

Prof. Dr. Boas Pucker  
**PBPM-BP-07**

# Availability of slides

- All materials are freely available (CC BY) - after the lectures:
  - eCampus: PBPM0 - Plant Biochemistry, Physiology and Molecular Biology (LEC)
  - GitHub: <https://github.com/bpucker/teaching/PBPM>
- Questions: Feel free to ask at any time
- Feedback, comments, or questions: [pucker\[a\]uni-bonn.de](mailto:pucker[a]uni-bonn.de)



My figures and content can be re-used in accordance with CC BY 4.0, but this might not apply to all images/logos. Some figure were constructed using bioRender.com.

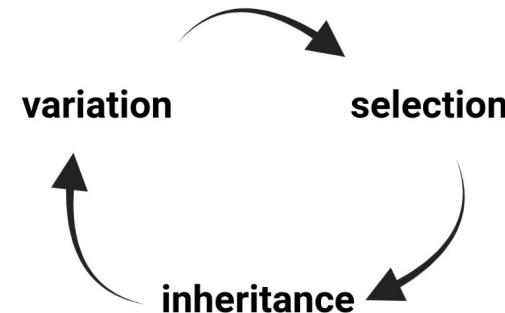
- Transformation methods
- Transformation examples: purple tomato, golden rice
- Genome editing: CRISPR/Cas9
- Legal implications (EU)
- Plant SynBio (iGEM)

“Nothing in biology makes sense except in the light of evolution”

Theodosius Dobzhansky

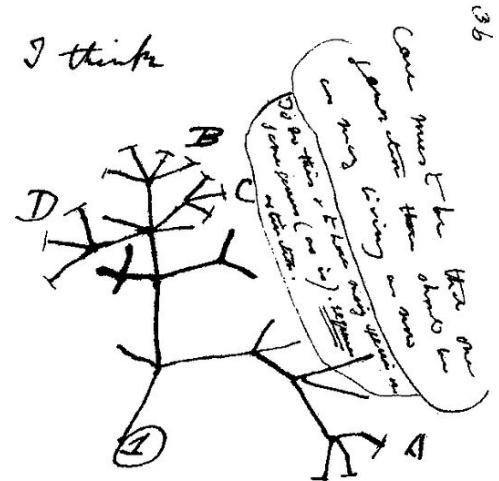
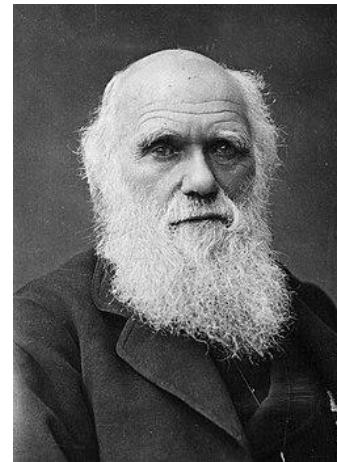
# Evolutionary concepts

- More offspring than needed
- Genetic basis of traits responsible for fitness
- Selection based on fitness



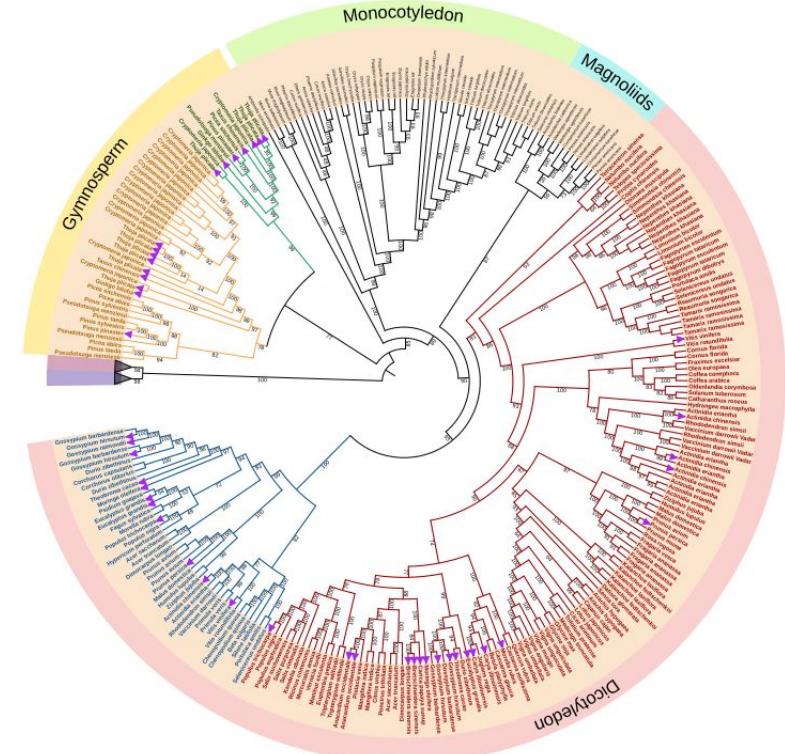
# Charles Darwin: I think

- Developed a theory to explain observations
- Most influential theory in biology?
- Theory vs. hypothesis
- Proving a theory?



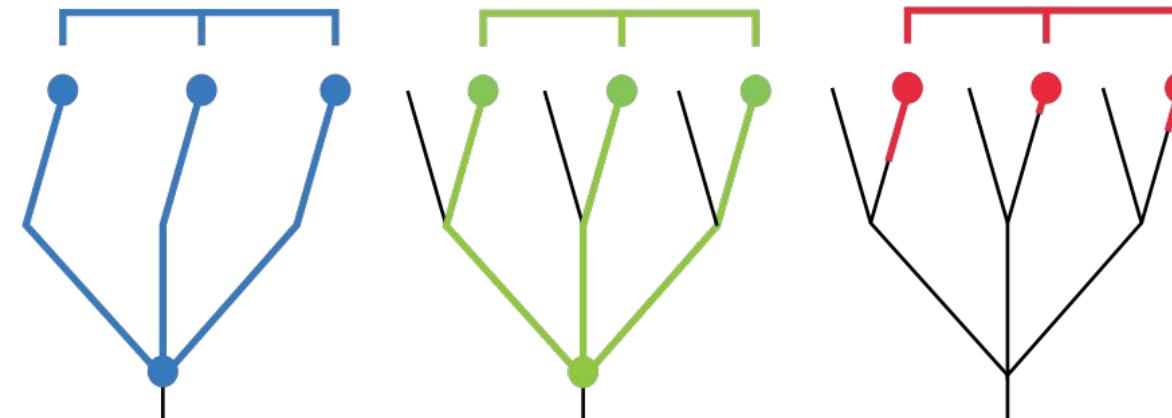
Thus between A & B. various  
sorts of relation. C & B. the  
first generation, B & D  
rather greater distinction.  
Thus genera would be  
formed. - binary relation

# Concept of a tree



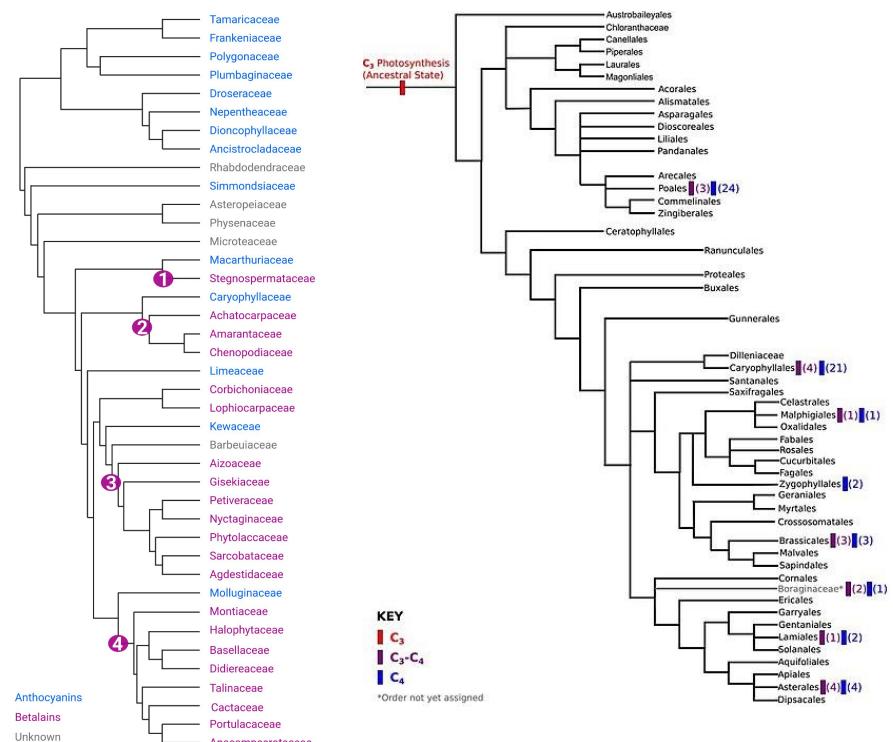
# Monophyletic, paraphyletic, and polyphyletic

- Monophyletic = all members of group have common ancestor in the group and all descendants of that ancestor are in group
- Paraphyletic = Group members have common ancestor in the group, but not all descendants of that ancestor are included in group
- Polyphyletic = Group members do not have common ancestor in the group

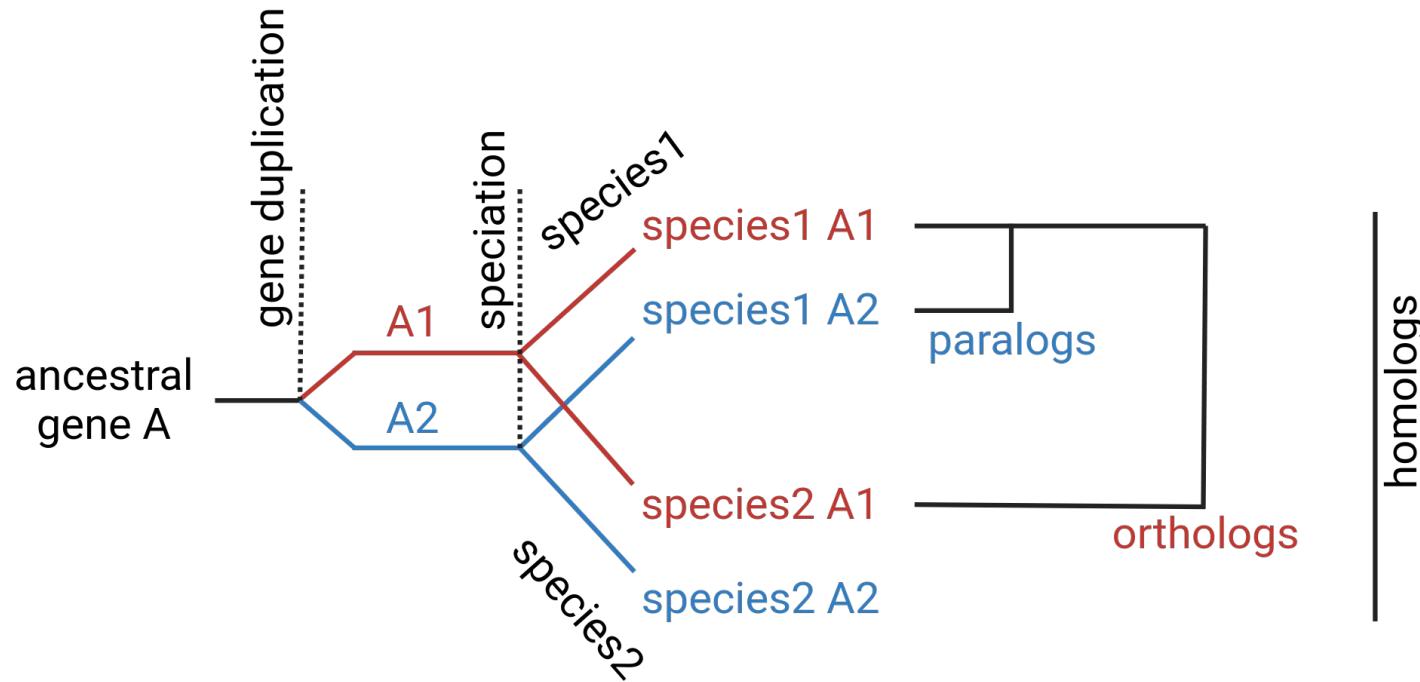


# Convergent and independent evolution

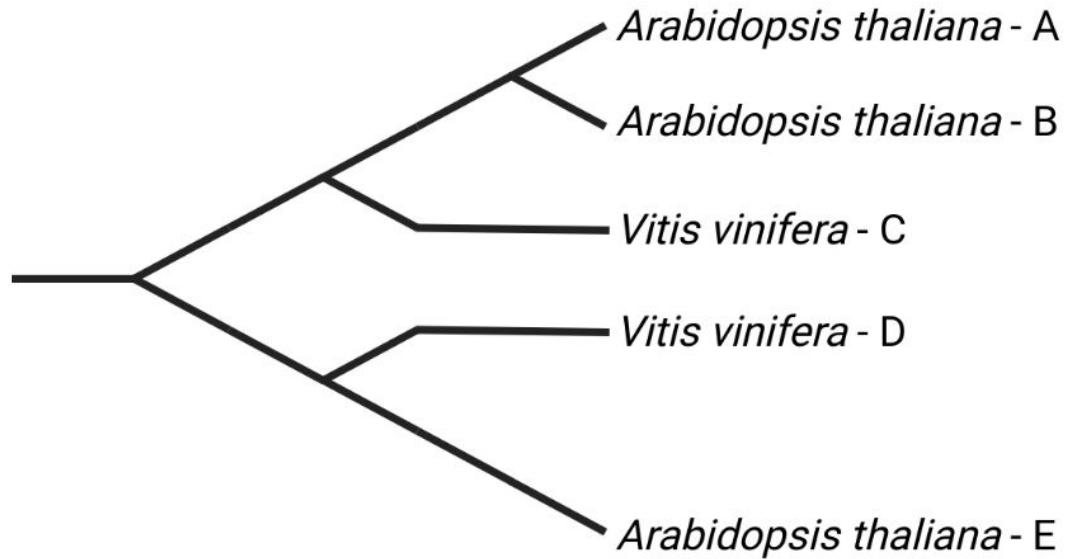
- Trait appears multiple times independently
- Same genes can be involved and modified in similar ways
- Examples: pigment biosynthesis in Caryophyllales; photosynthesis type



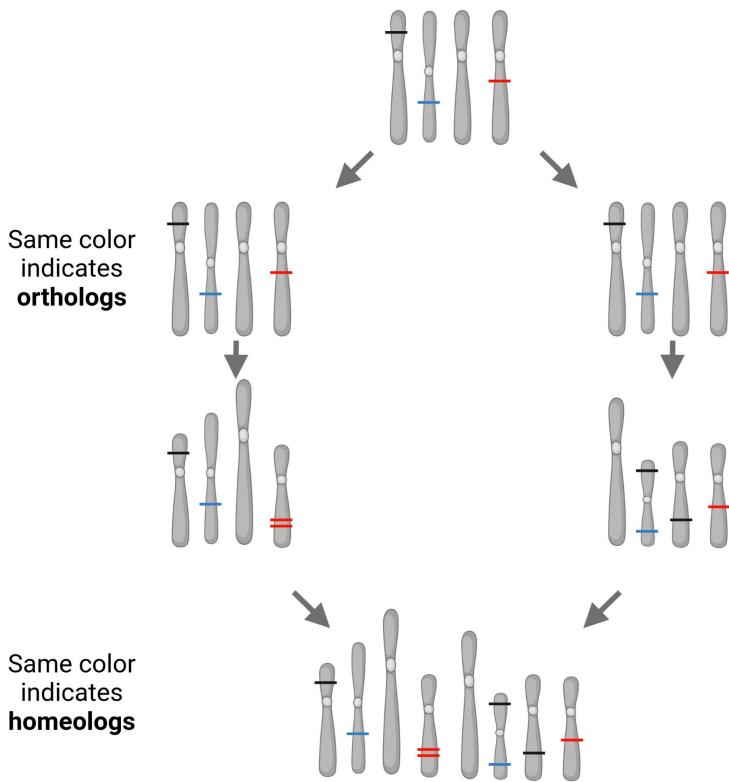
# Homologs, paralogs, orthologs



# Which are orthologs/paralogs?



# Homeologs

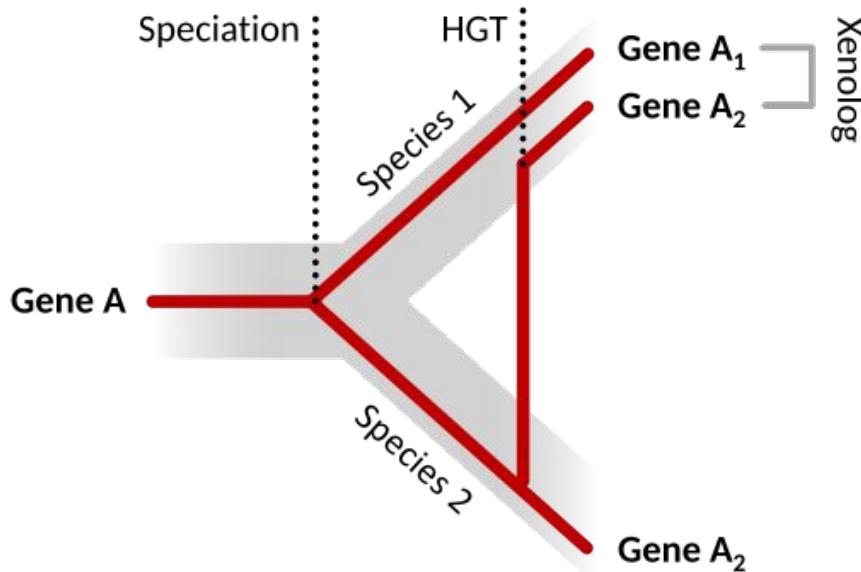


Speciation

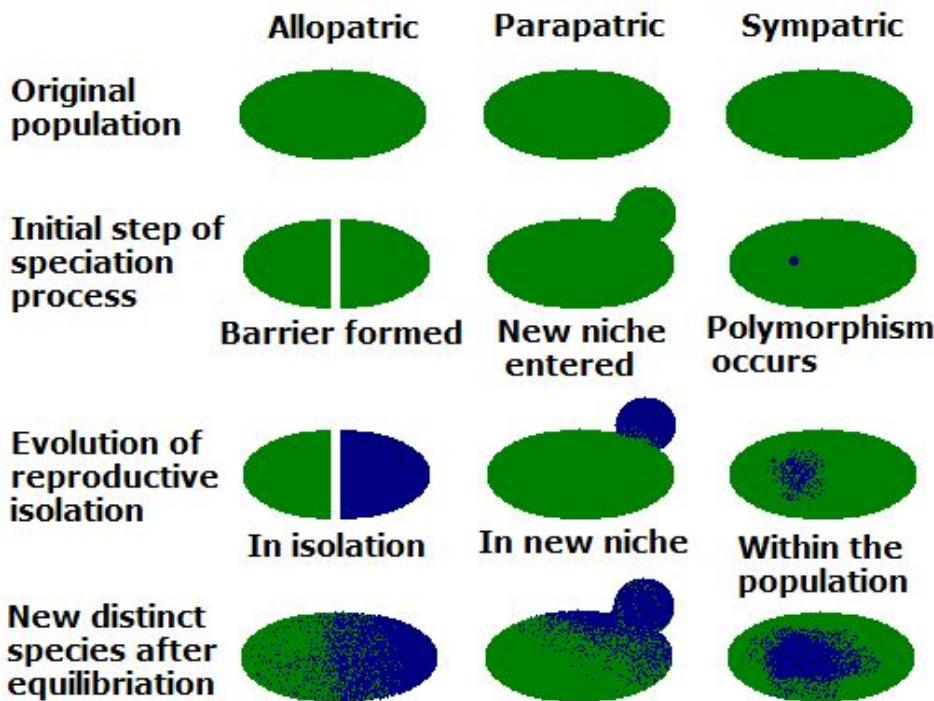
Gene duplications,  
translocations, and  
rearrangements

Hybridization leads to  
polyploidization

- Genes can be moved between species (horizontal gene transfer)
- Relationships turn from tree into networks
- Bacteria/viruses can transfer genes between plant species



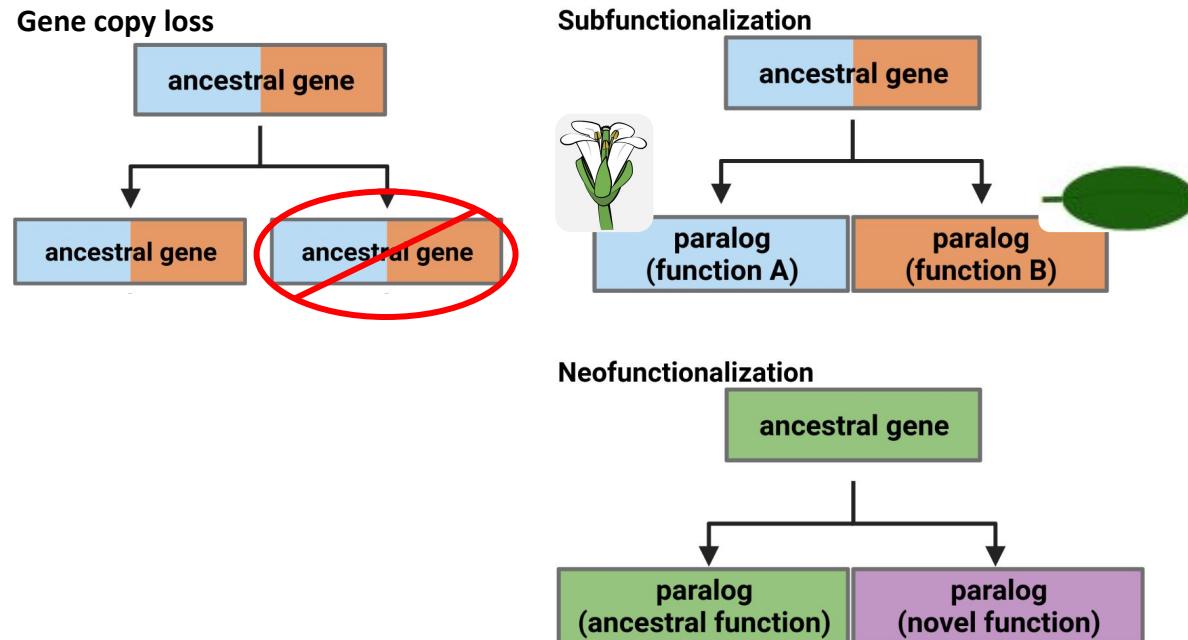
- sympatric = speciation due to genetic changes
- parapatric = speciation in a niche
- allopatric = physical separation of population leads to speciation



# Gene duplication

# Importance of gene duplications in plants

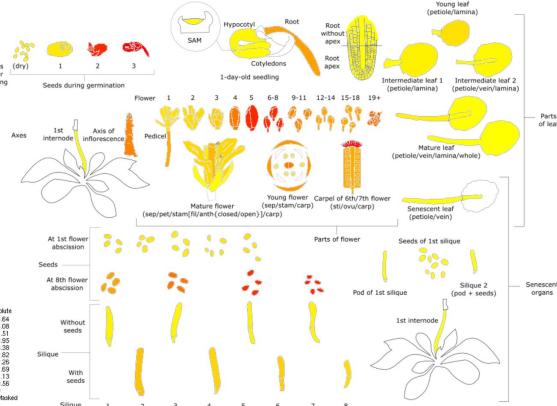
- 65% of annotated plant genes have a copy
- Most copies do not persist during evolution
- Gene duplications are crucial for metabolic novelty in specialized plant metabolism



# Example: SG7 MYB expression

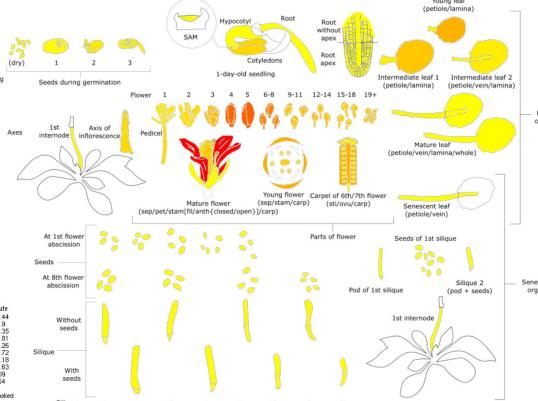
AT2G47460/AT2G47460

Klepikova Arabidopsis Atlas eFP Browser at bar.utoronto.ca  
Klepikova et al. 2016. Plant J. 88:1058-1070



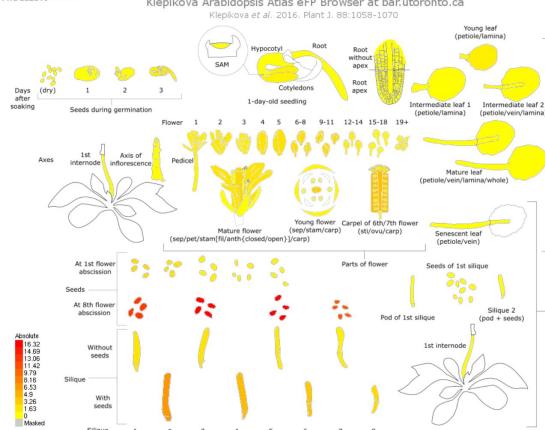
AT5G8030/AT5G8030

Klepikova Arabidopsis Atlas eFP Browser at bar.utoronto.ca  
Klepikova et al. 2016. Plant J. 88:1058-1070



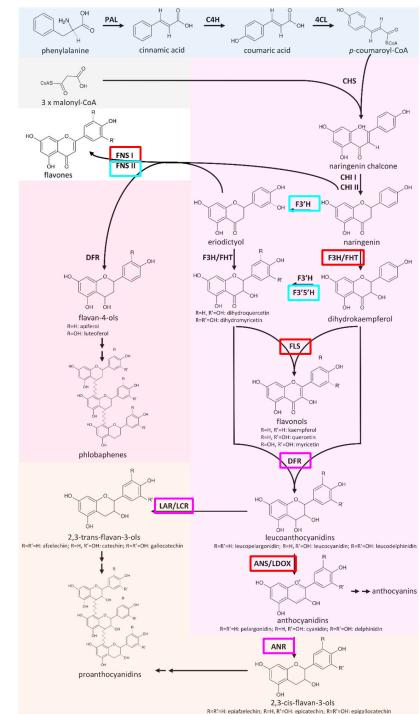
AT3G02610/AT3G02610

Klepikova Arabidopsis Atlas eFP Browser at bar.utoronto.ca  
Klepikova et al. 2016. Plant J. 88:1058-1070



# Importance of gene duplications in plants

- 65% of annotated plant genes have a copy
- Most copies do not persist during evolution
- Gene duplications are highly relevant in specialized metabolism:
  - FNS I / F3H / FLS / ANS (2-ODDs)
  - FNS II / F3'H / F3'5'H (CYPs)
  - DFR / LAR / ANR (SDRs)



# Gene duplication mechanisms

- Tandem duplication: copies are located next to each other on same chromosome
- Segment duplication: like tandem duplications, but multiple adjacent genes are also duplicated
- Whole genome duplication (WGD): all genes are duplicated

Tandem gene duplication



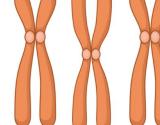
Segmental gene duplication



Whole genome duplication

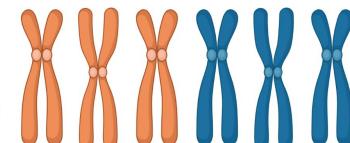


Chr1 Chr2 Chr3

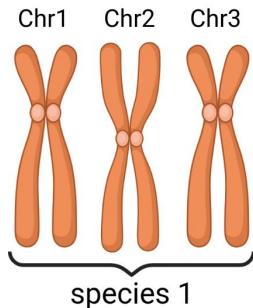


→  
whole  
genome  
duplication

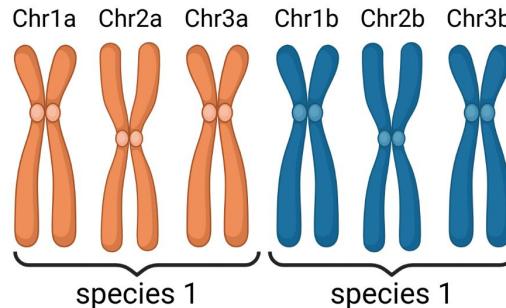
Chr1a Chr2a Chr3a Chr1b Chr2b Chr3b



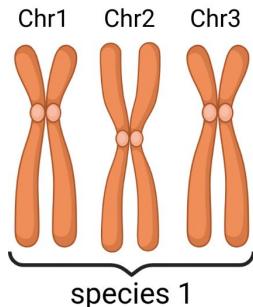
# WGD: autopolyploidy vs. allopolyploidy



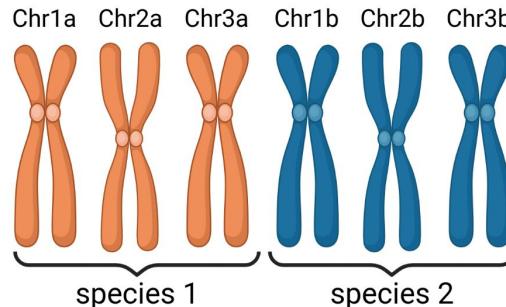
whole genome duplication



autopolyploidy



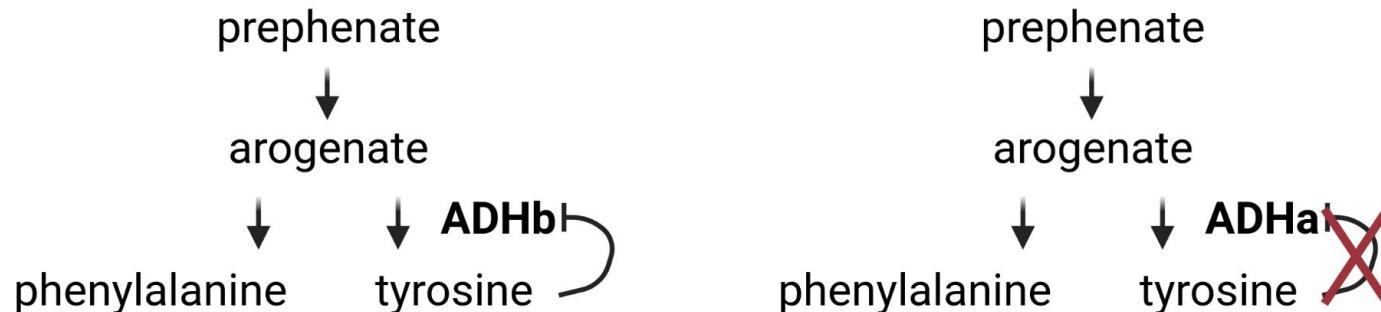
whole genome duplication



allopolyploidy

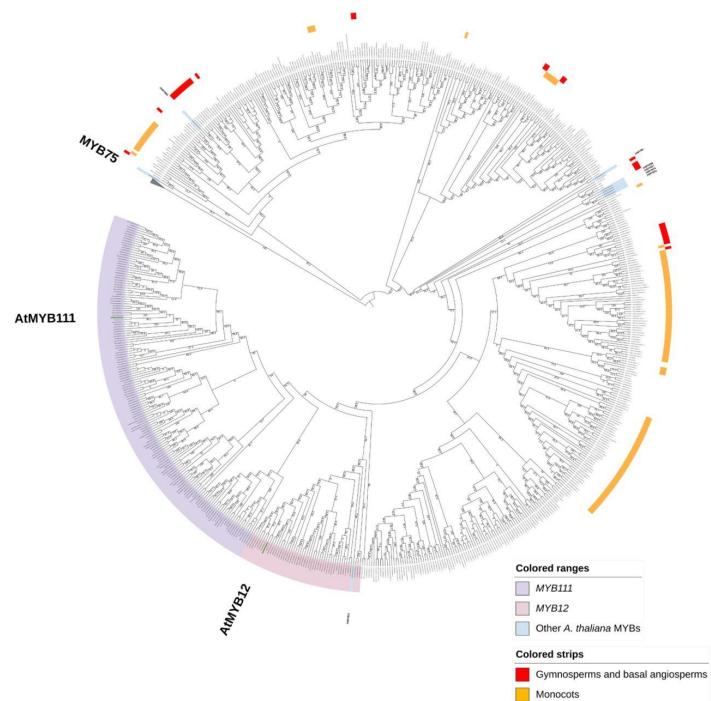
## Example: ADH in Caryophyllales

- ADH was duplicated in Caryophyllales
  - ADHa = feedback-resistance enzyme
  - ADHb = feedback-inhibited enzyme

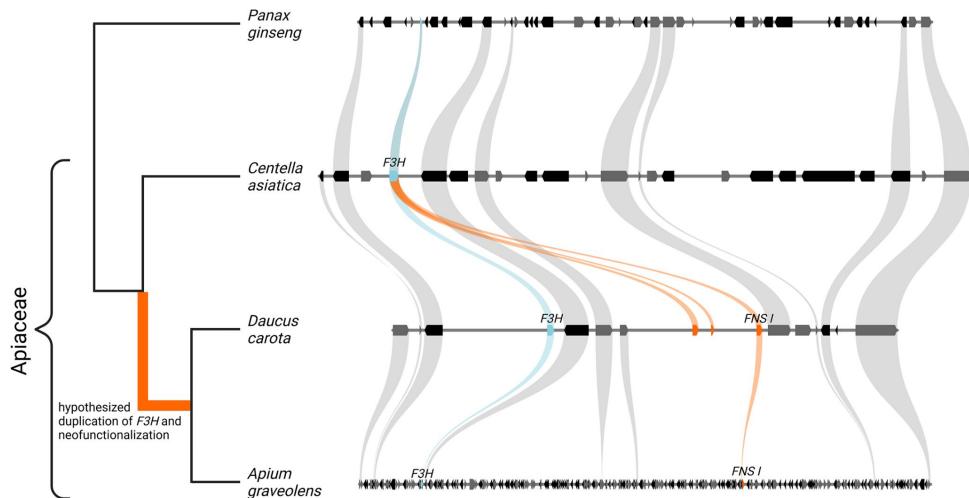
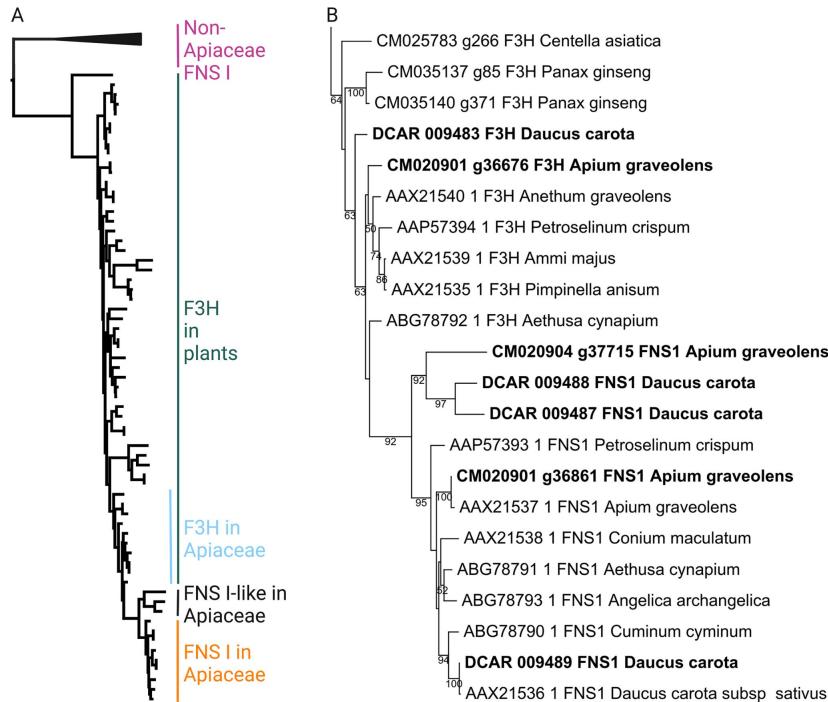


## Example: Flavonol regulating MYBs

- AtMYB111 and AtMYB12 are regulators of the flavonol biosynthesis
- Deep duplication of a flavonol-regulating MYB in angiosperms



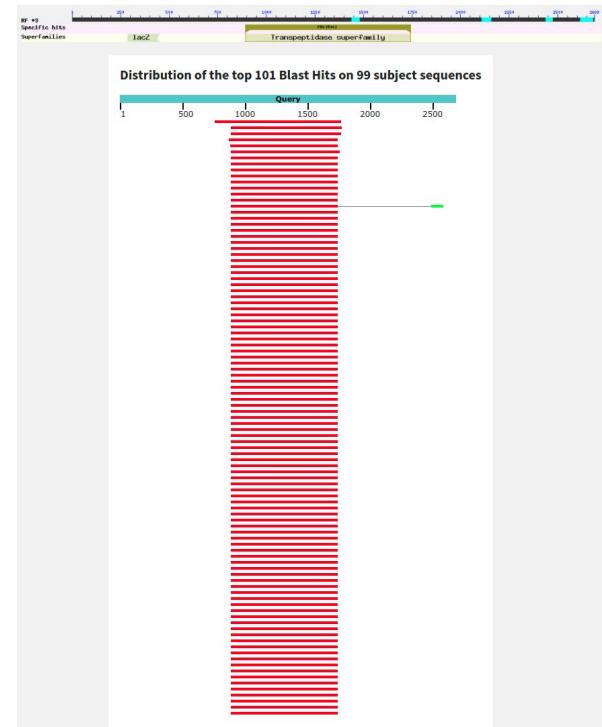
# Example: *F3H/FNS1* in Apiaceae



# Tools in phylogenetics

# Collecting sequences: BLAST

- BLAST = Basic Local Alignment Search Tool
- Identification of short stretches of sequence similarity
- Quick search against a large database
- Sequence similarity suggest homology i.e. common ancestors and similar function



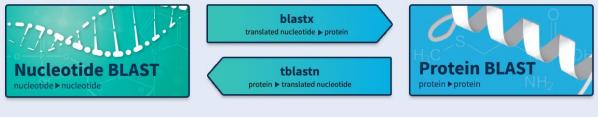
# Online vs. local BLAST

NCBI website	Local BLAST
Convenient to use	Requires some command line knowledge
Requires no local computational resources	Computational resources needed
Transfer of query sequence	Query kept secret
Large databases available	Download of databases required

**Basic Local Alignment Search Tool**  
 BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance.



#### Web BLAST



#### BLAST Genomes

Search

Human    Mouse    Rat    Microbes

```
USAGE
blastn [-h] [-help] [-import_search_strategy filename]
        [-export_search_strategy filename] [-task task_id] [-db database_name]
        [-taxid taxid] [-taxids taxids] [-taxidlist filename]
        [-negative_glist filename] [-negative_seqidlist filename]
        [-taxids taxids] [-negative taxids taxids] [-taxidlist filename]
        [-import_search_strategy filename] [-task task_id] [-db database_name]
        [-db_soft_mask filtering_algorithm] [-db_hard_mask filtering_algorithm]
        [-subject subject_input_file] [-subject_loc range] [-query input_file]
        [-gapopen float_value] [-gapextend float_value] [-gapcost float_value]
        [-gapopen open_penalty] [-gapextend extend_penalty]
        [-perc_identity float_value] [-qcov_hsp_perc float_value]
        [-max_hsps int_value] [-max_hsps_per_query float_value]
        [-xdrop_gap final float_value] [-searchsp int_value] [-penalty penalty]
        [-reward reward] [-no_greedy] [-min raw_gapped_score int_value]
        [-threshold float_value] [-length int_value] [-dust dust_options]
        [-filtering_db filtering_database]
        [-window_nasker taxid window_nasker taxid]
        [-uniprot_uniref uniprot_uniref] [-softmasking soft_masking]
        [-ungapped] [-colliding limit int_value] [-best_hit_overhang float_value]
        [-best_hit_score_edge float_value] [-subject_bestsht]
        [-subject_length int_value] [-query_length int_value]
        [-use_index boolean] [-index_name string] [-icase_masking]
        [-query_loc range] [-strand strand] [-parse_deflines] [-output format]
        [-query_length int_value] [-subject_length int_value]
        [-line_length line_length] [-html] [-sorthits num_sort_hits]
        [-sortby sort_by] [-max_target_seqs num_sequences]
        [-num_threads int_value] [-remote] [-version]
```

DESCRIPTION
 Nucleotide-Nucleotide BLAST 2.9.0+

## BLAST alternatives

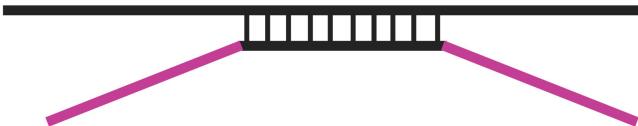
- HMMER: <http://hmmer.org/>
  - Dedicated to the identification of specific domains
- DIAMOND: <https://github.com/bbuchfink/diamond>
  - Much faster, but higher memory requirements



# Local vs. global alignment

- Local alignment:
  - highlights small stretches with high similarity
  - random hit possible
  - better for strong sequence length differences
  
- Global alignment:
  - only works by overall high similarity
  - similar sequence length required
  - full length similarity supports similar function

Local alignment



Global alignment



- MAFFT

## Multiple Sequence Alignment

MAFFT (Multiple Alignment using Fast Fourier Transform) is a high speed multiple sequence alignment program.

We have recently changed the default parameter settings for MAFFT. Alignments should run much more quickly and larger DNA alignments can be carried out by default.  
Please click the 'More options' button to review the defaults and change them if required.

**Important note:** This tool can align up to 500 sequences or a maximum file size of 1 MB.

STEP 1 - Enter your input sequences

Enter or paste a set of

AUTOMATIC

sequences in any supported format:

Or upload a file:  No file selected.

[Use a example sequence](#) | [Clear sequence](#) | [See more example inputs](#)

STEP 2 - Set your Parameters

OUTPUT FORMAT

Pearson/Fasta

The default settings will fulfill the needs of most users.

[More options...](#) (Click here, if you want to view or change the default settings.)

STEP 3 - Submit your job

Be notified by email (Tick this box if you want to be notified by email when the results are available)

MAFFT v7.453 (2019/Nov/8)  
<https://mafft.cbrc.jp/alignment/software/>  
 MBE 30:772-780 (2013), NAR 30:3059-3066 (2002)

High speed:  
 % mafft in > out  
 % mafft --retree 1 in > out (fast)

High accuracy (for <-200 sequences x <-2,000 aa/nt):  
 % mafft --maxiterate 1000 --localpair in > out (% linsi in > out is also ok)  
 % mafft --maxiterate 1000 --genafpair in > out (% einsi in > out)  
 % mafft --maxiterate 1000 --globalpair in > out (% ginsi in > out)

If unsure which option to use:  
 % mafft --auto in > out

```
--op # :      Gap opening penalty, default: 1.53
--ep # :      Offset (works like gap extension penalty), default: 0.0
--maxiterate # : Maximum number of iterative refinement, default: 0
--clustalout : Output: clustal format, default: fasta
--reorder :   Outorder: aligned, default: input order
--quiet :     Do not report progress
--thread # :  Number of threads (if unsure, --thread -1)
--dash :      Add structural information (Rozewicki et al., submitted)
```

- Letters representing amino acid residues are arranged to highlight similarities
- Gaps (-) are inserted to show differences

Input FASTA file with sequences

```
>seq1
MDTEKYMEKWIDQGHALFPEEDQ
>seq2
MDTDIKYSEKWIQGHSLFPDDQ
>seq3
METEKYMEKWIQGHSIFPEDQ
>seq4
MDTEIKYMEKWIQGHALFPDDQP
```



Alignment of sequences

seq1	MDTE-KYM <b>E</b> KWIDQGH <b>A</b> LFP <b>E</b> EDQ-
seq2	MDTD <b>I</b> KY <b>S</b> E <b>K</b> W <b>I</b> -Q <b>G</b> H <b>S</b> LF <b>P</b> D <b>D</b> Q-
seq3	METE-KYM <b>E</b> KW <b>I</b> -Q <b>G</b> H <b>S</b> IF <b>P</b> E-D <b>D</b> Q-
seq4	MDTE <b>I</b> KY <b>M</b> E <b>K</b> W <b>I</b> -Q <b>G</b> H <b>A</b> LF <b>P</b> D <b>D</b> QP

# Alignment formats: CLUSTAL vs. FASTA

AthCHS	MVMAGASSLDEIRQAQRADGPAGILAIGTANPENHVLOAEYPDYYFRITNSEHMTDLKEK
BvuCHS	M---ATPSVQEIRDAQRSNGPATILAIGTANPANEMYQAEYPDFYFRVTKSEHMSSELQK
	* . : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : *
AthCHS	FKRMCDKSTIRKRHMHLTEEFLKENPHMCAYMAPSLDTRQDIVVVEVPKLGKEAAVKAIK
BvuCHS	FKRMCDKSMIKKRYMHVTQELLEENPHMCDYNASSLNTRQDILATEVPKLGKEAAVKAIK
	***** * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : *
AthCHS	EWGQPKSKITHVVFCCTSGVDMPGADYQLTLKLGLRPSVKRLMMYQQGCFAAGTVLRIAK
BvuCHS	EWGQPRSKITHVIFCTTSGVDMPGADYQLTLKLGLRPSVKRFMLYQQGCYAGGTVLRLAK
	***** * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : *
AthCHS	DIAENNARGLVVCSEITAVTFRGSDTHLDLSVGQALFSDGAAALIVGSDPDTSVGEK
BvuCHS	DIAENNARGLVVAEITVICFRGPTETHLDMSGQALFGDGAGAVIVGADLDESI-ER
	* : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * :
AthCHS	PIFEMVSAAQTILPDSGAIKGHLREVGLTFHLLKDVPGLISKNIKSLDEAFKPLGIDD
BvuCHS	PIFQLWAAQTILPDSEGAIKGHLREVGLAFHLLKDVPGLISKNIKALVEAFKPIGIDD
	*** : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : *
AthCHS	WNSLFWIAHPGGPAILDQVEIKLGLKEEKMRATRHLSEYGNMSSACVLFILDEMRRSA
BvuCHS	WNSIFWAAHPGGPAILDQVESKLGKQDKLSTTRHVLSEFGNMSSACVLFILDEMRRKRM
	*** : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : *
AthCHS	KDGVATTGEGLEWGVLFGFGPGLTVETVVLHSVPL--
BvuCHS	KEGMATTGEGLEWGVLFGFGPGLTVETVMLHSVPIAN
	* : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : * : *

```
>AthCHS
MVMAGASSLDEIRQAQRADGPAGILAIGTANPENHVLOAEYPDYYFRITNSEHMTDLKEK
FKRMCDKSTIRKRHMHLTEEFLKENPHMCAYMAPSLDTRQDIVVVEVPKLGKEAAVKAIK
EWGQPKSKITHVVFCCTSGVDMPGADYQLTLKLGLRPSVKRLMMYQQGCFAAGTVLRIAK
DIAENNARGLVVCSEITAVTFRGSDTHLDLSVGQALFSDGAAALIVGSDPDTSVGEK
PIFEMVSAAQTILPDSGAIKGHLREVGLTFHLLKDVPGLISKNIKSLDEAFKPLGIDD
WNSLFWIAHPGGPAILDQVEIKLGLKEEKMRATRHLSEYGNMSSACVLFILDEMRRSA
KDGVATTGEGLEWGVLFGFGPGLTVETVVLHSVPL--
```

```
>BvuCHS
M---ATPSVQEIRDAQRSNGPATILAIGTANPANEMYQAEYPDFYFRVTKSEHMSSELQK
FKRMCDKSMIKKRYMHVTQELLEENPHMCDYNASSLNTRQDILATEVPKLGKEAAVKAIK
EWGQPRSKITHVIFCTTSGVDMPGADYQLTLKLGLRPSVKRFMLYQQGCYAGGTVLRLAK
DIAENNARGLVVAEITVICFRGPTETHLDMSGQALFGDGAGAVIVGADLDESI-ER
PIFOLWAAQTILPDSEGAIKGHLREVGLAFHLLKDVPGLISKNIKALVEAFKPIGIDD
WNSIFWAAHPGGPAILDQVESKLGKQDKLSTTRHVLSEFGNMSSACVLFILDEMRRKRM
KEGMATTGEGLEWGVLFGFGPGLTVETVMLHSVPIAN
```

- Removal of alignment columns with low abundance
- Columns with low abundance are not informative
- Reduced alignment lengths makes tree construction easier

Alignment of sequences

seq1	MDTE - KYMEKWII DQGHALFPEEDQ -
seq2	MDTDI KYSEKWI - QGHSLFPD - DQ -
seq3	METE - KYMEKWII - QGHSIFPE - DQ -
seq4	MDTEIKYMEKWII - QGHALFPD - DQP



Trimmed alignment of sequences

seq1	MDTE - KYMEKWIQGHALFPEDQ
seq2	MDTDI KYSEKWIQGHSLFPDDQ
seq3	METE - KYMEKWIQGHSIFPEDQ
seq4	MDTEIKYMEKWIQGHALFPDDQ

- Construction of phylogenetic trees is based on sequence alignments
- Similar sequences (recent shared ancestors) are grouped
- Tools: RAxML, FastTree, IQ-TREE, MEGA

Trimmed alignment of sequences

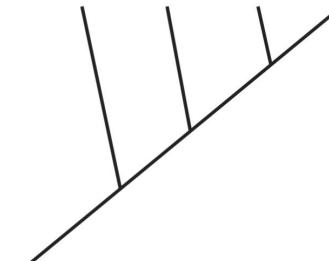
seq1	MDTE - KYMEKWIQGHALFPEDQ
seq2	MDTDI KYSEKWIQGHSLFPDDQ
seq3	METE - KYMEKWIQGHSIFPEDQ
seq4	MDTEIKYMEKWIQGHALFPDDQ

Distance calculation

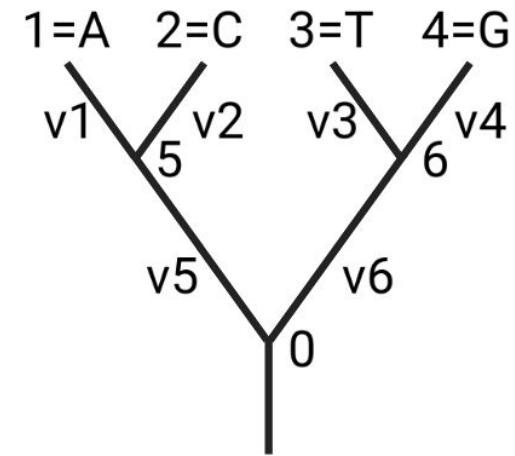
	seq1	seq2	seq3	seq4
seq1	.	5	3	2
seq2	.	.	6	3
seq3	.	.	.	5
seq4	.	.	.	.



seq2 seq3 seq1 seq4

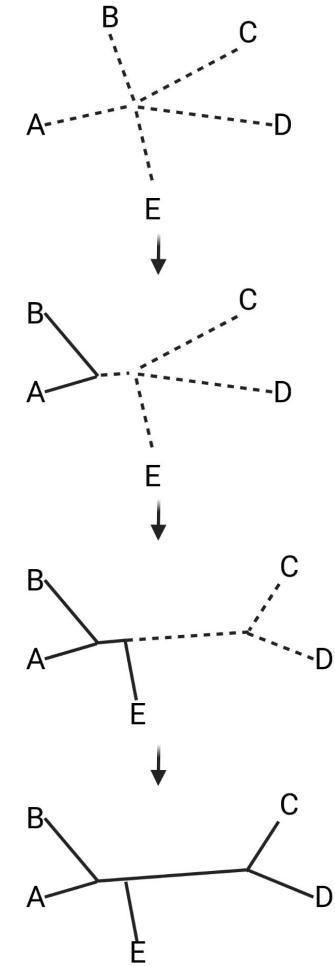


- Finding the tree with highest likelihood to explain the observed sequences
- Number of possible trees depends on number of sequences
- Often impossible to explore all possible trees (very computationally demanding)
- Tool: RAxML-NG, IQ-TREE



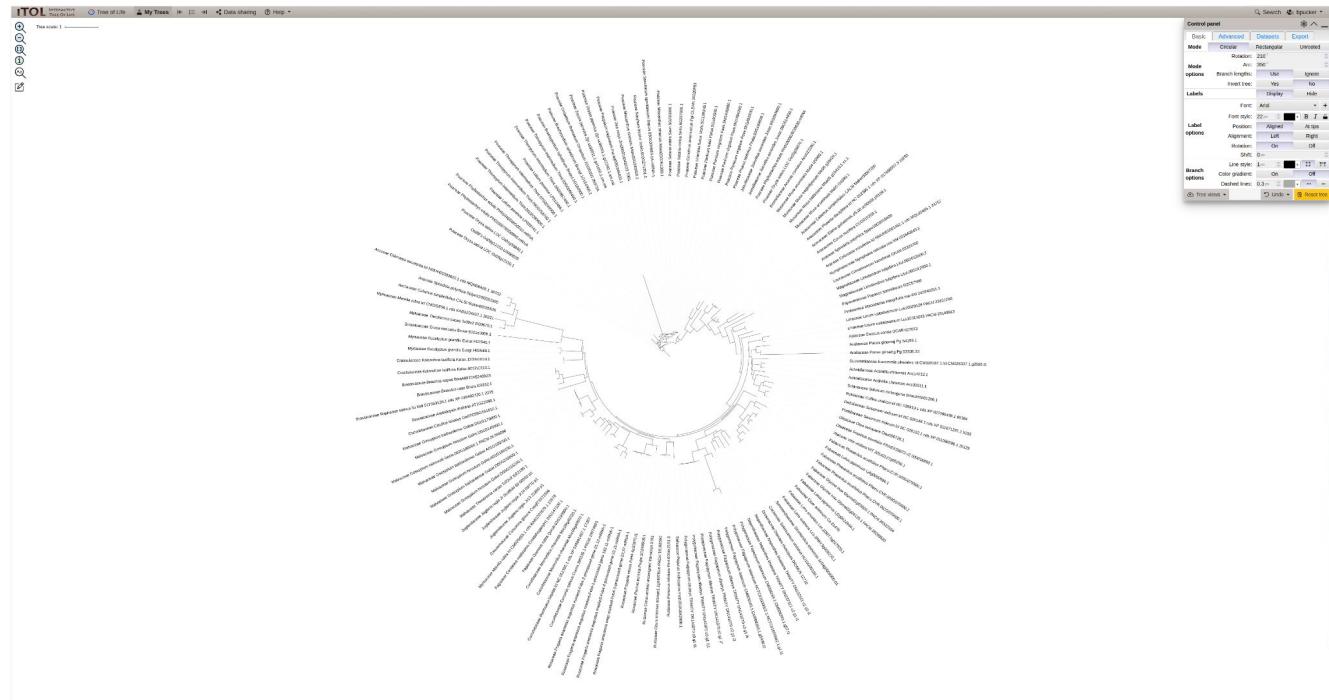
# Neighbour Joining (NJ) tree

- Combination of sequences (groups of sequences) with smallest distance
- Quick and computationally cheap approach
- Less reliable than maximum likelihood (ML) trees
- Tool: MEGA



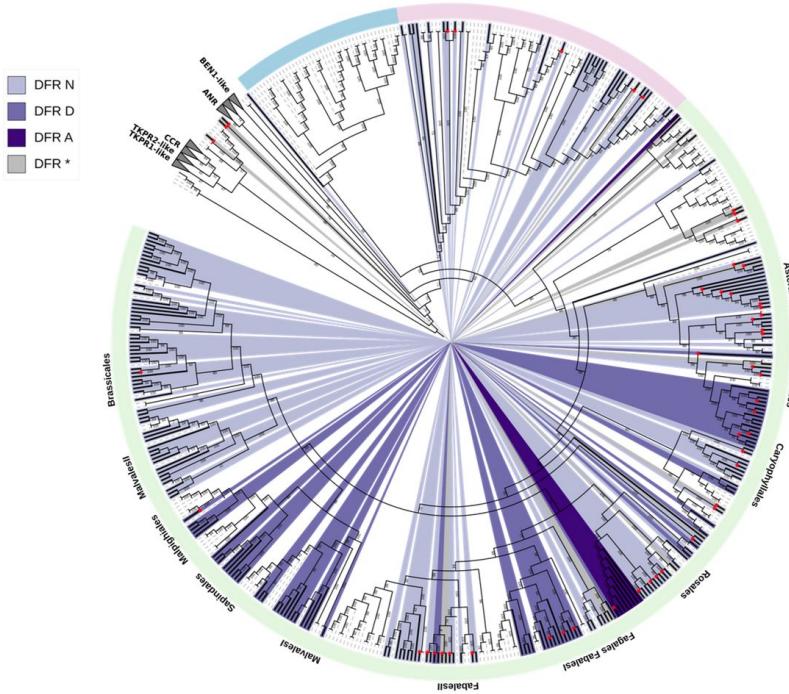
# Tree visualization

- iTOL:  
<https://itol.embl.de>
- FigTree:  
<https://tree.bio.ed.ac.uk/software/figtree/>



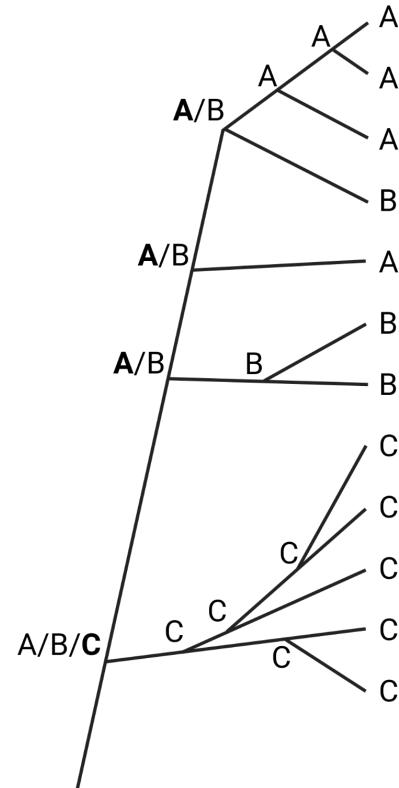
## Example: DFR evolution

- Evolution of *DFR* in land plants
- Multiple clades visible
- Sequences of closely related species cluster in lineages

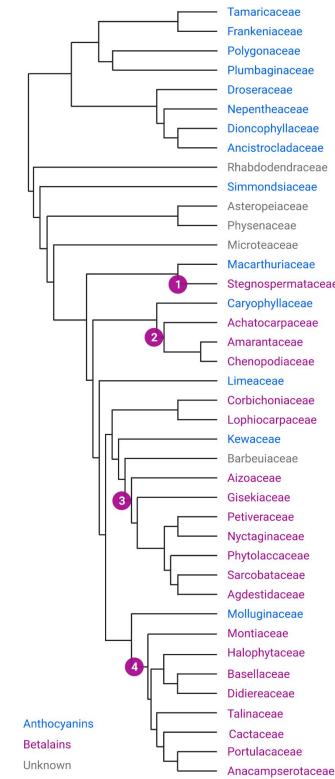
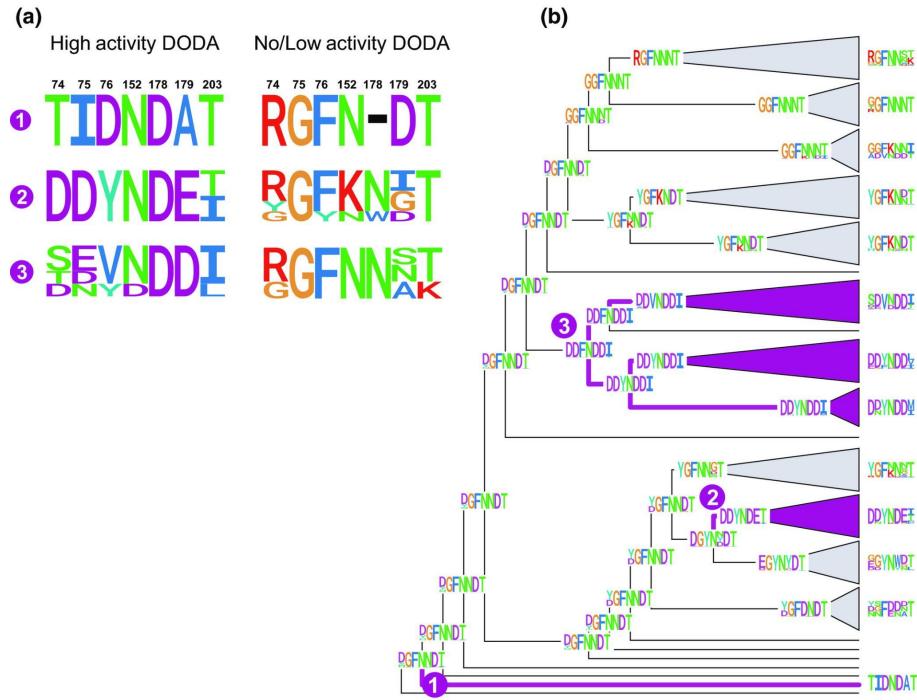


# Ancestral state reconstruction

- Identify the state of a certain sequence/trait in a last common ancestor
- Multiple states are possible in some cases
- Ancestral states are often predicted with a probability

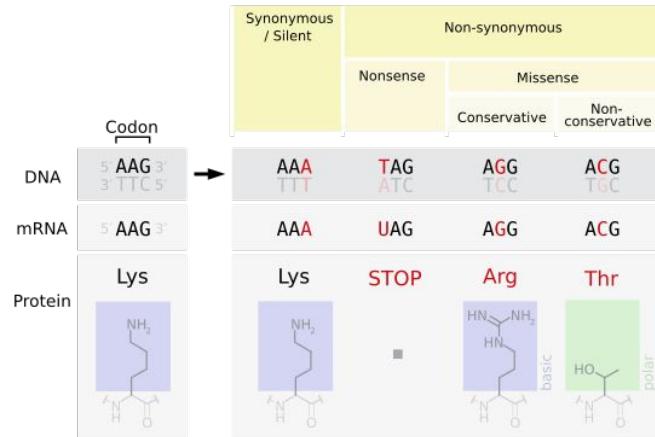


# Example: ancestral state reconstruction



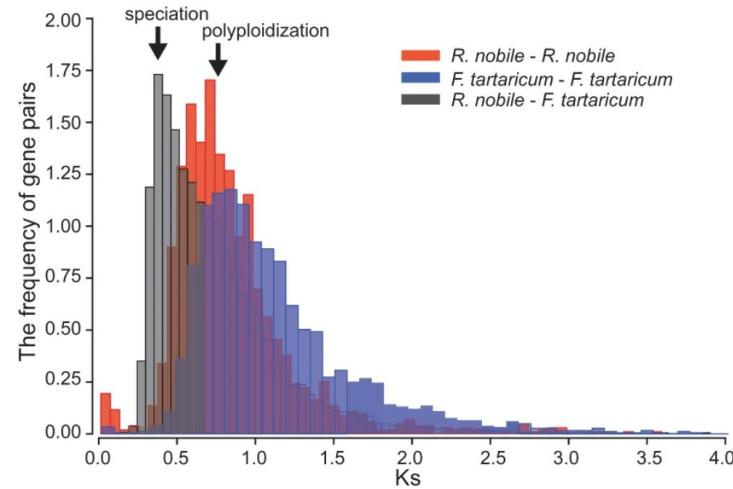
# Selection pressure analysis

- Nonsynonymous ( $K_a$ ) are compared against synonymous ( $K_s$ ) substitutions in protein-coding genes
- Interpretation:
  - $K_a/K_s < 1 \rightarrow$  purifying selection
  - $K_a/K_s = 1 \rightarrow$  neutral evolution
  - $K_a/K_s > 1 \rightarrow$  positive selection
- Requires orthologous protein-coding sequences, correct codon alignment, and sufficient evolutionary divergence for reliable estimates



# Dating evolutionary events

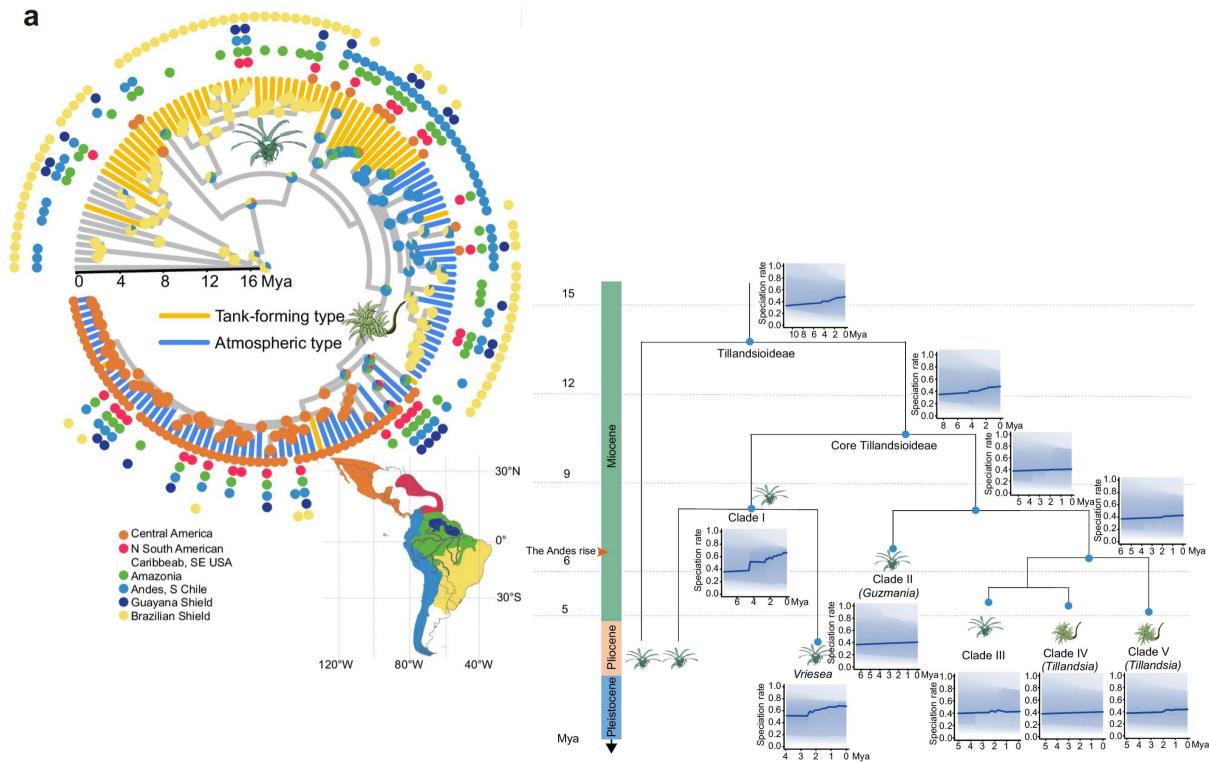
- Ks distribution peak reveal timing of duplication events (e.g., WGDs)
- Low vs. high Ks values indicate recent vs. ancient evolutionary events
- Fossil calibrations provide absolute time anchors
- Limitations: rate heterogeneity, misidentified orthologs/paralogs, and incomplete fossil records



# Evolution of plant traits

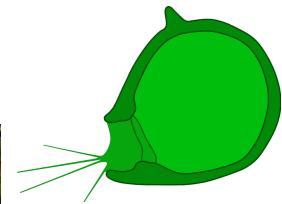
# Example: Tillandsioideae

- Tank formation to capture water and nutrients
- Absorptive trichomes on tank forming leaves
- Development of air roots



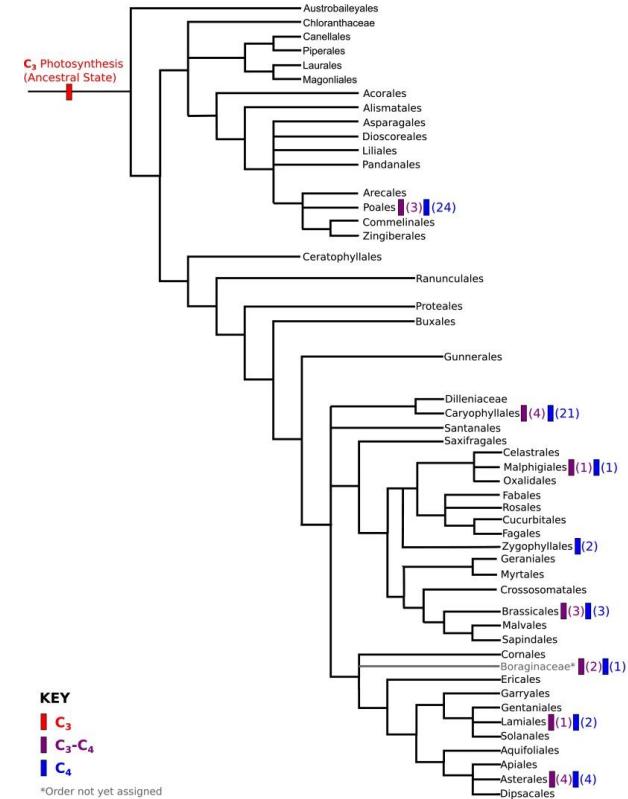
## Example: Carnivory in Lentibulariaceae

- Large and diverse carnivorous plant family (400+ species)
- Carnivory evolved to compensate for nutrient-poor habitats (low N+P)
- *Utricularia* (bladderworts): suction traps
- *Pinguicula* (butterworts): sticky leaf traps
- *Genlisea* (corkscrew plants): lobster-pot traps



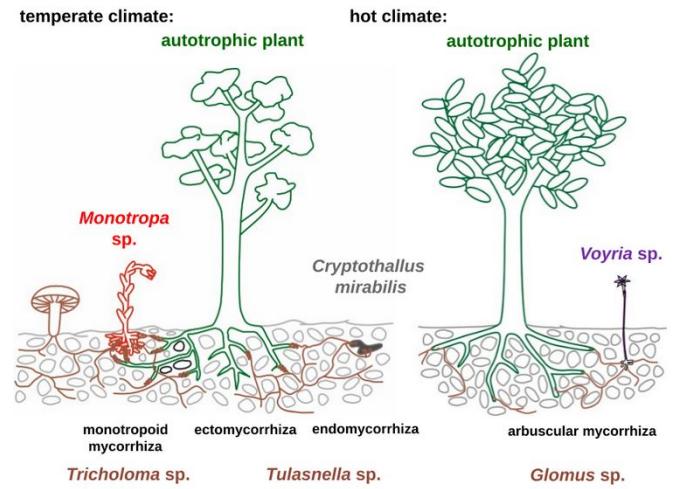
# Example: photosynthesis

- Different types of photosynthesis evolved multiple times independently
- C<sub>3</sub> was ancestral state
- C<sub>4</sub> could be more efficient under certain conditions



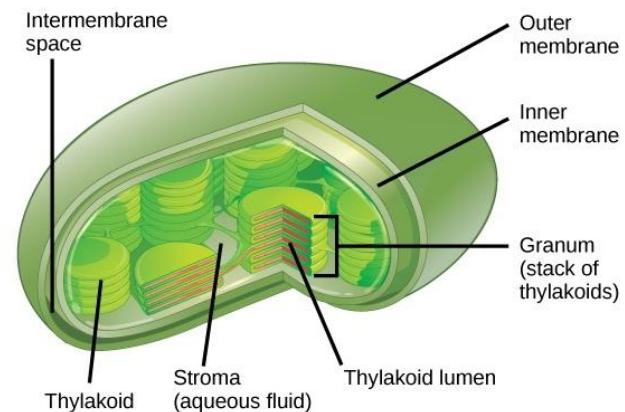
# Myco-heterotrophic plants

- Not all plants do photosynthesis  
(*Monotropha*)
- Myco-heterotrophic plants take nutrients through fungus from photoautotrophic plants

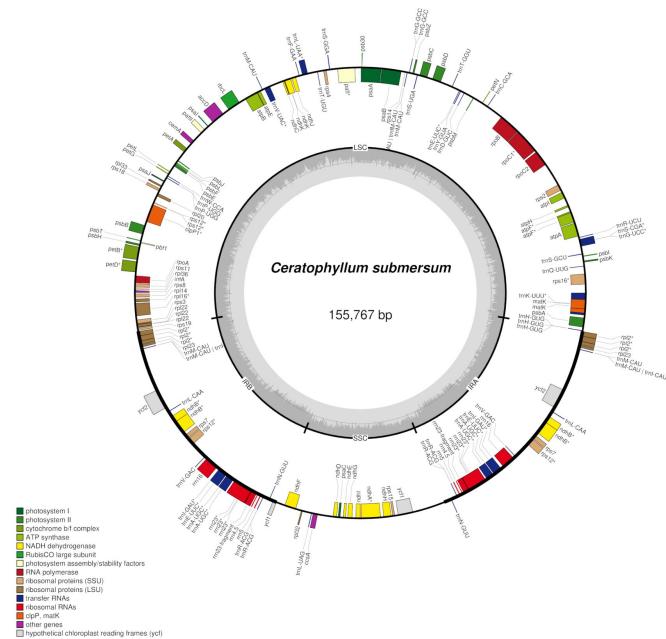


# Plastids and mitochondria

- Convert light energy into chemical energy through
- Stacks of thylakoids (grana) contain chlorophyll and other pigments essential for capturing light
- Contain circular DNA and 70S ribosomes; outer & inner membrane; prokaryotic gene structures
- Synthesize ATP, NADPH, and organic molecules like glucose, fueling plant metabolism

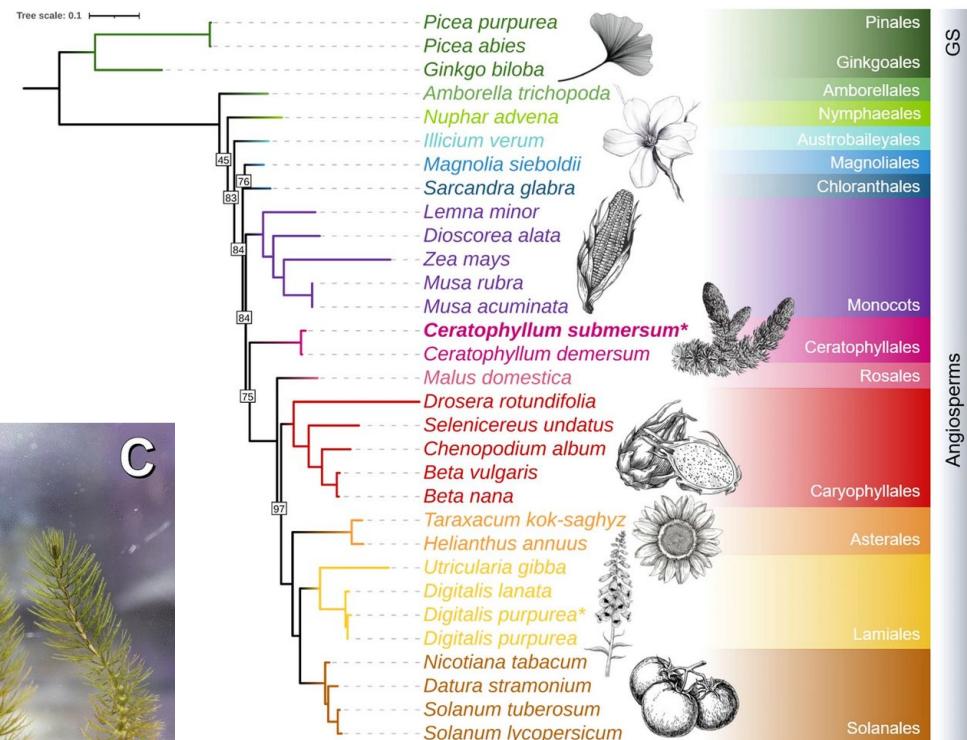


- Circular DNA found in chloroplasts and other plastids
- Typically encodes ~100–120 genes (photosynthesis, transcription/translation machinery, and metabolic pathways)
- Structure and gene order are relatively conserved across land plants
- Inherited maternally in most plants (exceptions exist)
- Target for transplastomic engineering, offering high gene expression and reduced risk of gene flow

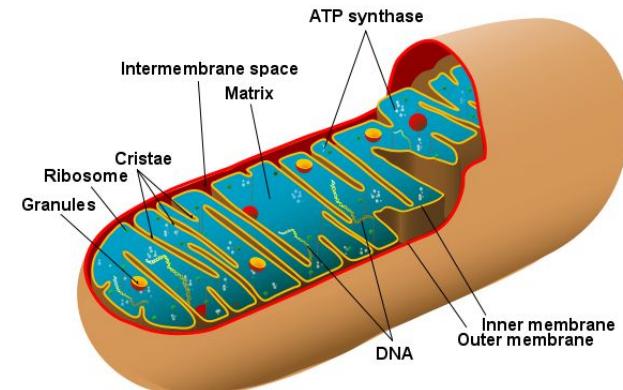


# Plastome as basis for phylogenetics

- *Ceratophyllum submersum* samples taken from local pond
- Long read sequencing resolved plastome
- Phylogenetic placement based on plastid marker genes

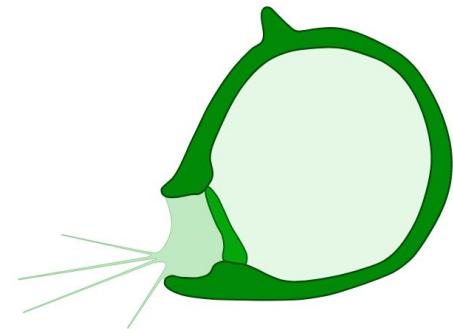


- Generate ATP through oxidative phosphorylation and electron transport chain
- Outer membrane and highly folded inner membrane (cristae) enabling compartmentalized metabolism
- Circular DNA (mtDNA) and 70S-like ribosomes
- Citric acid cycle,  $\beta$ -oxidation of fatty acids, and amino acid metabolism



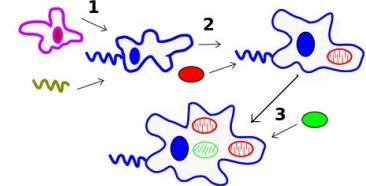
- Mitochondrial genome: large, multipartite, and highly recombinogenic DNA molecule
- Plant mitochondrial genomes are largest among eukaryotes
- Encodes core components of oxidative phosphorylation, ribosomal RNAs, tRNAs, and factors essential for mitochondrial function
- Sequence evolution is slow, yet genome structure changes rapidly
- Usually maternally inherited and valuable for studying plant evolution

- Endophytes are microorganisms (fungi/bacteria) living inside plant tissues without causing disease
- Can enhance plant growth, nutrient uptake, defence, and stress tolerance
- Important in sustainable agriculture, biocontrol, and biofertilization strategies



# Evidence for endosymbiont theory

- Mitochondria and chloroplasts retain their own circular DNA
- Genes lack introns (in many cases) and show operon-like arrangements
- Organelles contain 70S ribosomes and initiate prokaryotic translation
- Inner membranes contain bacterial-type lipids, distinct from eukaryotic membranes
- Dual-membrane structure reflects engulfment of a prokaryote by a host cell



- Evolutionary concepts: homologs, paralogs, orthologs
- Gene duplications: loss, subfunctionalization, neofunctionalization
- Phylogenetic tools: alignment, tree inference, tree visualization
- Evolution of traits in plants: carnivory, photosynthesis, pigmentation
- Plastids and mitochondria: endosymbionts

# Time for questions!

# Questions

1. What is a paralog and ortholog?
2. What is a monophyletic group?
3. What can happen to gene copies after a duplication event?
4. Provide an example for a gene duplication event in plants and explain its importance!
5. Which tool can quickly compare a sequence against a database?
6. Which tool can infer a phylogenetic tree?
7. What is a plastome?
8. Which theory explains the presence of chloroplasts in plants?