# Pan-genome analysis of *Theobroma cacao* reveals new genes and provides new insights into the diversity of the species.

X. Argout[1,2], G. Droc[1,2], O. Fouet[1,2], A. Lemainque[3], B. Rhoné[1,2], G. Loor[4], C. Lanaud[1,2]

[1] CIRAD, UMR AGAP Institut, Montpellier, France
[2] AGAP, Univ. Montpellier, CIRAD, INRA, Montpellier SupAgro, Montpellier, France
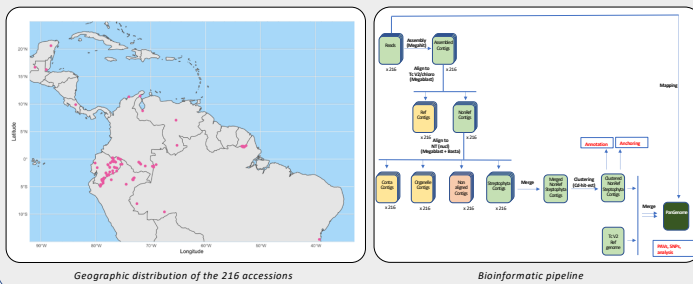[3] Genoscope, Institut de biologie François-Jacob, Commissariat à l'Energie Atomique (CEA), Université Paris-Saclay, Evry, France
[4] Instituto Nacional de Investigaciones Agropecurias, INIAP, Quito, Ecuador

While two reference genomes of *T. cacao* are available (B97-61/B2 and Matina1-6), a reference sequence cannot capture the entire gene content of a species owing to structural variants. In the recent years, high-throughput resequencing data of *T. cacao* genotypes have provided tools to discover allelic variants in the species or characterize the level of expression of genes, but much of the genotype-specific information is often lost by direct mapping of short sequence reads onto a single reference genome. For *T. cacao*, very little is known about the complete genomic content across the species.
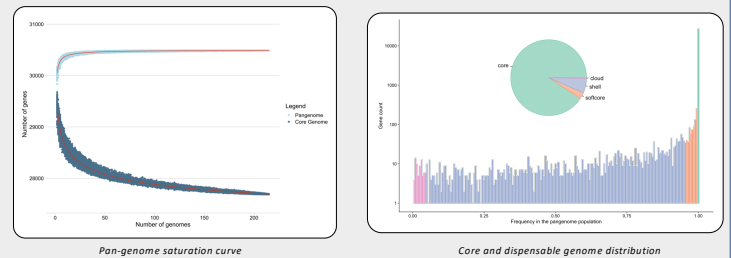To address this issue, we fully re-sequenced (58X) a collection of cocoa tree representative of the current known genetic diversity and built a pan-genome of *Theobroma cacao*. We comprehensively investigated gene PAVs inside the diversity of the species and captured new protein coding genes not included in the reference genomes. We also examined how the processes of diversification of the *Theobroma cacao* genetic groups have shaped their gene content and investigated putative protein-coding genes under selection during the domestication process of the Criollo group.

## ❖ Plant material and pan-genome construction



*Geographic distribution of the 216 accessions*



*Bioinformatic pipeline*

## ❖ Selection of gene PAVs during T. cacao diversification



*Genetic diversity of the population based on gene PAVs*



*Distribution of gene PAVs among the 15 groups*



*Gene space specificity for the 15 groups*



*Intersection between the gene space of the 15 groups*

**Key findings** :
- 15 genetic groups
- Accessions of Iquitos and Napo ancestry have the highest gene number
- Gene content of the Criollo group shows gene loss during domestication process
- Very few genes are specific to a genetic group
- Iquitos gene space overlap almost all genetic groups
- Criollo is most closely related to Caquetá group and Caquetá group contains more than 95% of the genes found in Criollo group

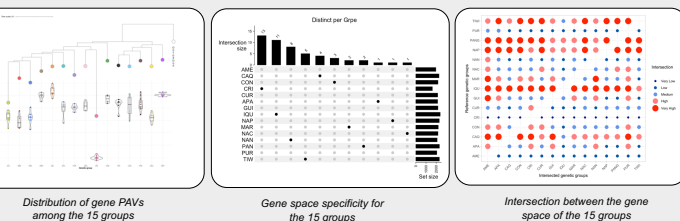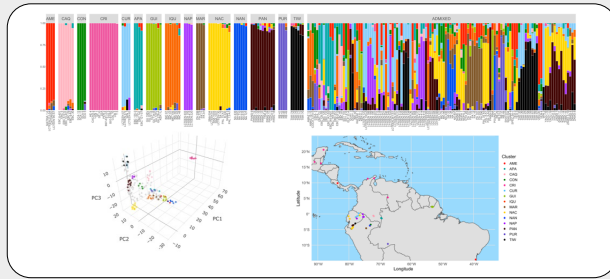## ❖ Assembly, Gene prediction, Core and dispensable genome

- The <u>non reference genome</u> dataset resulting from the assembly of 216 *T. cacao* accessions comprised **48.2 Mbp** and **1,407** candidate protein-coding gene models.
- The <u>*Theobroma cacao* pan-genome</u>, including the B97-61/B2 reference and non-reference genome, has a total size of **372.9 Mbp** and contained **30,489** protein-coding gene models.
- A vast majority of the genes, **27,687 (90,8%)**, are core genes, shared by all accessions while **2802** are variable genes.



*Pan-genome saturation curve*



*Core and dispensable genome distribution*



*Gene Ontology (GO) enrichment of protein-coding genes of the variable genome*

## ❖ Protein-coding genes under selection during Criollo domestication



*Geographic distribution of putative Criollo wild relative population*



*Gene frequencies comparison between Criollo accessions and accessions native from the putative center of origin of the domesticated Criollo.*

**Key findings** :
- 71 genes selected during Criollo domestication
- "Defense response" enrichment in non selected genes

## ❖ Conclusion / recommendation

- The *Theobroma cacao* pan-genome adds depth and completeness to the reference genomes, and is useful for future biological discovery
- The understanding of gene PAV, which contribute to trait variation, can support applications for the development of new varieties

Contact : xavier.argout@cirad.fr