



**HAL**  
open science

## The role of transposable elements on gene expression regulation in maize and its response to drought

Maud Fagny, Véronique Jamilloux, Johann Joets, Clémentine Vitte

### ► To cite this version:

Maud Fagny, Véronique Jamilloux, Johann Joets, Clémentine Vitte. The role of transposable elements on gene expression regulation in maize and its response to drought. 22nd national congress on transposable elements, Jul 2019, Lyon, France. hal-04330779

**HAL Id: hal-04330779**

**<https://hal.inrae.fr/hal-04330779v1>**

Submitted on 8 Dec 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

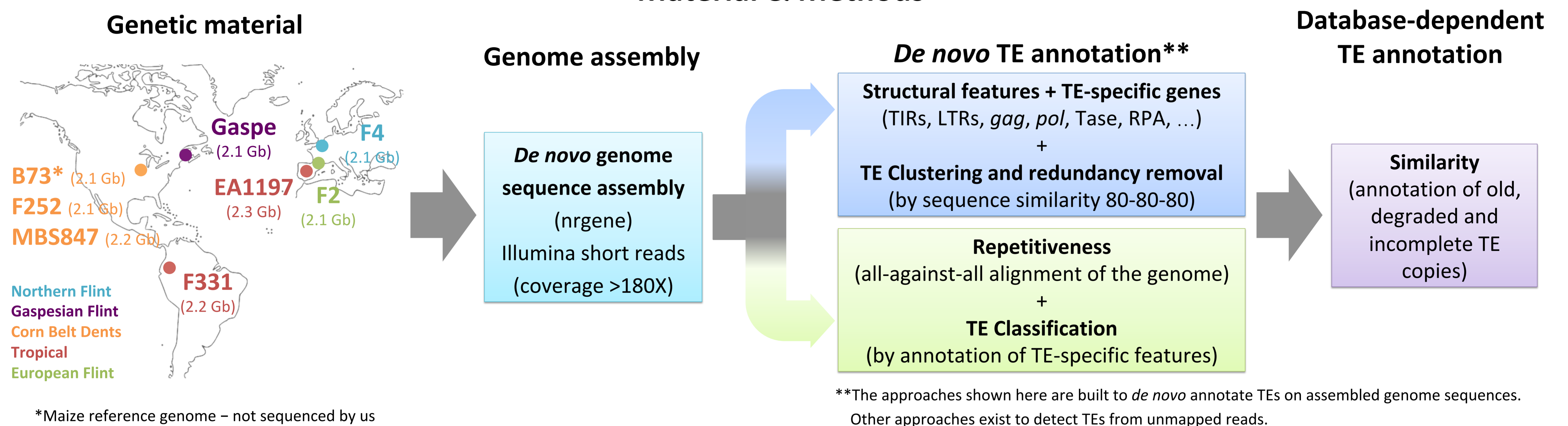
Maud Fagny<sup>1</sup>, Véronique Jamilloux<sup>2</sup>, Johann Joets<sup>1</sup>, Clémentine Vitte<sup>1</sup>

<sup>1</sup>Génomique Quantitative et Evolution – Le Moulon, INRA, Université Paris-Saclay, CNRS, AgroParisTech, Université Paris-Sud, Gif-sur-Yvette, France. <sup>2</sup>Unité de recherche Génomique Info, INRA, Versailles, France

## Background

Transposable elements (TEs) are major constituents of plant genomes and contribute to their dynamics and evolution. Existing softwares packaging tools using TE structural features detection and repetitiveness approaches have proven useful for TE annotation of small genomes. Attempts have been made to adapt them for large genomes analysis (Jamilloux *et al.*, 2016 - *Proceedings of the IEEE*). However, this remains challenging, mainly due to the amount of genomic data to analyze and to the structural complexity of these genomes. In particular, tools based on repetitiveness can be run only on a subset of the genome because of the required amount of memory and computation times. Tools based on structural features have been successfully applied on wheat (Wicker *et al.*, 2018 - *Genome Biology*) and maize (Anderson *et al.*, 2019 – *bioRxiv*). However, these studies were based on a succession of in-house scripts, and no pipeline has emerged so far to annotate all TE classes in large genomes. Consequently, while half of angiosperm species have a genome size above 1.6 Gb (Pellicer *et al.*, 2018 - *Genes*), only a handful of large plant genomes have been annotated so far. To better characterize the potential role of TEs in maize local adaptation, we aim at characterizing TE polymorphisms in seven maize inbred lines of contrasted origins. With its 2.1 Gb genome, we use maize as a model for TE annotation of an average angiosperm genome.

## Material & Methods



## Test of *de novo* TE annotation tools on the B73 genome

	Structural features	Repetitiveness	Similarity
Examples tested (not tested)	<i>LTRharvest, TEA, detectMITE, SINE-finder, HelitronScanner</i> ( <i>SINE_scan, MGEScan, LTR_retriever, ...</i> )	<i>REPET-TEdenovo</i> ( <i>RepeatScout, RepeatModeler</i> )	<i>RepeatMasker, REPET-TEannot</i> ( <i>TE-HMMER</i> )
All genomes	Independency from TE database	+	-
	Detection of old, incomplete TE copies	--	+/-
	Accuracy of TE boundaries	+	-
	Detection of highly variable TE superfamilies (Helitrons)	+	--
Large genomes	Detection of nested TE	++	+/-
	Sensitivity to low copy-number TE families	+	-
	Specificity (low false positive rate)	--	+/-
	Memory usage	+	--
	Computation time (32 cores, 96GB RAM) for maize B73 genome (2.1Gb)	Weeks	Months

Table 1: Comparison of different *de novo* TE annotation approaches

- Structural-based tools are relatively fast to run on large genomes, but have high false-positive rates and necessitate additional filtering steps.
- Repeats-based tools have a lower false-positive rate, but are slow and must be run on a subset of the genome (<300Mb), hampering the detection of low copy-number TEs, which may play an important role in adaptation.

## Conclusions

Finding the right accuracy/computing time balance to *de novo* annotate TEs in large plant genomes is challenging. A few softwares exist that integrate both structure-based and repetitiveness-based tools, such as *REPET-Tedenovo* (Quesneville *et al.*, 2015 - *PLoS Comp Bio*) or *PiRATE* (Berthelier *et al.*, 2018 - *BMC Genomics*). But their application to large genomes is hampered by the long computation time (up to months), the high amount of memory necessary to run the repetitiveness-based tools on the whole genome, and the absence of filtering steps to filter out the false positives generated by structural features-based tools.

## Future directions

### New pipeline based on structural features

- Package several existing structure-based tools
- Provide automatized, tunable filtering steps
- Include automatized TE clustering steps

### Modified integrated pipelines

- Modify data storage structures
- Improve memory usage for repetitiveness-based tools
- Add tunable filters to clean up structure-based results