



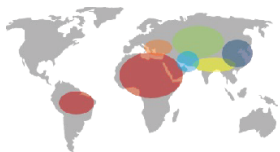
Bioinformatics tools for identifying traits in PGR

Luciana Gaccione

Unlocks the diversity in PGRs

Genomics

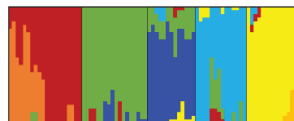
Collection of resources from sample regions



Storage in genebank with passport information



Regional characterization of genomic information



Screen for target traits

Phenomics - Metabolomics



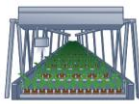
Satellite



Unmanned aerial vehicle



Unmanned ground vehicle



Gantry system



Facility phenotyping



Seed phenotyping



Root phenotyping



Microphenotyping



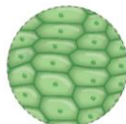
Population



Individual



Organ



Tissue



Cell

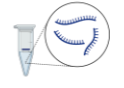


Metabolite

Transcriptomics



RNA extraction



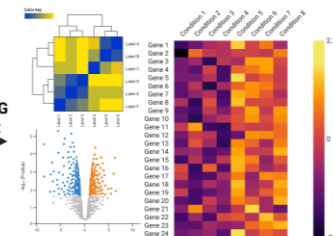
Sequencing



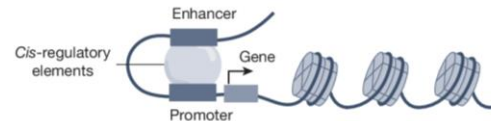
Reads count



DEGs
GO and KEGG
enrichment



Epigenomics



Chromatin
openness



H3K27ac



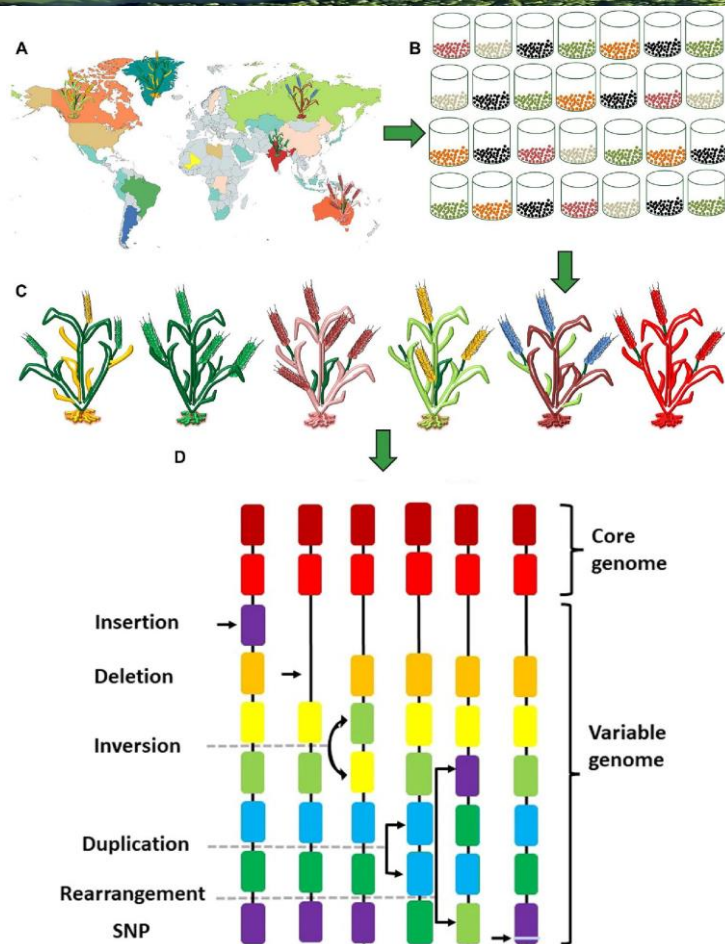
H3K4me3



Reference



Unlocks the diversity in PGRs



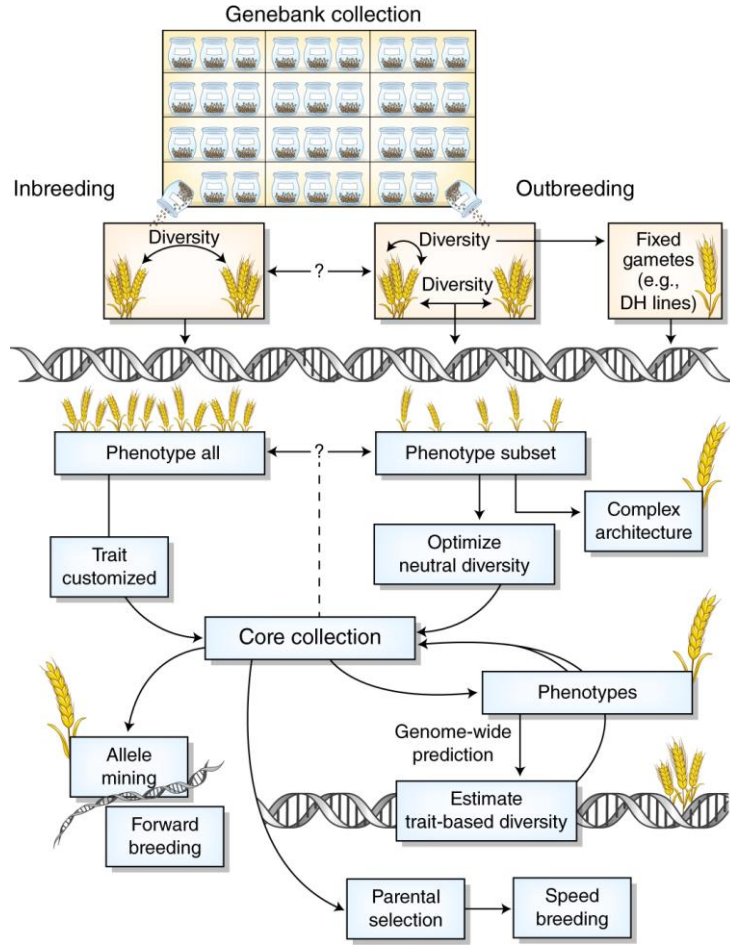
Omics technologies represent a leap forward for the conservation, management and characterization of PGRs.

Current PGR management does not involve the routine use of omic tools to trace accessions during seed regeneration or vegetative propagation and identifying traits.

A gap between what is technically possible with modern omics, and what is actually implemented in routine PGRs conservation practices

Bioinformatics pipelines can help bridge the gap

Informed selection for PGRs pre-breeding programs



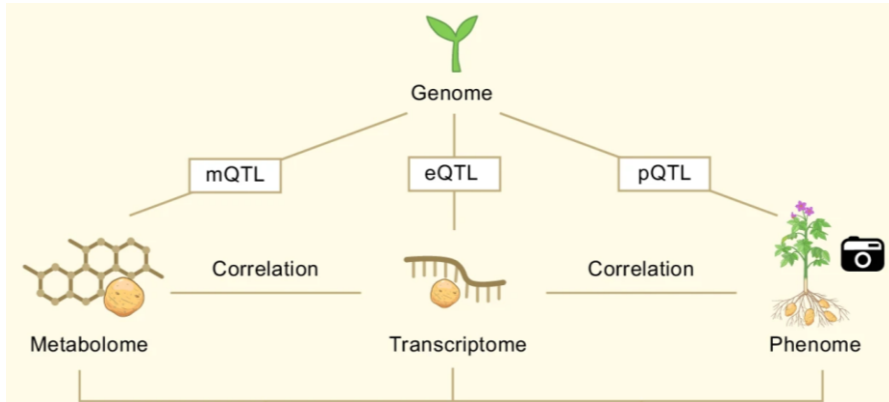
Genotypic and phenotypic information enables the informed selection of the most promising genetic resources.

Depending on the genetic architectures of the traits, **entire collections** can be phenotyped, or **core collections** maximizing genetic diversity.

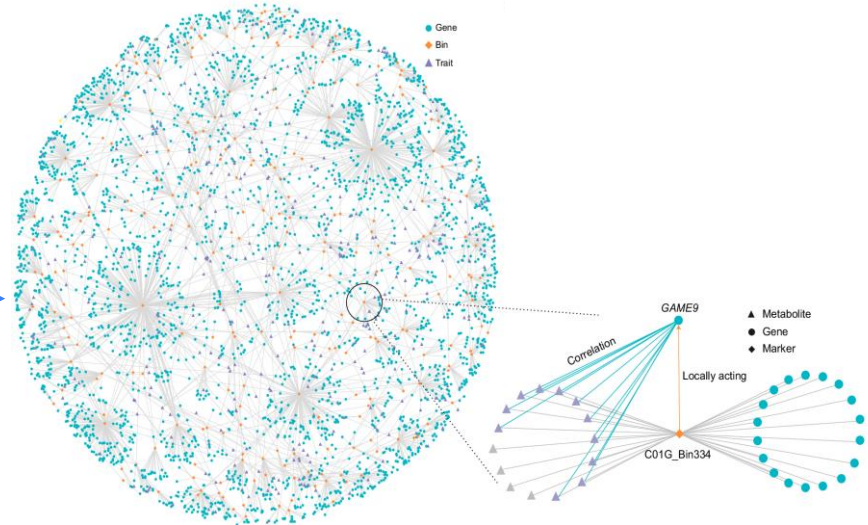


Genotypes with the highest breeding values enter pre-breeding programs and are used in genetic studies to elucidate the molecular basis of beneficial traits of PGRs.

Integrated multi-omics is more than the sum of its parts



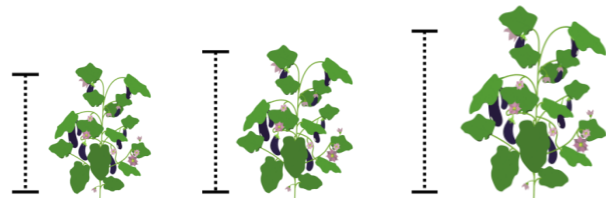
Multi-omics refers to the integrated study of various "omics" fields, including genomics, transcriptomics, proteomics, metabolomics, and epigenomics.



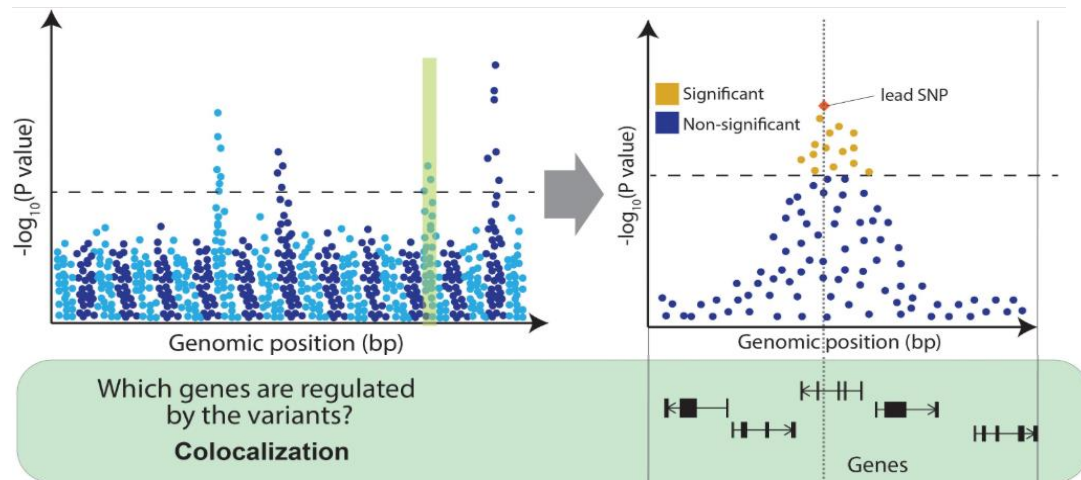
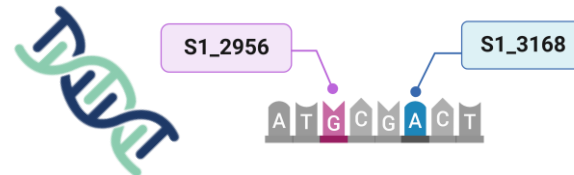
This approach is particularly valuable for unraveling intricate **networks and pathways associated with PGRs traits**

Traditional GWA approaches

Phenotype



Genotype





Linkage disequilibrium decay

What size should my sliding window be to define QTLs?



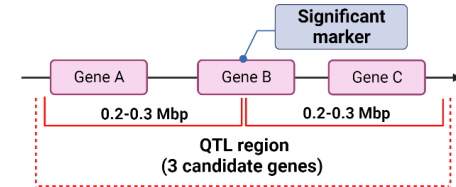
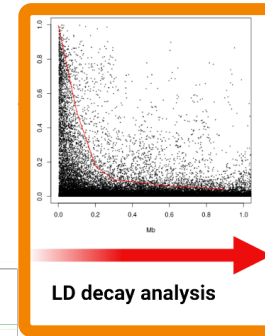
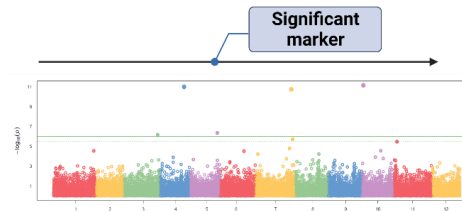
A simple rule of thumb is to set the window where linkage disequilibrium decays to the genome-wide background

Linkage equilibrium: haplotype frequencies in a population have the same value that they would have if the genes at each locus were combined at random.

Linkage disequilibrium: Non-random association of alleles at different loci in a given population

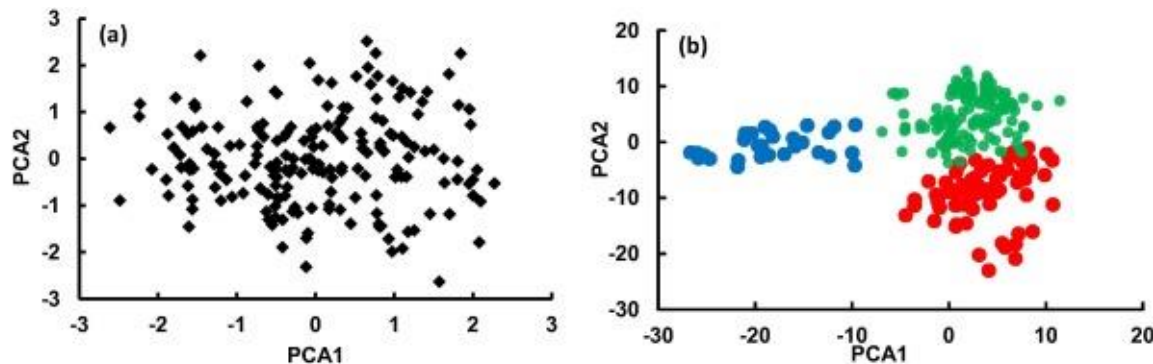
How LD decay is calculated?

1. **Get R^2 value** for each marker vs all the others (on the same chromosome or at a certain distance) using **plink**
2. **Use PopLDdecay**



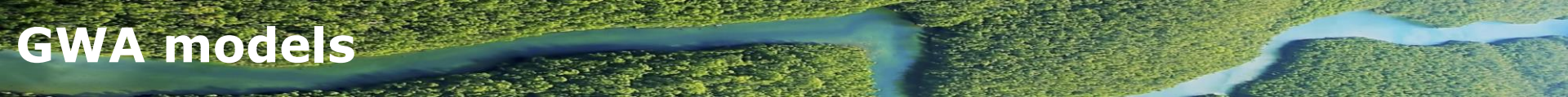
Population Structure and spurious associations

- The problem arises when cases disproportionally represent a genetic subgroup, **resulting (apparently) associated with the trait.**
- Geographic or growth habit, etc. lead to having a specific genetic variation and an effect on the end-use of association analysis



Main methods:

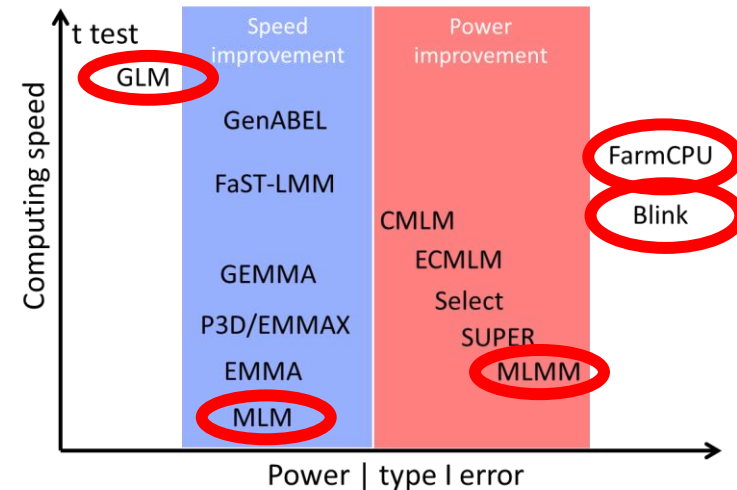
- Principal component analysis (PCA)
- Structured association (SA)
- Kinship analysis



GWA models

Multi-locus mixed models

- **MLMM** (Multi-Locus Mixed Model) is an extension of the standard mixed linear model (MLM) used in GWAS. The goal is to better handle **polygenic traits** by including **multiple associated loci** as cofactors, rather than testing SNPs one by one. The most significant SNPs as **pseudo-QTNs** (covariates), then removes non-significant ones.
- **FarmCPU** (Fixed and Random Model Circulating Probability Unification) **separates marker testing (fixed effect model) from population structure/kinship control (random effect model)**, iteratively updating pseudo-QTNs to reduce confounding. Iteratively updates both parts, so pseudo-QTNs and tested SNPs keep adjusting each other.
- **BLINK** (Bayesian-information and Linkage-disequilibrium Iteratively Nested Keyway) improves over FarmCPU by **replacing the random effect with LD-based clustering** and BIC model selection, making it faster and more powerful for large datasets.



Traditional GWA workflow

Phenotype

phenotyping
field data collection/metabolites
measurements
.txt

outliers' removal
inti
.no_outliers.txt

normalization
bestNormalize
.normalized.txt

BLUEs/BLUPs/mean
inti
.final.txt

test for associations
GAPIT and selected model
.pvalue_per_tested_marker.txt

determine significance threshold
Bonferroni correction/FDR/qvalue
significant_SNPs.txt

define QTLs
according to LD decay
bedtools
.QTL_regions.txt

Genotype

SNPs/indels calling
GATK/DeepVariant
.g.vcf/.vcf

joint vcf files
GATK/GLnexus
.merged.vcf

hard filters
GATK/bcftools
.merged.hardfilters.vcf

filtering for GWA
bcftools
.GWAfilters.vcf

imputation
Beagle/mean
.final.vcf

sequencing
GBS/resequencing
.fq.gz

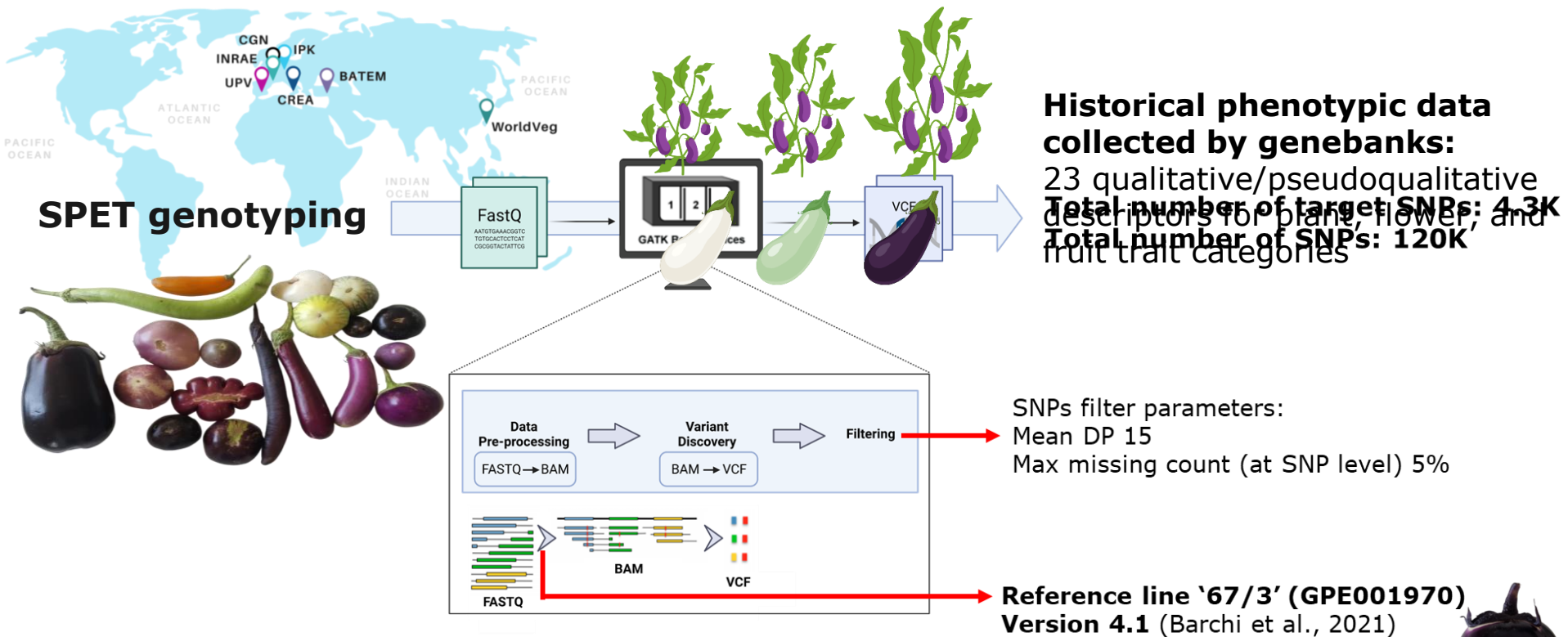
reads cleaning
trimmomatic/fastp
.gbz

alignment to a
reference genome
BWA
.bam

bam sort/index
samtools
.bam

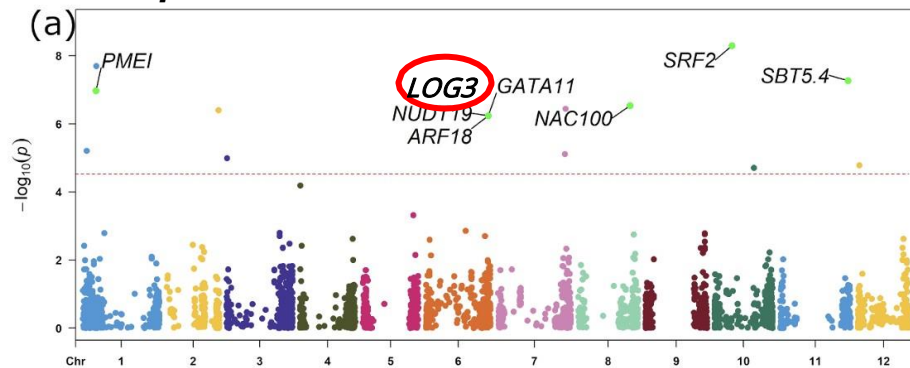
Clara Parabricks (NVIDIA® Parabricks®) is a GPU-based software suite for performing analysis of next generation sequencing (NGS) DNA and RNA data. It delivers results at blazing fast speeds and low cost.

GWA: Case study - SPET data and historical GWAS in eggplant

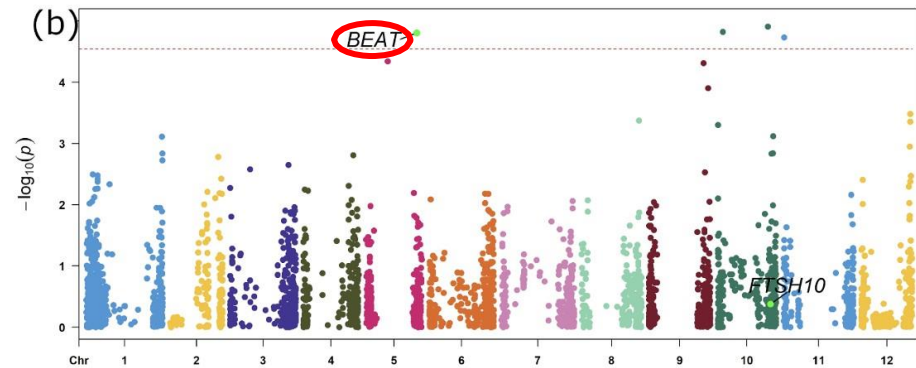


GWA: Case study – SPET data and historical GWAS in eggplant

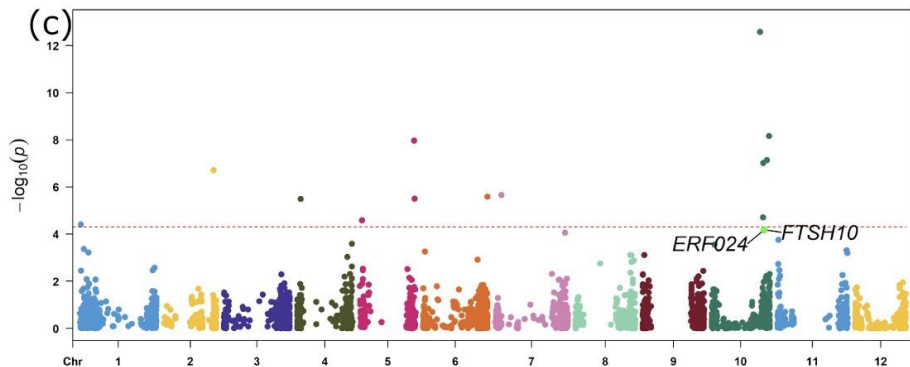
Leaf prickles



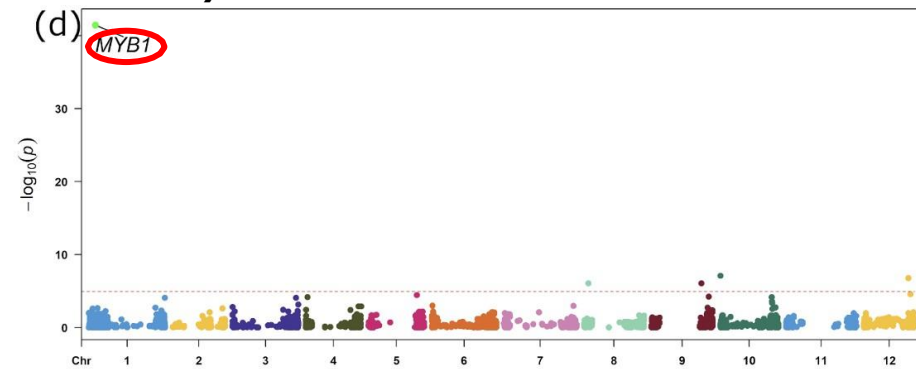
Corolla color



Fruit color

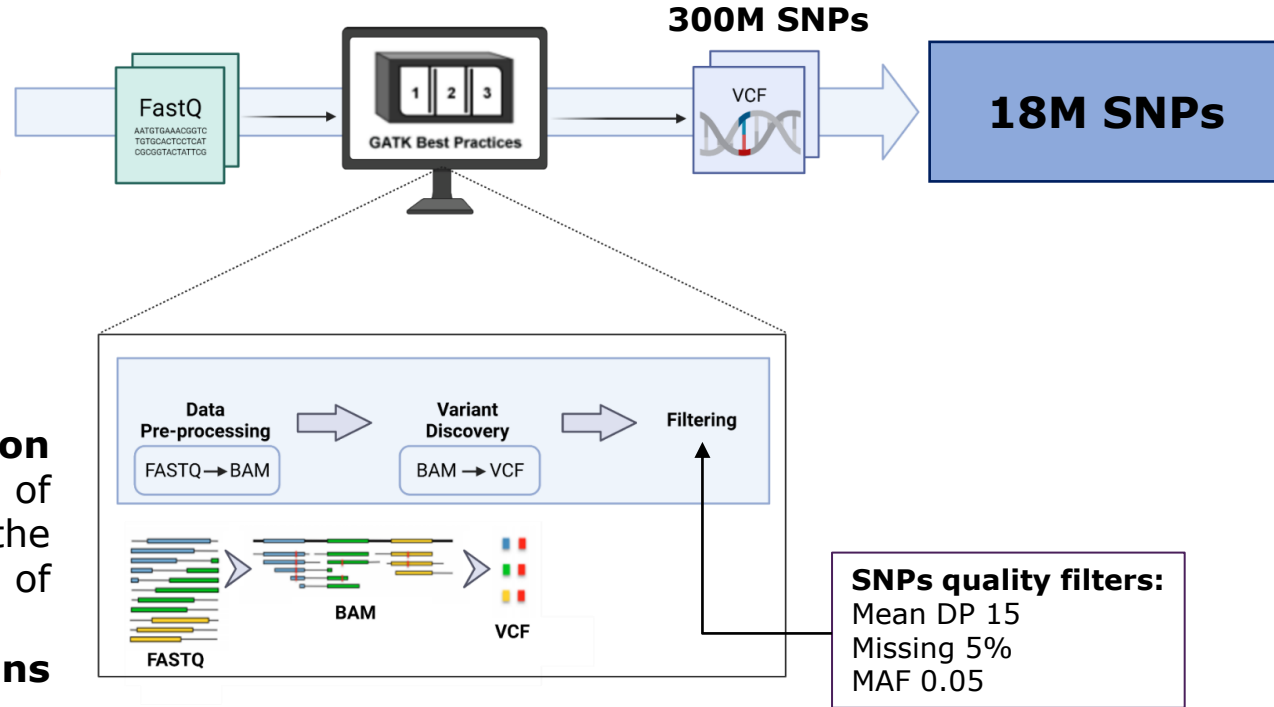


Anthocyanin distribution



Candidate genes circled in red have been confirmed by different models, confirming the quality of genebank data

GWA: Case study – Resequencing and GWA in pepper



- The **G2P-SOL core collection of *Capsicum* spp.**, consisting of 423 accessions representing the genetic variability of a panel of 10,083 accessions.
- **393 *C. annuum* accessions** used for GWAS

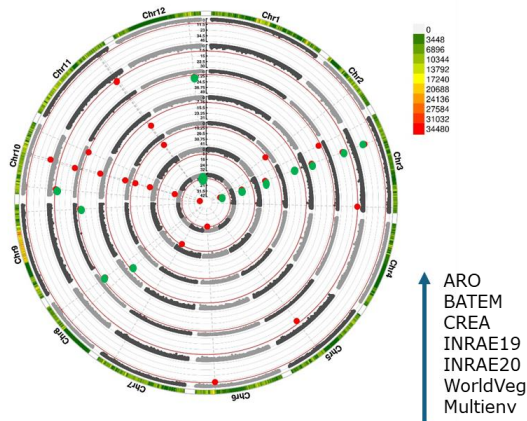
The genome sequence of the *C. annuum* var. *annuum* Zhangshugang was used as reference

GWA: Case study – Pepper: Agronomic traits

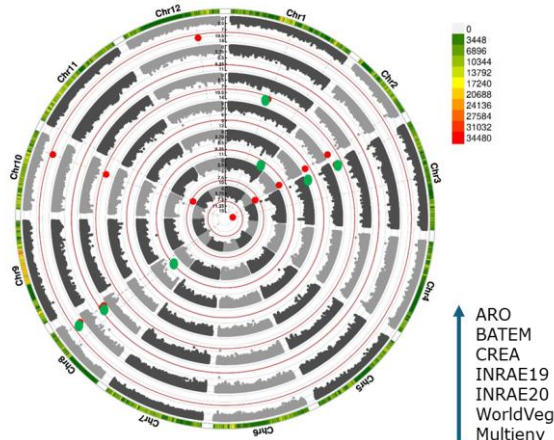
- **207 significant traits** based on pepper descriptors from the International Plant Genetic Resources Institute (IPGRI)
- **112 QTL regions** based on the linkage disequilibrium
- **Up to 6 independent field trials across different years and locations:** ARO in Israel (2020), BATEM in Türkiye (2020–2021), CREA in Italy (2019), INRAE in France (2019–2020), WorldVeg in Taiwan (2020–2021)



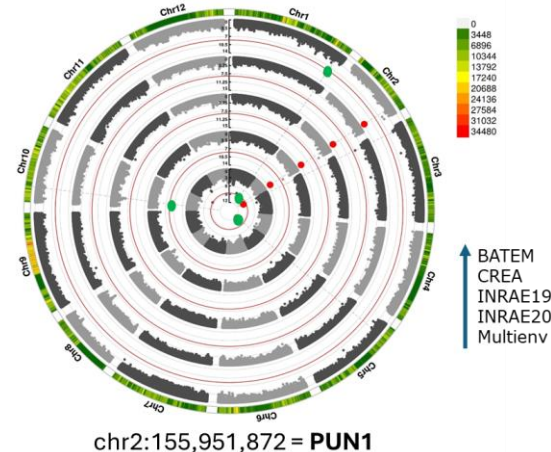
Fruit shape index (FSH_I)



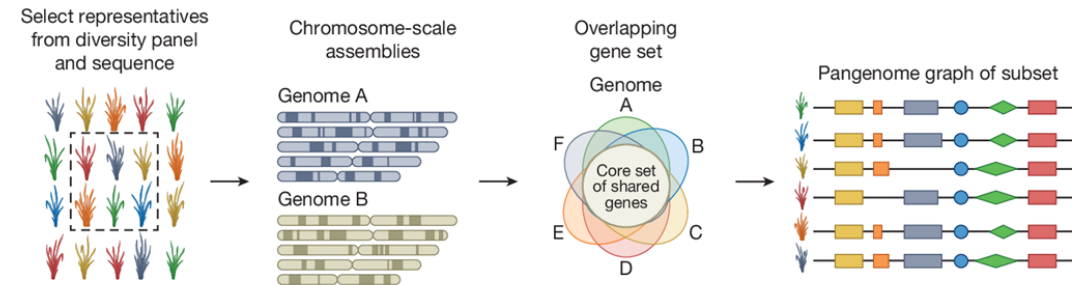
Fruit weight (FW_e)



Fruit pungency (FP)

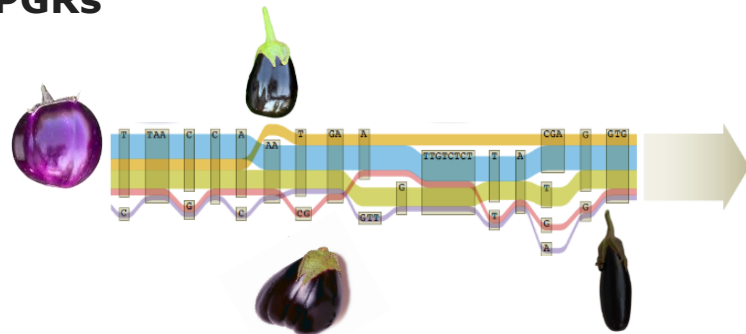
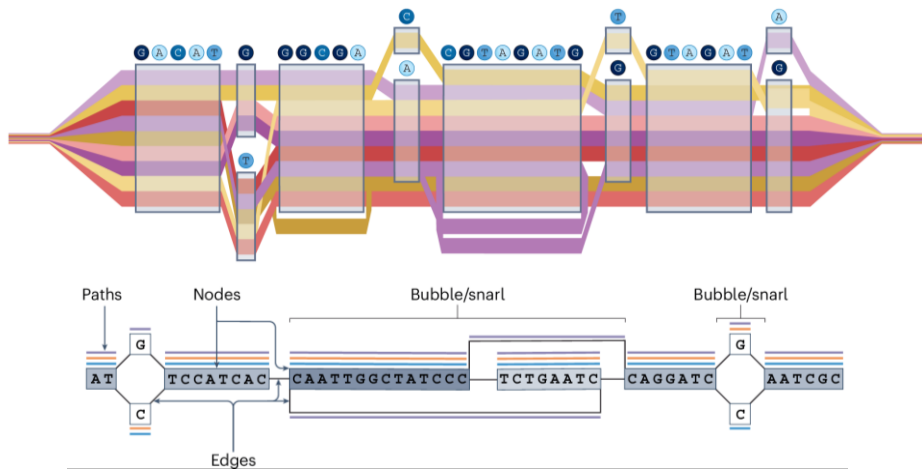


Pangenomics and Pangenome-wide association studies (Pan-GWA)

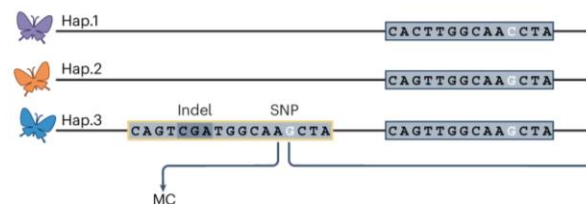
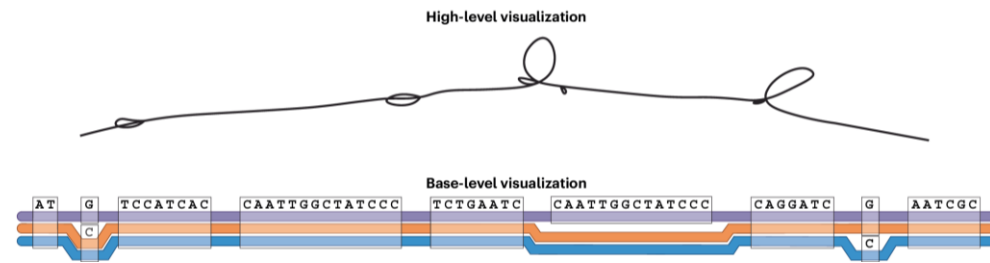
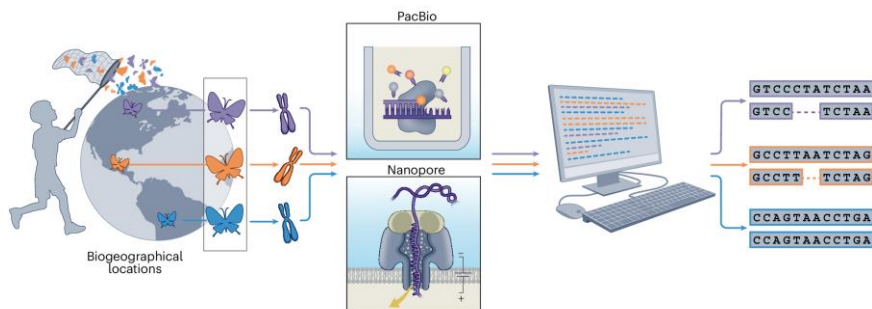


Pangenomics and pangenome graphs have emerged as transformative tools, reshaping the way we study genetic diversity and uncover the complexity of genomes across populations.

There is a need of **new tools and pipeline to handle graph structures for pangenome-wide associations analysis for identifying traits in PGRs**



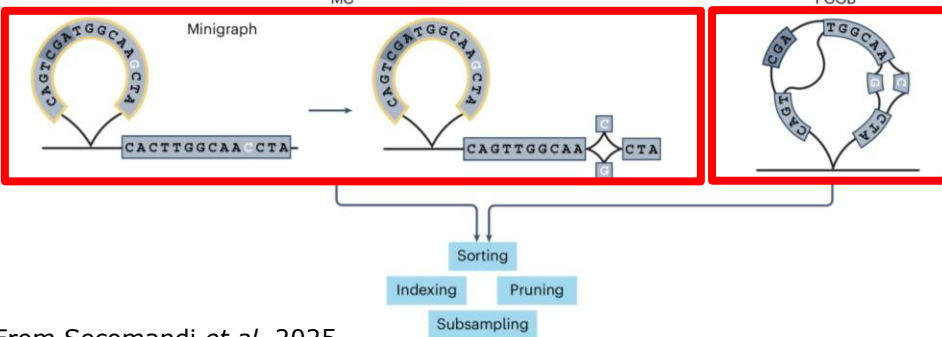
Pangenome graphs building pipeline



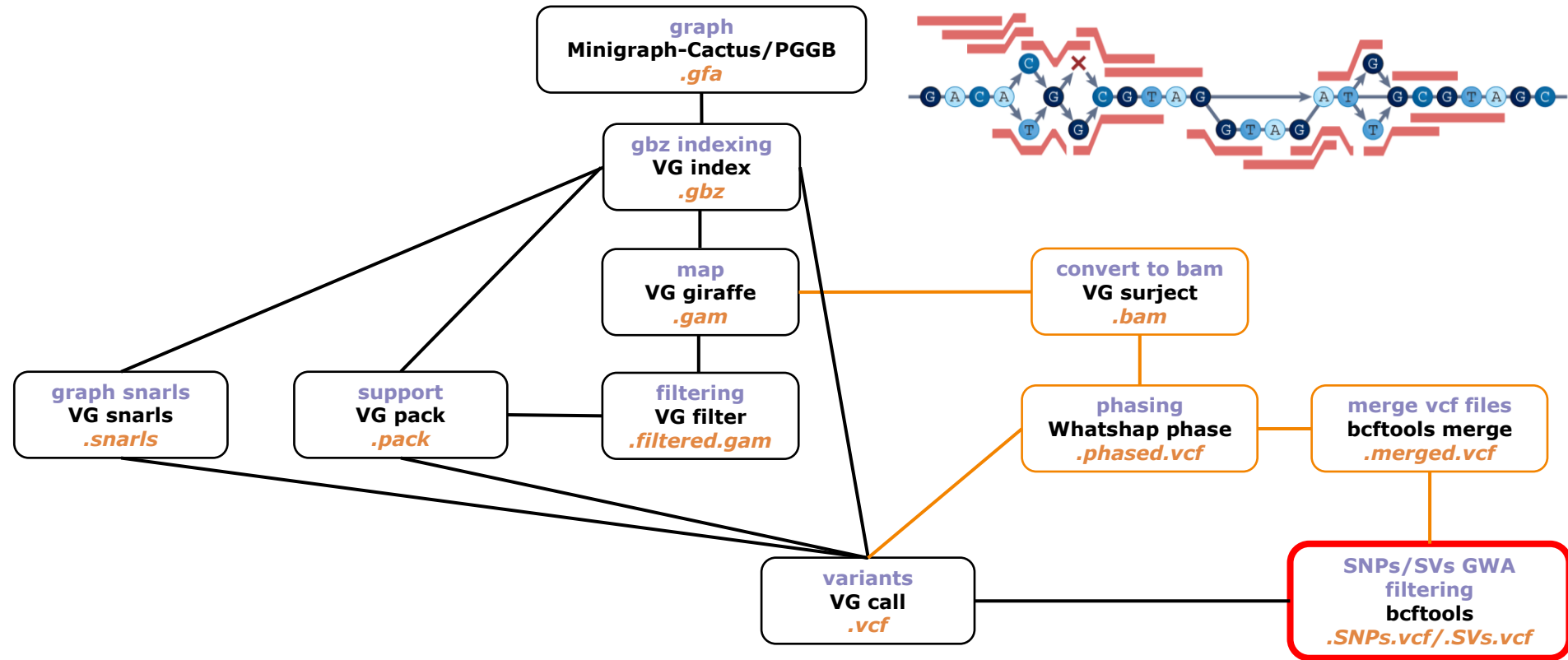
	<i>PG-SMA</i>	<i>PG-SM</i>
Tool	PGGB	Minigraph-Cactus
Length	2,934,734,179	2,540,649,469
Node count	141,826,052	33,186,357
Edge count	194,712,053	45,701,072
Average degree	2.75	2.75
Minimum degree	1	1
Maximum degree	1,073	29
Paths count	480	396
Steps count	3,473,585,945	635,992,212

Nodes maximum degree are the number of edges connected to a node.

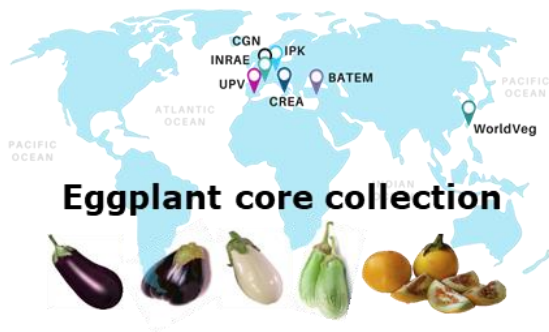
The higher the maximum degree, the more complex the graph structure, indicating highly connected nodes.



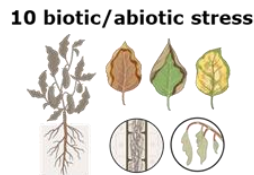
Pangenome-wide association studies (Pan-GWAS)



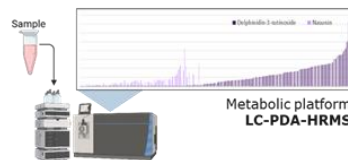
Pan-GWAS: Case study - Eggplant



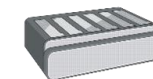
***S. melongena* pangenome graph "PG-SM"**



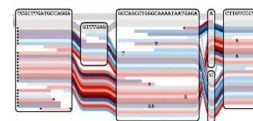
162 semi-polar fruit metabolites



Oxford Nanopore sequencing



Assembly + Scaffolding with Hi-C data (Yahs and Juicebox) or reference-based (RagTag)



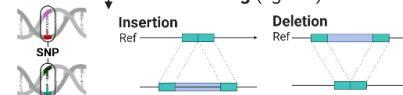
Align short reads to the graph (vg giraffe)

33 *S. melongena* chromosome-scale genome

Pangenome graph construction (Minigraph-Cactus)

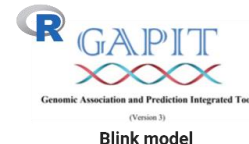
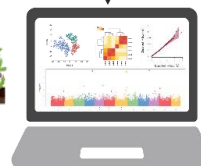


Variant calling (vg call)



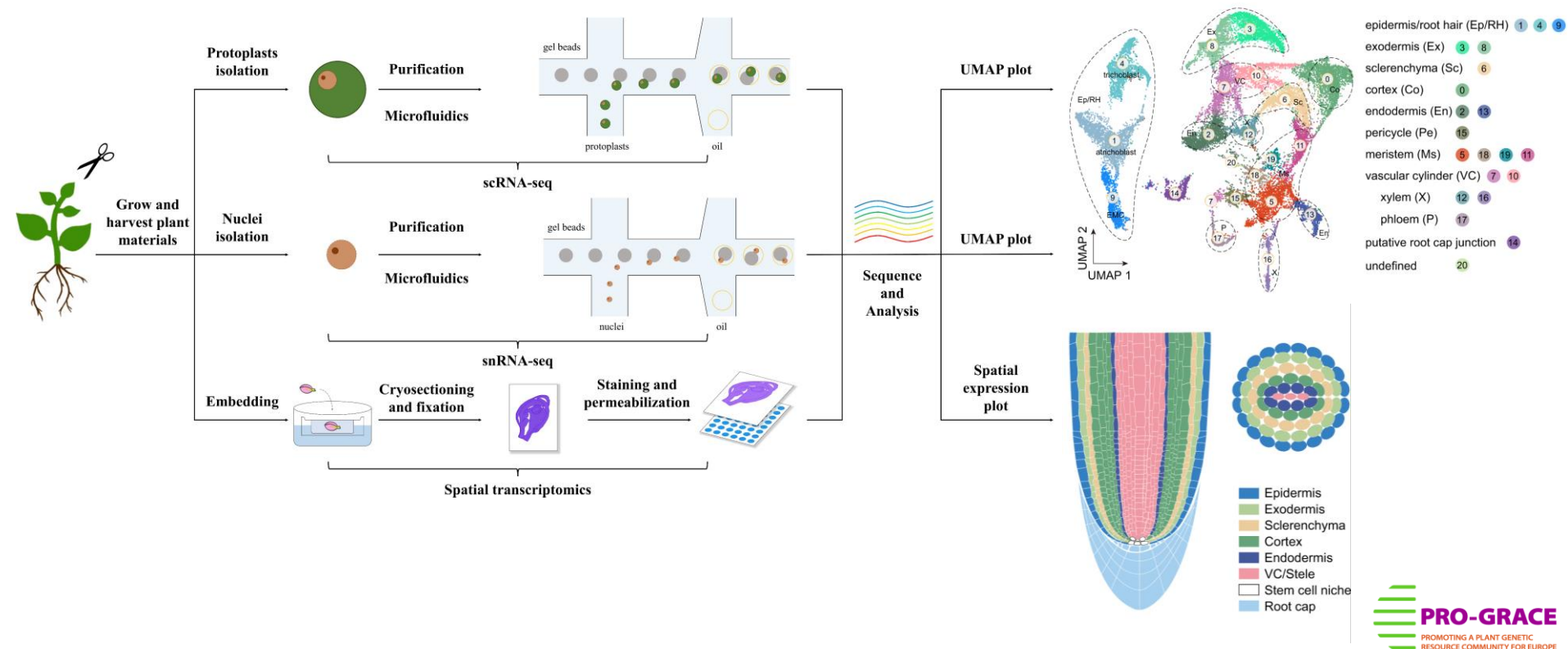
GWA analyses

Multi-environment approach



Single cell and spatial omics

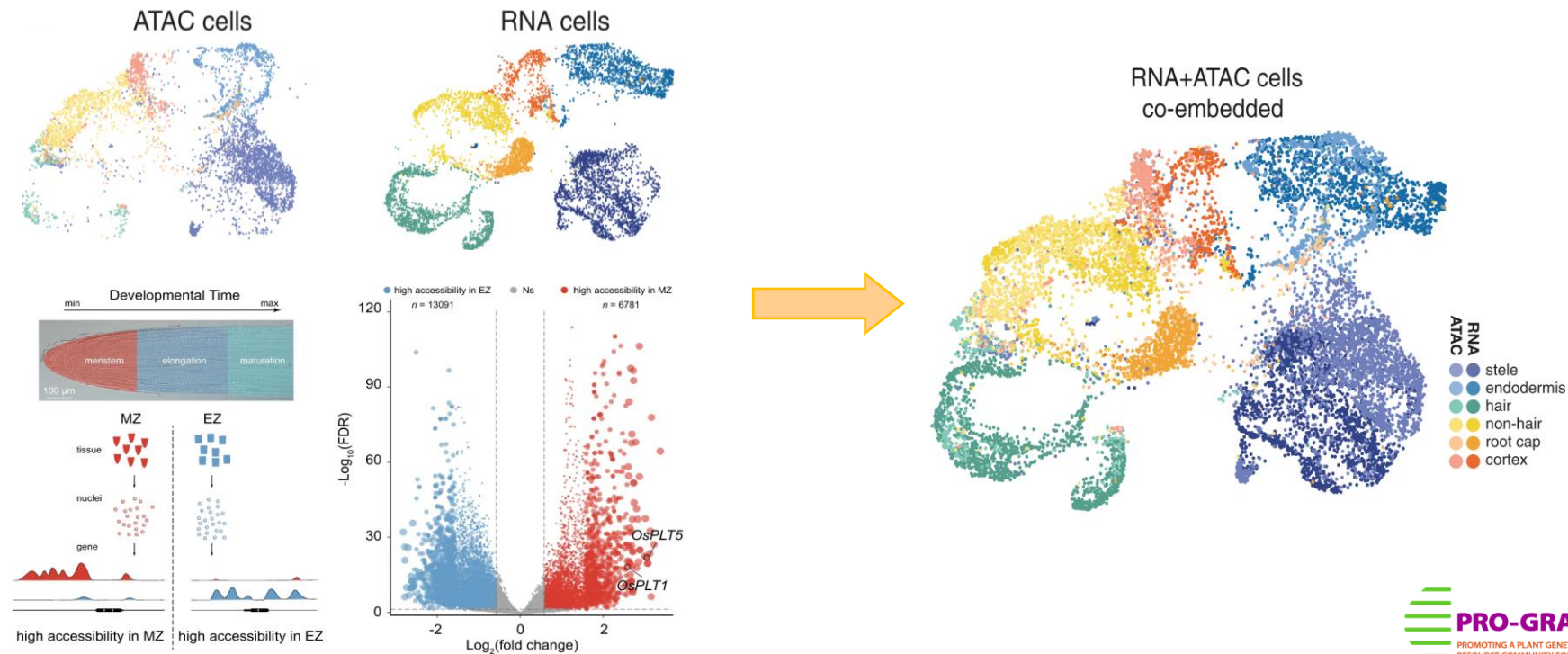
Single cell omics allow to study **cell type-specific transcriptomic** and **chromatin-accessible landscapes**



Single cell and spatial omics

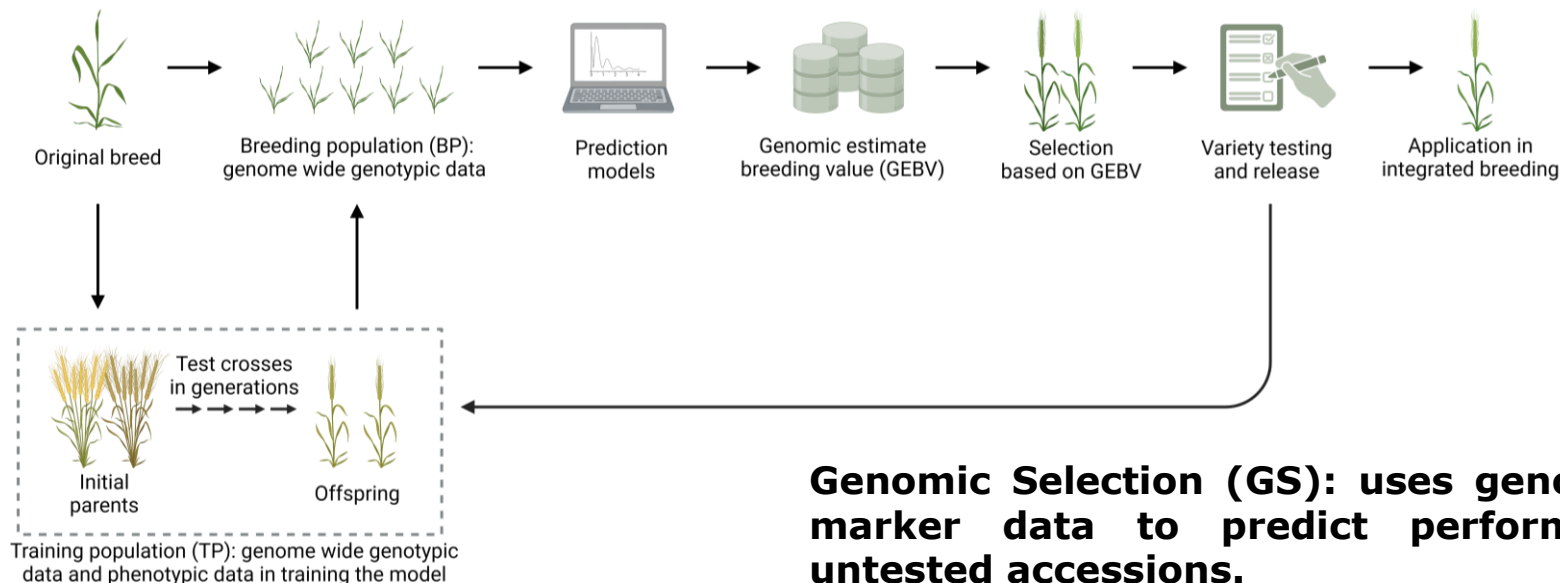
Integration of single cell omics possible with different tools:

- **Signac** and **Seurat** R packages
- **Sanpy** - Single-Cell Analysis in Python



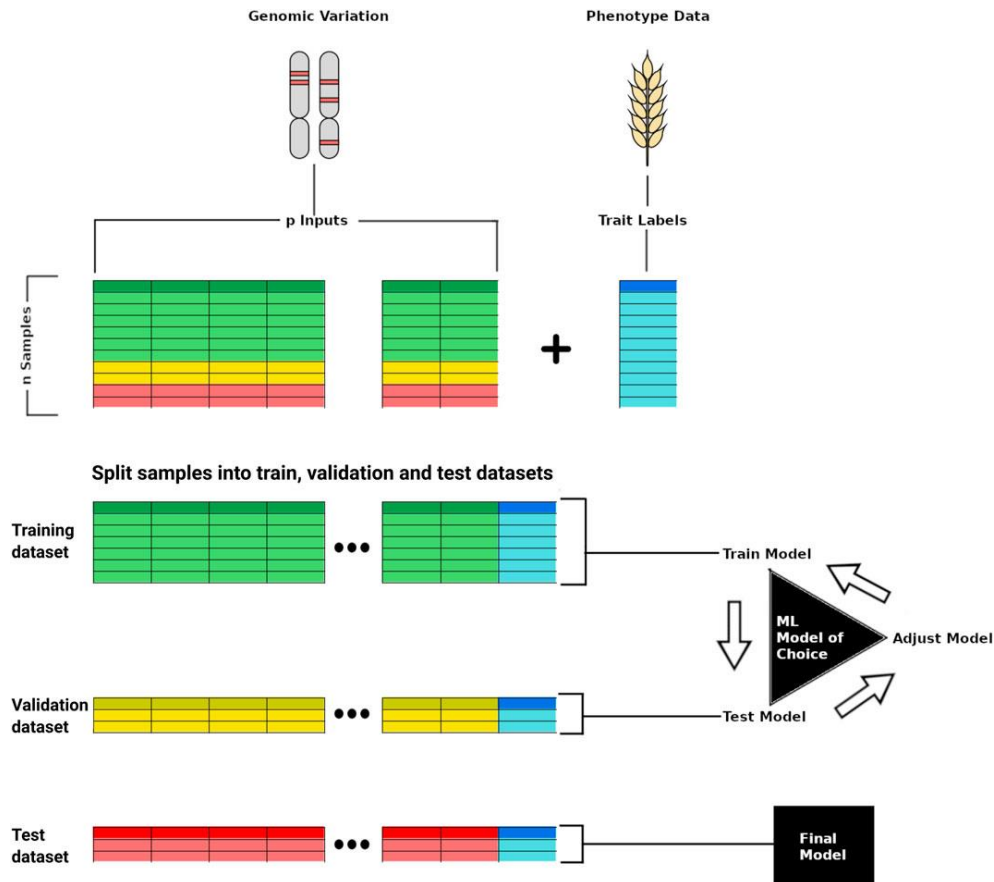
Beyond GWAS: Unlocking the Power of Genomic Selection

- GWAS and Pan-GWAS identify trait-marker associations across diverse accessions.
- Single-cell omics helps us understand in which cell types and contexts those variants act.
- **Translating markers into predictive breeding power** requires new strategies.



Genomic Selection (GS): uses genome-wide marker data to predict performance of untested accessions.

Genomic Selection Workflow in Practice

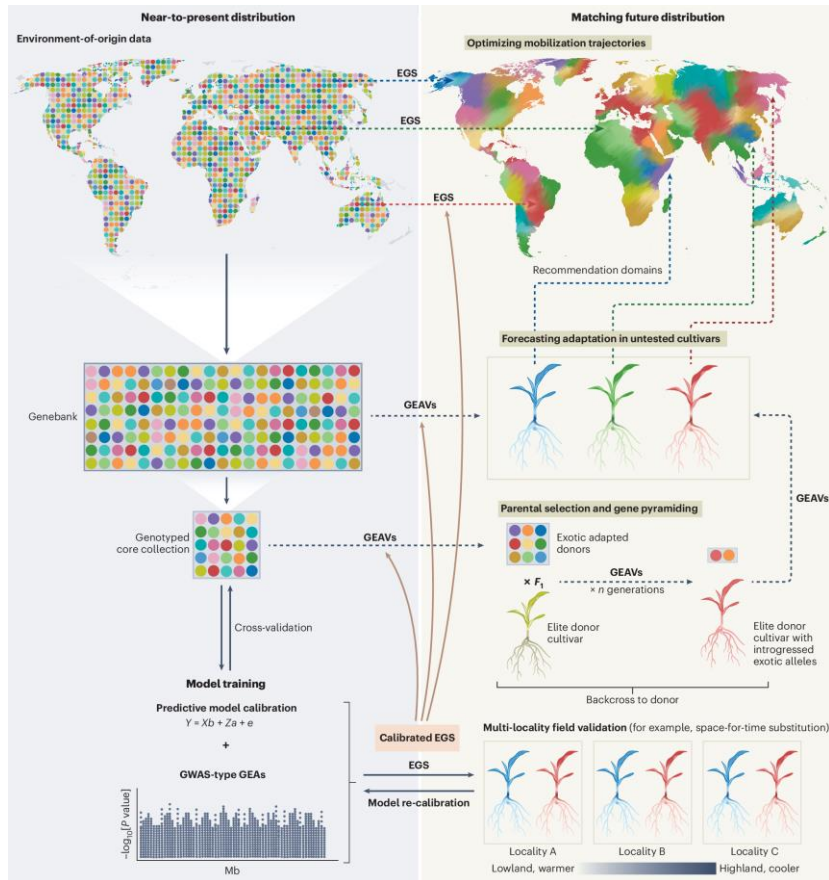


- Train Genomic and phenotypic data is split into training, validation, and test datasets.
- and adjust ML models to optimize prediction.



Select the final model for genomic prediction.

Take home messages



- **PGRs are reservoirs of untapped alleles** for resilience, adaptation, and nutrition.
- **Multi-omics, GWAS, Pan-GWAS and single-cell omics expand our power to identify and interpret trait associations.**
- **Genomic Selection** leverages genome-wide variation but **predicts breeding values early**.



The future of PGRs pre-breeding: **an AI-enabled multi-omics system that turns PGRs diversity into predictive power for future crops.**



**Thank you
for your attention**