



# A journey from classical to modern cassava breeding: Integrating genomic selection for faster and greater genetic gains

Danilo E. Moreta, Xiaofei Zhang, & Cassava Team  
Cassava Team Meeting  
Montería, Colombia | 21<sup>st</sup> November 2023

# Acknowledgements

**T**ogether  
**E**veryone  
**A**chieves  
**M**ore

## **Bioversity-CIAT Cassava Team**

Jonathan Newby + Research sub team + Management sub team (Ximena, Zulma, Oriana)

### **Breeding sub team**

Xiaofei Zhang  
Sean Fenstemaker  
Luis Fernando Delgado  
Lizbeth Pino  
Sandra Salazar  
Nelson Morante  
Jorge Iván Lenis  
Camilo Vargas  
Ana María Ensuncho

### **Genetics sub team**

Winnie Gimode  
Adriana Bohórquez  
Camilo Sánchez  
Luz Andrea Gómez  
Vianey Barrera (former)  
Janeth Gutiérrez  
Carmen Bolaños  
Carlos Ordóñez

### **Quality lab sub team**

Thierry Tran  
Luis Fernando Londoño  
María Alejandra Ospina  
Jorge Luna  
Jhon Larry Moreno  
Christian Duarte

### **Seed systems sub team**

Roosevelt Escobar  
Adriana Núñez  
Natalia Canacuan  
Auradela Ríos  
Carlos Dorado

### **Crop protection sub team**

Wilmer Cuellar  
Juan Manuel Pardo  
María Isabel Gómez

### **Field workers**

# Outline

## **PART I:** Background & context

- Transition from traditional to modern cassava breeding
  - Overview of the traditional cassava breeding method
    - Recurrent phenotypic selection
  - Genomic selection (GS) as a breeding tool
    - What it is? Why implement? How it works?

## **PART II:** Applications & plans

- Incorporating GS in CIAT cassava breeding pipeline
  - Work plan 2023-24, status, challenges & opportunities,
  - Discuss ideas & plans to optimize cassava GS models
  - A glance at the future: Hybrid cassava breeding

# Part I

## Transition from traditional to modern cassava breeding

# The traditional & long cassava breeding cycle

Selections based on phenotypic values

Improve population

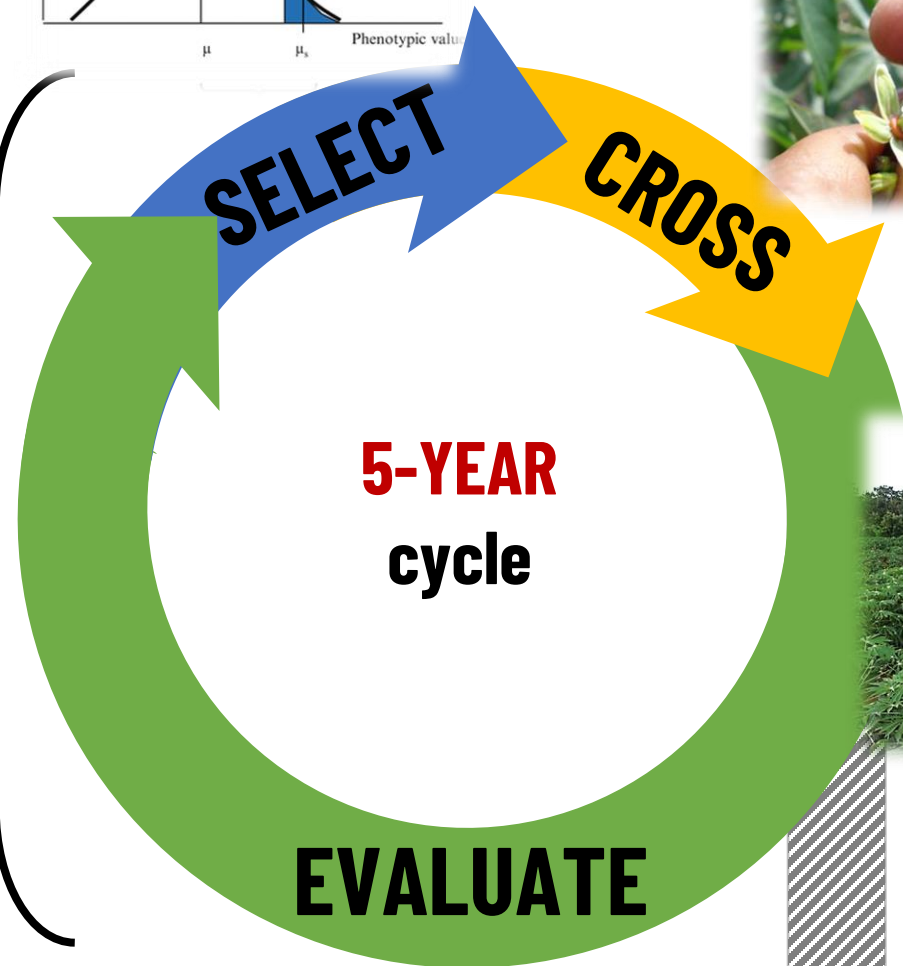
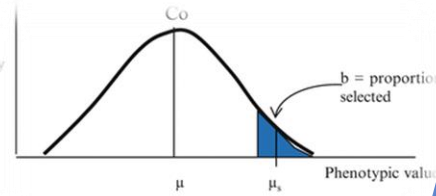


Image credits:  
G. Acquah, 2015; N. Morante; & D. E. Moreta

Release product (cultivar)

# A shorter breeding cycle with genomic selection

Selections based on genomic estimated breeding values (GEBVs)

$$GEBV = f(\text{DNA})$$

Improve population

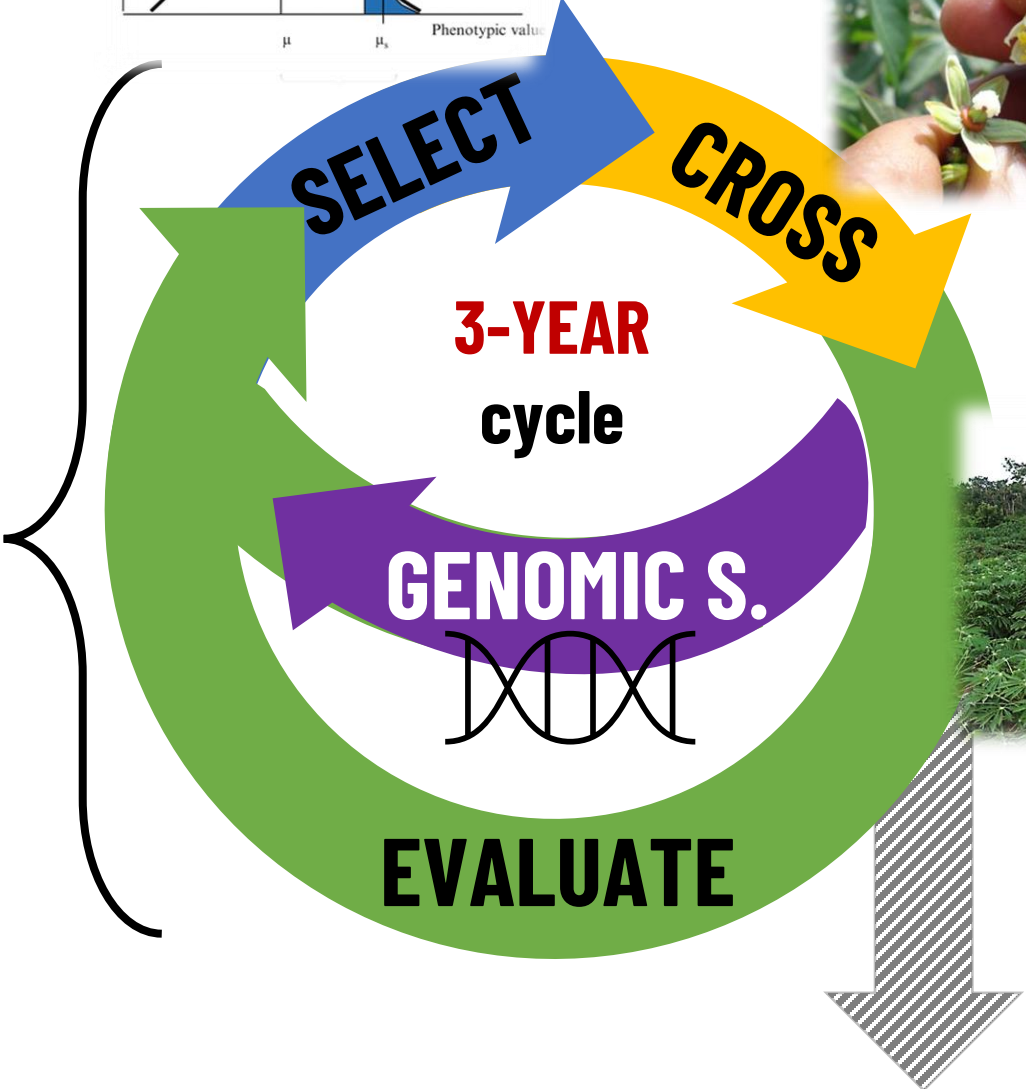
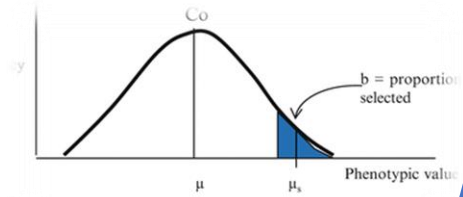


Image credits:  
G. Acquaah, 2015; N. Morante; & D. E. Moreta

Release product (cultivar)

# Genomic selection: The next (r)evolution in plant breeding

$$\Delta G = \frac{i \cdot r \cdot \sigma_a}{t}$$

$\Delta G$ : Genetic gain per unit time  
 $i$ : Intensity of selection  
 $r$ : Accuracy of selection  
 $\sigma_a$ : Additive genetic variation  
 $t$ : Duration of the breeding cycle

## Impacts of GS on the breeder's equation parameters



Introducing genomic information

Add genotyping to increase selection accuracy

$r$

Indirect impact  
 $i \ 1/t$

Reduce phenotyping effort with genotyping

$i \ 1/t$

Management of genetic diversity

Better choice of parents to optimize crossbreeding or preserve genetic diversity

$r \ \sigma_a$

Source: Fugerey-Scarbel et al., 2021

# Accelerated cassava breeding: The major benefit of genomic selection

*Genomic Selection* (GS, aka genomic prediction, genome-wide selection, genome-wide prediction, etc.) is a genomics-assisted breeding tool.

**F<sub>1</sub>**: seedling nursery

**F<sub>1</sub>C<sub>1</sub>**: cloned seedling nursery

**SRT**: single row trial

**PYT**: preliminary yield trials

**AYT**: advanced yield trials

**SIT**: seed increase trial

**TPY**: training population yield trials

**GWP**: genome-wide prediction

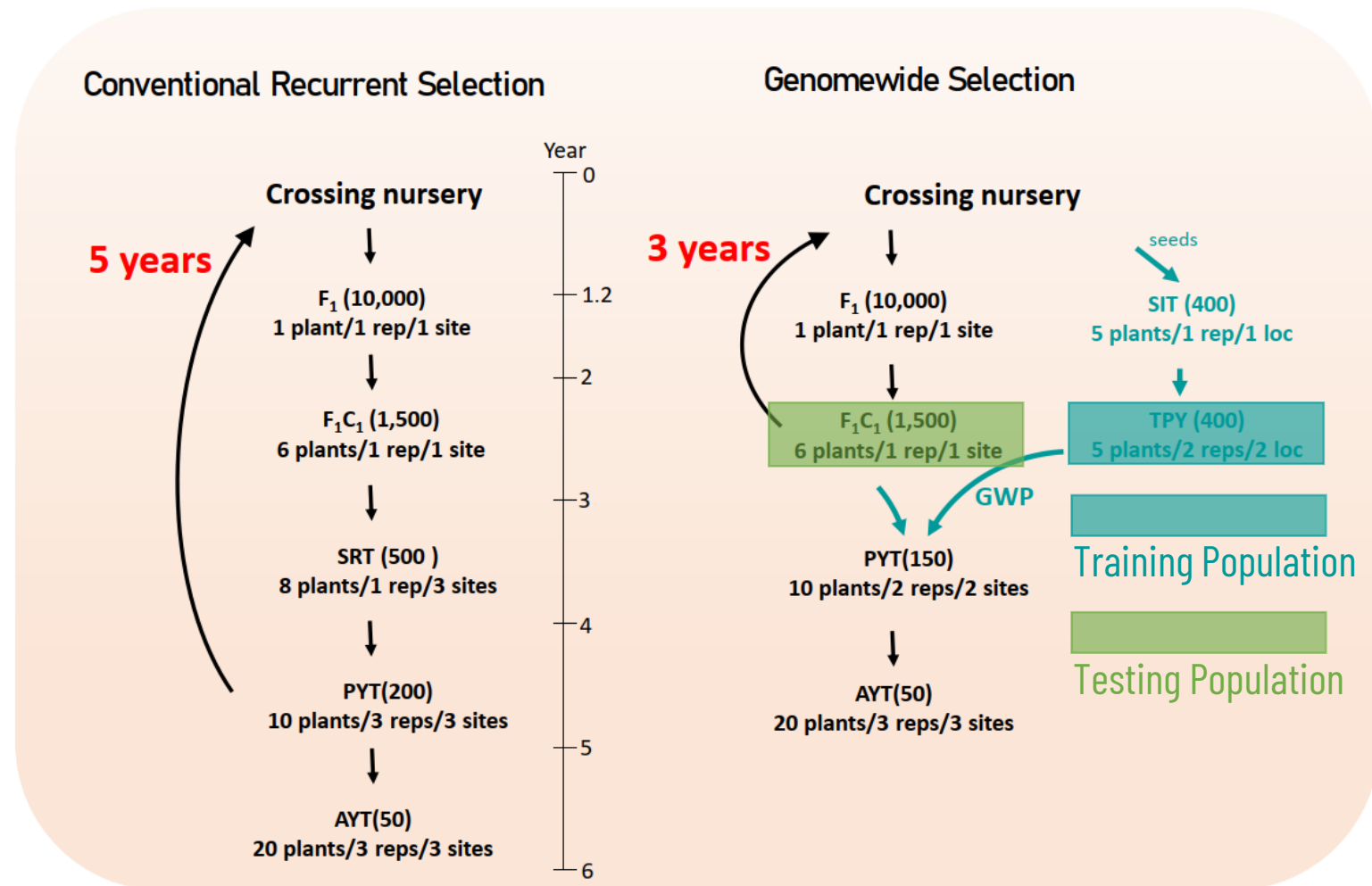


Figure credits: Cassava Breeding Team (Alliance Bioversity-CIAT)



# How genomic selection works?

$$GEBVs = f(\text{DNA})$$

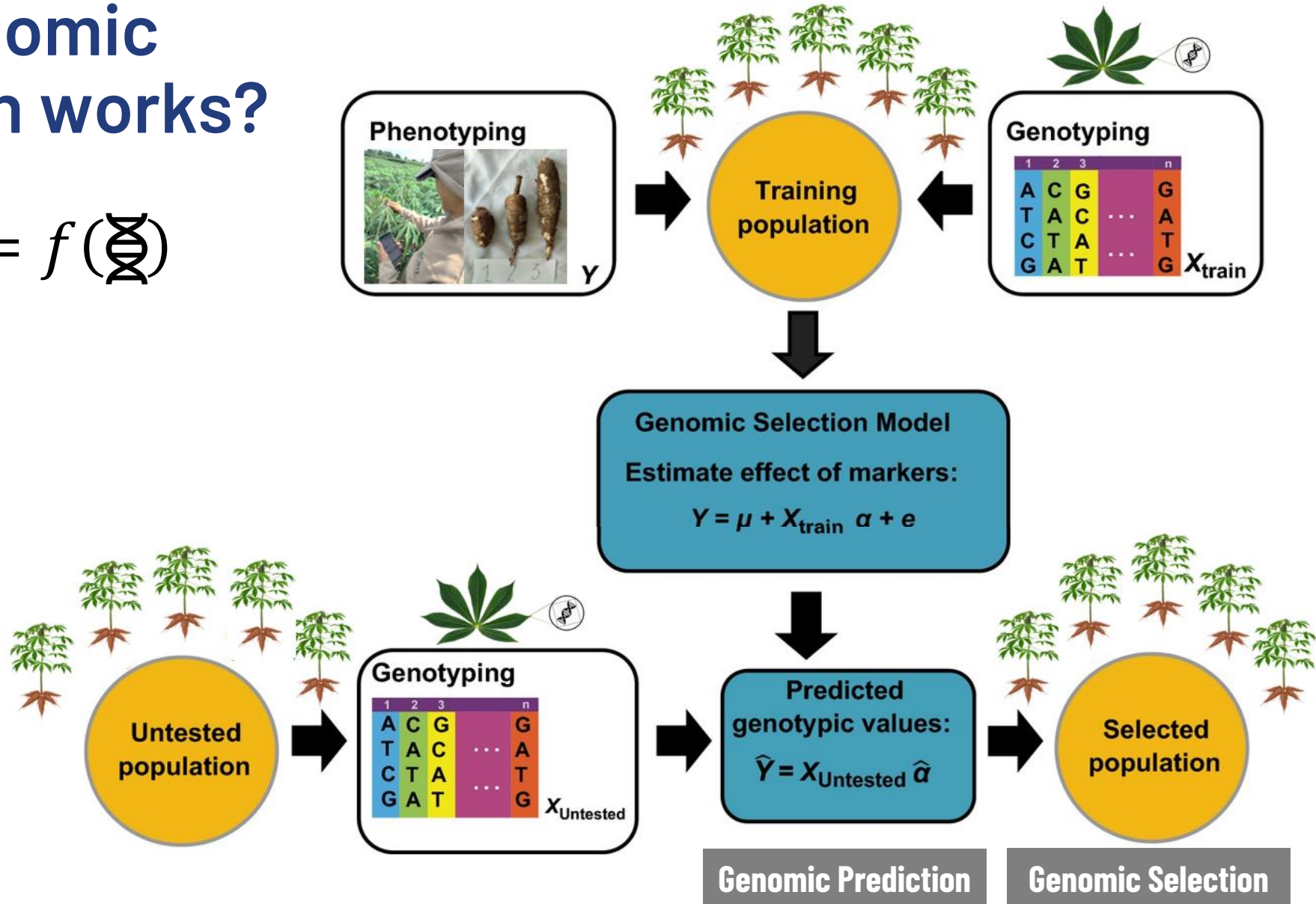


Figure adapted from Zhao et al., 2015 | Photo credits: S. Salazar

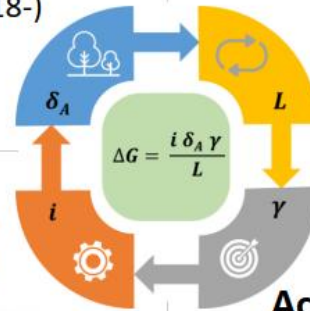
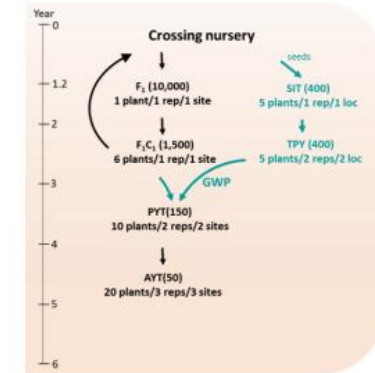
# GS as a breeding tool to accelerate cassava genetic gains

## Genetic Diversity

- New traits, e.g., CBSD res., CMD res., good cooking quality et al.
- Sequencing of progenitors (2020-)
  - Hybrid Breeding (2018-)

## Duration of Selection Cycle

- Flower Inducing (2016-)
- Genomewide Prediction (2019-)



## Intensity

- High throughput phenotyping

## Accuracy

- CassavaBase, Fieldbook & Barcode (2018-)
- Quality control and MAS (2020-)
- TPE, ≥ 2 Environments (2020-)
- ≥ 5 Checks, BLUP and GBLUP (2020-)
- Selection Index (2012-)
- NIRS & Image Analysis (2012-)
- Stage&Gate System (2020-)
- Operational Excellence (2019-)

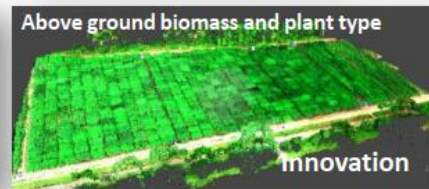
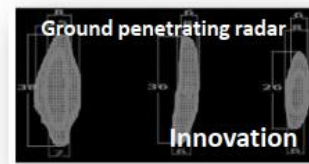


Figure credits: X. Zhang

# Why implement genomic selection?

$$GEBV = f(\text{DNA})$$

- Shorten lengthy breeding cycles
- Increase breeding efficiency: greater  $\Delta G$  per unit time
- Seed shortages
- Discard poor clones earlier
- Reduce manual labor (phenotyping)
- Hard-to-measure traits



# What factors affect GS accuracy?

- Marker density
- Size & composition of training population
- Number of QTLs
- Heritability
- Linkage disequilibrium (LD)
- Model used

$$Accuracy = corr(phenotype, prediction)$$

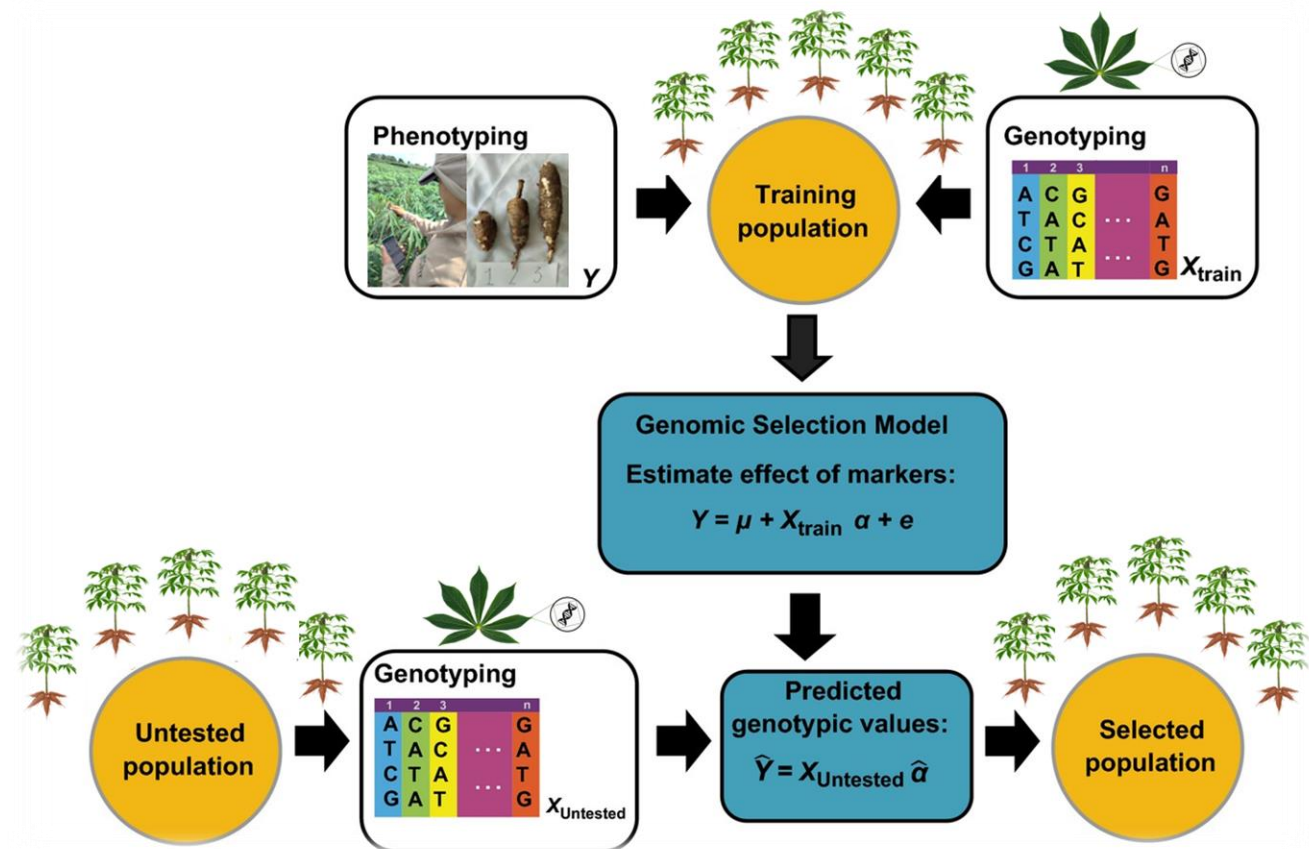



Figure adapted from Zhao et al., 2015 | Photo credits: S. Salazar

# Genomic selection models

*All models are approximations and hence wrong, but some are useful. ~ George Box*

The basic genetic model:  $P = G + E + (G \times E)$


$$GEBV_{GS} = f(\text{DNA})$$

## Ridge-regression best linear unbiased prediction (RR-BLUP)

Assumption: Infinitesimal model of genetic architecture (all markers have an equal effect)

**Step 1:** Estimate marker effects in training population (TP)

**Step 2:** Use marker effects & genotypes of selection candidates to predict GEBVs

## Genomic best linear unbiased prediction (G-BLUP)

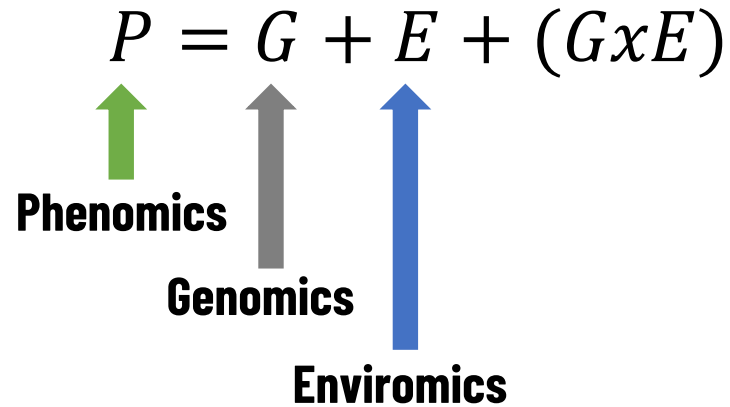
**Step 1:** Use markers to quantify genomic relationships

**Step 2:** Use genetic relatedness to TP of unevaluated to predict GEBVs

## Bayesian models

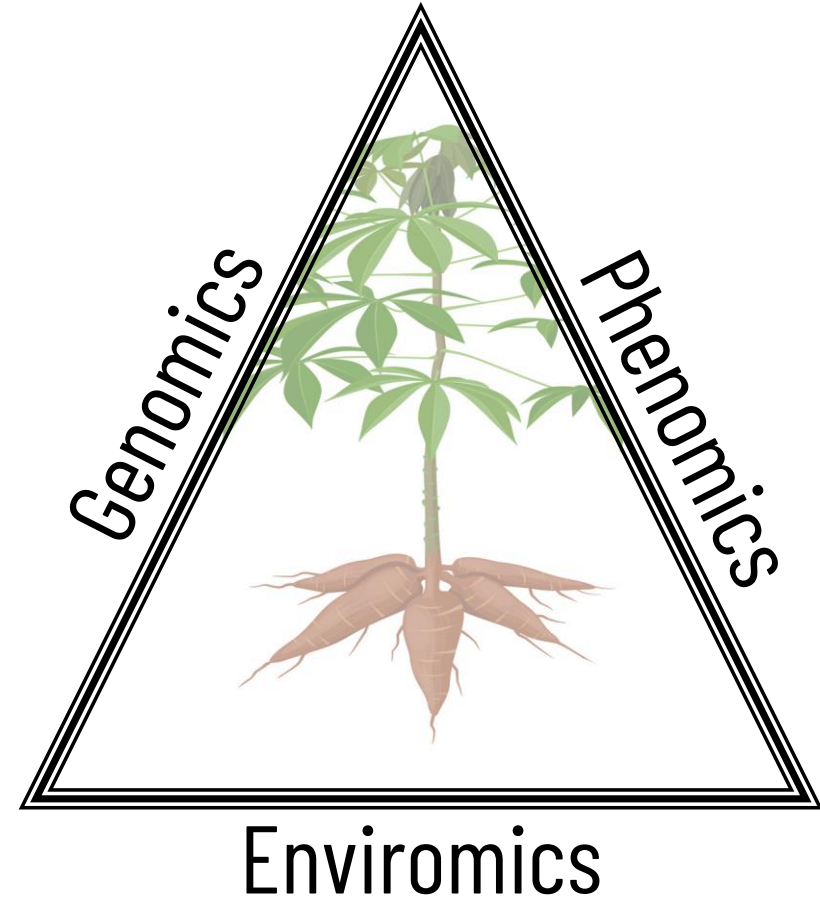
- Better model marker effects of differing sizes (Hayes, 2007)
- Separate variance estimated for each marker (Meuwissen et al., 2001)

# Modern plant breeding beyond molecular markers. Integrating and connecting disciplines

$$P = G + E + (G \times E)$$


Phenomics  
Genomics  
Enviromics

GS will not replace traditional breeding & field work, instead it will help optimize the system.



## **Part II**

# **Incorporating GS in CIAT cassava breeding pipeline**

# First steps to implementing GS in the cassava breeding pipeline at CIAT

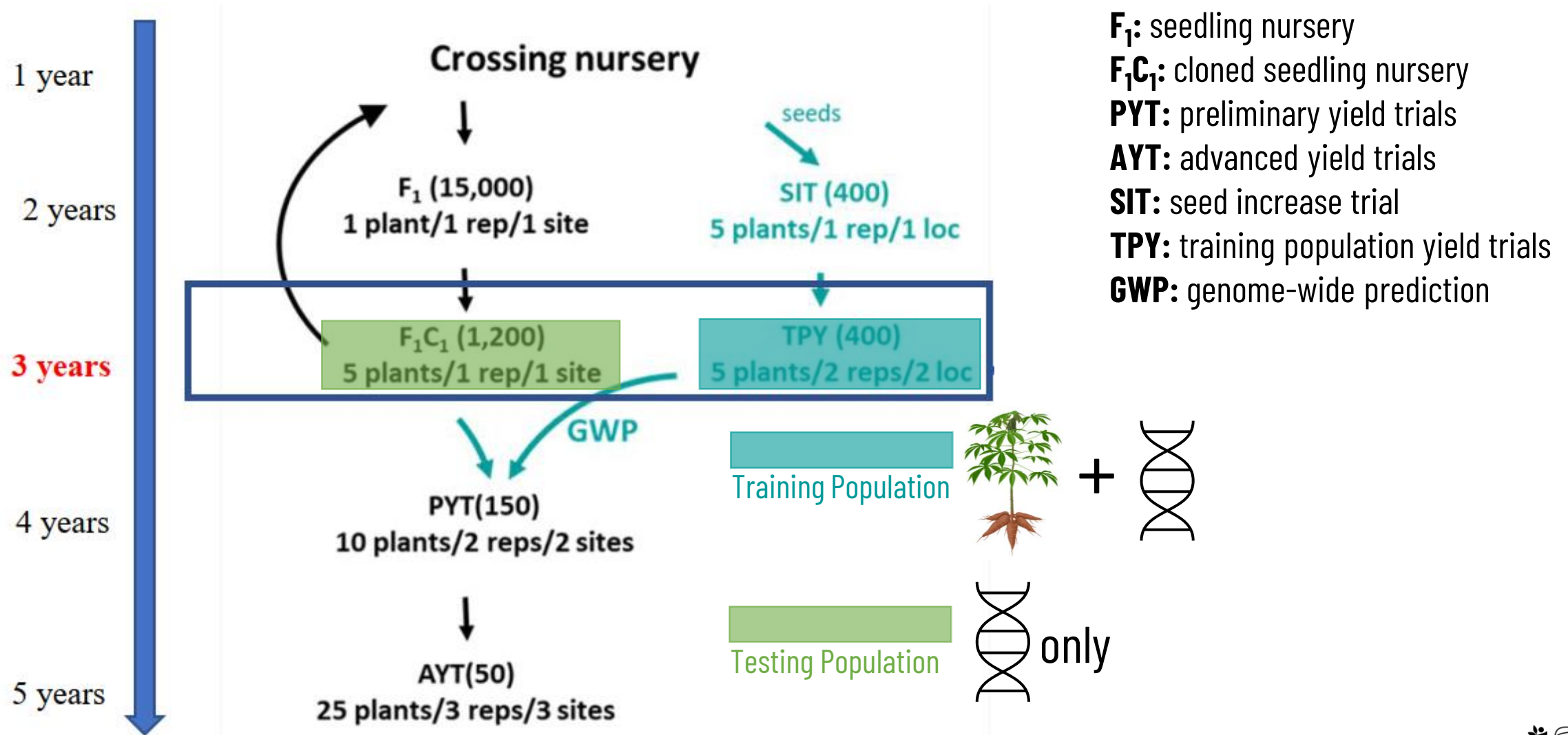


Figure credits: Cassava Breeding Team (Alliance Bioversity-CIAT)



# Training & testing (breeding) populations

## C1.1 (Cycle 1, Cohort 1) | 2020-2021



**Germplasm:** Full-sibs

**Trial names:** 2021DVGST (n = 399)      2020DVF1C (n = 651)

**Locations:** 3 locs

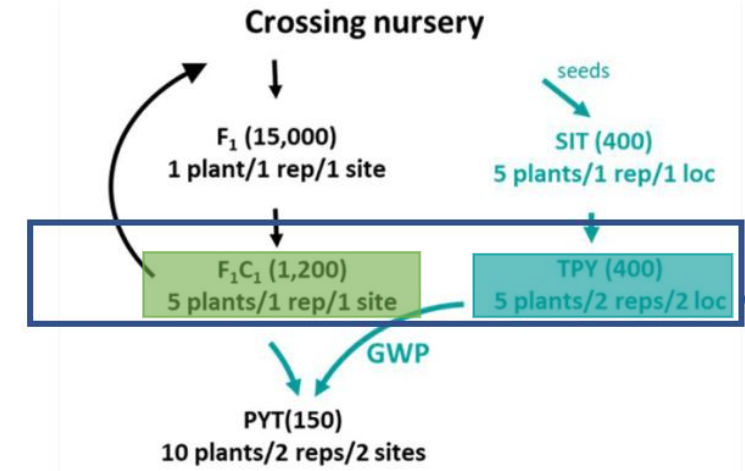
2 in north coast\*: Repelón  
Momil

Ag traits + Pest & disease resistance

1 in Palmira

Only quality traits

-----  
\*Semi-arid & sub-humid environments



# Training & testing (breeding) populations

## C1.2 (Cycle 1, Cohort 2) | 2021-2022

  
Training Population

  
Testing Population

**Germplasm:** Full & half-sibs

**Trial names:** 2022DVGST (n = ~873) 2021DMF1C (n = ~823)  
2021CQF1C (n = ~672)

**Locations:** 3 locs

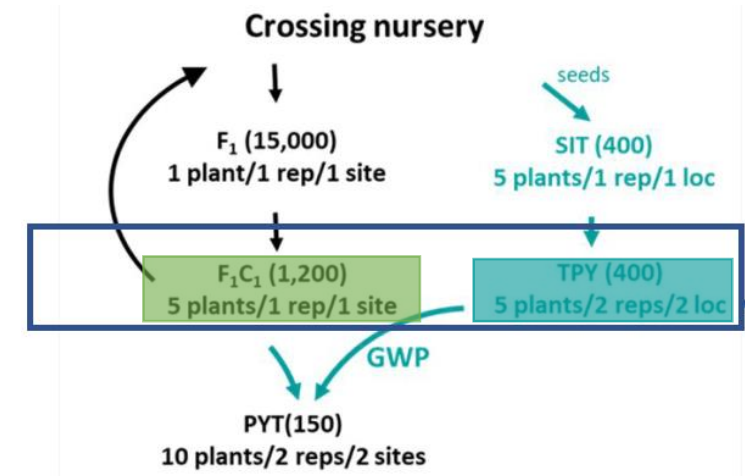
2 in north coast\*: Santo Tomás  
Momil

1 in Palmira

Ag traits + Pest & disease resistance

Only quality traits

-----  
\*Semi-arid & sub-humid  
environments



**GS for waxy cassava (COMING SOON!)**

n = ~200 clones

Ingredion agreement

# Genotyping the training & testing (breeding) populations

## C1.1 (Cycle 1, Cohort 1) | 2020-2021

$$\text{Phenotype} = \text{Genotype} + \text{Environment}$$

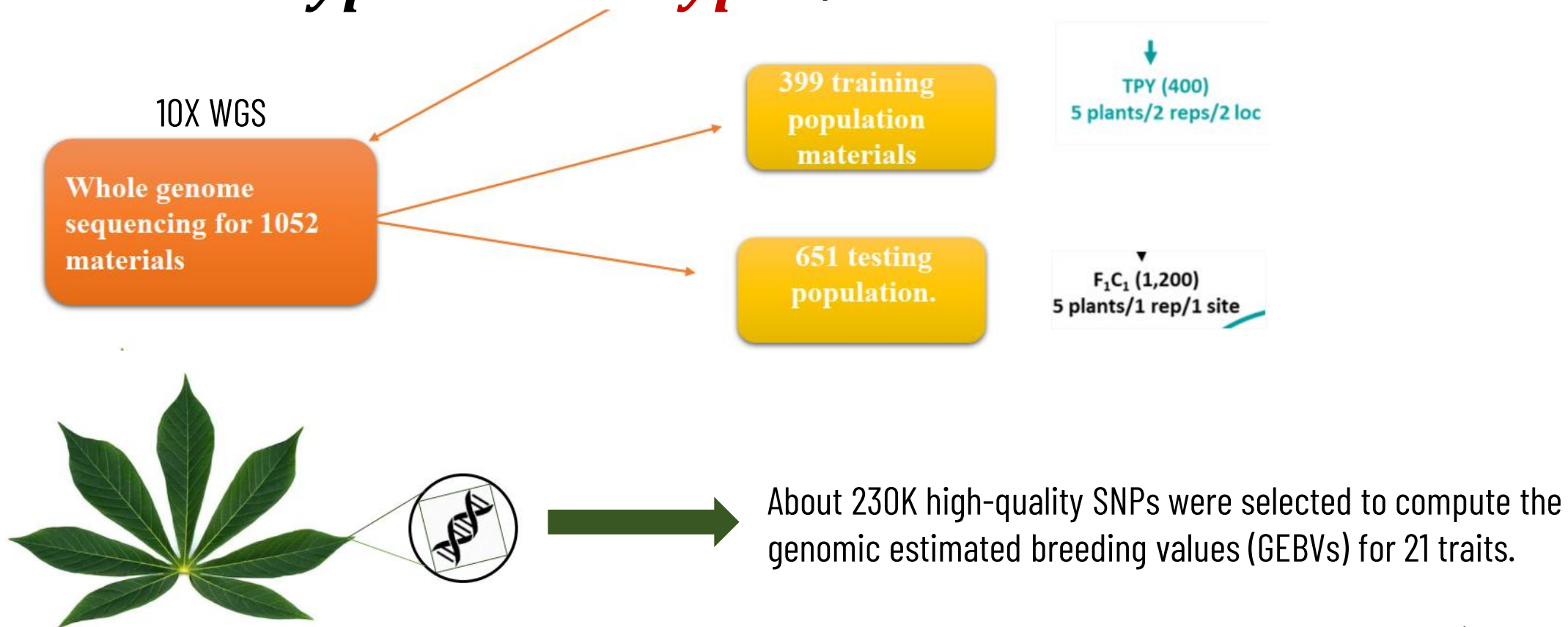


Figure credits: C. Vargas

# Prediction ability

## C1.1 (Cycle 1, Cohort 1) | 2020-2021

$$\text{Accuracy} = \text{corr}(\text{phenotype}, \text{prediction})$$

### Cross-Validation

How well do we predict individuals without phenotypes?

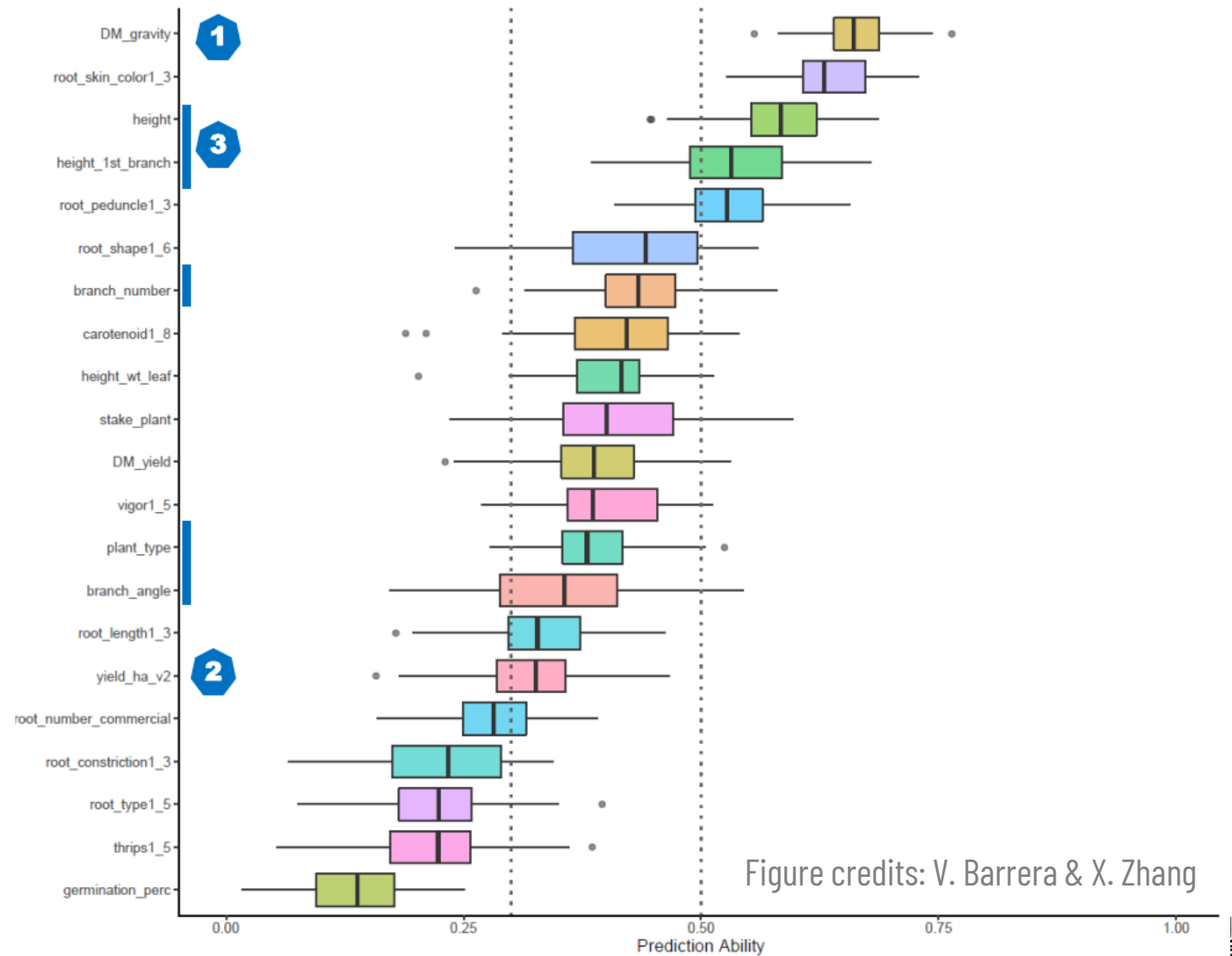
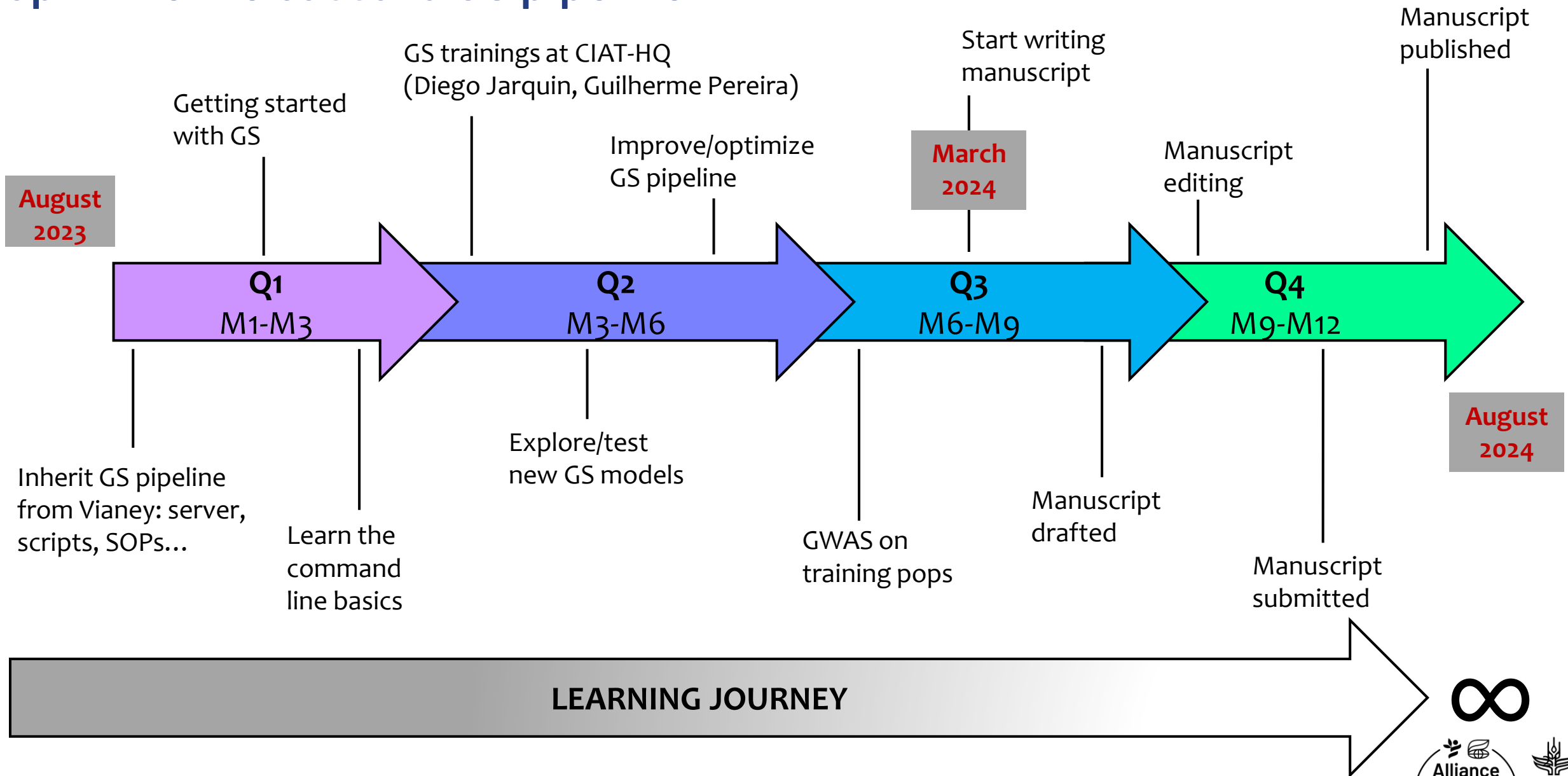


Figure credits: V. Barrera & X. Zhang

# My timeline/roadmap to implement & optimize the cassava GS pipeline

Q: quarter of a year (3-month period)

M: month



# Status: Trying to reproduce what others did

## 1 Examining VCF statistics in R

1.1 Setting up the R environment

## 2 Variant based statistics

2.1 Variant Quality

2.1.1 Phred quality score Q

2.1.2 Variant mean depth

2.1.3 Variant missingness

2.1.4 Minor allele frequency (MAF)

## 3 Individual based statistics

3.1 Mean depth per individual

3.2 Proportion of missing data per individual

3.3 Heterozygosity and inbreeding coefficient per individual

## VCFtutorial\_statsVisualization

*dEmc<sup>2</sup>*

Cali, CO | 17<sup>th</sup> October 2023

GUIDE/REFERENCE: [https://speciationgenomics.github.io/filtering\\_vcfs/](https://speciationgenomics.github.io/filtering_vcfs/)

*Notes<sub>dEmc<sup>2</sup></sub>*:

### DNA sequencing info

*10x whole genome sequencing*

VCF file source: [Group1\\_11.chromosome01.vcf](#)

$n_{samples} = 395$

$n_{SNPs} = 588,247$  (original/raw)

2021: Group 1-8

### Group

Genom

## 1 Set directory & paths

2 Read datasets

3 Preparing genotype data

3.1 Random sampling

3.2 Numerical scale

3.3 Imputation

4 Geno-Pheno match

4.1 Kinship

5 Prediction function

6 Run prediction

Learning by doing!  
Running demos & generating reports

## GP\_GS\_demo: rrBLUP

*dEmc<sup>2</sup>*

Cali, CO | 25<sup>th</sup> October 2023

GUIDE/REFERENCE: [GS\\_GEBVs\\_Prediction\\_Tutorial\\_dEmc2.jpynb](#)

*Notes<sub>dEmc<sup>2</sup></sub>*:

### Working files

1.  $Genotypes_{SNPs} = gs\_2023.reduced.txt$  | Should contain training (gsT) + breeding pops (gsB)

2.  $Phenotypes_{BLUPs} = GS\_2023\_second\_cohort\_blups.csv$  | Should be only gsT

3.  $Accessions_{names} = group12\_group12\_accessions.csv$  | Use to match genos & phenos

**NEED TO TRACE BACK ALL THESE FILES IN VIANEY'S FOLDERS. They :**

# Operational challenges & opportunities



## Challenges

Unorganized files on server  
(hard to track and find)

Lack of repository for GS  
scripts and related

IT micromanagement & lack  
of communication with  
researchers



## Opportunities




Systematic organization of files on server

Feed/update **Cassava2050 GitHub repository**  
(reproducible research)

Creation of a **server committee**: “Cassava byters”  
(Sean, Xiaofei, Winnie, Camilo, Danilo)

# Ideas to optimize (make it more accurate) the cassava GS models

## Phenos & environment

- Plug environmental covariates: weather & soil 
- Account for GxE effects   $\times$  
- Combine populations to increase power
- Include secondary traits (easier to predict?)
  - Relative yield
  - Total disease resistance
  - Harvest index, etc.
- High-throughput phenotyping for better model training (ongoing collab. w/ Mike Selvaraj)

$$P = G + E + (G \times E)$$

*Let's not miss these parts!*

$GEBV_{GS} = f(\text{DNA})$

## DNA markers

- Try different TP sizes & number of markers
- Dual purpose of TP: predictions + association (GWAS)
- Special treatment to some markers
  - Most significant & GWAS hits as fixed effects

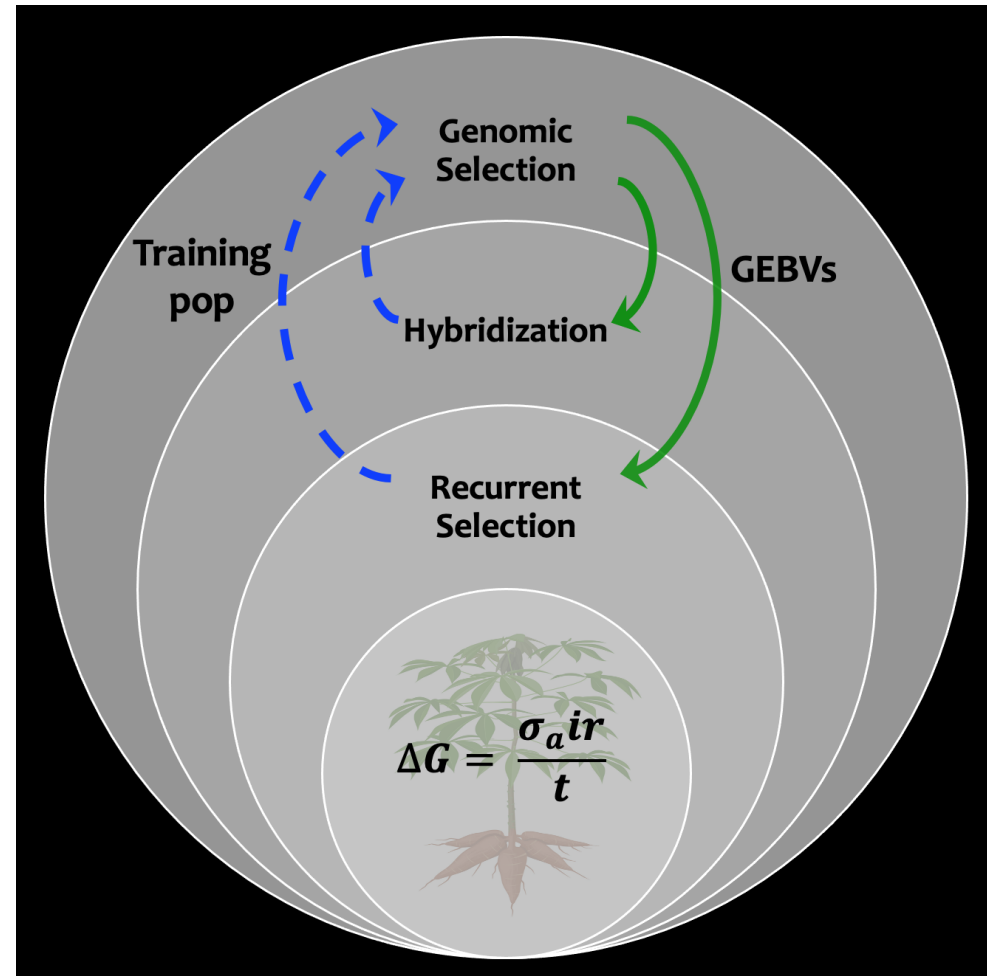


# Moving up cassava improvement & production to the next level through hybrid breeding

## How can GS help achieve this?

**Inbreds:** Combining ability

**Hybrids:** Performance



**Efficacy of  $\Delta G$  per unit time and cost**



**Thanks!**