

Genomics & Bioinformatics

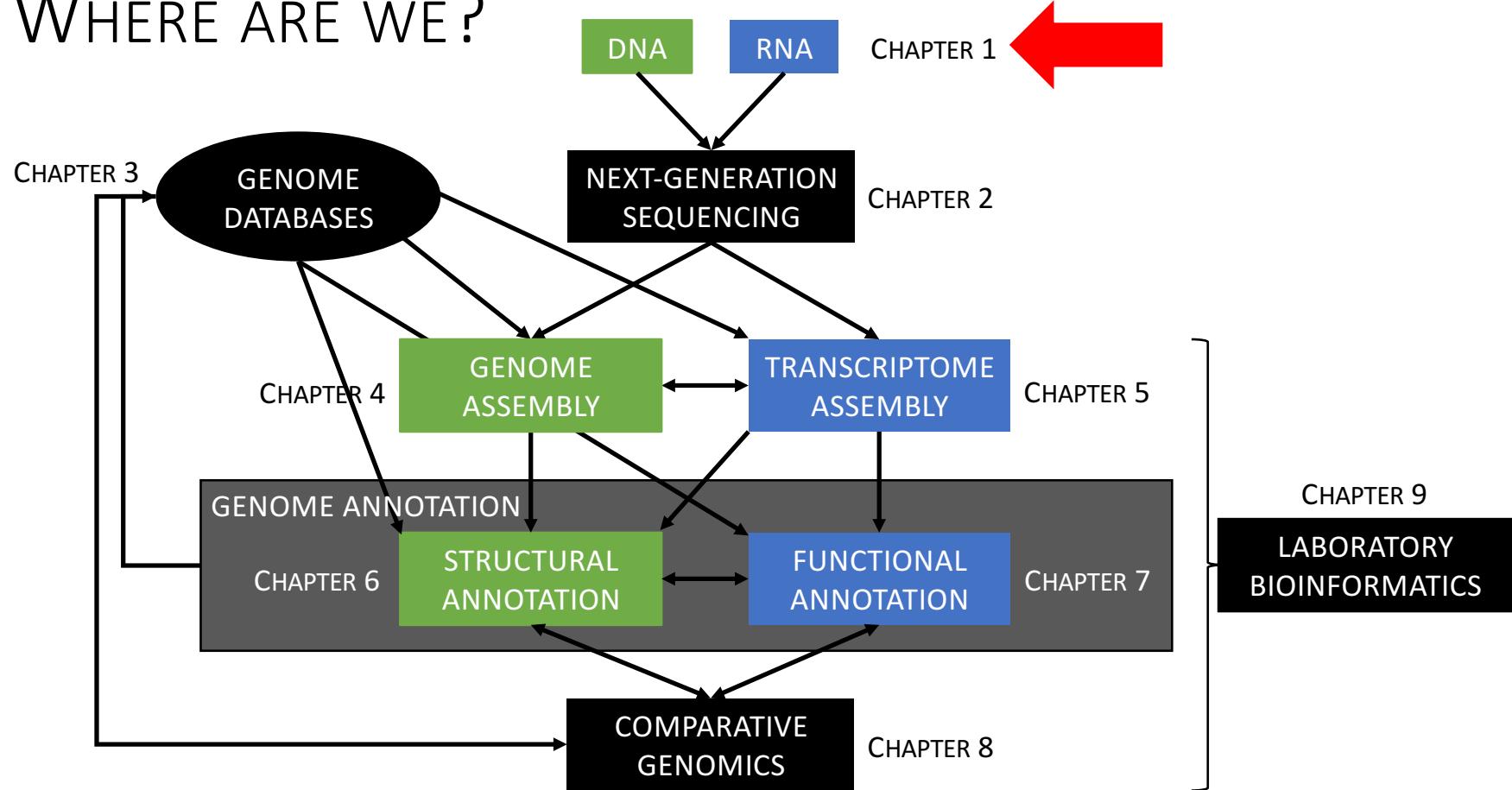


BIOL 497, 597

Boise State University

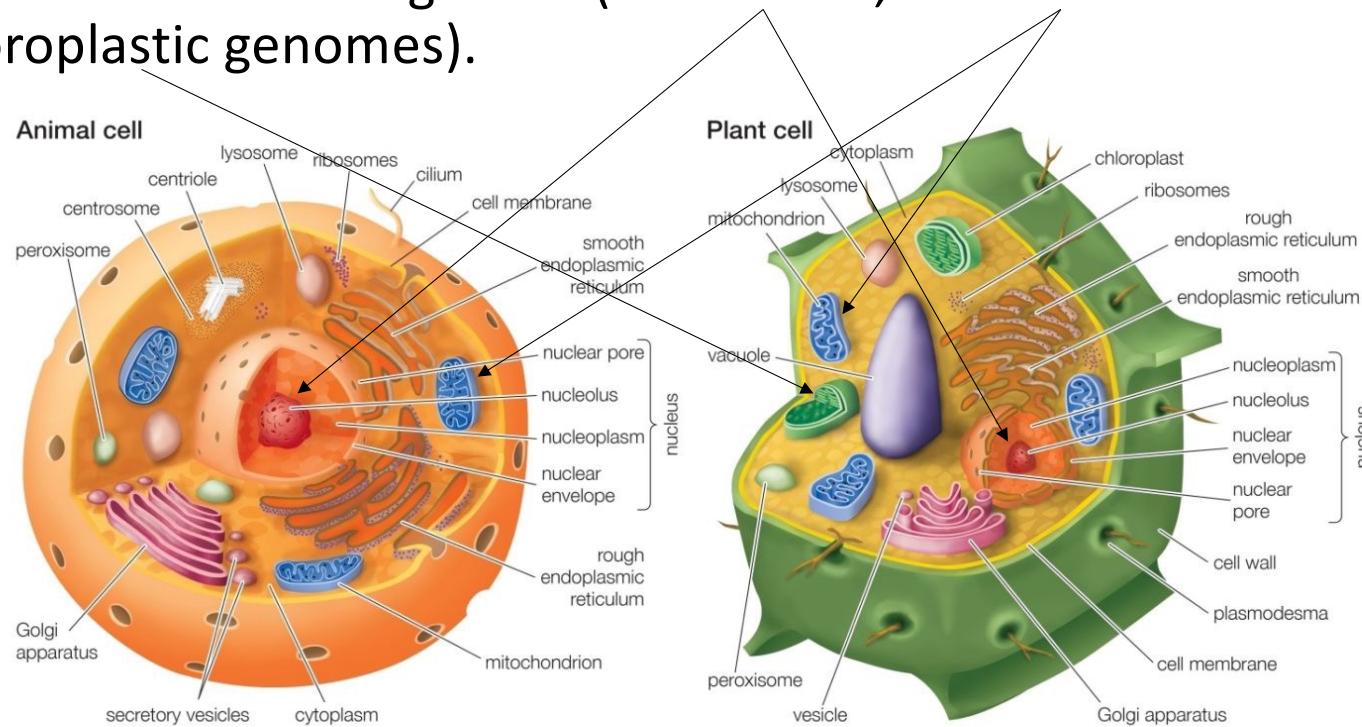
Spring 2022

WHERE ARE WE?



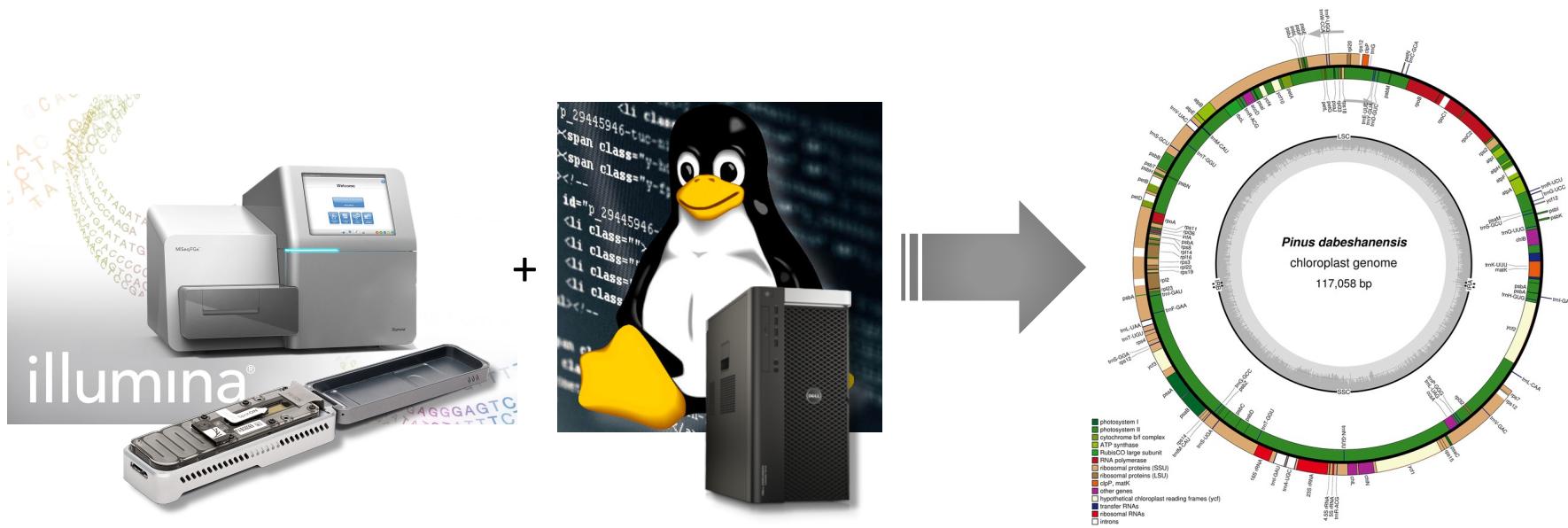
WHAT IS GENOMICS?

- **Genomics:** The study of genomes, the complete set of genetic material within an organism (i.e. nuclear, mitochondrial and chloroplastic genomes).



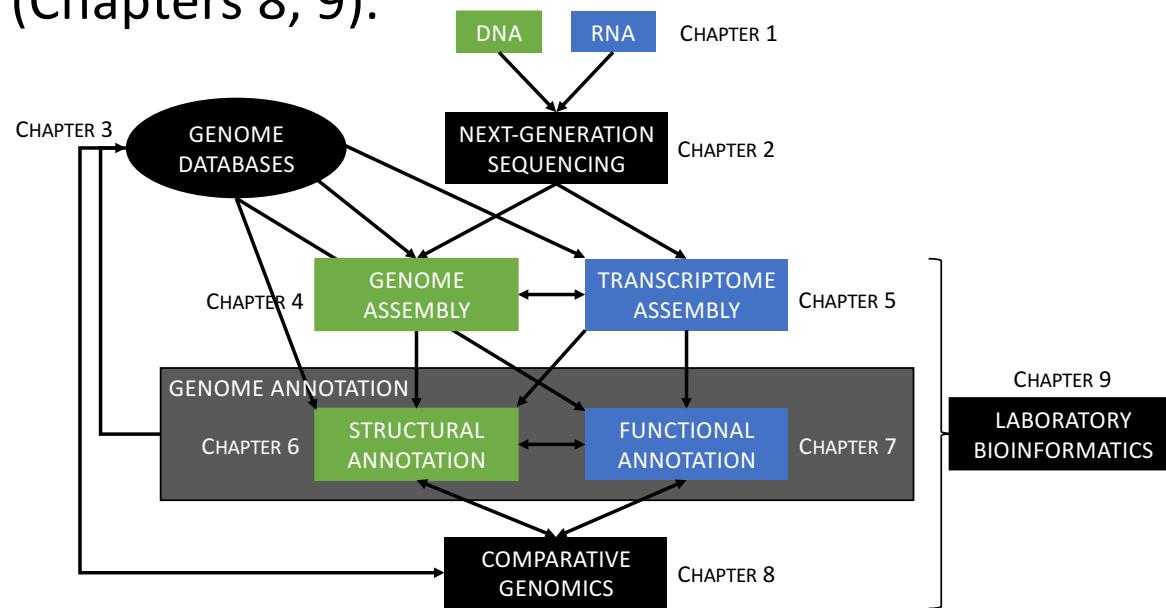
WHAT IS GENOMICS?

- **Genomics** uses an approach combining **next-generation sequencing** (hereafter NGS; Chapter 2 and Mini report 1) and **bioinformatics** (Chapter 9) to **assemble and annotate entire genomes**.

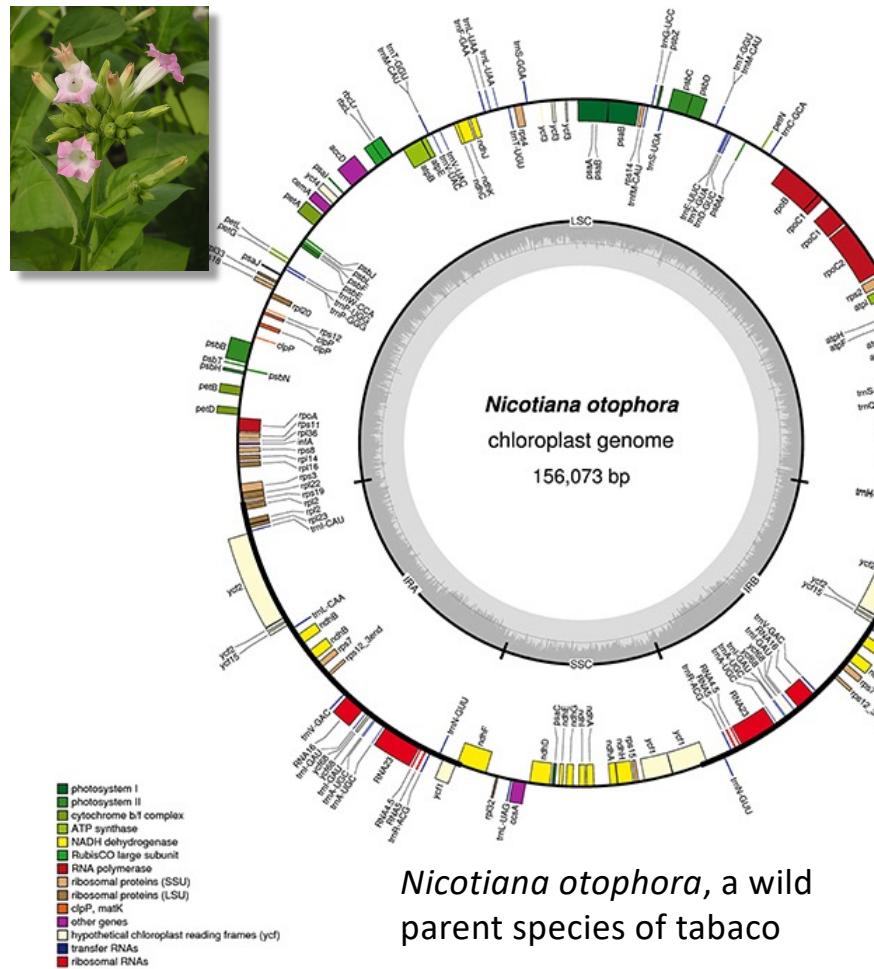


WHAT IS GENOMICS?

- Ultimately genomics aims at **inferring genomes' structures and functions** (Mining genome; Chapter 9).
- Such objectives are usually achieved by applying a comparative approach (Chapters 8, 9).



STRUCTURE AND FUNCTION OF CHLOROPLAST GENOME



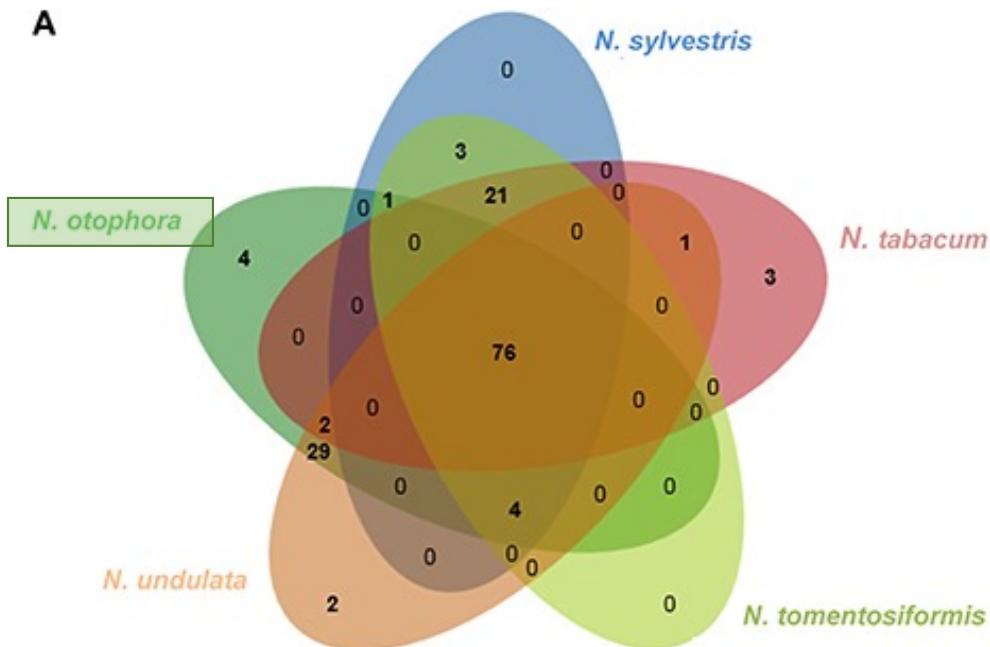
Category	Group of genes	Name of genes
Self-replication	Large subunit of ribosomal proteins	<i>rpl2, 14, 16, 20, 22, 23, 32, 33, 36</i>
	Small subunit of ribosomal proteins	<i>rps2, 3, 7, 8, 11, 12, 14, 15, 16, 18, 19</i>
	DNA dependent RNA polymerase	<i>rpoA, B, C1, C2</i>
	rRNA genes	<i>RNA</i>
	tRNA genes	<i>trnA-UGC, C-GCA, D-GUC, E-UUC, F-GAA, f-m-CAU, G-UCC, H-GUG, I-CAU, L-CAA, M-CAU, N-GUU, P-GGG, P-UGG, Q-UUG, R-ACG, R-UCU, S-GCU, S-GGA, S-UGA, T-GGU, T-UGU, V-GAC, V-UAC, W-CCA, Y-GUA</i>
Photosynthesis	Photosystem I	<i>psaA, B, C, I, J</i>
	Photosystem II	<i>psbA, B, C, D, E, F, H, I, J, K</i>
	NadH oxidoreductase	<i>ndhA, B, C, D, E, F</i>
	Cytochrome b6/f complex	<i>petA, B, D, G, L, N</i>
	ATP synthase	<i>atpA, B, E, F, H, I</i>
	Rubisco	<i>rbcL, rbcR</i>
Other genes	Translational initiation factor	<i>infA</i>
	Maturase	<i>matK</i>
	Protease	<i>clpP</i>
	Envelop membrane protein	<i>cemA</i>
	Subunit Acetyl-CoA-Carboxylate	<i>accD</i>
	c-type cytochrome synthesis gene	<i>ccsA</i>
	Conserved Open reading frames	<i>ycf1, 2, 3, 4, 15, 68</i>

Asaf et al. (2016) Front. Plant Sci.

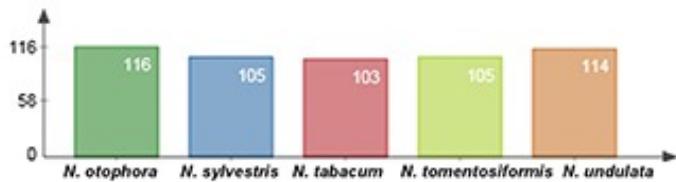
COMPARATIVE GENOMICS

Coding genes

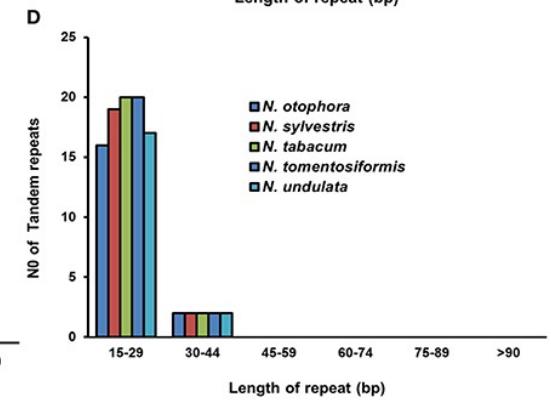
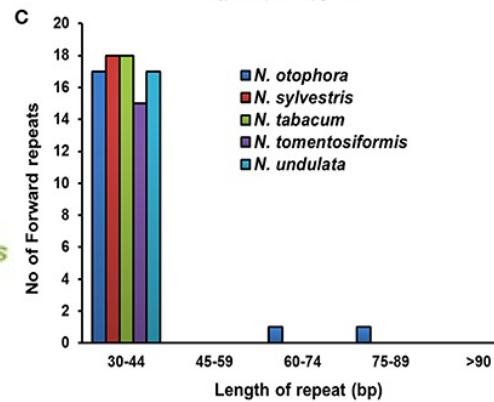
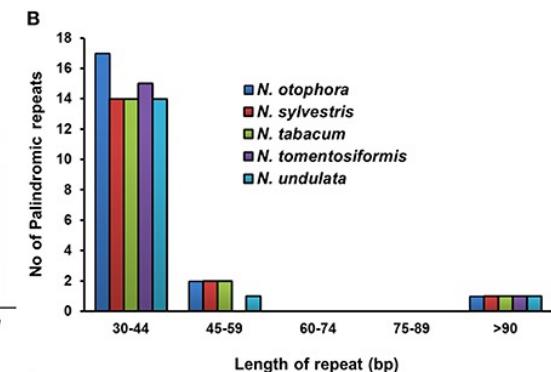
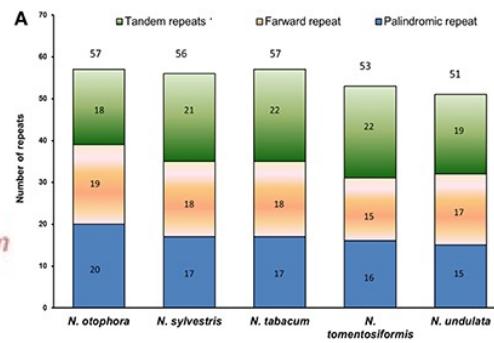
A



B

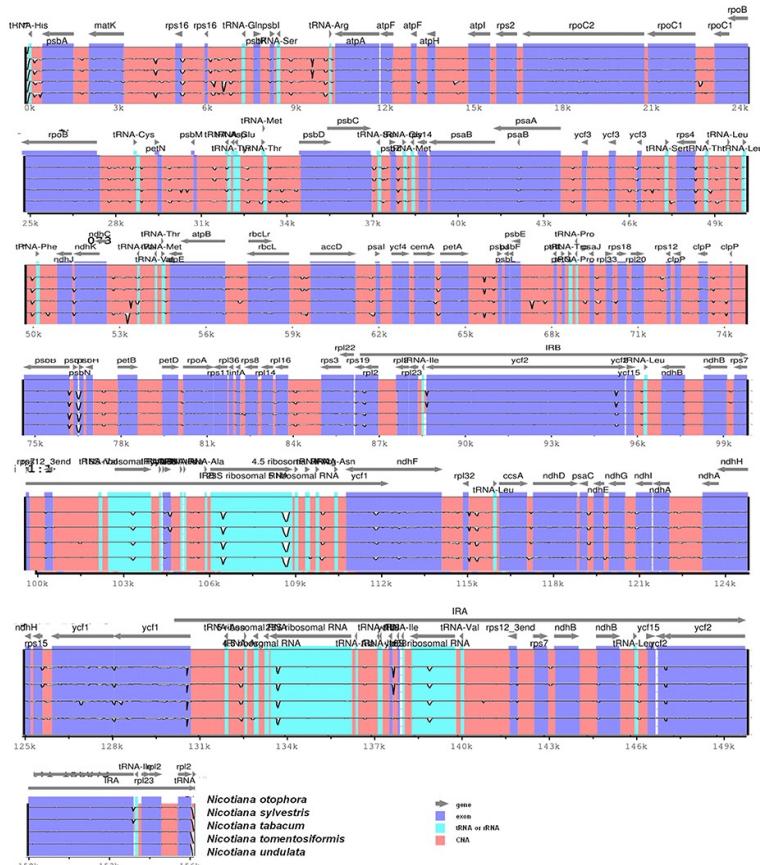


Repeat sequences

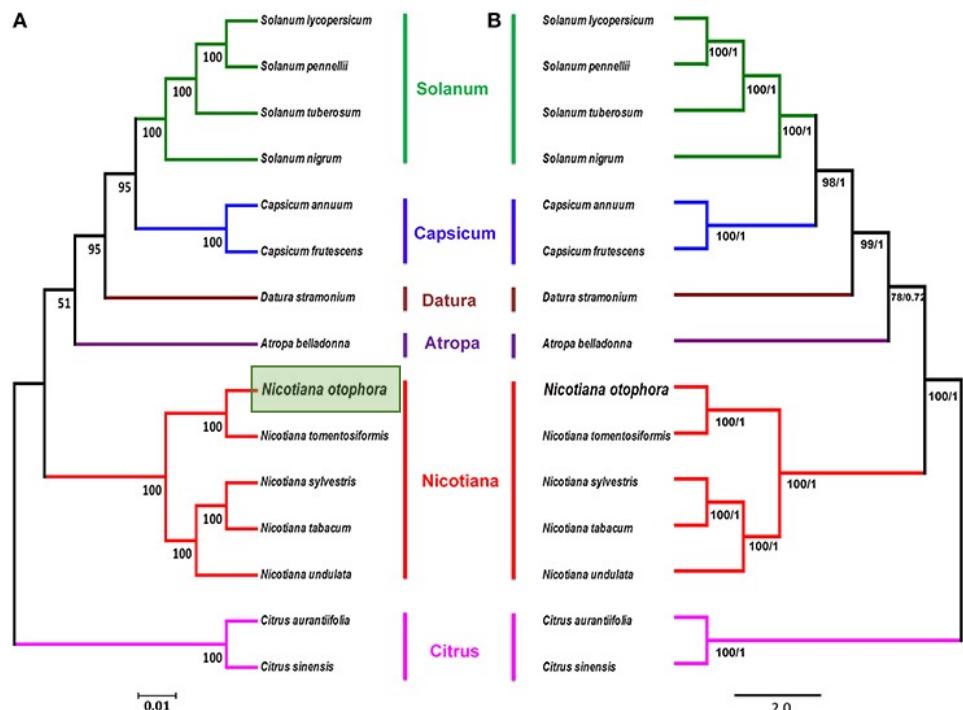


COMPARATIVE GENOMICS

Genomes alignment

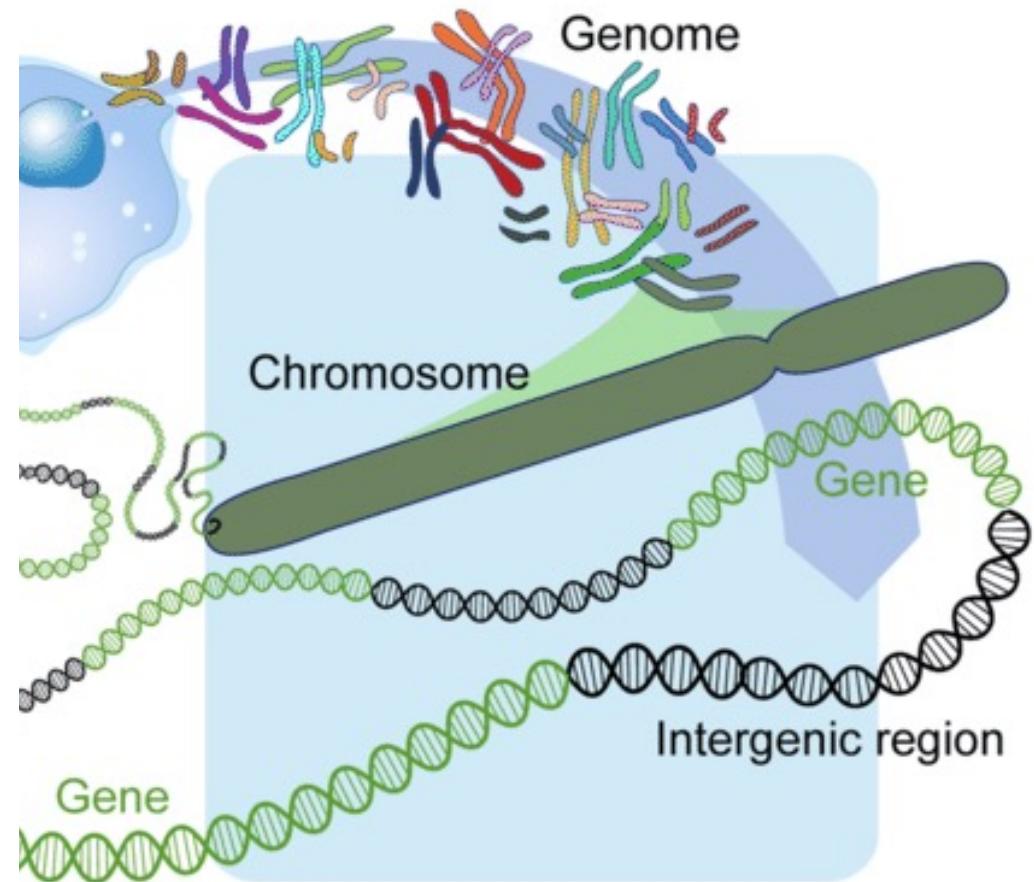


Phylogenetic relationships



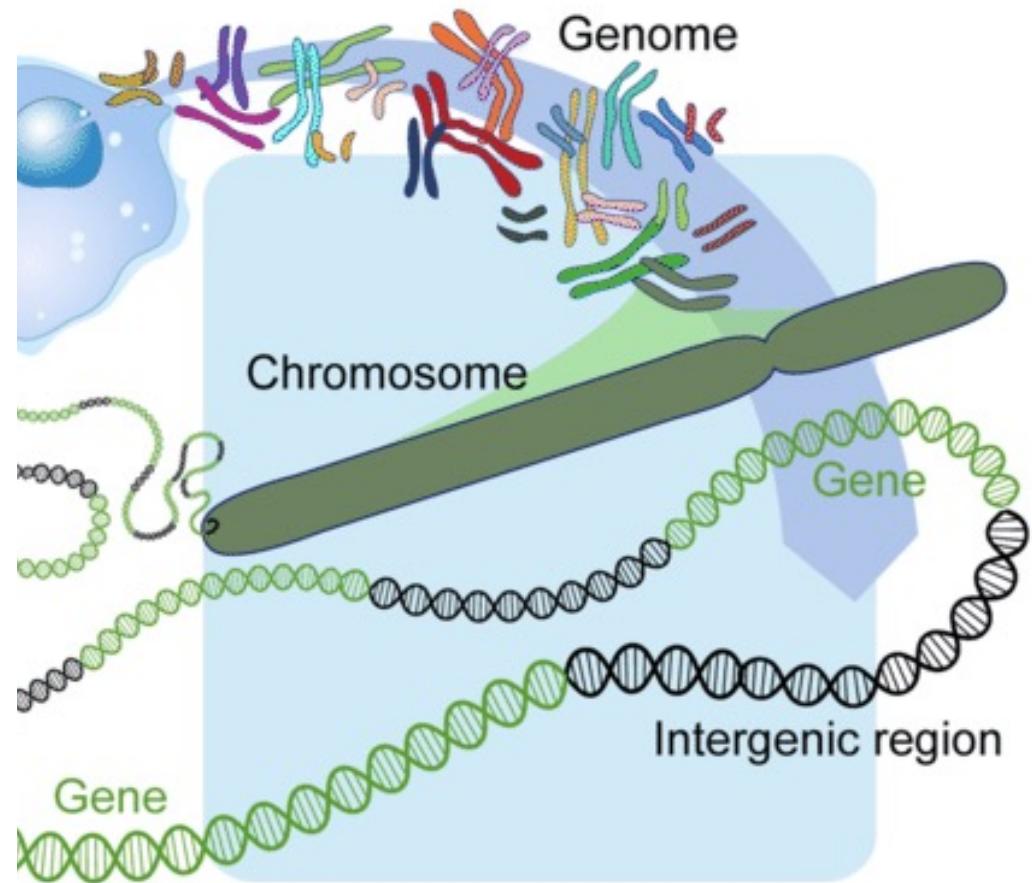
WHAT IS A GENOME?

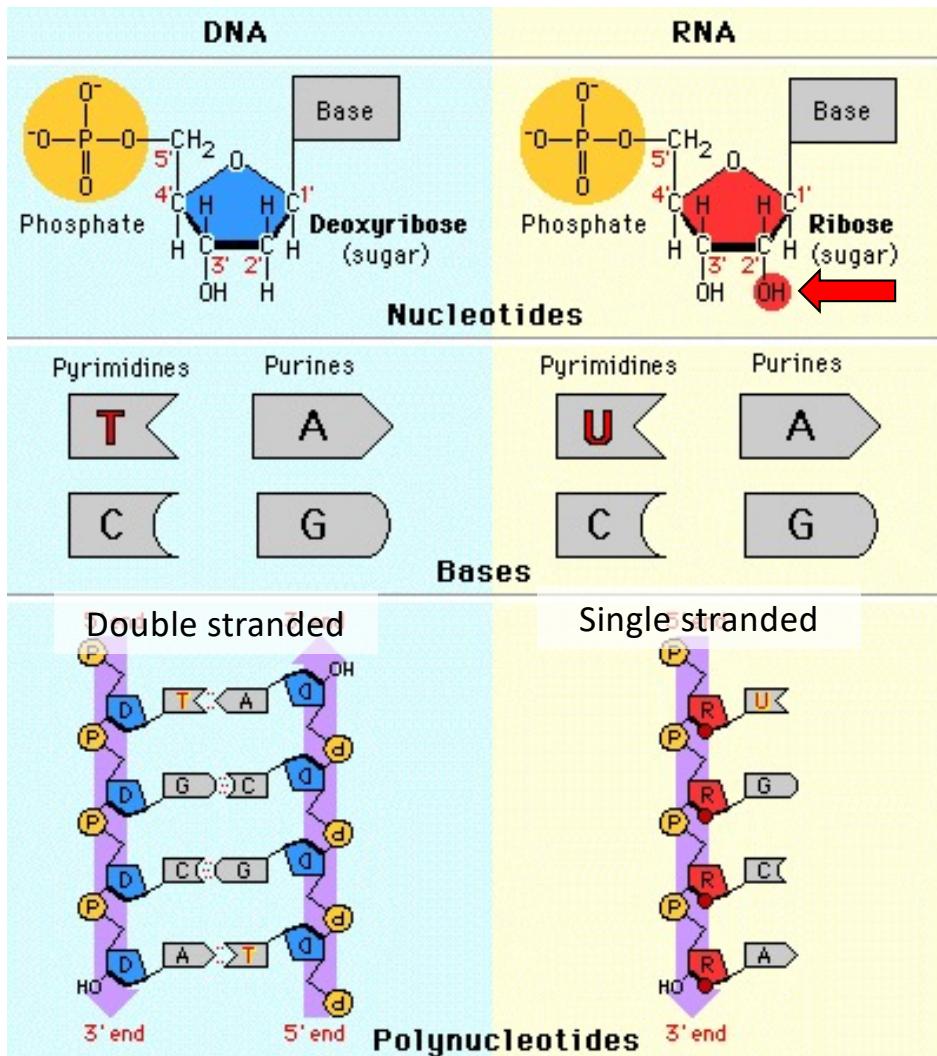
- A **genome** is the entire genetic complement of a living organism.
- Each organism possesses a genome containing the biological information needed to construct and maintain a living example of that organism.



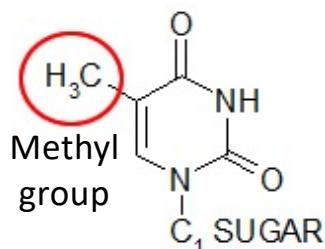
WHAT IS A GENOME?

- Most genomes are made of **DNA** (deoxyribonucleic acid), but few viruses have **RNA** (ribonucleic acid) genomes.
- DNA and RNA are polymeric molecules made up of chains of monomeric subunits called nucleotides.

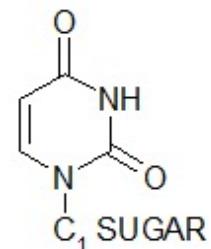




- The **OH bond** in RNA ribose makes molecule more reactive, especially in alkaline conditions.



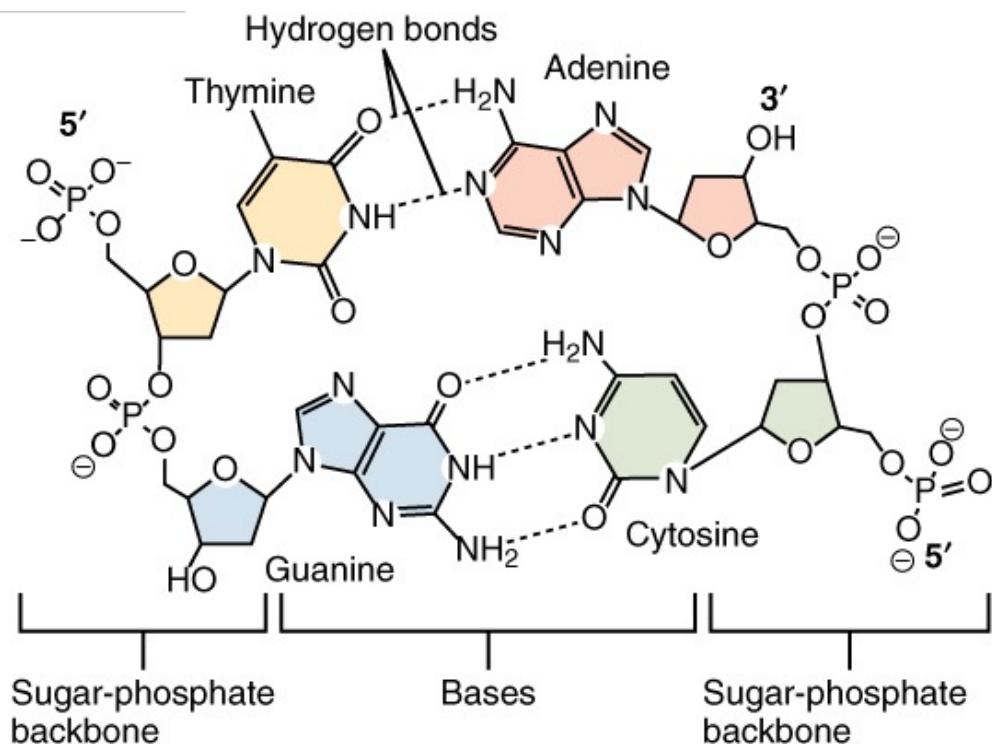
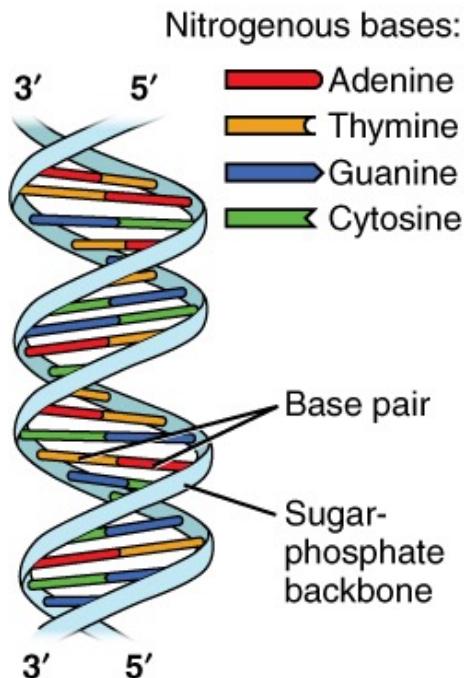
Thymine (DNA)



Uracil (RNA)

- DNA is more resistant to enzymatic attacks (due to its helix structure) compared to RNA.

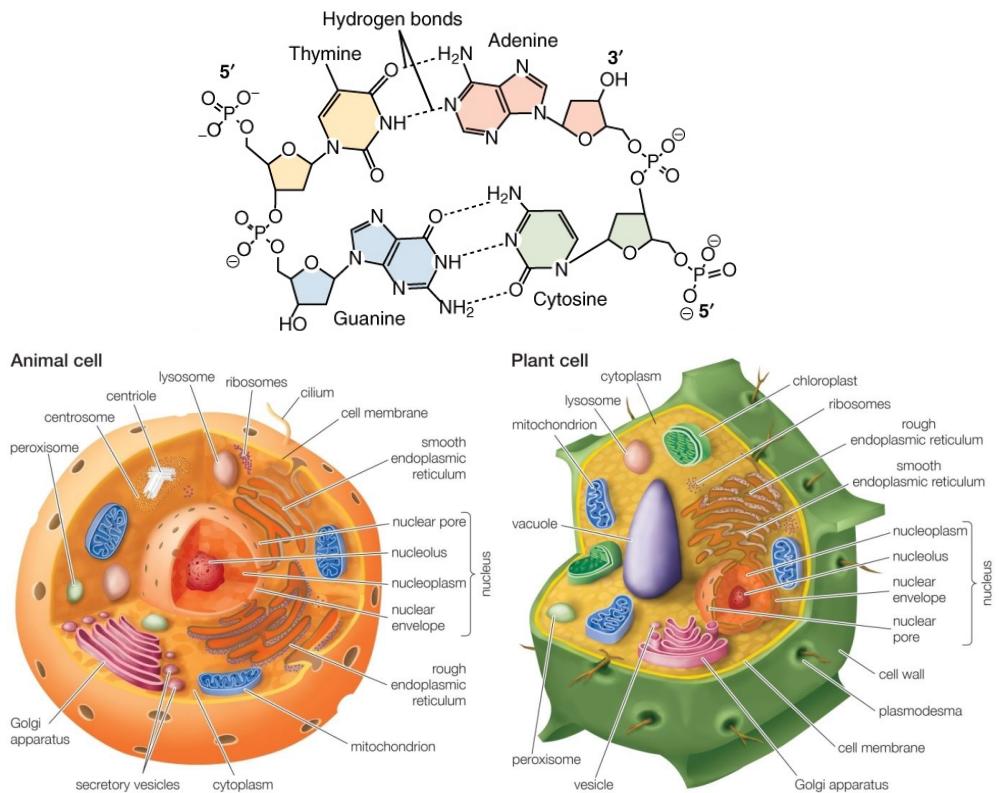
DNA – HYDROGEN BONDS



T-A: 2 hydrogen bonds
G-C: 3 hydrogen bonds

DNA – HYDROGEN BONDS

- Mitochondrial and chloroplastic genomes are enriched in AT.
- Nuclear genome is enriched in GC.
- On average, plastid DNA GC-content is ~37%, whereas nrDNA GC-content is ~41%.
- These genome structural properties can be used to filter reads in bioinformatics pipelines and study gene trafficking between genomes.



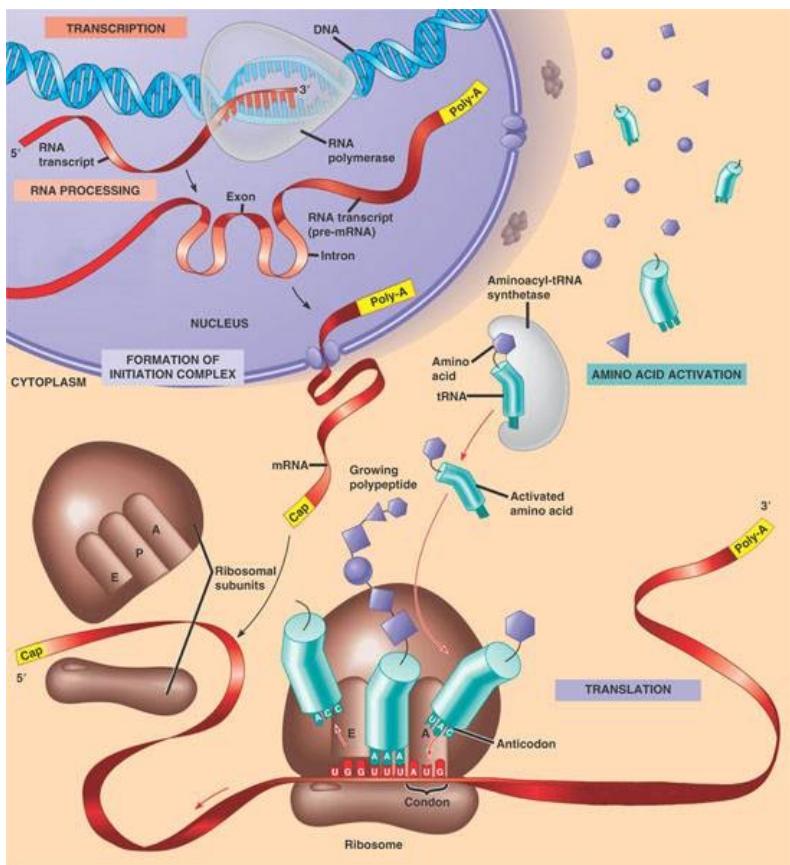
DNA VS. RNA – MODES OF TISSUE PRESERVATION

- The structural and chemical differences between DNA and RNA molecules must be considered to select a storage strategy.
- **What is the best strategy to preserve tissue for genomic/transcriptomic analyses?**
- Which tissue is best suited for the study? When should it be sampled?

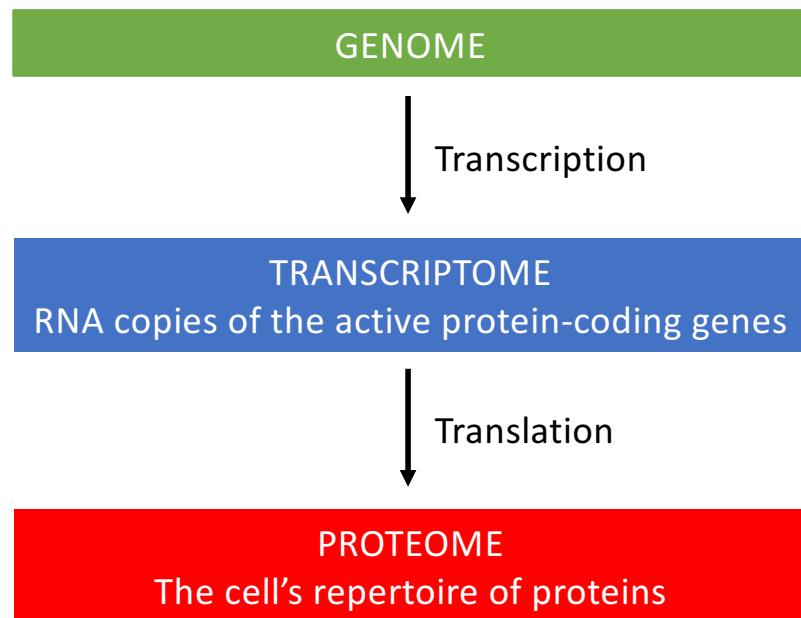


Important for Lab. report

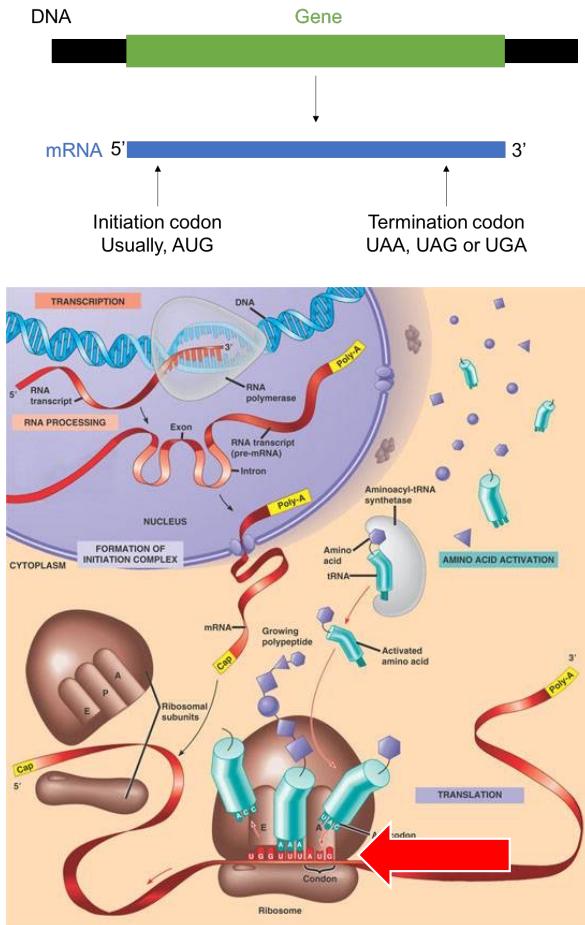
GENOME EXPRESSION



- The genome is a store of biological information, but on its own it is unable to release the information to the cell.
- This process is done through a complex series of biochemical reactions referred to as genome expression.



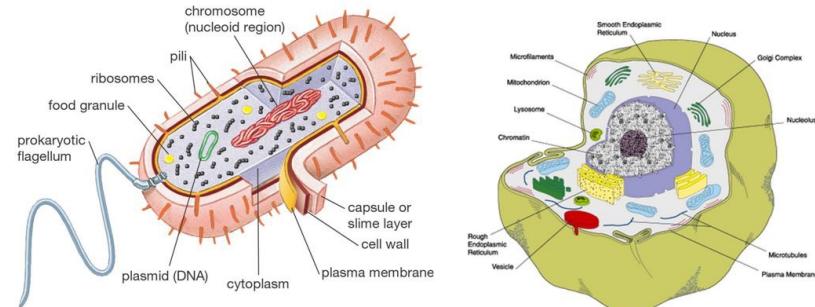
GENETIC CODE – LINKING TRANSCRIPTOME & PROTEOME



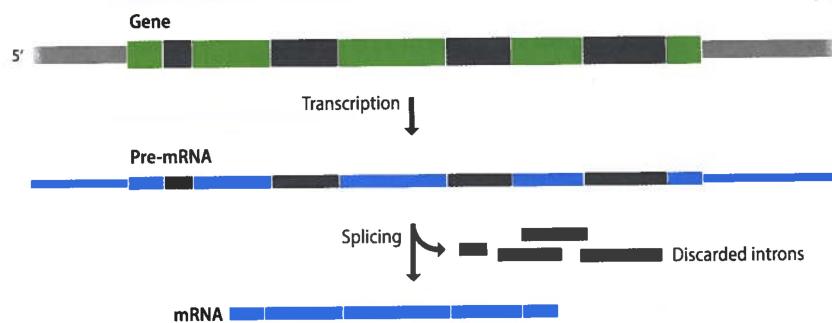
- The genetic code specifies **how the nucleotide sequence of an mRNA is translated into the amino acid (AA) sequence of a protein.**
- Different genetic codes depending on lineages.
- **20 amino acids:** high redundancy (except Met & Trp)
- **Start codon:** Met (sometimes Trp in cpDNA).

UUU	Phe	UCU	Ser	UAU	Tyr	UGU	Cys
UUC		UCC		UAC		UGC	
UUA	Leu	UCA		UAA	Stop	UGA	Stop
UUG		UCG		UAG		UGG	Trp
CUU		CCU		CAU	His	CGU	
CUC	Leu	CCC		CAC		CGC	
CUA		CCA	Pro	CAA	Gln	CGA	
CUG		CCG		CAG		CGG	Arg
AUU		ACU		AAU	Asn	AGU	Ser
AUC	Ile	ACC		AAC		AGC	
AUA		ACA		AAA	Lys	AGA	
AUG	Met	ACG	Thr	AAG		AGG	Arg
GUU		GCU		GAU	Asp	GGU	
GUC		GCC		GAC		GGC	
GUA	Val	GCA	Ala	GAA	Glu	GGC	Gly
GUG		GCG		GAG		GGG	

VARIETIES OF GENOME ORGANIZATION



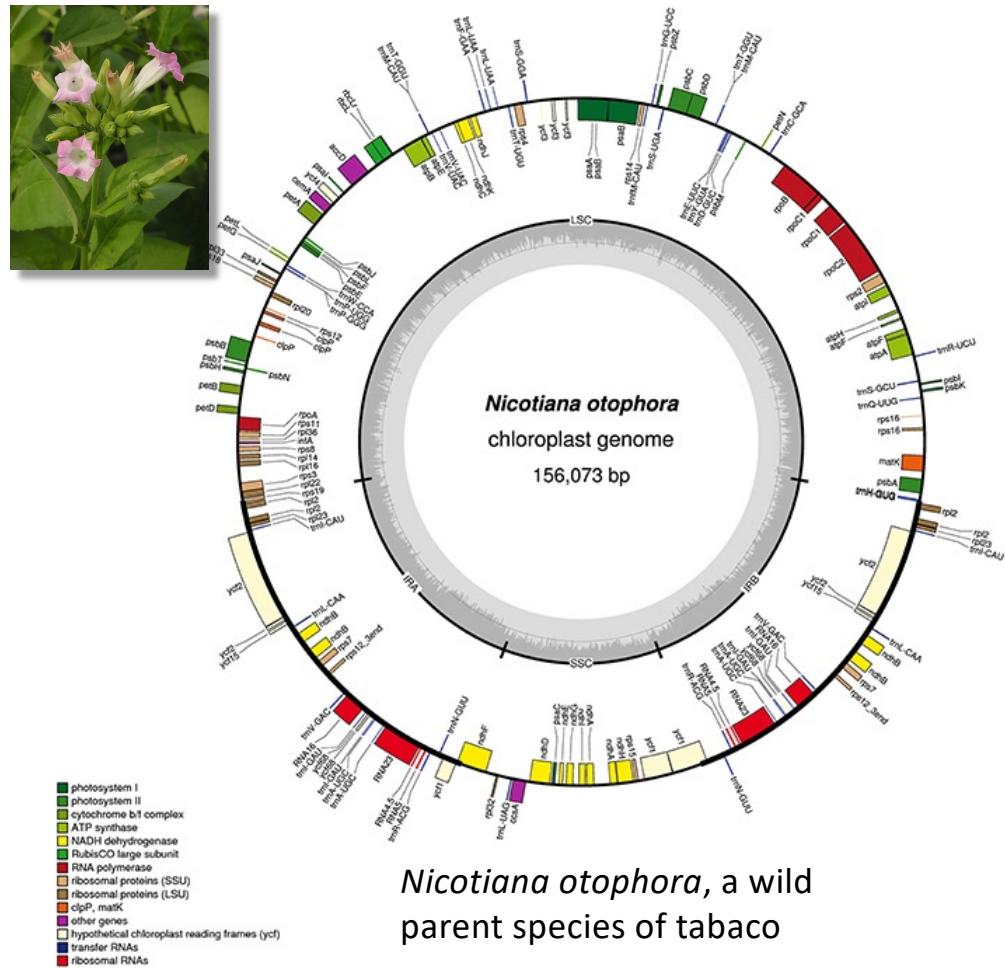
Gene annotation is “easier” in prokaryotes due to the lack of introns



	Prokaryotes	Eukaryotes
DNA	DNA is naked DNA is circular Usually no introns	DNA bound to protein DNA is linear Usually has introns
	No nucleus No membrane-bound	Has a nucleus Membrane-bound
	70S ribosomes	80S ribosomes
	Binary fission	Mitosis and meiosis
	Single chromosome (haploid)	Chromosomes paired (diploid or more)
	Smaller (~1–5 µm)	Larger (~10–100 µm)



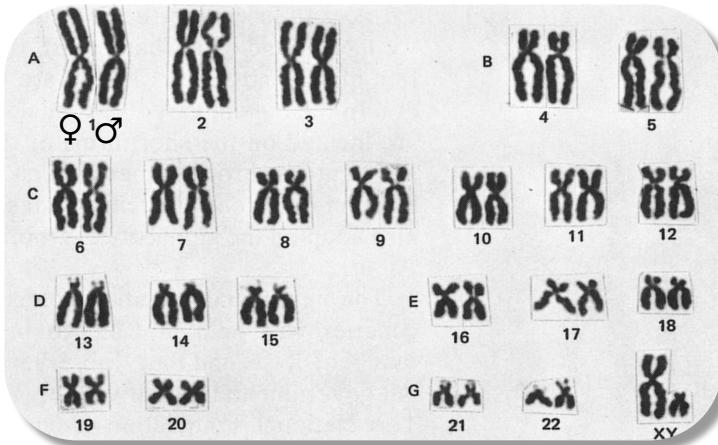
STRUCTURE AND FUNCTION OF CHLOROPLAST GENOME



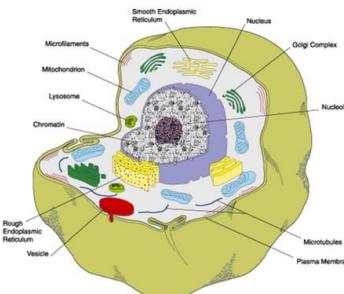
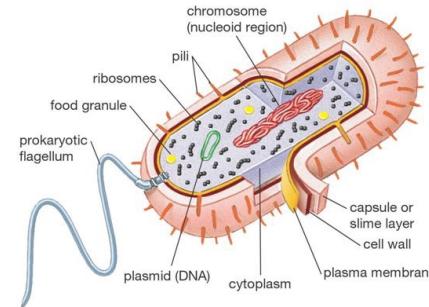
Although chloroplast genomes are of prokaryote origin several genes have introns!

VARIETIES OF GENOME ORGANIZATION

Human: $2n=2x=46$



This difference impacts on genome complexity (heterozygosity) and influence bioinformatics (chromosome reconstruction and phasing)

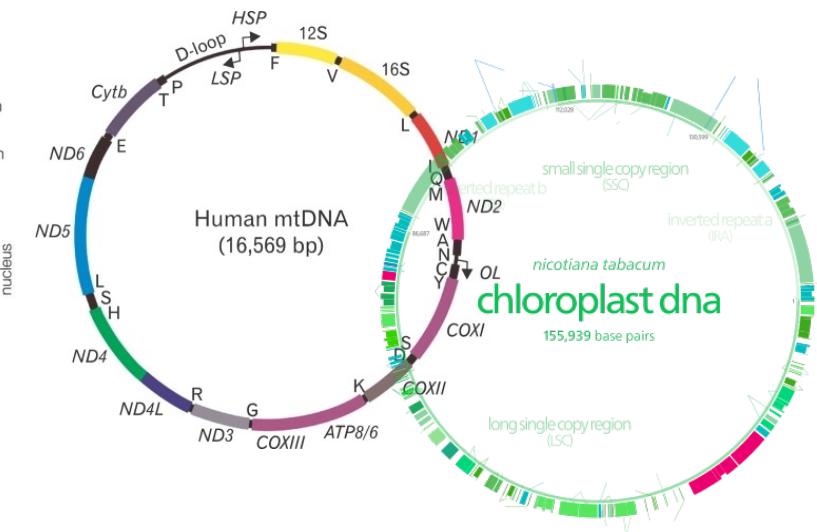
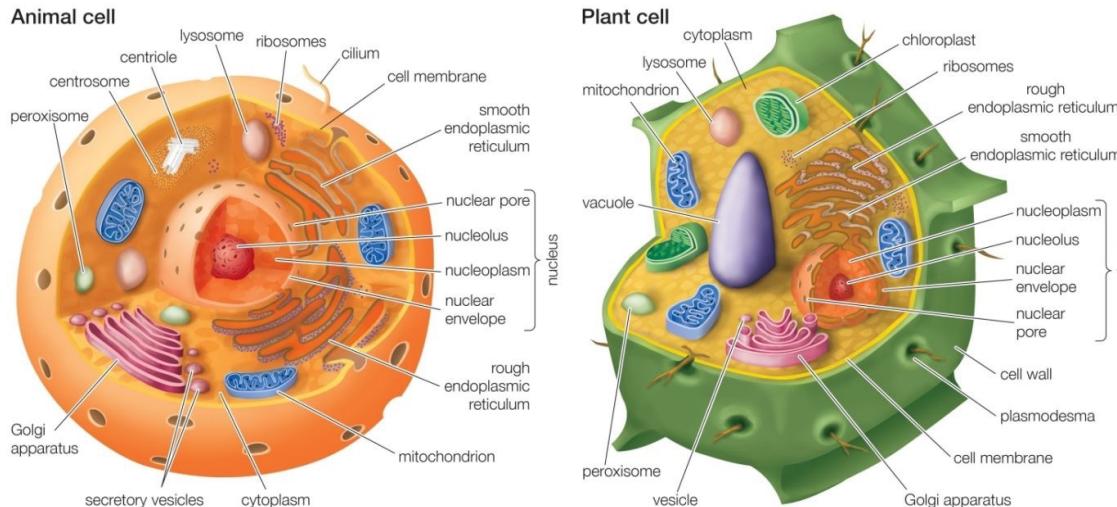


	Prokaryotes	Eukaryotes
DNA	DNA is naked	DNA bound to protein
	DNA is circular	DNA is linear
	Usually no introns	Usually has introns
Organelles	No nucleus	Has a nucleus
	No membrane-bound	Membrane-bound
	70S ribosomes	80S ribosomes
Reproduction	Binary fission	Mitosis and meiosis
	Single chromosome (haploid)	Chromosomes paired (diploid or more)
Average Size	Smaller (~1–5 µm)	Larger (~10–100 µm)



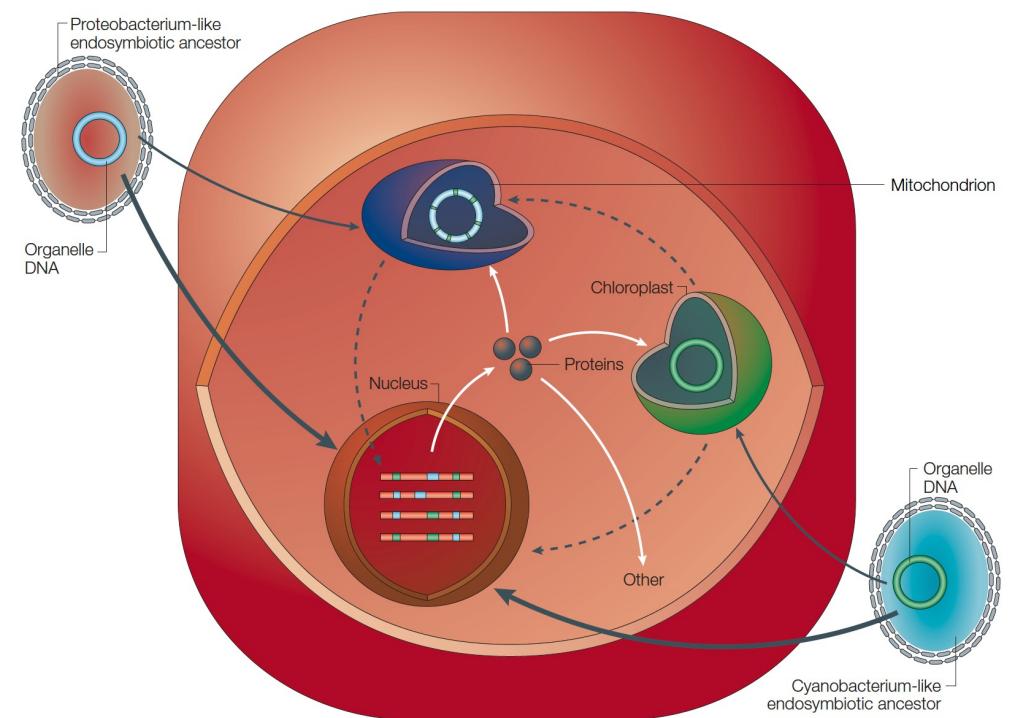
VARIETIES OF GENOME ORGANIZATION

- Eukaryotic cells also contain organelles with their own genomes.
 - These organelles contain additional DNA usually in the form of single closed or circular molecules; un-complexed with histones, like the DNA of prokaryotes, but some genes have introns.



ORGANELLE GENOMES FORGE EUKARYOTIC CHROMOSOMES

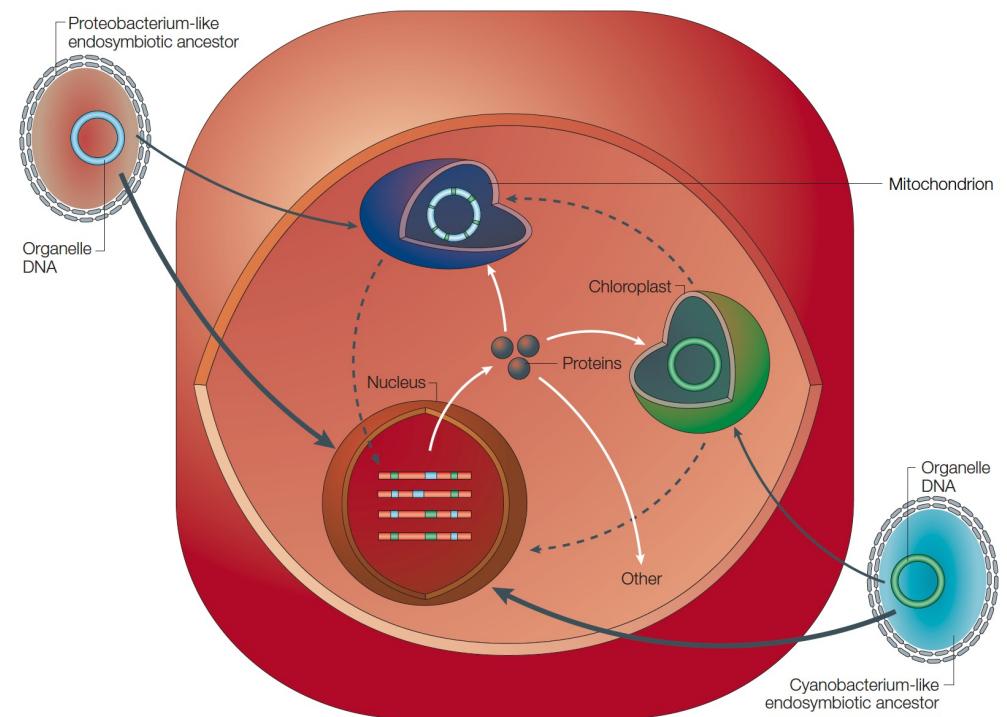
- Mitochondria and plastids were once free-living prokaryotes.
- **Most of the genes** originally present in the ancestors' genomes **have been transferred to the nucleus**, with few being retained in the organelles.
- These organelles heavily dependent on nuclear genes and import >90% of their proteins from the cytoplasm (e.g. RuBisCO).
- **Endosymbiotic gene transfer is not homogenous across lineages therefore complicating genome assembly.**



Timmis et al. (2004) *Nature Reviews*

ORGANELLE GENOMES FORGE EUKARYOTIC CHROMOSOMES

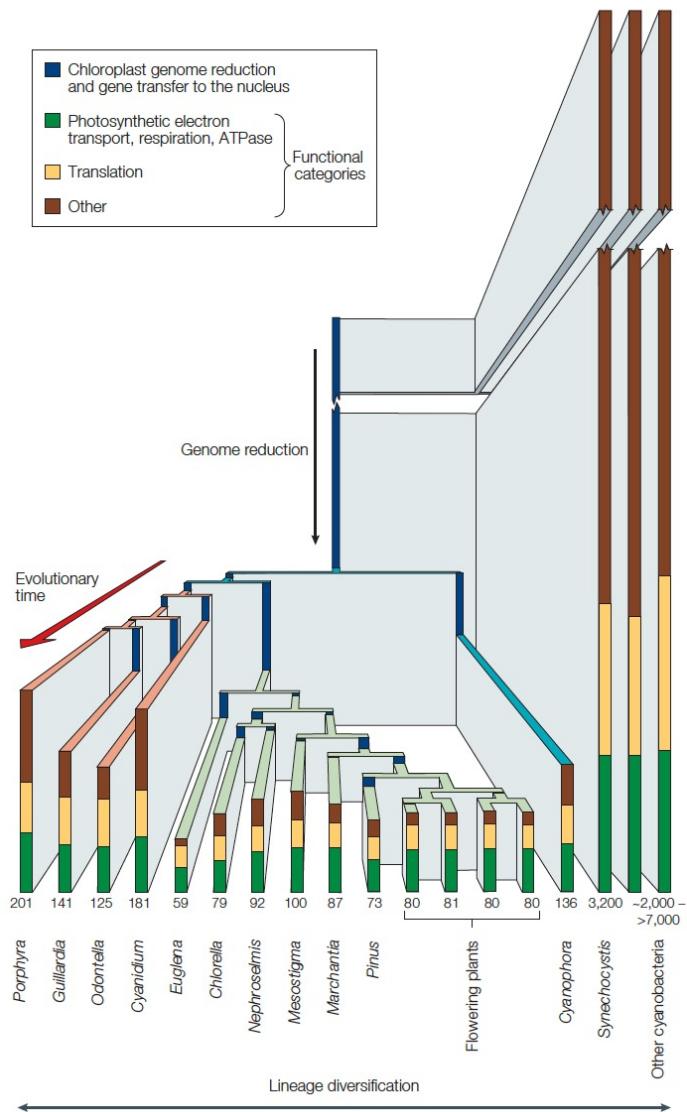
- Gene transfers into the nucleus is tricky to assess when conducting genome assemblies.
- **Tendency for genes to increase their GC-contents when they are transferred from organelles to the nucleus.**



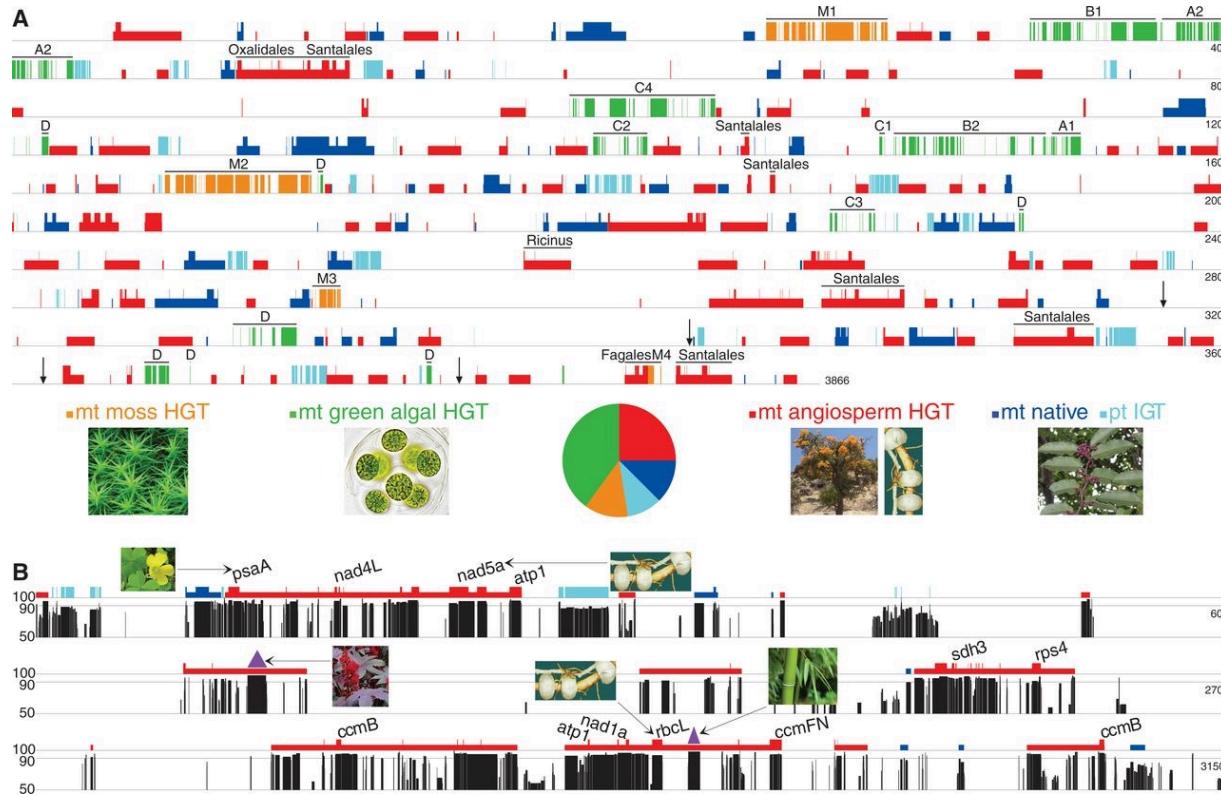
Timmis et al. (2004) *Nature Reviews*

REDUCTION OF SIZE OF CHLOROPLAST GENOMES OVER TIME

- Ancestor of plastids was a free-living cyanobacterium and therefore must have possessed several thousand genes as did its contemporaries.
- **Plastids have relinquished most of their genes to the genome of their host cell.**
- This gene relocation process occurred massively at the onset of endosymbiosis and continued in parallel during algal diversification.
- We are still at the infancy of unraveling this process, but NGS provides a boost.



HORIZONTAL GENE TRANSFERS IN PLANTS' MITOCHONDRIAL GENOMES

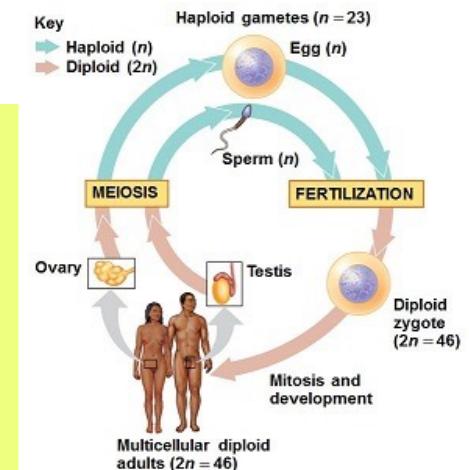
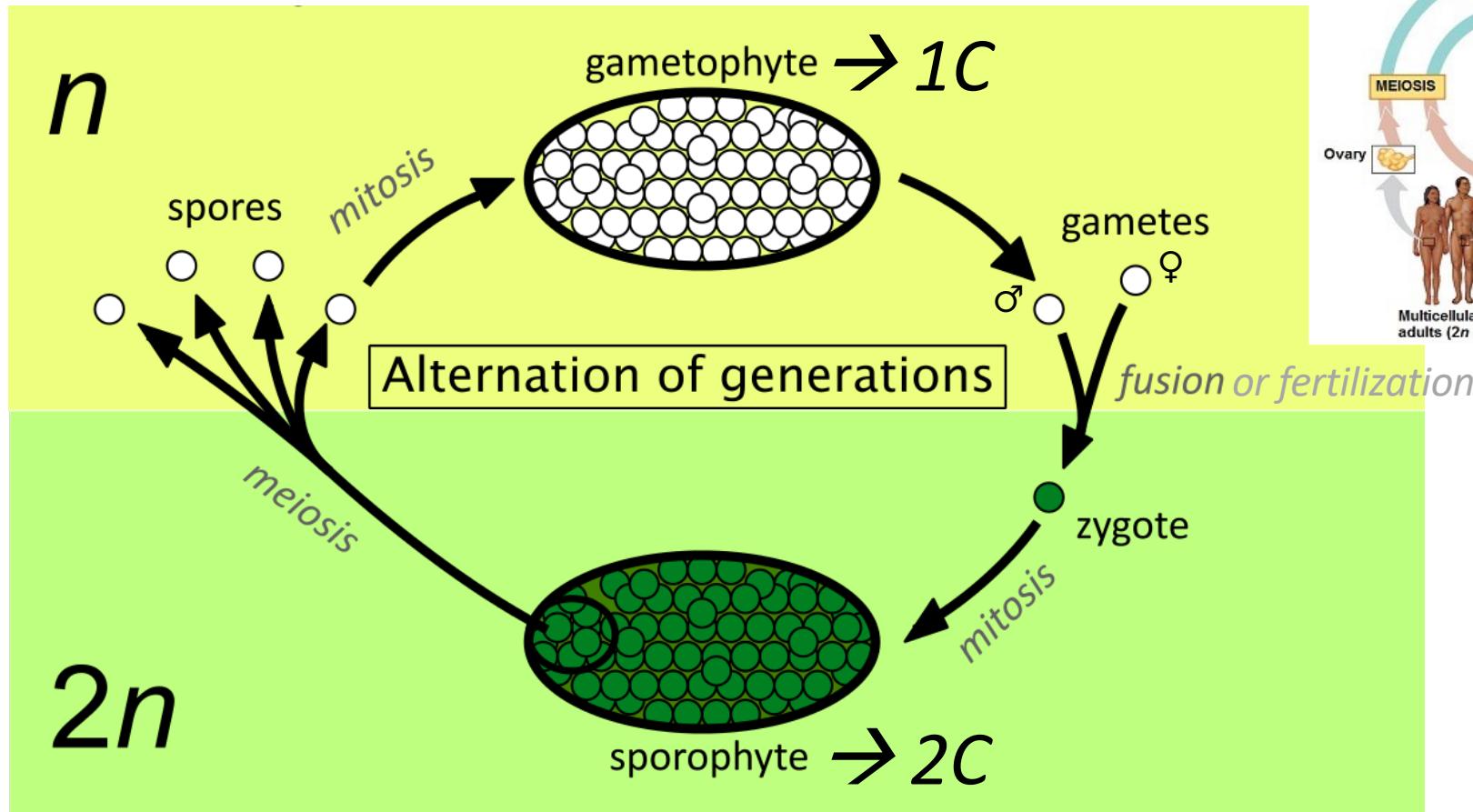


- Plant mitochondrial genomes (mtDNA) are not as conserved as those of animals.
- In plants, mtDNA genomes greatly vary in sizes and number of genes due to horizontal gene transfers (e.g. entire mtDNA genomes in *Amborella*)



Rice et al. (2013) *Science*

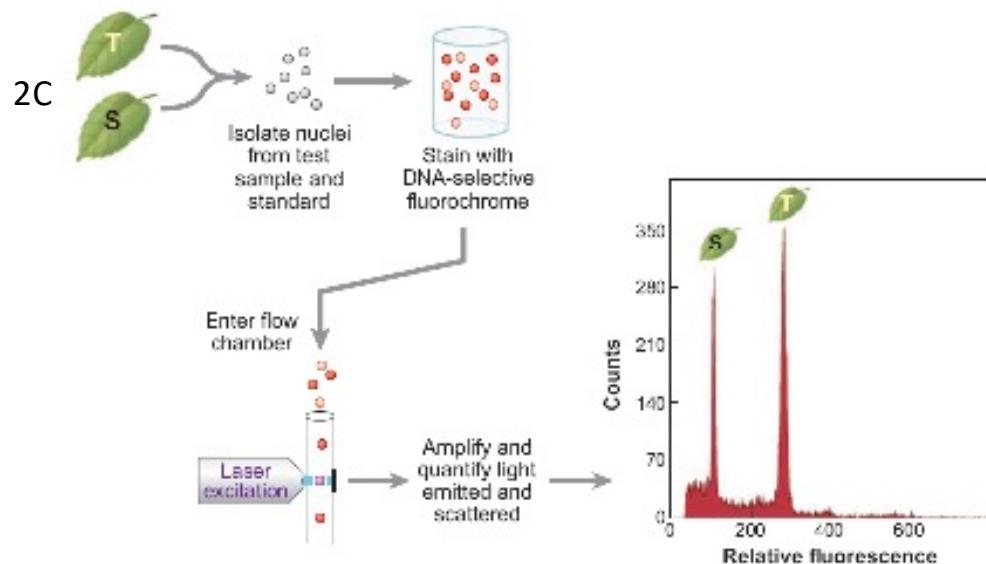
QUANTIFYING GENOME SIZE – LIFE CYCLE



1C = DNA content of the gametophytic (n) set of chromosomes

QUANTIFYING GENOME SIZE – FLOW CYTOMETRY

- Flow cytometry is the best strategy to estimate genome size (and ploidy).
- It is key to design sequencing strategy (e.g. how many sequence data should be produced to assemble genome).



$$S \text{ } 2C \text{ value (pg)} = T \text{ } 2C \text{ value} \times S \text{ } 2C \text{ mean peak position} / T \text{ } 2C \text{ mean peak position}$$

↓

$$S \text{ Nb bp} = \text{mass in pg} \times 0.978 \times 10^9$$

or

$$1 \text{ pg} = 978 \text{ Mbp}$$

C= DNA content of the haploid set of chromosomes

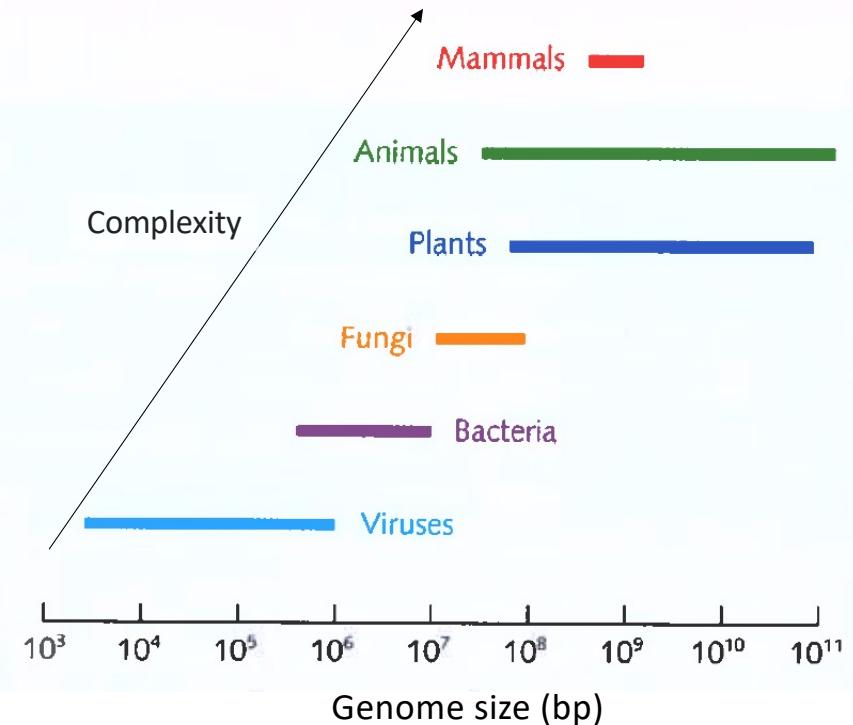
S= Sample

T= Reference (= sample of known genome size used as standard to estimate size of sample)

Doležel *et al.* (2007) *Nature Protocols*

GENOME SIZES & ORGANISMS' COMPLEXITY

- Large variation in genome sizes across the Tree of Life!
- **Broad scale: correlation between genome size (= amount of DNA per cell) and the complexity of an organism.**
- Fine scale: this correlation does not necessarily hold true within closely related species (especially in plants with processes of whole genome doubling).
- The number of genes is usually used as a proxy to reflect an organism's complexity.



GENOME SIZES & ORGANISMS' COMPLEXITY

Species	Genome size (Mb)	Coding (%)	Approx. number of genes	Gene density (kb/gene)
<i>Escherichia coli</i>	4.64	88	4485	1.03
Yeast	12.5	70	6000	2.1
Pufferfish	365	15	23000	10
<i>Arabidopsis thaliana</i>	115	29	23000	6
Human	3289	1.3	23000	143



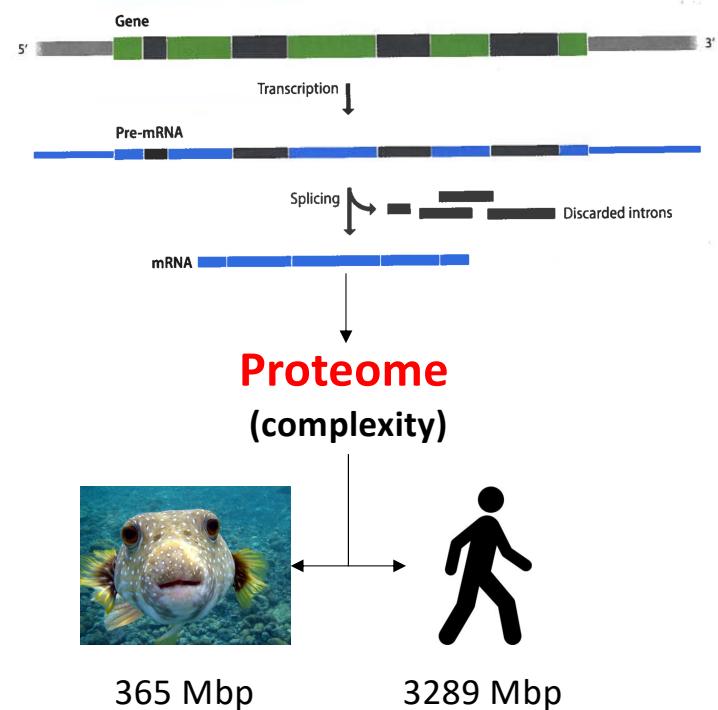
Pufferfish and human exhibit different levels of complexity, but they share the same number of genes

→ Where does complexity take place?



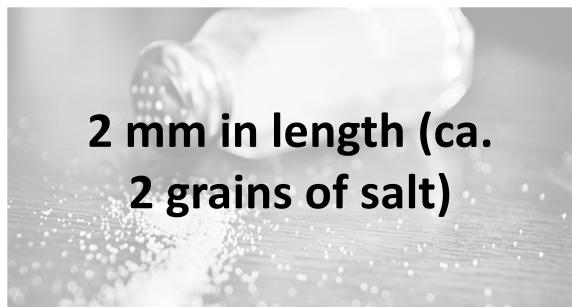
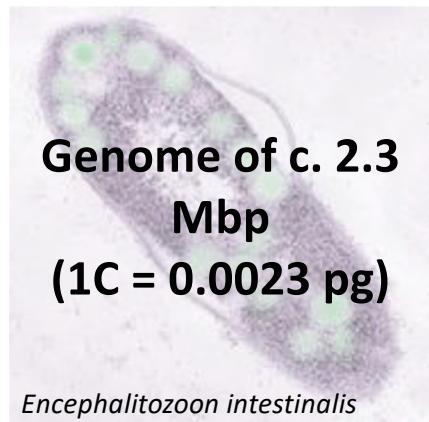
GENOME SIZES & ORGANISMS' COMPLEXITY

- Number of genes is not a good proxy of an organism's "complexity".
- **Complexity takes place at the **proteome** level, but it is orchestrated at the **transcriptome** level** via two main processes:
 - ✓ Alternative splicing.
 - ✓ RNA editing.
- **Alternative splicing is a process during gene expression that results in a single gene coding for multiple proteins.** In this process, particular exons of a gene may be included within or excluded from the final mRNA. For instance, in mammals, billions of antibodies arise from <100 exons.

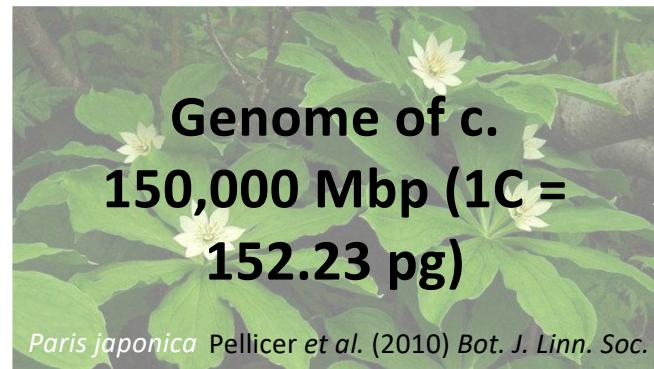


VARIATION OF GENOME SIZES – EUKAROYOTIC GENOMES

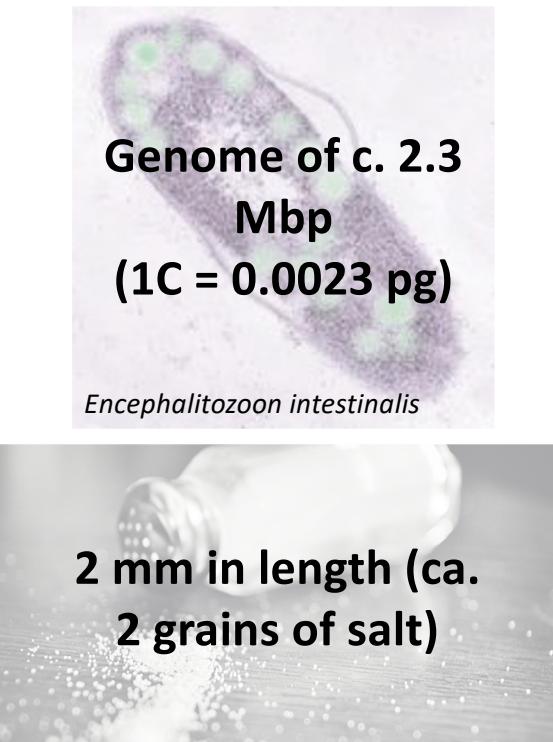
**66,000-fold
variation**



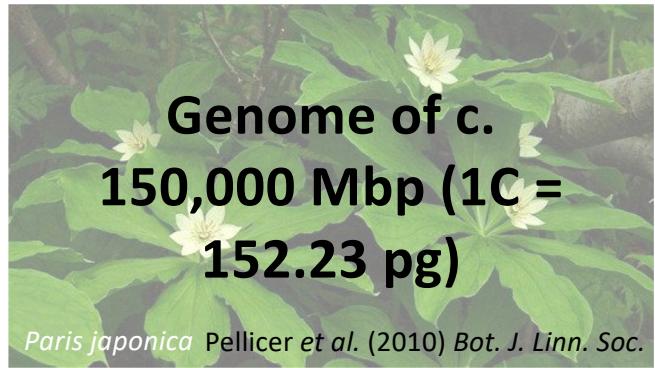
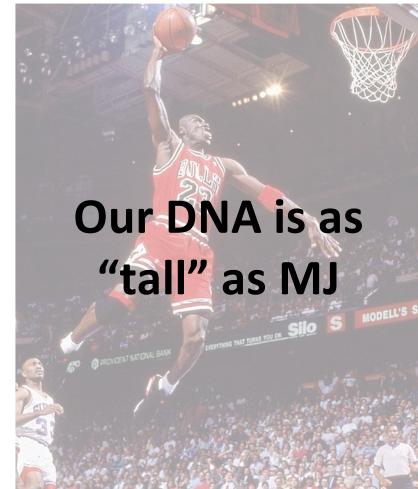
Mbp = mega base pairs = 1,000,000 bp
1 nucleotide = c. 0.34 nm



VARIATION OF GENOME SIZES – EUKARYOTIC GENOMES



**66,000-fold
variation**

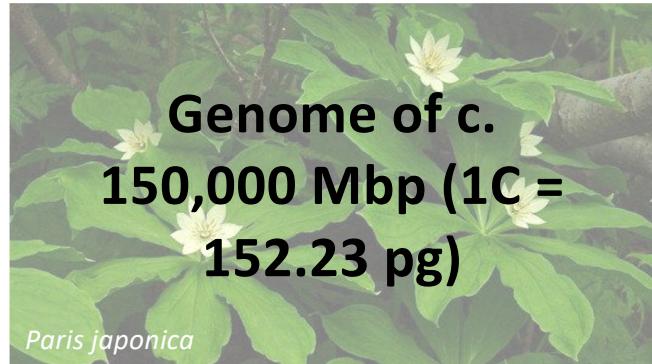


96 m in length (Big Ben)



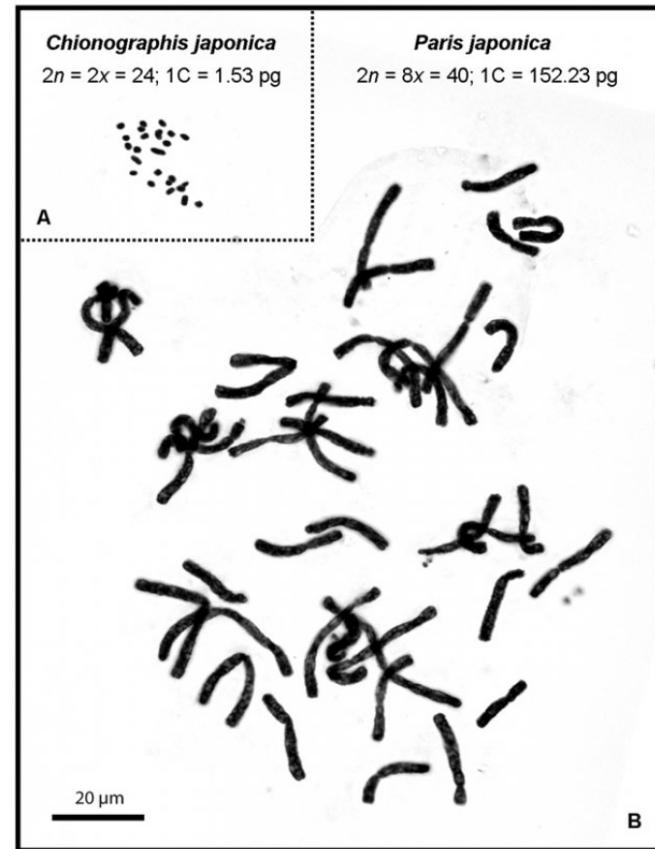
Mbp = mega base pairs = 1,000,000 bp
1 nucleotide = c. 0.34 nm

VARIATION OF GENOME SIZES – EUKAROYOTIC GENOMES

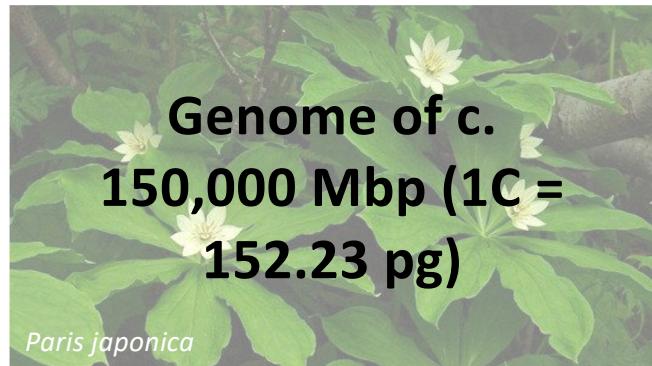


The “obese” genome of *P. japonica* is associated with:

1. Polyploidization event(s) or whole genome doubling (8x).
2. Extreme karyological rearrangements (> 10-fold difference in chromosome length with other species in family).

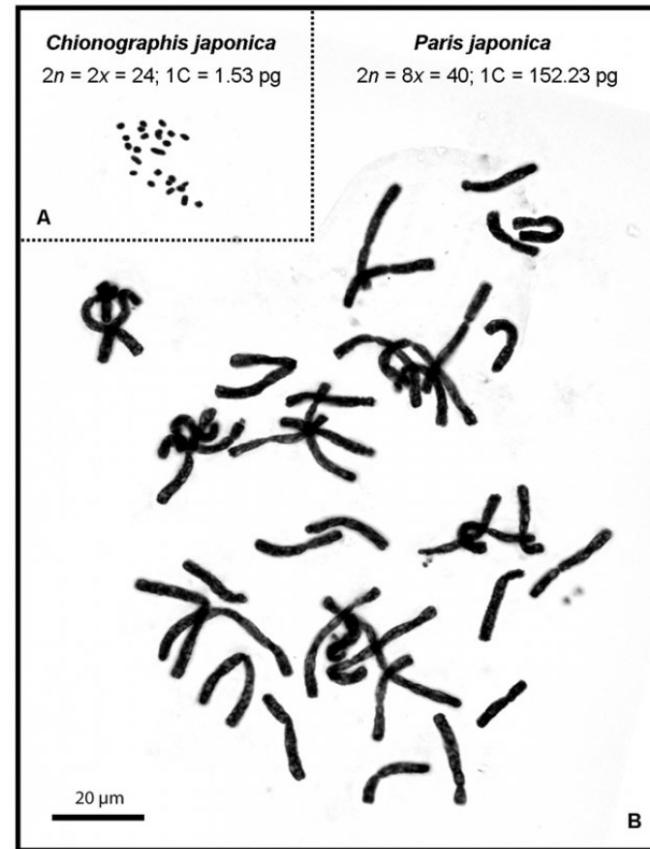


VARIATION OF GENOME SIZES – EUKAROYOTIC GENOMES



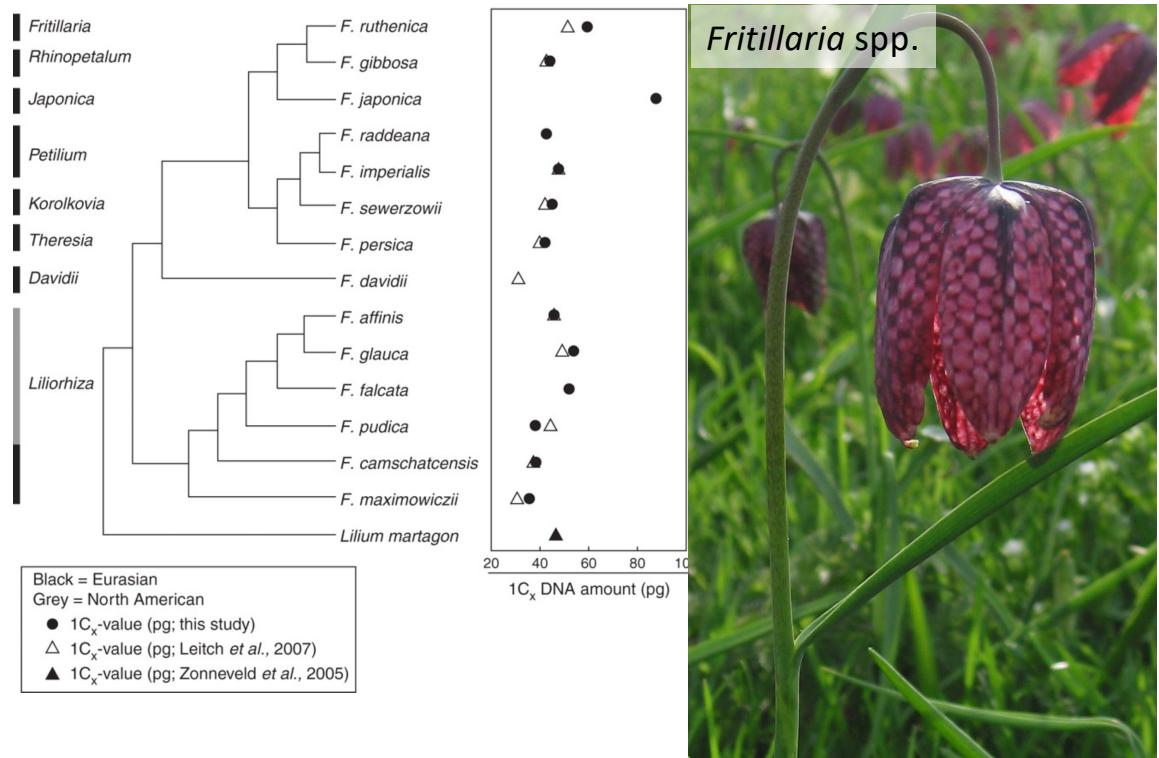
Basic knowledge on chromosome numbers and genome sizes are essential to design any genomic project, especially to develop a sequencing strategy.

Rule of thumb: Aim for sequencing depth between 100-150x to assemble an eukaryote nuclear genome (but this depends on level of heterozygosity).



VARIATION OF GENOME SIZES – EUKAROYOTIC GENOMES

- Large variation of genome sizes in *Fritillaria*, most likely connected to transposable elements.
- **Phylogenetic evidence can't predict species genome sizes.**
- First step in designing a genomic project is therefore to estimate the organism's genome size and ploidy level.

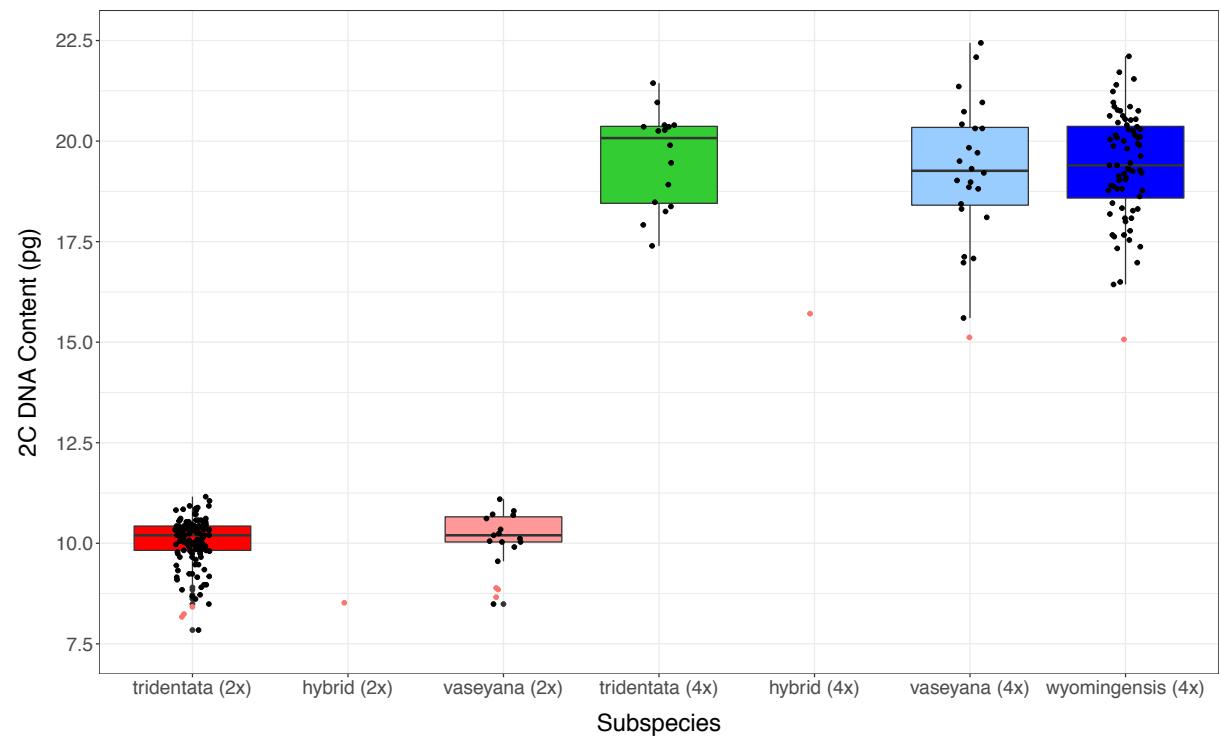


Ambrozova *et al.* (2011) *Ann. Bot.*

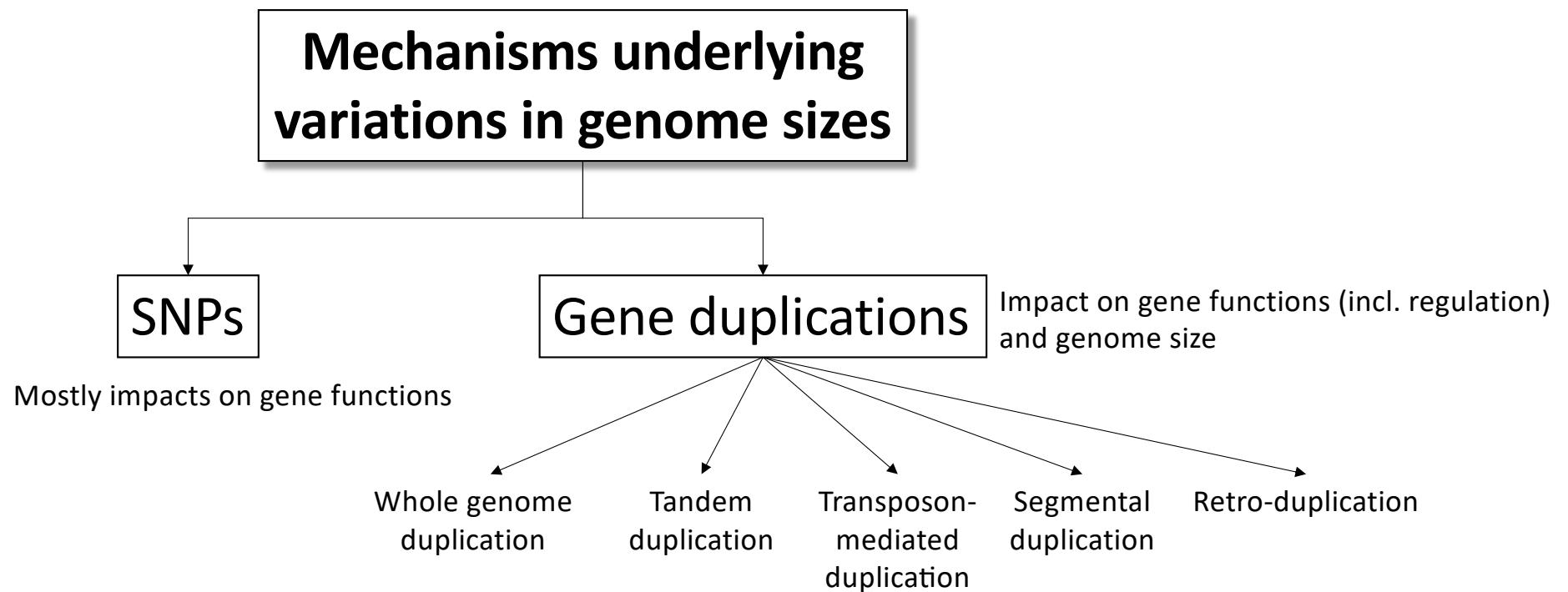
VARIATION OF GENOME SIZES – EUKAROYOTIC GENOMES



Sagebrush genome is 3x bigger than human genome (ca. 9Gbp)



HOW DO GENOMES DIFFER?



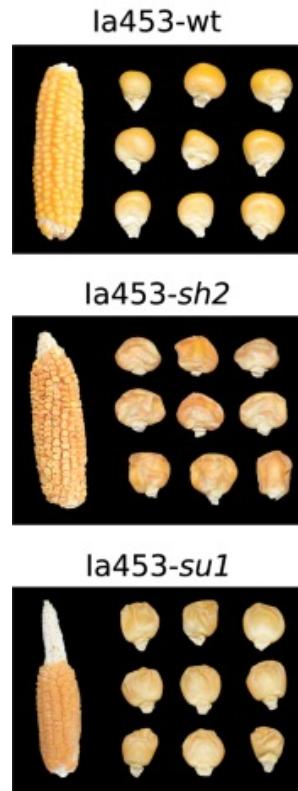
HOW DO GENOMES DIFFER?

Single nucleotide polymorphisms (SNPs)

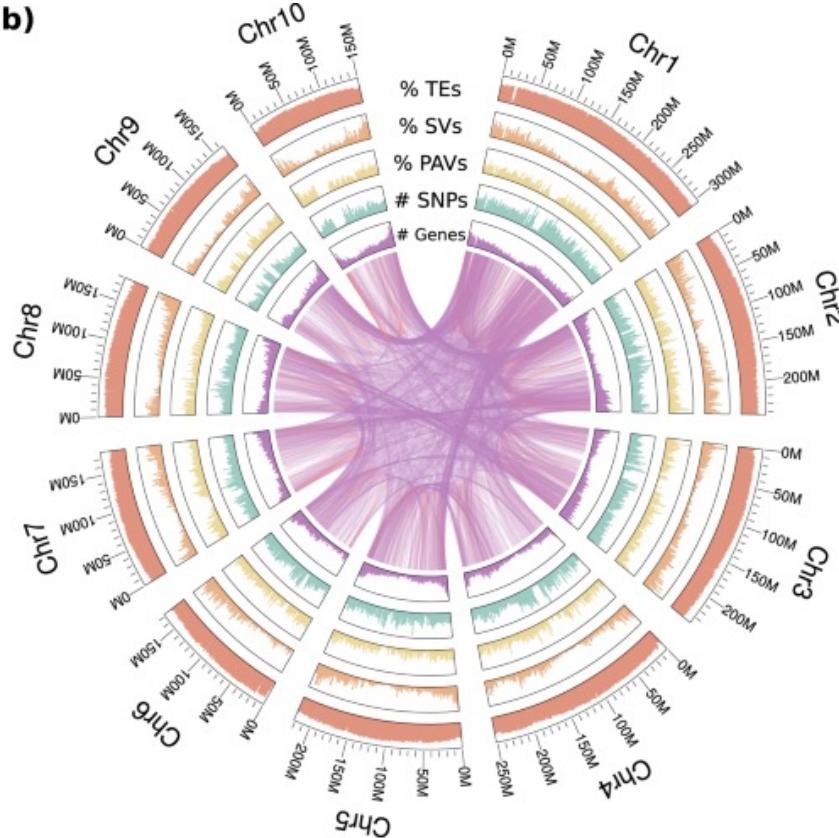
- The most common type of genetic variation.
- Each SNP represents a difference in a single nucleotide.
 - ✓ For instance, a C is replaced by a T.
- In the human genome, SNPs occur on average once in every 300 nucleotides.
 - ✓ There are roughly 10 million SNPs in our genome.

EXAMPLE OF SNPs AND GENOME FEATURES IN CORN

a)



b)

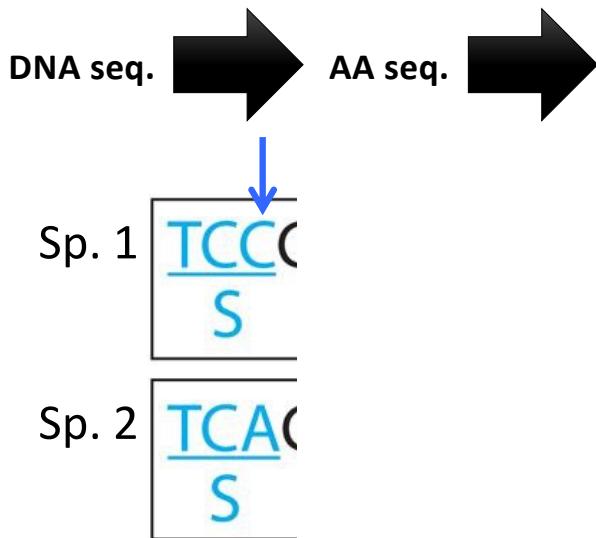


<https://www.nature.com/articles/s41467-021-21380-4>

HOW DO GENOMES DIFFER?

Impact of SNPs on coding genes

- Synonymous vs. non-synonymous substitutions



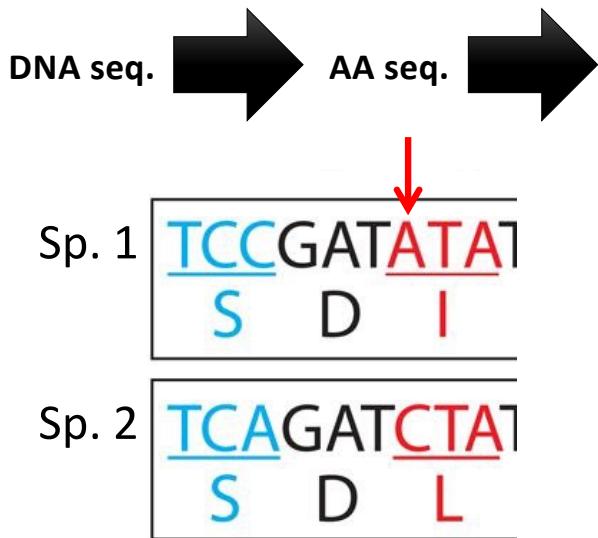
Synonymous substitution – mutation that does not alter amino acid seq.

		Second Letter					
		U	C	A	G		
1st letter	U	UUU Phe UUC UUA Leu UUG	UCU Ser UCC UCA UCG	UAU Tyr UAC UAA Stop UAG	UGU Cys UGC UGA Stop UGG	U C C G A G	
	C	CUU Leu CUC CUA CUG	CCU Pro CCC CCA CCG	CAU His CAC CAA CAG	CGU Arg CGC CGA CGG	U C C A A G	
	A	AUU Ile AUC AUA Met AUG	ACU Thr ACC ACA ACG	AAU Asn AAC AAA AAG	AGU Ser AGC AGA AGG	U C C A A G	
	G	GUU Val GUC GUA GUG	GCU Ala GCC GCA GCG	GAU Asp GAC GAA GAG	GGU Gly GGC GGA GGG	U C C A A G	

HOW DO GENOMES DIFFER?

Impact of SNPs on coding genes

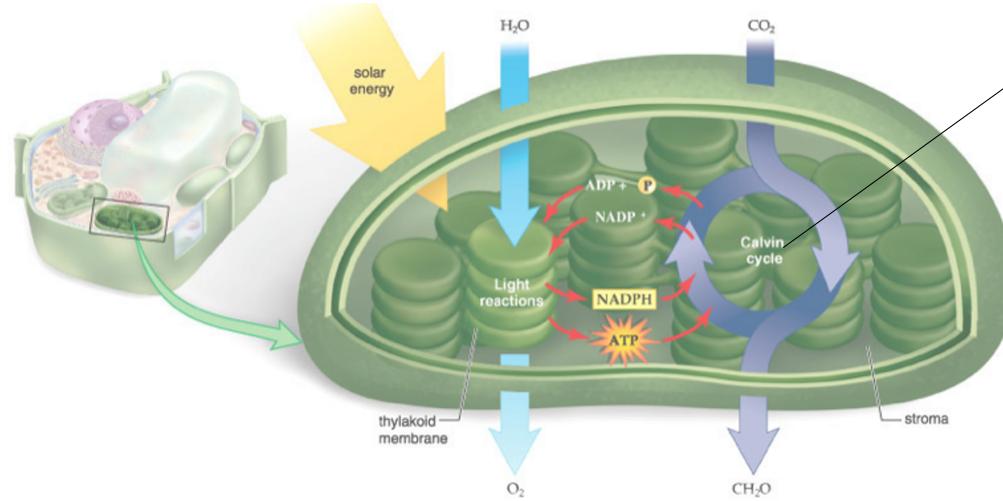
- Synonymous vs. non-synonymous substitutions



Non-synonymous substitution – mutation altering the amino acid seq.

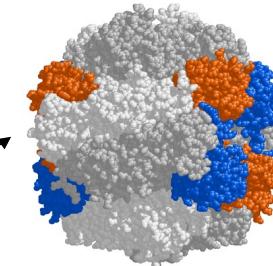
		Second Letter					
		U	C	A	G		
1st letter	U	UUU Phe UUC UUA Leu UUG	UCU Ser UCC UCA UCG	UAU Tyr UAC UAA Stop UAG	UGU Cys UGC UGA Stop UGG	U C C G A G	
	C	CUU Leu CUC CUA Leu CUG	CCU Pro CCC CCA CCG	CAU His CAC CAA Gln CAG	CGU Arg CGC CGA CGG	U C C A A G	
	A	AUU Ile AUC AUA Met AUG	ACU Thr ACC ACA ACG	AAU Asn AAC AAA Lys AAG	AGU Ser AGC AGA AGG	U C C A A G	
	G	GUU Val GUC GUA GUG	GCU Ala GCC GCA GCG	GAU Asp GAC GAA Glu GAG	GGU Gly GGC GGA GGG	U C C A A G	

EVOLUTION OF CARBON FIXATION IN FACE OF CLIMATE CHANGE

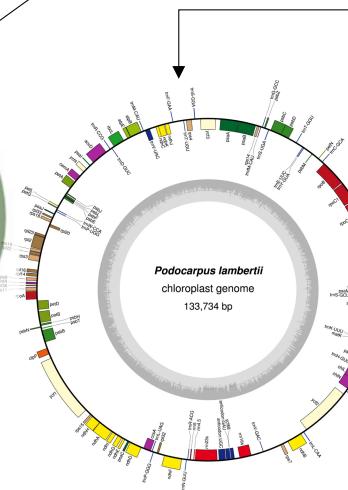


Carbon
fixation

RuBisCO relies on genes located in nuclear and chloroplast genomes

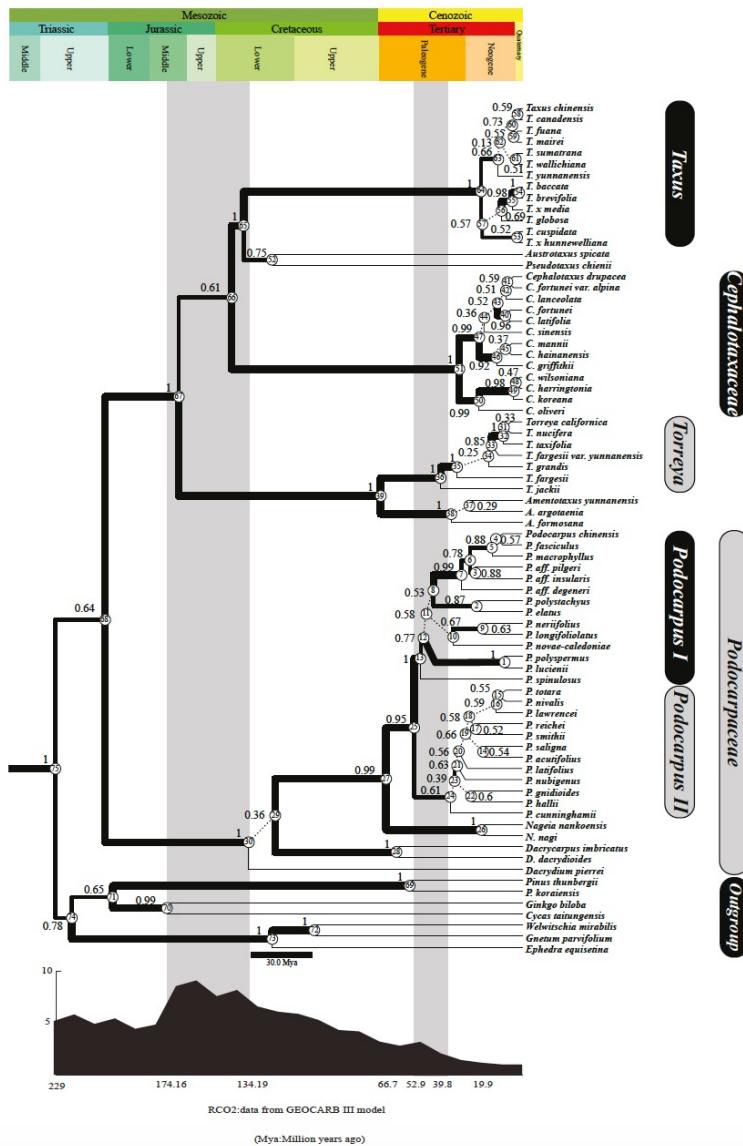


Ribulose-1,5-bisphosphate carboxylase/oxygenase (RuBisCO)

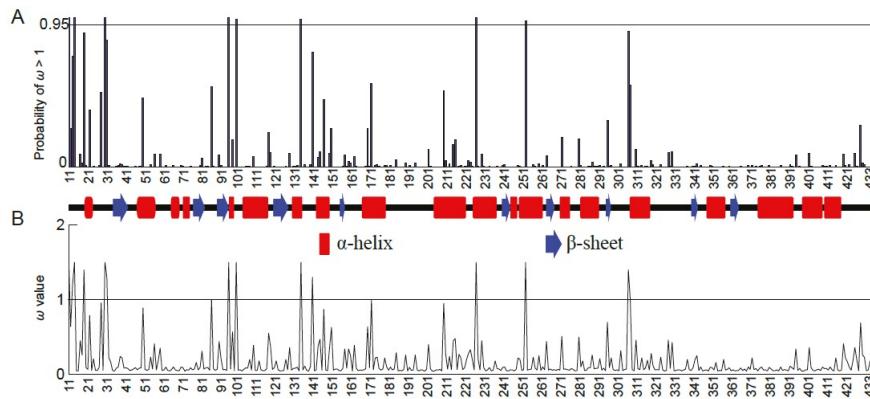


Consists of two subunits:

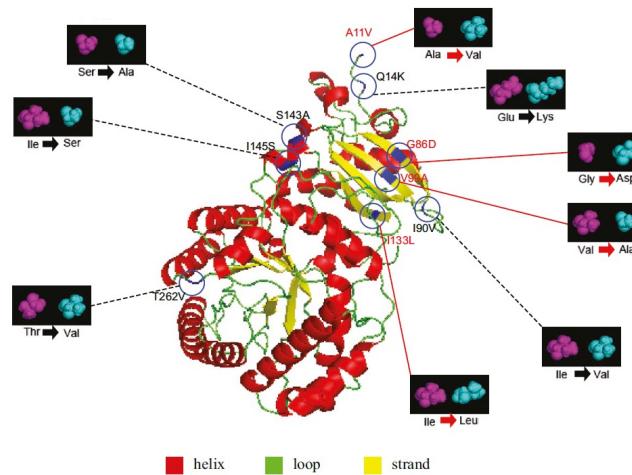
- Large (about 55,000 Da) subunit encoded by chloroplastic *rbcL*.
- Small (about 13,000 Da) subunit encoded by several nuclear genes.



Look for signatures of natural selection along sequence



Positively selected sites evolving in ancestor of Podocarpaceae

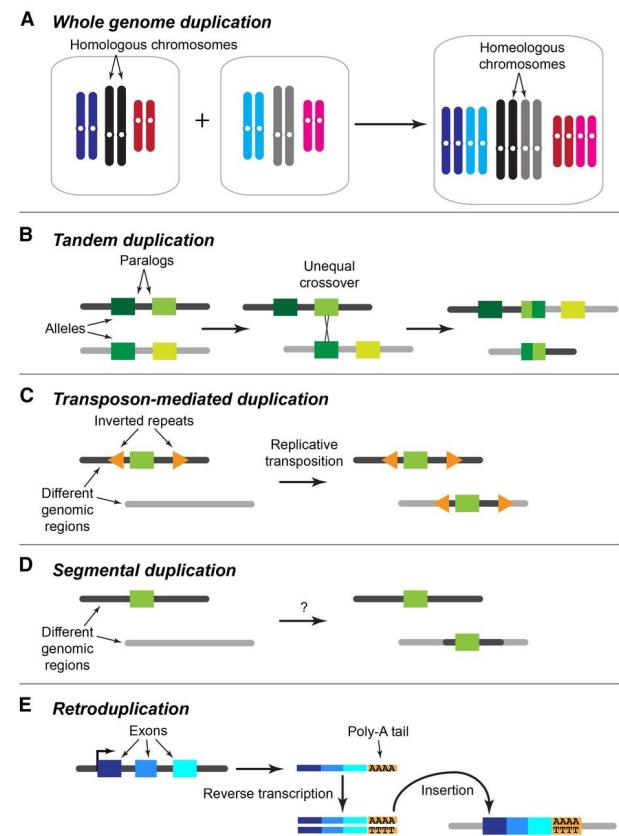


Sen et al. (2011) Biology Direct

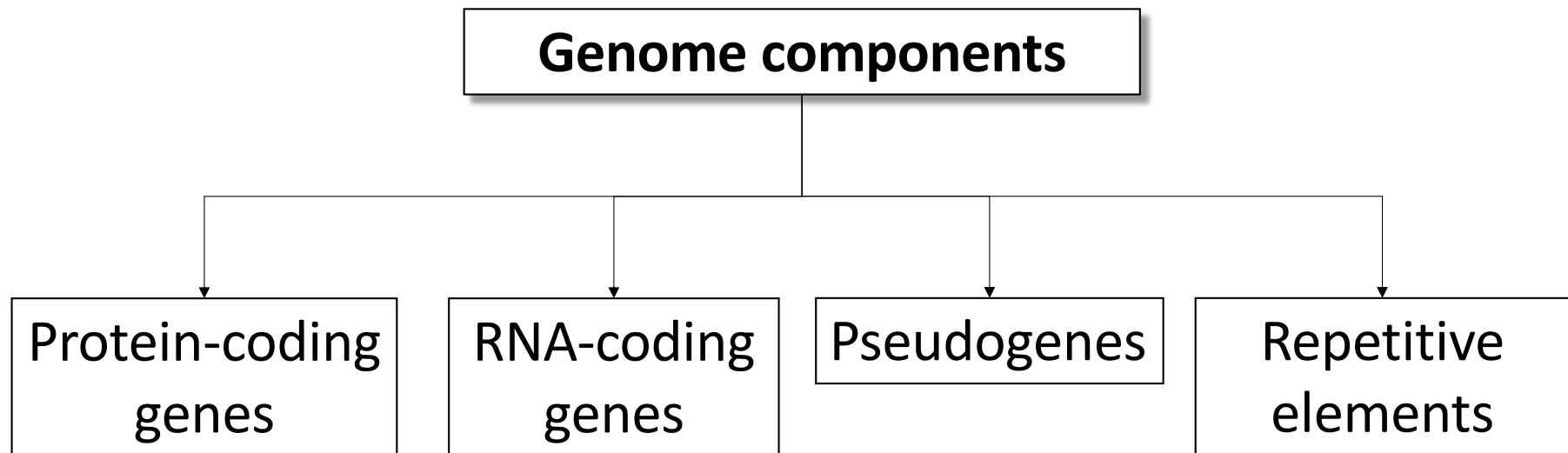
HOW DO GENOMES DIFFER?

Gene duplications

- This process drives the evolution of organisms.
- **Plant diversity has arisen largely following the duplication and adaptive specialization of pre-existing genes (e.g. flower).**
- On average, **65% of annotated genes in plant genomes have a duplicate copy**. Most copies were derived from whole genome doubling (WGD).
- In animals, WGD is less frequent, but there are several cases of large-scale segmental duplications. For instance, this process explains the difference in genome sizes between the human and chimpanzee genomes.



WHAT IS "IN" A GENOME?



WHAT IS "IN" A GENOME?

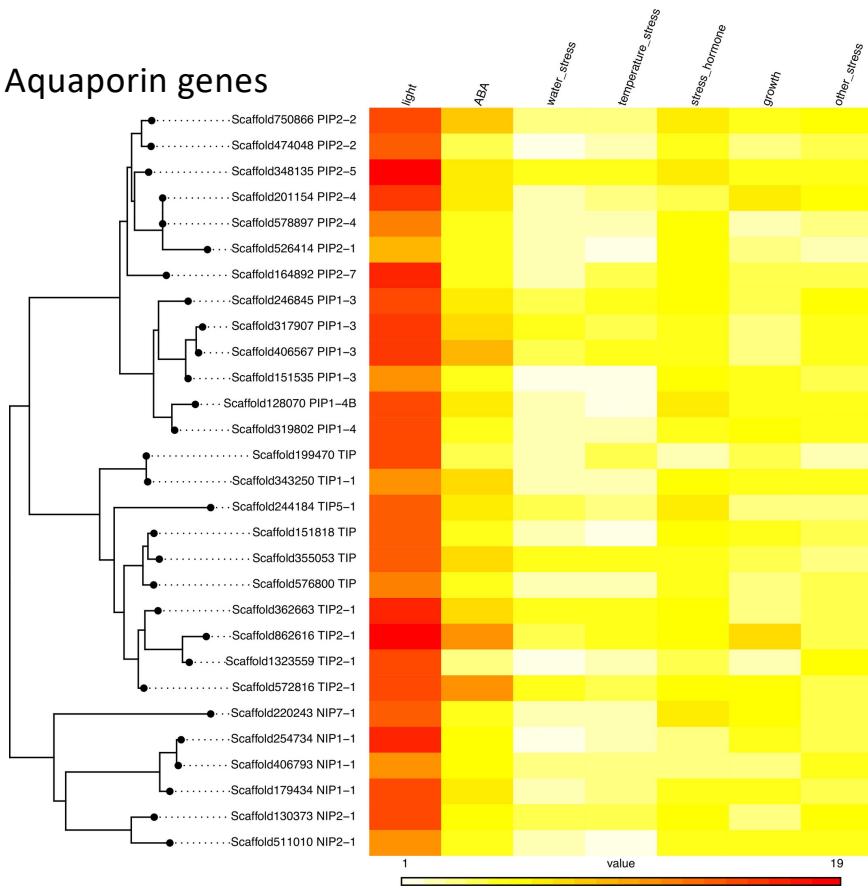
Protein-coding genes

- Small fraction of genomes. For instance, >2-3% in the human genome representing ca. 23,000 genes.
- Many protein-coding genes appear in multiple copies (= paralogues), either identical or diverging into families obtained through gene duplication.

Important for Lab. report

EXAMPLE OF PARALOGUES IN SAGEBRUSH

Aquaporin genes



Aquaporin genes in sagebrush genome

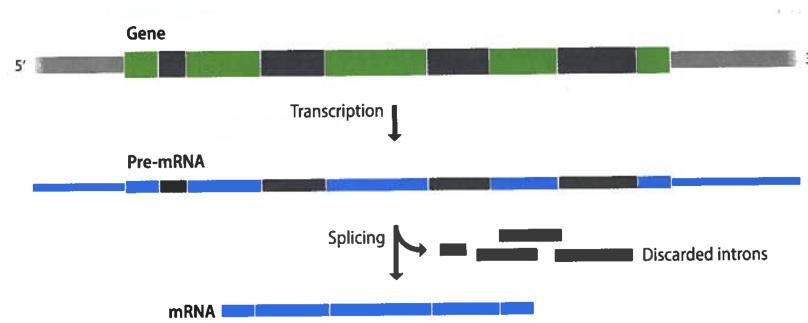
- Mining genome unveiled ca. 50 genes.
- Phylogenetic analyses and protein reconstructions allowed identifying that most of these genes are paralogues, most likely evolving through multiple rounds of whole genome doubling.
- Promoter sequence analysis suggest that they serve different purposes.

<https://onlinelibrary.wiley.com/doi/full/10.1002/ece3.8245>

WHAT IS "IN" A GENOME?

Protein-coding genes

- RNA-Seq helps identifying coding genes in genome sequence, but this process is made difficult due to splicing in eukaryotes.



WHAT IS "IN" A GENOME?

Genes coding for RNAs

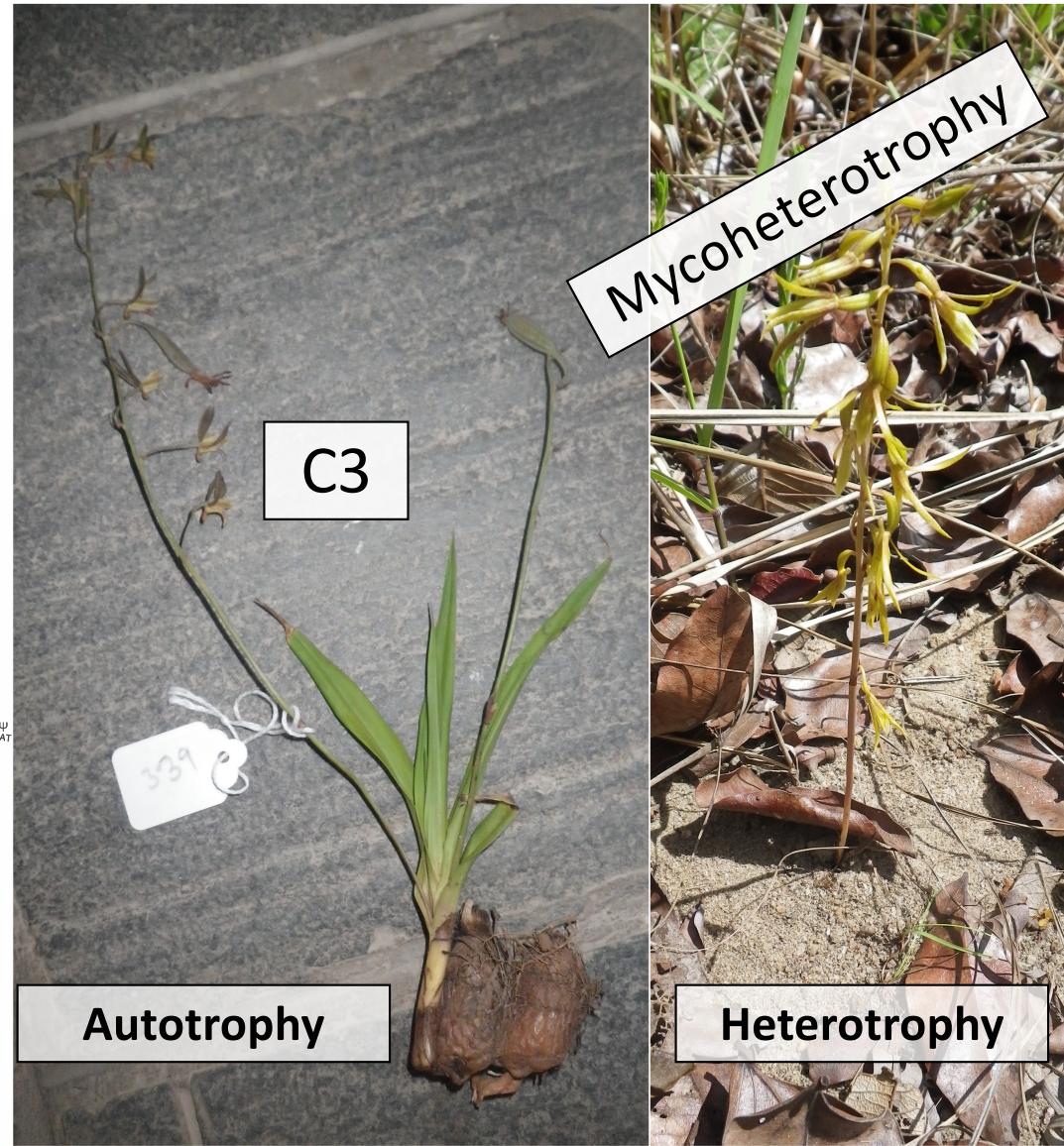
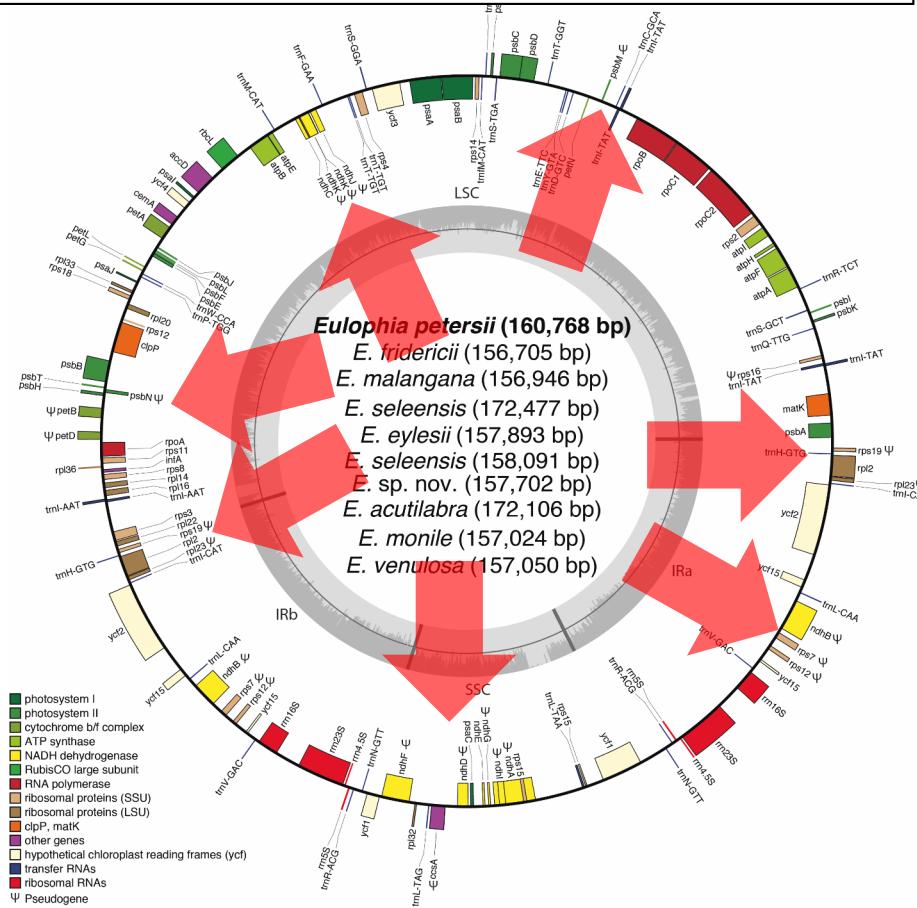
- About 3000 genes coding for RNAs (exclusive of the mRNAs translated to proteins).
- It is becoming clear that the RNA-ome is much richer than previously suspected.
- Except for RNAs involved in the machinery of protein synthesis (e.g. tRNAs and the ribosome itself), most non-coding RNAs are involved in gene expression (e.g., miRNA).

WHAT IS "IN" A GENOME?

Pseudogenes

- Degenerated genes that have mutated so far from their original sequences that the polypeptide sequence they encode will not be functional.
- First step before “losing” genes that are not anymore under selection.

Pseudogenes associated with shifts from autotrophy to heterotrophy in orchids



WHAT IS "IN" A GENOME?

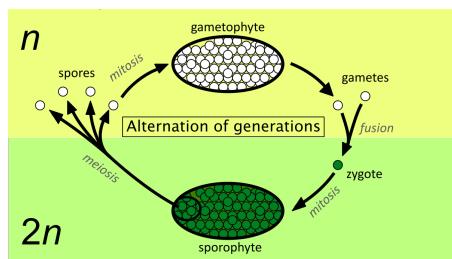
Repetitive elements of unknown function

- In humans, **long and short interspersed elements account for 21% and 13%** of the genome.
- Even more highly repeated sequences – **minisatellites** (about 10-100 base pairs) and **microsatellites** (mostly 2-4 base pairs) – may appear in hundreds of thousands of copies **totaling up to 15%** of the genome. These regions are widely used for population genetic analyzes.
- **Polypliody and the accumulation of repetitive DNA sequences** (often derived from retrotransposons) **are the main factors driving the diversification of genome sizes.**

SEQUENCING STRATEGY

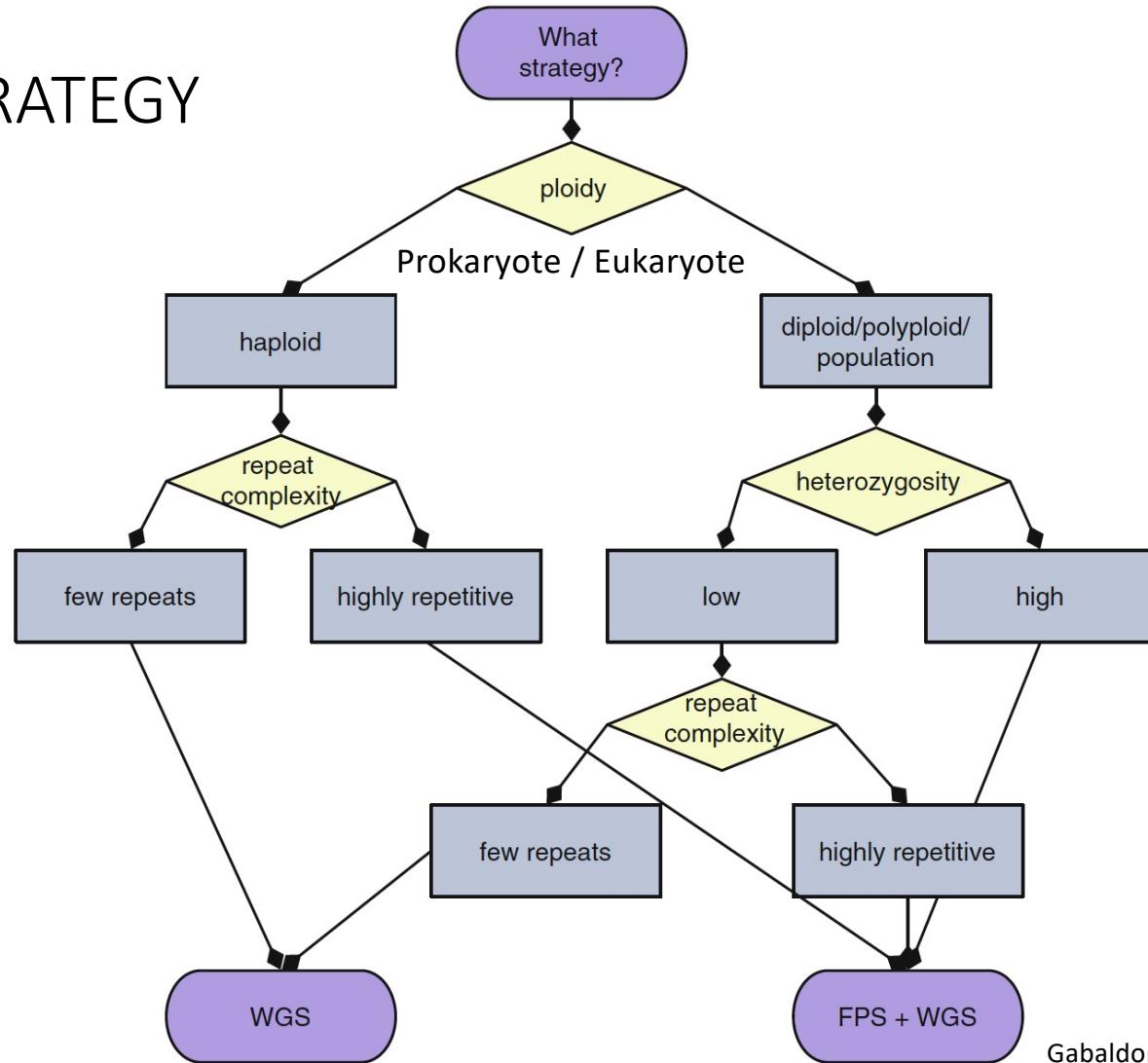
Other aspects to consider:

- Test species monophyly.
- If species is a polyploid complex, sample diploids.
- If heterozygosity is still too high, consider sampling gametophyte (n).



FPS: fosmid pool sequencing (cloning)

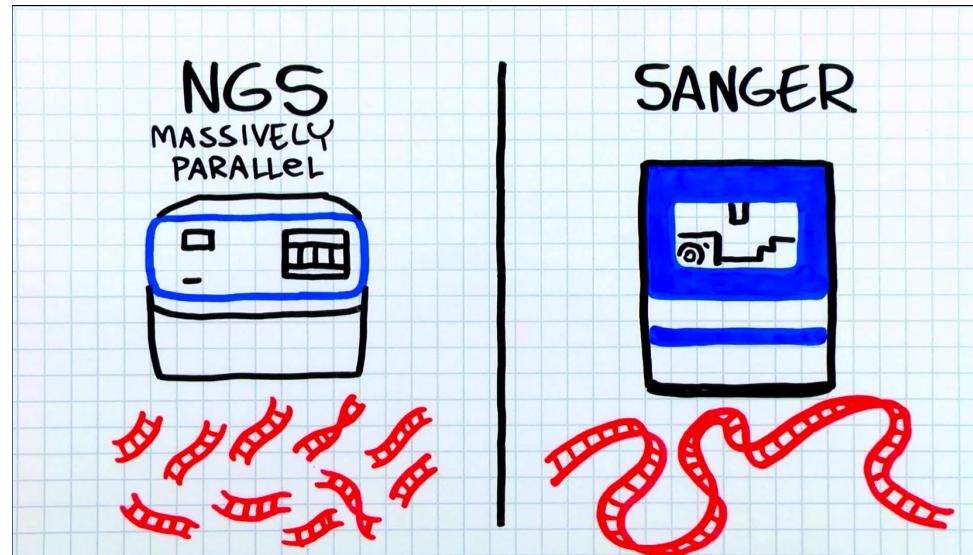
WGS: whole-genome shotgun sequencing



Gabaldon & Alioto (2016)

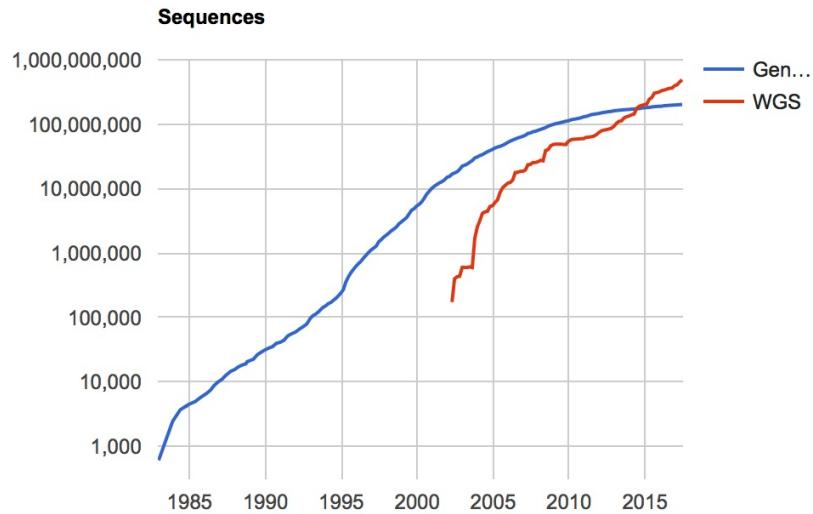
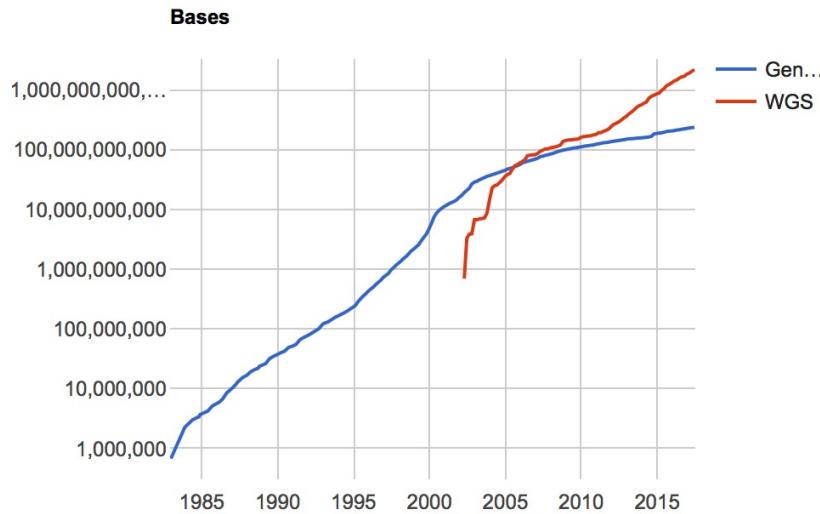
WHAT IS NGS?

- **NGS** is the term applied to methods enabling **thousands of millions of DNA fragments to be sequenced in parallel in a single experiment.**
- NGS outperformed Sanger sequencing, which was able to sequence only individual DNA fragments, each fragment obtained by a different PCR.



WHAT IS NGS?

- NGS methods enable the vast amounts of data needed to assemble an entire genome sequence to be obtained much more rapidly and cheaply than the Sanger sequencing method.
- The emergence of NGS technologies (between 2005-2011) had an impact on the availability of DNA sequences on GenBank.



Mini Report 3

HOW MANY GENOMES ARE PUBLISHED?



Genome > Genome Information by Organism

283,756 genomes published

Organism name (common or scientific) or Accession (Assembly, BioProject or replica)

Search

[Download Reports from FTP site](#)

Overview (33434); [Eukaryotes \(4998\)](#); [Prokaryotes \(126593\)](#); [Viruses \(12995\)](#); [Plasmids \(10991\)](#); [Organelles \(11179\)](#)

Filters Columns Download

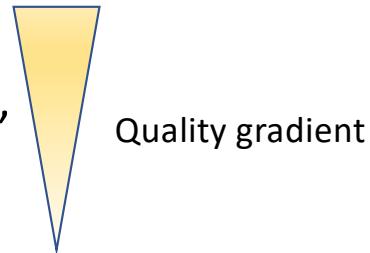
#	Organism Name	Organism Groups	Size(Mb)	Chromosomes	Organelles	Plasmids	Assemblies
1	Pinus lambertiana	Eukaryota;Plants;Land Plants	27,602.7	-	-	-	1
2	Picea glauca	Eukaryota;Plants;Land Plants	24,633.1	-	1	-	3
3	Pinus taeda	Eukaryota;Plants;Land Plants	22,103.6	-	-	-	1
4	Triticum aestivum	Eukaryota;Plants;Land Plants	15,344.7	-	-	-	13
5	Pseudotsuga menziesii	Eukaryota;Plants;Land Plants	14,673.2	-	-	-	1
6	Triticum dicoccoides	Eukaryota;Plants;Land Plants	10,495	14	-	-	2
7	Hordeum vulgare	Eukaryota;Plants;Land Plants	9,788.86	-	-	-	7
8	Rana catesbeiana	Eukaryota;Animals;Amphibians	6,250.35	-	1	-	1
9	Locusta migratoria	Eukaryota;Animals;Insects	5,759.8	-	-	-	1
10	Onychorinus afer	Eukaryota;Animals;Mammals	4,444.09	-	1	-	1

FEEDBACK

BUT NOT ALL GENOMES ARE OF THE SAME QUALITY...

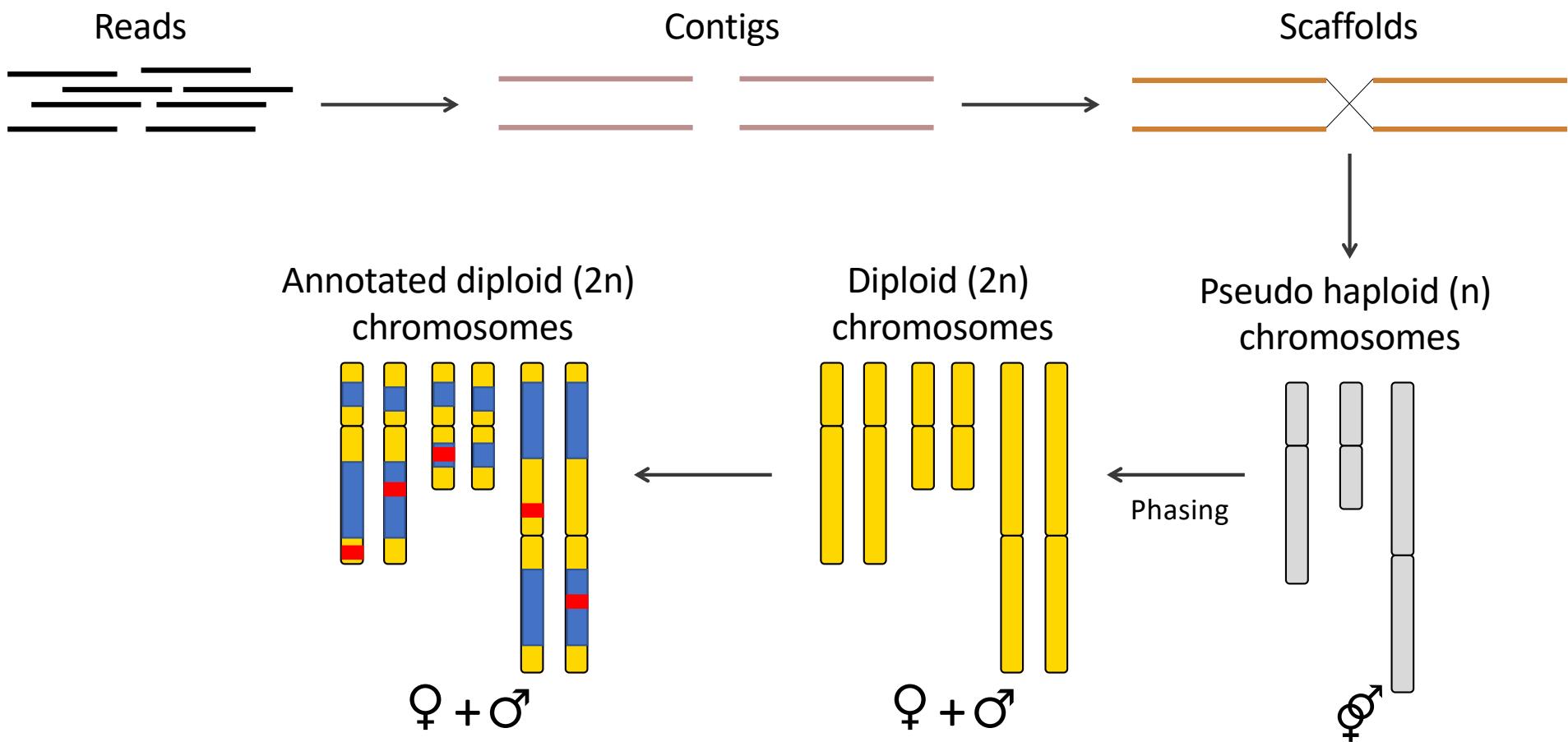
- Not all genomes are assembled and annotated at the same depth.
- Genome assemblies could be submitted as:

- i. Complete,
- ii. Chromosome,
- iii. Scaffold,
- iv. Contig.



- The quality of the submitted genomes depend on the complexity of the genome (i.e. amount of repeated DNA, level of ploidy), but also on the aims of the researchers submitting the genome (why producing a fully annotated genome when we do not need it?)

Genome assembly and annotation workflow



WHAT DOES IT MEAN TO SEQUENCE A GENOME?

- Ideally, a **draft genome would represent the complete nucleotide base sequence for all chromosomes** in the species of interest, a ‘physical map’ of its genetic content.
- There are four major complications with the concept of a “genome sequence”:
 1. There is not one true sequence for a species due to individual SNPs.
 2. It is essentially impossible to sequence and assemble all nucleotides in the genome due to repetitive DNA.
 3. There will always be some degree of error in the characterized genome sequence: sequencing and assembly errors.
 4. Every genome assembly is the result of a series of assembly heuristics and should accordingly be treated as a working hypothesis.

Conclusions: Many current genome assemblers produced useful assemblies, containing a significant representation of their genes and overall genome structure. However, the high degree of variability between the entries suggests that there is still much room for improvement in the field of genome assembly and that approaches which work well in assembling the genome of one species may not necessarily work well for another.

Assemblathon 2: evaluating *de novo* methods of genome assembly in three vertebrate species

Keith R Bradnam^{1*†}, Joseph N Fass^{1†}, Anton Alexandrov³⁶, Paul Baranay², Michael Bechner³⁹, Inanç Birol³³, Sébastien Roisvert^{10,11}, Jarrod A Chapman²⁰, Guillaume Chenuic^{7,9}, Raven Chikhi^{7,9}, Hamidreza Chitsaz⁶, Docking³³, Richard Durbin³⁴, Dent Earl⁴⁰,

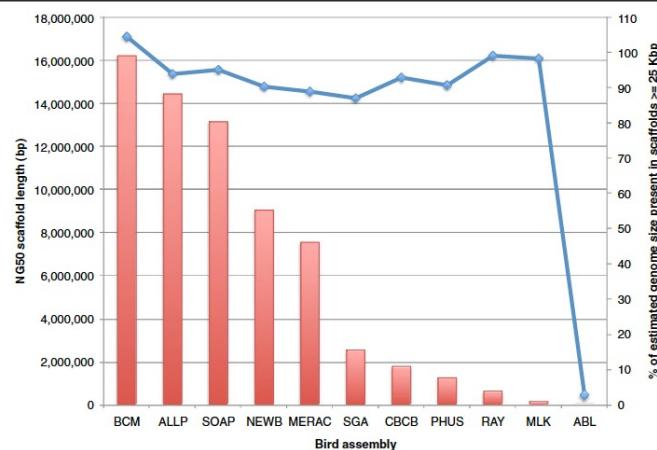


Figure 4 NG50 scaffold length distribution in bird assemblies and the fraction of the bird genome represented by gene-sized scaffolds. Primary Y-axis (red) shows NG50 scaffold length for bird assemblies: the scaffold length that captures 50% of the estimated genome size (~1.2 Gbp). Secondary Y-axis (blue) shows percentage of estimated genome size that is represented by scaffolds ≥ 25 Kbp (the average length of a vertebrate gene).

This field, especially bioinformatics, is exploding and methods for *de novo* genome assembly are improving every day!