
This is the **accepted version** of the article:

Castanera, Raúl; Vendrell-Mir, Pol; Bardil, Amélie; [et al.]. «Amplification dynamics of miniature inverted-repeat transposable elements and their impact on rice trait variability». *The Plant Journal*, (April 2021). DOI 10.1111/tpj.15277

This version is available at <https://ddd.uab.cat/record/241555>

under the terms of the  **CC BY-NC-ND** license

DR JOSEP CASACUBERTA (Orcid ID : 0000-0002-5609-4152)

Article type : Original Article

The amplification dynamics of MITEs and their impact on rice trait variability.

Raul Castanera^{1*}, Pol Vendrell-Mir¹, Amélie Bardil¹, Marie-Christine Carpentier², Olivier Panaud², Josep M. Casacuberta^{1*}

¹ Centre for Research in Agricultural Genomics CSIC-IRTA-UAB-UB, Campus UAB, Edifici CRAG, Bellaterra, 08193 Barcelona, Spain.

² Laboratoire Génome et Développement des Plantes, UMR CNRS/UPVD 5096, Université de Perpignan Via Domitia, 52 Avenue Paul Alduy, 66860, Perpignan Cedex, France.

* Corresponding authors

Running title

Impact of MITEs on rice trait variability

Keywords

MITE, transposable elements, GWAS, traits, rice, transposition, genetic factor

Corresponding authors details:

- Raúl Castanera

raul.castanera@cragenomica.es

Centre for Research in Agricultural Genomics CSIC-IRTA-UAB-UB, Campus UAB, Edifici CRAG, Bellaterra, 08193 Barcelona, Spain.

- Josep M. Casacuberta

josep.casacuberta@cragenomica.es

This article has been accepted for publication and undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/TPJ.15277](https://doi.org/10.1111/TPJ.15277)

This article is protected by copyright. All rights reserved

Centre for Research in Agricultural Genomics CSIC-IRTA-UAB-UB, Campus UAB, Edifici
CRAG, Bellaterra, 08193 Barcelona, Spain.

Accepted Article

Summary

Transposable elements (TEs) are a rich source of genetic variability. Among TEs, Miniature Inverted-repeat Transposable Elements (MITEs) are of particular interest as they are present in high copy numbers in plant genomes and are closely associated with genes. MITEs are deletion derivatives of class II transposons, and can be mobilized by the transposases encoded by the latter through a typical cut-and-paste mechanism. However, MITEs are typically present at much higher copy numbers than class II transposons.

We present here an analysis of 103,109 Transposon Insertion Polymorphisms (TIPs) in 738 *O. sativa* genomes representing the main rice population groups. We show that an important fraction of MITE insertions has been fixed in rice concomitantly with its domestication. However, another fraction of MITE insertions is present at low frequencies. We performed MITE TIP-GWAS to study the impact of these elements on agronomically important traits and found that these elements uncover more trait associations than SNPs on important phenotypes such as grain width. Finally, using SNP-GWAS and TIP-GWAS we provide evidences of the replicative amplification of MITEs.

Introduction

Transposable Elements (TEs) are an essential component of plant genomes. Their capacity to amplify and create new genetic variability by insertion/excision, and the possibility that their multiple copies offer for recombination, make them a rich source of genomic variants that can be selected through evolution (Tenailon *et al.*, 2010). TE-induced mutations include gene knock-outs, but also the induction of gene epigenetic silencing or changes in gene regulation by inactivating enhancers or repressors upon insertion or by adding new regulatory elements contributed by the TE (Lisch, 2013).

Miniature Inverted-repeat Transposable Elements (MITEs) are short non-coding TEs thought to be deletion derivatives of class II cut-and-paste TEs (Feschotte *et al.*, 2002). However, differently to the class II TEs they derive from, MITEs are found in high-copy number in plant genomes (Chen *et al.*, 2014). The apparent contradiction between a conservative cut-and-paste mechanism of transposition and MITEs high copy number was highlighted short after MITE discovery (Feschotte *et al.*, 2002; Casacuberta and Santiago, 2003), but the mechanism by which MITEs amplify is still obscure.

MITEs are tightly associated with plant genes (Santiago *et al.*, 2002; Lu *et al.*, 2012; Benjak *et al.*, 2009), and examples of these elements potentially altering gene expression have accumulated over the years (Santiago *et al.*, 2002; Lu *et al.*, 2012; Naito *et al.*, 2009; Yang *et al.*, 2005; Xu *et al.*, 2020; Zheng *et al.*, 2019; Yin *et al.*, 2020). Moreover, it has recently been shown that MITEs frequently contain Transcription Factor Binding Sites (TFBS) in plants, which may allow their mobilization to alter transcriptional networks by rewiring new genes (Morata *et al.*, 2018), and that they can induce structural variability through aberrant transposition events (Chen *et al.*, 2020).

In the last few years, different methods for analyzing TE insertion polymorphisms (TIPs) in resequenced genomes have been developed (Vendrell-Mir *et al.*, 2019), and LTR-retrotransposon (LTR-RT) TIPs have recently been used to perform GWAS to study LTR-RT dynamics in rice (Carpentier *et al.*, 2019a), and the genetic basis of agronomic traits in rice and tomato (Akakpo *et al.*, 2020; Domínguez *et al.*, 2020). These recent publications show that LTR-RT TIPs can allow discovering associations not seen with conventional GWAS strategies based on SNPs and highlight the importance of analyzing the fraction of the genetic variability TEs account for. MITEs have not been used in TIP-GWAS as they were considered as more difficult to analyze due

to their small size and high copy number (Domínguez *et al.*, 2020). However, recent benchmarking efforts of TIP prediction tools have allowed us to propose efficient approaches to analyze MITE insertions (Vendrell-Mir *et al.*, 2019). Here we used these approaches to identify MITE TIPs in 1,059 rice varieties and used them to perform TIP-GWAS on different agronomic traits. We show that, indeed, MITE TIPs reveal an additional fraction of the genetic variability associated to important crop traits and may allow discovering the underlying causal genes for some traits. In addition, we used GWAS to study MITE dynamics in rice highlighting the impact of MITEs on rice domestication and breeding and suggesting that MITEs amplify by a replicative mechanism from a reduced number of MITE copies.

Results

1. TE annotation and selection of rice varieties

Rice, as all other crops, has been continuously improved since its domestications, giving rise to different groups of varieties frequently linked to different geographical origins. Japonica and Indica are the two major types of rice cultivars that originated from two distinct origins of domestication. Japonica varieties can be further subdivided in tropical and temperate. Two additional groups of varieties have been described, the aromatic group, closely related to the Japonica varieties, and the Aus/Boro group of varieties that have been shown to originate from a third domestication (Carpentier *et al.*, 2019b). The availability of assembled genomes corresponding to the different major groups of rice varieties allowed us to perform a TE annotation taking into account the potential intraspecific differences. We performed a stringent TE annotation using the REPET package (Flutre *et al.*, 2011). As MITEs are particularly prevalent in plants (Chen *et al.*, 2014), we complemented this annotation with a targeted annotation of MITEs based on MITE-hunter (Han *et al.*, 2010) and the already published annotations in the P-MITE database (Chen *et al.*, 2014). We used this pipeline to annotate three assembled rice reference genomes belonging to the three major rice subgroups of varieties (Indica, Japonica and Aus), and clustered all TE consensus sequences to obtain a global TE library of the species. We identified a total of 821 complete TE families belonging to all major TE orders, excluding incomplete and chimeric elements (Supporting Figure 1, Supporting dataset S1). In spite of the stringency of the approach used to build the TE library, a RepeatMasker (<http://www.repeatmasker.org>) annotation with the 821 consensus sequences on the

Nipponbare reference genome (MSU7) (Kawahara *et al.*, 2013) showed highly congruent results in comparison to recently published rice TE annotations. More specifically, 95% of the annotation overlaps rice 6.9.5.liban annotation (Ou *et al.*, 2019) (78% overlap in the opposite direction), and the classification agreement at the order level between reciprocal consensus of the two libraries was 99%. However, the annotation used here has fewer LTR-RT consensus (112 vs 389), probably due to the high stringent approach that did not retained partial, low-copy number and degenerated LTR-RT families, and has much more MITE consensus (400 vs 173) due to the dedicated tools used.

In order to study the TE dynamics in rice, and in particular that of MITEs, we took advantage of the availability of resequencing data of 3000 rice varieties to look for TE insertion polymorphisms (TIPs) in the different varieties of the TEs annotated in the reference genomes. A recent benchmarking exercise indicated that PopoolationTE2 (Kofler *et al.*, 2016) is a good tool for this purpose, but also that coverage is a key factor than can limit the detection of TIPs (Vendrell-Mir *et al.*, 2019). For this reason, we selected the 1,059 rice genomes sequenced at 15x or more (Supporting Table S1) and subsampled all genomes to 15x in order to be able to perform an unbiased comparisons of TIP abundance between varieties.

We found an average of 40,568 insertions per variety. We filtered the insertions using a zygosity cut-off of 0.7 as recently recommended to avoid false positive calls (Vendrell-Mir *et al.*, 2019). After applying this filter, we obtained an average of 18,463 insertions per variety.

The number of TIPs per variety is variable (Figure 1), but in general CAAS varieties had a significantly lower number of TIPs than IRIS varieties (Supporting Figure S2), despite the identical processing pipeline followed to sequence the selected varieties (Li *et al.*, 2014). Although we could not identify the reason for such a difference, we decided to use only the IRIS data, which reduced the number of rice varieties analyzed to 738.

2. TE dynamics in rice varieties

We detected a total of 103,109 TIPs in the 738 rice varieties corresponding to 32,449 retrotransposons and 70,660 DNA transposons of all orders. PCA analysis and phylogenetic trees performed based on predicted TIPs are congruent with the previously defined groups of rice varieties based on SNPs (Wang *et al.*, 2018) (Figure 1A, Figure 1B). The analysis of TIPs shows that most TE families have similar mean number of TIPs in the four main varietal groups

(Japonica, Indica, Aus and Aro), with a Coefficient of Variation between varietal groups (CV) lower than 20% in 82% of the families, and with a median of 9.5% (Figure 2A). However, a few number of families are much more variable, the CV reaching a maximum of 128.2%. Up to 150 TE families have a coefficient of variation between varietal groups higher than 20% and have potentially experienced a varietal group-specific amplification (Figure 2A). Figure 2B shows an example of one of these families where the number of insertions per variety is clearly higher in Indica and Aus varieties compared with Japonica and aromatic varieties. Among the different TE orders, LTR-retrotransposons are the only order significantly enriched among the group-specific amplification families (Fisher's exact test p-value = 0.0062, Table 1). This result is consistent with recent reports suggesting that rice LTR-RT have been actively transposing *in agro* (Carpentier *et al.*, 2019a). However, a group-specific amplification is not the only indicator of activity, as a recent activity in all groups of varieties would not result in changes in the CV. An analysis of the total number of different insertions of each family with respect to the mean number of insertions per variety shows that families belonging to all TE orders have been recently active (not shown). This is particularly clear for MITEs and LTR-RTs (Figure 2C), where small and medium-size families respectively seem to have been particularly active recently.

In order to get more insight on the dynamics of the different types of TEs we analyzed the TIP frequency distribution for each type of TEs in a subset of 382 varieties that correspond to traditional varieties and are a good representation of the variability of the species (Gutaker *et al.*, 2020). As shown in Figure 3A, the frequency distribution is different for different orders of TEs. In particular, and as previously shown (Carpentier *et al.*, 2019a), our study reveals that most LTR-RT TIPs are present at very low frequency (LF), which could suggest both a recent activity and a high turnover of LTR-RT insertions. On the contrary, although there is an important number of MITE TIPs present at low frequency, there is a significant number of them that appear to be fixed in the population, presenting a "U-shape" frequency spectrum with an excess of high-frequency (HF) derived mutations at the expense of middle-frequency variants. This suggests that, whereas most LTR-RT insertions are recent, rice has retained an important fraction of older MITE insertions. This can also be seen when analyzing the fraction of TIPs corresponding to insertions present in rice (*Oryza sativa* ssp. *Japonica* cv. Nipponbare) and nine different *Oryza* genomes (Zhou *et al.*, 2020) (Figure 3B). Whereas less than 10% of both *gypsy* and *copia* LTR-RT insertions present in rice (Nipponbare) are also present in *Oryza glumaepatula*, more than 30% of the MITE insertions are also present in this relatively distant *Oryza* genome (Figure 3B). This data

clearly suggests that an important fraction of MITE insertions have been fixed during rice evolution. As MITEs tend to concentrate close to genes (Casacuberta and Santiago, 2003; Feschotte *et al.*, 2002) these fixed insertions may have played an important role in the evolution of rice genes. On the other hand, an analysis of the frequency within rice varieties of the insertions shared between Nipponbare rice and *O. glumaepatula*, shows that these relatively old insertions have been maintained in the whole population (Figure 3C), confirming that, as already proposed (Feschotte *et al.*, 2002), most MITE insertions are highly stable, in spite of the possibility of being excised by related transposases (Chen *et al.*, 2020).

An analysis of the chromosomal distribution of MITE TIPs shows that it is similar to that of *copia* LTR-RTs (and different from that of *gypsy* LTR-RTs) and also follows that of genes (Figure 4A and Supporting Figure S3). This is expected as *copia* LTR-RTs and MITEs are known to be closely associated to plant genes (Casacuberta and Santiago, 2003). However, a comparison of the distribution of TIPs present at low- and high-frequency (1st and 4th frequency quartiles, respectively, that represent frequency cutoffs of 6.6% and 85.5% for MITEs, 0.2% and 25.3% for *copia* LTR-RTs and 1.3% and 30.8% for *gypsy* LTR-RTs), which allow starting to discriminate between the role of target specificity and selection in shaping TE distribution, shows clear differences between the two TE groups. Whereas only low frequency *copia* LTR-RTs insertions follow gene distribution, MITE insertion association with genes is stronger for high frequency insertions (Figure 4, Supporting Figure S4). This suggests that whereas in general both *copia* LTR-RTs and MITEs seem to target genic regions for integration, LTR-RTs are progressively cleaned away from these regions whereas MITEs inserted close to genes are more frequently maintained. A detailed analysis of the non-reference insertions with respect to genes shows that MITEs concentrate in the 5' and 3' proximal regions (<500 nt) of, and in particular in the 5' upstream regions, and that this was found for both high frequency and low frequency insertions (Figure 4B). Therefore, MITEs could have had an important impact on gene variability in the recent evolution of rice. As shown in Figure 3B, although a significant fraction of MITE insertions are present in the rice wild relative *Oryza rufipogon*, more than 60% of them seem specific of domesticated rice. We analyzed the MITE insertions that are present at HF in rice varieties and are absent from *O. rufipogon* as they could potentially be linked to domestication mutations, and found that the fraction of insertions close to genes is significantly higher than for the whole MITE insertions here described (63% vs 54%, Fisher exact test $p=0.027$). Among the 834 genes that have a MITE insertion present at HF in domesticated rice and absent from wild rice, there are

some already characterized rice genes with important functions during development or under stress conditions and that may be related to important agronomic traits (Supporting Table S2).

3. Association of MITE insertions with trait variability.

Genome-wide association studies (GWAS) using LTR-RT TIPs as a genotype, instead of SNPs, have recently been performed in rice and tomato (Domínguez *et al.*, 2020; Akakpo *et al.*, 2020). These studies showed that these TIPs can reveal additional genetic associations with traits that are not seen with SNPs. Here we wanted to explore the potential of MITE TIPs for TIP-GWAS. We reasoned that their small size, high copy number, and close association with genes, could make MITEs particularly suited for this purpose. We performed GWAS using MITE TIPs as genotype and the phenotypic data available for the 451 Indica varieties (Mansueto *et al.*, 2017; Jackson, 1997). We used Indica varieties as they are present in a higher number in our dataset, which provide us with more power for the GWAS analyses. We obtained significant associations between specific MITE insertions and the eight phenotypes analyzed (grain width, length and weight, salt injury at EC12 and EC18, flowering time, leaf length and panicle shattering Supporting Table S3). Some of the associations obtained were coincident with peaks obtained using SNPs as genotype, suggesting that both were in Linkage Disequilibrium (LD) with the causal mutation (that could also be the SNP or the MITE TIP themselves). However, in many cases the MITE TIP-GWAS revealed different regions allowing to explore additional genetic regions potentially linked to trait variability (Supporting Table S3). In order to evaluate the potential of TIPs to reveal additional associations with respect to SNPs, we analyzed the LD of all MITE TIPs with surrounding SNPs (located at 100 kb upstream and downstream from each MITE TIP) and compared it to the LD among those SNPs. Figure 5A shows that most MITE TIPs are in lower LD with surrounding SNPs than SNPs with other SNPs. As we used more SNPs than TIPs, this may somehow influence this analysis. For this reason, we repeated the analysis subsampling the SNPs dataset to have a comparable number between SNPs and TIPs. As it is shown in Figure 5A, the clear difference of LD between TIP-SNPs and SNP-SNPs is not the result of a different number of variants between the two datasets. As LD may vary along chromosomes, we also performed a local comparison of the LD of MITE TIPs with SNPs located close by with the LD among those SNPs. Figure 5B shows that the majority of MITE TIPs (89%) show a mean r^2 for their linkage to the surrounding SNPs that is lower than the mean of the r^2 of the SNPs to other

SNPs located in the same region. This strongly suggests that MITE TIPs can reveal additional genetic variability with respect to SNPs and highlights their potential for GWAS.

Indeed, among the 51 TIPs associated with eight different traits, 29 are not in LD with surrounding SNPs ($r^2 < 0.2$), and 21 of those do not have any surrounding SNP (100 kb upstream and downstream) significantly associated with the trait. On the other hand, we have detected TIP associations with four traits for which there is not a single significant association with SNPs (resistance to salt injury EC12 at and EC18, flowering time and leaf length) (Supporting Figure S5). Interestingly 27 out of the 51 TIPs associated with traits are located within or close to previously described QTLs for the traits, which gives more support to the associations found.

As an example, we present the GWAS analysis of the grain width phenotype. Whereas the SNP-GWAS revealed a single genomic locus linked to grain width, the MITE TIP-GWAS revealed nine different regions with significant association TIPs (Figure 6, Supporting Table S3). The different number of associations obtained with TIPs and SNPs was not related to the different number of genetic data used in both analyses (153,744 SNPs and 34,023 TIPs) as a random subsampling of SNPs to a number similar to that of MITE TIPs did not modified the result obtained (Supporting Figure S6).

The only SNP association obtained was found at the well-characterized GSE5 gene of the GW5 locus, which is associated with rice grain size variation (Duan *et al.*, 2017). The leading SNP was located at 3.9 Kb of the GSE5 gene. Interestingly, the strongest association using MITE TIPs, GW-MITE 19 ($p\text{-value} = 1.02e^{-23}$) (Figure 6) corresponded to a MITE insertion located 350 nt upstream of the GSE5 gene. A 367bp insertion located at this position was previously described as being in strong LD with a deletion located 4,500 bp upstream the GSE5 gene, and the co-occurrence of these two structural variants constitute one of the three main haplotypes of GSE5 in cultivated rice ($GSE5^{\text{DEL1+IN1}}$) (Duan *et al.*, 2017). However, our data suggests that, in spite of the MITE insertion being closer to the GS5 gene than the deletion, the MITE insertion is not the causal mutation responsible of grain width, as the few varieties carrying the MITE but not the deletion have narrow grains (Supporting Figure S7).

In addition to the association related to the GSE5 gene, we detected eight more regions with MITE insertions associated with grain width. An analysis of the leading TIP shows that in most cases the TIP is in low LD with the surrounding SNPs (Figure 6). Interestingly, several of these TIP are located close to genes or within QTLs already characterized as linked to grain phenotypes. The GW-MITE 5 insertion is located at 33 Kb of the OsARG gene, which contributes to grain yield

variation in rice (Ma *et al.*, 2013), and five additional MITE TIPs are within a previously characterized QTLs linked to grain yield (GW-MITE 3 at qGY-3 (Mao *et al.*, 2003) and GW-MITE 22 at QTARO QTL-450 (Yonemaru *et al.*, 2010)), embryo length (GW-MITE 1 at qEML-2 (Dong *et al.*, 2003)), grains per panicle (GW-MITE 4 at gpp4 (Xiao *et al.*, 1996)), grain shape (GW-MITE 22 at qLW-6 (Yan *et al.*, 2003)) and 1000 kernel weight (GW-MITE 21 at QKw5 (Li *et al.*, 1997)), which reinforces the relevance of the associations found. A detailed analysis of the MITE insertion positions showed that four out of these nine significant MITE TIPs were located inside genes or in their close vicinity (1000bp upstream or downstream) (Supporting Table S3). We anticipate that these genes could be good candidates to be linked to the grain width trait.

4. Genetic factors linked to MITE amplification

In addition to study the genetic basis of agronomic traits, GWAS can also be used to study the genetic determinants of TE activity, as recently done for LTR-RTs in rice (Carpentier *et al.*, 2019a). Here we use this approach to study the genetic determinants of MITE amplification. This is particularly relevant as although MITEs are known to be mobilized by transposases of class II related elements (TIR-TEs), a canonical cut-and-paste mechanism does not usually result in the high-copy numbers that MITEs frequently attain in eukaryotes, and particularly in plant genomes (Chen *et al.*, 2014).

We performed a GWAS analysis on 451 rice varieties belonging to Indica subspecies using a SNP matrix of 153,744 SNPs as genotype, and the MITE copy number of 400 MITE families as phenotype, running one association analysis per TE family. In order to directly check for the potential role of transposase-encoding TEs related to MITE families in their amplification, we analyzed 18 MITE families showing significant sequence similarity with TIR-TEs from which they probably derive (Supporting Table S4). We identified peaks of SNPs significantly associated with the copy number of seven of those 18 MITE families. None of these SNPs peaks had a TIR-TE insertion with significant sequence similarity with the MITE family at less than 100kb. Interestingly, five of the associated SNPs peaks (for four families) had an insertion of a MITE of the same family at less than 100kb (Supporting Table S4). As a complementary approach, we performed TIP-GWAS using transposase-coding TIR-TE TIPs (4,898) or MITE TIPs (34,023) as genotype. TIR-TE TIP-GWAS did not reveal any association of a particular TIR-TE with the copy number of the related MITE family. On the contrary, a TIP-GWAS performed with MITE TIPs

revealed significant associations of a MITE insertion with the MITE family copy number for 11 out of 18 MITE families, including the four families for which we detected an associated SNP peaks with a MITE of the same family located close. As an example, Figure 7 shows the analysis of the genetic factors associated with the copy number of the MITE family SE260300235fam318_632. The SNP-GWAS revealed two different SNP peaks strongly associated with the MITE family copy number (Figure 7A). MITE TIP-GWAS revealed several peaks of MITE TIPs significantly associated with the MITE family copy number (Figure 7B), and in all cases the leading TIP corresponded to a MITE of the same family (in green in Figure 7B). The leading TIPs of chromosomes 2 and 6 are closely linked to the corresponding leading SNPs of the associations detected by SNP-GWAS ($r^2 = 0.89$ and 0.98 , Figure 7E), and their presence strongly correlates with an increase of copy number of the MITE family (Figure 7D). On the contrary, the TIR-TE TIP-GWAS did not reveal any association of the TIR-TE family closely related to the MITE family (see Figure 7C), the hAT family DTX_comp_IRGSP_B_R1932, with the amplification of the MITE family (Figure 7F).

We extended this analysis to the 400 MITE families identified here and found significant association SNP peaks for 175 out of the 400 MITE families. The SNPs peaks associated with MITE family amplification were distributed along the 12 chromosomes without a particular enrichment in any genomic region (Supporting Table S5). None of these associated SNPs had a TIP related to a TIR-TE with similarity (even limited similarity at the TIR level) to the corresponding MITE family. On the contrary, up to 37% of these peaks (66) have a MITE TIP of the same family at less than 100kb (Supporting Table S6). In 96% of the cases, the presence of this MITE copy is linked to an increase of MITE copy number of the respective family. We used three assembled *Indica* varieties from the recently published platinum genomes (Zhou *et al.*, 2020) to analyze the insertions corresponding to the TIPs located close to the 66 SNP peaks associated to MITE copy number and in 48 out of 49 cases where the TIP is predicted to be present in one of the platinum genomes, we confirmed the presence of the MITE insertion in the genome assembly.

As a complementary approach, we performed GWAS using MITE TIPs or TIR-TE TIPs instead of SNPs as genotype. We detected associations of a MITE TIP with the copy number of 301 different MITE families (75% of them). In most of these cases (186) the most significant MITE TIP belongs to the same MITE family for which they are potentially affecting the copy number. The high number of associations with MITE insertions of the same family strongly suggests that this particular MITE copy is probably at the origin of the increase in copy number of the MITE family.

On the contrary, the TIR-TE GWAS gave only 51 associations between the presence of a TIR-TE TIP (39 different elements) and a MITE family copy number, but in this case only two of these 39 TIR-TE elements have significant similarity (even limited to the TIRs) with the corresponding MITE family. However, in these two cases the analysis of these TIPs suggested that they are due to related MITE insertions and not to TIR-TEs.

Most of the families for which we have detected the association of the presence of a particular MITE with its family copy number are relatively small families and the differences in copy number are also small. We reasoned that in big families the presence of older elements or elements arisen from different waves of amplification could mask the effect of GFs on recent amplifications. We therefore decided to repeat the TIP GWAS for the 20 biggest families (mean TIP number/variety of more than 50, and more than 500 TIPs at $MAF > 1\%$, to take into account both the number of insertions in each variety and the number of different insertions in the population) using as a phenotype only the number of TIPs present at a frequency lower than 10%. We detected particular MITEs of the same family associated with the number of recent insertions for 18 of these 20 families. In all cases, the presence of this particular MITE is correlated with an increase in recent insertions from 1.4 to 2.9 fold (Table 2).

In summary, our results show a positive association of the presence of particular MITE copies with MITE copy number, which may suggest that MITEs amplify by a replicative mechanism of few "master" MITE copies, as it is often the case for retrotransposons.

In order to check whether TIP-GWAS reliably identifies "master" copies of replicative TEs, we decided to look for the genetic determinants of rice LINE amplification, as "master" LINES have particular structural characteristics. Indeed, LINES are transcribed from an internal promoter located in their 5' end (Swergold, 1990), absent from the vast majority of elements due to 5' deletions that make them inactive (Farley *et al.*, 2004). This is the consequence of a transposition mechanism that frequently leads to the integration of elements truncated in 5', due in part to the low processivity of the RT and to the microhomology-facilitated recombination during integration (Martin *et al.*, 2005). Therefore, most newly transposed elements are incapable of expression and transposition, and their structure is different from that of the few "master" elements that are at the origin of each LINE family. It is therefore relatively straightforward to differentiate a "master" LINE from the rest of the copies of the same family.

We performed a GWAS analysis on 104 LINE families and we found 133 significant peaks corresponding to 79 LINE families (79% of the total families, Supporting Table S5). 57% of the

significant peaks (76 peaks from 51 families) had a TIP of the same family at less than 100kb of the corresponding SNP (median distance to SNP = 30.7 Kb, Supporting Table S6). In 99% of these cases, the varieties with the TIP have higher copy number than varieties without the TIP ($p < 0.05$, two-tailed Wilcoxon rank sum test), suggesting that, indeed, the presence of this particular copy of the TE could be at the origin of the increase in LINE copy number.

We compared the structure of the LINE copies identified as potentially responsible for the amplification of a LINE family with that of the rest of the copies. As an example, Figure 8A (top) shows the Manhattan plot corresponding to the family RIX_comp_MH63_B_R3107_Map6. The three major peaks in chromosomes 8, 11 and 12 identify positions contain a TIP of the same family at positions chr08:5,711,756, chr11:23,993,235 and chr12:1,525,695, respectively. These TIPs (hereafter referred as genetic factors, GF) are in strong linkage disequilibrium (LD) with the minor SNP allele (ie, GF1 and GF2) or with the major SNP allele (ie, GF3), which may explain the association of the SNP with the LINE family copy number. The association of these three LINE insertions with the copy number of this LINE family can also be seen in a TIP-GWAS performed using LINE TIPs as genotype (Figure 8C). Moreover, the presence of the three LINE copies is associated with an increase of copy number of this LINE family (Figure 8B), suggesting that indeed, the three LINE copies may correspond to "master" LINE elements at the origin of this LINE family. In order to validate the presence of these insertions and study the structure of these three potential "master" LINES, we used the sequence of the Indica varieties from the recently published platinum genomes (Zhou *et al.*, 2020). The element corresponding to GF-1 is present in the LIMA::IRGC 81487-1 genome and both GF-1 and GF-3 are present in the genome of the IR 64 variety. A comparison of these elements with all the other copies of the same family present in these genomes shows that they are the only complete ones, the other copies being truncated at 5' (Figure 8D). All these results suggest that the associations detected point to the active (or "master") elements that are at the origin of the most recent replicative amplification of these families.

These results confirm that GWAS is able to identify the "master" elements at the origin of a replicative amplification of a TE family and strongly suggest that the particular MITE copies found associated with MITE copy number are "master" elements whose replicative amplification is at the origin of MITE families.

5- Molecular characteristics of MITEs linked to MITE copy number

In order to characterize the MITE copies involved in MITE amplification, we compared the sequence of 20 of these MITEs present in the assembled Indica rice LIMA::IRGC 81487-1, as well as of their flanking sequences, with that of the rest of the members of their respective families. We analyzed potential differences in GC content and Minimum Free Energy Structure of these elements including flanks of different length (see methods), as well as TE content of the regions, but failed to detect significant differences between the MITE copies involved in MITE amplification and the rest of the copies of their respective families (Supporting Table S7).

In order to analyze the possible influence of different chromatin characteristics we used the data available for the Nipponbare reference genome to look for the presence in this genome of the MITE TIPs characterized here as associated to the copy number of their respective families. We were able to identify 43 of these MITEs in Nipponbare and we compared them with all the annotated copies of the corresponding MITE families in the whole genome (Table 3). The characteristics of the chromatin associated to these MITEs seem different than that of the rest of the copies of their respective families. In particular, they show a potential enrichment in transcription activation marks, with 33% more H3K4me3 and 14% more DNase I hypersensitive peaks coinciding with these specific insertions, as compared with the rest of the copies of their respective families. On the other hand, a recent analysis has characterized the meiotic recombination hotspots in Indica rice mapping them to the Nipponbare genome (Marand *et al.*, 2019). An analysis of the overlap of these recombination hotspots with MITE TIPs shows that 18.6 % of the 43 MITEs characterized here as associated to the copy number of their respective families overlap with these recombination hotspots whereas only 12.2 % of all MITEs overlap with these sites.

Discussion

MITEs and the evolution of rice genome

The analysis of 103,109 TIPs in 738 rice varieties, which include 382 traditional varieties that represent the variability of the species, allowed us to study the dynamics of the rice mobilome. Our data shows that TEs have been active during the recent evolution of rice and their insertions allow discriminating the different varietal group that have been defined based on SNPs.

Interestingly, different types of TEs show a different dynamics in rice. LTR-RTs insertions are present at very low frequency. This result is in agreement with a recent analysis showing that rice LTR-RTs have been active in agro (Carpentier *et al.*, 2019a), but contrasts with another recent study that found higher population frequencies for LTR-retrotransposons (Kou *et al.*, 2020). In line with what the authors of this latter study discuss, we think that these differences are likely due to the inclusion in that study of partial and truncated LTR-RTs and probably old elements (not included in the dataset used here), likely found at higher population frequencies. In contrast to LTR-RTs, MITE show a "U-shape" frequency distribution, with insertions found at low frequency but with an important fraction that seems almost fixed. Almost 40% of the insertions are present in the wild *O. rufipogon* and more than 30% in the relatively distant *O. glumaepatula* genome. This suggests that whereas most LTR-RT insertions are rapidly eliminated, MITE insertions are frequently retained. The smaller size of MITEs, as compared with LTR-RTs, could make their insertions less deleterious and more easily tolerated, and, as MITEs are preferentially found close to genes, the retention of MITEs could simply be due to the difficulty of eliminating them without affecting neighboring genes. However, the fixation of MITEs close to genes could also be the result of a positive selection of some MITE insertions. Our data confirms that MITEs are closely associated with rice genes, concentrating in gene proximal upstream and downstream regions. Many examples of MITE insertions in 5' and 3' of genes that alter gene expression in different plants have accumulated over the years, including insertions in promoters enhancing (Zheng *et al.*, 2019; Shimada *et al.*, 2018; Yin *et al.*, 2020; Xu *et al.*, 2020) or repressing (Mao *et al.*, 2015; Xu *et al.*, 2020) transcription, or repressing translation (Shen *et al.*, 2017). MITEs have been shown to frequently contain transcription factor binding sites in plants, including rice (Morata *et al.*, 2018), suggesting a potential impact on promoter evolution. On the other hand, the methylation of MITEs located upstream or downstream of genes can also repress or activate gene expression, as it has recently been shown for several genes linked with rice tillering (Xu *et al.*, 2020).

It is therefore possible, that some of the MITE insertions described here may actually modify gene expression, and that some of the insertions present at high frequency in rice may have been positively selected. In particular a number of MITE insertions appear to be present at HF in rice but are absent from *Oryza rufipogon*, the wild species from which rice was domesticated. Therefore, these insertions may have been selected concomitantly with rice domestication. An important fraction of these insertions (66%) are tightly associated with genes and may have altered their coding capacity or their expression. The fact that some of these genes have already been

characterized as responsible for important functions linked to rice development and stress responses, suggests that MITEs have played an important role in generating variability used in rice domestication. Future work will be needed to determine the extent of this impact.

Whereas an important fraction of MITEs are fixed in rice, a quarter of the MITE TIPs are present at a population frequency lower than 7% and are probably recent insertions. Up to 61 % of these recent insertions are tightly linked to genes and may therefore be involved in differences of gene expression among varieties and be at the origin of trait variability. In addition, MITEs, as TEs in general, can generate an important number of mutations in relatively short timeframe, which could make their insertions a good complement to SNPs for GWAS. LTR-RT TIPs have been recently used for GWAS and it was shown that they can uncover additional associations as compared with SNPs (Akakpo *et al.*, 2020; Domínguez *et al.*, 2020). Here we show that MITE TIPs can also reveal additional associations in GWAS. Indeed, our data shows that most MITE TIPs show a low LD with the surrounding SNPs, similarly to what has been previously found for TE variants in *Arabidopsis* (Stuart *et al.*, 2016). MITEs are short elements that should be better tolerated than LTR-RTs when inserting close to genes. This may make them a particularly suited source of variability for gene evolution but may also make them particularly suited for markers of gene variability, and therefore, particularly useful for GWAS.

The mechanism of MITE amplification

MITEs are a particular type of TEs as, although they are deletion derivatives of TIR-TEs which can mobilized them by a conservative cut-and-paste mechanism, they are frequently present in genomes at very high copy numbers. Class II TEs can increase in copy number and are even able to invade the genome of a species in a short period of time, as it has been clearly shown for the *Drosophila P* element (Clark and Kidwell, 1997). For the *P* element, as well as for other TEs, a link of transposition with DNA replication could explain the increase in copy number through an otherwise conservative cut-and-paste mechanism. Targeting insertions to replication origins may allow coupling P transposition and DNA replication (Spradling *et al.*, 2011). In plants, the maize Ac transposase is known to preferentially bind to freshly replicated hemimethylated DNA (Ros and Kunze, 2001), increasing the frequency of transposition associated with DNA replication.

Although these mechanisms may allow for an increase in copy number of class II TEs, MITEs are usually present at much higher copy numbers, and they resemble in this aspect replicative TEs such as retrotransposons. It has been suggested that MITE high copy number may be the result of an increased efficiency of transposition due to their promiscuity in using transposases (Feschotte *et al.*, 2003), to particular characteristics of the related transposase in some genomes (Guermónprez *et al.*, 2008) or to a higher mobilization efficiency with respect to the related autonomous transposons (Yang *et al.*, 2009). Recent genomic studies on the *mPing* MITE have linked its amplification with the presence of a particular copy of the related *Ping* TIR-TE, as this copy is present in varieties showing bursts of *mPing* amplification (Chen *et al.*, 2019). On the other hand, an analysis of RILs showing variation in *mPing* copy number, identified QTLs containing multiple copies of *Ping* TIR-TE suggesting a link between *Ping* copy number and *mPing* amplification (Chen *et al.*, 2020).

Here we have used SNP- and TIP-GWAS to look for genetic determinants of the amplification of 400 MITE families in rice and failed to uncover a link between the presence of particular TIR-TEs and MITE amplification, although, the mobile nature of TIR-TEs may make it difficult to reveal this association. Interestingly our analyses revealed a positive correlation between the presence in the genome of particular MITE copies and the copy number of the correspondent family, as if only few MITE copies were capable of amplifying. This is what commonly happens with some replicative TEs, such as LINEs, where the replication of one or few "master elements" is at the origin of a whole family. In fact, the phylogenetic analyses of MITE populations in different plant genomes accumulated in the last 20 years (Santiago *et al.*, 2002; Naito *et al.*, 2009; Xin *et al.*, 2019; Lu *et al.*, 2012; Feschotte *et al.*, 2003) are compatible with the replicative amplification of MITEs, and a replicative transposition mechanism independent of related transposases was proposed long-time ago (Izsvák *et al.*, 1999). Izsvák and co-workers proposed that the high secondary structure of MITEs could allow them to fold back to form a stem-loop ssDNA molecule that might detach from the chromosome during replication, in a mechanism similar to that of certain bacterial transposons that move through single-strand transposition associated with DNA replication (Lavatine *et al.*, 2016). The fact that few MITE copies are at the origin of a MITE family could suggest that the amplification is an efficient but rare event that leads to a substantial increase in copy number at each amplification event. Alternatively, some MITE copies could be more prone to amplification due to their location in the genome or their particular sequence or structural characteristics. The analysis reported here did not reveal major differences between

MITE copies positively correlated with an increase of copy number of a MITE family and the rest of the elements of the same family. However, we detected slight enrichment in chromatin marks associated with active transcription at these elements as well as some enrichment in recombination hotspots compared with the rest of the elements of the family. Origins of replication are closely associated with genes and are enriched in active transcription epigenetic marks (Costas *et al.*, 2011), and there are clear links between recombination and replication (Syeda *et al.*, 2014). More work will be needed to clarify to what extent these differences may allow some MITE copies to be amplified, and to what extent the amplification of MITEs is linked to DNA replication.

Conclusions

The results presented here show that an important fraction of MITE insertions is present at high frequency in rice while being absent from its wild ancestor, suggesting that they have been fixed concomitantly with domestication. On the other hand, another fraction of MITE insertions is present at low frequencies among rice varieties, which shows that MITEs have also transposed after domestication. MITEs concentrate close to genes and have generated gene variability during rice domestication and breeding. We used MITE TIPs as genetic information for GWAS and we show that they uncover more associations with rice traits than SNPs. We also used this approach, together with SNP-GWAS to shed light on the still unknown mechanism of MITE amplification. Our results suggest that these elements amplify by a replicative mechanism.

Methods

Data source

We used the fastq files from the 1,059 rice varieties sequenced at a minimum coverage of 15X from the 3000 rice genomes project (Li *et al.*, 2014), and three Platinum Indica assemblies (LIMA::IRGC 81487-1, KHAO YAI GUANG::IRGC 65972-1 and LARHA MUGAD::IRGC 52339-1) (Zhou *et al.*, 2020).

Reconstruction of TE family consensus

We used TEdenovo from the REPET package (Flutre *et al.*, 2011) to build TE consensus sequences from the Japonica Nipponbare genome (Sasaki and Project, 2005), Indica MH63

(Zhang *et al.*, 2016) and Aus N22 (Stein *et al.*, 2018). TE consensus were classified at the order level using PASTEC (Hoede *et al.*, 2014), and only those classified as “complete“ elements were retained. We concatenated the three datasets and removed redundancy by obtaining a library of centroids at 80% identity using VSEARCH (Rognes *et al.*, 2016). We run MITE-Hunter (Han and Wessler, 2010) in the three rice genome assemblies and combined the predictions with the MITE dataset from PMITE (Chen *et al.*, 2014) and followed the same clustering approach to obtain the MITE library. We concatenated the two libraries (generic TEs and MITEs) to obtain the complete rice TE consensus (Supporting Dataset S1). Repeatmasker (<http://www.repeatmasker.org/>) was used to identify genomic regions with similarity to TE consensus.

Detection of TE insertions from sequencing reads

We used PopoolationTE2 (Kofler *et al.*, 2016) with the mode “separate” to detect TE insertions using the genome of Nipponbare as reference, and we discarded predictions below zygosity of 0.7 as previously recommended (Vendrell-Mir *et al.*, 2019). In order to avoid the bias caused by differential sequencing depth, we randomly sampled every accession to 15X prior running the tool using seqtk (<https://github.com/lh3/seqtk>). To identify non-reference insertions, we intersected all detected insertions with the regions annotated by RepeatMasker and selected all the insertion points that were further than 25bp to any annotated TE in the rice Nipponbare reference genome used here as a reference.

Estimation of TIP insertion frequencies

Nipponbare genome was split into 500bp (for MITE analyses) and 1Kb (for analyses of the rest of TE) windows. The results of PopoolationTE2 were transformed into bed format. Predicted insertions from different TE orders were separated and intersected with the genome windows to obtain TIPs. TIPs from the 1,059 varieties were concatenated in one file and windows were collapsed to remove redundancy. TIPs of the same order from adjacent windows were merged to obtain a set of high-confidence TIPs. In a second iteration, we intersected all the insertions of the 1,059 varieties, this time excluding the zygosity filter, with the positions of the high-confidence TIP dataset to score the presence/absence status of each TIP in each variety and obtain the final TIP matrices. TIP insertion frequencies were calculated by obtaining the proportion of varieties that contain each insertion. For MITE and LTR-RTs, we considered as Low-frequency (LF) the 25% TIPs with lower frequency (frequency cut-off: MITE = 6.6%, *gypsy* LTR-RT = 1.3%, *copia*

LTR-RT = 0.2 %), and High-frequency (HF) the 25% TIPs with higher frequency (frequency cut-off: MITE = 85.5%, *gypsy* LTR-RT = 30.8%, *copia* LTR-RT = 25.3% %).

Population structure analyses

Principal component analysis was performed using `prcomp` function from the *stats* R package (The R Foundation for Statistical Computing, 2011). TIP-based Neighbor-joining tree was built using *ape* library from the R package (Paradis *et al.*, 2004). We used the R function `dist.gene` (from the *ape* package) to estimate distances from the genotype matrix

SNP and TIP-GWAS

We used the LFMM 2 R package (Caye *et al.*, 2019) to obtain genotype-phenotype associations applying latent factor mixed models in Indica varieties to correct for population structure ($K = 4$). For the SNP-GWAS we used a SNP matrix obtained from the 173 K SNP matrix obtained by Carpentier *et al.* (Carpentier *et al.*, 2019a) in `lfmm` format (153,744 after $MAF > 1\%$). For TIP-GWAS, we used TIP matrices as genotype. For both TIP GWAS and SNP GWAS we used a MAF threshold of 1%. Independent Bonferroni corrections were used to set significance p-value thresholds for SNP-GWAS and TIP-GWAS. Phenotypic traits were downloaded from SNP-seek database (<https://snp-seek.irri.org/>).

Patterns of linkage disequilibrium of MITE TIPs and SNPs.

The software `ngsLD` (Fox *et al.*, 2019) was used to estimate the pairwise linkage disequilibrium (r^2) between MITE TIPs and SNPs located close by (at a distance of 100 Kb upstream and downstream) and between SNPs and other SNPs located close by. For each MITE TIP the mean TIP-SNP r^2 was compared to the mean SNP-SNP r^2 of all SNPs in the corresponding 200Kb window. The number of comparisons where the TIP-SNP average r^2 was higher than SNP-SNP was calculated. We considered a r^2 of 0.2 as a cut-off to define TIPs in low LD with SNPs, as previously done by (Domínguez *et al.*, 2020).

Identification of genetic factors

The results of SNP-GWAS using TE family copy number as phenotype and SNPs as genotype were manually inspected to identify peaks with significant associations (Bonferroni-adjusted p

value threshold = $3.2e^{-07}$). We looked for TIPs of the same family located at < 100 Kb, a cutoff based on the patterns of linkage disequilibrium in Indica rice (Mather *et al.*, 2007). Statistical significance between family copy number between varieties carrying or not the genetic factor was assessed using two-tailed Wilcoxon rank sum test (p-value cutoff = 0.05).

Validation of Genetic Factors

Windows containing the genetic factors were extended 4000bp upstream and downstream and extracted from the reference genome (IRGSP, Nipponbare). The regions were mapped to three Platinum Indica assemblies present in our dataset (LIMA::IRGC 81487-1, KHAO YAI GUANG::IRGC 65972-1 and LARHA MUGAD::IRGC 52339-1) using minimap2 (Li, 2018). By screening the best mapping hits we identified the orthologous regions in the three assemblies. Using the family consensus as query for a BLASTN search (Altschul *et al.*, 1990) (cutoff e^{-20}), we verified the presence of the corresponding TIP in each orthologous region. Schematic representation of all genomic BLAST hits obtained using the consensus sequence as query was performed using Sushi R library (Phanstiel *et al.*, 2014).

Sequence analysis of Genetic factors

The sequence of every MITE in LIMA::IRGC 81487-1 genome was extracted including 100, 250 and 500 bp of upstream and downstream sequences. GC content was determined with bbmap (Bushnell, 2014) and the minimum free energy with ViennaRNA package (Lorenz *et al.*, 2011).

Data availability

The datasets generated during and/or analyzed during the current study are available in the following repositories: TE consensus, raw insertions detected by PopoolationTE2, TIP matrices and all the necessary input files and code to reproduce the analyses (TIP detection, GWAS and LD analyses) are available in Zenodo (10.5281/zenodo.4058696); Raw sequencing data of the 1059 rice genomes is available as part of the 3,000 rice genomes project in SRA accession PRJEB6180; H3K4me3, H3K9ac, H3K27me3 and DNase I hypersensitivity sites were obtained from GEO database (accession numbers GSM489075, GSM489079, GSM489083, and GSM655033). Indica rice crossover hotspots were obtained from https://github.com/plantformatics/rice_rec_rates.

Acknowledgements

RC is recipient of a Juan de la Cierva-formation contract from the Spanish Ministerio de Ciencia y Innovación. This work was supported by grants from the Spanish Ministerio de Economía, Industria y Competitividad (AGL2016-78992-R/FEDER) and Ministerio de Ciencia y Innovación (PID2019-106374RB-I00 /AEI / 10.13039/501100011033) to JMC. Computing resources have been provided by the Red Española de Supercomputación at the Pirineus machine (RES activity BCV-2019-2-0006). We thank Marie Mirouze (Université de Perpignan), and Miguel Perez-Enciso (CRAG, Barcelona) for helpful discussions. The authors have no conflict of interest.

Authors' contributions

JMC and RC designed the analysis. RC, PV-M and AB performed TIP detection and analyzed the data. RC performed GWAS analyses with the collaboration M-C C and OP. JMC and RC wrote the manuscript with collaborations of PV-M and AB. All authors revised and approved the manuscript.

Conflicts of interest

The authors declare that they have no conflict of interest

Supporting material legends

Supporting Figure S1. Percentage of the TE library occupied by each TE order. DHX = Helitrons, DTX = TIR-TE transposons, RIX = LINEs, RLX = LTR-retrotransposons, RSX = SINEs.

Supporting Figure S2. Insertion detection bias between CAAS and IRIS varieties. Number of PopulationTE2 filtered insertions per variety in 1059 coverage-homogenized genomes.

Supporting Figure S3. Genome wide density plot of MITE and LTR-retrotransposon TIPs.

Supporting Figure S4. Genome wide density plot genes, MITEs and LTR-retrotransposons according to their High (Q1) or Low (Q4) frequency on the population.

Supporting Figure S5. SNP-GWAS of seven agronomic traits. Bonferroni cut-off = $3.2e-07$. Positions annotated over leading significant SNPs indicate the absolute position starting from the first nucleotide of chromosome 1.

Supporting Figure S6. SNP-GWAS of grain width phenotype using a genotype matrix of 35,340 SNPs. MAF > 1%. Bonferroni cut-off = 1.4e-06

Supporting Figure S7. Grain width of varieties carrying MITE-1 insertion (chr05:5364000-5365000) and the deletion located 4,500 bp upstream the GSE5 gene described in (Duan et al. 2017).

Supporting Dataset S1. Consensus sequences of the 821 TE families.

Supporting Table S1. Description of the 1,059 varieties used in this study.

Supporting Table S2. List of High-frequency MITE insertions present in *O. sativa* ssp. *Japonica* cv. Nipponbare and absent in *O. rufipogon*, present in genic regions.

Supporting Table S3. Table S3. Significant associations between MITE insertions and seven agronomic traits. In green, MITEs in low LD with surrounding SNPs ($r^2 < 0.2$). The presence of a surrounding significant SNP associated with the trait is shown in column 7.

Supporting Table S4. List of 18 MITE families and their linked transposases showing significant similarity on the 5' and 3' extremities, along with the results from the SNP-GWAS and TIP-GWAS using TIR-TEs and MITEs.

Supporting Table S5. GWAS leading SNPs associated with family copy number for 428 TE families.

Supporting Table S6. Information related to the GWAS leading SNPs and TIPs associated with LINE family copy number.

Supporting Table S7. GC content, Minimum Free Energy Structure and TE content (250 Kb region) of GF and non-GF loci of 25 MITE families present in LIMA::IRGC 81487-1 Genome.

References

- Akakpo, R., Carpentier, M.C., Ie Hsing, Y. and Panaud, O.** (2020) The impact of transposable elements on the structure, evolution and function of the rice genome. *New Phytol.*, **226**, 44–49.
- Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J.** (1990) Basic local alignment search tool. *J. Mol. Biol.*, **215**, 403–410.
- Benjak, A., Boué, S., Forneck, A. and Casacuberta, J.M.** (2009) Recent amplification and impact of MITEs on the genome of grapevine (*Vitis vinifera* L.). *Genome Biol. Evol.*, **1**, 75–84. Available at: <http://gbe.oxfordjournals.org/content/1/75.abstract>.
- Bushnell, B.** (2014) *BBMap: a fast, accurate, splice-aware aligner*, Available at: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3572422&tool=pmcentrez&rendertype=abstract>.
- Carpentier, M.C., Manfroi, E., Wei, F.J., et al.** (2019a) Retrotranspositional landscape of Asian rice revealed by 3000 genomes. *Nat. Commun.*, **10**, 24. Available at: <https://doi.org/10.1038/s41467-018-07974-5>.
- Carpentier, M.C., Manfroi, E., Wei, F.J., et al.** (2019b) Retrotranspositional landscape of Asian rice revealed by 3000 genomes. *Nat. Commun.*, **10**, 24.
- Casacuberta, J.M. and Santiago, N.** (2003) Plant LTR-retrotransposons and MITEs: Control of transposition and impact on the evolution of plant genes and genomes. *Gene*, **311**, 1–11.
- Caye, K., Jumentier, B., Lepeule, J. and François, O.** (2019) LFMM 2: Fast and accurate inference of gene-environment associations in genome-wide studies. *Mol. Biol. Evol.*, **36**, 852–860.
- Chen, J., Hu, Q., Zhang, Y., Lu, C. and Kuang, H.** (2014) P-MITE: A database for plant miniature inverted-repeat transposable elements. *Nucleic Acids Res.*, **42**, D1176–D1181.
- Chen, J., Lu, L., Benjamin, J., Diaz, S., Hancock, C.N., Stajich, J.E. and Wessler, S.R.** (2019) Tracking the origin of two genetic components associated with transposable element bursts in domesticated rice. *Nat. Commun.*, **10**, 641.
- Chen, J., Lu, L., Robb, S.M.C., Collin, M., Okumoto, Y. and Stajich, J.E.** (2020) Genomic diversity generated by a transposable element burst in a rice recombinant inbred population. *Proc. Natl. Acad. Sci.*, **117**, 26288–26297.
- Clark, J.B. and Kidwell, M.G.** (1997) A phylogenetic perspective on P transposable element

evolution in *Drosophila*. *Proc. Natl. Acad. Sci. U. S. A.*, **94**, 11428–11433.

Costas, C., La Paz Sanchez, M. De, Stroud, H., et al. (2011) Genome-wide mapping of *Arabidopsis thaliana* origins of DNA replication and their associated epigenetic marks. *Nat. Struct. Mol. Biol.*, **18**, 395–400. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/21297636>.

Domínguez, M., Dugas, E., Benchouaia, M., Leduque, B., Jiménez-Gómez, J.M., Colot, V. and Quadrana, L. (2020) The impact of transposable elements on tomato diversity. *Nat. Commun.*

Dong, Y., Tsuzuki, E., Kamiunten, H., Terao, H. and Lin, D. (2003) Mapping of QTL for embryo size in rice. *Crop Sci.*, **43**, 1068–1071.

Duan, P., Xu, J., Zeng, D., et al. (2017) Natural Variation in the Promoter of GSE5 Contributes to Grain Size Diversity in Rice. *Mol. Plant*, **10**, 685–694.

Farley, A.H., Luning Park, E.T. and Kazazian, H.H. (2004) More active human L1 retrotransposons produce longer insertions. *Nucleic Acids Res.*, **32**, 502–510.

Feschotte, C., Jiang, N. and Wessler, S.R. (2002) Plant transposable elements: Where genetics meets genomics. *Nat. Rev. Genet.*, **3**, 329–341. Available at: <http://dx.doi.org/10.1038/nrg793>.

Feschotte, C., Swamy, L. and Wessler, S.R. (2003) Genome-wide analysis of mariner-like transposable elements in rice reveals complex relationships with Stowaway miniature inverted repeat transposable elements (MITEs). *Genetics*, **163**, 747–758.

Flutre, T., Duprat, E., Feuillet, C. and Quesneville, H. (2011) Considering transposable element diversification in de novo annotation approaches. *PLoS One*, **6**.

Fox, E.A., Wright, A.E., Fumagalli, M. and Vieira, F.G. (2019) NgsLD: Evaluating linkage disequilibrium using genotype likelihoods. *Bioinformatics*, **35**, 3855–3856.

Guermónprez, H., Loot, C. and Casacuberta, J.M. (2008) Different strategies to persist: The pogo-like Lem1 transposon produces miniature inverted-repeat transposable elements or typical defective elements in different plant genomes. *Genetics*, **180**, 83–92.

Gutaker, R.M., Groen, S.C., Bellis, E.S., et al. (2020) Genomic history and ecology of the geographic spread of rice. *Nat. Plants*, **6**, 492–502.

Han, Y. and Wessler, S.R. (2010) MITE-Hunter: A program for discovering miniature inverted-repeat transposable elements from genomic sequences. *Nucleic Acids Res.*, **38**.

Hoede, C., Arnoux, S., Moisset, M., Chaumier, T., Inizan, O., Jamilloux, V. and Quesneville,

- H. (2014) PASTEC: An automatic transposable element classification tool. *PLoS One*, **9**.
- Izsvák, Z., Ivics, Z., Shimoda, N., Mohn, D., Okamoto, H. and Hackett, P.B. (1999) Short inverted-repeat transposable elements in teleost fish and implications for a mechanism of their amplification. *J. Mol. Evol.*, **48**, 13–21.
- Jackson, M.T. (1997) Conservation of rice genetic resources: The role of the International Rice Genebank at IRRI. *Plant Mol. Biol.*
- Kawahara, Y., la Bastide, M. de, Hamilton, J.P., et al. (2013) Improvement of the oryza sativa nipponbare reference genome using next generation sequence and optical map data. *Rice*, **6**, 3–10.
- Kofler, R., Gómez-Sánchez, D. and Schlötterer, C. (2016) PoPoolationTE2: Comparative Population Genomics of Transposable Elements Using Pool-Seq. *Mol. Biol. Evol.*, **33**, 2759–2764.
- Kou, Y., Liao, Y., Toivainen, T., Lv, Y., Tian, X., Emerson, J.J., Gaut, B. and Zhou, Y. (2020) Evolutionary genomics of structural variation in Asian rice (*Oryza sativa*) domestication. *Mol. Biol. Evol.* Available at: <https://doi.org/10.1093/molbev/msaa185>.
- Lavatine, L., He, S., Caumont-Sarcos, A., Guynet, C., Marty, B., Chandler, M. and Ton-Hoang, B. (2016) Single strand transposition at the host replication fork. *Nucleic Acids Res.*, **44**, 7866–7883.
- Li, H. (2018) Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics*, **34**, 3094–3100.
- Li, J.Y., Wang, J. and Zeigler, R.S. (2014) The 3,000 rice genomes project: New opportunities and challenges for future rice research. *Gigascience*, **3**, 8. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4035671/>.
- Li, Z., Pinson, S.R.M., Park, W.D., Paterson, A.H. and Stansel, J.W. (1997) Epistasis for three grain yield components in rice (*Oryza sativa* L.). *Genetics*, **145**, 453–465.
- Lisch, D. (2013) How important are transposons for plant evolution? *Nat. Rev. Genet.*, **14**, 49–61. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/23247435>.
- Lorenz, R., Bernhart, S.H., Höner zu Siederdisen, C., Tafer, H., Flamm, C., Stadler, P.F. and Hofacker, I.L. (2011) ViennaRNA Package 2.0. *Algorithms Mol. Biol.*, **6**.
- Lu, C., Chen, J., Zhang, Y., Hu, Q., Su, W. and Kuang, H. (2012) Miniature inverted-repeat transposable elements (MITEs) have been accumulated through amplification bursts and play important roles in gene expression and species diversity in *oryza sativa*. *Mol. Biol. Evol.*, **29**,

1005–1017.

- Ma, Xuefeng, Cheng, Z., Qin, R., et al.** (2013) OsARG encodes an arginase that plays critical roles in panicle development and grain production in rice. *Plant J.*
- Mansueto, L., Fuentes, R.R., Borja, F.N., et al.** (2017) Rice SNP-seek database update: New SNPs, indels, and queries. *Nucleic Acids Res.*, **45**, D1075–D1081.
- Mao, B.B., Cai, W.J., Zhang, Z.H., Hu, Z.L., Li, P., Zhu, L.H. and Zhu, Y.G.** (2003) Characterization of QTLs for harvest index and source-sink characters in a DH population of rice (*Oryza sativa* L.). *Yi Chuan Xue Bao*, **30**, 1118–1126.
- Mao, H., Wang, H., Liu, S., Li, Z., Yang, X., Yan, J., Li, J., Tran, L.S.P. and Qin, F.** (2015) A transposable element in a NAC gene is associated with drought tolerance in maize seedlings. *Nat. Commun.*, **6**, 8326.
- Marand, A.P., Zhao, H., Zhang, W., Zeng, Z., Fang, C. and Jianga, J.** (2019) Historical meiotic crossover hotspots fueled patterns of evolutionary divergence in rice. *Plant Cell*, **31**, 645–662. Available at: <http://www.plantcell.org/content/31/3/645.abstract>.
- Martin, S.L., Li, W.L.P., Furano, A. V. and Boissinot, S.** (2005) The structures of mouse and human L1 elements reflect their insertion mechanism. *Cytogenet. Genome Res.*, **110**, 223–228.
- Mather, K.A., Caicedo, A.L., Polato, N.R., Olsen, K.M., McCouch, S. and Purugganan, M.D.** (2007) The extent of linkage disequilibrium in rice (*Oryza sativa* L.). *Genetics*, **177**, 2223–2232.
- Morata, J., Marín, F., Payet, J. and Casacuberta, J.M.** (2018) Plant lineage-specific amplification of transcription factor binding motifs by miniature inverted-repeat transposable elements (MITEs). *Genome Biol. Evol.*, **10**, 1210–1220.
- Naito, K., Zhang, F., Tsukiyama, T., Saito, H., Hancock, C.N., Richardson, A.O., Okumoto, Y., Tanisaka, T. and Wessler, S.R.** (2009) Unexpected consequences of a sudden and massive transposon amplification on rice gene expression. *Nature*, **461**, 1130–1134. Available at: <http://dx.doi.org/10.1038/nature08479>.
- Ou, S., Su, W., Liao, Y., et al.** (2019) Benchmarking transposable element annotation methods for creation of a streamlined, comprehensive pipeline. *Genome Biol.*, **20**, 275.
- Paradis, E., Claude, J. and Strimmer, K.** (2004) APE: Analyses of phylogenetics and evolution in R language. *Bioinformatics*, **20**, 289–290.
- Phanstiel, D.H., Boyle, A.P., Araya, C.L. and Snyder, M.P.** (2014) Sushi.R: Flexible,

quantitative and integrative genomic visualizations for publication-quality multi-panel figures. *Bioinformatics*, **30**, 2808–2810.

Rognes, T., Flouri, T., Nichols, B., Quince, C. and Mahé, F. (2016) VSEARCH: A versatile open source tool for metagenomics. *PeerJ*, **2016**.

Ros, F. and Kunze, R. (2001) Regulation of Activator/Dissociation transposition by replication and DNA methylation. *Genetics*.

Santiago, N., Herráiz, C., Ramón Goñi, J., Messeguer, X. and Casacuberta, J.M. (2002) Genome-wide analysis of the Emigrant family of MITEs of *Arabidopsis thaliana*. *Mol. Biol. Evol.*, **19**, 2285–2293.

Sasaki, T. and Project, I.R.G.S. (2005) The map-based sequence of the rice genome. *Nature*, **436**, 793–800. Available at: <https://doi.org/10.1038/nature03895>.

Shen, J., Liu, J., Xie, K., Xing, F., Xiong, F., Xiao, J., Li, X. and Xiong, L. (2017) Translational repression by a miniature inverted-repeat transposable element in the 3' untranslated region. *Nat. Commun.*, **8**, 14651.

Shimada, T., Endo, T., Fujii, H., Nakano, M., Sugiyama, A., Daido, G., Ohta, S., Yoshioka, T. and Omura, M. (2018) MITE insertion-dependent expression of CitRKD1 with a RWP-RK domain regulates somatic embryogenesis in citrus nucellar tissues. *BMC Plant Biol.*, **18**, 166.

Spradling, A.C., Bellen, H.J. and Hoskins, R.A. (2011) *Drosophila* P elements preferentially transpose to replication origins. *Proc. Natl. Acad. Sci. U. S. A.*, **108**, 15948–15953.

Stein, J.C., Yu, Y., Copetti, D., et al. (2018) Genomes of 13 domesticated and wild rice relatives highlight genetic conservation, turnover and innovation across the genus *Oryza*. *Nat. Genet.*, **50**, 285–296.

Stuart, T., Eichten, S.R., Cahn, J., Karpievitch, Y. V., Borevitz, J.O. and Lister, R. (2016) Population scale mapping of transposable element diversity reveals links to gene regulation and epigenomic variation. *Elife*, **5**, e20777.

Swergold, G.D. (1990) Identification, characterization, and cell specificity of a human LINE-1 promoter. *Mol. Cell. Biol.*, **10**, 6718–6729.

Syeda, A.H., Hawkins, M. and McGlynn, P. (2014) Recombination and Replication. *Cold Spring Harb. Perspect. Biol.*, **6**, 1–14.

Tenaillon, M.I., Hollister, J.D. and Gaut, B.S. (2010) A triptych of the evolution of plant transposable elements. *Trends Plant Sci.*, **15**, 471–478.

The R Foundation for Statistical Computing (2011) R Development Core Team. *R A Lang. Environ. Stat. Comput.* Available at: <http://www.r-project.org/>.

Vendrell-Mir, P., Barteri, F., Merenciano, M., González, J., Casacuberta, J.M. and Castanera, R. (2019) A benchmark of transposon insertion detection tools using real data. *Mob. DNA*, **10**, 53. Available at: <https://doi.org/10.1186/s13100-019-0197-9>.

Wang, W., Mauleon, R., Hu, Z., et al. (2018) Genomic variation in 3,010 diverse accessions of Asian cultivated rice. *Nature*, **557**, 43–49.

Xiao, J., Li, J., Yuan, L. and Tanksley, S.D. (1996) Identification of QTLs affecting traits of agronomic importance in a recombinant inbred population derived from a subspecific rice cross. *Theor. Appl. Genet.*, **92**, 230–244.

Xin, Y., Ma, B., Xiang, Z. and He, N. (2019) Amplification of miniature inverted-repeat transposable elements and the associated impact on gene regulation and alternative splicing in mulberry (*Morus notabilis*). *Mob. DNA*, **10**, 27.

Xu, L., Yuan, K., Yuan, M., Meng, X., Chen, M., Wu, J., Li, J. and Qi, Y. (2020) Regulation of Rice Tillering by RNA-Directed DNA Methylation at Miniature Inverted-Repeat Transposable Elements. *Mol. Plant*, **13**, 851–863.

Yan, C.J., Liang, G.H., Chen, F., Li, X., Tang, S.Z., Yi, C.D., Tian, S., Lu, J.F. and Gu, M.H. (2003) Mapping quantitative trait loci associated with rice grain shape based on an indica/japonica backcross population. *Acta Genet. Sin.*, **30**, 711–716.

Yang, G., Lee, Y.H., Jiang, Y., Shi, X., Kertbundit, S. and Hall, T.C. (2005) A two-edged role for the transposable element Kiddo in the rice ubiquitin2 promoter. *Plant Cell*.

Yang, G., Nagel, D.H., Feschotte, C., Nathan Hancock, C. and Wessler, S.R. (2009) Tuned for transposition: Molecular determinants underlying the hyperactivity of a stowaway MITE. *Science (80-.)*, **325**, 1391–1394.

Yin, S., Wan, M., Guo, C., et al. (2020) Transposon insertions within alleles of BnaFLC.A10 and BnaFLC.A2 are associated with seasonal crop type in rapeseed. *J. Exp. Bot.*, **71**, 4729–4741.

Yonemaru, J. ichi, Yamamoto, T., Fukuoka, S., Uga, Y., Hori, K. and Yano, M. (2010) Q-TARO: QTL annotation rice online database. *Rice*, **3**, 194–203.

Zhang, J., Chen, L.-L., Sun, S., et al. (2016) Building two indica rice reference genomes with PacBio long-read and Illumina paired-end sequencing data. *Sci. Data*, **3**, 160076. Available at: <https://doi.org/10.1038/sdata.2016.76>.

Zheng, X., Zhu, K., Sun, Q., et al. (2019) Natural Variation in CCD4 Promoter Underpins

Species-Specific Evolution of Red Coloration in Citrus Peel. *Mol. Plant*, **12**, 1294–1307.

Zhou, Y., Chebotarov, D., Kudrna, D., et al. (2020) A platinum standard pan-genome resource that represents the population structure of Asian rice. *Sci. Data*, **7**, 113.

Accepted Article

Figure legends

Figure 1. TIP-based population structure and TIP content in 738 rice varieties. A) Neighbor-joining tree based on a presence/absence matrix of 103,109 insertion polymorphisms. Colors denote the nine SNP-based sub-populations described by Wang et.al (2018), comprising 4 Indica (XI) subpopulations (1A, 1B, 2 and 3) as well as accessions with admixture components (adm), 3 Japonica (GJ) subpopulations (trp, sbtrp and trop) as well as accessions with admixture components (adm), the Aus varieties (cA) and the aromatic varieties which include Basmati varieties (cB). B) PCA based on a presence/absence matrix of 103,109 insertion polymorphisms, using the same color code. C) Reference insertions per variety grouped by subpopulation. D) non-reference insertions per variety grouped by subpopulation.

Figure 2. Signatures of subspecies-specific activity in rice TE families.

A) Distribution of inter-subspecies Coefficient of Variation in the average family insertion number. Dotted line represents a CV of 20%. B) Example of a *gypsy* LTR-retrotransposon family (RLX_comp_MH63_B_R175_Map20) showing signs of amplification in Indica and Aus subspecies after their split with Japonica and Aro. C) Mean insertion per family *versus* the total number of family TIPs in the population.

Figure 3. TIP population frequencies and their conservation in wild *Oryza* species. A) TIP frequencies in 382 rice traditional varieties. Number of TIPs per order are: 4,113 Helitrons, 45,837 MITEs, 18,592 TIR-TEs, 5,026 *copia* LTR-RTs, 18,742 *gypsy* LTR-RTs, 2,402 LINEs and 2,000 SINEs. B) Percentage of TIPs present in *O. sativa* ssp. Japonica cv. Nipponbare and in 6 wild *Oryza* species. A dendrogram representing the phylogenetic relationships between the *Oryza* species analyzed is shown. Branch lengths do not represent real phylogenetic distances. Divergence times were obtained from <http://www.timetree.org/>. C) Population frequency of the TIPs conserved between *O. sativa* ssp. Japonica cv. Nipponbare and *O. glumaepatula* in the traditional rice varieties.

Figure 4. Genome-wide distribution of TIPs and their impact on genes.

A) Density plot showing the chromosomal distribution of High-frequency (HF) and Low-Frequency (LF) TIPs of MITEs, *gypsy* and *copia* LTR-RTs and genes (Chr05, full genome shown

in Supporting Figure S4). B) Number of non-reference MITE insertions per variety. B) Number of non-reference HF and LF MITE insertions per 100Kb of genomic feature. Asterisks denote significant differences between each group and “intergenic” group.

Figure 5: Patterns of linkage disequilibrium of MITE TIPs and SNPs.

A) Distribution of MITE TIP-SNPs (blue) and SNP-SNP (in red, full SNP dataset; in green, a subsample of 35,340 SNPs) linkage disequilibrium at ± 100 Kb (r^2). B) Local comparison of linkage disequilibrium (r^2) between TIP-SNP and SNP-SNP. TIP-SNP r^2 was calculated for each TIP and all surrounding SNPs at ± 100 Kb. SNP-SNP r^2 was calculated for all pairs of SNPs located in the same region. Each mean TIP-SNP r^2 was compared to every mean SNP-SNP r^2 at ± 100 Kb (mean of 104 comparisons per TIP). The percentile represents the proportion of comparisons where the mean TIP-SNP r^2 is higher than mean SNP-SNP.

Figure 6. Association of MITE transposon insertions with rice grain width.

A) Rice grain width SNP-GWAS using 451 rice Indica varieties. B) Rice grain width TIP GWAS using MITEs insertion polymorphisms as genotype (34,023 TIPs at MAF > 1%). Red line represents the Bonferroni-adjusted significant threshold (Panel A: $3.2e^{-07}$, panel B = $1.47e^{-06}$). Significant MITE TIPs overlapping known seed QTLs are marked in the manhattan plots. In green, significant MITE TIPs in low LD with surrounding SNPs ($r^2 < 0.2$). C) Violin plots showing the rice grain width in the different subsets of varieties carrying (present) or lacking (absent) the MITE TIPs in low LD, and not having any significant SNP at (\pm)100Kb. Grain width of the varieties carrying each of the two alleles of the leading SNP in this trait (chr05:5361195) are also included in the analysis as another group (T-SNP and A-SNP). Differences between means were tested using Student's t-test (p -value cutoff = 0.05). More information about the significant MITE-TIPs is shown in Supporting Table S3.

Figure 7. Identification of MITE copies as main genetic factors responsible for their amplification in Indica rice.

A) SNP-GWAS analysis of MITE family SE260300235fam318_632 using insertion numbers as phenotype and SNPs as genotype (451 Indica varieties). B) TIP-GWAS analysis using insertion numbers as phenotype and all MITE TIPs as genotype. C) TIR-TE TIP-GWAS analysis using insertion numbers as phenotype and all TIR-TE TIPs as genotype. In green, TIPs from the family

SE260300235fam318_632 (panel B) and from the distantly related TIR-TE DTX_comp_IRGSP_B_R1932 (panel C). Red lines represent the Bonferroni-adjusted significant threshold (Panel A: $3.2e^{-07}$, Panel B: $1.47e^{-06}$, Panel C: $9.3e^{-06}$). D) Mean MITE insertion numbers of varieties containing or not the associated MITEs. Control are all the varieties that do not contain any of the two MITEs. The number of varieties present in each category are indicated upstream each violin. P-values correspond to two-tailed Wilcoxon rank sum tests comparing each category vs the control category. E) Genotypes of the leading SNPs in all Indica varieties and in the subsets of varieties carrying the two associated MITEs. Distances between each MITE and the leading SNP is reported in Kb. F) Schematic representation of the nucleotide conservation between the consensus sequences of the full TIR-TE element DTX_comp_IRGSP_B_R1932 and its related MITE SE260300235fam318_632. ORF = Open Reading Frame.

Figure 8. Identification of full-length copies as main genetic factors responsible for LINE amplification in Indica rice.

A) SNP-GWAS and TIP-GWAS (451 Indica varieties) analysis of LINE family RIX_comp_MH63_B_R3107_MAP6 using insertion numbers as phenotype. SNPs (upper Manhattan plot) or LINE TIPs (lower Manhattan plot) are used as genotype. Red lines represent the Bonferroni-adjusted significant threshold (Panel A: $3.2e^{-07}$, Panel B: $3.03e^{-05}$). In green, LINE TIPs of the target family (RIX_comp_MH63_B_R3107_MAP6). GF = Genetic factor) Mean insertion numbers of varieties containing or lacking the putative genetic factors. The number of varieties present in each category are indicated upstream each violin. P-values correspond to two-tailed Wilcoxon rank sum tests comparing each category vs the no-GF category. C) Genotypes of the leading SNPs in all Indica varieties and in the subsets of varieties carrying the three different genetic factors. Distance between each GF and the leading SNP is reported in Kb. D) Identification of Genetic Factors in the assembled genomes of LIMA::IRGC 81487-1 and IR64. In red, family consensus sequence. In blue, schematic representation of all genomic BLAST hits obtained using the consensus sequence as query.

Table legends

Table 1. Number of TE-families per TE order showing group-specific amplification. The number of families showing a coefficient of variation among varietal groups (CV) of more than 20% are shown with respect to the total number of families for each specific TE order. Only the LTR-RT order is significantly enriched (Fischer's exact test) in families that show group-specific amplification (shown in bold)

Table 2. Families with individual MITE TIPs significantly associated with recent amplification in indica population.

Table 3. Overlap between MITE TIP genetic factors, chromatin marks and indica recombination hotspots.

Tables

Table 1

Order	CV > 20 %	N. families	%	pvalue
Helitrons	1	29	3.45	0.0708
TIR-TEs	34	168	20.24	0.5843
MITEs	62	400	15.50	0.1094
LINEs	19	104	18.27	1,00
LTR-RTs	34	112	30.36	0.0062
SINEs	0	8	0.00	0.6171

Table 2

Family	TIPs/var	Total TIPs	Leading TIP	MAF *	Present**	Absent**	Fold increment
MITE_MH63fam8_344	825	2087	chr07_129000_129500	0,03	18	8	2,2
MITE_MH63fam47_235	462	3603	chr07_26029000_26030500	0,03	42	21	2,1
MITE_N22fam34_480	431	2609	chr02_9131000_9132500	0,07	23	14	1,7
MITE_MH63fam6_341	427	3055	chr05_20497500_20498000	0,03	29	18	1,6
MITE_MH63fam73_259	276	1002	chr07_25096000_25097500	0,04	9	4	2,3
MITE_MH63fam72_365	214	971	chr05_24965500_24966000	0,02	28	10	2,9
MITE_MH63fam106_364	203	1050	chr07_290500_291000	0,04	15	10	1,6
MITE_MH63fam14_237	180	1751	chr03_29812500_29813500	0,09	17	12	1,4
MITE_Oryza1fam20_279	169	2137	chr07_6848500_6850000	0,03	26	15	1,7
MITE_MH63fam50_219	146	1199	chr05_21075000_21075500	0,04	17	8	2,0
MITE_N22fam30_347	142	790	chr04_20639500_20640500	0,04	11	5	2,2
MITE_MH63fam29_244	141	1810	chr12_16399500_16400500	0,08	26	16	1,6
MITE_N22fam5_230	127	1097	chr01_28203000_28204000	0,05	15	8	1,9
MITE_MH63fam13_234	126	1188	chr04_8247500_8250000	0,04	11	6	2,0
MITE_MH63fam32_236	91	2190	chr07_25306000_25306500	0,03	41	23	1,8
MITE_MH63fam51_257	87	699	chr03_14054000_14055000	0,08	8	5	1,6
MITE_SE260500112fam219_340	86	2367	chr09_18380000_18381000	0,04	22	14	1,5
MITE_SE260500111fam211_334	64	1000	chr10_21262000_21262500	0,07	11	7	1,5

*Leading TIP

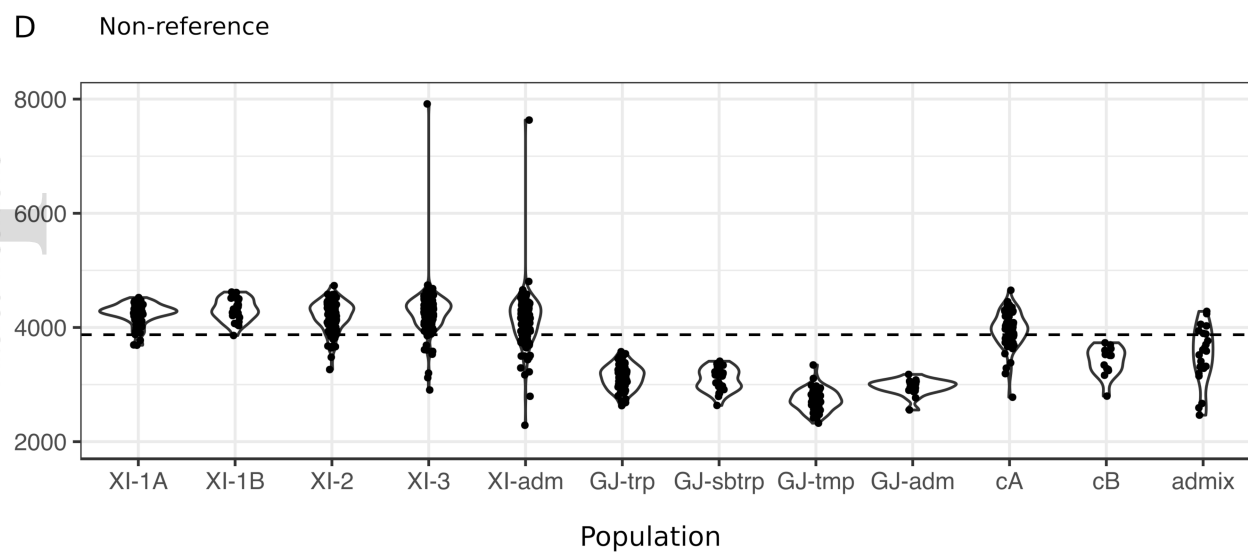
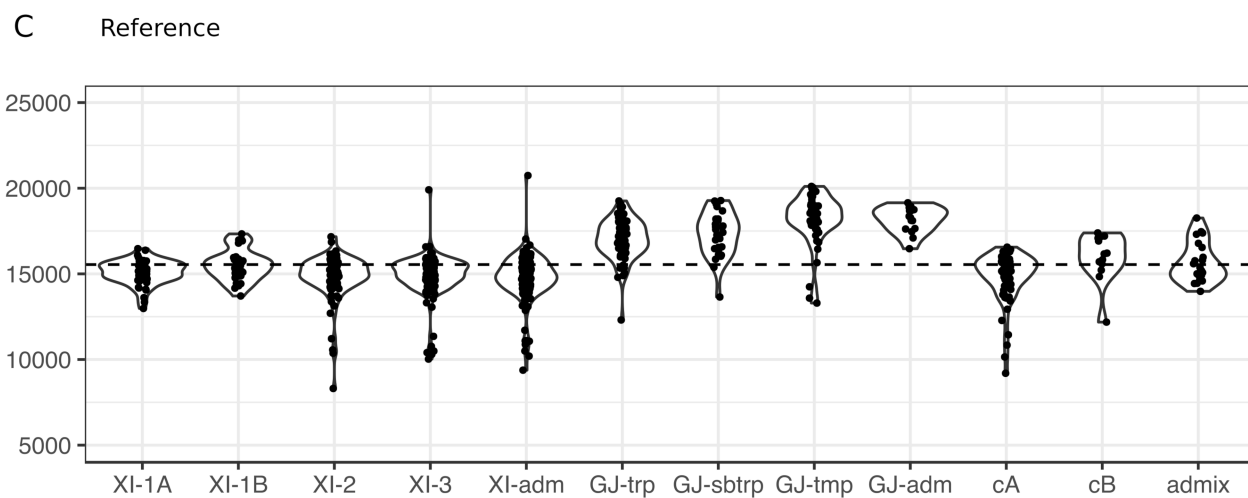
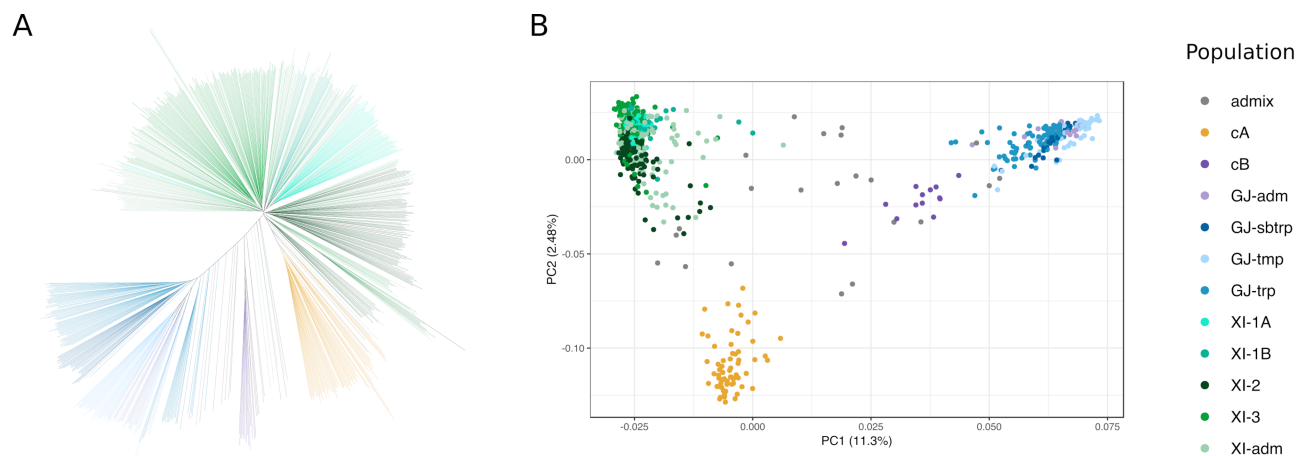
** Mean TIPs/var at a frequency in population lower than 10% in varieties with the leading TIP present or absent

Table 3

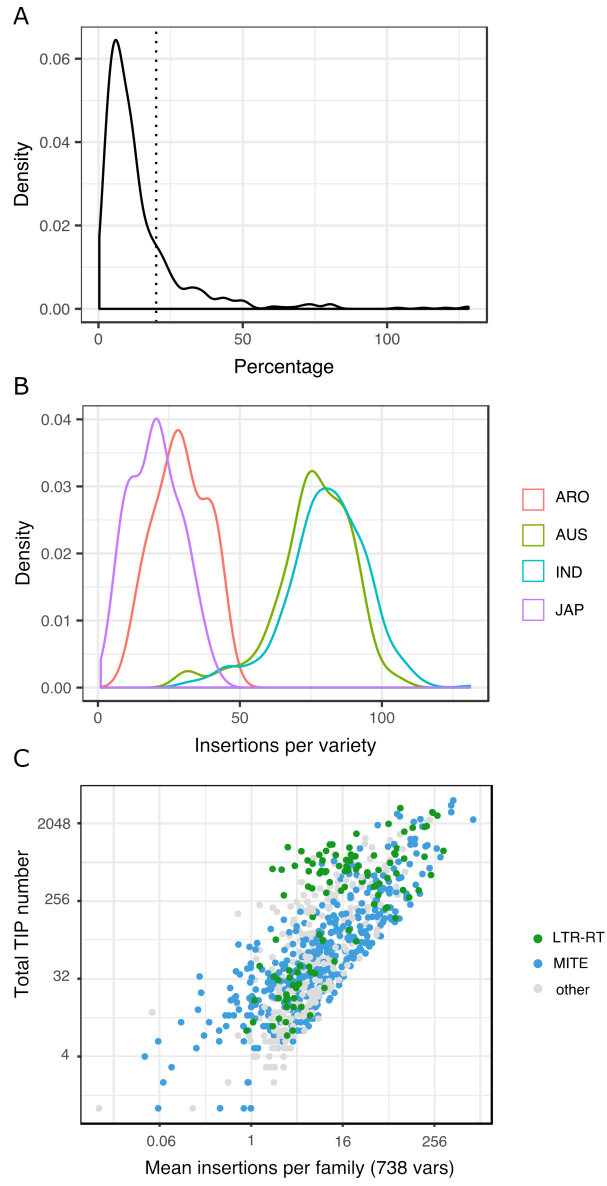
	H3K4me3	H3K9ac	H3K27me3	DH	Rec.hotspot
MITE-GF *	16,3	23,3	11,6	16,3	18,6
MITE-All **	14,7	17,4	10,9	14,3	12,2

*43 MITEs

** 71,796 MITEs longer than 200bp

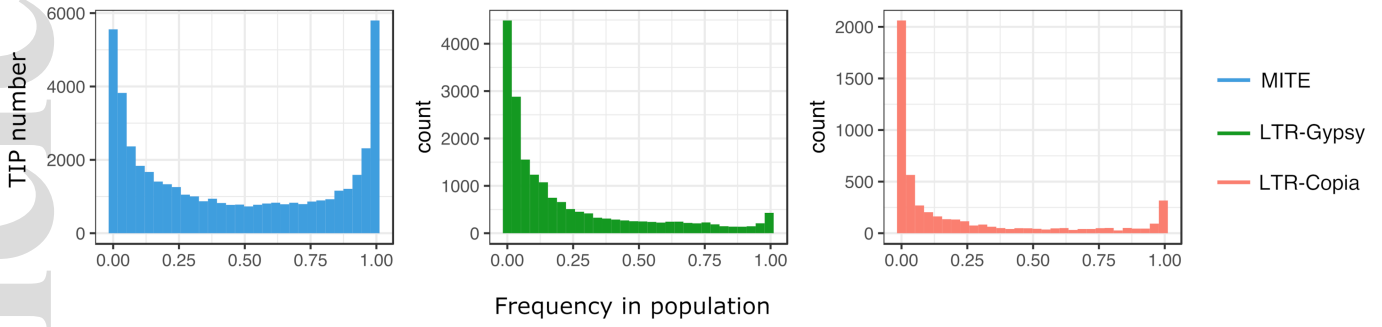


tpj_15277_f1.png

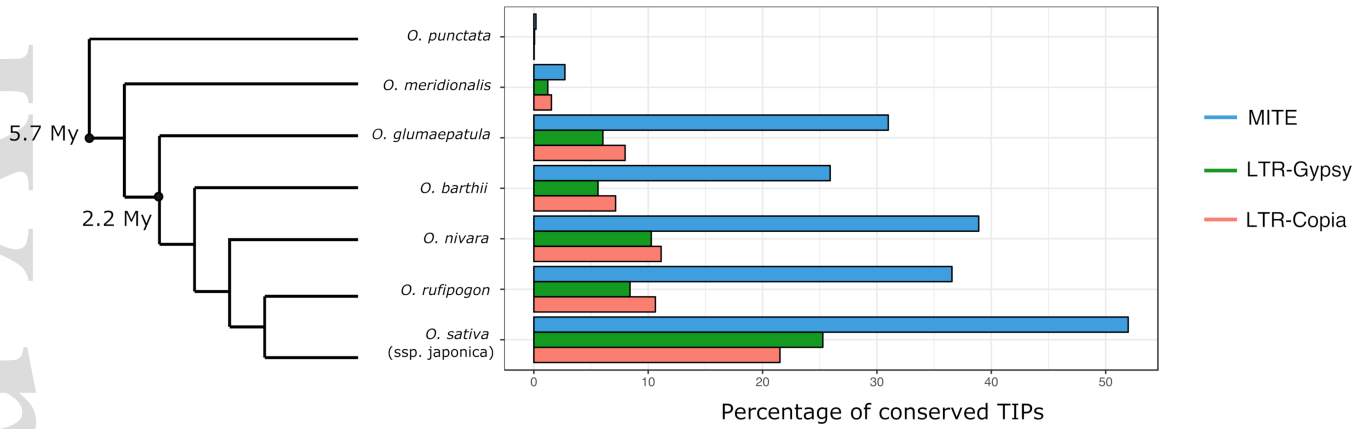


tpj_15277_f2.png

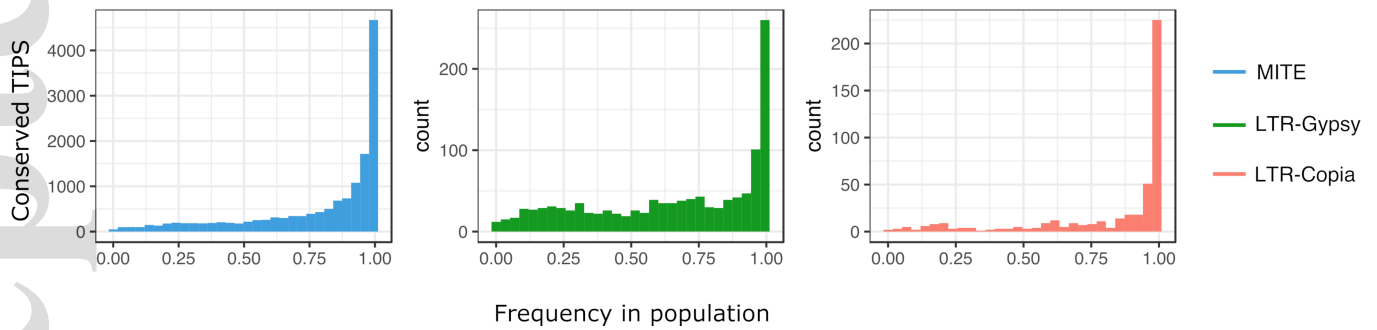
A



B

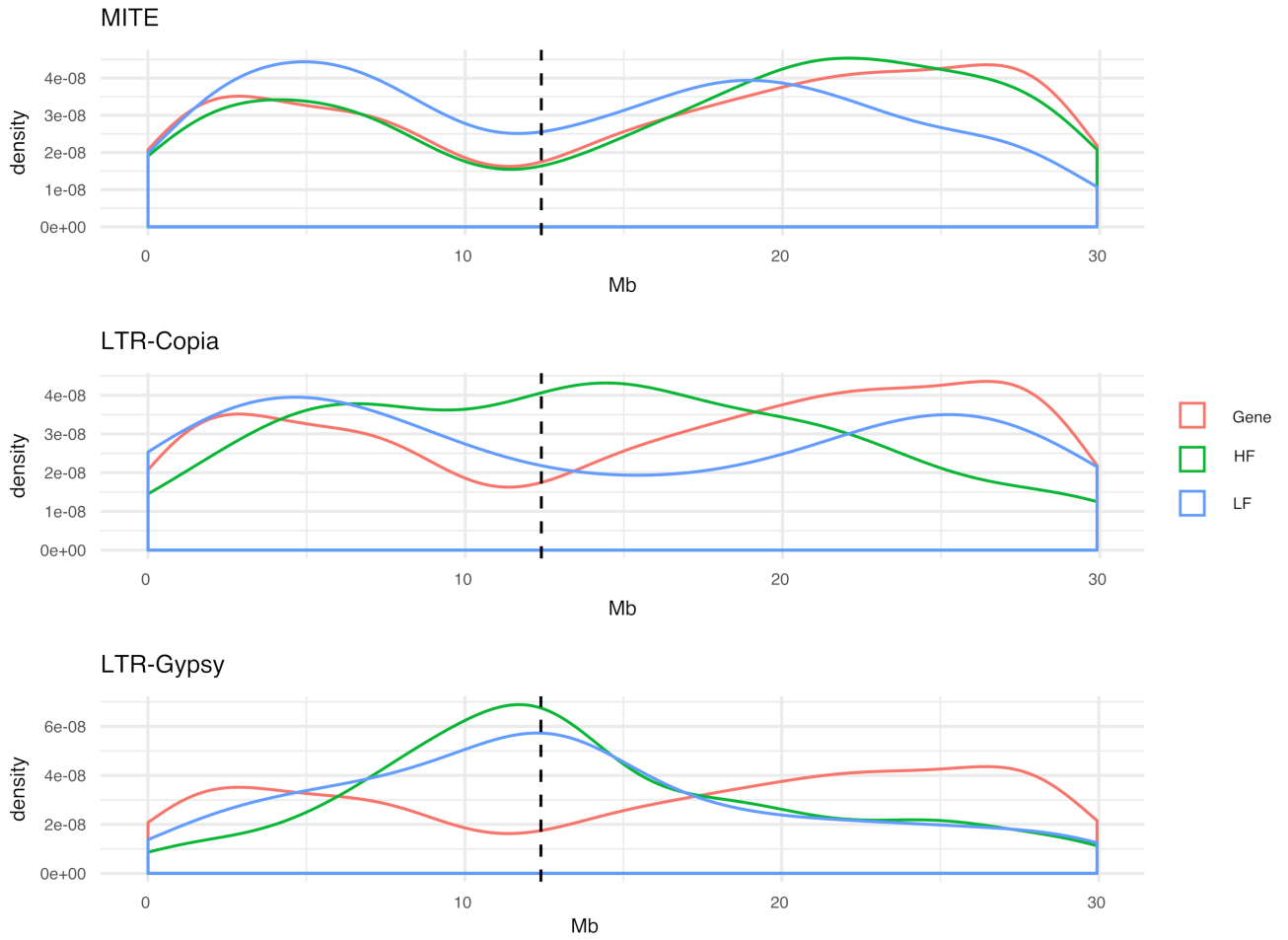


C

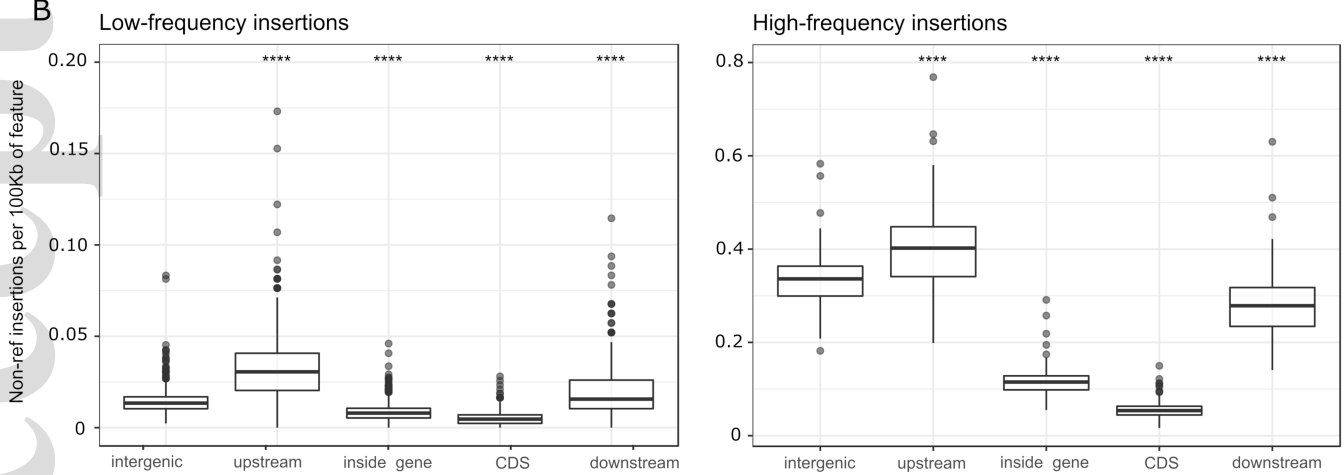


tpj_15277_f3.png

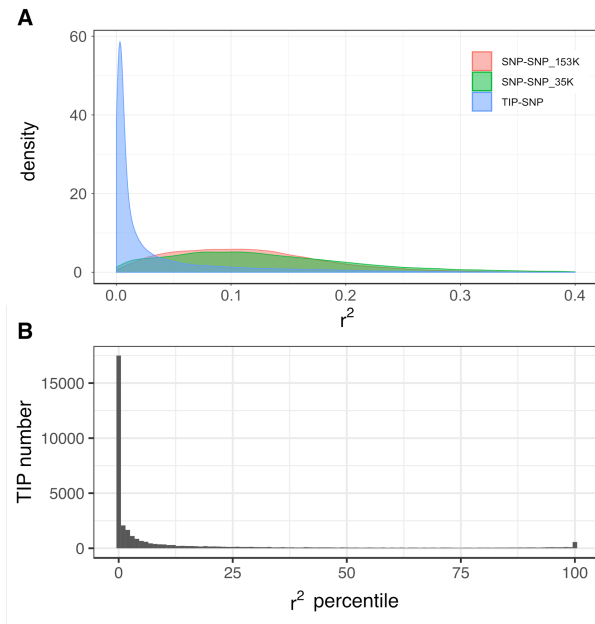
A



B

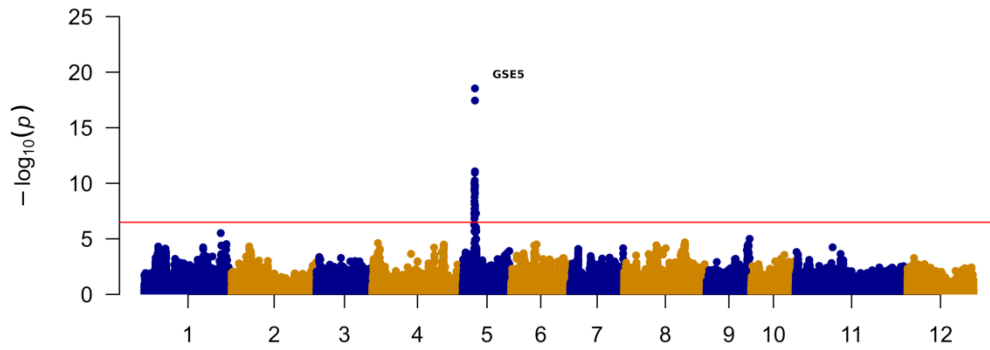


tpj_15277_f4.png

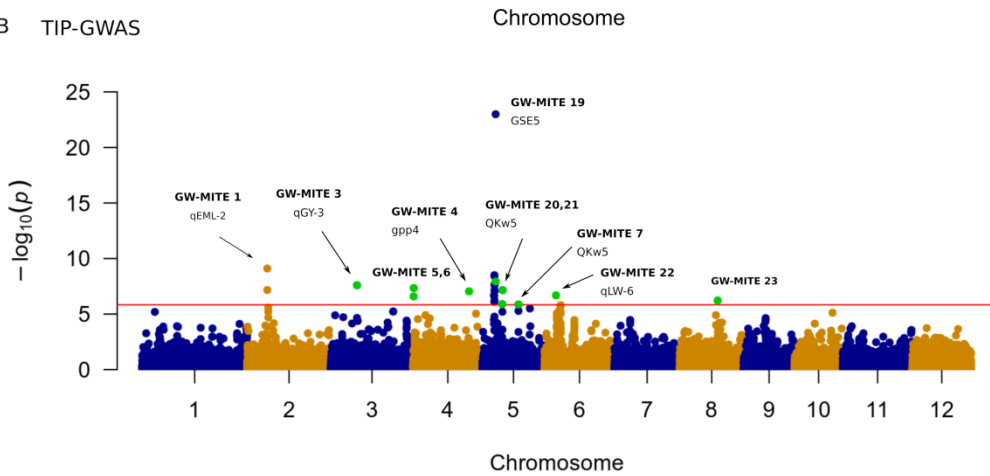


tpj_15277_f5.png

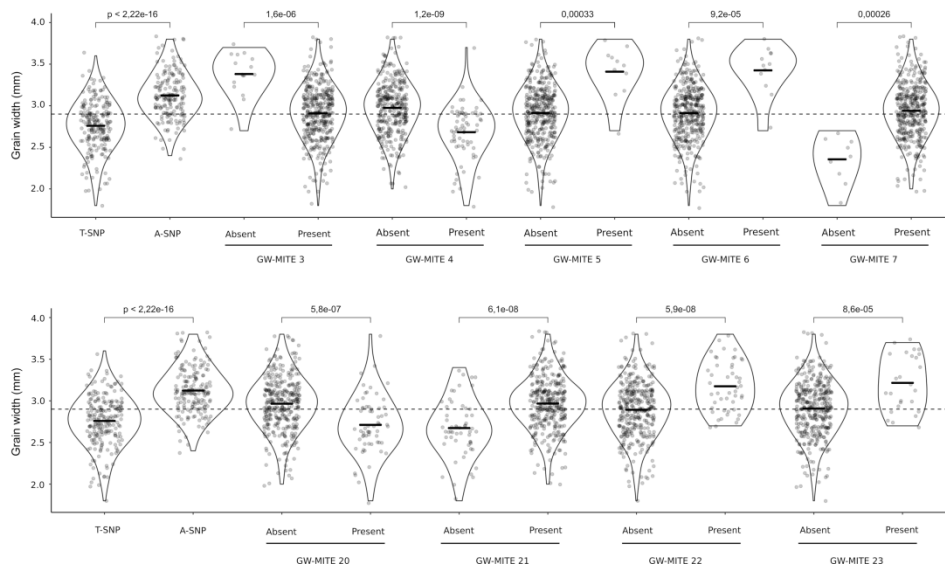
A SNP-GWAS



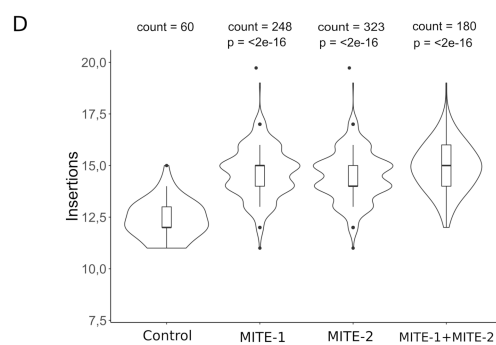
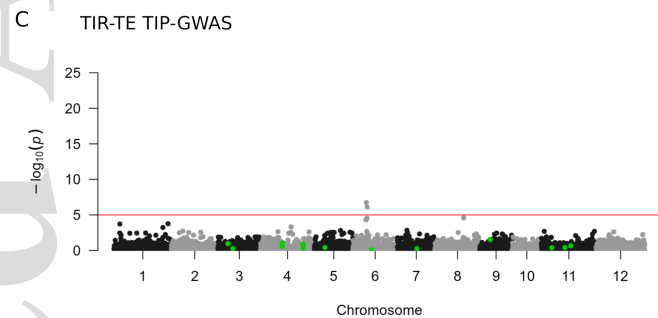
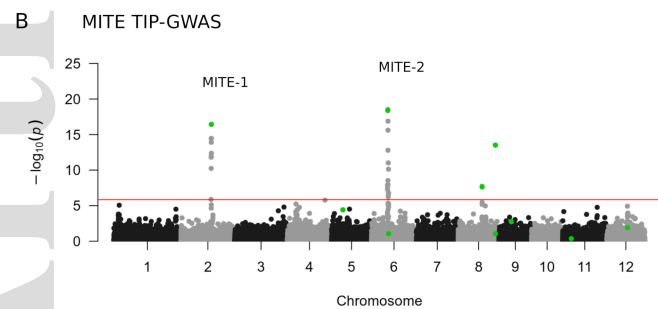
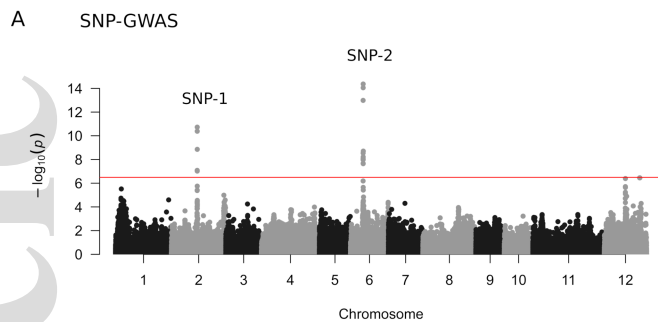
B TIP-GWAS



C



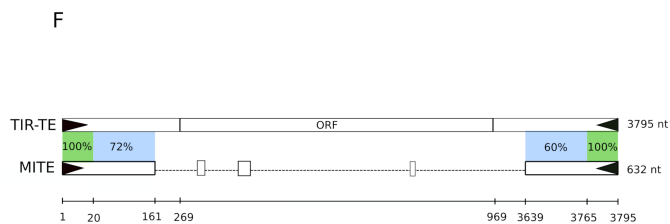
tpj_15277_f6.png



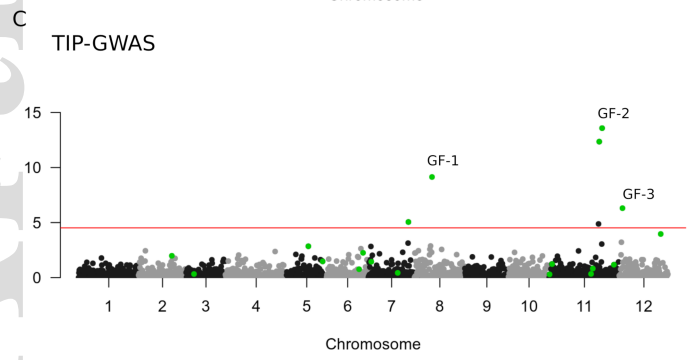
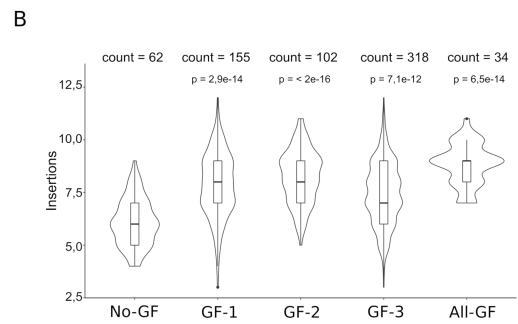
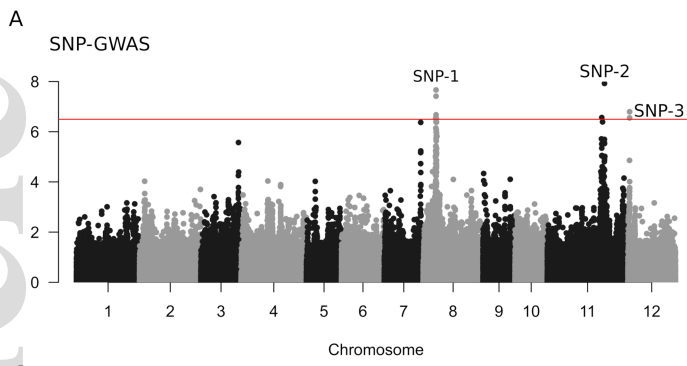
E

Allele	SNP-1	MITE-1 (30.5 Kb)
A	238	230
A/G	9	7
G	196	2

Allele	SNP-2	MITE-2 (85.7 Kb)
C	314	313
C/G	5	4
G	127	0



tpj_15277_f7.png

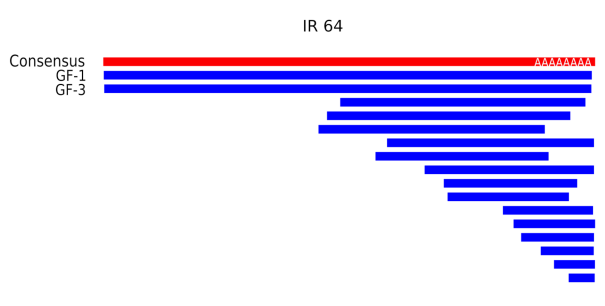
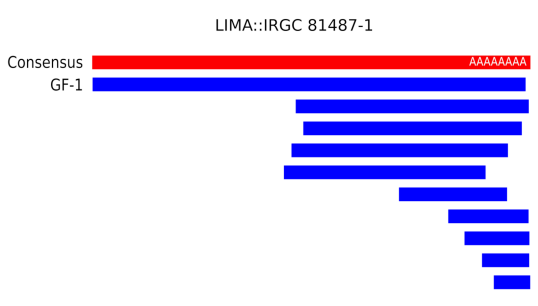


D

Allele	SNP-1	GF1 (9.7 Kb)
A (REF)	298	3
A/G	9	9
G	143	143

Allele	SNP-2	GF2 (72.7 Kb)
C (REF)	348	2
C/T	3	3
T	98	96

Allele	SNP-3	GF3 (71.6 Kb)
A (REF)	317	308
A/T	11	10
T	124	0



tpj_15277_f8.png