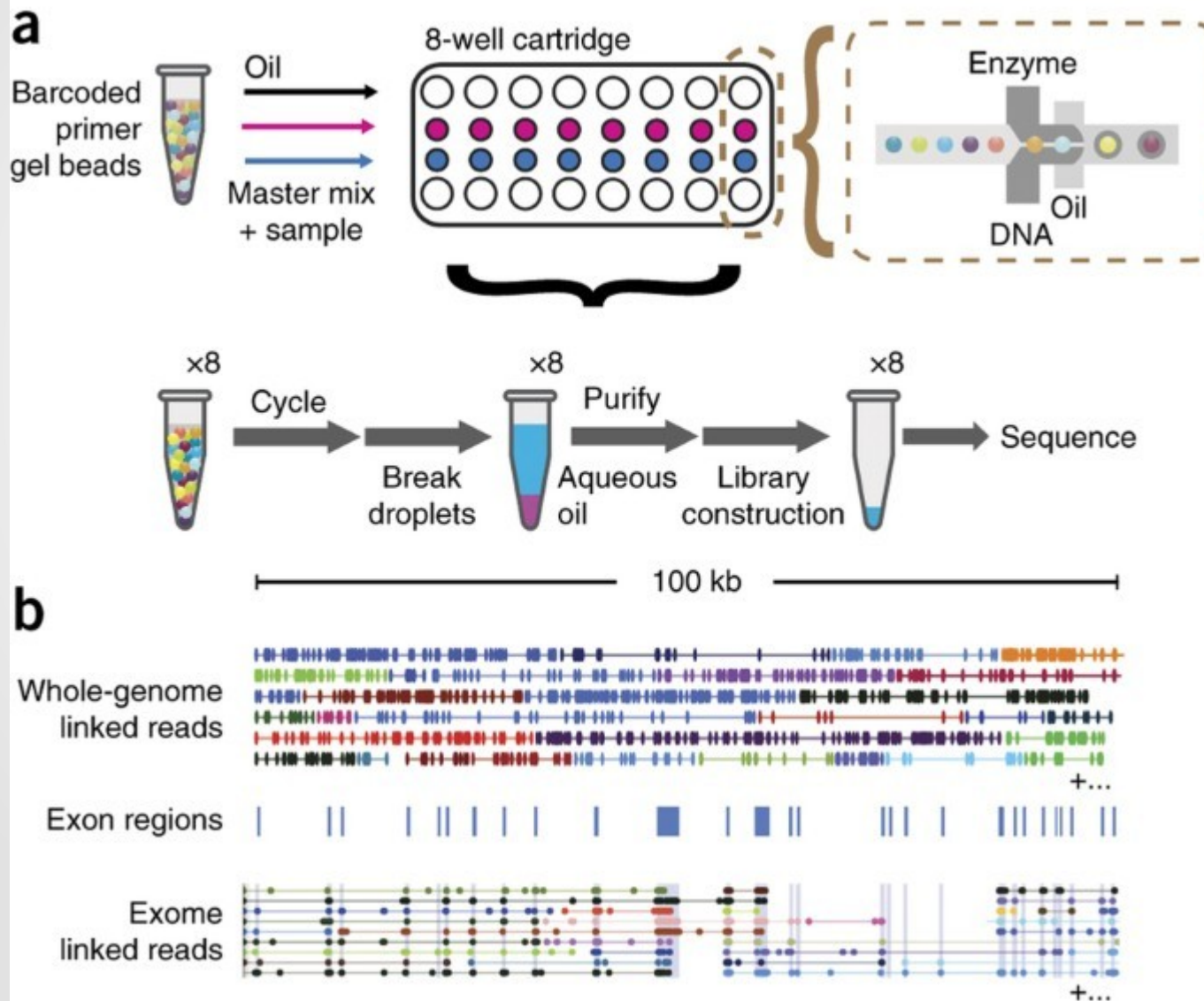


## Mate Pair Sequencing





# Special types of libraries: linked reads



Microfluidics based:

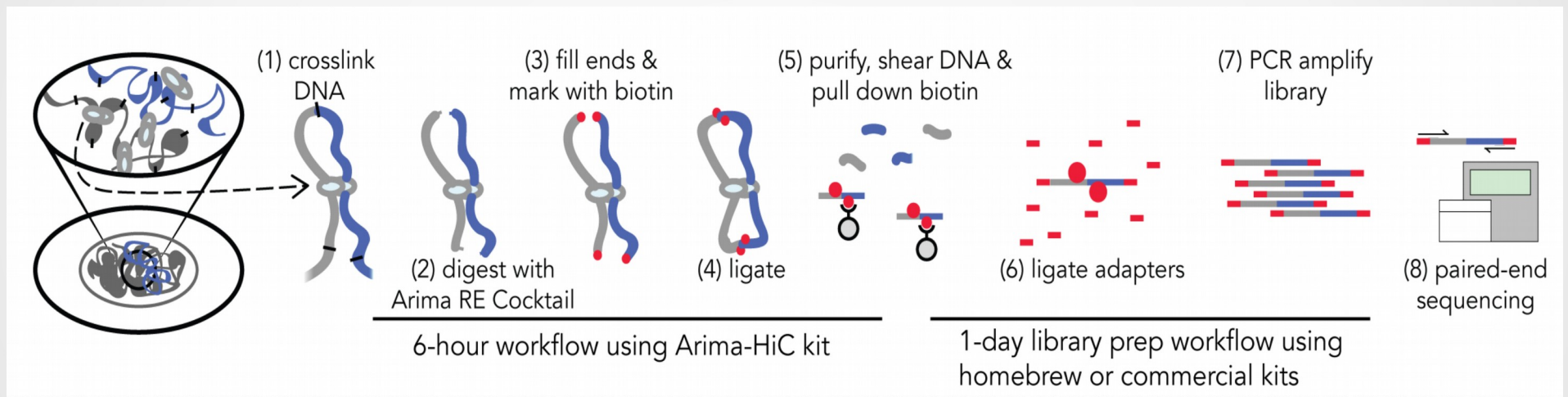
- 10X Genomics

Transposase-based

- Tell-Seq
- stLFR

Zheng et al, 2016

# Special types of libraries: HiC



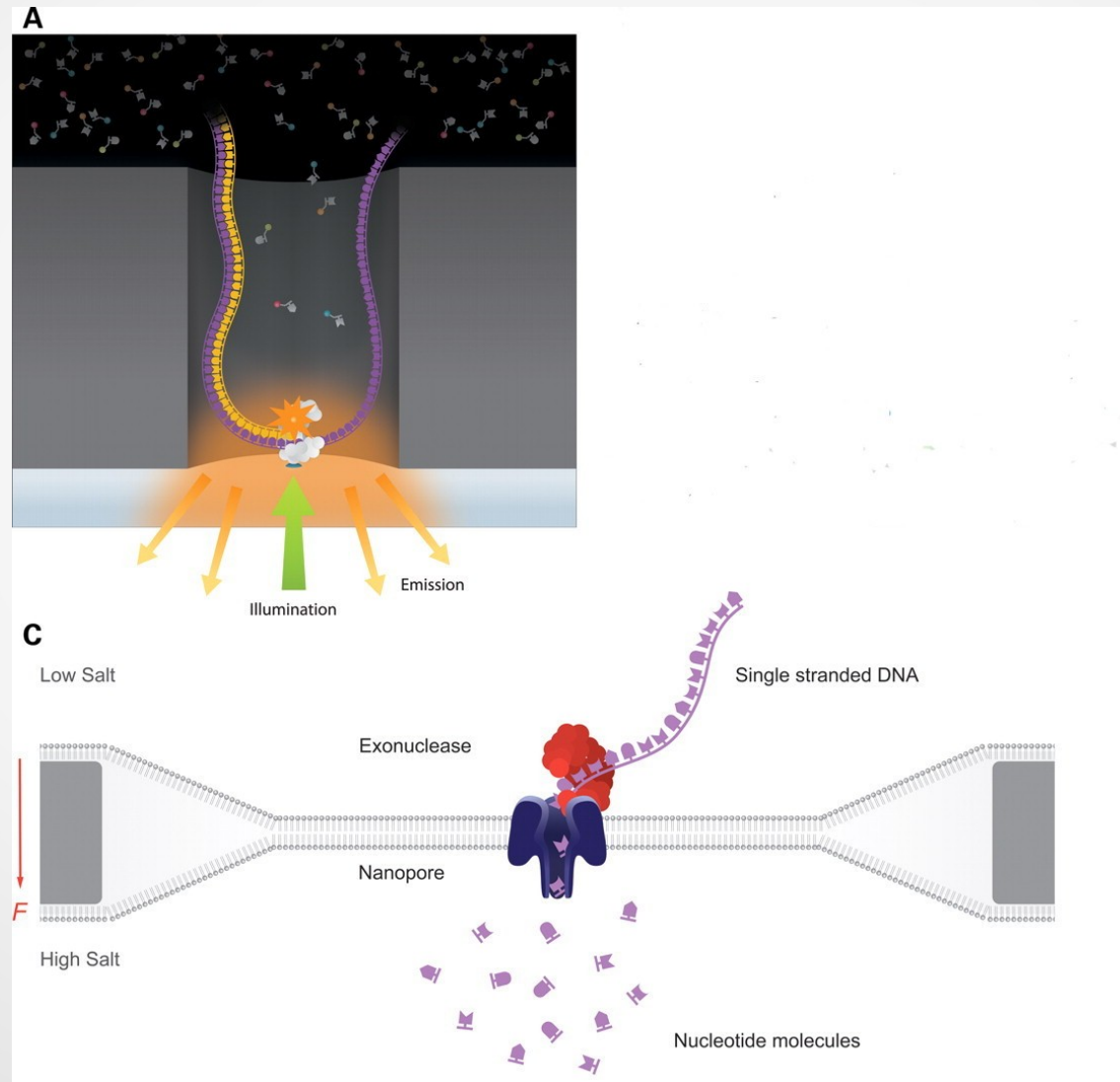
HiC-sequencing allows to scaffold even very fragmented draft genomes to the chromosome level.

# Sequencing technologies: third generation

## PacBio vs Nanopore

### PacBio

- Shorter reads
- Less errors
- More expensive



### Nanopore

- Longer reads
- More(?) errors
- Cheaper

Waiting  
forthcoming  
release of  
new chemistry  
soon!

# Assembly strategies and assemblers

## Short read based Sequencing

- Linked reads + HiC
- Overlapping PE reads + HiC

## Assemblers

Supernova + 3D-DNA/Salsa  
w2rap + 3D-DNA/Salsa

## Long read based Sequencing

- PacBio HiFi + HiC
- Nanopore + Illumina + HiC
- Pacbio + Nanopore + HiC

## Assemblers

HiFiasm

MaSuRCA/Flye/Canu + Medaka + 3D-DNA/Salsa  
Falcon/Flye/Canu + Arrow + 3D-DNA/Salsa

For highly repetitive genomes an optical mapping (Bionano) could be added as intermediate step before HiC-scaffolding

# Assembly of mtDNA from WGS

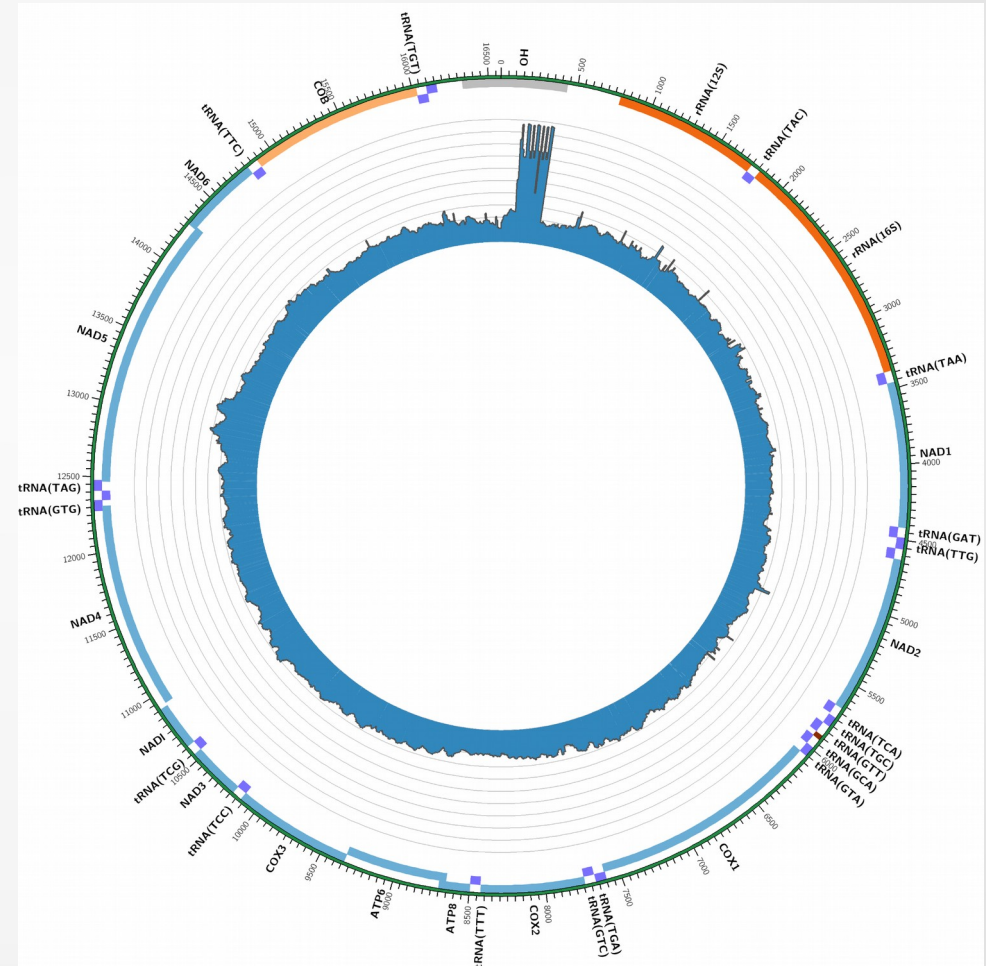
## Main issue

dramatic (~100-1000x) difference in coverage between nuclear and mtDNA genome.

## Solution

Independent assembly from downsampled reads (~1-5x coverage of nuclear genome) using different tools.

Tool	Function
MitoZ	assembly
Mitos	annotation



Tomarovsky, 2021



# Quality control of the assembly

## Assembly quality metrics

- Number of C-scaffolds should be equal to number of chromosome pairs
- Number of breaks introduced by HiC-scaffolding
- N50 and L50 (not informative for chromosome length assemblies)
- Number of Ns in the assembly
- Number of unplaced scaffolds
- BUSCO metrics

and so on ...

**N50** - maximal length of contig/scaffold in the assembly for which all sequences of such length and longer encompass no less than 50% of the assembly.

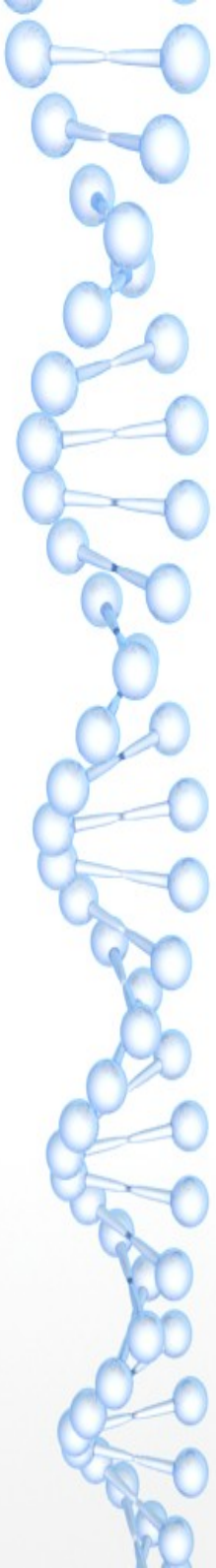
# BUSCO

## Benchmarking Universal Single-Copy Orthologs

- Assesses number of conservative genes present in the assembly
- specific databases for different taxa

Example for 4 mustelid species (Mammalian database, 9226 BUSCOs)

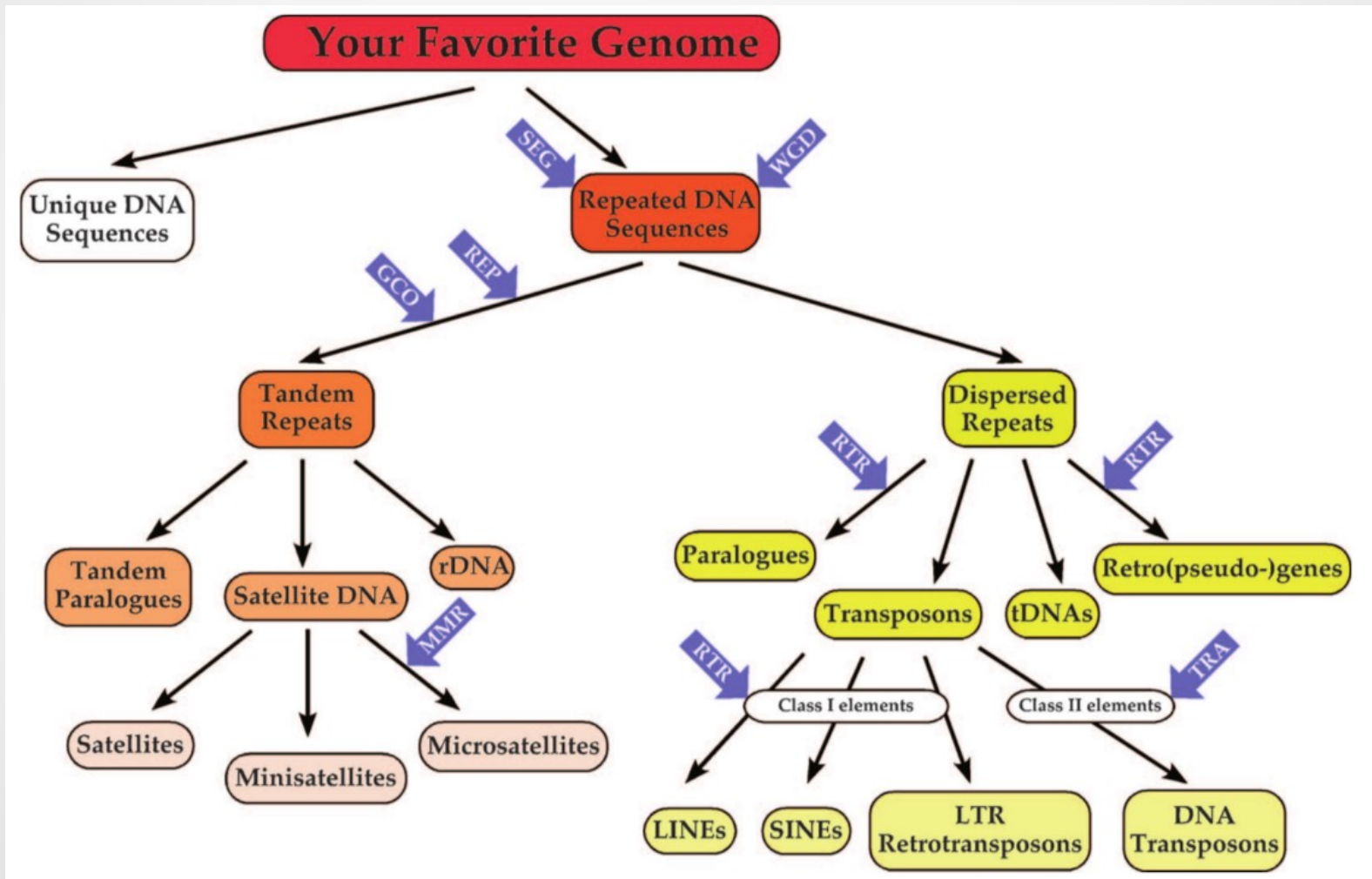
Species	Complete, %	Complete and single copy, %	Complete and duplicated, %	Fragmented, %	Missing, %
<i>Mustela nigripes</i>	96.2	94.0	2.2	1.0	2.8
<i>Mustela putorius furo</i>	94.0	92.8	1.2	1.2	4.8
<i>Enhydra lutris</i>	96.4	95.5	0.9	0.8	2.8
<i>Pteronura brasiliensis</i>	95.2	94.1	1.1	1.4	3.4



## III. Genome Projects Analysis



# Repeat types



Richard et al, 2008

# Major types of mobile elements

## Class I: retrotransposons

Reverse Transcriptase (RT) copies genome into DNA

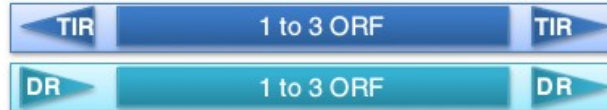
### Long Terminal Repeat (LTR) elements

(*copia*, *gypsy*, etc)



Double stranded DNA is integrated using a transposase-related integrase

### DIRS elements (*DIRS*, *Ngaro*, *viper*)



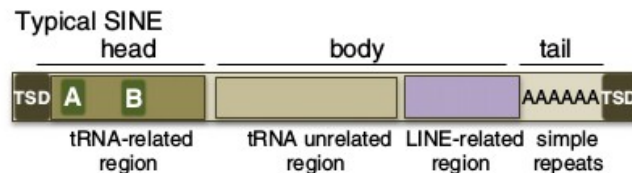
DNA is integrated using a tyrosine recombinase

### Non-LTR elements (*LINE*, *Penelope*, etc)



Integration and priming of RT reaction is mediated by an endonuclease

### Non-autonomous retroelements (*SINE*, *SVA*, etc)



## Class II: DNA transposons

DNA genome itself serves as the template for transposition

### Terminal Inverted Repeat (TIR) elements

(*Tc1-Mariner*, *PiggyBac*, *P*, etc)



DNA is cleaved from donor site and integrated at target site using transposase

### Helitron elements



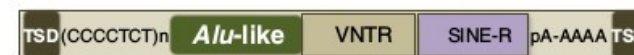
Single strand DNA intermediate is replicated from donor site by a rolling-circle mechanism

### Maverick elements

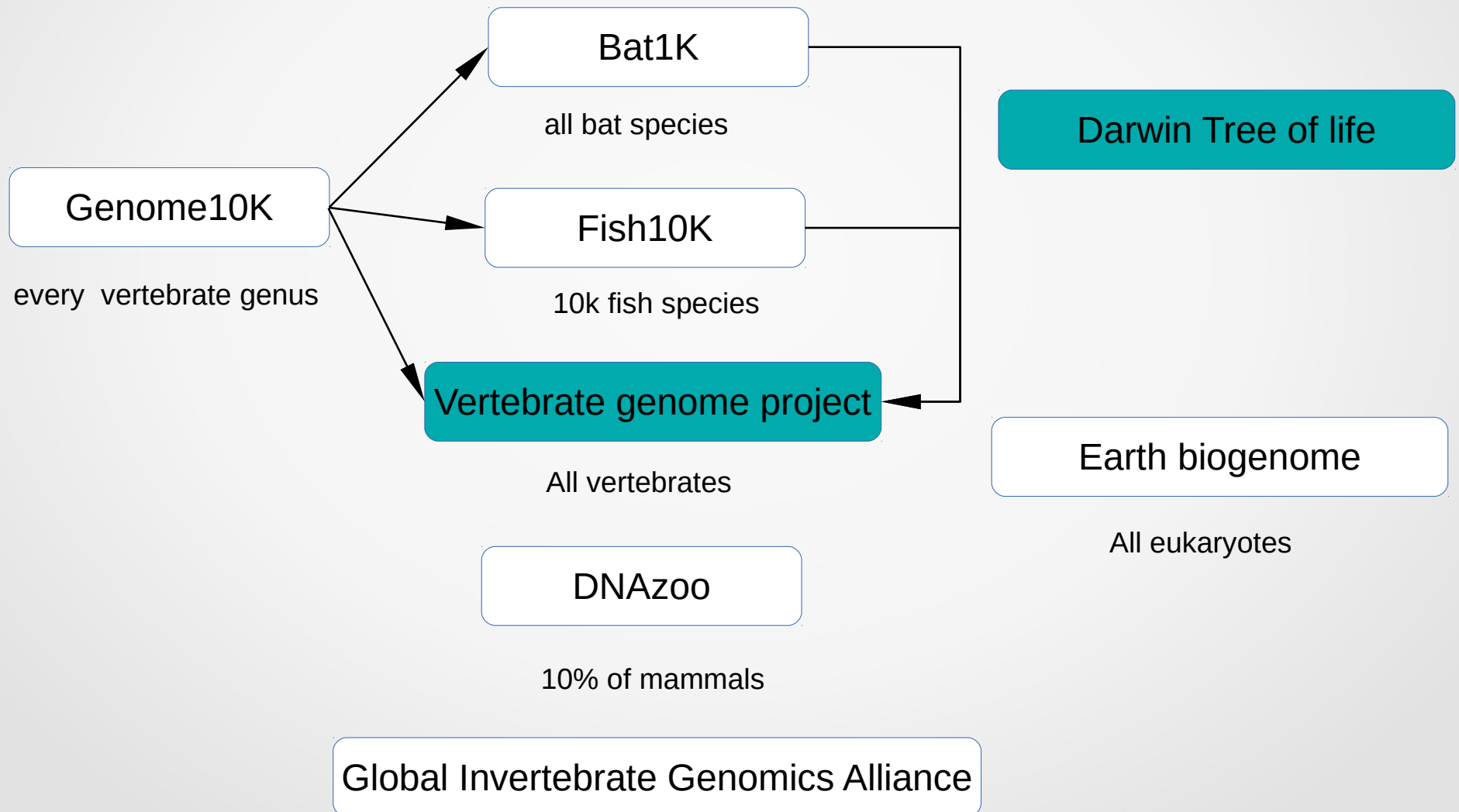


Both viral-like DNA polymerase and retroviral-like integrase are thought to be involved in double stranded DNA integration

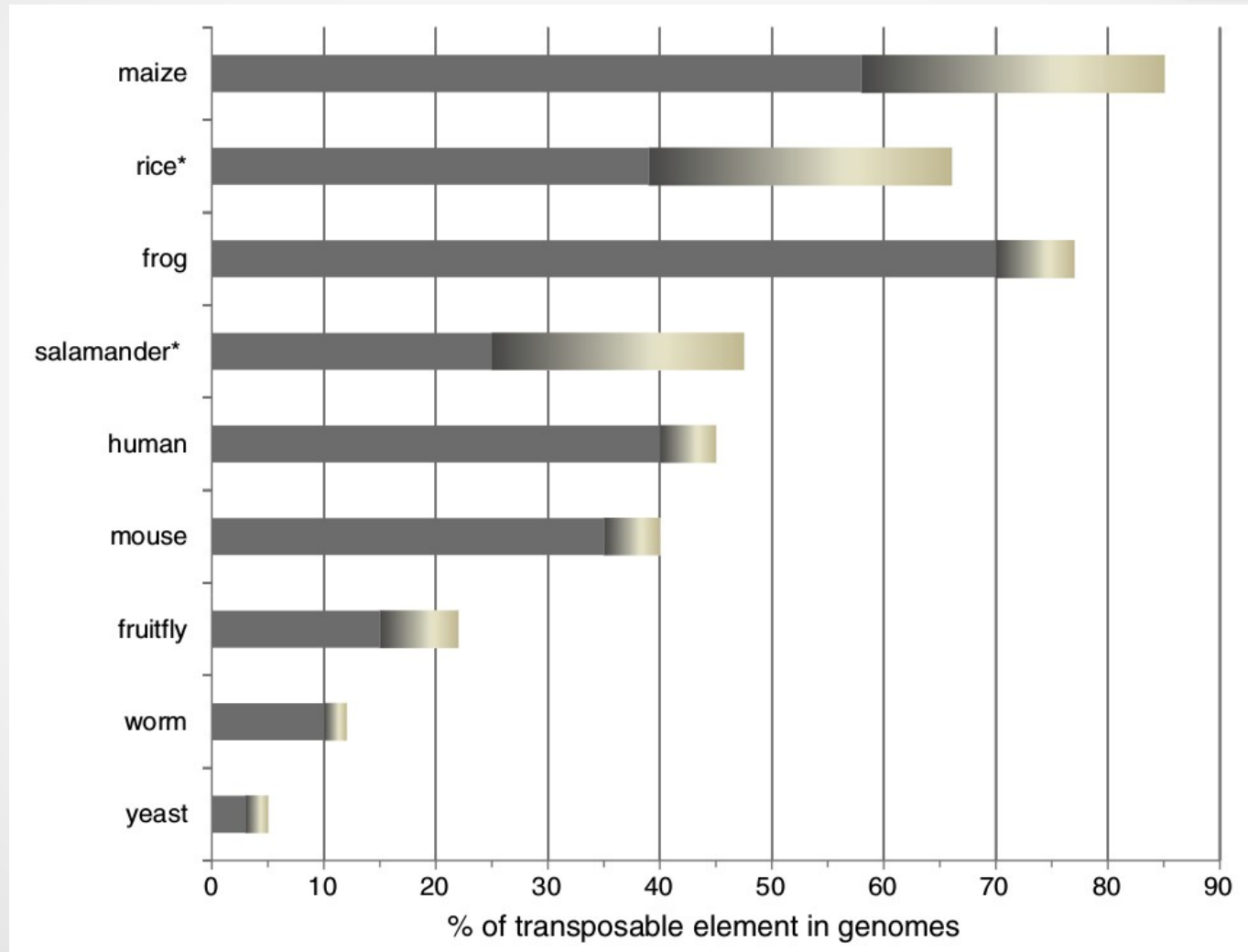
### SVA



# Big multigenome projects and their aims



# Presence of mobile elements in eukaryotic genomes



Chénais et al, 2012

# Major types of mobile elements

## Class I: retrotransposons

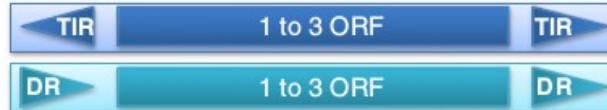
Reverse Transcriptase (RT) copies genome into DNA

### Long Terminal Repeat (LTR) elements (*copia*, *gypsy*, etc)



Double stranded DNA is integrated using a transposase-related integrase

### DIRS elements (*DIRS*, *Ngaro*, *viper*)



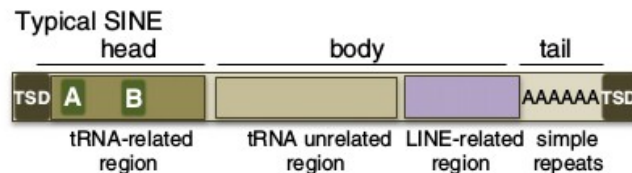
DNA is integrated using a tyrosine recombinase

### Non-LTR elements (*LINE*, *Penelope*, etc)



Integration and priming of RT reaction is mediated by an endonuclease

### Non-autonomous retroelements (*SINE*, *SVA*, etc)



## Class II: DNA transposons

DNA genome itself serves as the template for transposition

### Terminal Inverted Repeat (TIR) elements (*Tc1-Mariner*, *PiggyBac*, *P*, etc)



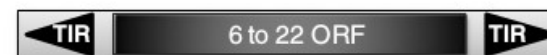
DNA is cleaved from donor site and integrated at target site using transposase

### Helitron elements



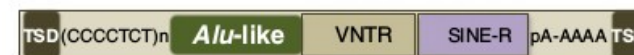
Single strand DNA intermediate is replicated from donor site by a rolling-circle mechanism

### Maverick elements

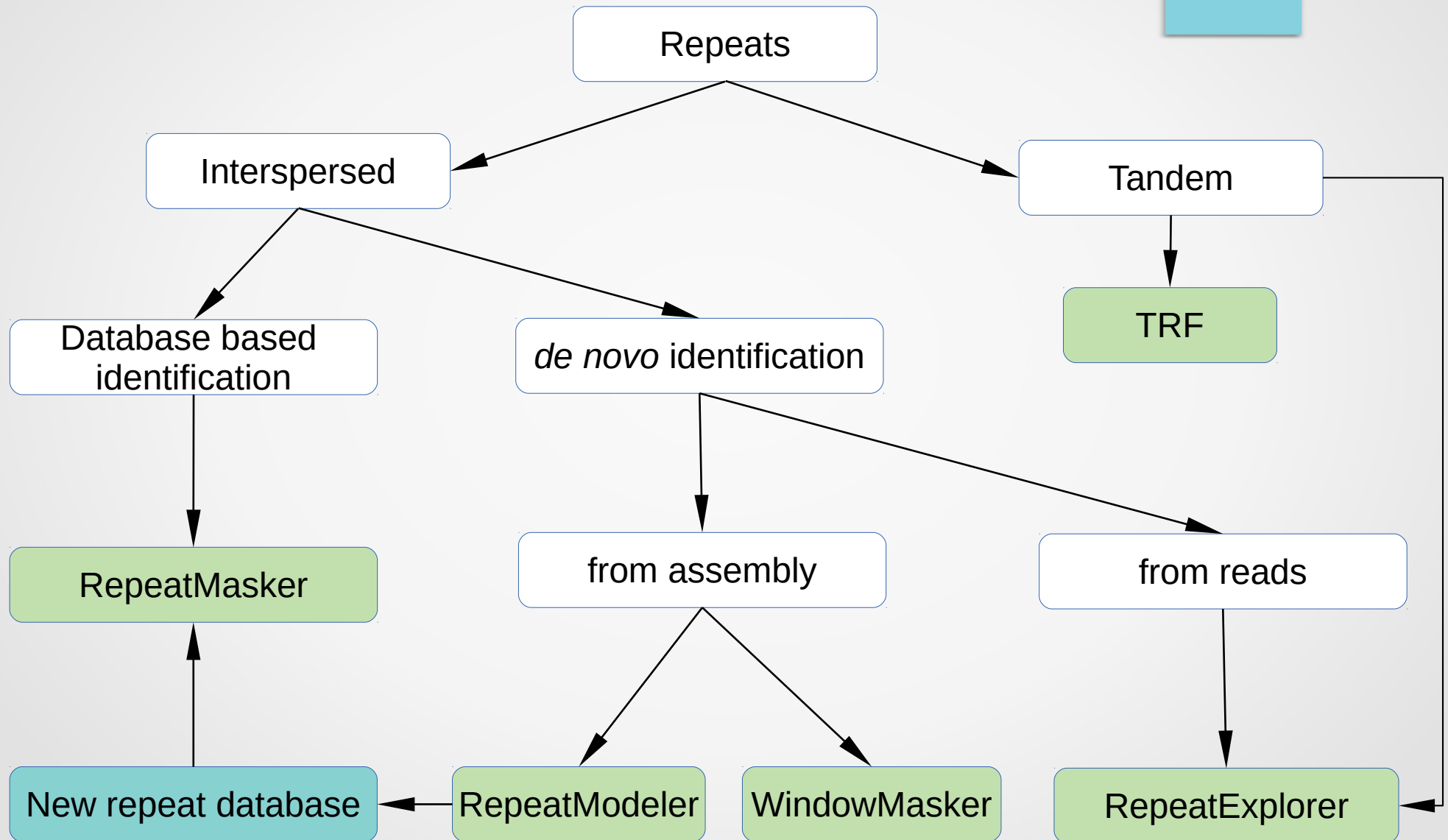


Both viral-like DNA polymerase and retroviral-like integrase are thought to be involved in double stranded DNA integration

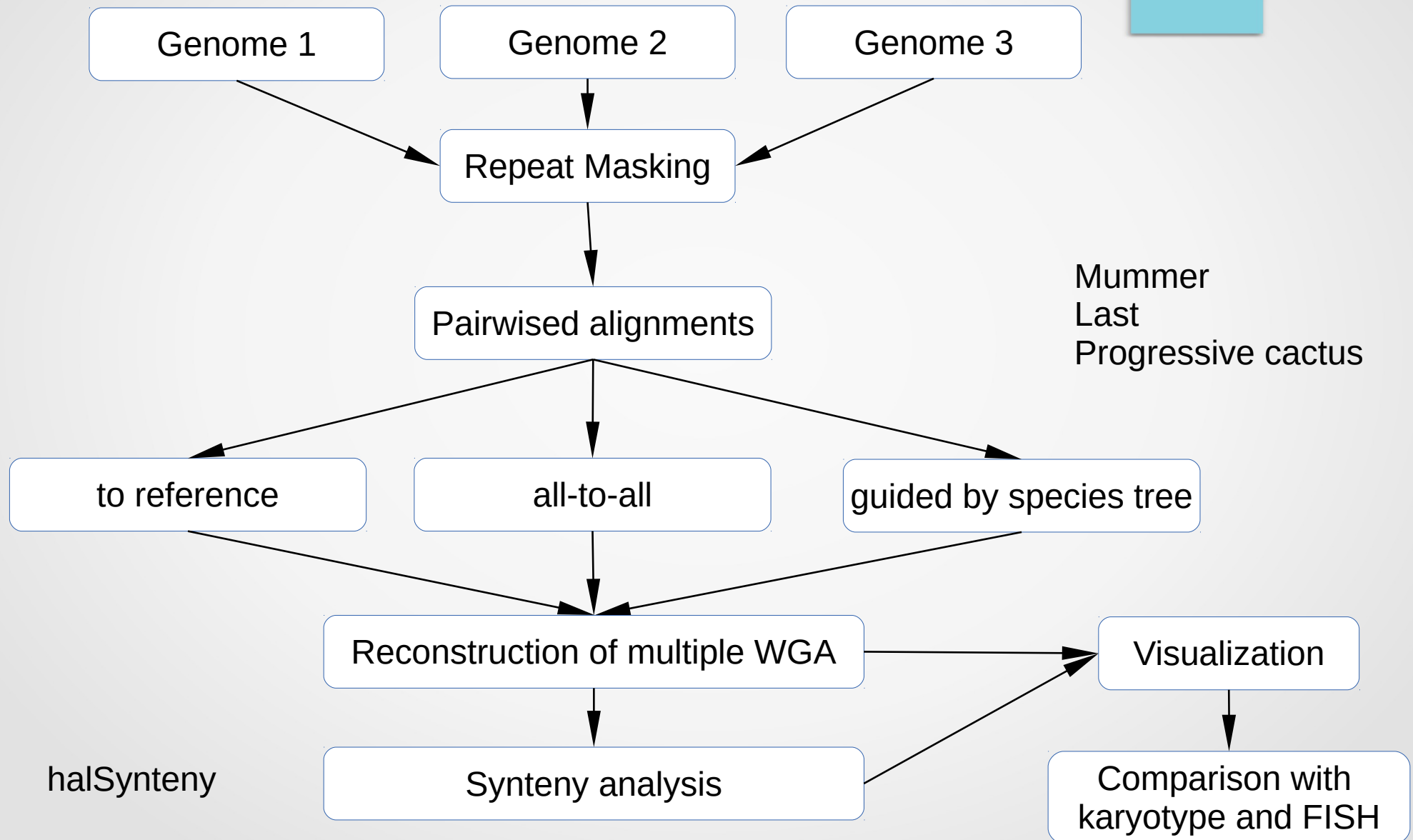
### SVA



# Tools for repeat identification

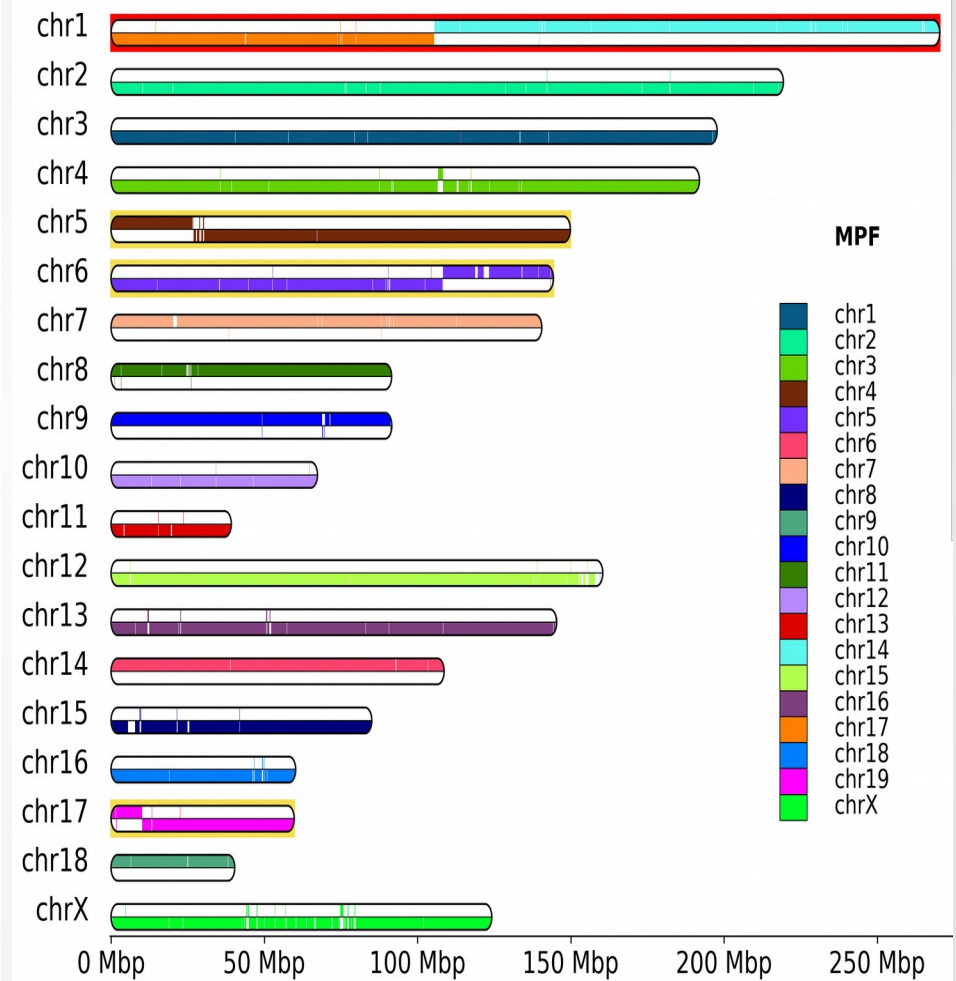
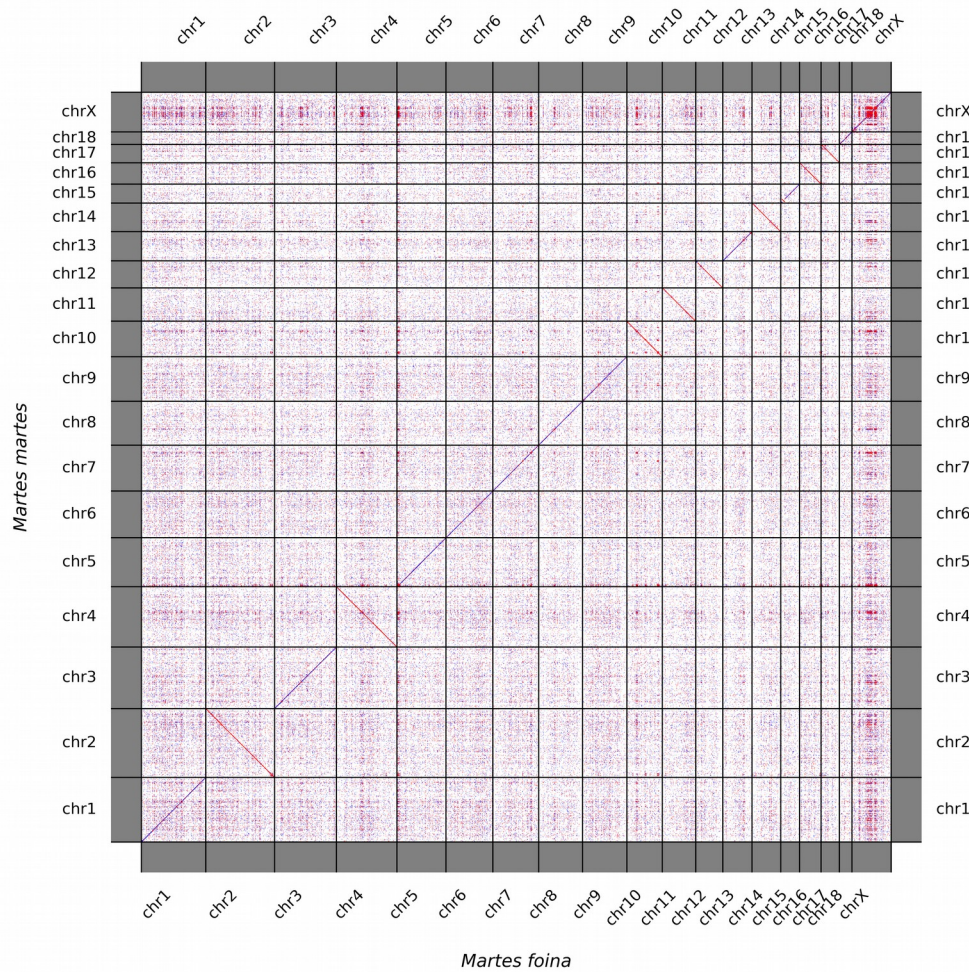


# Whole genome alignment (WGA)





# Whole genome alignment and synteny blocks





# Protein-coding gene annotation

C  
o  
m  
b  
i  
n  
a  
t  
i  
o  
n

## I. RNAseq-based

## II. homology-based

CESAR2.0

Comparative annotation toolkit

Exonerate

## III. de novo predictions

Augustus

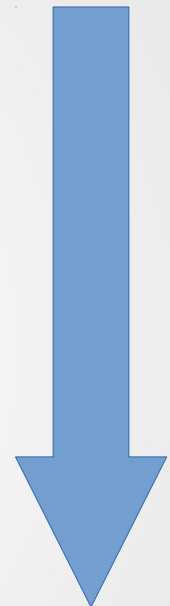
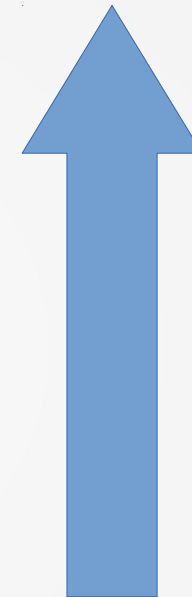
Metaeuk

## IV. Hybrid approach

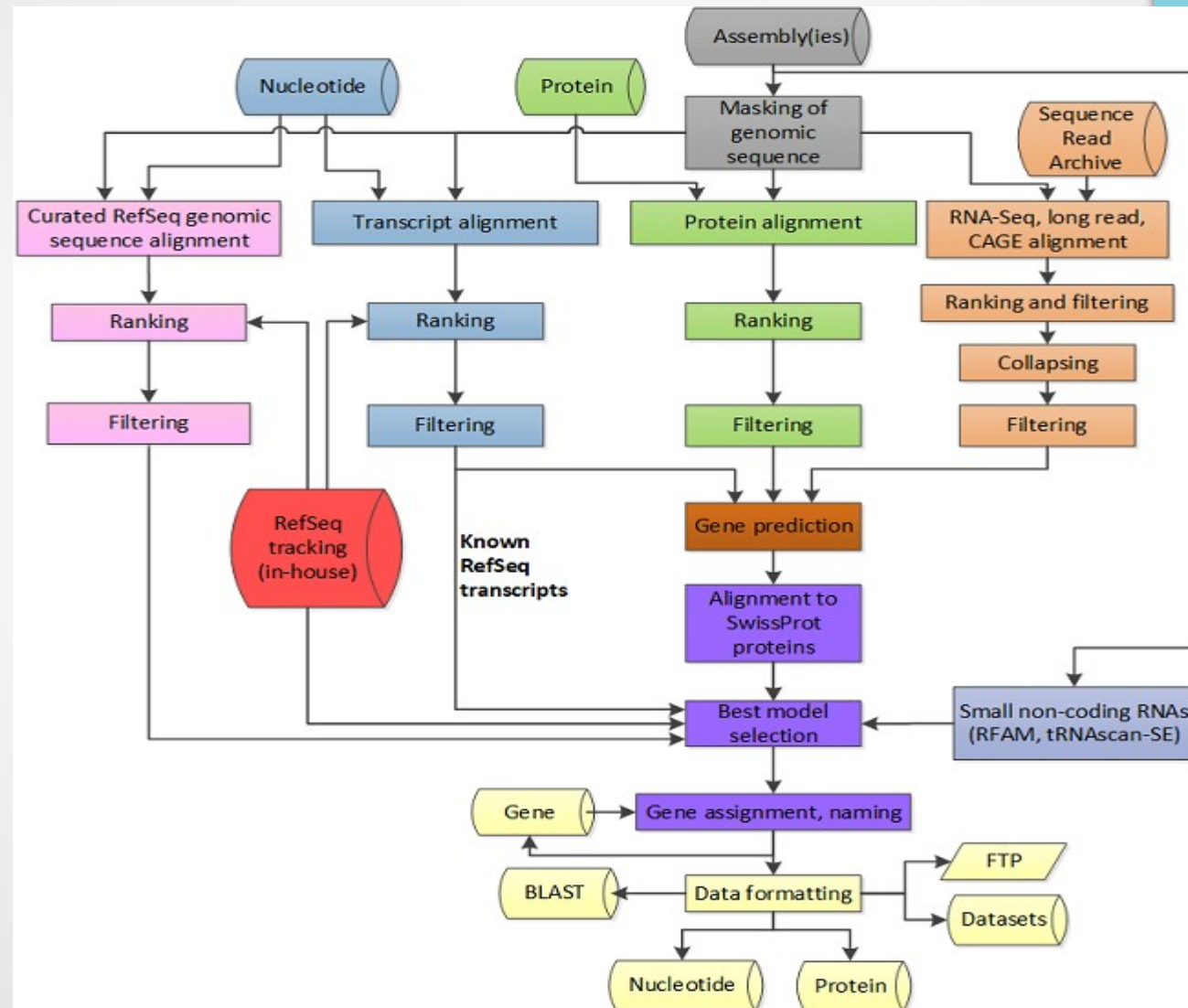
MAKER2

Quality

Price

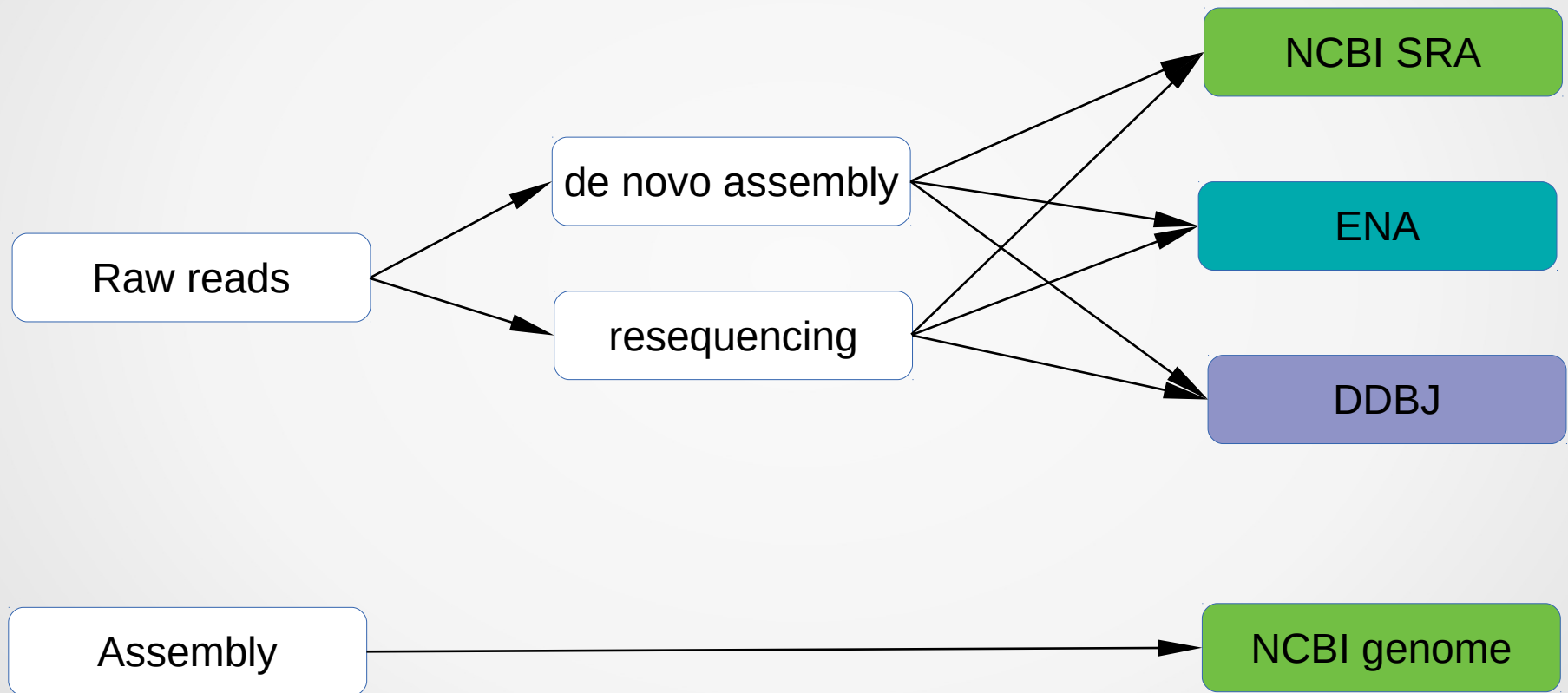


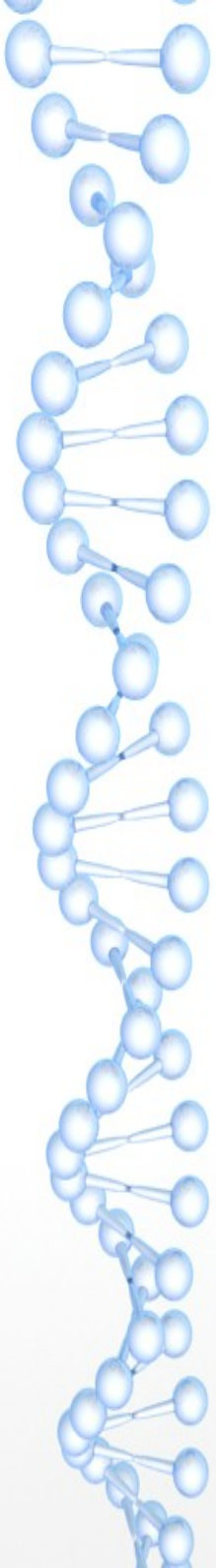
# The NCBI Eukaryotic Genome Annotation Pipeline



[https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/process/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/process/)

# Data sharing





End of module III