

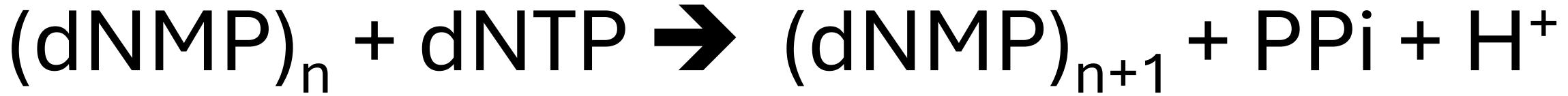


wellcome
connecting
science

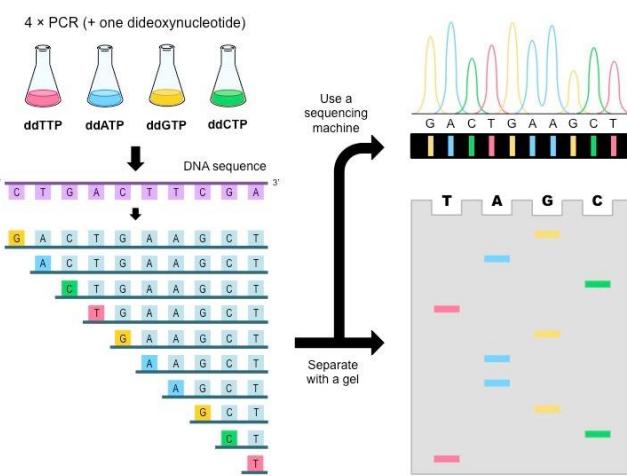
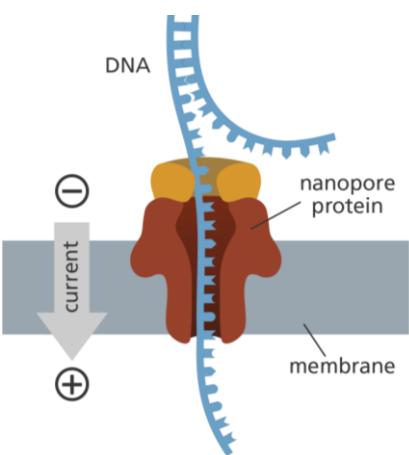
Mitochondrial Assembly. Genome skimming strategy.

Mg. Agustín Baricalla
Universidad Nacional de San Luis.

May 2025

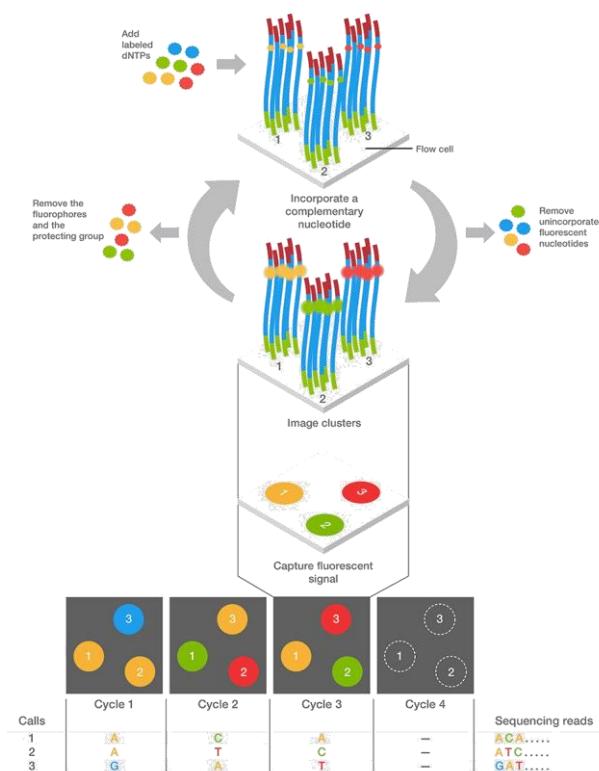


DNA chain
(Nanopore)



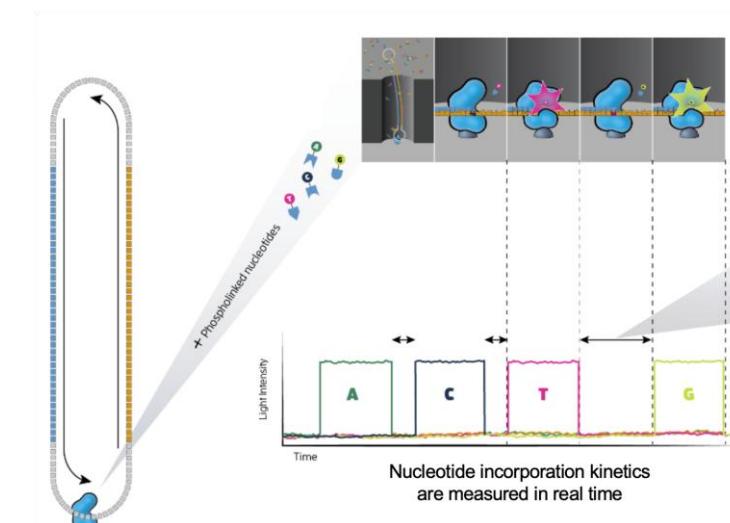
Nucleotide
(modified dye)

Illumina,
Sanger,
Pacbio



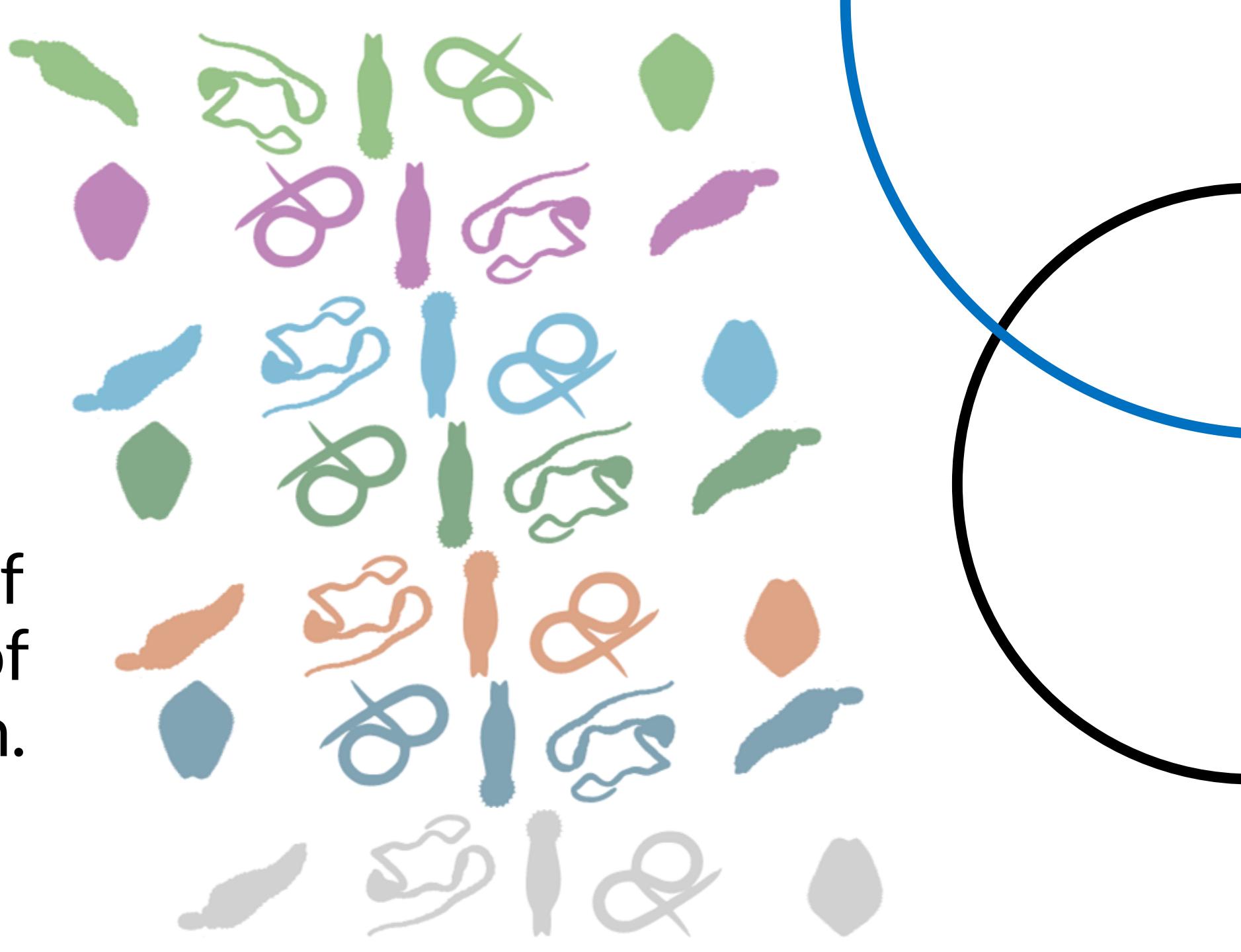
Pyrosequencing

Ion Torrent

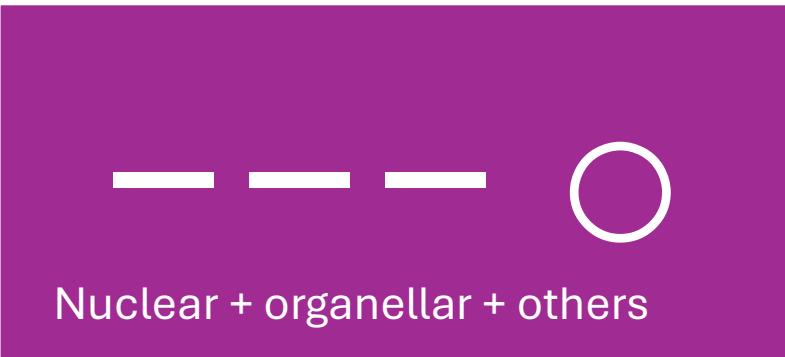
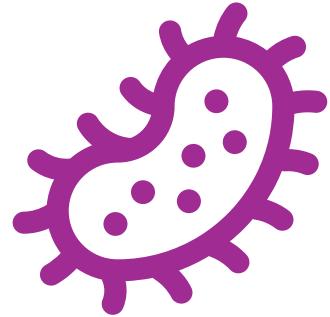


Step I.
The sample.

Previous
knowledge of
the biology of
the organism.

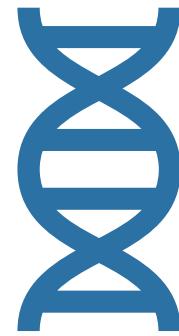


Preliminary considerations - Sample

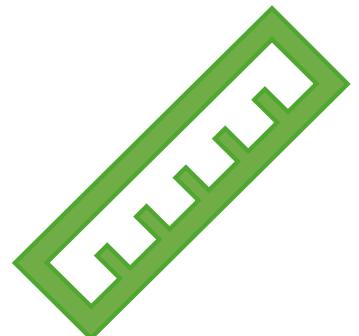


Extraction
Inhibitors (!)

Ploidy?
XY | XO | ZW (?)



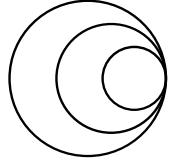
QC



Length



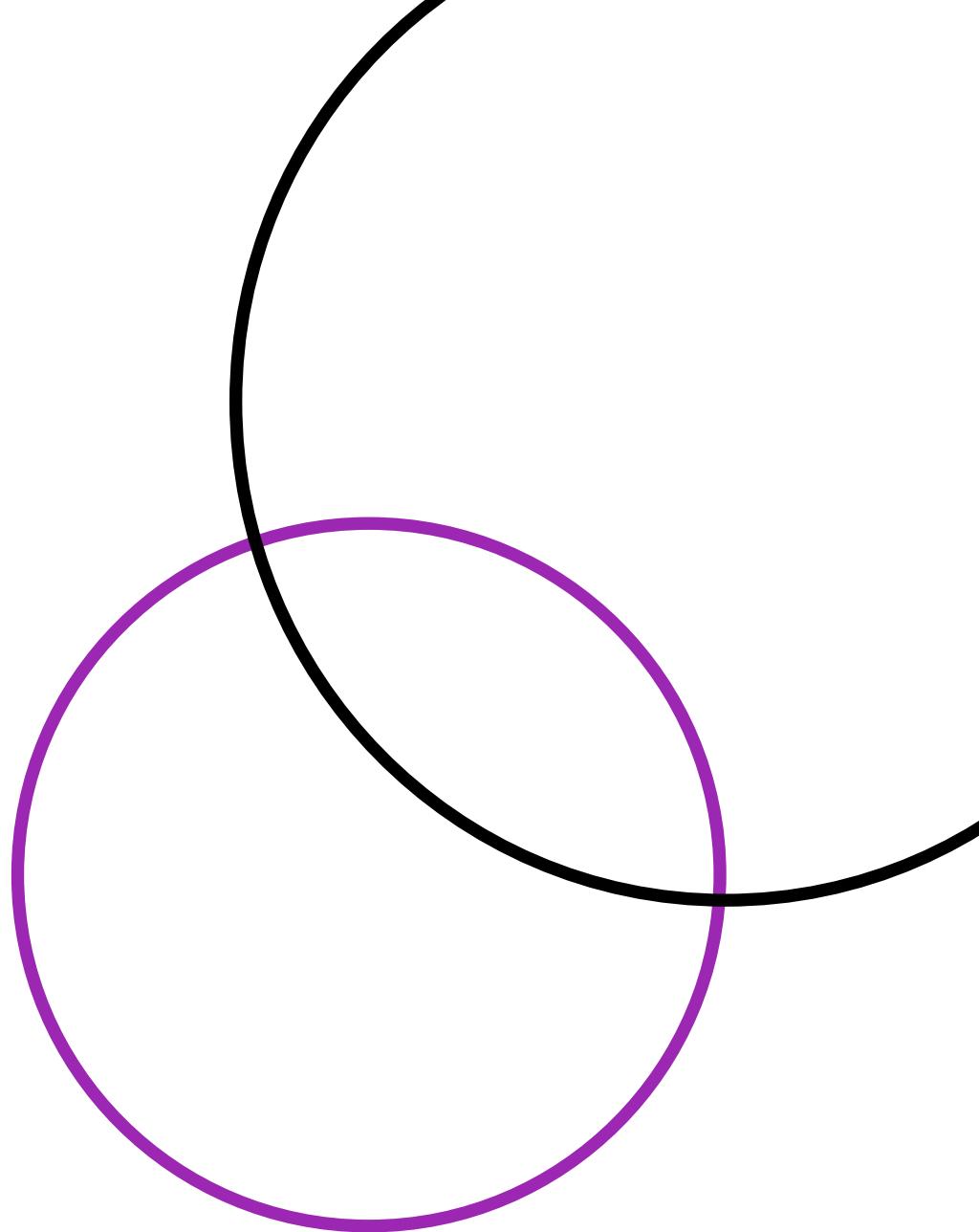
Sequencing



wellcome
connecting
science

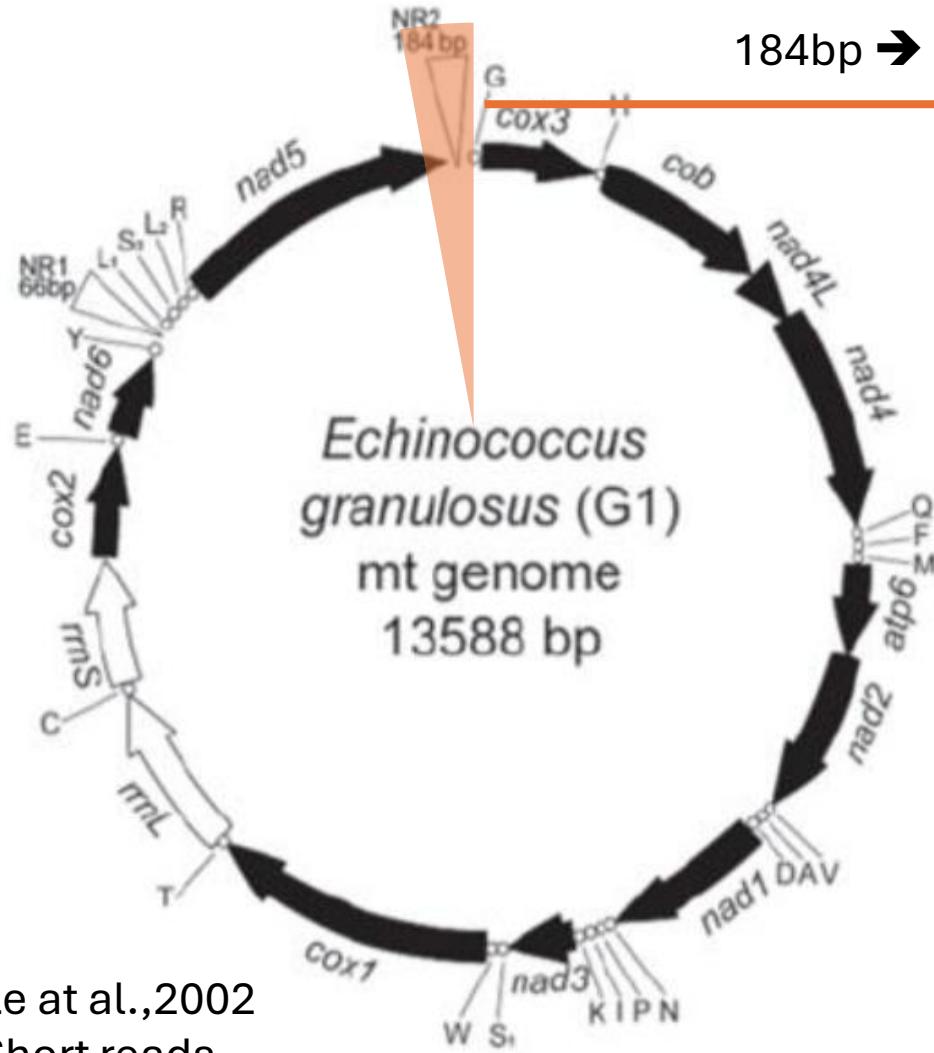
Mitochondrial genome assembly

May 2025

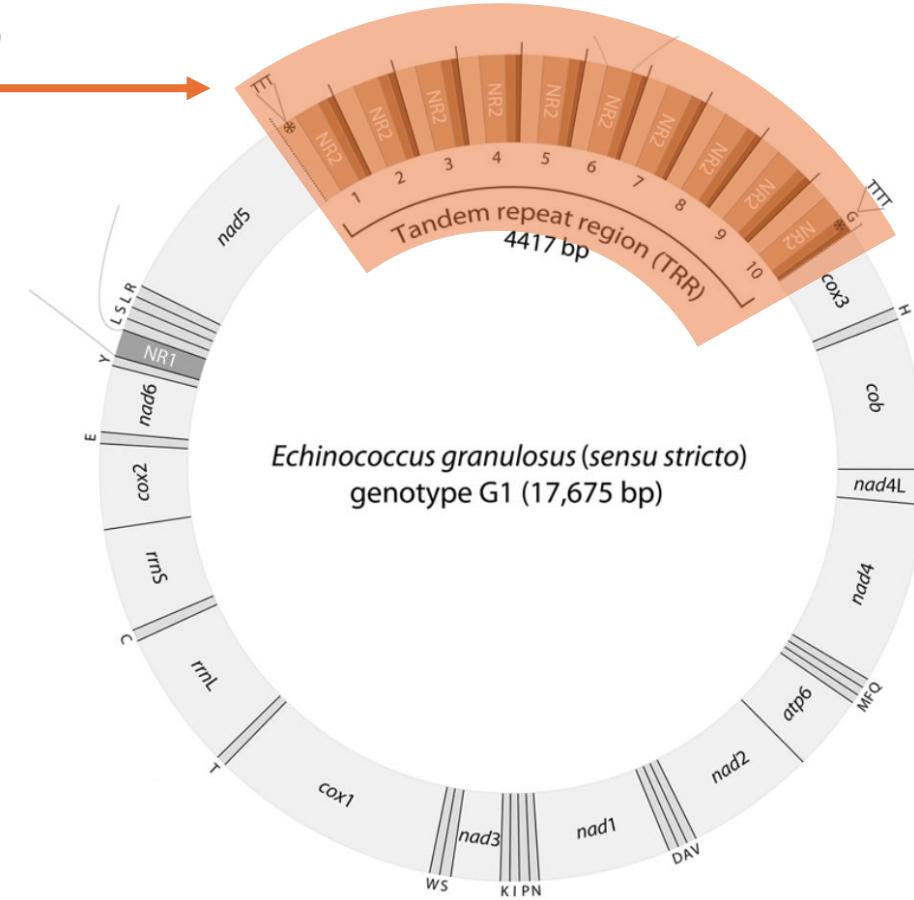


Why we choose long reads for mitochondrial assembly?

mtDNA *Echinococcus granulosus*:

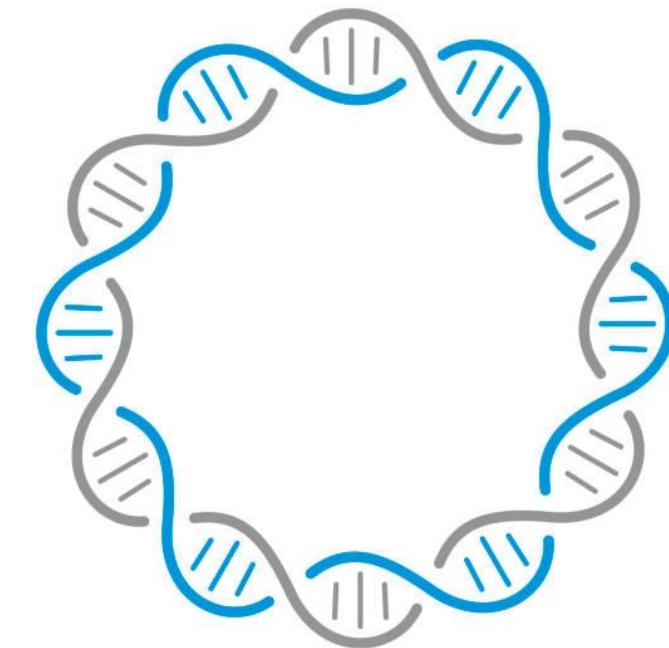
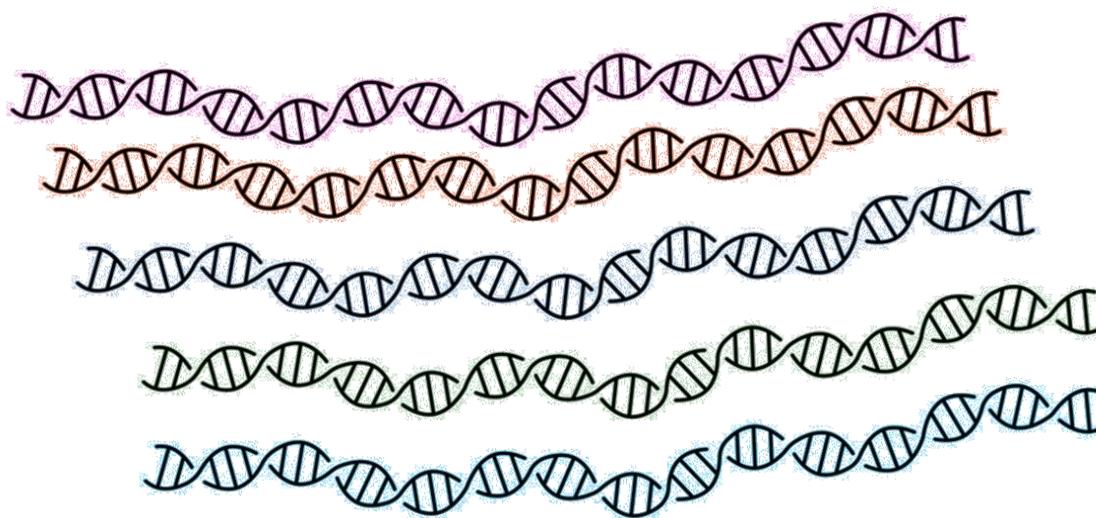


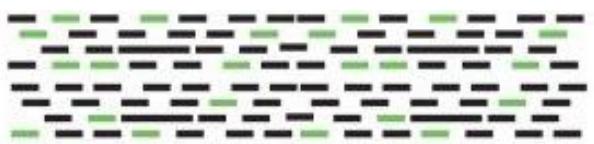
Le et al., 2002
Short reads



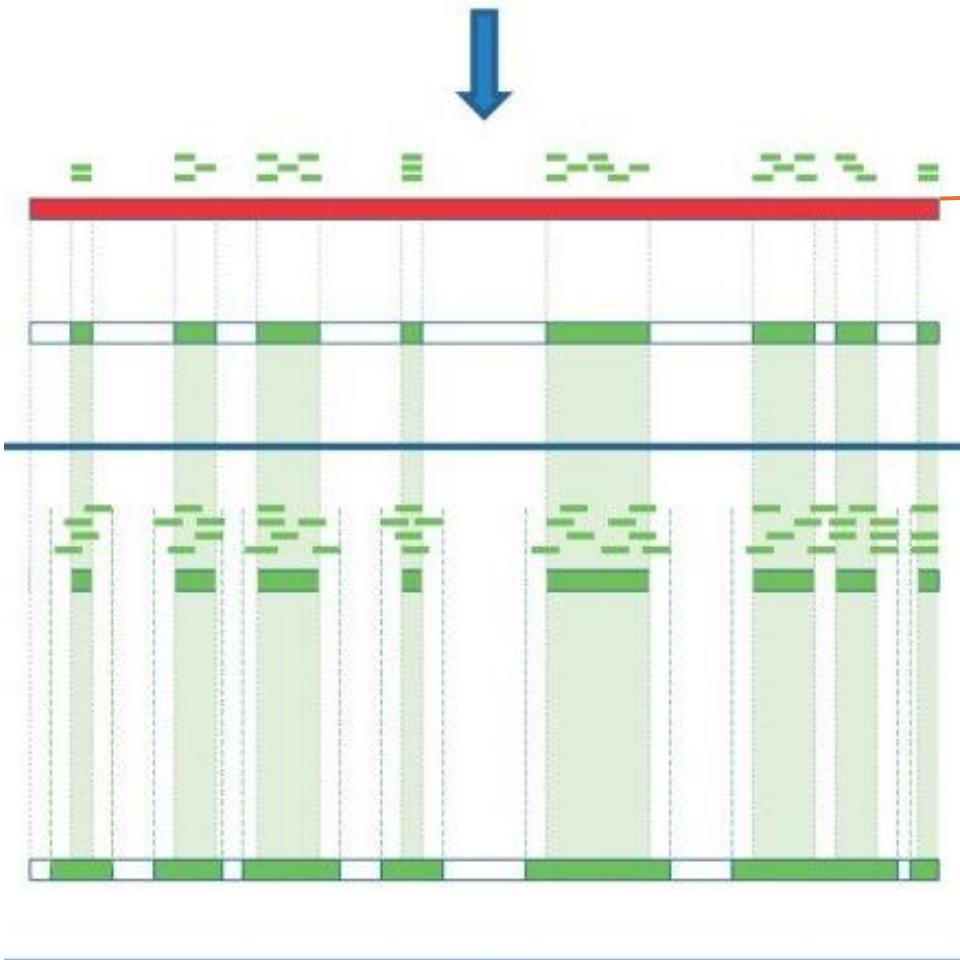
Kinkar et al., 2019
Long -reads

Mitochondrial genome assembly strategies





Reads:
Nuclear
Mitochondrial



**Mapping
Alignment
Extension**

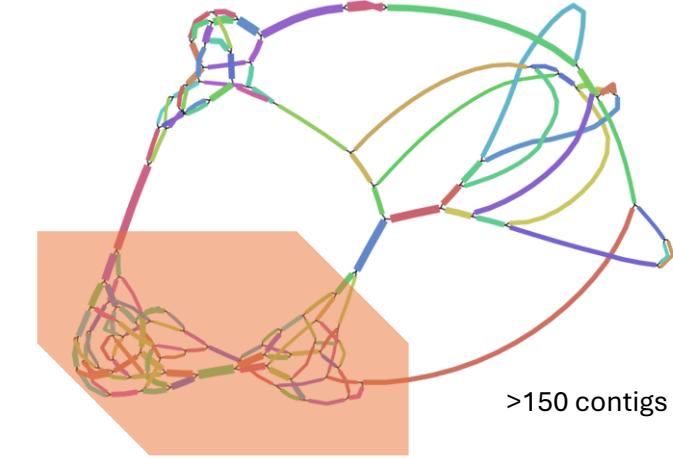


**Re-Mapping
Re-Alignment
Extension**

Consensus

Closely-related
reference
Genes or
mitogenomes

MitoBim

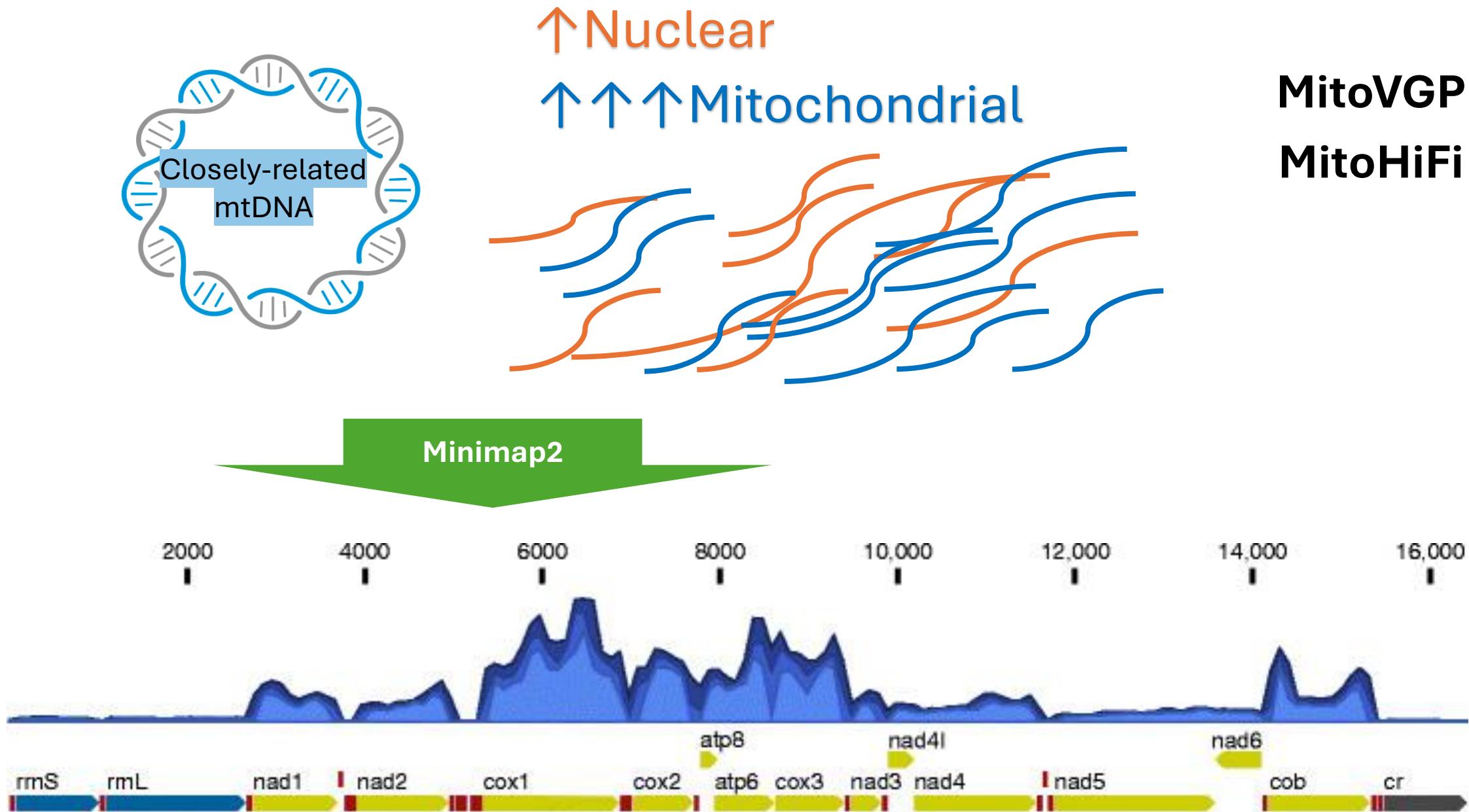


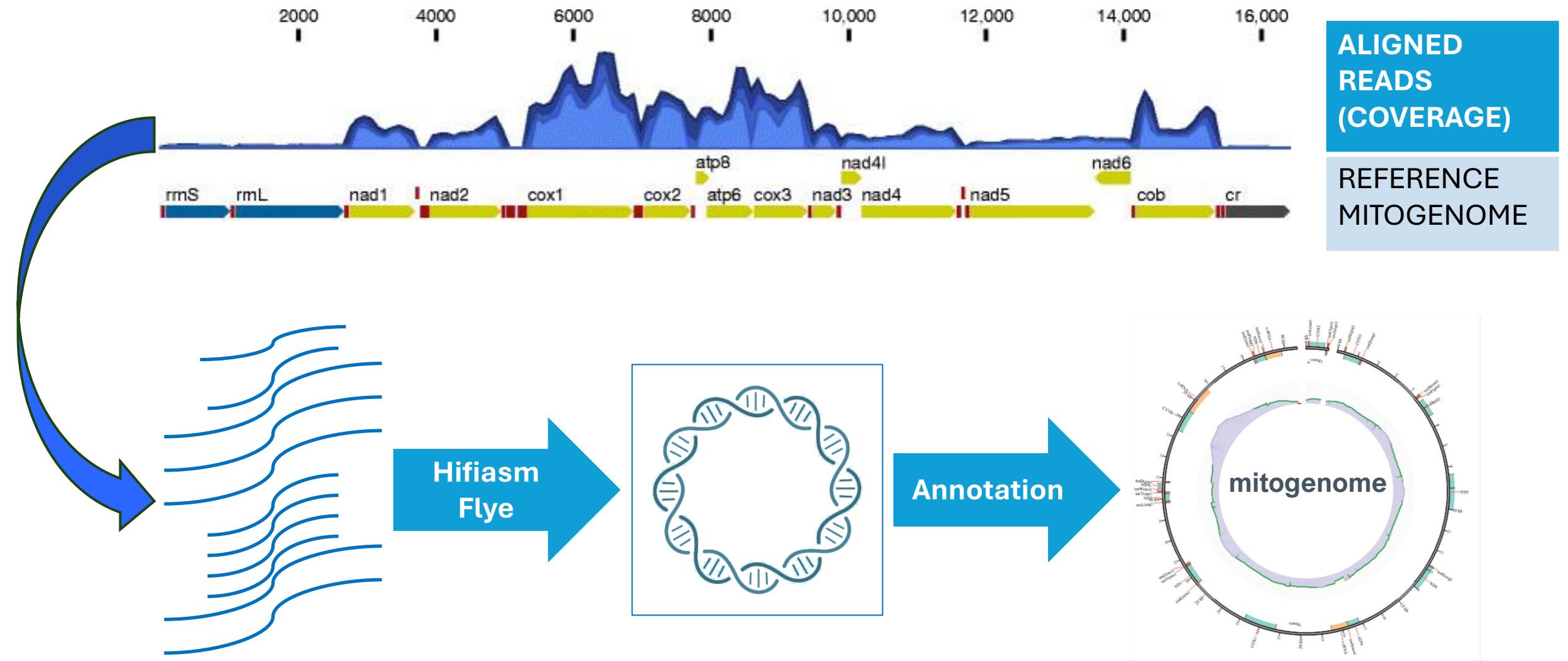
Orange: complex repetitive region

Bandage visualization of mtDNA

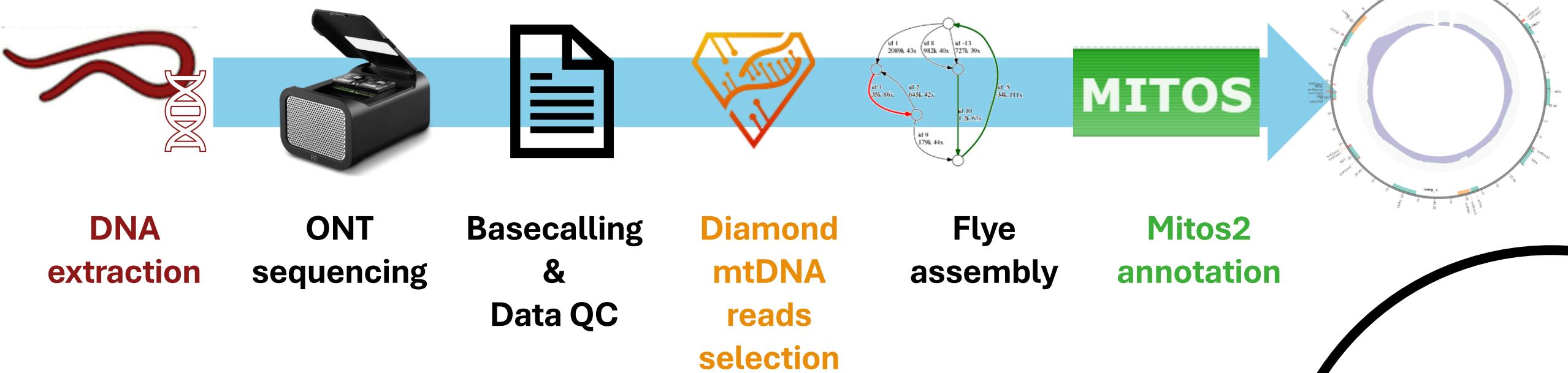


Highly repetitive sequences
or tandem repeats



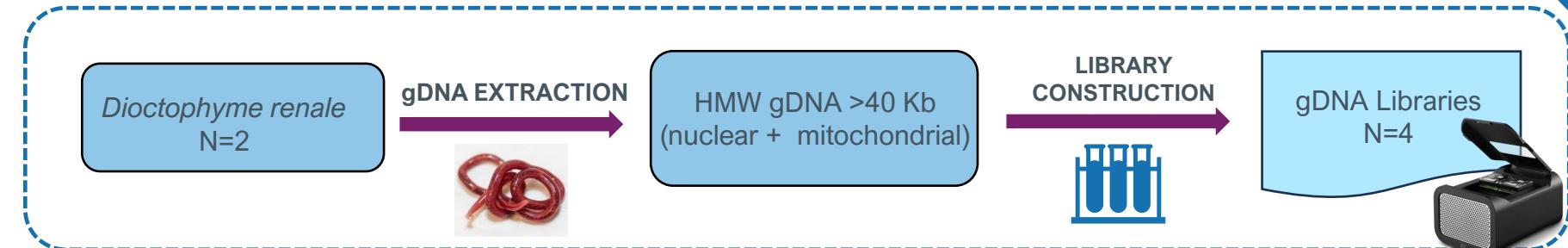


Pipeline

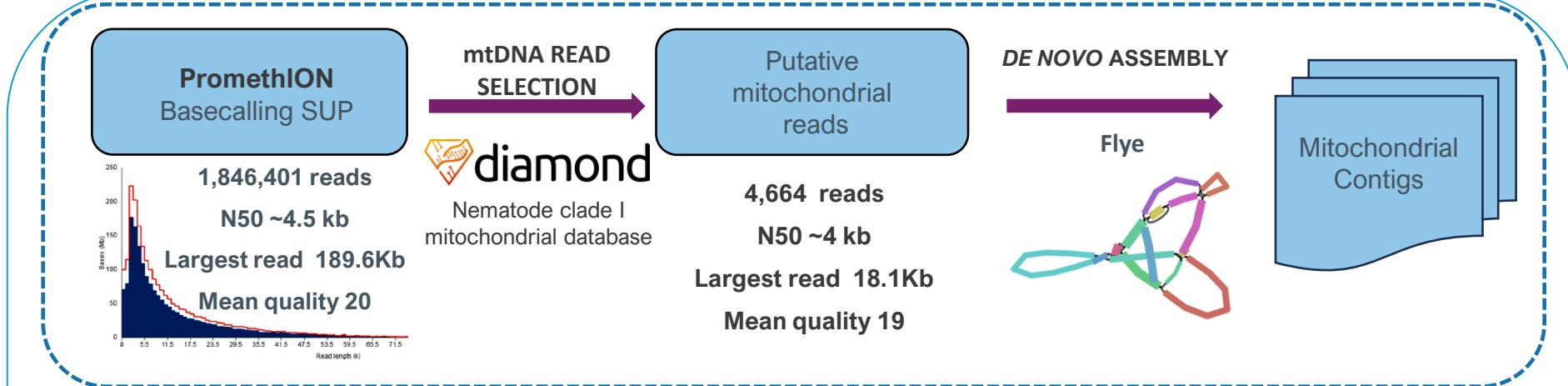


Workflow

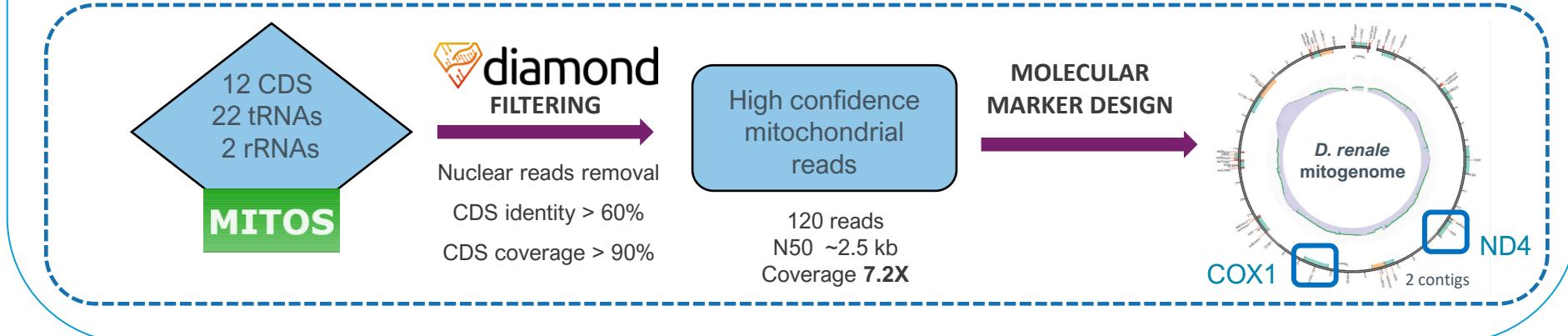
Sampling Sequencing



QC Assembly

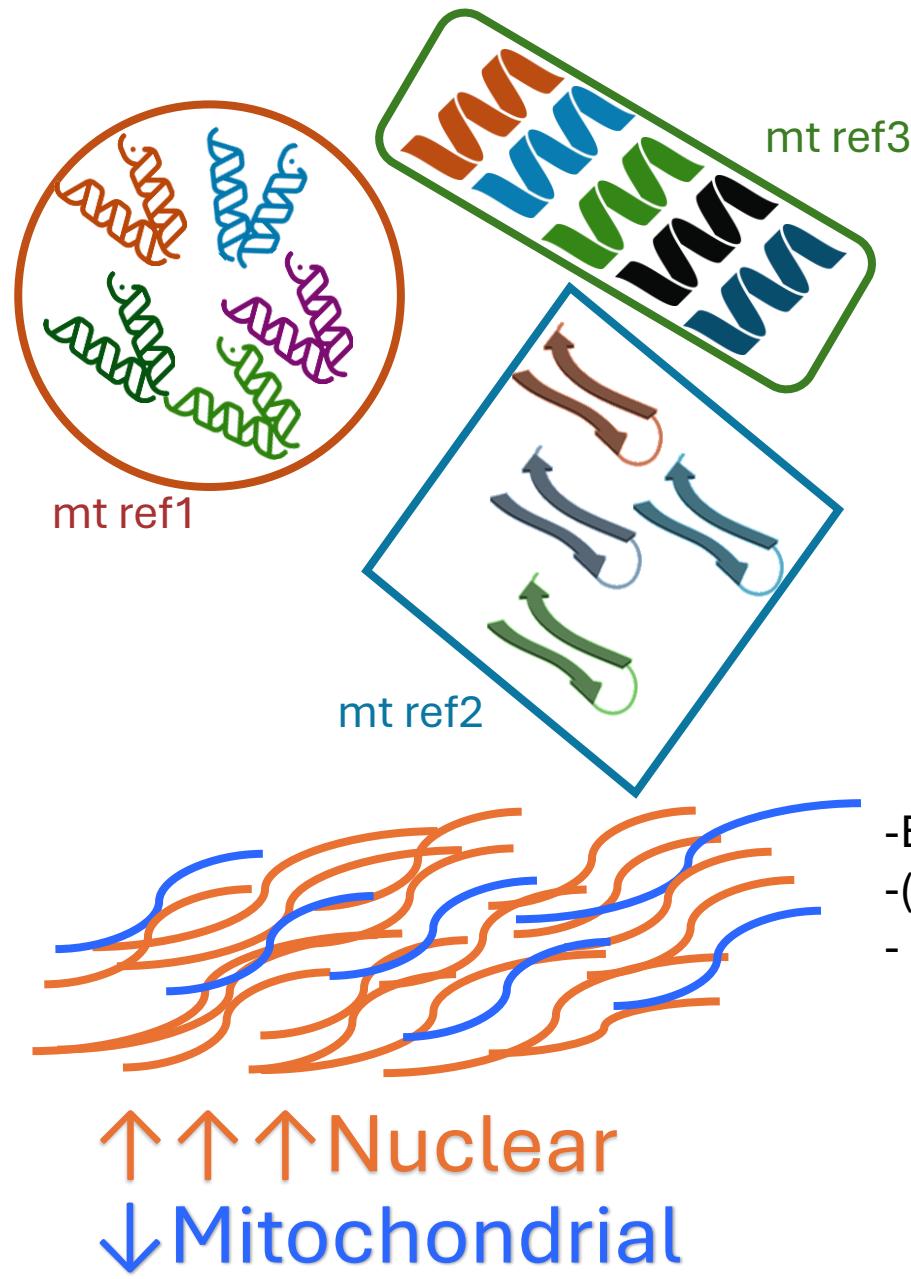


Annotation Molecular markers design

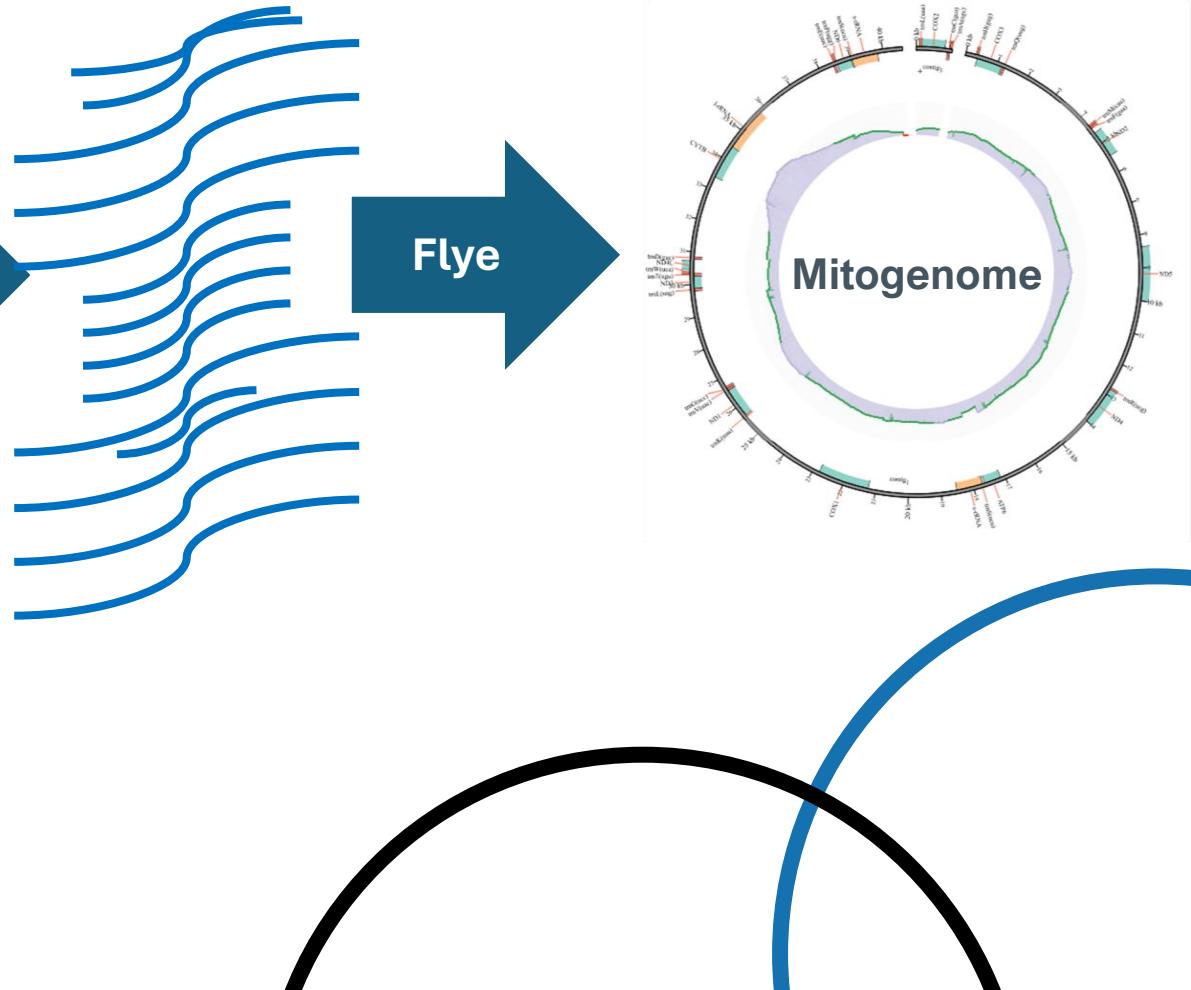


BIOINFO

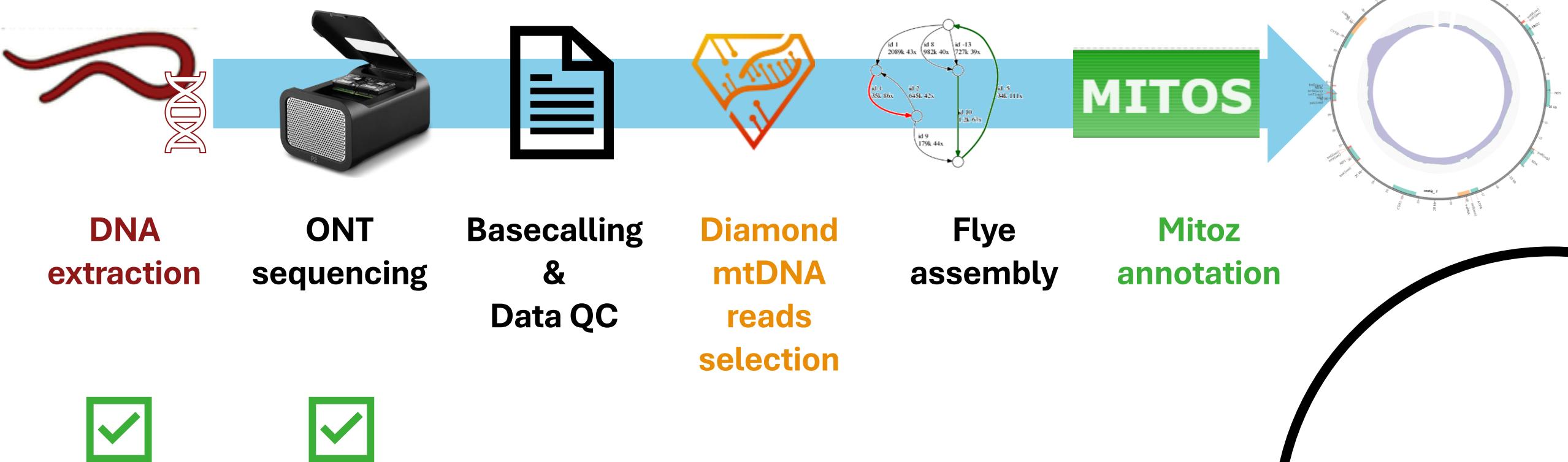
Mitochondrial genome *de novo* Lab10 - pipeline



- Based on proteins
- (Non) Closely related
- Reads selection



Pipeline



Readings files format

FASTQ format

- Simple format of unaligned sequence reads
- Sequences and quality score associated with each base.
- File extension : .fastq o .fq
 - (Compressed fastq.gz | fq.gz)

Readings files format

FASTQ format

@IL16_4408:3:5:17860:13258
CTGGCTCACATACAGGCCAGTATAAAGCGTCTCCTTTAAA
+
HHHE

@IL16_4408:3:14:13276:1210
GTAGAAAAACAATATAGAAGCCTTAGGACAGAAAACCCTAA
+
HHHHHHHHHHHHHHHHHDHHHHHHHBGHHHHHCDCCF>

'+' sign for parity →

Read ID = '@' sign + sequence identifier →

Read sequence →

Read Quality by base (ASCII code) →

Nanoplot - QC

NanoPlot

-t 8

--fastq minion.fq

--dpi 300

--N50

-o ./nanoplot

--huge

Program (BD creation)

Threads

Fastq file name

Image quality

To show the N50 value mark

Output folder name

Big data file

Nanoplot - QC

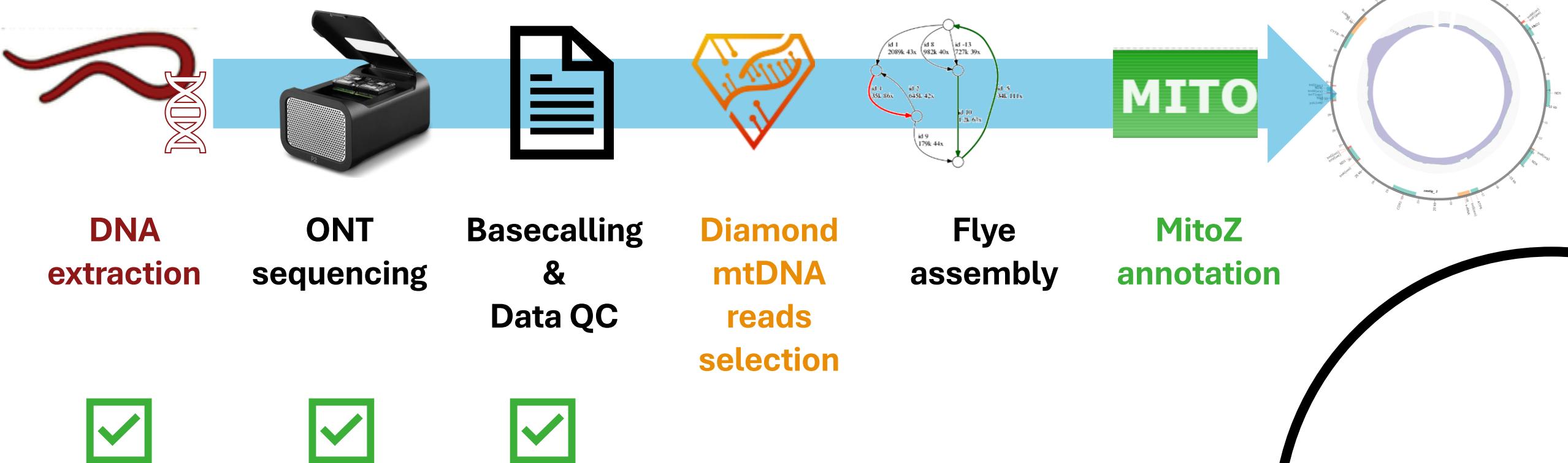
Summary statistics

Metrics	
number_of_reads	26190
number_of_bases	686824347.0
median_read_length	12433.0
mean_read_length	26224.7
read_length_stdev	35320.1
n50	52343.0
mean_qual	18.8
median_qual	21.7
longest_read_(with_Q):1	560959 (13.3)
longest_read_(with_Q):2	540464 (14.5)
longest_read_(with_Q):3	531365 (15.8)
longest_read_(with_Q):4	421500 (18.3)
longest_read_(with_Q):5	373488 (23.4)
highest_Q_read_(with_length):1	32.5 (5063)
highest_Q_read_(with_length):2	31.9 (9176)
highest_Q_read_(with_length):3	30.9 (5644)
highest_Q_read_(with_length):4	29.7 (43859)
highest_Q_read_(with_length):5	29.5 (7944)
Reads >Q10:	26190 (100.0%) 686.8Mb
Reads >Q15:	24320 (92.9%) 632.1Mb
Reads >Q20:	18101 (69.1%) 469.3Mb
Reads >Q25:	1461 (5.6%) 21.5Mb

Read lengths vs Average read quality plot using dots



Pipeline



Diamond BLAST

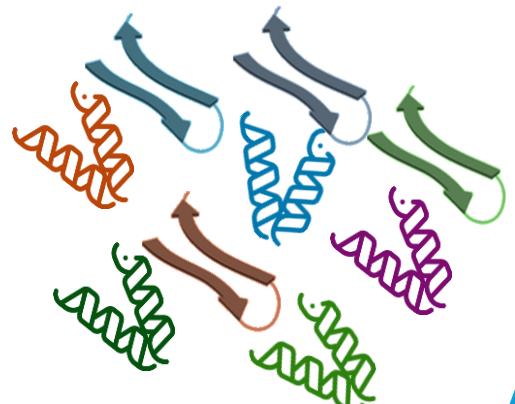


DIAMOND is a sequence aligner for protein and translated DNA searches vs Protein database.

Key features:

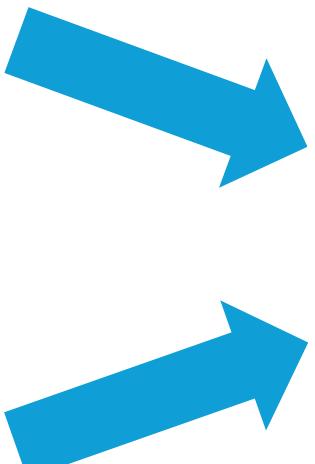
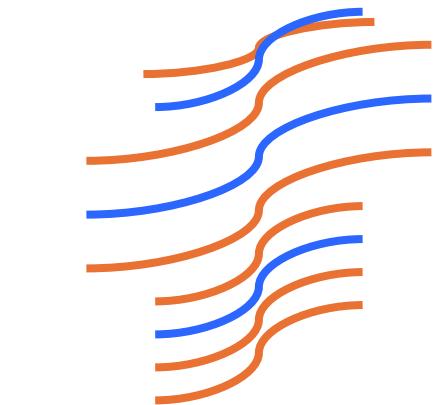
- Pairwise alignment of proteins or DNA.
- Faster than BLAST.
- Allows alignments with reading frame changes for analysis of long reads (error rate tolerance).
- Various **output formats**, including **tabular** and **fasta/q**.

Diamond BLAST

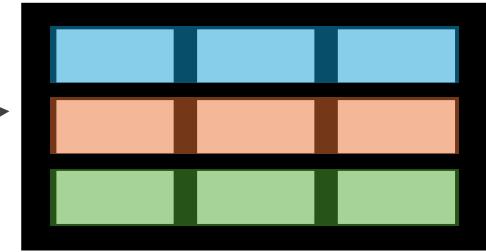


PROTEINS
FASTA

NUCLEIC ACIDS
FASTA/Q



PROTEIC
DATABASE

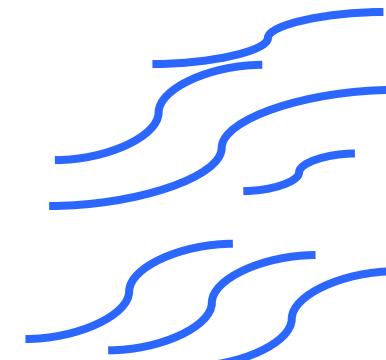


TABULAR FILE



FASTA/Q
MATCHING
SEQUENCES

NUCLEIC ACIDS
FASTA/Q





Diamond - outformat 6

qseqid	sseqid	qlen	slen	pident	length	mismat	gapopen	qstart	qend	sstart	send	evalue	bitscore	nident	positive	ppos
300262.1_14095	contig_1_ATP6_len=168_17134_17639_+	14095	168	99.4	168	1	0	12145	11642	1	168	3.24E-100	309	167	168	100
38576.1_8924	contig_1_COX1_len=520_21238_22799_+	8924	520	100	422	0	0	68	1333	1	422	5.02E-283	838	422	422	100
599553.1_8749	contig_1_CYTB_len=372_33381_34500_+	8749	372	100	370	0	0	2671	3780	3	372	1.82E-246	732	370	370	100
599553.1_8749	contig_1_ATP6_len=168_17134_17639_+	8749	168	99.4	168	1	0	7962	8465	1	168	1.99E-100	309	167	168	100
455962.1_8005	contig_1_ND4_len=393_12828_14008_-	8005	393	100	259	0	0	4122	3346	135	393	5.77E-152	467	259	259	100
						...										
65512.1_7635	contig_1_ND1_len=299_25659_26558_+	7635	299	100	116	0	0	1955	1608	1	116	3.22E-64	211	116	116	100

Diamond - outformat 6



Qseqid	sseqid	qlen	slen	pident	length	mismatch	gapopen	qstart	qend	sstart	send	evalue	bitscore	nident	positive	ppos
65512. 1_7635	ND1	7635	299	100	116	0	0	1955	1608	1	116	3.22E -64	211	116	116	100
65512. 1_7635	COX1	7635	520	100	422	0	0	6378	5113	1	422	7.27E -286	838	422	422	100

Diamond tabular



qseqid	Query sequence ID.	qend	End of alignment in the query sequence.
sseqid	ID of the subject sequence (the one in the database).	sstart	Beginning of the alignment in the subject sequence.
qlen	Length of the query sequence.	send	End of alignment in the subject sequence.
slen	Length of the subject sequence.	evalue	E-value, which indicates the number of alignments expected by chance with a score equal or better. Lower values indicate higher significance.
pident	Percentage of identity between the query sequence and the subject sequence.	bitscore	Bit score of the alignment, a measure of the size of the alignment.
length	Length of alignment.	nident	Number of identities (identical base pairs) in the alignment.
mismatch	Number of misalignments (differences) in the alignment.	positive	Number of positive base pairs (similar amino acids) in the alignment.
gapopen	Number of gap openings in the alignment.	ppos	Percentage of positive positions in the alignment.
qstart	Beginning of the alignment in the query sequence.		

Parameter highlights



qseqid	ID of the query sequence.	eval	E-value, which indicates the number of alignments expected by chance with a score equal or better. Lower values indicate higher significance.
sseqid	ID of the subject sequence (the one in the database).	bitscore	Alignment bit score, a measure of alignment size.

Database preparation

```
diamond makedb  
--threads 4  
--db cladoi.dmnd  
--in proteins_mt.fa
```

Program (BD creation)

Threads to be used

Database name

Protein file for indexing

Diamond Search

```
diamond blastx  
--threads 4 --db cladoi.dmnd  
--query minion.fq.gz  
--query-gencode 5 --ultra-sensitive  
--out mt.outfmt6.tsv --outfmt 6 --header  
--unal 0 --alfmt fastq  
--al minion_drenale.fq
```

Diamond Search

```
diamond blastx  
--query minion.fq.gz  
--db cladoi.dmnd  
--query-gencode 5  
--ultra-sensitive  
--out mt.outfmt6.tsv  
--unal 0  
--alfmt fastq  
--al minion_drenale.fq
```

Program (BLAST X / P)

File input (**sequenced**)

Database

Codon usage table (**5-invertebrates, 9-
Platyhelminthes**)[#]

Sensitivity *

Name **Output table**

Non-aligned reads are not saved

Output format(fastA/Q)

Output file name

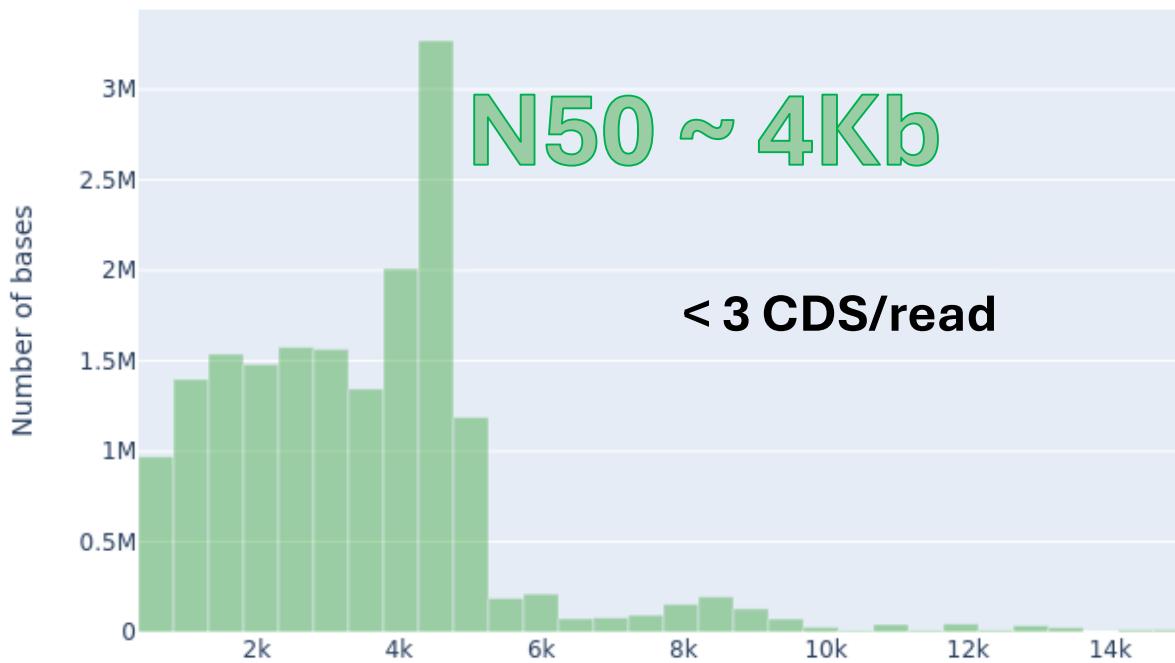
- * <https://github.com/bbuchfink/diamond/wiki/3.-Command-line-options#sensitivity-modes>
- # <https://www.ncbi.nlm.nih.gov/Taxonomy/Utils/wprintgc.cgi>

Alternative genetic codes

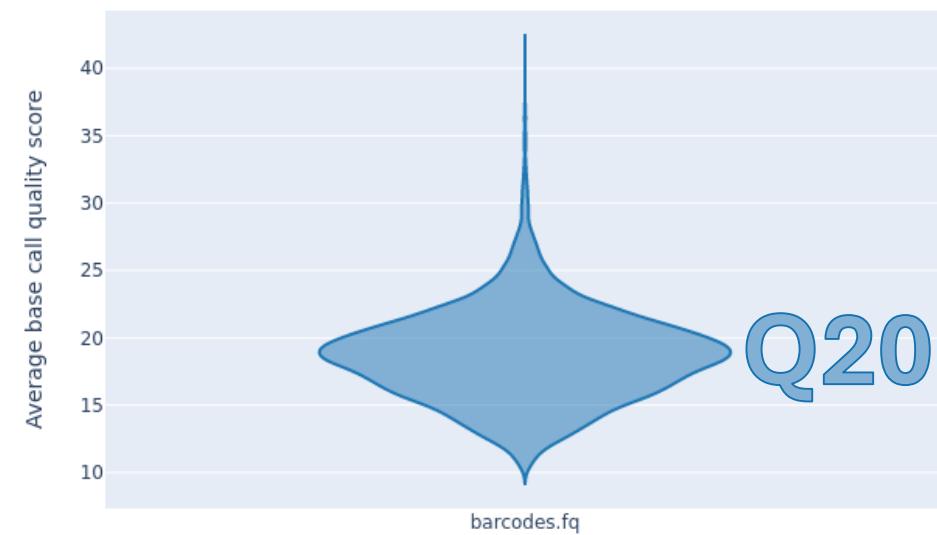
		2. nucleotide					
		U	C	A	G		
U	UUU } F	UCU }	UAU }	UGU }	U		
	UUC }	UCC }	UAC }	UGC }	C		
	UUA }	UCA }	UAA St, Q ₆	UGA St, W _{1,2,3,4,5}	A		
	UUG }	UCG }	UAG St, Q ₆	UGG W	G		
C	CUU }	CCU }	CAU }	CGU }	U		
	CUC }	CCC }	CAC }	CGC }	C		
	CUA }	CCA }	CAA }	CGA }	A		
	CUG L, T ₅ , S ₇	CCG }	CAG }	CGG }	G		
A	AUU }	ACU }	AAU }	AGU }	U		
	AUC }	ACC }	AAC }	AGC }	C		
	AUA I, M _{2,3,4,5}	ACA }	AAA K, N ₁	AGA }	A		
	AUG M	ACG }	AAG K	AGG }	G		
G	GUU }	GCU }	GAU }	GGU }	U		
	GUC }	GCC }	GAC }	GGC }	C		
	GUA }	GCA }	GAA }	GGA }	A		
	GUG }	GCG }	GAG }	GGG }	G		

Read Length & Quality

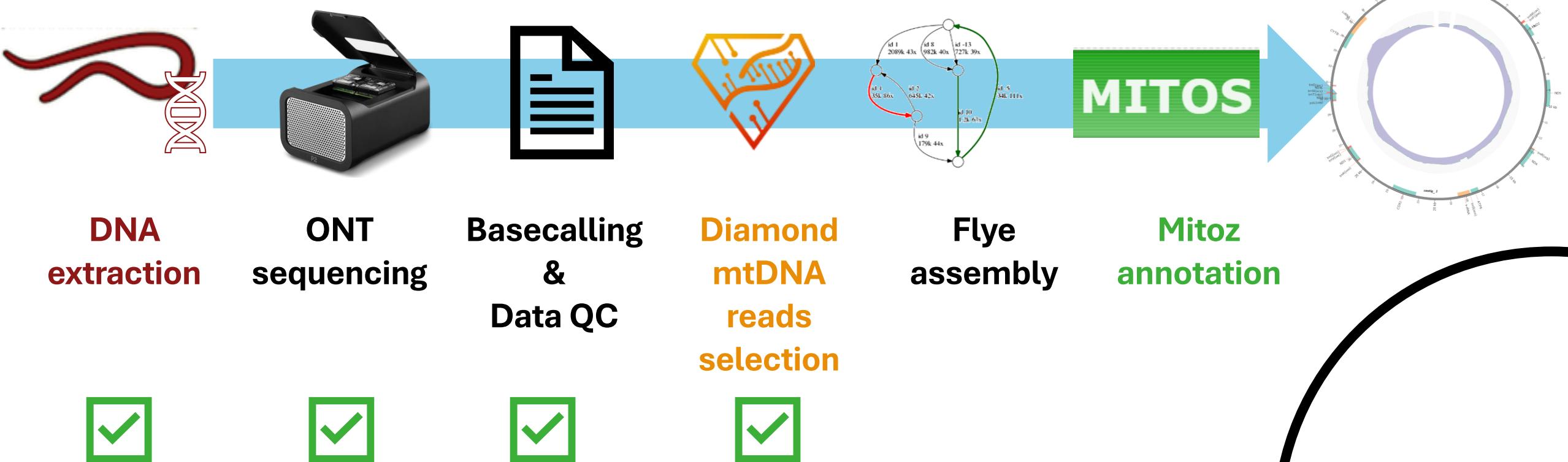
Weighted histogram of read lengths



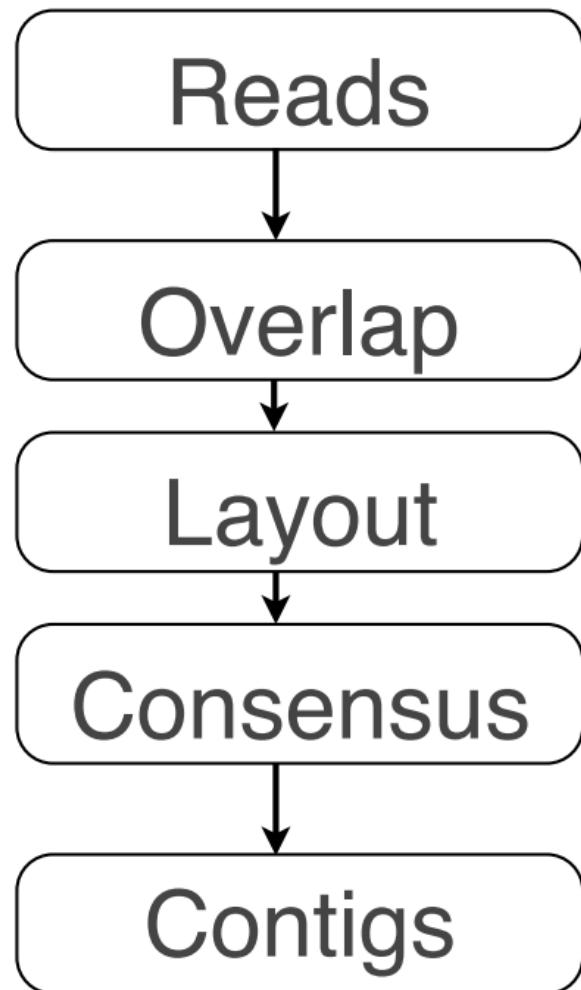
Comparing average base call quality score



Pipeline



Long Read Assembly Pipeline

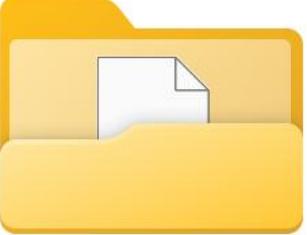
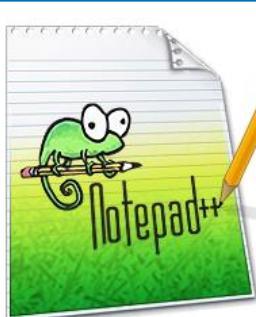
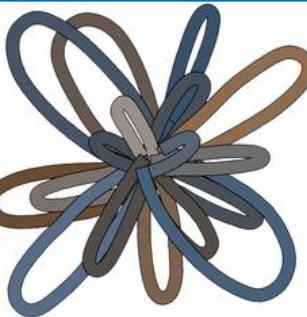
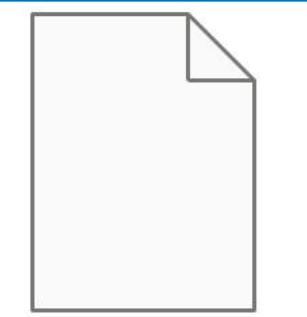
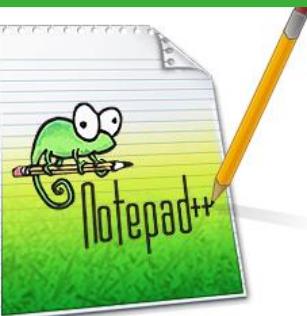


Build overlap graph

Bundle stretches of the overlap graph into contigs

Pick most likely nucleotide sequence for each contig

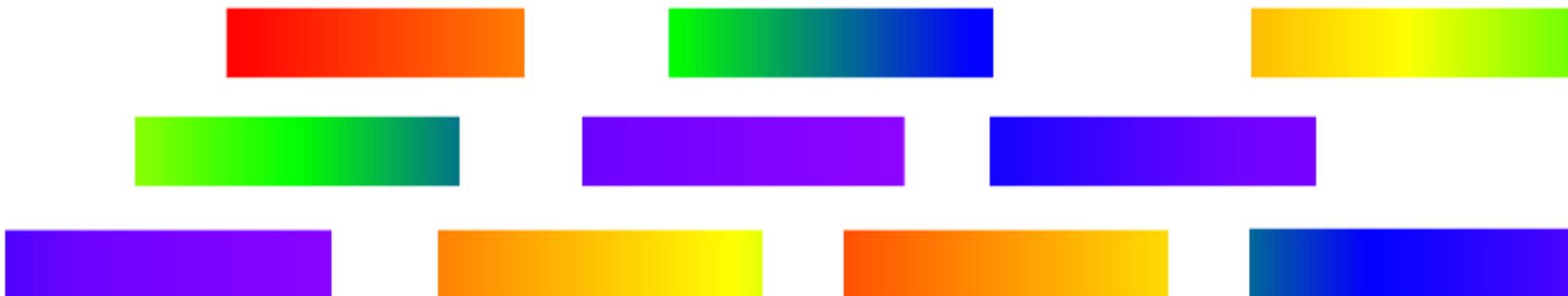
Flye

First assembly	Repeat resolution	Contigger	Final polish
			
00-assembly	20-repeat	30-contigger	40-polishing
			
assembly	assembly_graph	assembly_graph.gv	assembly_info
			
fly	params		

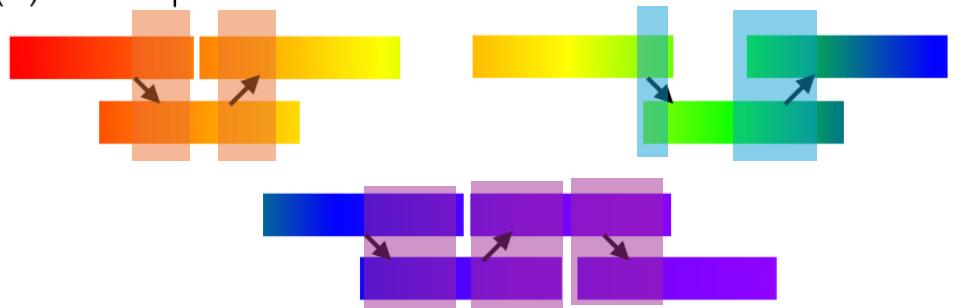
GENOME



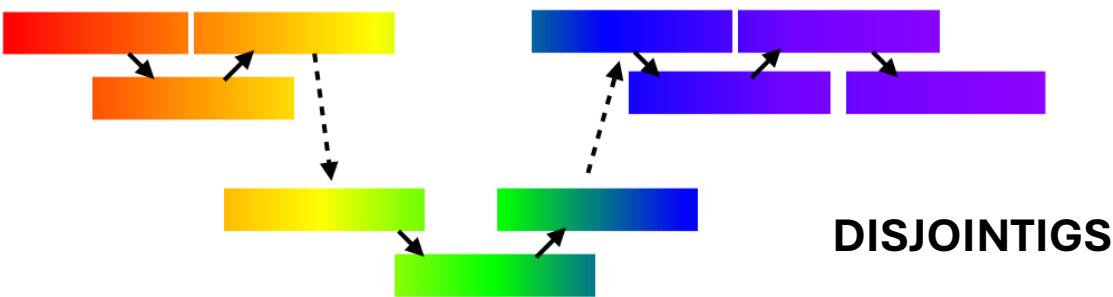
READS



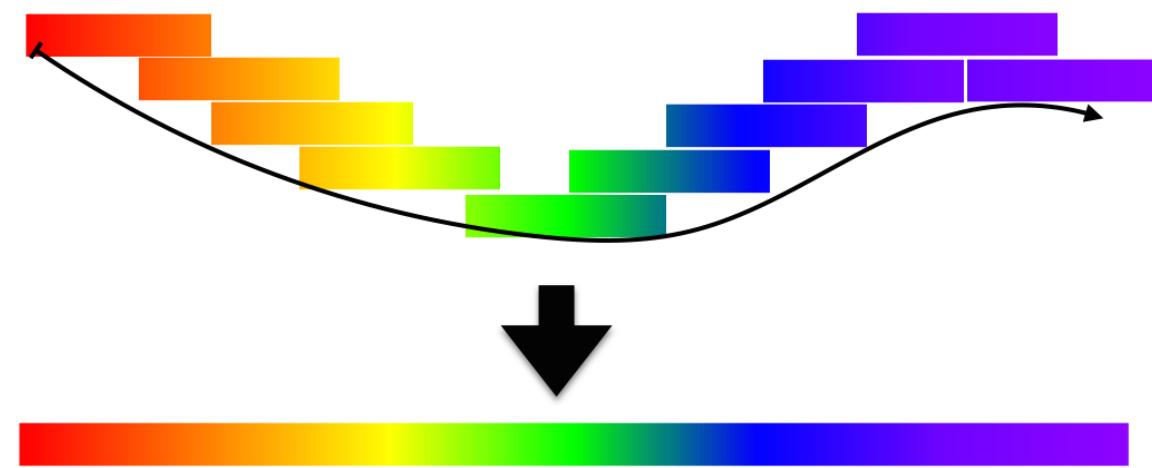
(1) Overlap

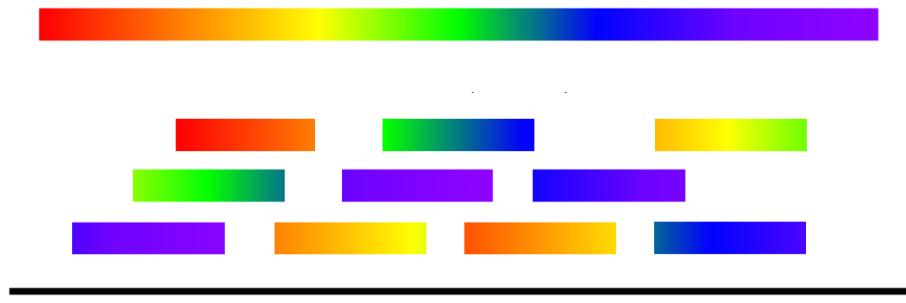


(2) Layout

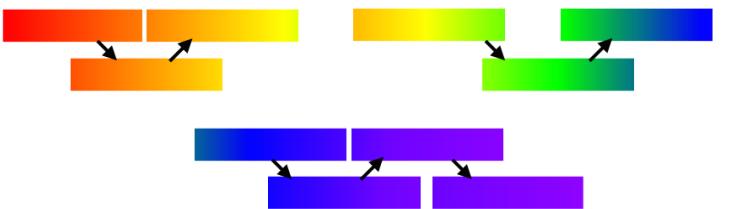


(3) Consensus

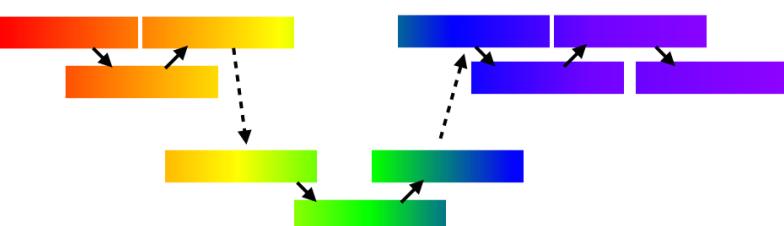




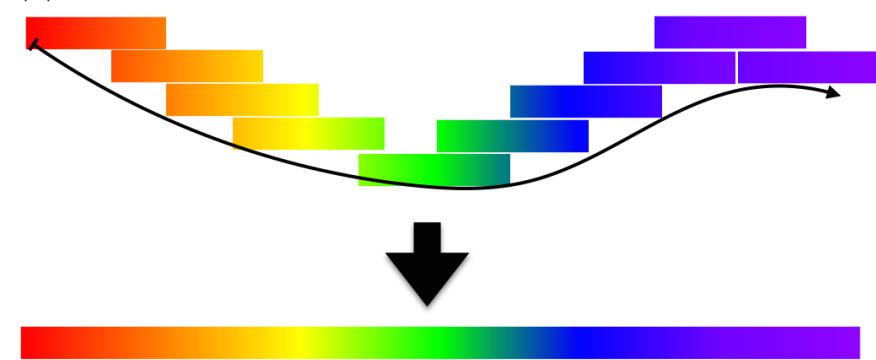
(1) Overlap



(2) Layout



(3) Consensus



Assembly with Flye

```
flye  
--nano-raw drenale.fq  
-t 4  
--meta  
--keep-haplotypes  
-o ./flye
```

Program

Input file (**sequenced**)

Threads

For metagenomes or **with variable coverage**

Maintain haplotypes
(do not collapse alternative haplotypes)

Name **Output folder**

Sequencing methodology-dependent options:

--nano-raw Nanopore regular reads (<20% error)

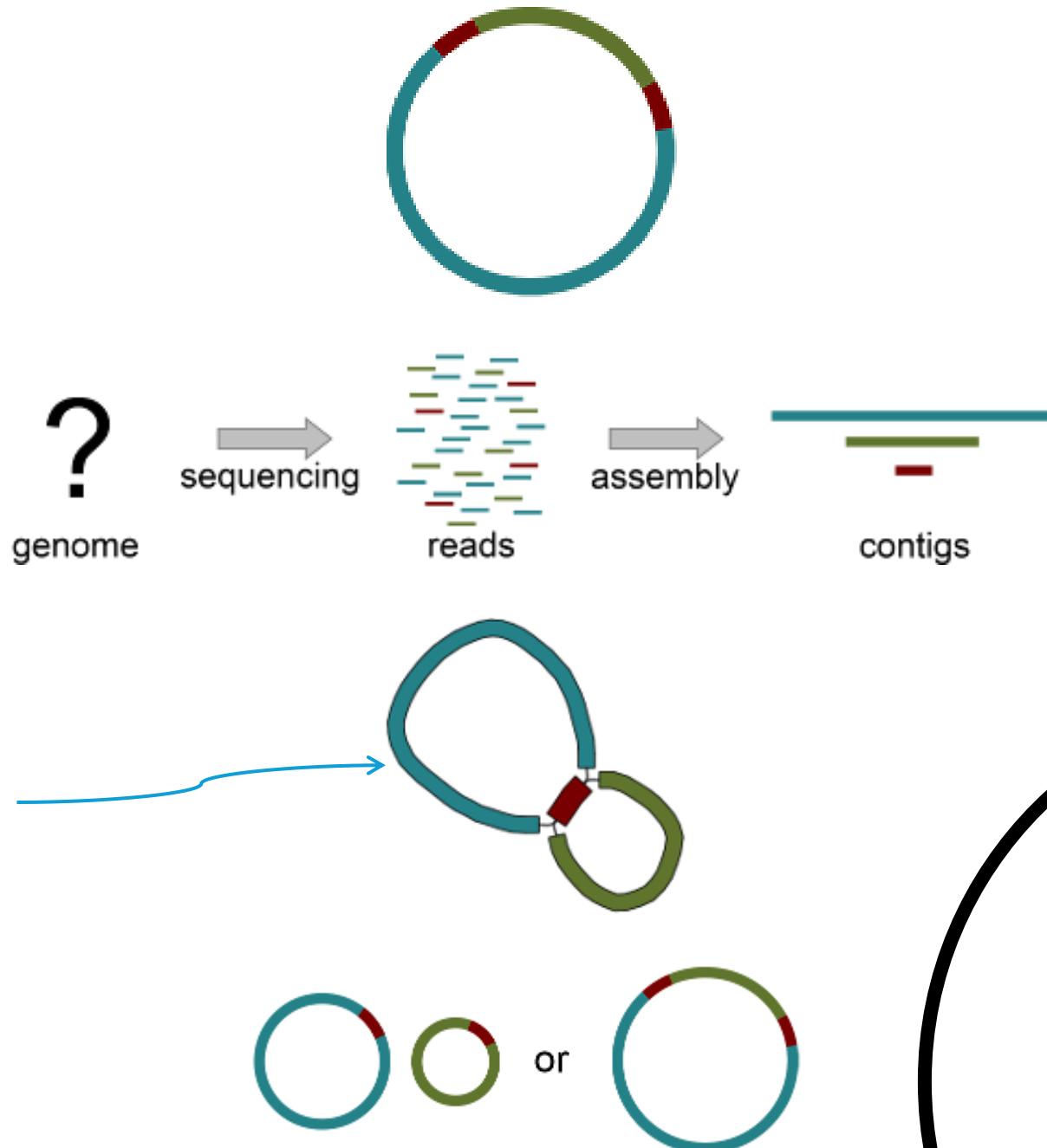
--nano-hq Nanopore high-quality reads: SUP or **Q20** (<3% error)

OUR CASE

Similar options are designed for Pacbio:
(<https://github.com/mikolmogorov/Flye/blob/flye/docs/USAGE.md#inputdata>)

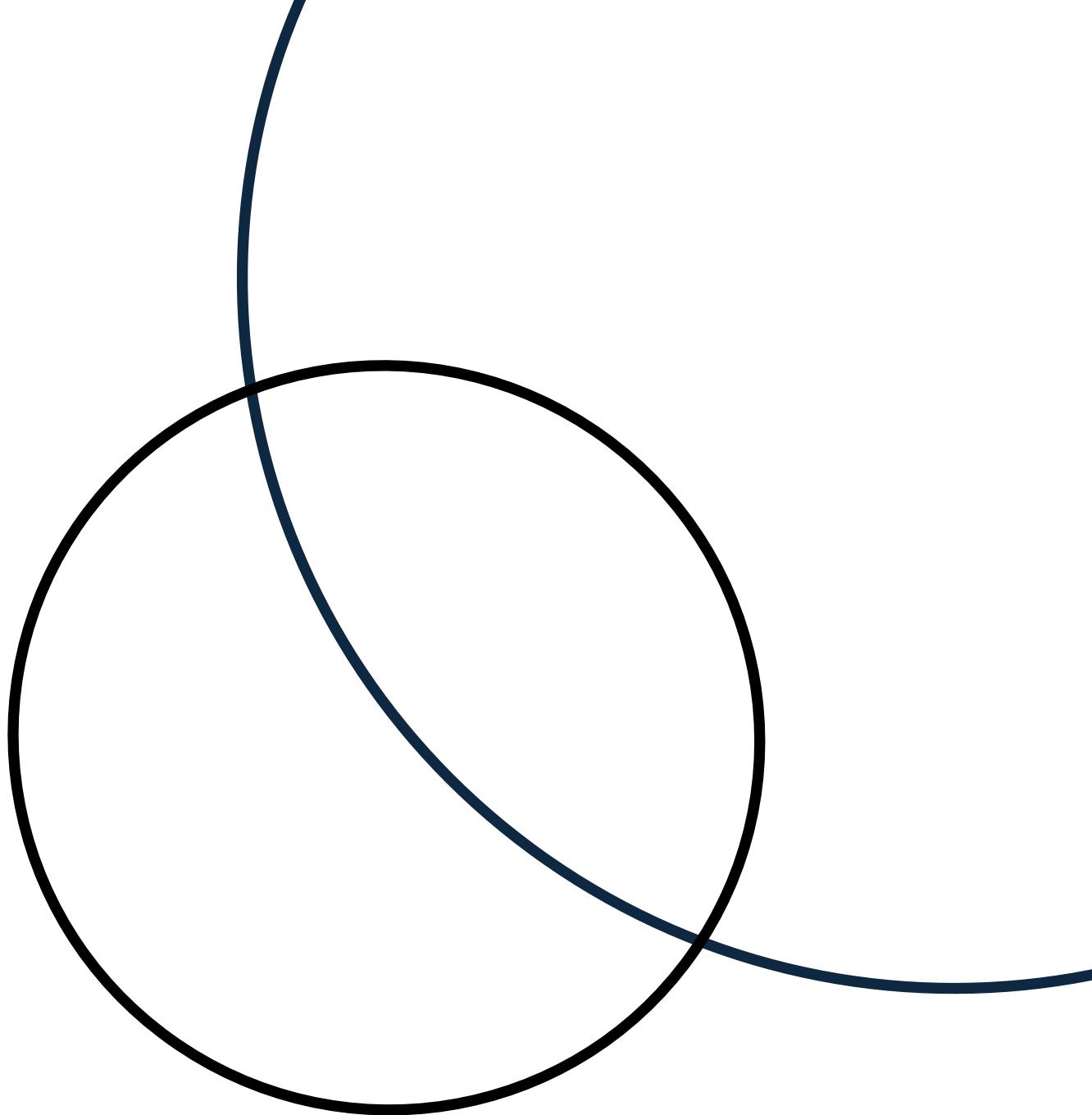
Assembly QC

- assembly_info.txt :
 - Number of Contigs
 - Total length of assembly
- Visualize GFA in Bandage

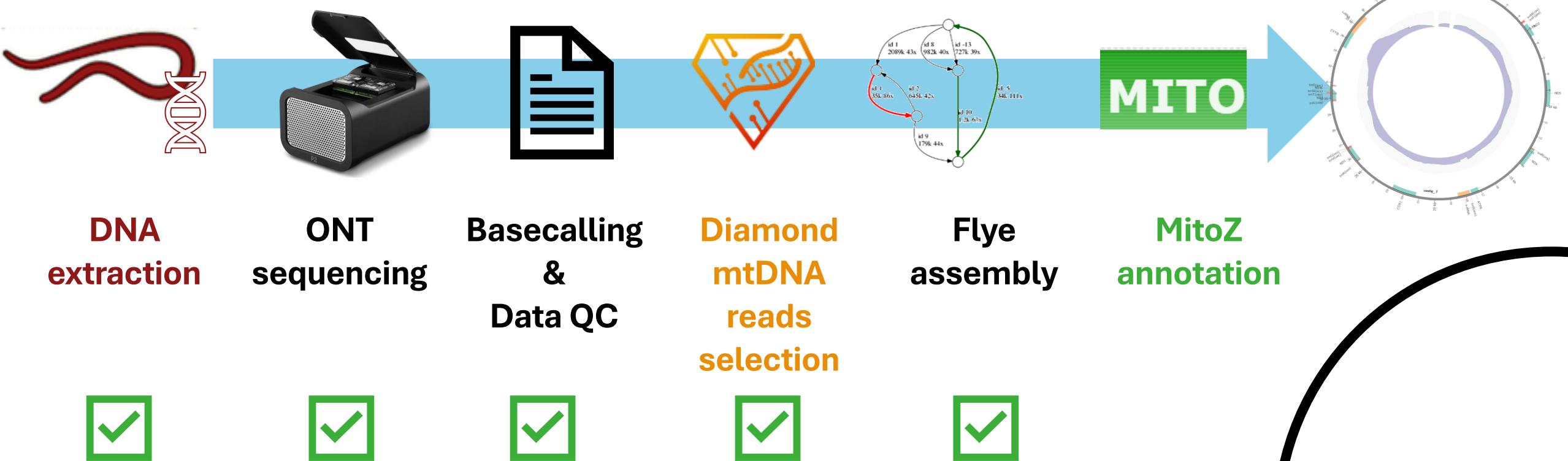


Any questions?

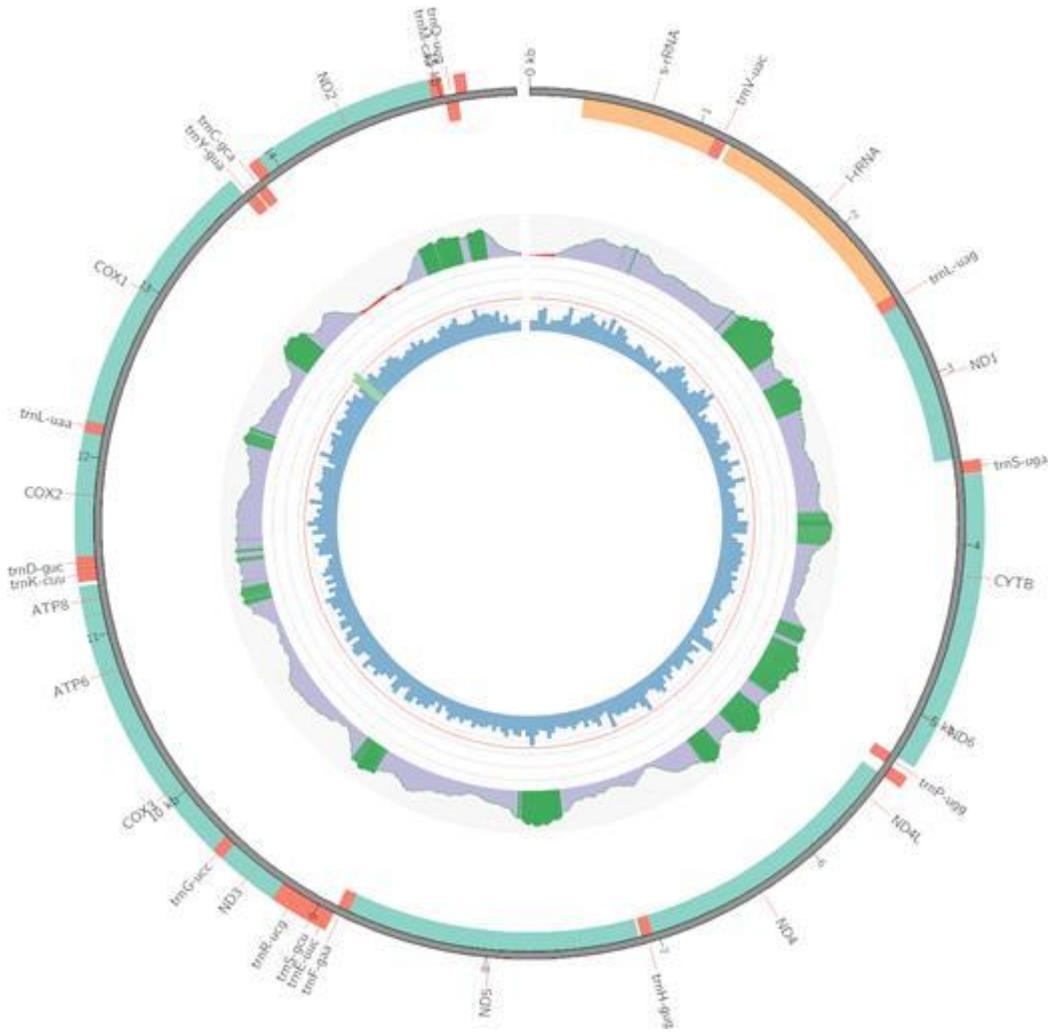
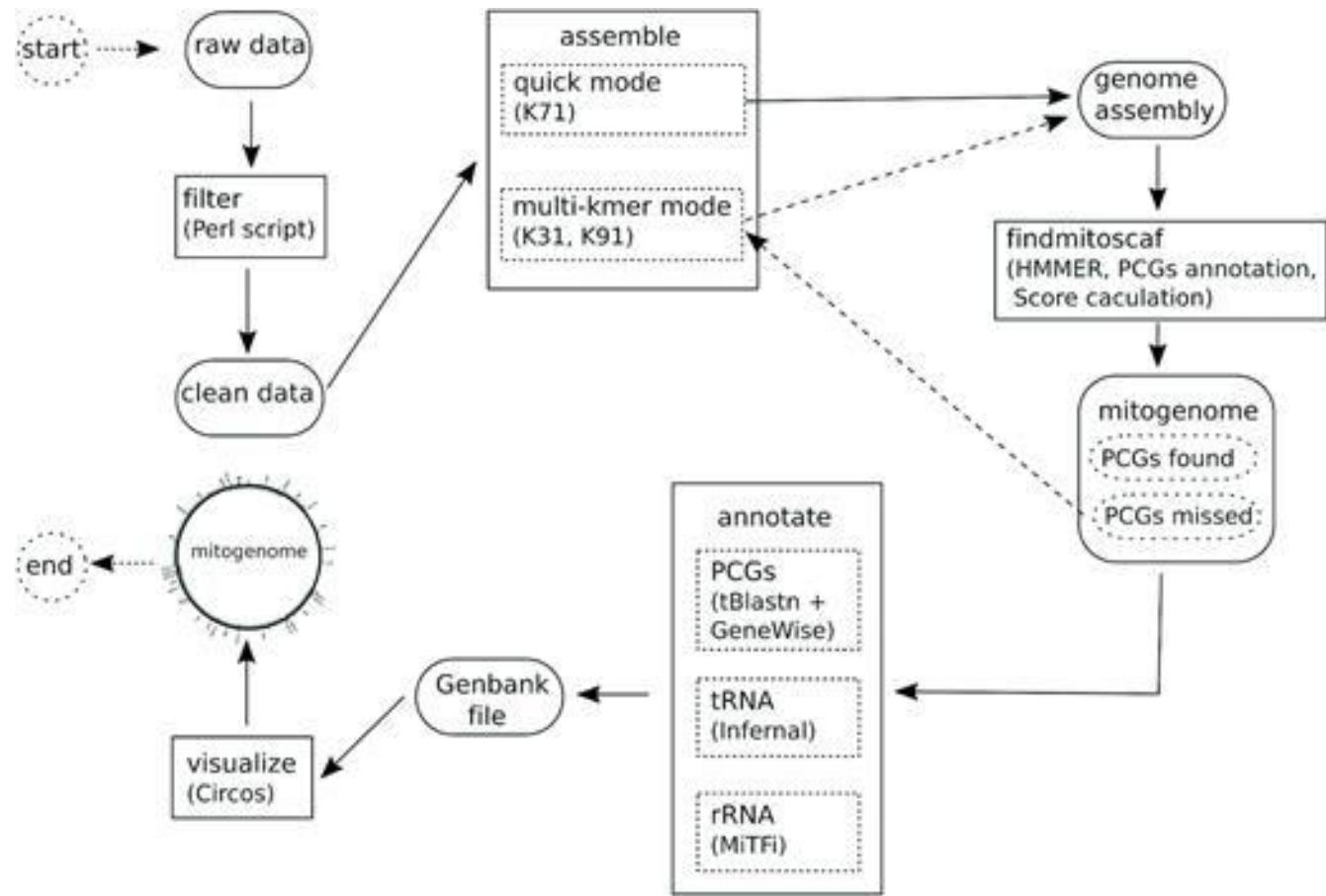
Please contact wellcomeconnectingscience.org
for more information.



Pipeline



MITOZ



MITOZ

```
mitoz annotate  
--workdir ./  
--fastafiles cox1.fasta  
--clade Nematoda  
--genetic_code 5  
--outprefix mitoz_drenale
```

Program

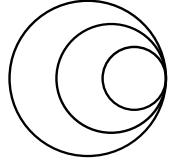
Working directory name

Fasta file with mitogenome assembly

Clade

Clade specific genetic code

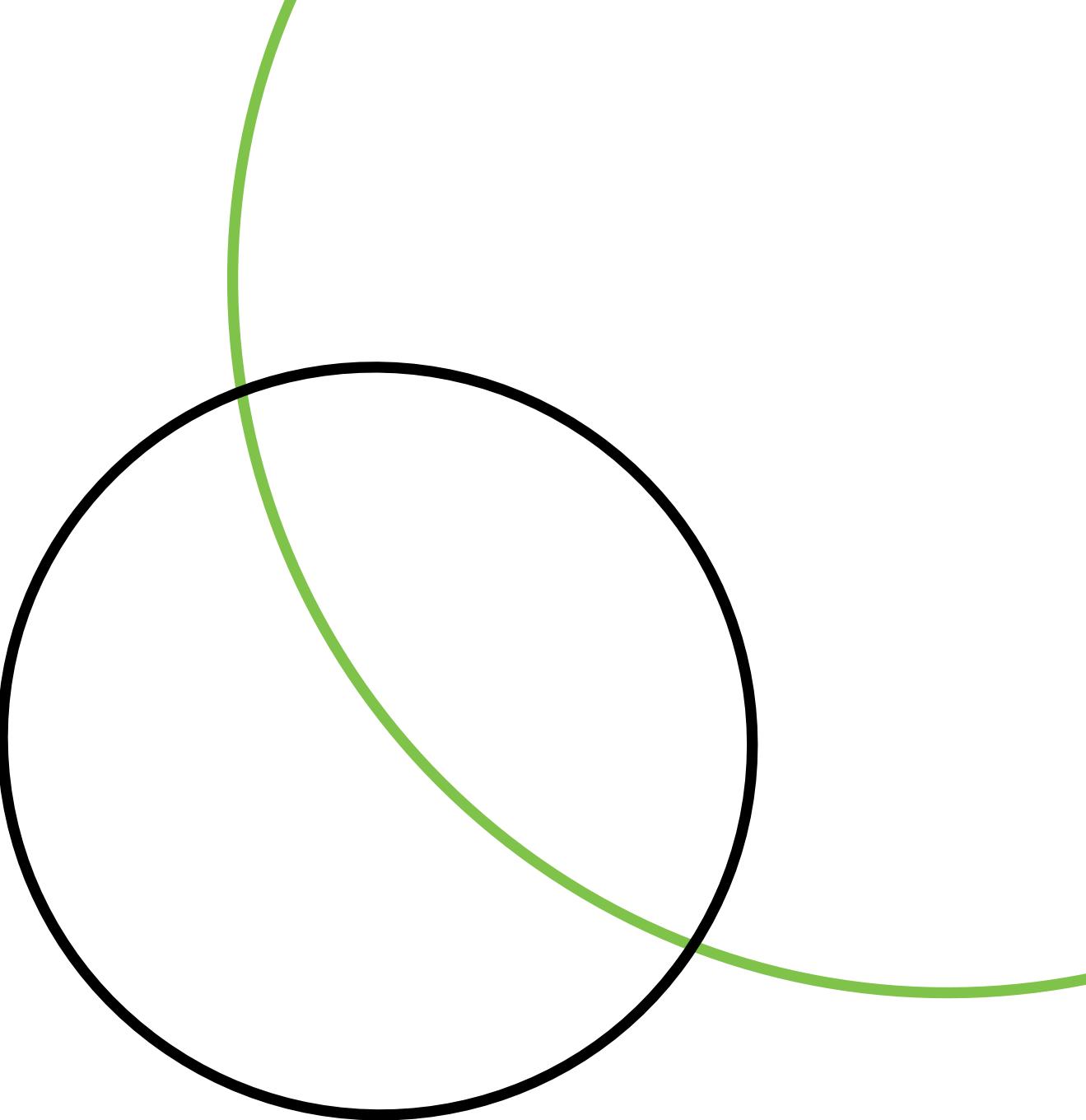
Name **Output folder**



wellcome
connecting
science

thanks!

Please contact wellcomeconnectingscience.org
for more information.



May 2025

Dioctophyme renale

Caracterización mitogenómica

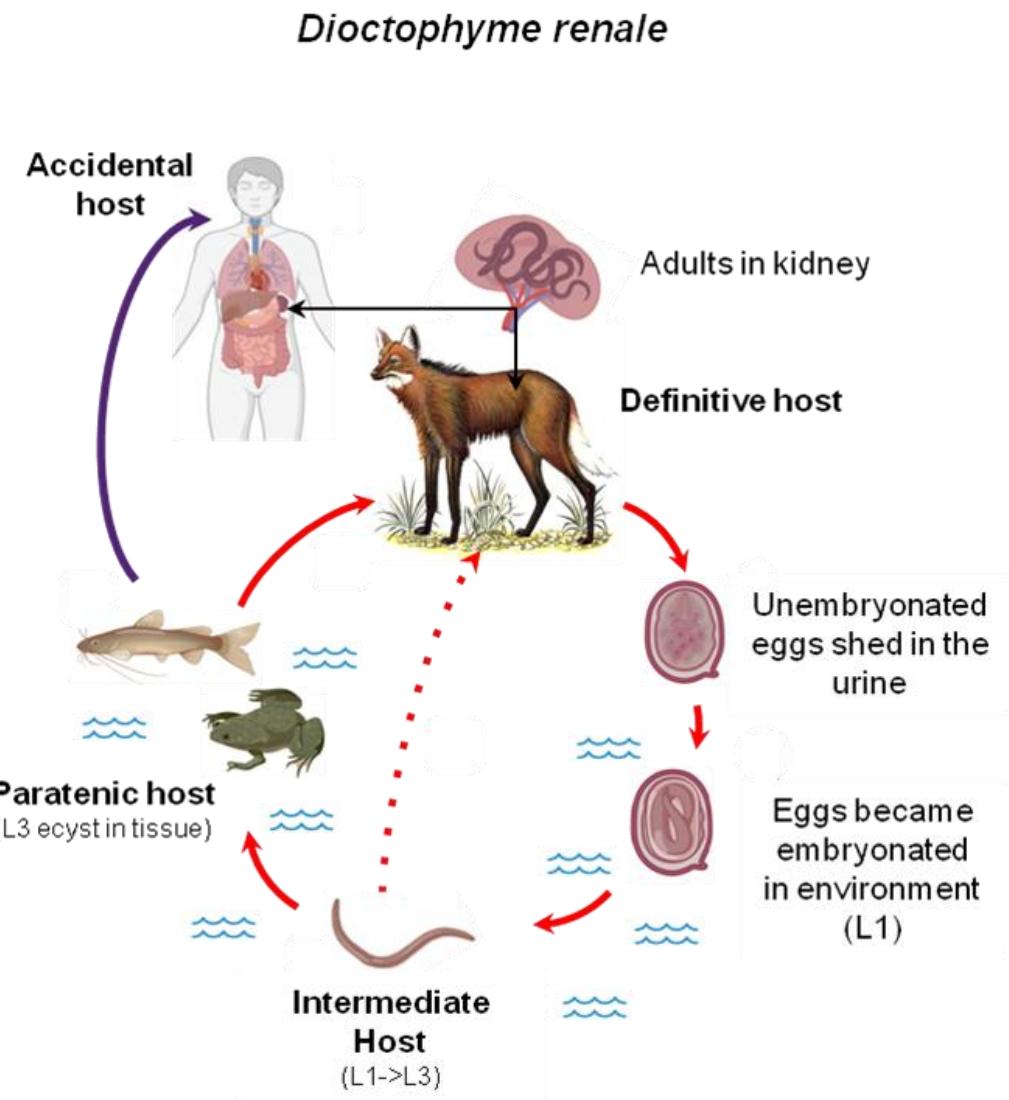


Dioctophyme renale

“Giant kidney worm”

- **Nematode parasite**
- Considered the **largest parasitic nematode** of terrestrial vertebrates
- Health problem in **dogs and threatened wildlife** living near aquatic environments
- High risk of causing **infections in human** populations in riparian areas
- **Hard to sample.**
or roadkill wildlife (*degraded DNA*)
- Little molecular information on this organism
(no **genome** nor **transcriptome**)

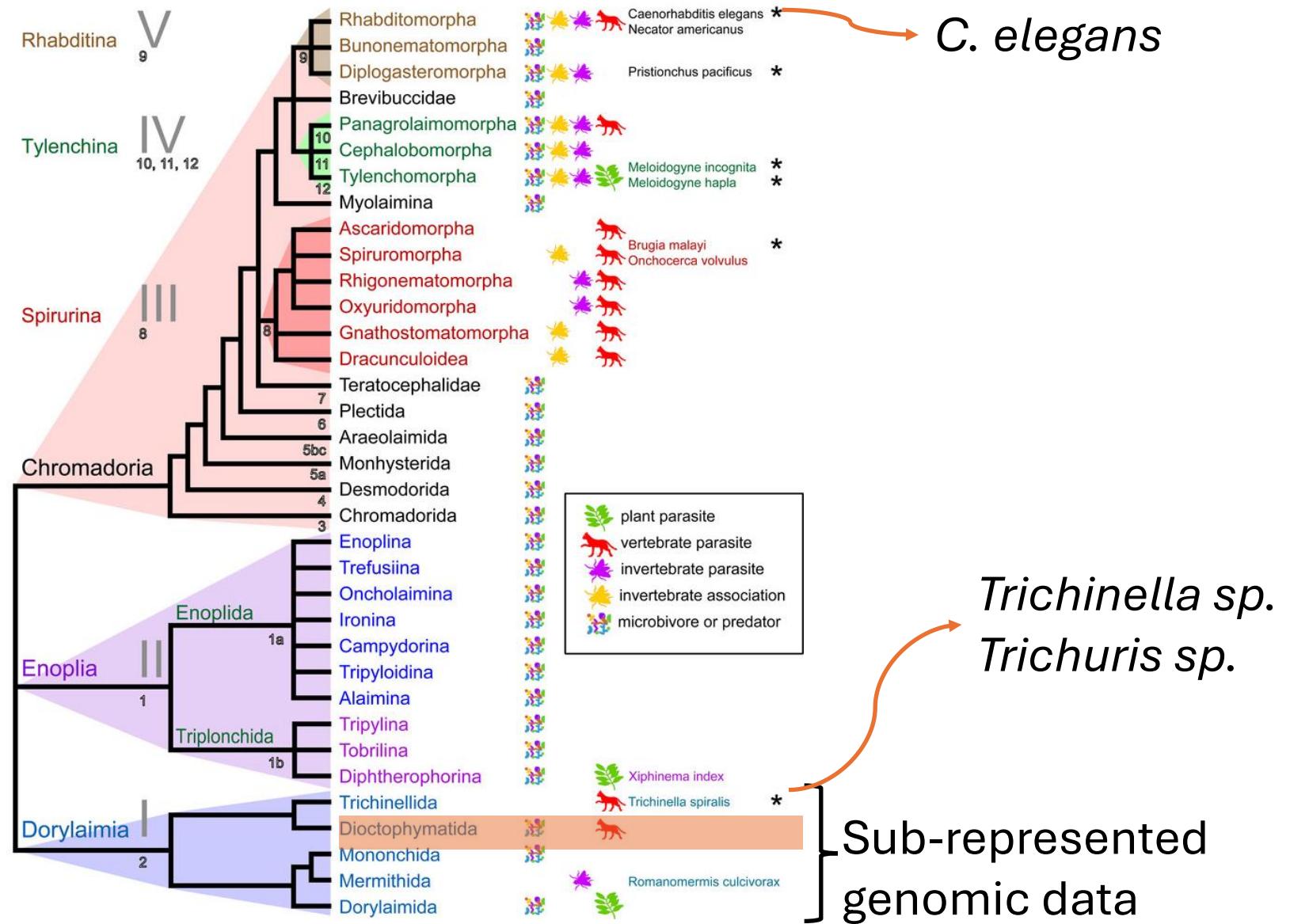
surgeries



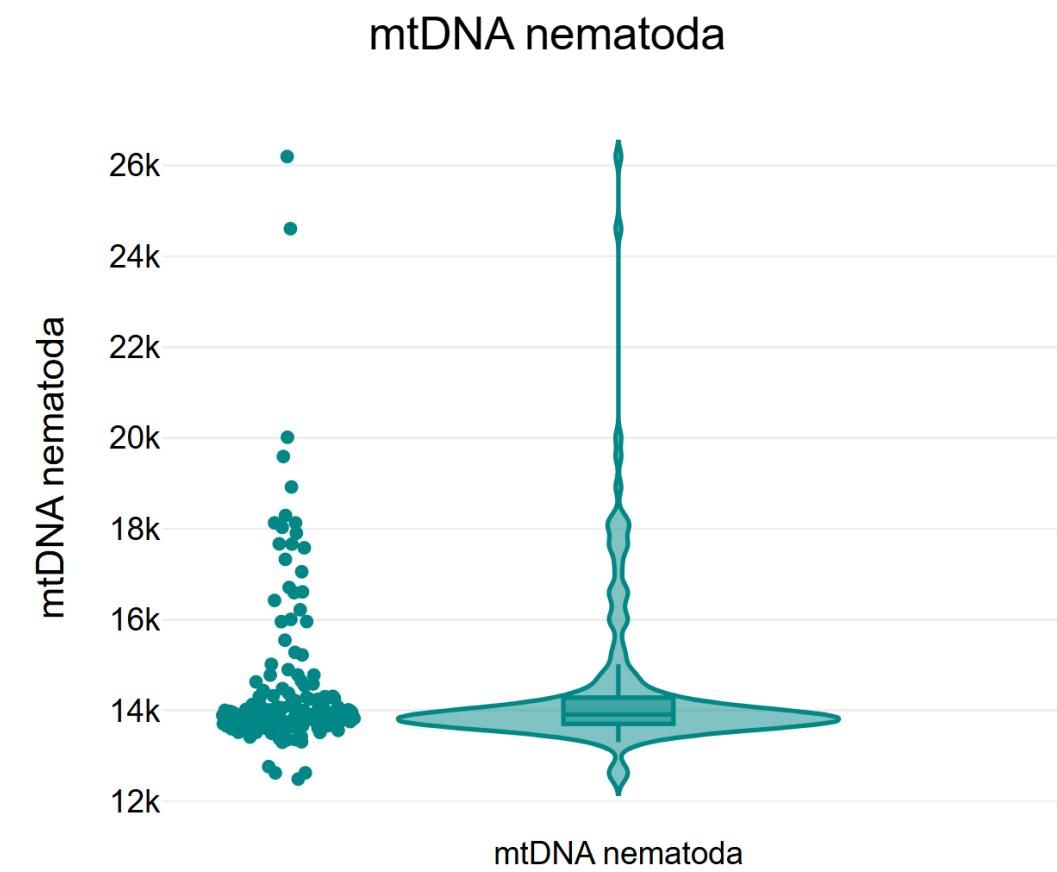
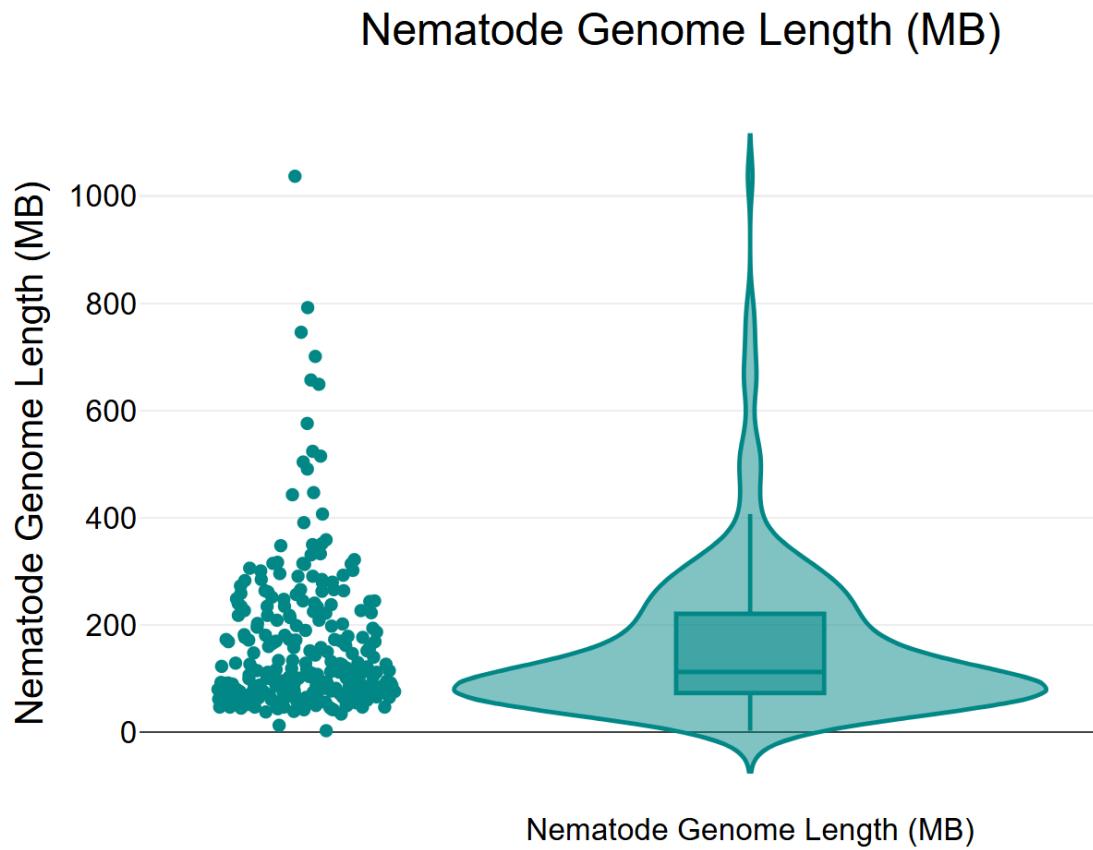
Phylogenetic Relationships - Nematodes

mtDNA has valuable characteristics :

- high mutation and substitution rates
- infrequent genetic recombination
- maternal transmission
- high copy number (!)
- easy accessibility.



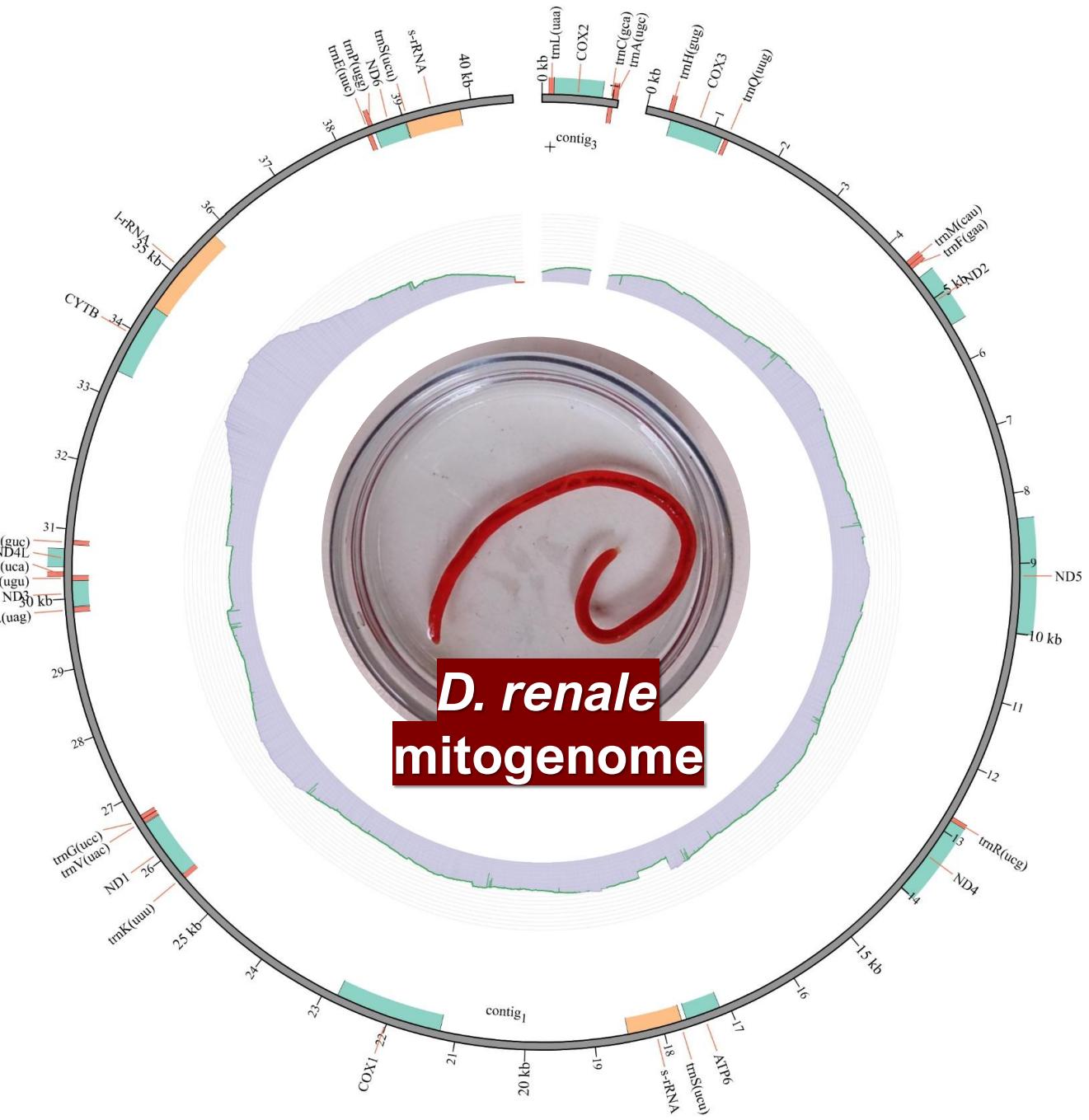
(Mito)Genomes lengths in Nematodes

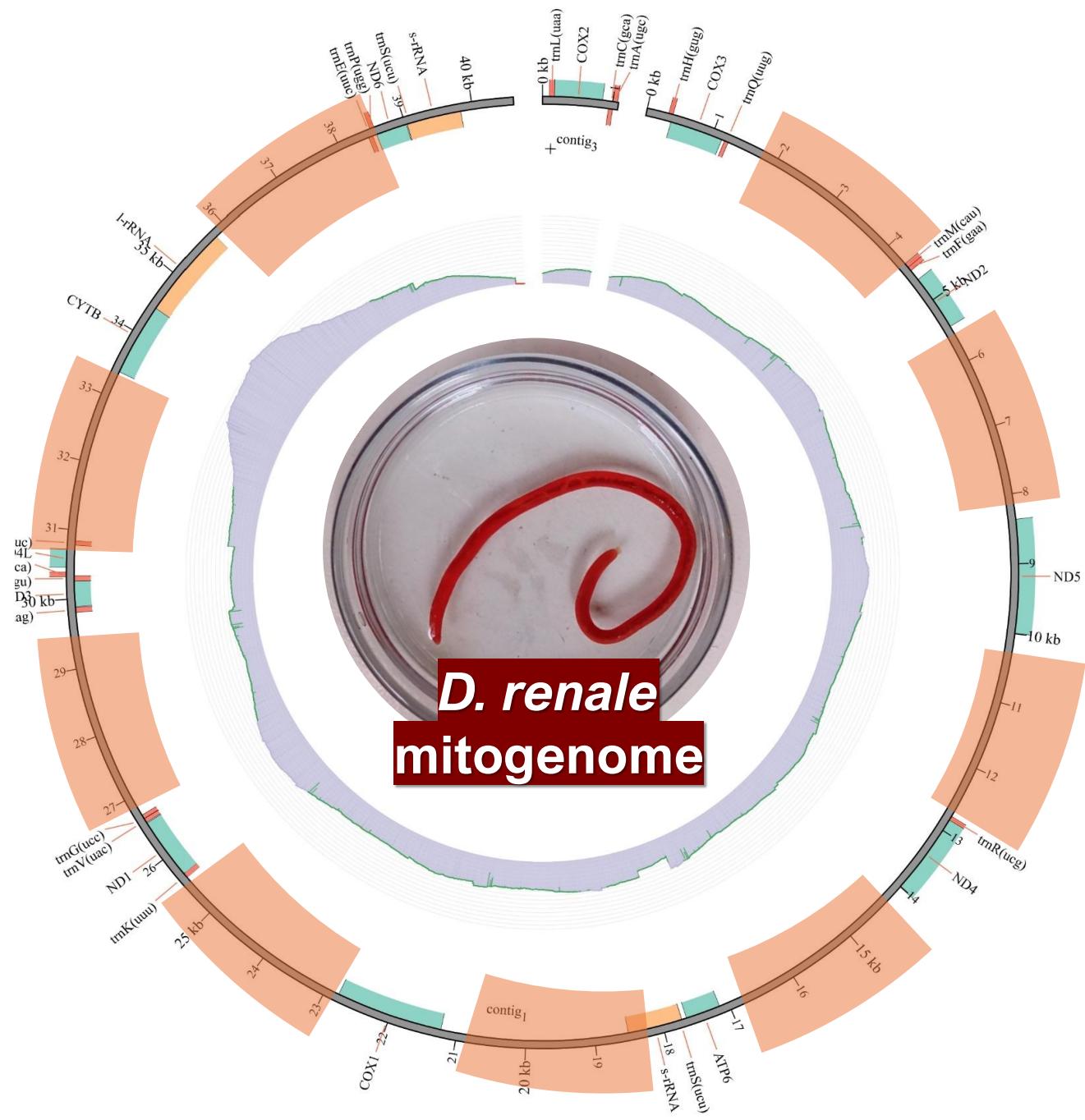
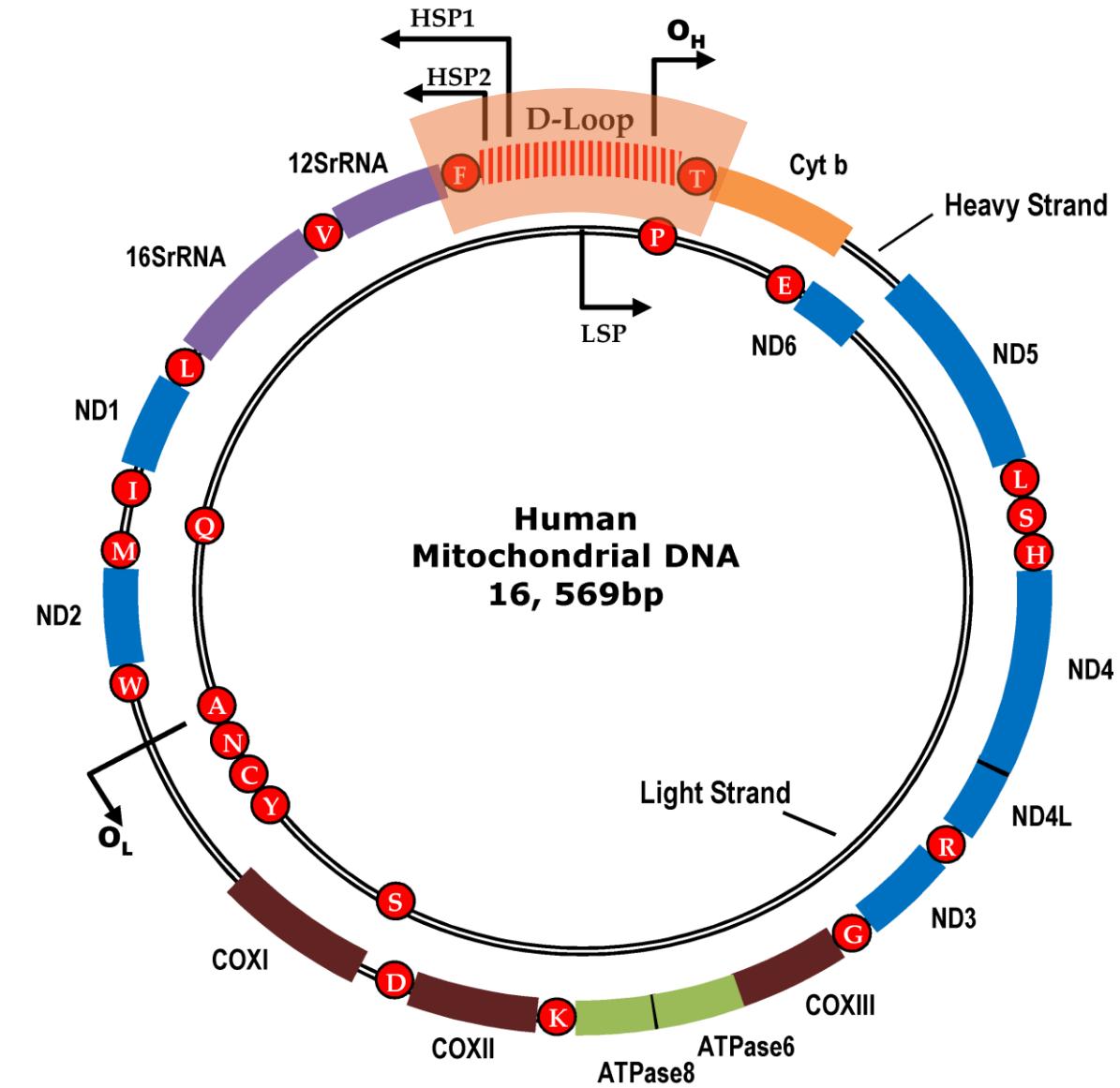


ATP6
COX1
COX2
COX3
CYTB
NAD1
NAD2
NAD3

NAD4
NAD4L
NAD5
NAD6

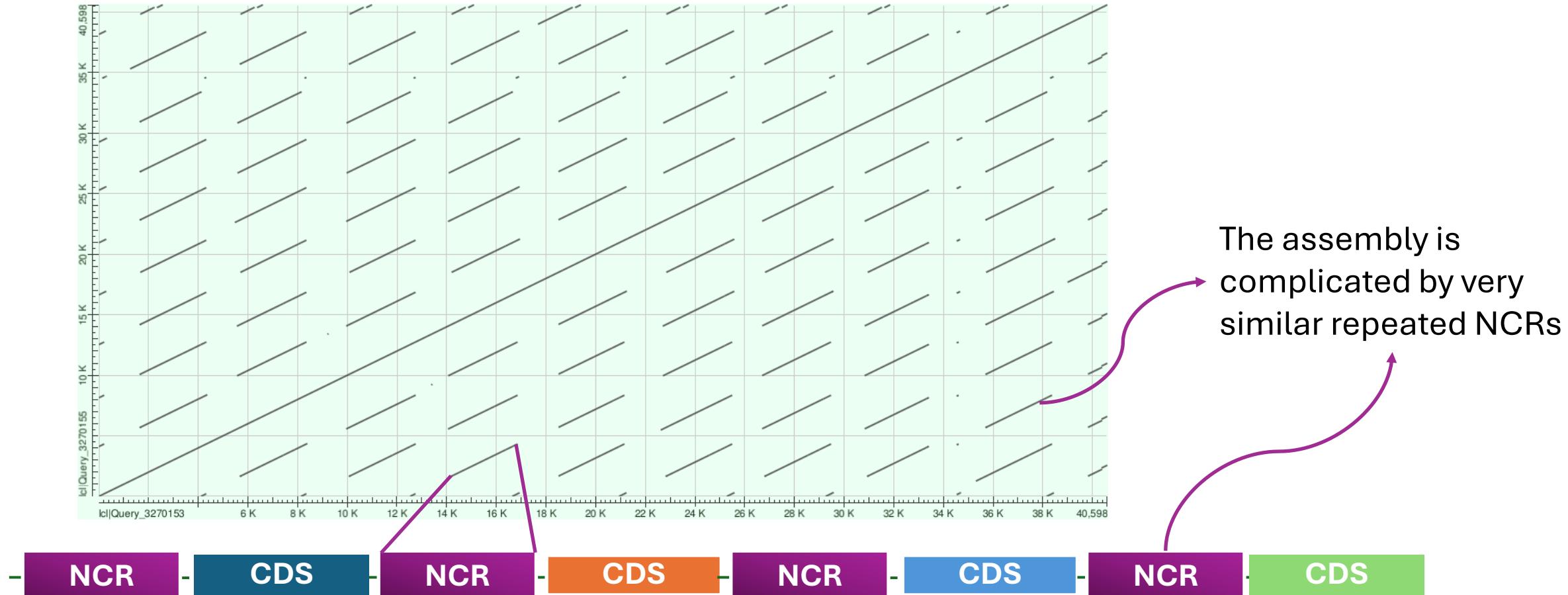
rRNA-L
rRNA-S
22 x tRNA



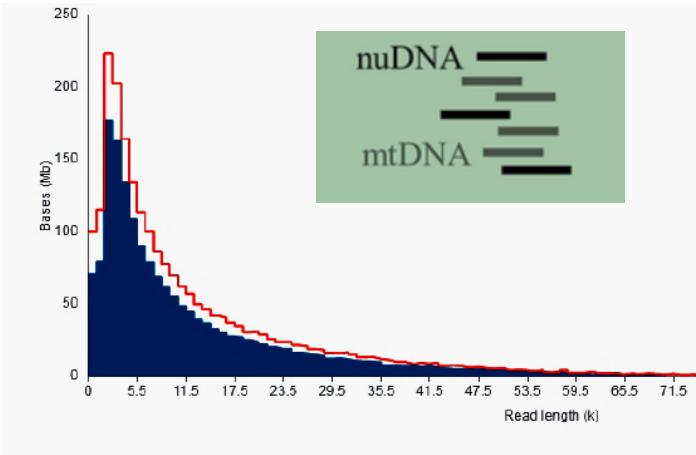


Non-Coding regions (NCR)

Each CDS is flanked by an NCR

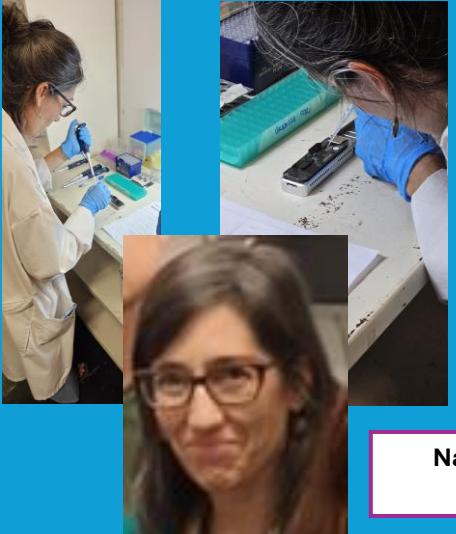


Conclusions



- ❖ Genome skimming bioinformatics allowed us to **assemble and annotate** a new mitochondrial genome
 - ✓ **Without a reference genome**
 - ✓ **No mitochondrial separation required**
 - ✓ **No PCR amplification needed**
 - ✓ All mitochondrial proteins were annotated with **82-99% AA identity** to other phylogenetically related nematodes
 - ✓ All mitochondrial **rRNA** and **tRNA** were found
 - ✓ Two mitochondrial **molecular markers** were designed and implemented in phylogeography studies
 - ✓ Unique nematode mitochondrial **features** were identified

Genomics and Bioinformatics LAB10. iB3-UBA-CONICET



Natalia Macchiaroli
Researcher



Lucas Arce
Fellowship



Prof. Gisela Franchini
INBIOLP-UNLP-CONICET



Kevin Calupiña
Fellowship



Ines Sananez
Fellowship



Marina Ingravidì Fellowship



Juan Arrabal
Fellowship



Prof. Laura Kamenetzky
Principal Investigator



Melisa Magallanes
Fellowship



Prof. Enrique Caviedes-Vidal
Principal Investigator



Instituto de Biociencias,
Biotecnología y Biología traslacional

INIBIOLP
Instituto de Investigaciones Bioquímicas de La Plata
"Profesor Doctor Rodolfo R. Brenner"

CONICET

I M I B I O - S L