

Prof. Dr. Boas Pucker

Long Read Genomics

- **Read Mapping and Variant Calling**

Availability of slides

- All materials are freely available (CC BY) - after the lectures:
 - GitHub: <https://github.com/bpucker/LRG>
- Questions: Feel free to ask at any time
- Feedback, comments, or questions: [pucker\[a\]uni-bonn.de](mailto:pucker[a]uni-bonn.de)



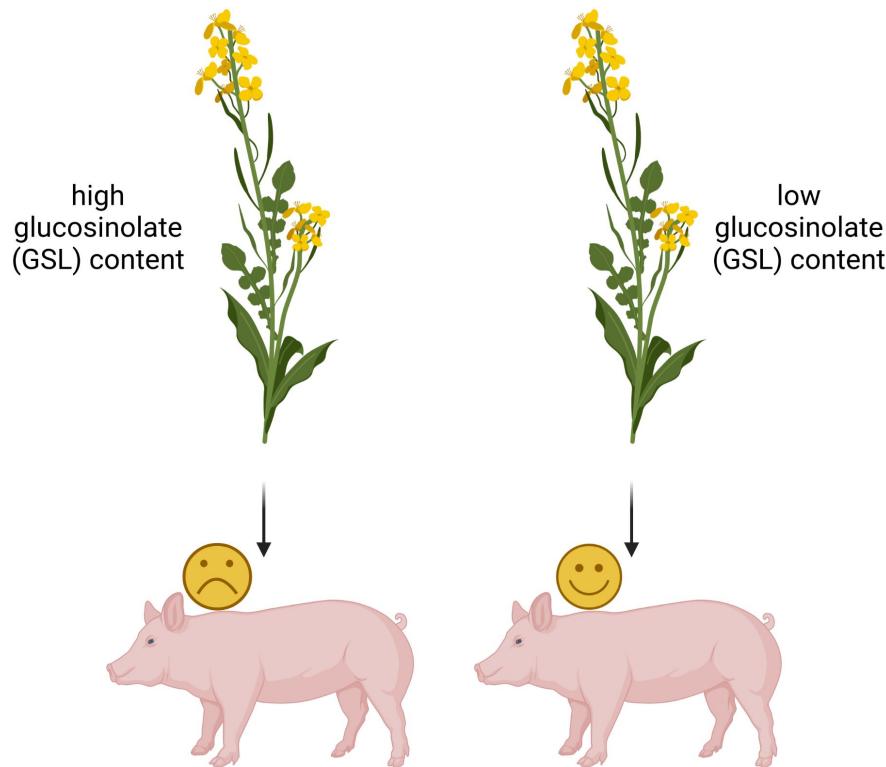
My figures and content can be re-used in accordance with CC BY 4.0, but this might not apply to all images/logos.

Repetition

- Genome sequence assembly
- Assembly quality evaluation (completeness, contiguity, correctness)
- Structural annotation (GeMoMa, BRAKER3)
 - ab initio
 - hint-based
- Functional annotation (KIPES, InterProScan5)

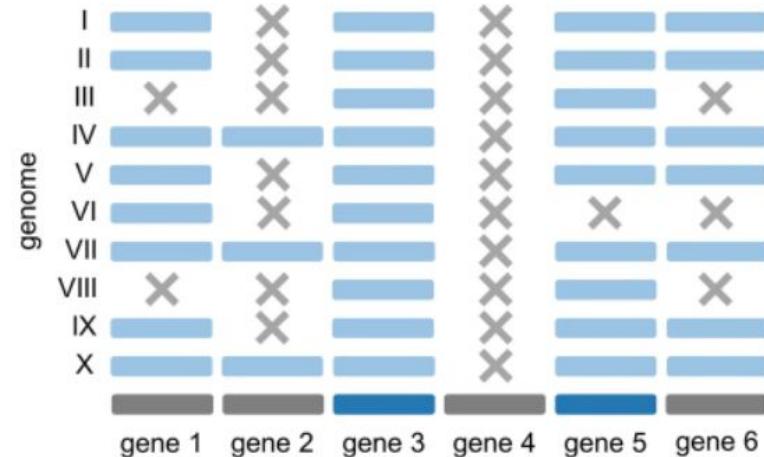
Read Mapping & Variant Calling

Question: What determines differences between plants?



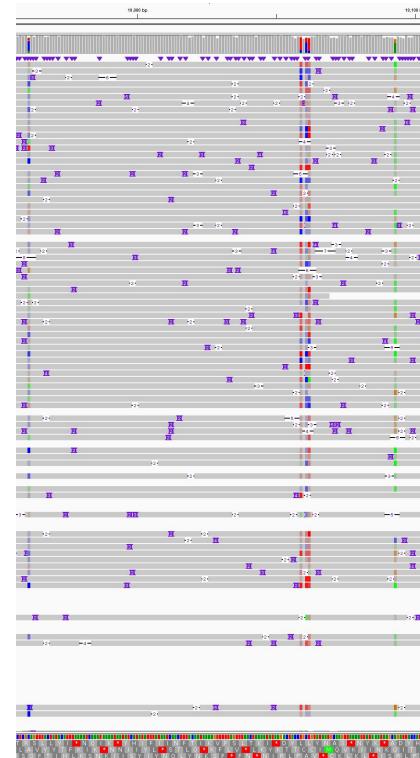
Intraspecific variation

- Different accessions of the same species have genomic differences
- Genotype determines phenotype (traits)
- Specific traits are relevant for breeding or biotechnological/agronomical applications



Differences between samples and reference

- Differences are identified against a reference genome sequence (classic approach)
- Reads of a sample are aligned against the reference (mapping)
- Differences are identified by tools (variant calling)



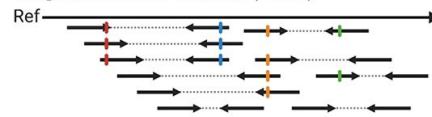
Types of sequence variants

- SNVs vs. SNPs
 - Single Nucleotide Variants = no information about their effect
 - Single Nucleotide Polymorphisms = no detrimental impact
- MNVs
 - Multiple Nucleotide Variants
- InDels
 - Insertions: additional bases in sample compared to reference
 - Deletions: loss of bases in sample compared to reference
- Inversions
 - Orientation of sequence differs between sample and reference
- Tandem duplications

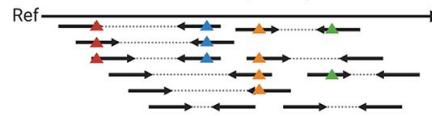
Read mapping vs. de novo genome sequence assembly

A NGS variant calling

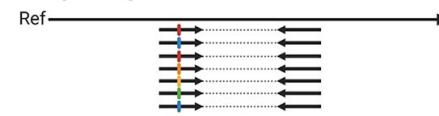
Single nucleotide variants (SNVs)



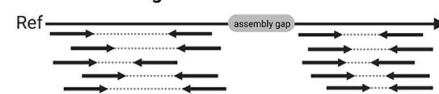
Small insertion/deletions (InDels)



Collapsed repeats



Inaccessible regions

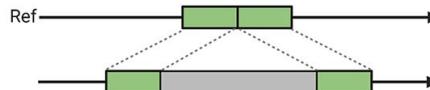


B long read variant calling

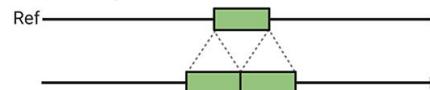
Deletion



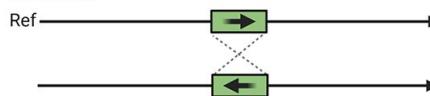
Insertion



Tandem duplication

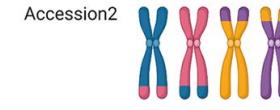
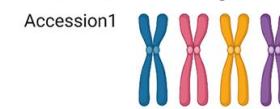


Inversion

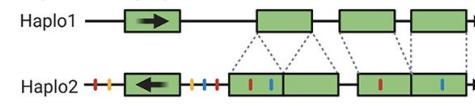


C *de novo* assembly

Chromosomal rearrangements



Separated haplotypes

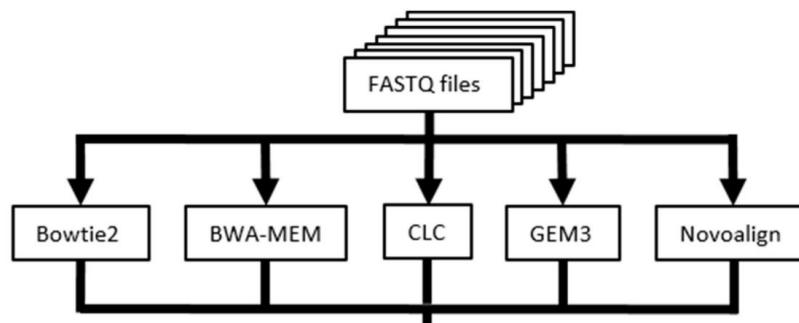


Challenges of read mappings

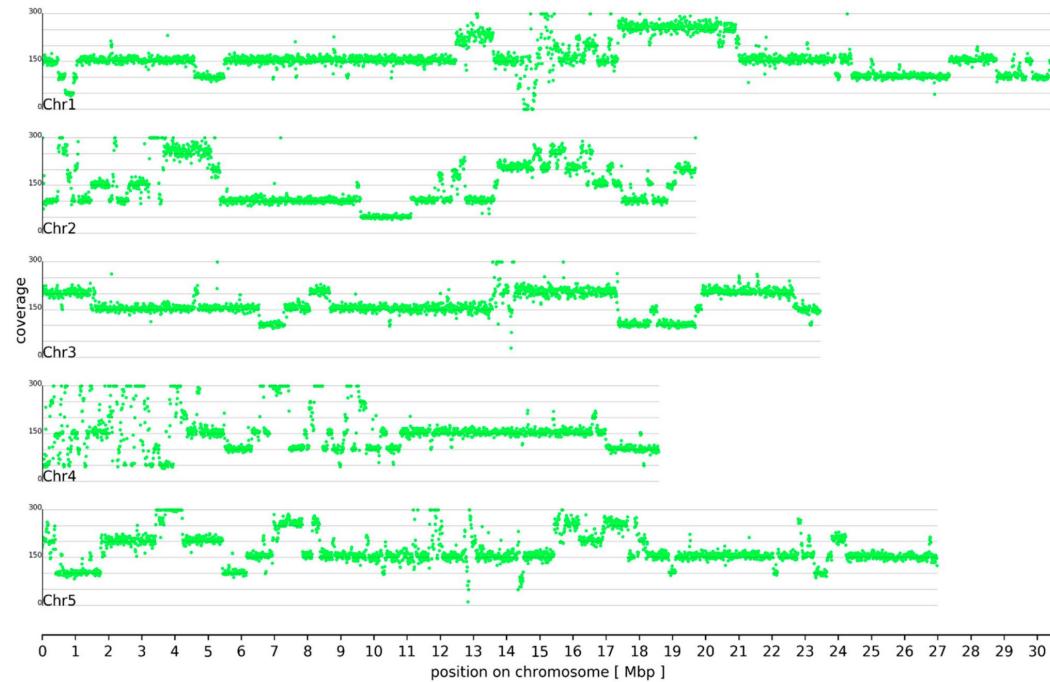
- Speed / computational costs
 - BLAST would be too slow
- Specific assignment of reads / repeats
 - (especially challenging in polyploid species)
- Splitting around InDels or across introns

Short read mappings

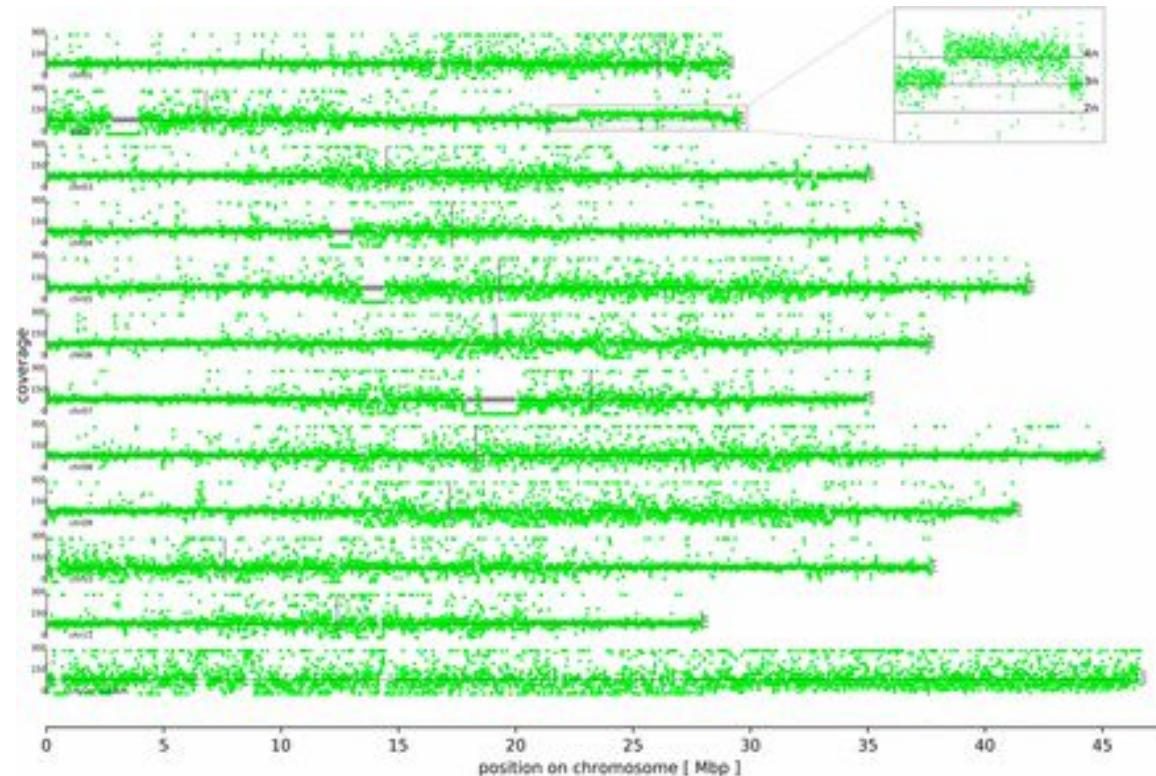
- Large number of tools for mappings
- DNA seq read mappers: BWA MEM, bowtie2
- RNA-seq read mappers (split read): STAR, HiSAT2



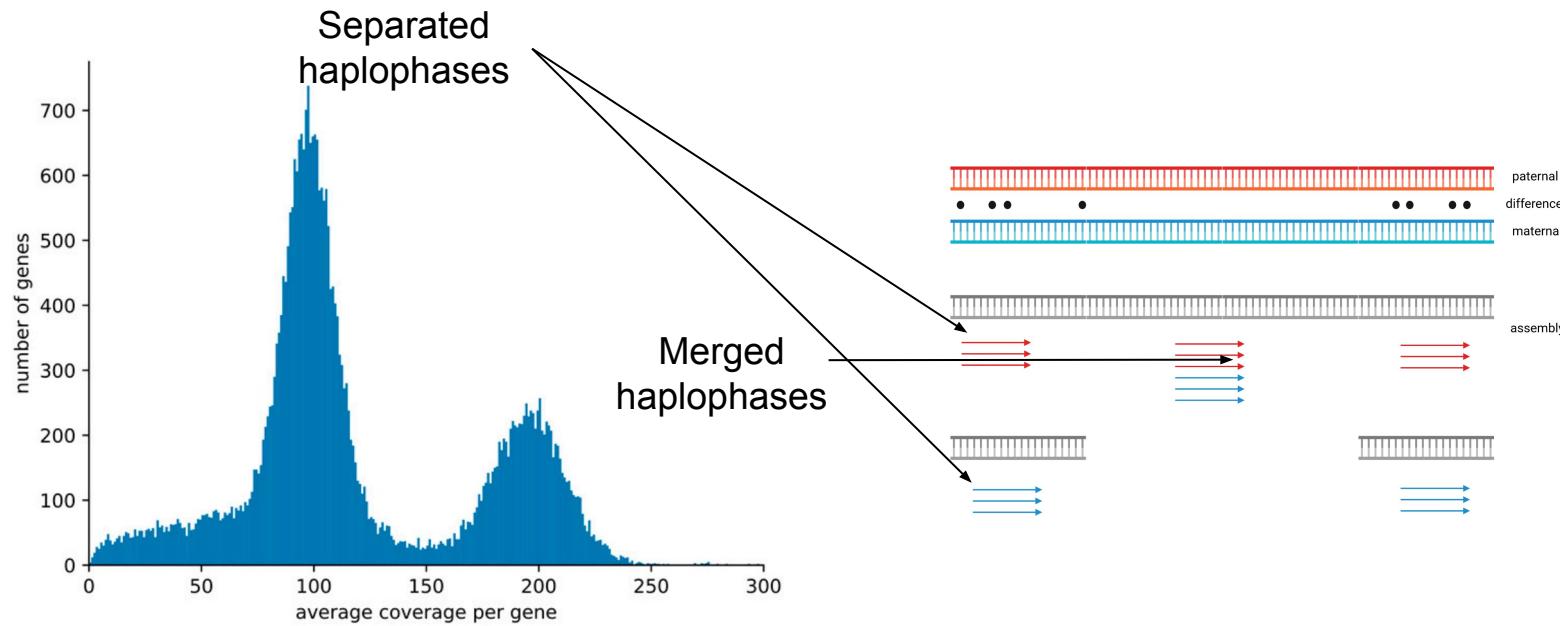
Coverage analysis: genomic copy numbers



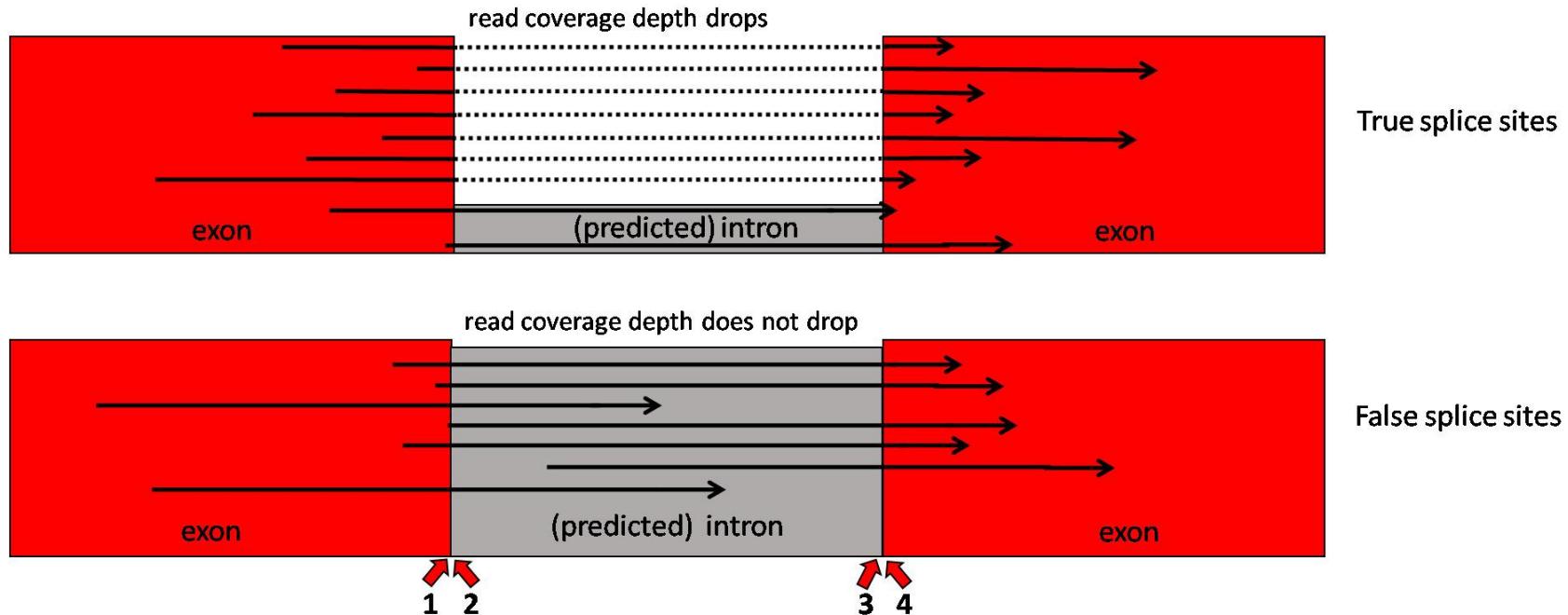
Coverage analysis: segmental duplication



Coverage analysis: haplotype phasing / ploidy



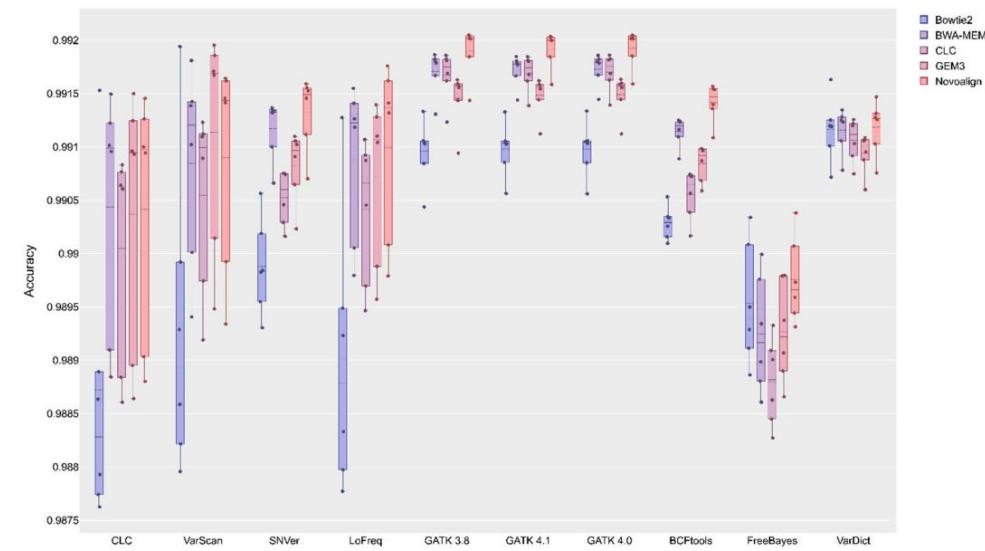
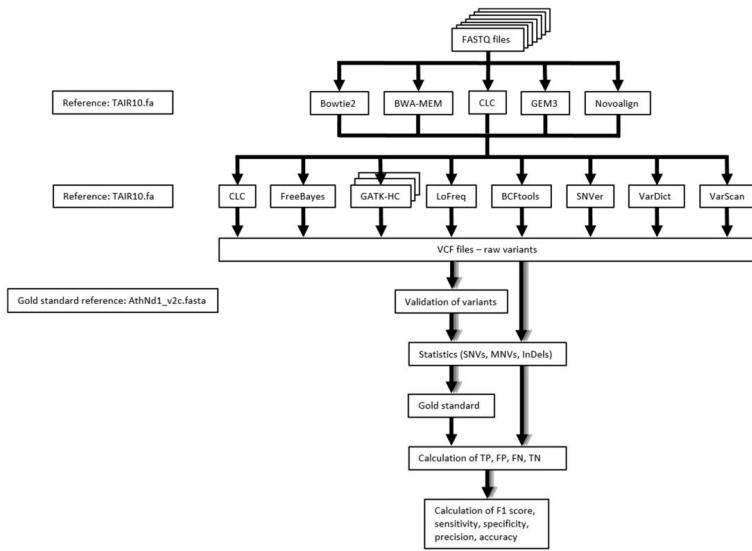
Coverage analysis: splice sites



Short read variant callers

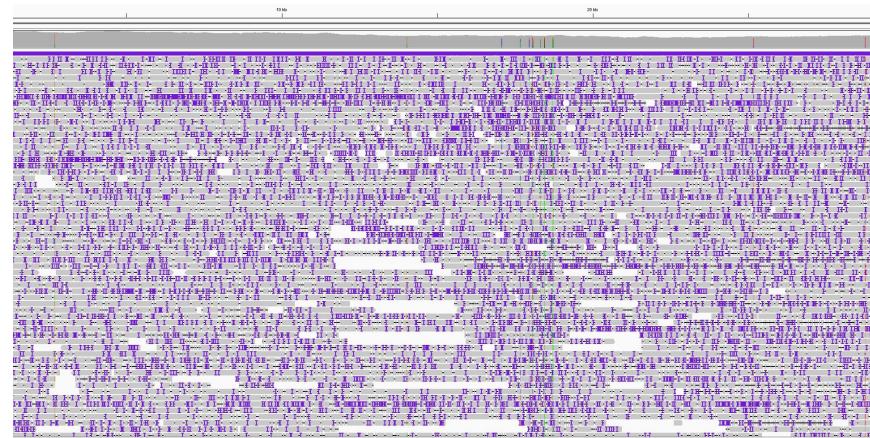
- Millions of differences between samples and reference
- Automatic detection with tools required
- Tools:
 - Genome Analysis Tool Kit (GATK): <https://gatk.broadinstitute.org/hc/en-us>
 - Samtools: <http://www.htslib.org/>
 - VarDict: <https://anaconda.org/bioconda/vardict>
 - VarScan: <https://github.com/dkoboldt/varscan>
 - SNVer: <http://snver.sourceforge.net/>
 - FreeBayes: <https://github.com/freebayes/freebayes>

Benchmarking of variant callers



Long read mapping

- Noisy alignments due to high rate of sequencing errors; InDels are particularly abundant (purple ticks)
- Rapid development of new tools
- Benchmarking studies are missing and would require frequent updates



Long read mapping tools

- Collection of long read tools: <https://long-read-tools.org/>
- Over 600 tools available
- Examples:
 - BWA-SW: <http://bio-bwa.sourceforge.net/>
 - Minimap2: <https://github.com/lh3/minimap2>
 - GraphMap: <https://github.com/isovic/graphmap>
 - NGMLR: <https://github.com/philres/ngmlr>
 - Winnowmap2: <https://github.com/marbl/Winnowmap>

Long read variant calling

- Rapid development of new tools and lack of benchmarking studies
- Tools:
 - SVIM2: <https://github.com/eldariont/svim>
 - Longshot: <https://github.com/pjedge/longshot>
 - NanoCaller: <https://github.com/WGLab/NanoCaller>
 - Sniffles2: <https://github.com/fritzsedlazeck/Sniffles>

Variant Call Format (VCF)

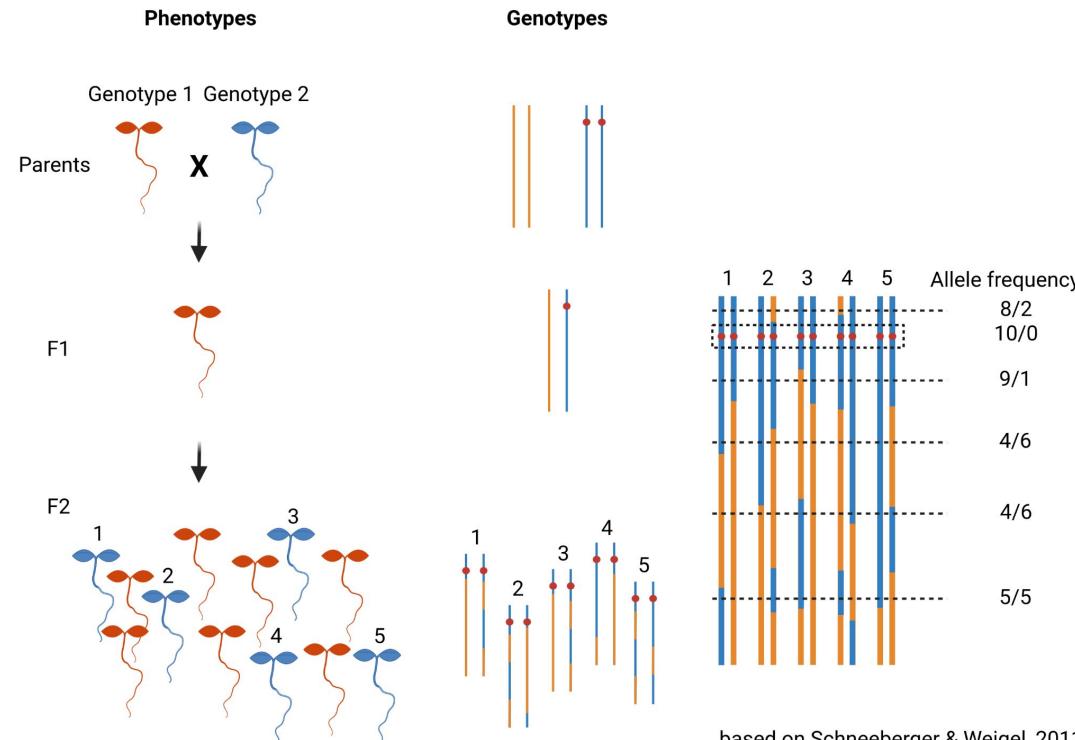
- Chromosome: name of sequence in the reference
- Position: position on the specified sequence
- ID (not relevant in plant biology): variants in humans have IDs
- Reference allele: nucleotide(s) in the reference sequence at the specified position
- Alternative allele: nucleotide(s) in the sample at the same position
- Quality: tool specific value that can be used for filtering
- Filter status: specifies filter name if variant was filtered out
- Info: collection of information
- Format: explains the fields in the following sample columns
- DATA_SET1, DATA_SET2, DATA_SET3, ...: one column per sample

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	DATA_SET	
R105	10929	.	A	G	58.54	PASS	AC<1;AF<0.500;AN<2;BaseRankSum<-1.242;DP<8;ExcessHet<=3;FS<0;0.000;MLEAC<1;MLEAF<0;580;HO<44;0.1;MRankSum<=2.109;OD<7.33;ReadPosRankSum<0.513;SOR<0.307	GT:AD:DP:GQ:PL	0/1:6;2:8;66:66,0,239	
R105	10934	.	A	G	58.54	PASS	AC<1;AF<0.500;AN<2;BaseRankSum<-1.060;DP<8;ExcessHet<=3;FS<0;0.000;MLEAC<1;MLEAF<0;580;HO<44;0.1;MRankSum<=2.109;OD<7.33;ReadPosRankSum<-1.465;SOR<0.307	GT:AD:DP:GQ:PL	0/1:6;2:8;66:66,0,239	
R105	10955	.	G	T	61.64	PASS	AC<1;AF<0.500;AN<2;BaseRankSum<-1.068;DP<7;ExcessHet<=3;FS<0;0.000;MLEAC<1;MLEAF<0;580;HO<46;16;MRankSum<=1.981;OD<8.81;ReadPosRankSum<-1.465;SOR<0.446	GT:AD:DP:GQ:PL	0/1:5;2:7;69:69,0,204	
R105	10962	.	A	G	61.64	PASS	AC<1;AF<0.500;AN<2;BaseRankSum<-1.068;DP<8;ExcessHet<=3;FS<0;0.000;MLEAC<1;MLEAF<0;580;HO<48;11;MRankSum<=1.981;OD<8.81;ReadPosRankSum<-1.465;SOR<0.446	GT:AD:DP:GQ:PL	0/1:5;2:7;69:69,0,204	
R105	13234	.	T	C	32.64	MQ_filter:0.00	filter	AC<1;AF<0..500;AN<2;BaseRankSum<-1.062;DP<35;ExcessHet<=3.0103;FS<2.782;MLEAC<1;MLEAF<0..500;HO<48;11;MRankSum<=1.981;OD<8.81;ReadPosRankSum<-1.465;SOR<0.446	GT:AD:DP:GQ:PL	0/1:5;2:7;69:69,0,204
R105	13240	.	T	C	70.64	MQ_filter:0.00	filter	AC<1;AF<0..500;AN<2;BaseRankSum<-1.062;DP<35;ExcessHet<=3.0103;FS<2.782;MLEAC<1;MLEAF<0..500;HO<48;11;MRankSum<=1.981;OD<8.81;ReadPosRankSum<-1.465;SOR<0.446	GT:AD:DP:GQ:PL	0/1:5;2:7;69:69,0,204
R105	13295	.	T	A	183.64	PASS	AC<1;AF<0..500;AN<2;BaseRankSum<-1.979;DP<19;ExcessHet<=3;FS<0;0.000;MLEAC<1;MLEAF<0;580;HO<43;94;MRankSum<=1.631;OD<9.67;ReadPosRankSum<1.186;SOR<0.3126	GT:AD:DP:GQ:PL	0/1:13;6;19;99:191,0,413	
R105	13538	.	T	C	61.64	PASS	AC<1;AF<0.500;AN<2;BaseRankSum<1.359;DP<27;ExcessHet<=3;FS<0;0.000;MLEAC<1;MLEAF<0;580;HO<46;51;MRankSum<=1.847;OD<9.769;ReadPosRankSum<1.189;SOR<0.769	GT:AD:DP:GQ:PL	0/1:21;4:25;69:69,0,660	
R105	14511	.	C	T	39.64	PASS	AC<1;AF<0..500;AN<2;BaseRankSum<-0.539;DP<13;ExcessHet<=3.0103;FS<2.522;MLEAC<1;MLEAF<0..500;HO<43;42;MRankSum<=1.771;OD<9.42;ReadPosRankSum<-0.217;SOR<1.546	GT:AD:DP:GQ:PL	0/1:9;4:13;47:47,0,352	
R105	14542	.	G	A	15.64	PASS	AC<1;AF<0.500;AN<2;BaseRankSum<-1.125;DP<10;ExcessHet<=3.0103;FS<0;0.000;MLEAC<1;MLEAF<0;580;HO<43;30;MRankSum<=1.843;OD<9.44;ReadPosRankSum<0.319;SOR<1.022	GT:AD:DP:GQ:PL	0/1:10;3:15;59:59,0,209	
R105	14547	.	A	T	156.64	PASS	AC<1;AF<0..500;AN<2;BaseRankSum<-0.899;DP<43;ExcessHet<=3;FS<0;0.000;MLEAC<1;MLEAF<0;580;HO<46;40;MRankSum<=2.000;OD<9.92;ReadPosRankSum<0.303;SOR<1.022	GT:AD:DP:GQ:PL	0/1:11;4:15;59:59,0,1161	
R105	15712	.	C	T	60.64	MQ_filter:0.00	filter	AC<1;AF<0..500;AN<2;BaseRankSum<-1.209;DP<35;ExcessHet<=3.0103;FS<2.026;MLEAC<1;MLEAF<0;580;HO<39;58;MRankSum<=2.129;OD<9.08;ReadPosRankSum<=0.994;SOR<1.236	GT:AD:DP:GQ:PL	0/1:38;4;34;78:0,0,1247
R105	15719	.	C	T	64.64	MQ_filter:0.00	filter	AC<1;AF<0..500;AN<2;BaseRankSum<-0.166;DP<37;ExcessHet<=3.0103;FS<2.037;MLEAC<1;MLEAF<0;580;HO<34;75;MRankSum<=0.303;OD<7.23;ReadPosRankSum<=0.712;SOR<1.802	GT:AD:DP:GQ:PL	0/1:5;2:7;58:58,0,171
R105	16037	.	C	A	58.64	MQ_filter:0.00	filter	AC<1;AF<0..500;AN<2;BaseRankSum<-0.030;DP<37;ExcessHet<=3.0103;FS<2.043;MLEAC<1;MLEAF<0;580;HO<34;80;MRankSum<=0.366;OD<7.23;ReadPosRankSum<=0.722;SOR<1.802	GT:AD:DP:GQ:PL	0/1:5;2:7;58:58,0,171
R105	16042	.	C	T	58.64	PASS	AC<1;AF<0..500;AN<2;BaseRankSum<-0.030;DP<37;ExcessHet<=3.0103;FS<2.043;MLEAC<1;MLEAF<0;580;HO<34;80;MRankSum<=0.366;OD<7.23;ReadPosRankSum<=0.722;SOR<1.802	GT:AD:DP:GQ:PL	0/1:5;2:7;58:58,0,171	
R105	17718	.	A	G	79.64	MQ_filter:0.00	filter	AC<1;AF<0..500;AN<2;BaseRankSum<-0.722;DP<37;ExcessHet<=3.0103;FS<2.043;MLEAC<1;MLEAF<0;580;HO<29;81;MRankSum<=0.858;OD<8.42;ReadPosRankSum<0.099;SOR<0.041	GT:AD:DP:GQ:PL	0/1:13;5;18;87:87,0,392
R105	37909	.	A	G	51.64	MQ_filter:0.00	filter	AC<1;AF<0..500;AN<2;BaseRankSum<-0.549;DP<2.74;ReadPosRankSum<0.549;SOR<0.206	GT:AD:DP:GQ:PL	0/1:7;2:9;59:59,0,254
R105	45724	.	C	CTTATA	61.60	PASS	AC<1;AF<0.500;AN<2;BaseRankSum<-0.210;DP<7;ExcessHet<=3.0103;FS<0;0.000;MLEAC<1;MLEAF<0;580;HO<43;83;MRankSum<=0.967;OD<8.80;ReadPosRankSum<=0.210;SOR<0.44	GT:AD:DP:GQ:PL	0/1:5;2:7;69:69,0,204	

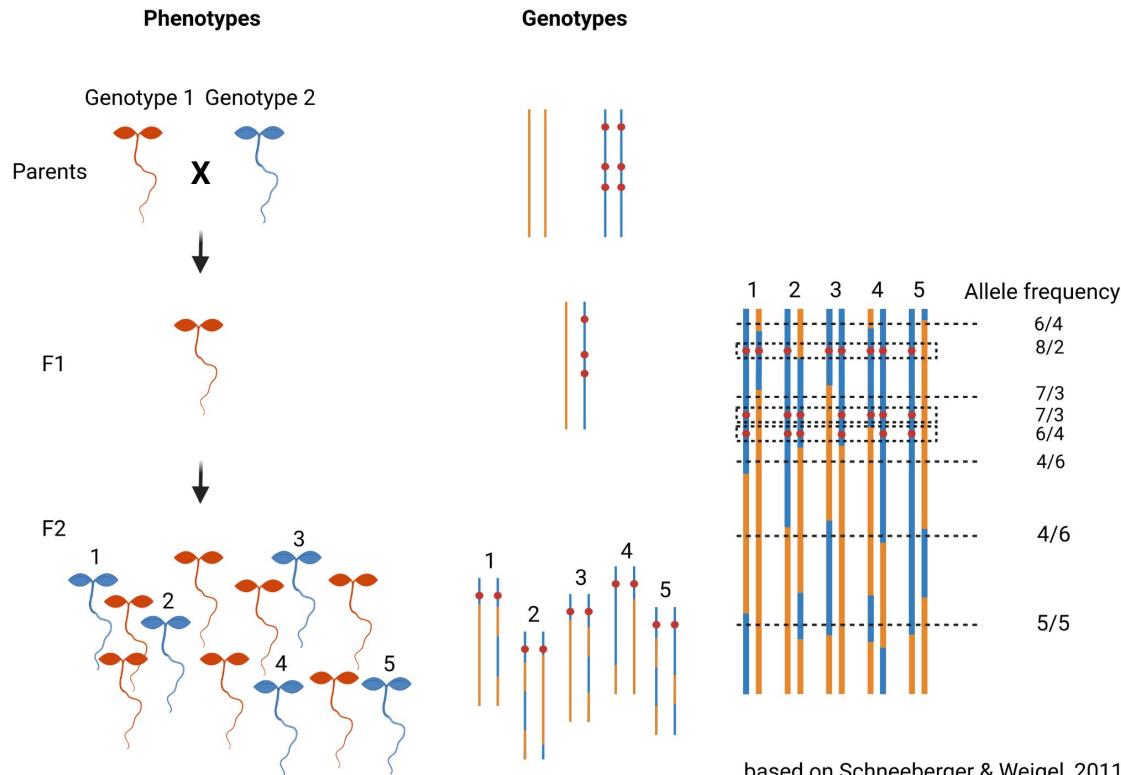
genomic Variant Call Format (gVCF)

- Generally similar to VCF file
- One entry for each position in the reference
- Not just variants, but also non-variants are listed

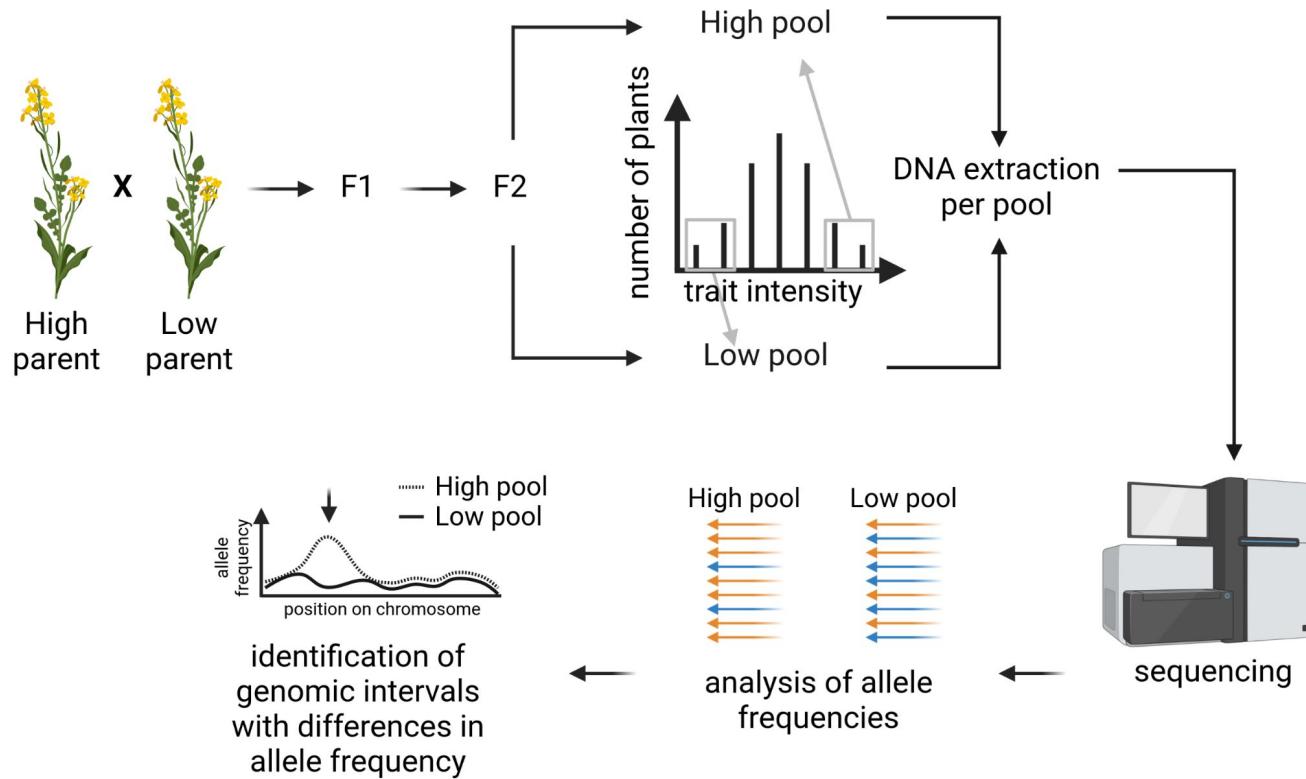
Mapping-by-sequencing (1)



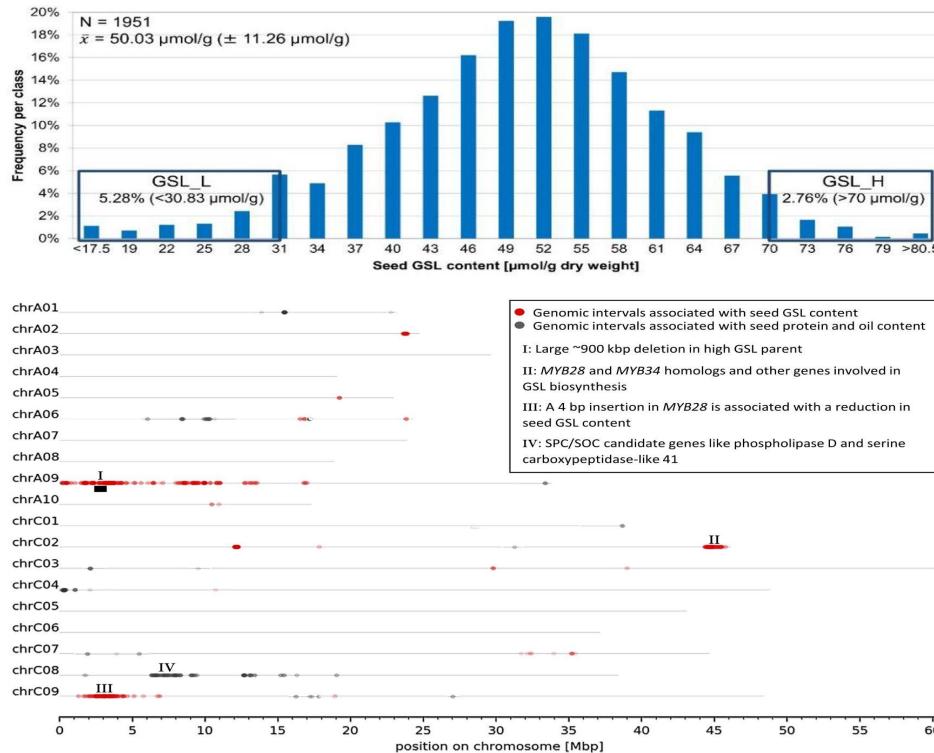
Mapping-by-sequencing (2)



Mapping-by-sequencing (3)



Mapping-by-sequencing (4)



Variant annotation

- Variants can have functional consequences
- Variants in coding sequences can change amino acids
- Variants outside CDS can have regulatory consequences

Variant Annotation - SnpEff

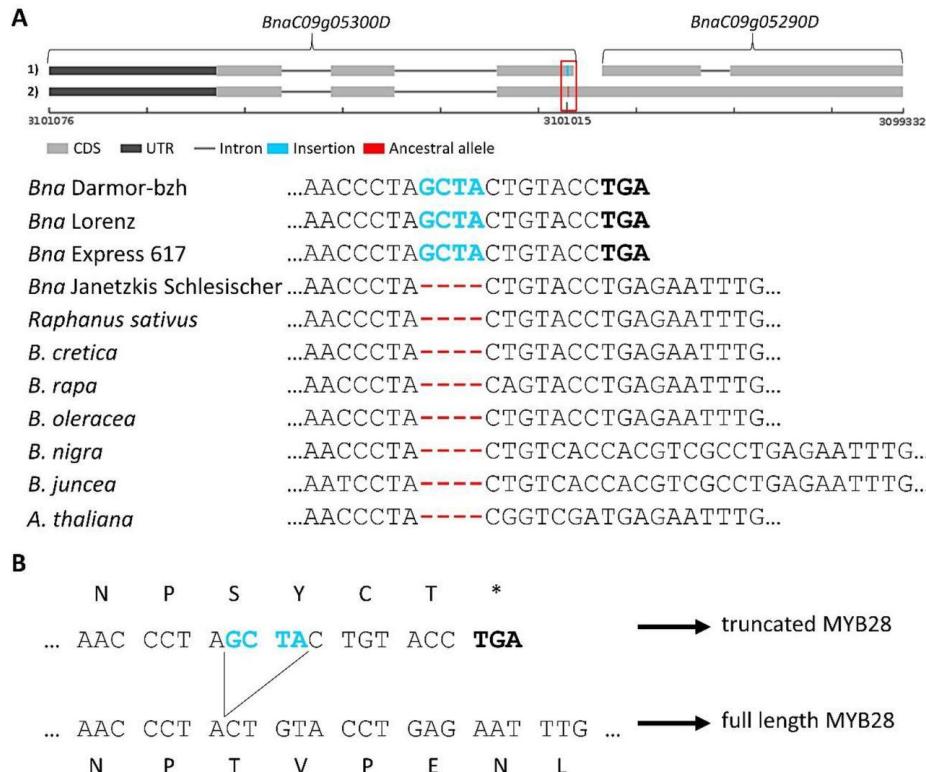
- High throughput annotation of SNVs and larger variants
- Reference sequence (FASTA) and structural annotation (GFF) required
- Freely available software: <http://pcingola.github.io/SnpEff/>
- Addition of one annotation field to VCF file

chrA01 179964 . T A . PASS	GSL_High,GSL_Low;ANN=A stop_gained HIGH Gene_BnaA01g00340D Gene_BnaA01g00340D transcript BnaA01g00340D Coding 5/9 c.863T>A p.Leu288* 863/1834 863/1834 288/610 GT:AD:DP:QQ:PL
chrA01 179967 . T C . PASS	GSL_High,GSL_Low;ANN=C mssense_variant MODERATE Gene_BnaA01g00340D Gene_BnaA01g00340D transcript BnaA01g00340D Coding 5/9 c.866T>C p.Phe289Ser 866/1834 866/1834 289/610 GT:AD:DP:QQ:PL
chrA01 179979 . C T . PASS	GSL_Low;ANN=T mssense_variant MODERATE Gene_BnaA01g00340D Gene_BnaA01g00340D transcript BnaA01g00340D Coding 5/9 c.878C>T p.Thr293Met 878/1834 878/1834 293/610 GT:AD:DP:QQ:PL
chrA01 179988 . A T . PASS	GSL_Low;ANN=T mssense_variant MODERATE Gene_BnaA01g00340D Gene_BnaA01g00340D transcript BnaA01g00340D Coding 5/9 c.887A>T p.Tyr296Phe 887/1834 887/1834 296/610 GT:AD:DP:QQ:PL
chrA01 180003 . G A . PASS	GSL_Low;ANN=A mssense_variant MODERATE Gene_BnaA01g00340D Gene_BnaA01g00340D transcript BnaA01g00340D Coding 5/9 c.902G>A p.Gly301Asp 902/1834 902/1834 301/610 GT:AD:DP:QQ:PL
chrA01 180024 . A G . PASS	GSL_High,GSL_Low;ANN=G mssense_variant MODERATE Gene_BnaA01g00340D Gene_BnaA01g00340D transcript BnaA01g00340D Coding 5/9 c.923A>G p.His308Arg 923/1834 923/1834 308/610 GT:AD:DP:QG:PL
chrA01 180039 . T C . PASS	GSL_High,GSL_Low;ANN=C mssense_variant MODERATE Gene_BnaA01g00340D Gene_BnaA01g00340D transcript BnaA01g00340D Coding 5/9 c.938T>C p.Ile313Thr 938/1834 938/1834 313/610 GT:AD:DP:QG:PL
chrA01 180048 . G A . PASS	GSL_High,GSL_Low;ANN=A stop_retained_variant LOW Gene_BnaA01g00340D Gene_BnaA01g00340D transcript BnaA01g00340D Coding 5/9 c.947G>A p.Ter316Ter 947/1834 947/1834 316/610 GT:AD:DP:QG:PL

SnpEff predictions

Type	Meaning	Example
SNP	Single Nucleotide Polymorphism	Reference = 'A', Sample = 'C'
Ins	Insertion	Reference = 'A', Sample = 'AGT'
Del	Deletion	Reference = 'AC', Sample = 'C'
MNP	Multiple Nucleotide Polymorphism	Reference = 'ATA', Sample = 'GTC'
MIXED	MNP and InDel	Reference = 'ATA', Sample = 'GTCAGT'

Example: Glucosinolate biosynthesis in rapeseed



Is there something that SnpEff might miss?

Variant Annotation - NAVIP

- NAVIP = Neighborhood-Aware Variant Impact Predictor
- Variants can have interacting effects e.g. compensate each other
- Considering multiple variants while predicting effects
- Tool freely available:
<https://github.com/bpucker/NAVIP>

wild type:	GAT TCA AGA AGA ATG
peptide:	D S R R M
variant 1:	GAT TCA TGA AGA ATG (STOP)
peptide:	D S * R M
variant 2:	GAT TCA AGC AGA ATG (amino acid substitution)
peptide:	D S S R M
combined:	GAT TCA TGC AGA ATG (amino acid substitution)
peptide:	D S C R M

A>T A>C
↓ ↓

Connected SNVs

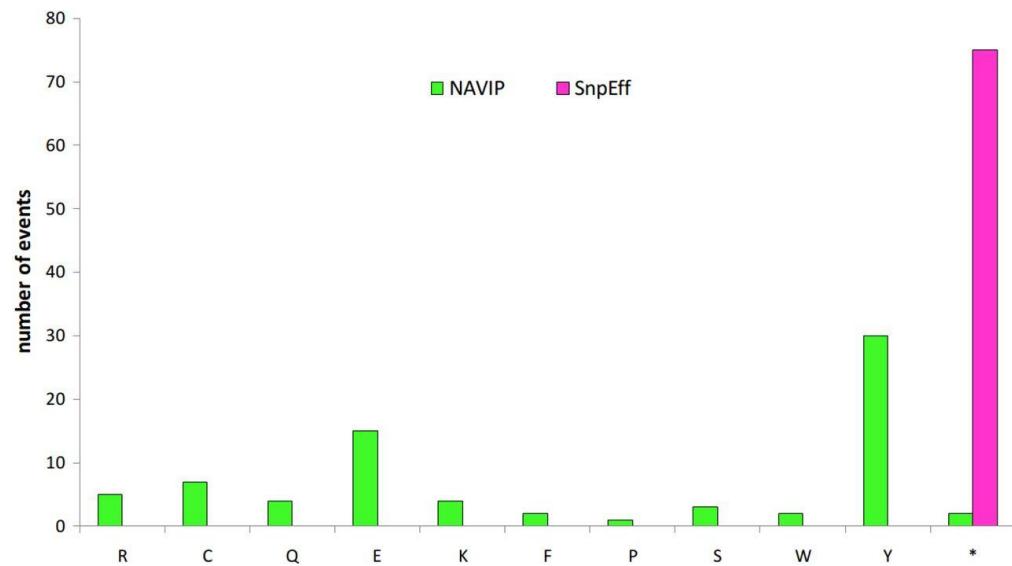
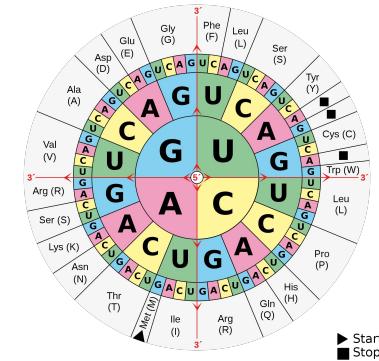
wild type:
peptide: GAT TCA **AGA** AGA ATG
D S R R M

A>T
A>C

variant 1:
peptide: GAT TCA **TGA** AGA ATG (STOP)
D S * R M

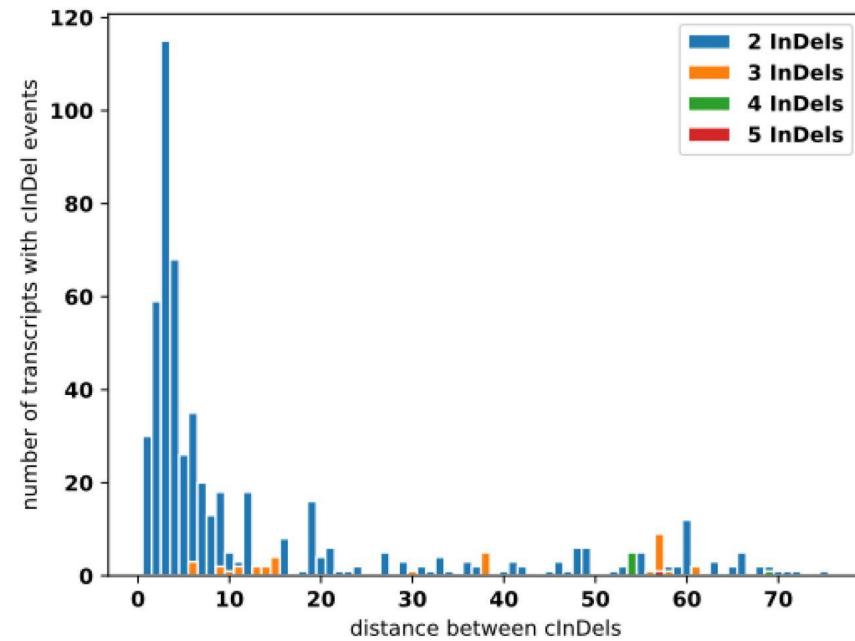
variant 2:
peptide: GAT TCA **AGC** AGA ATG (amino acid substitution)
D S **S** R M

combined:
peptide: GAT TCA **TGC** AGA ATG (amino acid substitution)
D S **C** R M



Compensating InDels

wild type:	GTG TAT CTG CGC ATT	T>TA	C>CGC
peptide:	V Y L R I		
variant 1:	GTG TAT ACT GCG CAT T	(frameshift)	
peptide:	V Y T A H		
variant 2:	GTG TAT CTG CGC GCA TT	(frameshift)	
peptide:	V Y L R A		
combined:	GTG TAT ACT GCG CGC ATT		
peptide:	V Y L A R I		

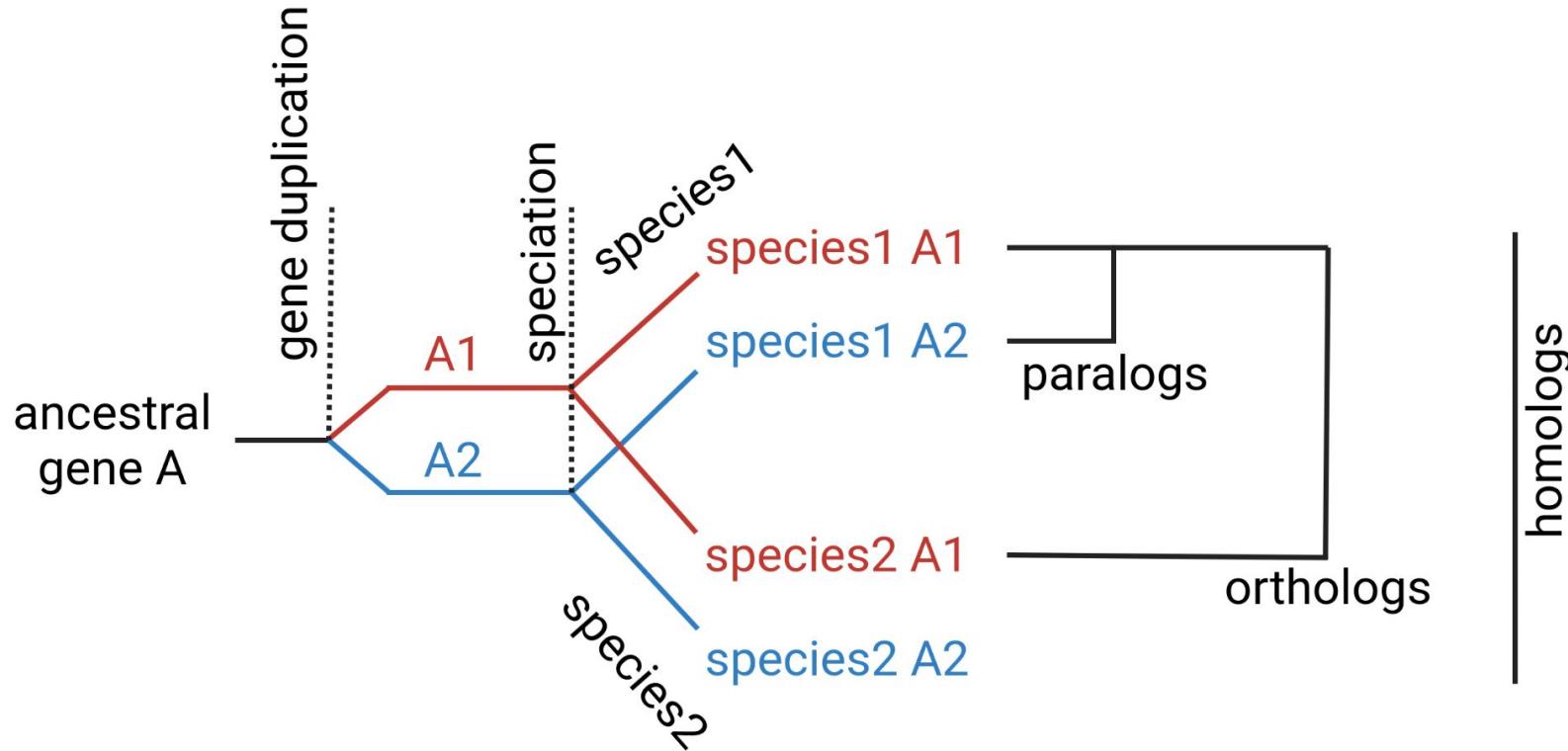


Evolutionary Genomic

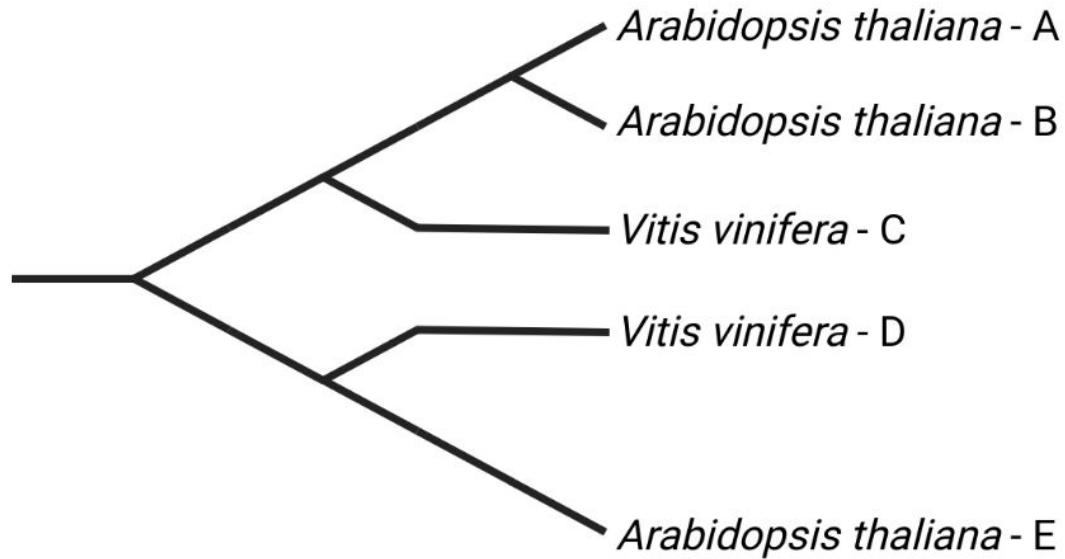
**“Nothing in biology makes sense
except in the light of evolution”**

Theodosius Dobzhansky

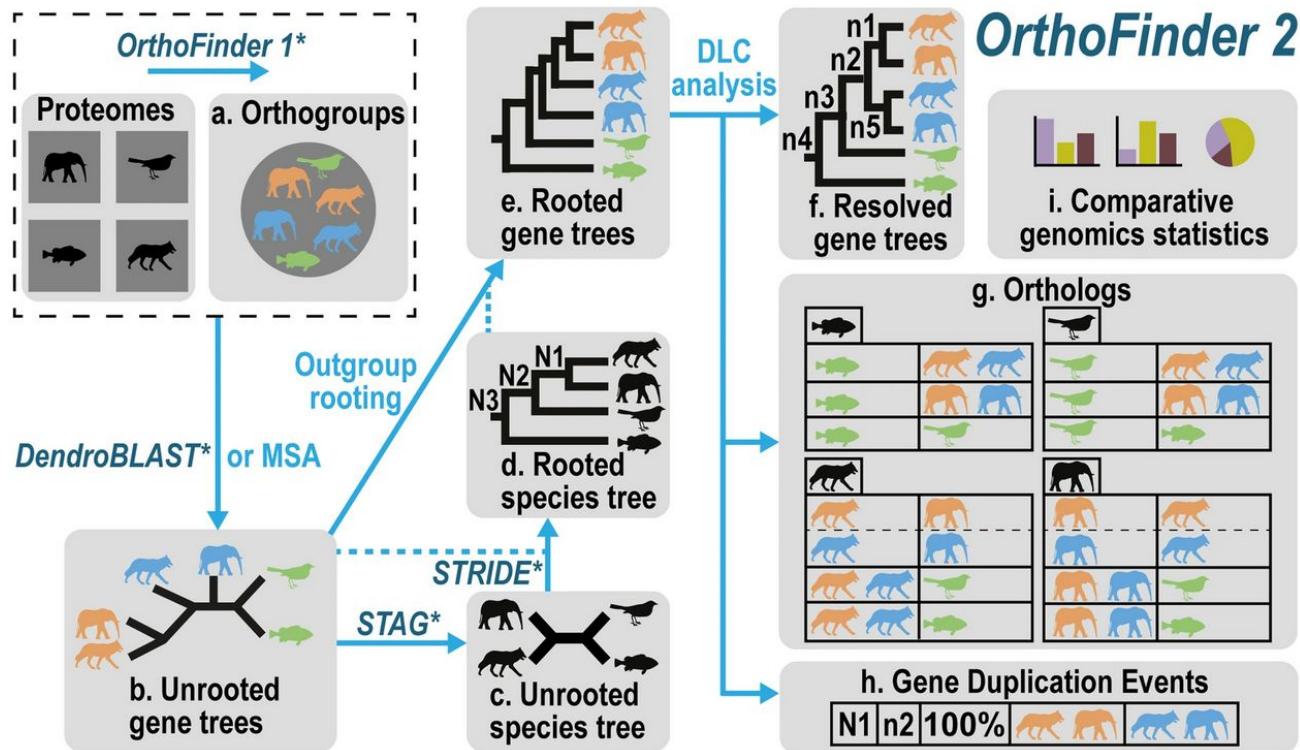
Paralogs & orthologs



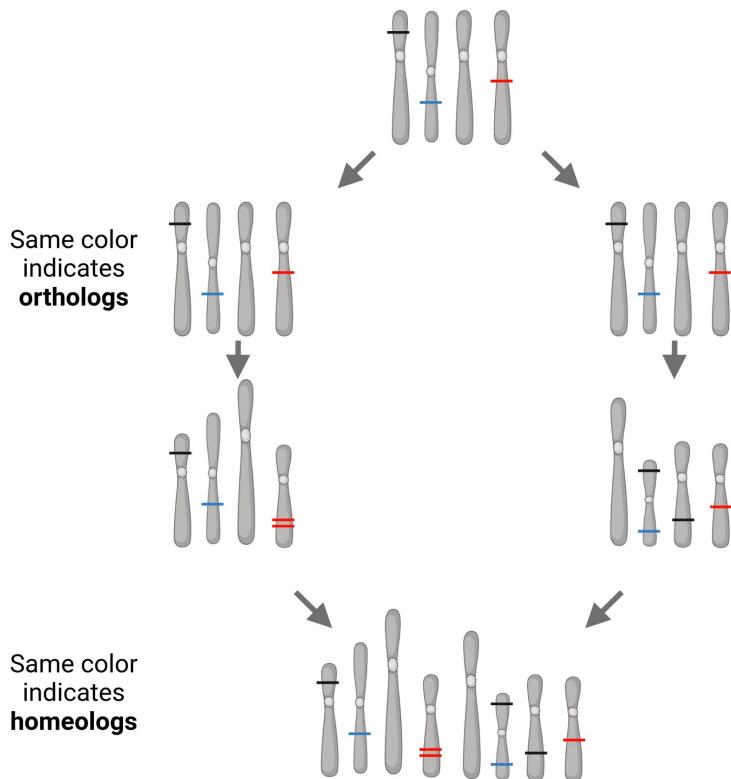
Which are orthologs/paralogs?



OrthoFinder2



Homeologs



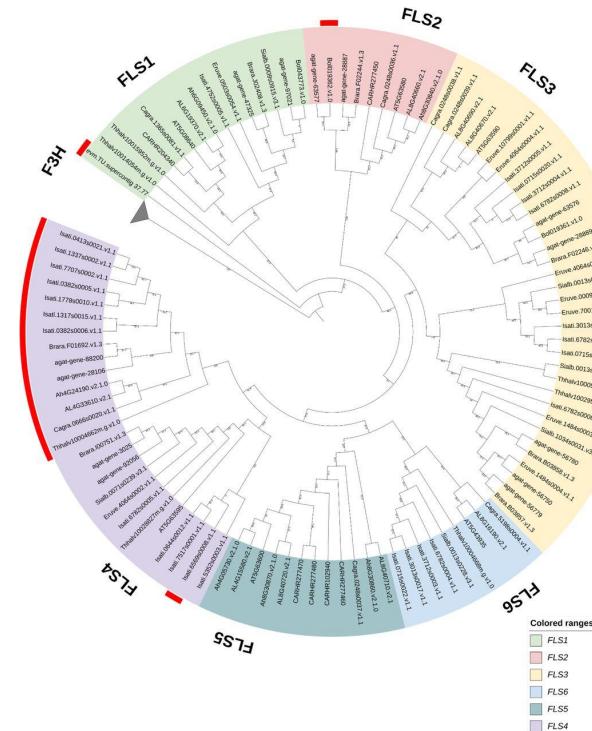
Speciation

Gene duplications,
translocations, and
rearrangements

Hybridization leads to
polyploidization

DupyliCate - investigation of gene duplications

- Python tool for investigation of gene copies
- Classification of duplicates
- Investigation of sub-/neofunctionalization at expression level



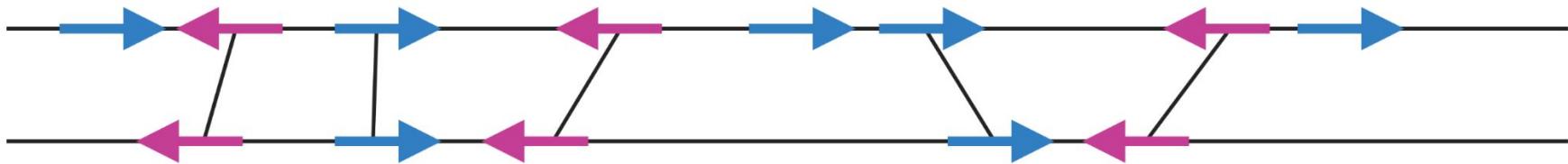
Comparative Genomics

Comparative genomics

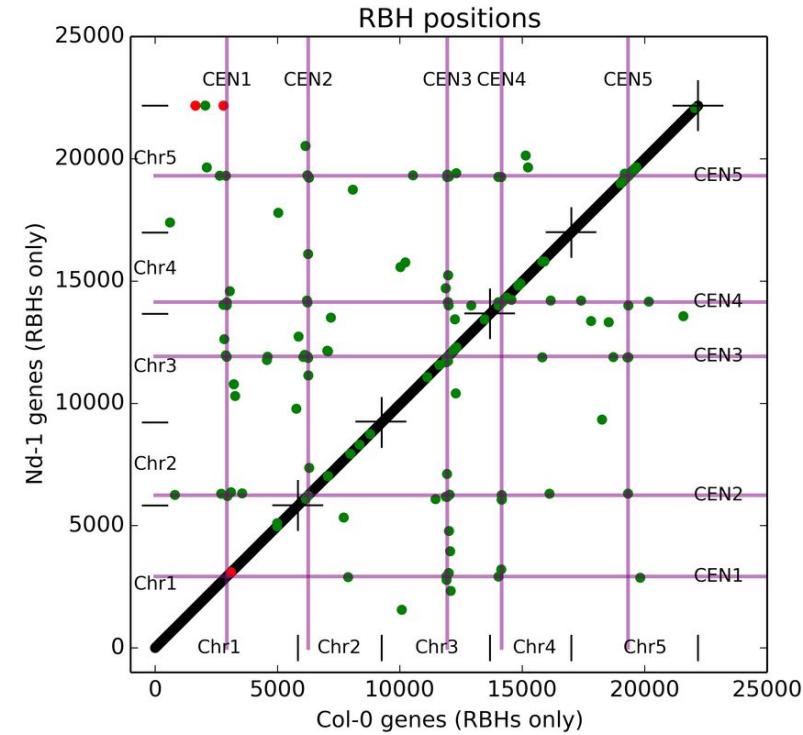
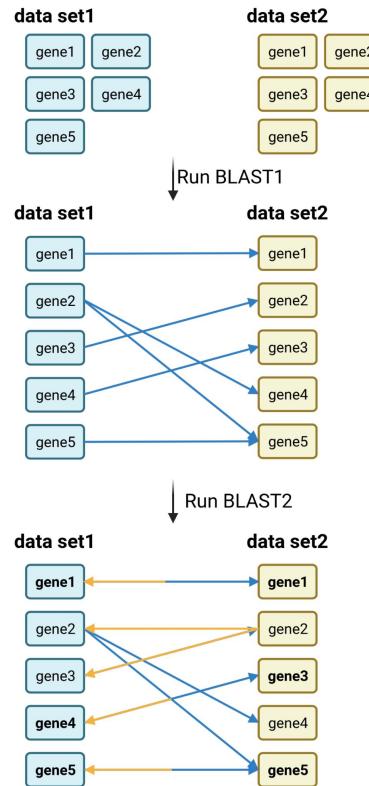
- Compare genome sizes
- Compare chromosome numbers
- Compare gene sets
- Compare gene positions & chromosome structures

Synteny

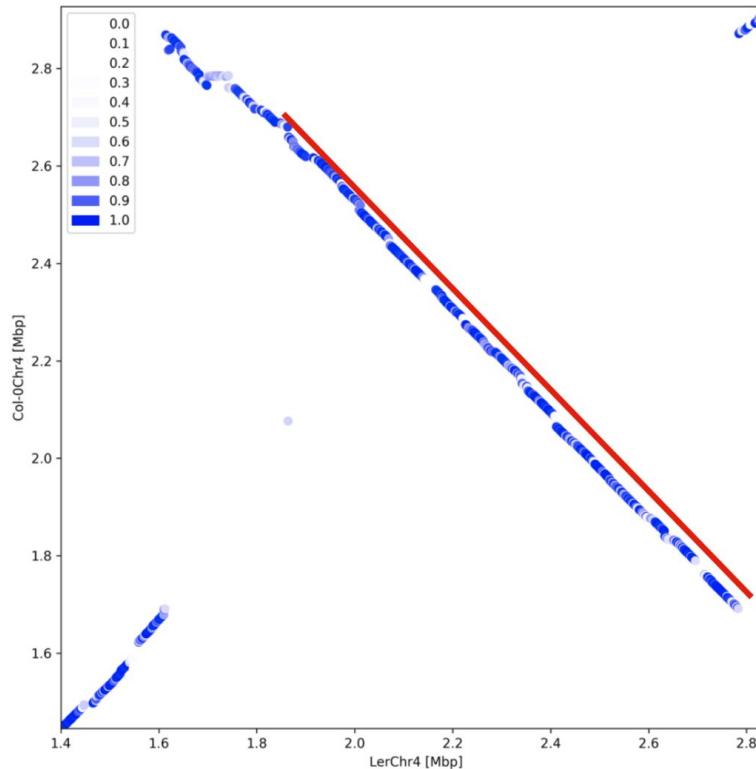
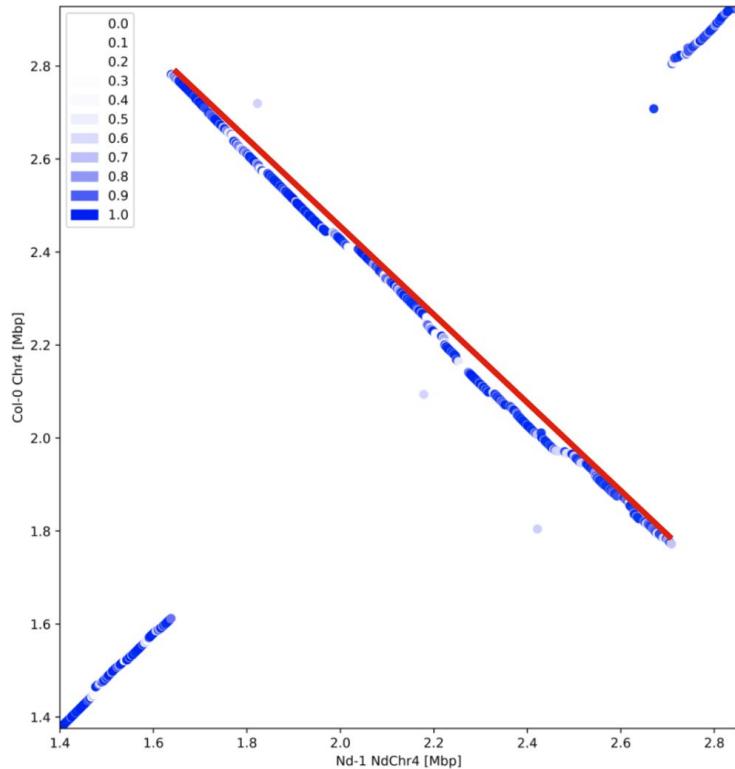
Synteny = Same order of genes in different genomes



Reciprocal Best BLAST hits (RBHs)



Dot plots



- 1 Extraction of mRNA sequences based on genomic positions of genes

species1 ——————

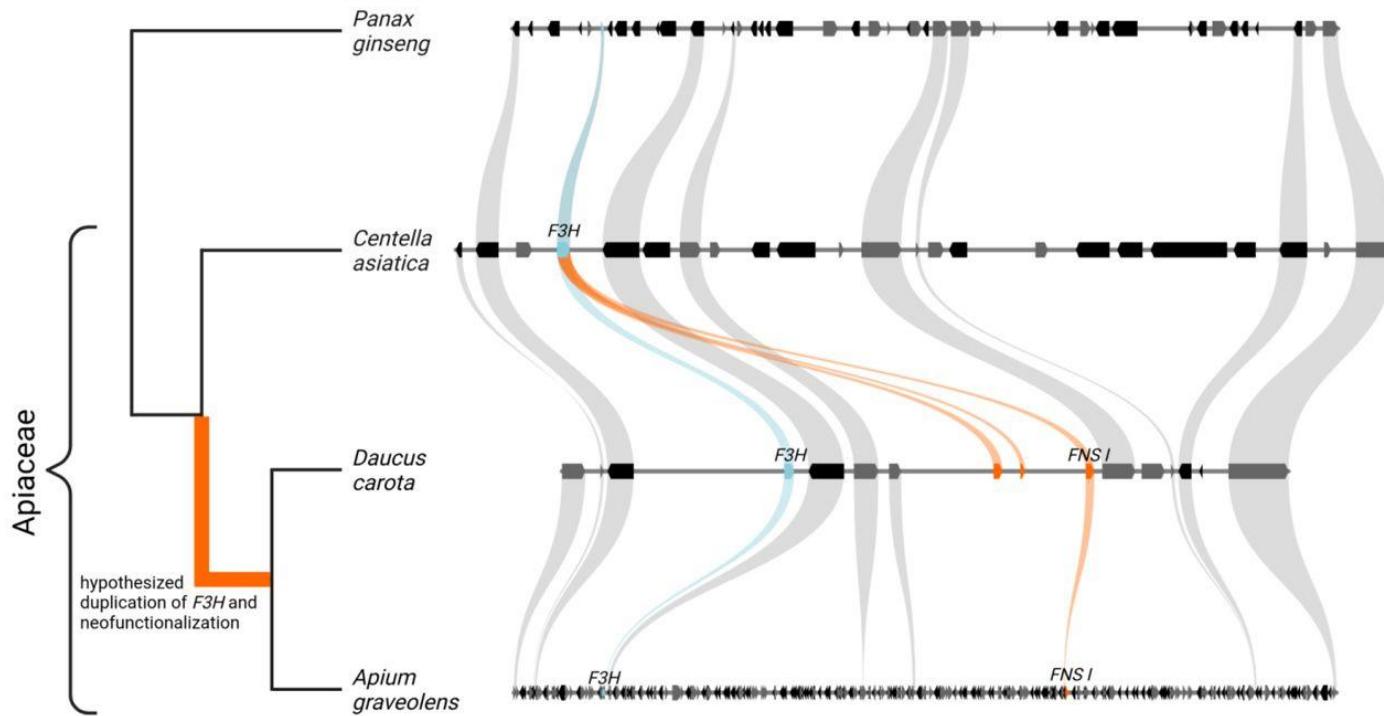
- 2 Comparison of concatenated mRNA sequences via BLAST

species2 ——————

- 3 Identification of syntenic blocks

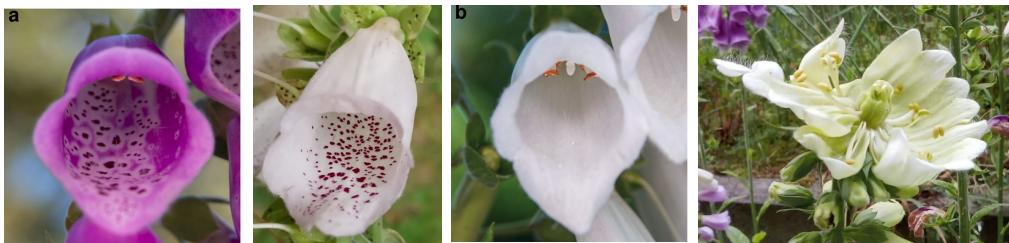


Evolution of FNS I in the Apiaceae

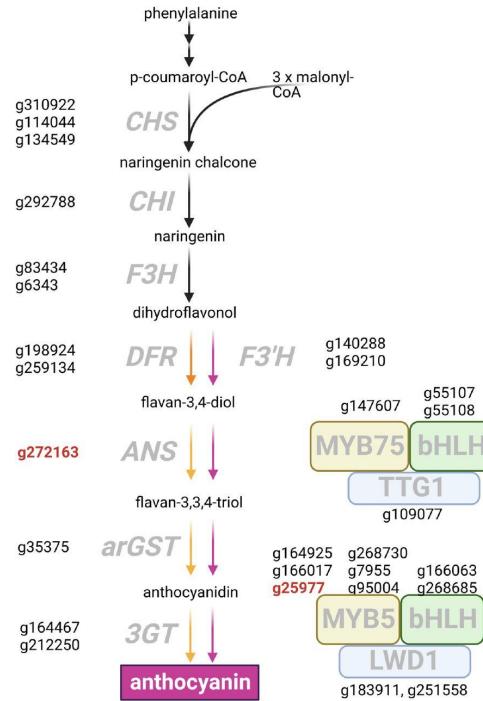
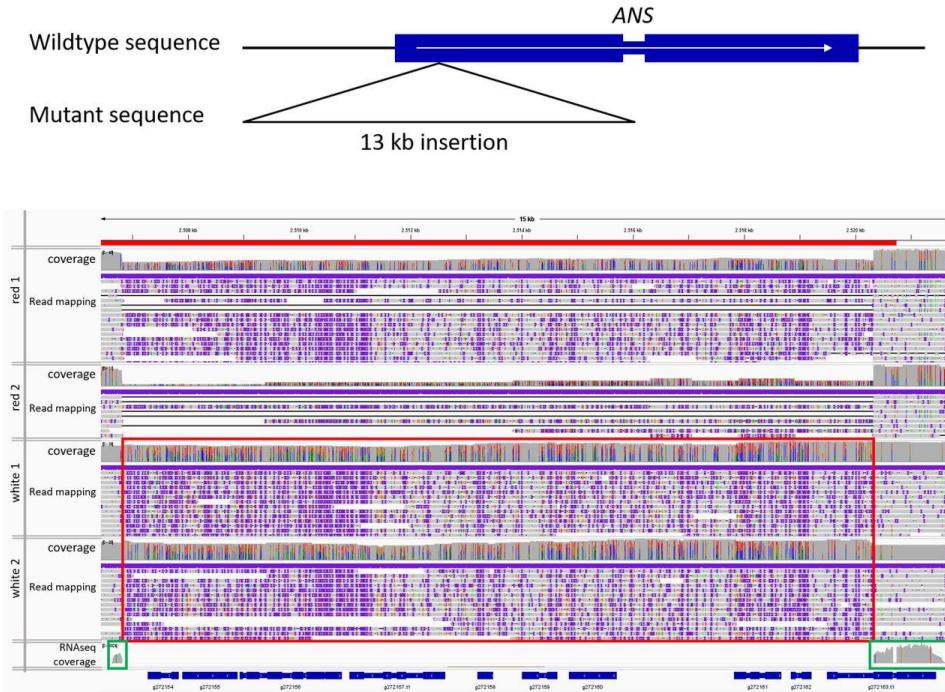


Recent Genomic Studies - Examples

Flower diversity in *Digitalis purpurea* I



Flower diversity in *Digitalis purpurea* L.



Discovery of the withanolide biosynthesis I

Withania



W. somnifera

Physalis

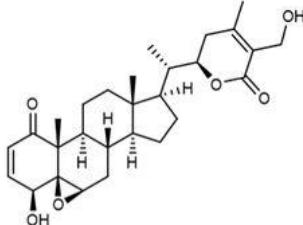


P. pruinosa

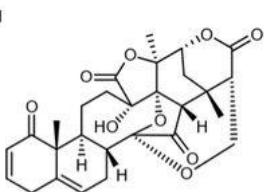
Datura



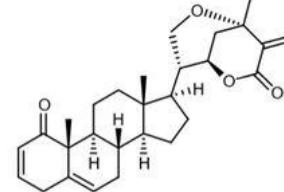
D. innoxia



Withaferin A

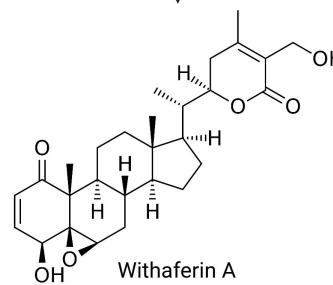
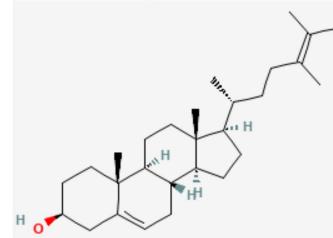


Physalin B



Withametelin

24-Methyldesmosterol

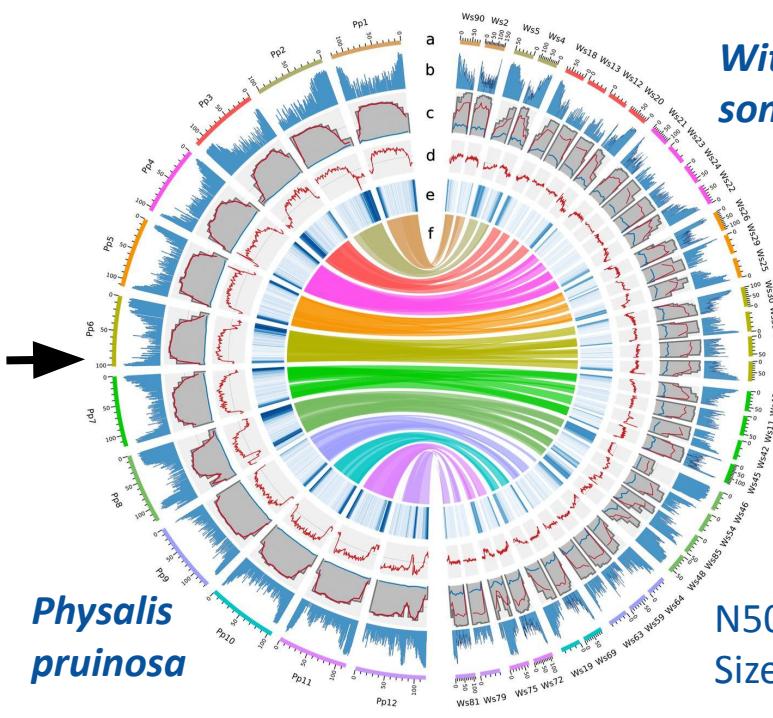


Withaferin A

Discovery of the withanolide biosynthesis II



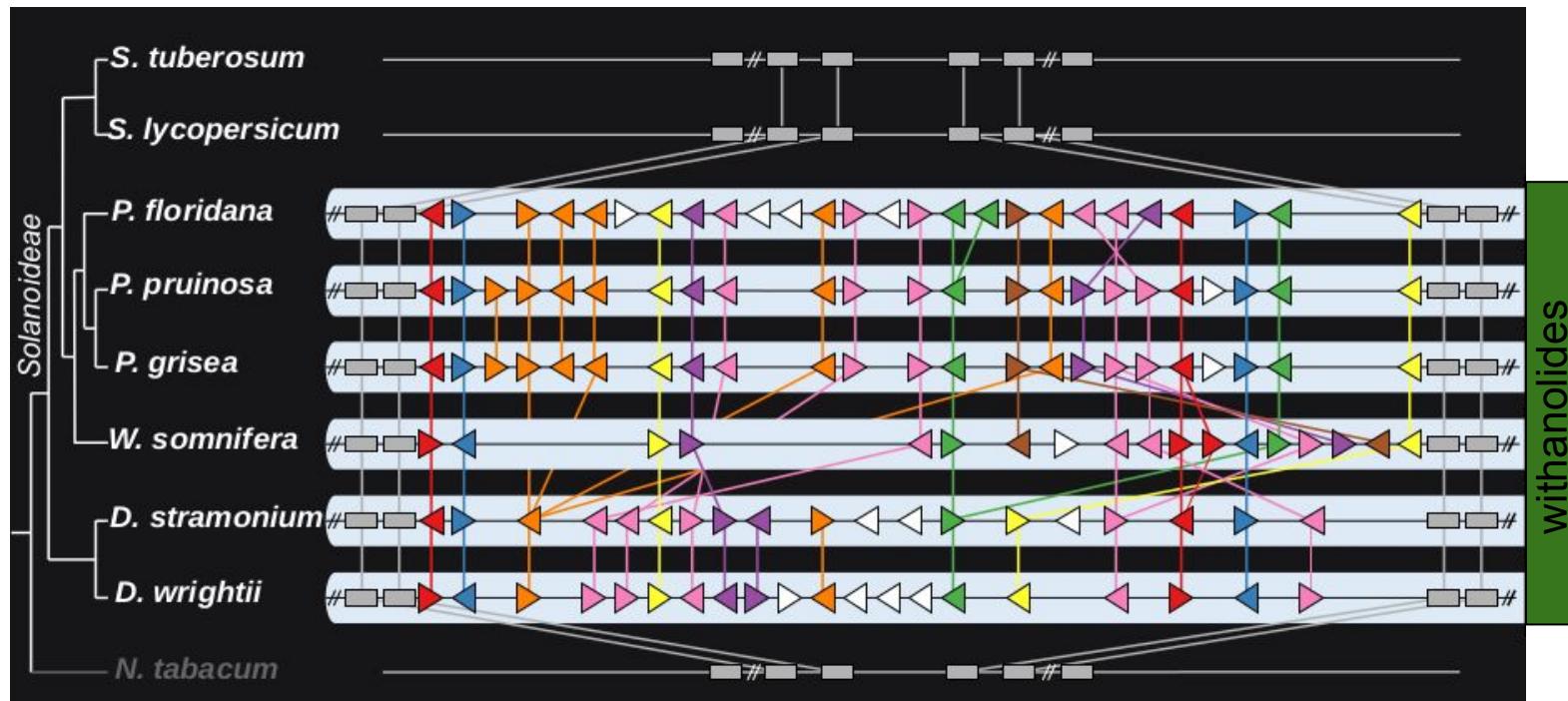
Withania produces especially interesting withanolides



*Withania
somnifera*

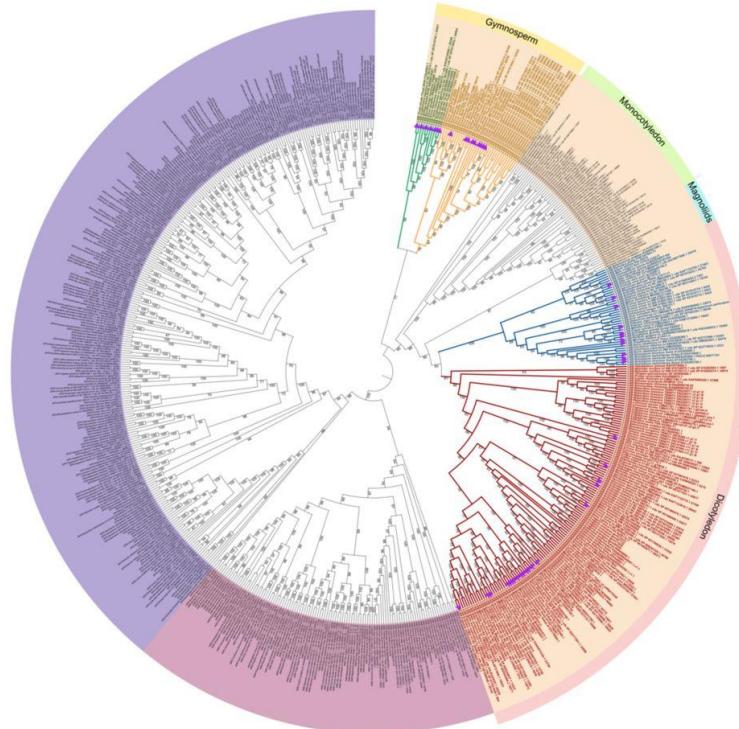
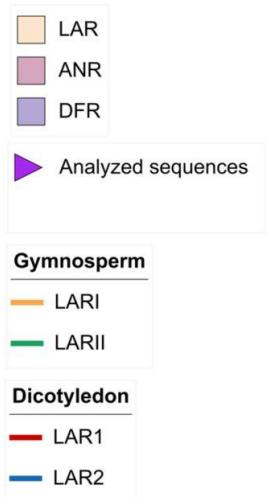
N50: 71 Mb
Size: 2.8 Gbp

Discovery of the withanolide biosynthesis III



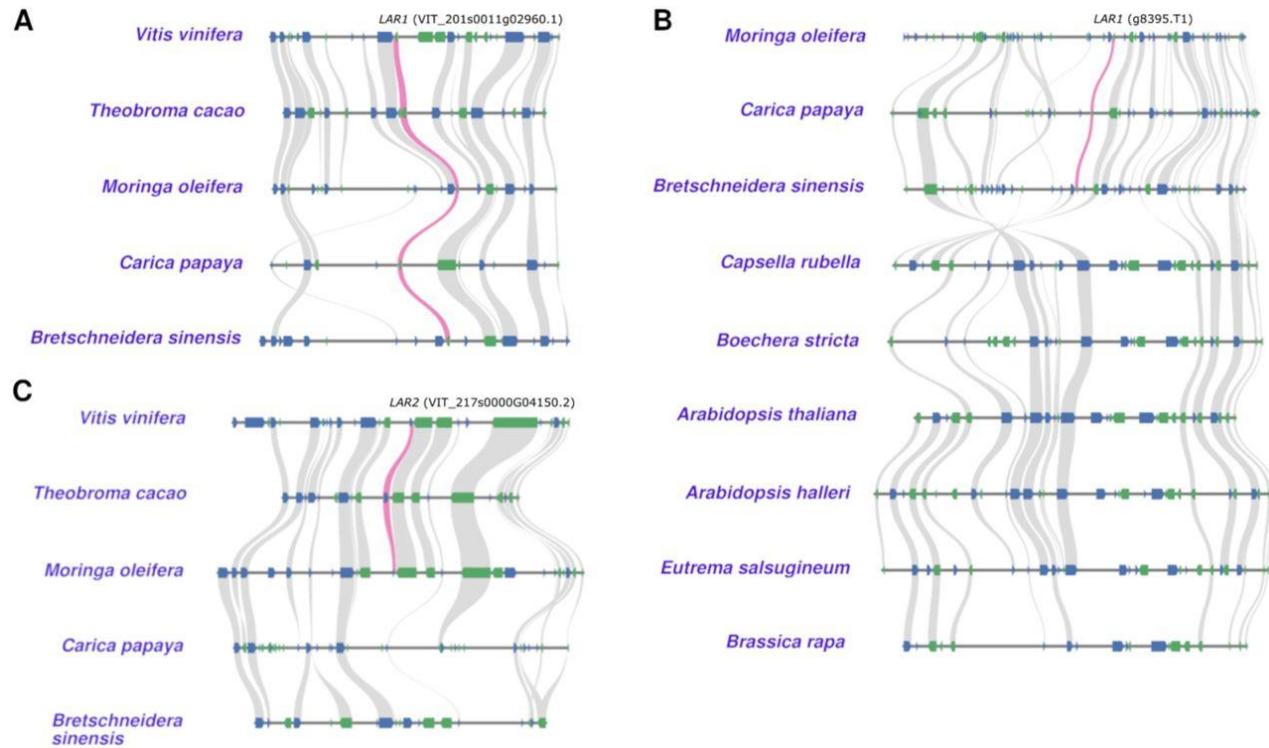
LAR evolution I

- Leucoanthocyanidin reductase (LAR)
- LAR was duplicated at basis of angiosperms
- Large numbers of genome sequences enables such large scale analyses

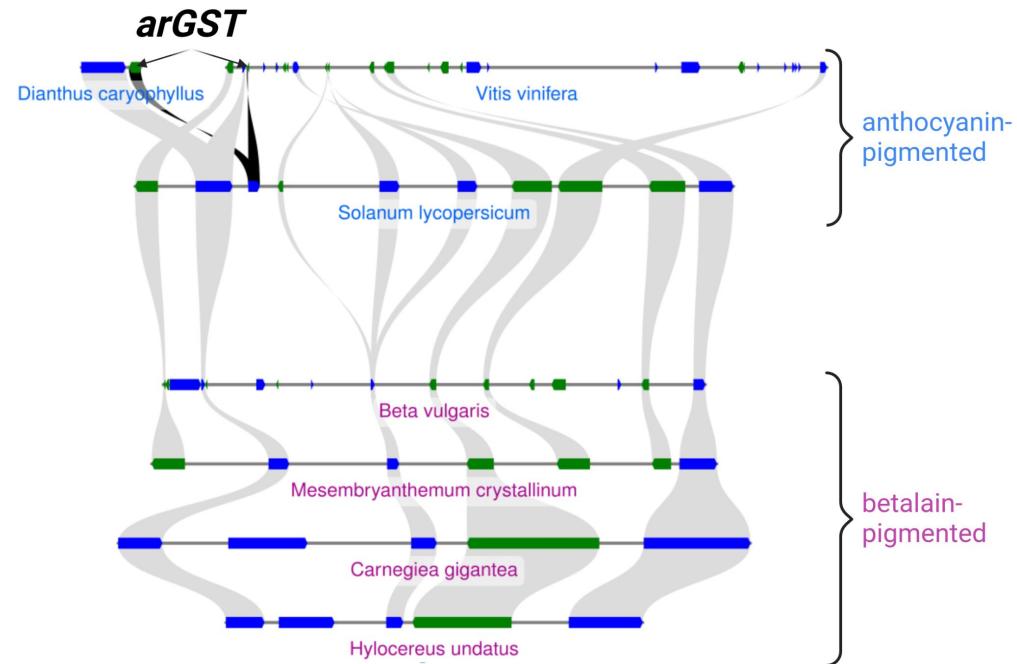
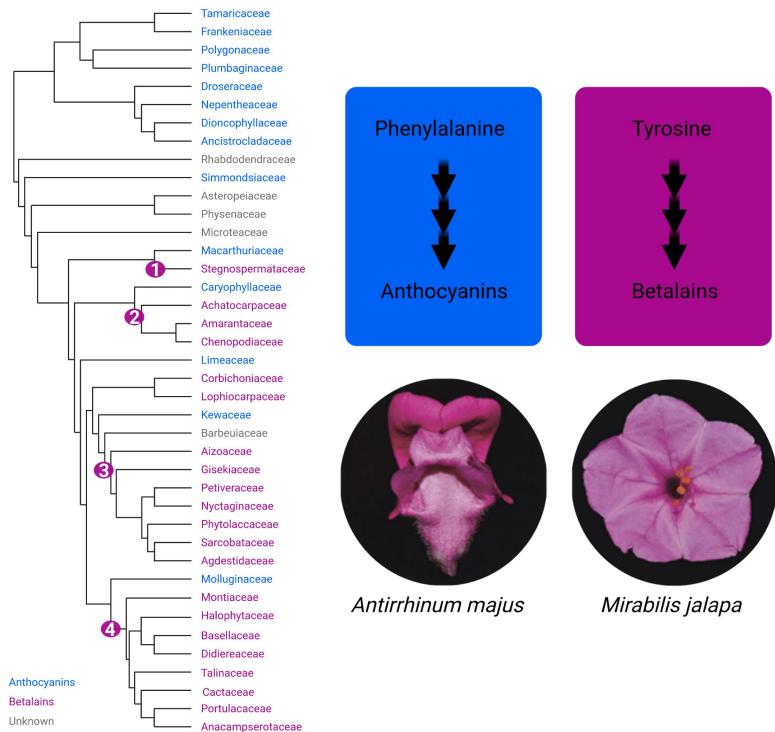


LAR evolution II

- LAR appears lost in Brassicales
- Synteny analyses allow to confirm this

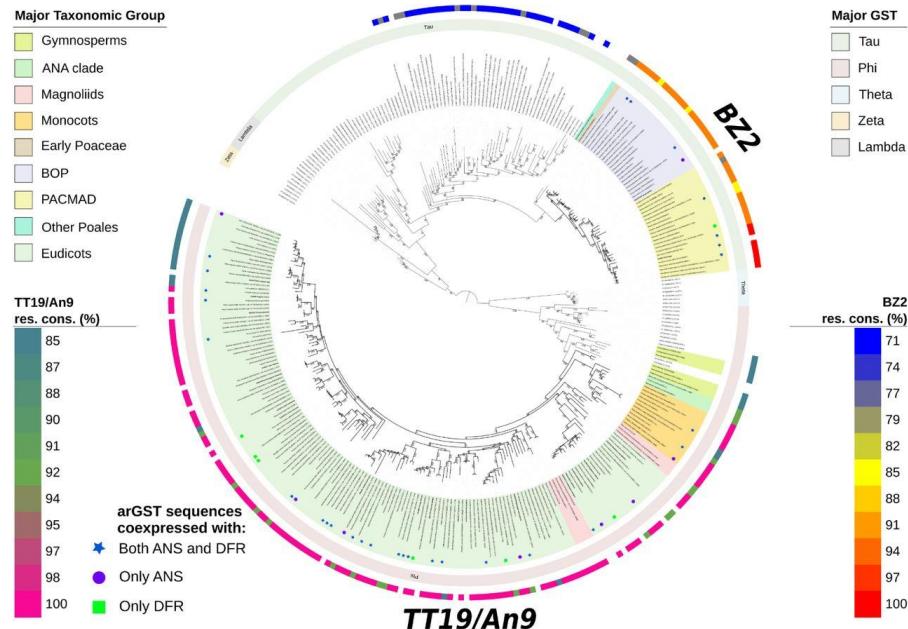


Anthocyanin loss in Caryophyllales

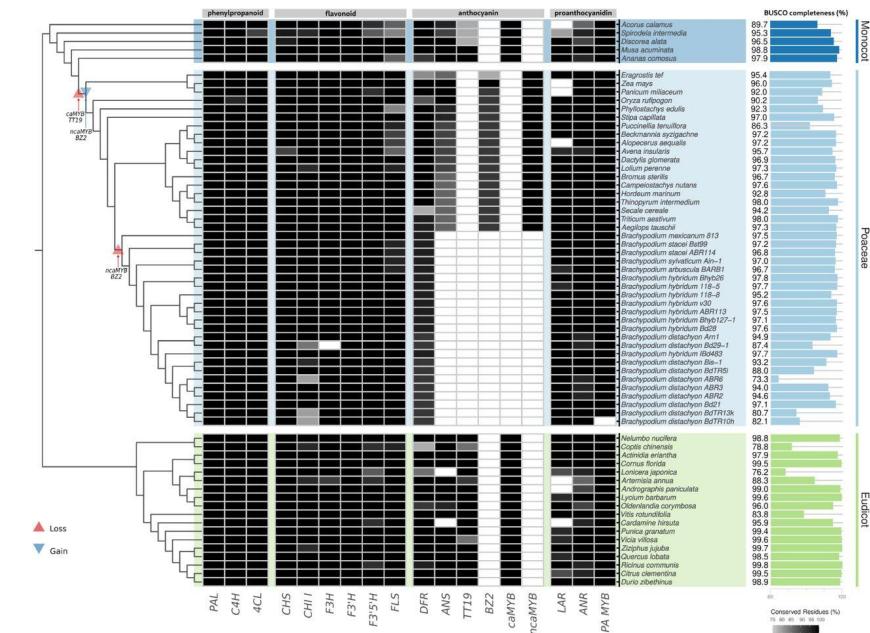
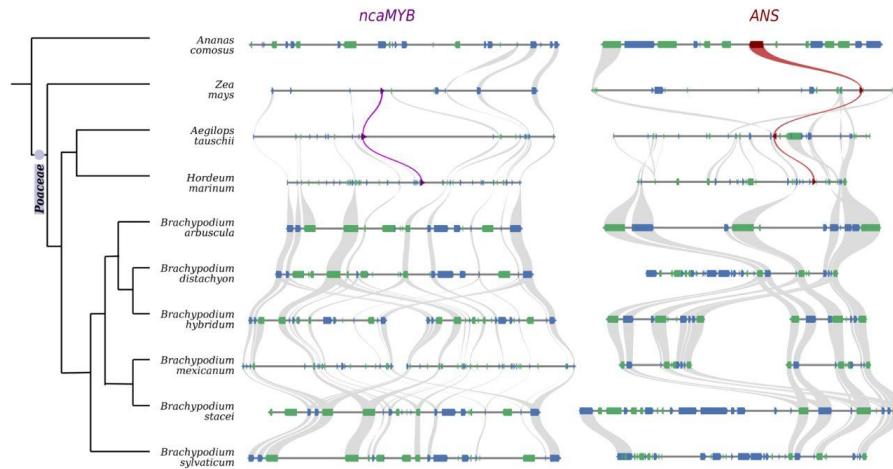


Anthocyanin loss in Poaceae I

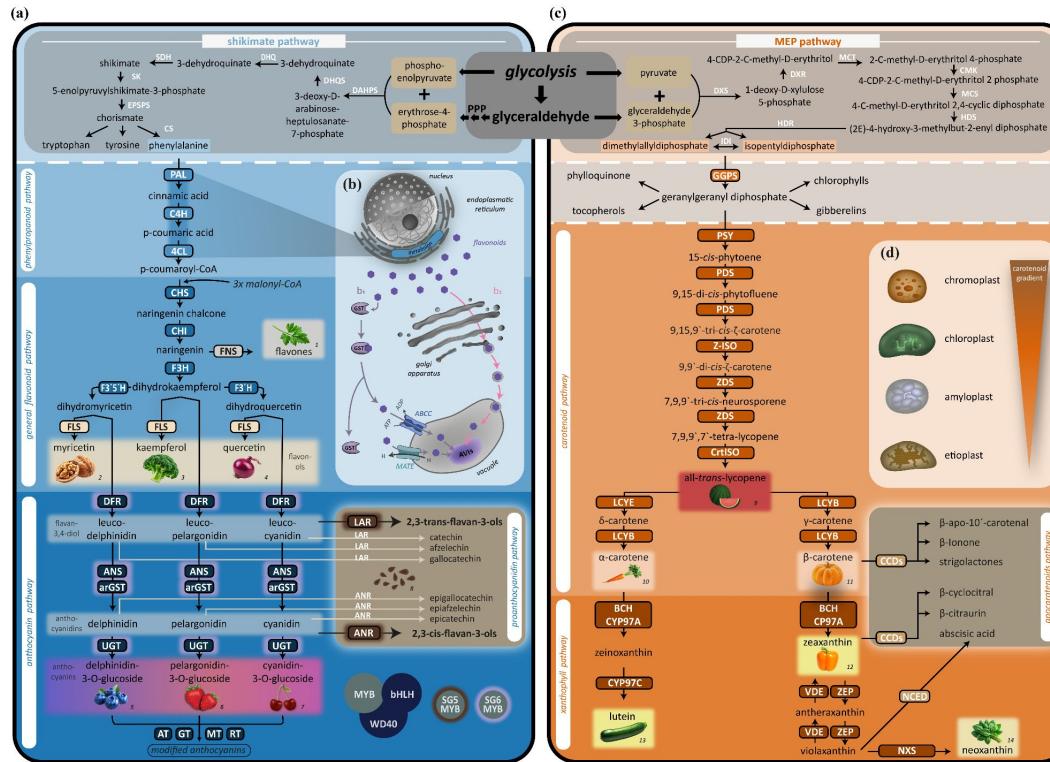
- Poaceae (true grasses) harbour numerous crops
- Anthocyanin biosynthesis gene arGST (TT19/An9) was lost
- Independent recruitment of a novel arGST in Poaceae (BZ2)



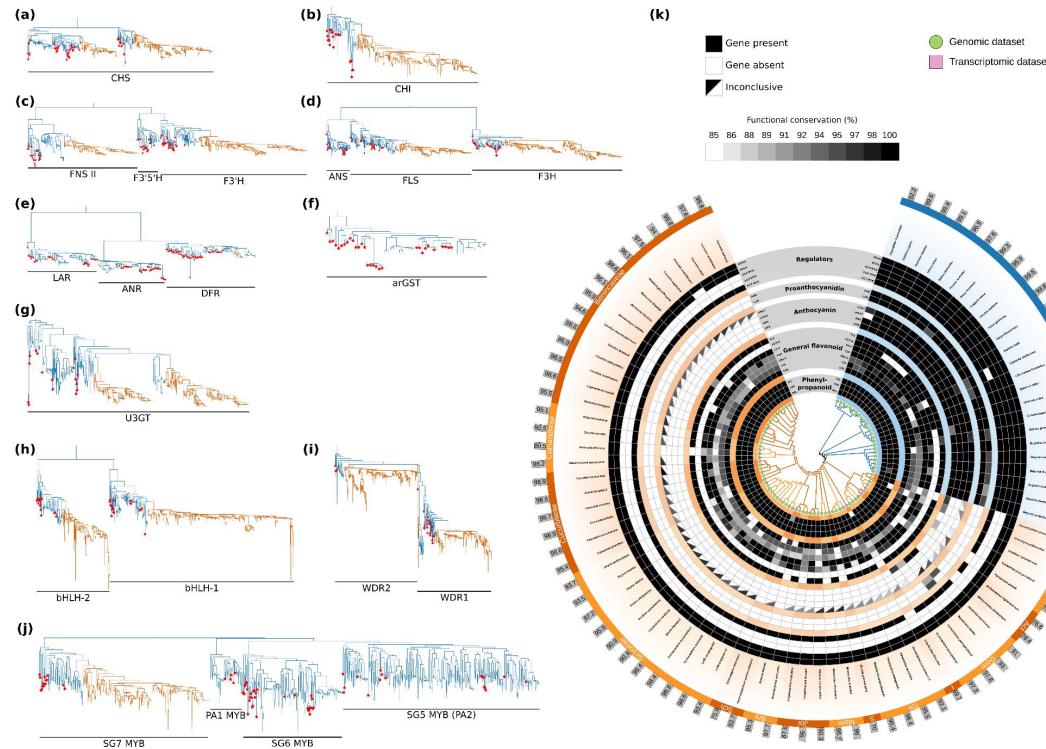
Anthocyanin loss in Poaceae II



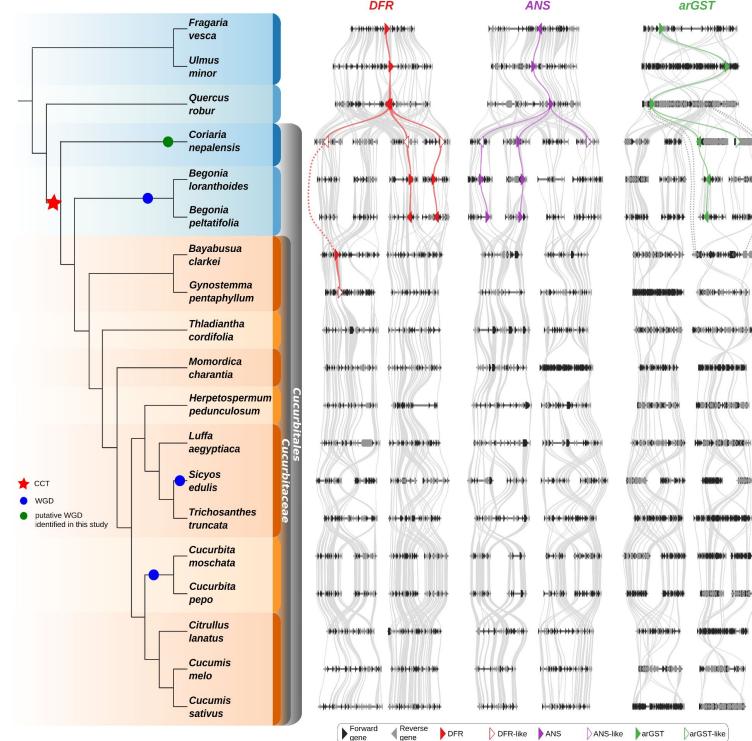
Anthocyanin loss in Cucurbitaceae I



Anthocyanin loss in Cucurbitaceae II



Anthocyanin loss in Cucurbitaceae III



Summary

- Read mapping (short reads vs. long reads)
- Variant calling (short reads vs. long reads)
- Evolutionary genomics
- Comparative genomics & synteny analysis
- Application examples

Time for questions!

Literature

- de Oliveira, J. A. V. S.; Choudhary, N.; Meckoni, S. N.; Nowak, M. S.; Hagedorn, M.; Pucker, B. (2025). Cookbook for Plant Genome Sequences. doi: [10.20944/preprints202508.1176.v2](https://doi.org/10.20944/preprints202508.1176.v2).
- Wolff, K.; Friedhoff, R.; Schwarzer, F.; Pucker, B. (2023). Data Literacy in Genome Research. *Journal of Integrative Bioinformatics*, 2023, pp. 20230033. doi: [10.1515/jib-2023-0033](https://doi.org/10.1515/jib-2023-0033).
- Pucker B, Irisarri I, de Vries J and Xu B (2022). Plant genome sequence assembly in the era of long reads: Progress, challenges and future directions. *Quantitative Plant Biology*, 3, E5. doi: [10.1017/qpb.2021.18](https://doi.org/10.1017/qpb.2021.18).