

1 Illuminating the transposon insertion landscape in plants using Cas9- 2 targeted Nanopore sequencing and a novel pipeline

3 Ilya Kirov^{1,2*}, Pavel Merkulov¹, Sofya Gvaramiya¹, Roman Komakhin¹, Murad Omarov¹, Maxim
4 Dudnikov^{1,2}, Alina Kocheshkova^{1,2}, Alexander Soloviev¹, Gennady Karlov¹, Mikhail Divashuk^{1,2}

5 ¹All-Russia Research Institute of Agricultural Biotechnology, Timiryazevskaya str. 42, Moscow 127550,
6 Russia;

7 ²Kurchatov Genomics Center of ARRIAB, All-Russia Research Institute of Agricultural Biotechnology,
8 Timiryazevskaya Street, 42, 127550 Moscow, Russia

9 Abstract

10 Transposable elements (TEs), which occupy significant portions of most plant genomes, are a major source
11 of genomic novelty, contributing to plant adaptation, speciation and new cultivar production. The often
12 large, complex genomes of plants make identifying TE insertions from short reads challenging, while
13 whole-genome sequencing remains expensive. To expand the toolbox for TE identification in plants, we
14 used the recently developed Cas9-targeted Nanopore sequencing (CANS) approach. Additionally, as no
15 current bioinformatics tools automatically detect TE insertions after CANS, we developed NanoCasTE, a
16 novel pipeline for target TE insertion discovery. We performed CANS of three copies of EVD
17 retrotransposons in wild-type *Arabidopsis thaliana* and obtained up to 40× coverage of the targets after
18 only a few hours of sequencing on a MinION sequencer. To estimate the ability to detect new TE
19 insertions, we exploited the *A. thaliana ddm1* mutant, which has elevated TE activity. Using CANS, we
20 detected 84% of these insertions in *ddm1* after generating only 4420 Nanopore reads (0.2× genome
21 coverage), and also unambiguously identified their locations, demonstrating the method's sensitivity.
22 CANS of pooled (~50 plants) *ddm1* plants captured >800 EVD insertions, especially in centromeric regions.
23 CANS also identified insertions of a Ty3/Gypsy retrotransposon in the genomes of two *Aegilops tauschii*
24 plants, a species with a large genome.

26 Introduction

27 Transposable elements (TEs) are a diverse group of genomic elements that can transpose across the
28 genome via RNA intermediates (Class I transposons, or retrotransposons) or via a “cut-and-paste”
29 mechanism (Class II, DNA transposons). TEs have been identified in almost all living organisms, where they
30 and their remnants can comprise up to approximately 93% of the genome (Rabanus-Wallace et al., 2021).
31 In plants, LTR retrotransposons are a major group of TEs with hundreds of ancient as well as very recent
32 insertions (Baduel et al., 2021). TE insertions (TEI) may alter gene expression, create a new transcriptional
33 repertoire and cause genomic reorganization (Chuong et al., 2017). The effects of TEs on genomic,
34 transcriptomic and proteomic activity have made them a major force in plant species diversity and
35 evolution (Lisch, 2013). TEIs are classified as genic (exonic and intronic) or intergenic (Sultana et al., 2017)
36 based on their locations in the genome. TEIs in or near genes are often more deleterious than single-
37 nucleotide polymorphism (SNPs) and may create null alleles. However, the true impacts of TEIs on the
38 functions of affected genes are more variable, ranging from complete inactivation to the generation of
39 novel transcriptional programs.

40 Progress in pangenome sequencing and genome assembly algorithms has helped unravel the genome-
41 wide TEI landscape in natural plant populations and germplasm (Domínguez et al., 2020; Baduel et al.,
42 2021). These studies demonstrated that TEIs are an important source of novelty for plant adaptation
43 (Domínguez et al., 2020). For example, recurrent insertions of TEs in *FLOWERING LOCUS C* (*FLC*; a key
44 determinant of flowering time variation in plants) in different *Arabidopsis thaliana* populations have
45 contributed to the local adaptation of this species (Baduel et al., 2021). Recent genome-wide association
46 study (GWAS) analysis using 602 resequenced genomes of wild and cultivated tomatoes revealed that at

47 least five of 17 analyzed traits were associated with TEI polymorphisms (Domínguez et al., 2020), including
48 fruit color and leaf morphology. Hence, ongoing TE activity has also been widely exploited by plant
49 breeders, as many novel and low-frequency alleles of the genes involved in agronomically important traits
50 originated via TEIs (Kobayashi et al., 2004; Jiang et al., 2009; Butelli et al., 2012; Domínguez et al., 2020).

51 TE transposition activity in plant cells is under strict epigenetic control, and most TEs are silenced under
52 non-stressed conditions (Slotkin and Martienssen, 2007; Nuthikattu et al., 2013). To overcome this
53 limitation, mutants that are defective in key genes of TE epigenetic control systems have been obtained.
54 These plants are fertile and have elevated TE activity, making them indispensable tools for the study of TE
55 biology (Mirouze et al., 2009; Tsukahara et al., 2009; Panda and Slotkin, 2020).

56 Identifying TEIs in the genome is challenging, primarily because most plant genomes possess a bulk of
57 repetitive DNA carrying many TE insertions. This ‘dark side’ of the genome is difficult to explore using the
58 short-read approach, a major tool for TEI identification in plants (Quadrana et al., 2016; Carpentier et al.,
59 2019; Chen et al., 2020; Domínguez et al., 2020; Baduel et al., 2021). Long-read sequencing using PacBio
60 and Nanopore technologies provides a dramatic increase in the resolution of the structures of repetitive
61 DNA sequences (Jung et al., 2019). While the whole-genome sequencing (WGS) approach theoretically
62 allows you to identify all TEIs in an individual plant, its application in plants with large genomes is still
63 quite expensive. In addition, WGS is even more difficult when the identification of TEIs in a plant
64 population is needed (Handsaker et al., 2011). Therefore, several approaches have been proposed to
65 capture specific genomic loci possessing TEIs without the need to sequence the rest of the genome (Li et
66 al., 2019). Enriched DNA samples are usually sequenced using short reads. Therefore, identifying TEIs that
67 have occurred in repeat-rich regions such as centromeres and in other TEs is challenging. These
68 approaches are mostly based on oligo-probe-mediated ‘fishing’ of target DNA fragments (Williams-Carrier
69 et al., 2010; Quadrana et al., 2016). A microfluidics-based method (Xdrop) was recently developed and
70 shown to be effective for capturing the integration sites of human papillomavirus 18 (Madsen et al., 2020).
71 While these approaches are effective, they require either a PCR amplification step, which may introduce
72 chimeric reads, or the use of specific equipment (for the Xdrop method).

73 The combination of TE capture with Nanopore sequencing has greatly improved integration site mapping
74 (Li et al., 2019; Madsen et al., 2020). A novel enrichment approach combining Cas9-mediated adaptor
75 ligation and long-read sequencing (CANS) was recently proposed (Gabrieli et al., 2018; Gilpatrick et al.,
76 2020; Stangl et al., 2020). The principle of this method is based on the finding that Nanopore sequencing
77 adapters preferentially ligate to Cas9-cleaved DNA sites. To decrease the ligation of random adapters, the
78 DNA is initially dephosphorylated. The advantages of this approach over oligo-based enrichment include
79 the lack of a PCR amplification step, the ability to detect DNA methylation and the shorter time needed
80 to complete the protocol. CANS has been applied to sequence target genes in humans (Gilpatrick et al.,
81 2020) and plants (López-Girona et al., 2020). The use of CANS to identify TE integration sites in human
82 cells was recently reported (McDonald et al., 2021). However, the application of CANS for TE integration
83 identification in plants has not yet been reported.

84 Here, we utilized CANS and developed a novel pipeline (NanoCasTE) for the identification of TEIs in plants
85 with small (*A. thaliana*, 157 Mb/1C) and large (*Aegilops tauschii*, 4.3 Gb/1C) genomes. Using the *A.*
86 *thaliana ddm1* mutant, which carries T-DNA and possesses elevated TE activity, we were able to detect
87 most (84%, 21 of 25 TEIs) of the genetically inherited insertions of EVD retrotransposons in an individual
88 plant with only 0.2× genome coverage obtained by CANS reads. This coverage is two orders of magnitude
89 lower than what is needed for TEI identification based on whole-genome sequencing. In addition, >800
90 new TEIs could be detected using a pool of *ddm1* plants, resulting in a high-density EVD integration map.
91 Finally, using *Ae. tauschii* as a model, we demonstrated that CANS is also well suited for TEI identification
92 in plants with large, complex genomes. The application of CANS to two *Ae. tauschii* plants resulted in the
93 capture of hundreds of insertions of a Ty3/Gypsy retrotransposon. Altogether, our results indicate that
94 CANS and NanoCasTE are effective tools to capture TEIs in diverse plant species, opening an avenue for

95 obtaining a deeper understanding of the TE integration landscapes in plant genomes and exploring the
96 potential of TE ‘mutagenesis’ in plant breeding.

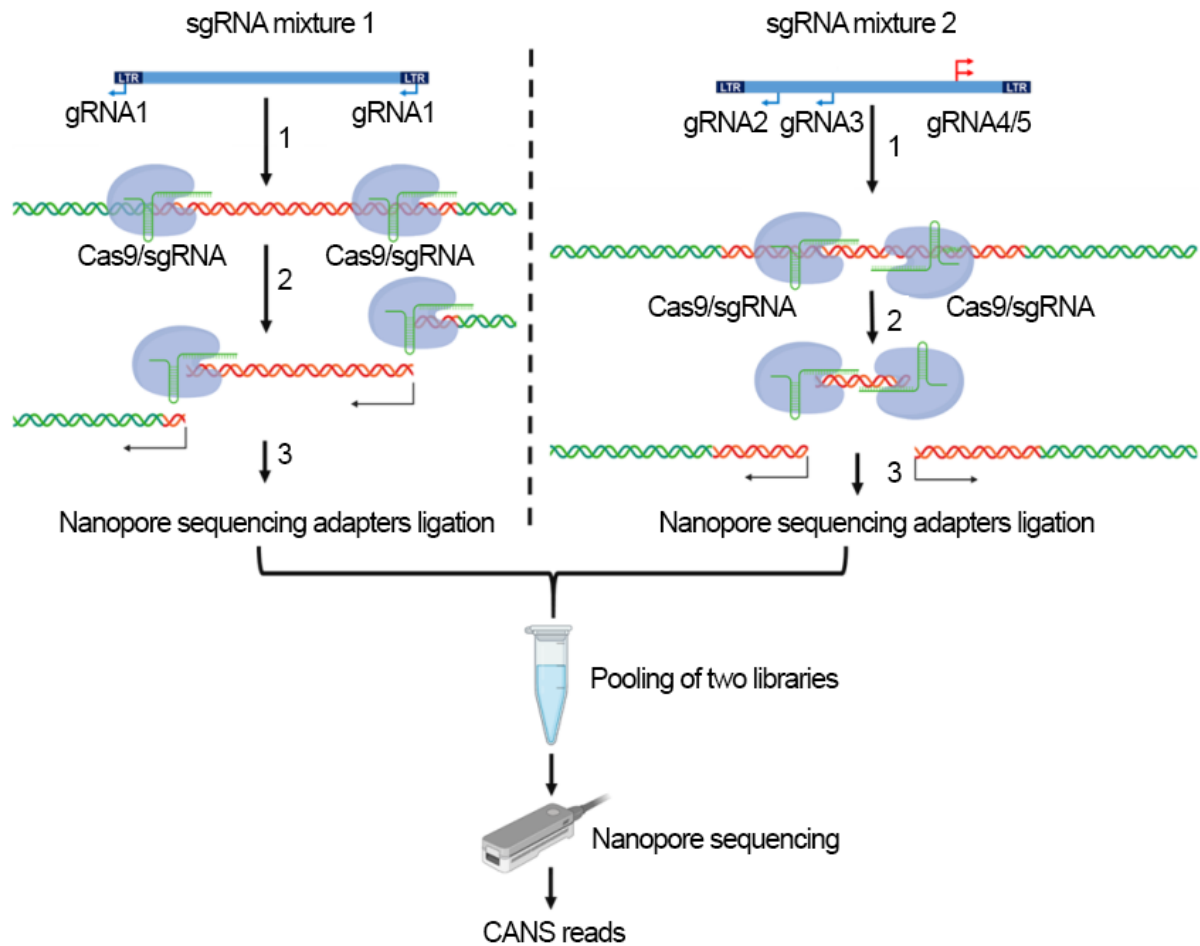
97

98 Results

99 Cas9-mediated targeted sequencing of three copies of the EVD retrotransposon in *A. thaliana*

100 To demonstrate that CANS is suitable for tracing individual transposon insertions in plants, we chose the
101 EVD5 (5333 bp) LTR retrotransposon (Mirouze et al., 2009), which can generate new insertions in some
102 *A. thaliana* mutants (e.g. *met1* and *ddm1* [*decrease in DNA methylation 1*]) (Tsukahara et al., 2009)). EVD5
103 belongs to the Evadé subfamily (ATCOPIA93 family) of Ty1/Copia retrotransposons, which includes other
104 EVD elements whose activity was not detected. Of these, EVD1 (AT1TE41575, 1548 bp) and EVD2
105 (AT1TE41580, 5336 bp) are densely located elements with high sequence similarity (~99%) to EVD5
106 (Mirouze et al., 2009).

107 We designed a set of seven guide RNAs located in LTRs and inter-LTR regions of EVD5 and EVD2. We
108 predicted that one of these gRNAs also recognizes the EVD1 copy. Because CANS is strand-specific, we
109 designed gRNAs to direct Nanopore sequencing to sequence (1) EVD retrotransposons *per se* and (2)
110 upstream and downstream sequences flanking retrotransposon insertions. The gRNAs were transcribed *in*
111 *vitro* as single-guide RNAs (sgRNAs). We carried out two runs on a MinION sequencer using seven and five
112 gRNA sets. The pilot experiment using all gRNAs as one mixture yielded 15,777 reads (144Mb total bases),
113 with 231 (1.4%) on-target reads detected after 4.4 hours of sequencing. Because of the many short reads
114 produced from the internal gRNA cutting site, for our second experiment, we decided to use five of the
115 most efficient sgRNAs divided into two mixtures to prevent the generation of short reads. The selected
116 gRNAs, designated gRNA1 to gRNA5, included one LTR gRNA (gRNA1, sgRNA mixture 1) and four gRNAs
117 (sgRNA mixture 2) located in the inter-LTR region (Figure 1).

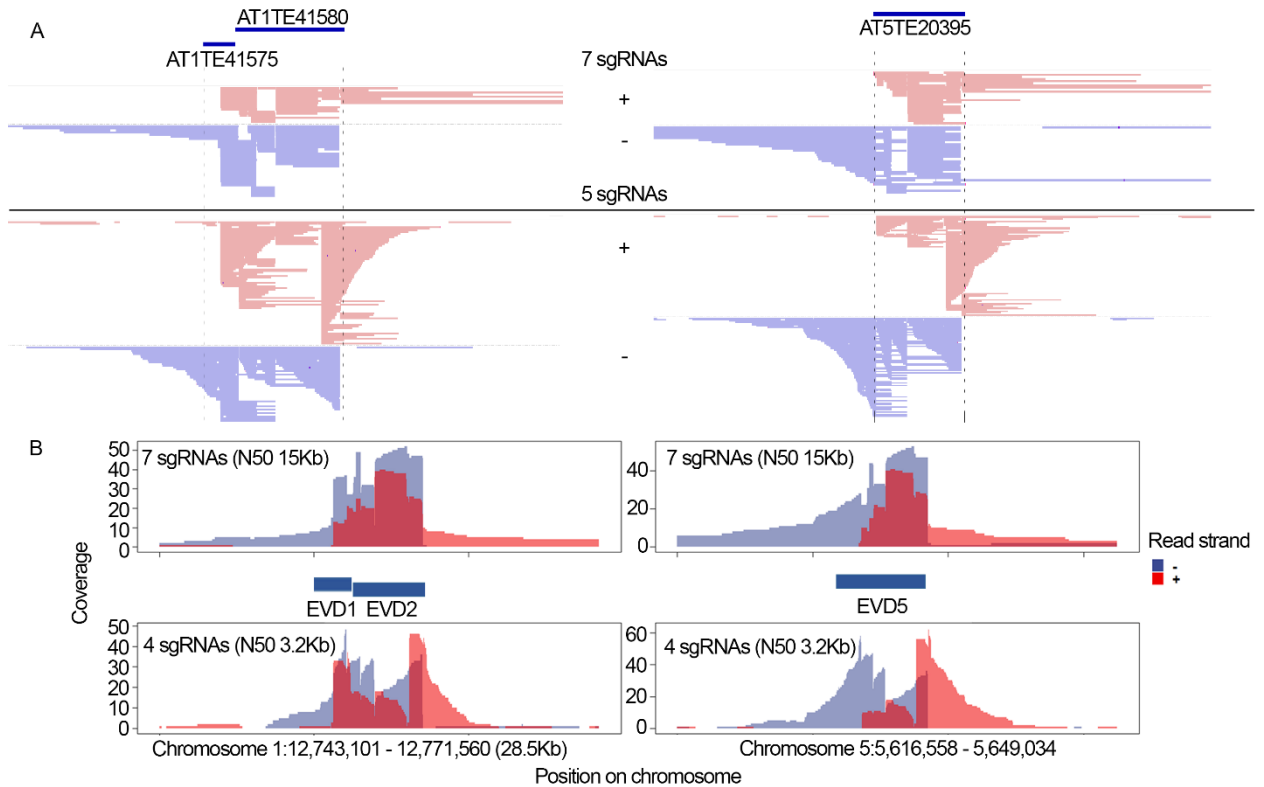


118

119 **Figure 1. Schematic view of CANS of LTR retrotransposon insertions using five gRNAs divided into two**
120 **mixtures.** (1), (2) and (3) show the RNP binding, target DNA cleavage and ONP library preparation steps
121 of the CANS protocol, respectively. Blue and red horizontal lines correspond to CANS reads mapped to
122 negative and positive strands, respectively, compared to the reference genome.

123

124 The experiment using five gRNAs yielded 88,000 reads (207Mb total bases), with 259 (0.3%) on-target
125 reads detected after 4 hours of sequencing on the MinION flow cell. Both CANS runs using Col-0 genomic
126 DNA resulted in up to 40× coverage of the target retrotransposons EVD5 and EVD1/2 (Figure 2A, 2B).
127 Mapped reads from the first and second runs covered 35,000 bp and 14,000 bp regions, respectively,
128 including target EVDs and flanking regions. The broad coverage of the flanking sequences demonstrates
129 the advantage of CANS for detecting TE insertions in the genome. Inspection of the gRNA locations and
130 the read distribution revealed that CANS had high strand specificity, with up to 82% of reads mapped to
131 strands with 3'-ends unprotected by Cas9 and possessing PAM (NGG) sequences. This is another
132 advantage of CANS, as it allows the coverage of the target region to be increased. We also noticed that
133 the bulk of reads from internal parts of EVD2 and EVD5 were unambiguously assigned to the
134 corresponding EVD copy, corroborating previous reports that cDNA ONP reads have high mappability
135 (Panda and Slotkin, 2020).



136

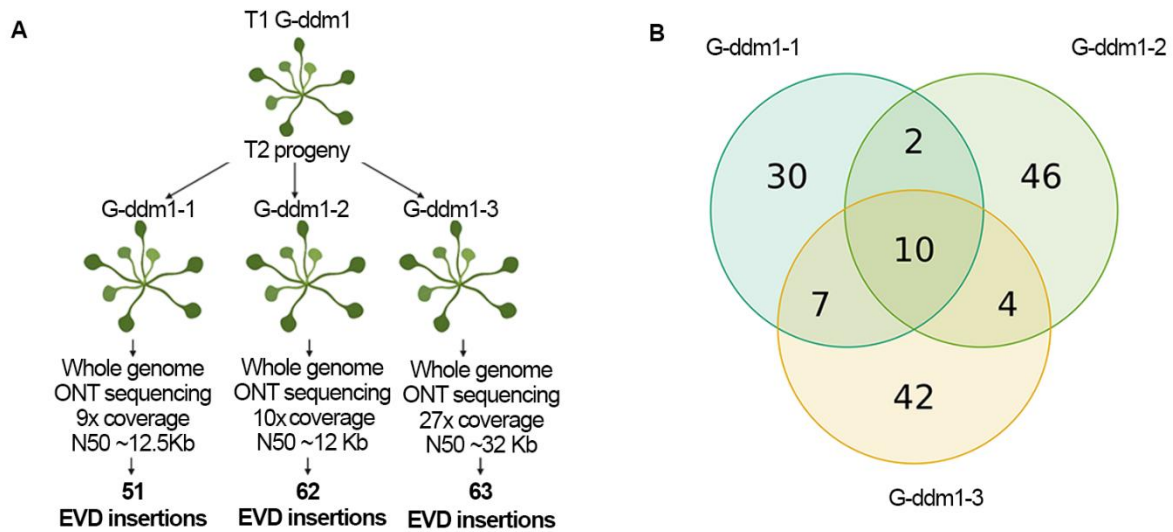
137 **Figure 2. Targeted enrichment of EVD loci following CANS using a set of 7 or 5 gRNAs.** (A) IGV view of
 138 Nanopore read alignment to EVD1 (AT1TE41575), EVD2 (AT1TE41580) and EVD5 (AT5TE20395)
 139 retrotransposon copies. (B) Coverage plot of EVD1, EVD2 and EVD5 using forward (red) and reverse (blue)
 140 reads obtained following CANS with 7 and 5 sgRNAs.

141 Thus, using CANS, we sequenced EVD1, EVD2 and EVD5 retrotransposon copies and their flanking
 142 sequencing with up to 40× coverage, demonstrating the advantage of this approach for detecting TE
 143 insertions.

144 [Detection of new EVD insertions in the genome using whole-genome Nanopore sequencing](#)

145 To assess whether new EVD insertions could be detected using CANS, we exploited the *A. thaliana ddm1*
 146 mutant, whose EVD5 activity was previously described (Tsukahara et al., 2009). We created transgenic T₂
 147 *ddm1* plants carrying T-DNA insertions of the EVD5 GAG open reading frame, a target for one of the five
 148 gRNAs used for CANS (G-*ddm1* plants). This allowed us to simultaneously detect EVD insertions and T-
 149 DNA integration sites using a single sgRNA set and to obtain information about the sufficiency of CANS
 150 sequencing depth. To identify all EVD insertions in a G-*ddm1* plant, we performed whole-genome
 151 Nanopore sequencing of one plant (plant G-*ddm1*-3). We collected 194,354 high-quality reads with N50
 152 ~32Kb, corresponding to 27× coverage (the assembled genome size of *A. thaliana* is 119 Mb) (Figure 3A).
 153 We designed a pipeline called NanoWgsTE to identify EVD and T-DNA insertion sites using whole-genome
 154 Nanopore sequencing data. Using this pipeline, we identified 63 EVD and two T-DNA insertions. However,
 155 37 TEIs (58%) were supported by only one read, which is significantly less than we expected based on
 156 Poisson distribution (p -value = 1.85e-05 with $\lambda = 13.5$). These results may point to the occurrence of
 157 somatic EVD insertions that are shared by a subset of cells in an individual plant. Because this would affect
 158 the interpretation of the results, we performed Nanopore sequencing of two additional G-*ddm1* plants
 159 (G-*ddm1*-1 and G-*ddm1*-2) that are siblings of G-*ddm1*-3 (progeny in the same family, Figure 3A). We
 160 obtained 9× and 10× coverage with N50 ~12Kb (Figure 3A). The NanoWgsTE pipeline revealed 51 and 62
 161 EVD insertions in G-*ddm1*-1 and G-*ddm1*-2 plants, respectively. Of these, 34 (67%) and 45 (70%) TEIs were
 162 supported by one read, corroborating the results of G-*ddm1*-3 sequencing. Next, to identify truly
 163 genetically inherited TEIs, we estimated the number of EVD TEIs shared by at least two G-*ddm1* plants
 164 and found 23 EVD insertions that were common to at least two sibling plants and 10 TEIs that were shared

165 by all three plants (Figure 3B). We then estimated the number of reads supporting these insertions (TEI
166 reads) according to the NanoWgsTE report and found a significantly higher number of TEI reads for
167 genetically inherited TEIs than for other TEIs (Supplementary Figure S1).



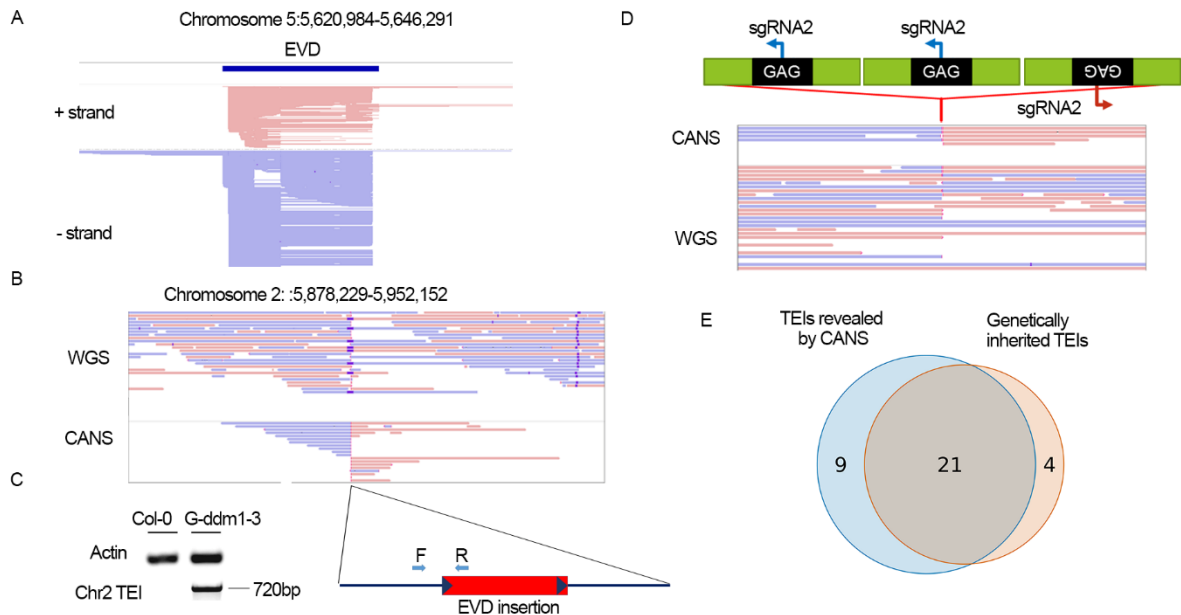
168

169 **Figure 3. Detection of new EVD and T-DNA insertions in G-ddm1 plants.** (A) Overview of whole-genome
170 Nanopore sequencing of three T₂ sibling G-ddm1 plants and the numbers of EVD insertions detected using
171 the NanoWgsTE pipeline. (B) Venn diagram showing the numbers of unique and common EVD insertions
172 in three G-ddm1 plants identified by whole-genome Nanopore sequencing.

173 To further validate the suitability of CANS using these results, we selected 21 genetically inherited TEIs
174 from a G-ddm1-3 plant. We also included four other TEIs of G-ddm1-3 that were supported by >2 TEI
175 reads, i.e., the number of TEI reads found for >86% (18 of 21) of truly genetically inherited TEIs. Thus,
176 using whole-genome Nanopore sequencing, we identified 25 TEIs in G-ddm1-3 plants that were
177 genetically inherited (Supplemental Data Set 1).

178 CANS of *ddm1* and the NanoCasTE pipeline

179 We carried out CANS using five gRNAs and generated 4420 reads with N50 ~11 Kb. Following read
180 mapping to the genome, we looked at the original EVD copy on chromosome 5 to estimate the sequence
181 coverage (Figure 4A). The EVD sequence had up to 400× coverage by 1170 reads (25% of all CANS reads)
182 (Supplementary Figure S2). It is worth noting that the reads spanning EVD sequences from LTR to LTR
183 could not be assigned to the certain novel EVD insertion, as all EVD copies had identical sequences.
184 Consequently, the observed EVD coverage primarily resulted from the alignment of Nanopore reads that
185 originated from different EVD insertions in the genome. To support the notion that the coverage we
186 obtained is sufficient for EVD TEI identification, we analyzed EVD-genome junctions and found that they
187 had ~8× coverage, which is sufficient for insertion detection (Figure 4A). As an additional control, we also
188 checked the CANS reads covering a hemizygous T-DNA insertion site carrying a single gRNA position
189 (gRNA2, Figure 4D) and found that even using a single gRNA, CANS generated ~4× coverage of the T-DNA
190 integration site. These observations indicate that the number of CANS reads obtained is sufficient for TEI
191 identification.



192

193 **Figure 4. EVD coverage and insertion detection in a G-ddm1-3 plant using CANS.** (A) IGV snapshot of
194 coverage of the original copy of EVD on Chromosome 5 using forward (red) and reverse (blue) ONT reads
195 obtained after CANS of G-ddm1-3 plant. (B) IGV browser snapshot of one of the 25 genetically inherited
196 EVD insertions detected by whole-genome Nanopore sequencing and CANS. (C) PCR validation of the EVD
197 insertion on Chromosome 2 (left) using primers based on the EVD and TEI flanking region (right). (D)
198 Inserted T-DNA structure (upper panel), gRNA2 positions and its coverage based on CANS and whole-
199 genome Nanopore sequencing reads. (E) Venn diagram of the numbers of genetically inherited and CANS-
200 detected EVD insertions in a G-ddm1-3 plant.

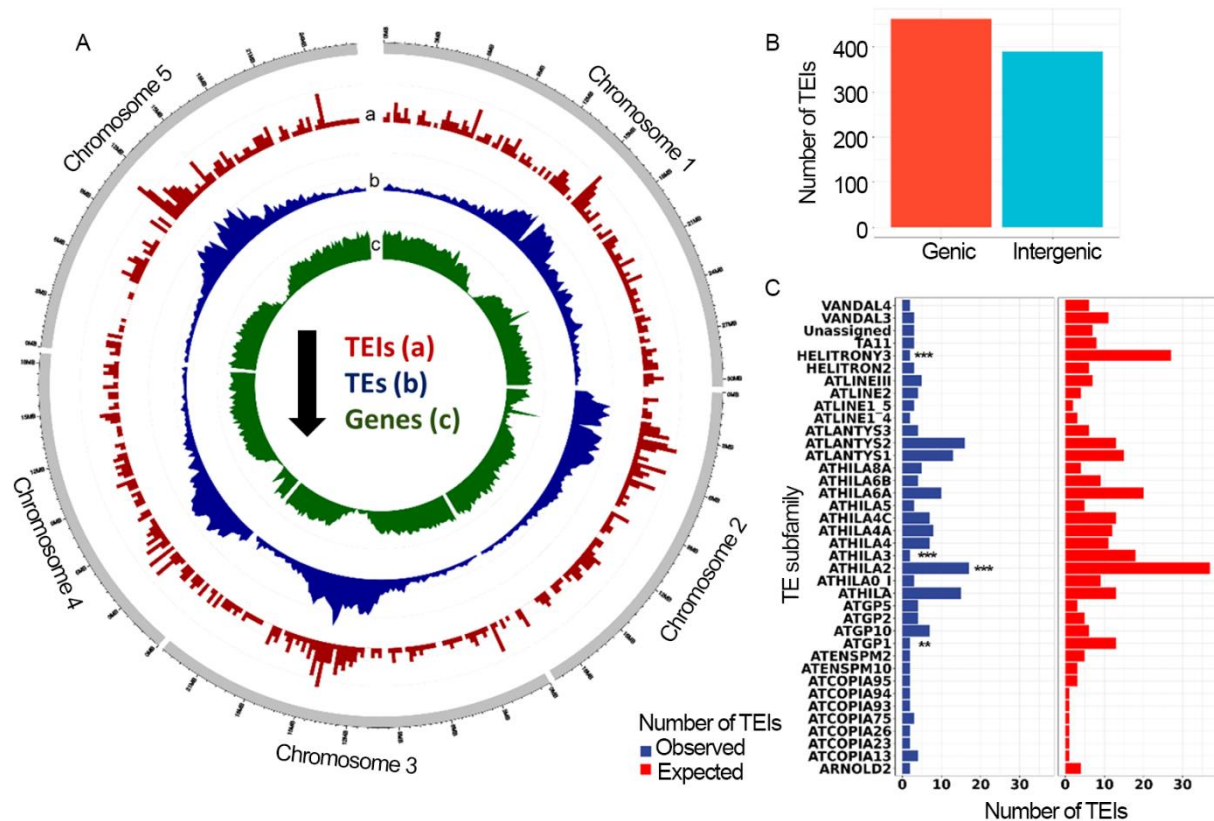
201

202 To automate the detection of TE and T-DNA insertions, we developed a novel pipeline called NanoCasTE
203 (see Methods, Figure 7). By applying NanoCasTE to the CANS dataset, we detected 30 EVD insertions. Of
204 these, 22 were high-confidence TEIs (TEIs supported by two or more reads). An example of one TEI located
205 on chromosome 2 that was homozygous according to the WGS data (see above) is shown in Figure 4B.
206 This TEI was also validated by PCR with primers designed based on EVD and insertion sites (Figure 4C).
207 Next, we estimated how many genetically inherited TEIs were identified by CANS and found that 21 of 25
208 (84%) genetically inherited TEIs were captured by CANS (Figure 4E). These results suggest that even the
209 low genome coverage (0.2 \times in our experiment) of CANS may be sufficient to capture most TEIs in *A.*
210 *thaliana*, highlighting the good sensitivity of the method. Additionally, the newly developed NanoCasTE
211 pipeline enables the rapid identification of novel TEIs.

212

213 Capture of EVD insertions in a *ddm1* population

214 Because EVD can generate new TE copies in individual *ddm1* plants, we asked whether we could detect
215 new insertions in a population of approximately 50 *ddm1* plants. To address this issue, we isolated high
216 molecular weight DNA from the pooled plant material and carried out CANS. After 2.5 hours of MinION
217 sequencing, we generated 18,246 Nanopore reads with N50 \sim 5.5 Kb. Using the NanoCasTE pipeline, we
218 detected 851 TEIs, including 29 high-confidence TEIs (supported by two or more reads), in the *ddm1*
219 population. Thirteen TEIs were PCR validated using DNA from individual plants, resulting in 1 to 6 plants
220 carrying a single TEI. Analysis of the chromosome-wide distribution of the TEIs showed that they tended
221 to be clustered in the pericentromeric regions of all chromosomes, while the density of TEIs in
222 chromosome arms was significantly reduced (Figure 5A).



223

224 **Figure 5. Genome organization of EVD insertions detected by CANS in a pooled (~50 plant) *ddm1***
 225 **sample. (A)** Chromosome distribution of TEIs detected by CANS (a), all *A. thaliana* TEs (Panda & Slotkin
 226 2020) (b) and genes (c) along the *A. thaliana* chromosomes. **(B)** Bar plot showing the number of genic
 227 and intergenic TEIs detected by CANS. **(C)** Bar plot showing the number of EVD insertions in members of
 228 different TE subfamilies of *A. thaliana*.

229 Classification of these TEIs revealed that 54% (462) of them are located in genic regions, including 403
 230 (87%) exonic and 59 (13%) intronic TEIs (Figure 5B). This value is expected by chance based on the total
 231 length of the genic regions in the assembled (119 Mb, TAIR10) *A. thaliana* genome (68.3 M, or 57% of the
 232 sequenced genome). Gene Ontology (GO) analysis revealed no enrichment of any GO categories in the
 233 set of genes with TEIs, suggesting that most EVD insertions in *ddm1* are randomly distributed among genes
 234 with different functions.

235 Among intergenic TEIs (389, 46%), 241 (28%) were located in the sequences of other TEs. This value is
 236 significantly higher (chi-squared test p -value = 0.0004) than expected by chance based on the total length
 237 of the annotated TEs in the *A. thaliana* genome (24.9 Mb, or 21% (Panda and Slotkin, 2020)). Moreover,
 238 20 TEIs had 2 or 3 different TEIs. We analyzed the number of TEIs in TEIs of different subfamilies and found
 239 that members of four subfamilies were significantly underrepresented (Fisher's exact test p -value <0.01)
 240 in a set of TEIs with detected TEIs, including HELITRONY3 (p -value = 6.78e-07), Athila2 (p -value = 0.004),
 241 Athila3 (p -value = 0.0002) and ATGP1 (p -value=0.004) (Figure 5C). This finding can be partly attributed to
 242 the differences in chromosome distribution between TEIs of these families and EVD TEIs. For example,
 243 most (1104 of 1399, 79%) of the members of HELITRONY3 are located in chromosome arms, where the
 244 EVD TEI frequency is much lower than in the pericentromeric regions.

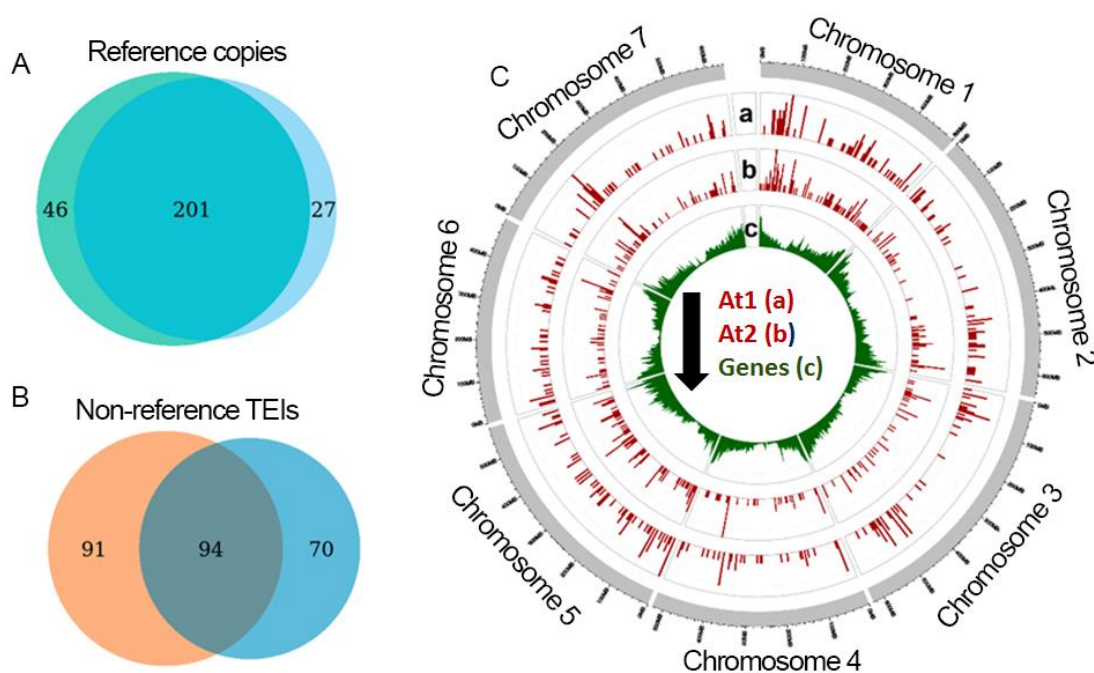
245 Thus, using CANS coupled with the NanoCasTE pipeline, we built a high-density physical map of EVD5
 246 insertion sites in the *A. thaliana* genome and showed that EVD insertions are biased toward
 247 pericentromeric regions and that EVD tends to mobilize frequently into other TEIs.

248

249

250 Targeting of LTR retrotransposons in *Ae. tauschii*

251 To examine whether NanoCasTE could be utilized to trace insertions of multicopy TEs in a plant with a
252 large genome, we used *Ae. tauschii* as a model. *Ae. tauschii*, a donor of the D-genome of bread wheat
253 (*Triticum aestivum*), has a relatively large genome (~4.3 Gb (Jia et al., 2013)). To choose a target for
254 NanoCasTE sequencing, we aimed to identify multicopy TEs with recent copies in the genome. We
255 selected a Ty3/Gypsy LTR retrotransposon (hereafter called Aty3_169) of the Retand lineage, with 327
256 BLAST hits per genome and a number of recent insertions that occurred <0.5 million years ago. This TE is
257 >13 Kb in length and is specifically found in the bread wheat D-genome (Supplementary Figure S3),
258 suggesting that its transposition activity occurred after the speciation of the A, B and D genome
259 progenitors of bread wheat. We synthesized five sgRNAs including 2 gRNAs (gRNA1, gRNA2) located on
260 LTR sequences in the both ends of the Aty_169 and three gRNAs located in the inter-LTR region (gRNA3-
261 5). We carried out CANS using two sgRNA mixtures: sgRNA Mixture 1 (gRNA1, gRNA2) and sgRNA Mixture
262 2 (gRNA3-5). Genomic DNA from two *Ae. tauschii* plants (At1 and At2) was sequenced using two MinION
263 flow cells (one flow cell per plant). In each experiment, we obtained ~70,000 reads. First, we evaluated
264 the coverage of the reference copies of Aty3_169 based on reads. For this, only high-quality mapped reads
265 (MQ > 40, no supplementary alignments allowed, <30 bp end clipping) were used. This analysis resulted
266 in 228 (69%) and 247 (75%) reference copies covered by two or more ONT reads, respectively. Of these,
267 201 reference copies were common between the two plants (Figure 6A).



268

269 **Figure 6. Insertion identification for a Ty3/Gypsy LTR retrotransposon (Aty3_169) of the Retand lineage**
270 **in the genomes of two *Ae. tauschii* plants. (A)** Numbers of reference insertions covered by two or more
271 ONT reads. **(B)** Numbers of non-reference insertions covered by two or more Nanopore reads. **(C)**
272 Chromosome distribution of At1 and At2 TEIs and annotated genes along the *Ae. tauschii* chromosomes.

273

274 To identify non-reference insertions, we applied the NanoCasTE pipeline to the At1 and At2 reads. We
275 identified 185 and 164 non-reference TEIs in the At1 and At2 genomes, respectively. Of these, 94
276 insertions were common to both datasets (Figure 6B). We designed primers based on 10 TEIs, and 7 of
277 them were validated by PCR.

278 To estimate any biases in the transposition of Aty3_169 in the *Ae. tauschii* genome, we evaluated the
279 chromosome distribution of the TEIs. In general, TEIs were randomly distributed along the chromosomes,

280 although a few regions of chromosomes 1, 3, 4 and 7 were highly enriched by Aty3_169 insertions (Figure
281 6C). In the future, it would be interesting to assess whether these hot spots of Aty3_169 insertions could
282 be attributed to the recent introgressions.

283

284 Discussion

285 Here, we demonstrated that CANS is a valuable tool for identifying TE insertions (TEIs) in both repetitive
286 (e.g. TE rich) and genic regions of plant genomes. CANS-based detection of TEIs has several advantages
287 over the existing method of TEI enrichment. First, the CANS protocol, including RNP complex formation,
288 DNA cleavage and sequencing library preparation, is rapid and can be completed in 4-5 hours. Second,
289 because a PCR amplification step is not needed, the resulting data can be used for epigenetic modification
290 calling, although higher coverage is required for this process (Gabrieli et al., 2018; Ni et al., 2021). Third,
291 long reads have high mappability to the genome, paving the way for TE insertion detection, even in repeat-
292 rich regions of the genome such as centromeres and heterochromatin. Fourth, utilizing NanoCasTE allows
293 TEIs to be identified easily using CANS data, even if hundreds of insertions are captured. Finally, CANS can
294 be applied to small as well as large plant genomes. Despite the many advantages of CANS for TEI detection
295 in plants, there are also several constraints to the method that need to be overcome in the future: (1) a
296 large amount of high-molecular-weight genomic DNA (>5 µg) is required; (2) deciphering the TE sequence
297 of the insertion is challenging and requires specific sets of overlapping sgRNAs; (3) only a small subset of
298 sequencing pores (5-10%) are sequencing at the same time, and therefore the sequencing efficiency is
299 lower than for whole-genome sequencing (90-98% of sequencing pores), and (4) the genotype of TEIs
300 (homozygote or heterozygote) cannot be determined using CANS.

301 While whole-genome Nanopore sequencing provides excellent resolution to potentially detect all TEIs in
302 a genome, it is expensive when the genome is relatively large. For example, in our hands and based on
303 previous reports using different species, it is challenging to generate more than 8 Gb of data per single
304 MinION flow cell from plant material (Lee et al., 2019; Schalamun et al., 2019; Dmitriev et al., 2021).
305 Obtaining only 4400 CANS reads (equivalent to 0.2× coverage of the *A. thaliana* genome) on a single flow
306 cell, we were able to detect 86% of the genetically inherited EVD insertions. Additionally, we also
307 unraveled the positions of T-DNA containing the GAG region of the EVD retrotransposon and carrying a
308 target site for one of the five sgRNAs. Thus, even a single sgRNA can be effective for insertion detection.
309 However, a higher number of sgRNAs and cleavage sites in genomic DNA may produce a higher proportion
310 of target reads, as we also observed in our pilot experiment with seven sgRNAs. The use of a higher
311 number of sgRNAs per target was also proposed to be beneficial for CANS of target genes (Gilpatrick et
312 al., 2020). This is an important property of the CANS approach that allows more targets to be included per
313 experiment without losing the efficiency of CANS. Based on this assumption, insertions from many TEs
314 can be captured in one experiment using specifically designed and pooled sgRNAs. This approach was
315 recently applied to detect TE insertions in the human genome (McDonald et al., 2021).

316 Here, we identified >800 TEIs in a small *ddm1* population after 2 hours of sequencing using only a part of
317 the capacity of MinION. This experiment provided a proof of concept that CANS is a valuable tool for
318 identifying TEIs in a population. A study on the distribution of TE insertions in the *A. thaliana* 1001
319 genomes population demonstrated that TEs are a major source of large-effect alleles that can be fixed in
320 the populations, contributing to the local adaptation of different accessions (Baduel et al., 2021). Indeed,
321 stress may induce TE activation and consequently the origination of new alleles. TEI-derived alleles have
322 spontaneously occurred in fields of crop species and have long been exploited in plant breeding programs
323 (Kobayashi et al., 2004; Jiang et al., 2009; Butelli et al., 2012; Domínguez et al., 2020). Although the
324 discovery of such alleles in crop fields using whole-genome sequencing is very challenging, CANS might
325 represent a valuable alternative to capture TEI-derived alleles in a population of cultivated plants.
326 However, prior knowledge of active TEs capable of transposing is required for such experiments.
327 Fortunately, this information is available for some plants. For example, we recently identified active LTR

328 retrotransposons in sunflower (Kirov et al., 2020a) and triticale (Kirov et al., 2020b). There are many other
329 examples of active plant TEs that can serve as potential targets for CANS-based capture of TEs. The
330 development of high-throughput tools that can capture TE insertions in a population represents a
331 milestone towards the acceleration of the plant breeding process and a more effective usage of natural
332 variation. Further increasing the coverage of sequencing using tools such as the PromethION sequencer
333 could greatly increase the sensitivity of this approach and might even allow the detection of TEs that have
334 occurred in a single plant within an entire population. Such results could then be used to design specific
335 primers based on the insertion site in order to identify all plants carrying the target TEs. The capture of
336 TEs occurring in individual plants in a population will also depend on the genome size of the species and
337 the number of sgRNAs and target loci in the genome. These parameters must be optimized for individual
338 species.

339 Estimating the integration preferences for distinct TEs is not a straightforward process, as post-integration
340 selection pressure has a strong effect on the current TEI landscape. Therefore, artificially inducing TE
341 activity (e.g. via stress) and exploiting mutant lines followed by TEI identification are the most direct ways
342 to study TE integration biases and targets. CANS can be a suitable tool for identifying TEs in such
343 experiments. Here, we showed that EVD has preferences for integrations in the pericentromeric regions
344 of *A. thaliana* chromosomes. It is important to note that the established EVD insertion distribution was
345 shaped by very recent events and that post-integration selection pressure had only a minor effect on this
346 process. While wild-type *A. thaliana* contains only a few EVD insertions, its close relative *Arabidopsis*
347 *lyrata* contains multiple EVD insertions, most of them in centromeric satellite regions (Tsukahara et al.,
348 2009). Thus, the centromere-biased distribution of EVD insertions is similar under selection pressure (*A.*
349 *lyrata*) and without post-integration selection (*A. thaliana ddm1* mutant). This can be explained by the
350 low gene density in the centromeric regions, which decreases the rate of gene insertions and relaxes the
351 selection pressure on EVD TEs.

352 Thus our pioneering work thus demonstrates that CANS is a valuable, cost-effective, rapid tool for
353 deciphering the mobilome landscapes of plant genomes. CANS can be utilized for TEI identification in
354 plants with small as well as large genomes and that TEs with few or multiple insertions can be detected
355 using this method and the NanoCasTE pipeline.

356

357 [Methods](#)

358 [Plant material and growth conditions](#)

359 Seeds of *ddm1* mutants (*ddm1-2*, F7 generation) were kindly provided by Vincent Colot (Institut de
360 Biologie de l'Ecole Normale Supérieure (IBENS), Paris, France). *Arabidopsis thaliana* Col-0 plants (wild type
361 and *ddm1* mutants) were grown in a light chamber for a month under 22°C and long-day conditions (16h
362 light/8h dark). For sequencing of pooled plants, *Arabidopsis* seeds were surface-sterilized with 75%
363 ethanol for 2 min, then washed with 5% sodium hypochlorite for 5 min. After this, the seeds were rinsed
364 with sterile distilled water 3 times. 0,1% agarose was added to each tube and the seeds were resuspended
365 and dripped on Petri dishes with MS nutrient medium (Murashige and Skoog, 1962), supplemented with
366 3% of sucrose (PanReac AppliChem, Germany) and 1% of agar (PanReac AppliChem, Germany). Plates with
367 seeds were sealed with Parafilm (Pechiney Plastic Packaging Company, USA) and kept in the dark at 4°C
368 for 3 days for vernalization and synchronous germination. Afterwards, dishes were transferred into a light
369 chamber with 16h day/8h night photoperiods for further growth.

370 Two *Aegilops tauschii* (K-1216 and K-112) accessions were kindly provided by Dr. E.D. Badaeva (Laboratory
371 of Genetic identification of plants, Vavilov Institute of General Genetics, Russian Academy of Sciences,
372 Moscow). The plants were grown in a greenhouse under natural conditions.

373

374 Plant transformation

375 To create G-ddm1 plants the nucleotide sequence of the GAG open reading frame of EVD5
376 retrotransposon (AtCOPIA93 (Mirouze et al., 2009)) was synthesized by the Synbio Technologies service
377 (USA). The GAG fragment was cloned at the restriction sites BamH1 and Sac1 into the pBI121 vector. The
378 *Agrobacterium* transformation of *A. thaliana* plants was carried out as described previously (Clough and
379 Bent, 1998). Primary transformants were selected based on their resistance to kanamycin (50 lg/ml). At
380 least 60 transgenic seeds of the T2 generation were selected and used for segregation analysis. Three
381 plants from one T2 family (G-ddm1) with 3:1 segregation of the marker for kanamycin resistance were
382 then transferred to pots with soil, grown three weeks more and used for DNA isolation for whole-genome
383 and CANS sequencing.

384

385 HMW DNA isolation and size selection

386 High molecular weight DNA was isolated from 200-500 mg of fresh and young leaves that were
387 homogenized in liquid nitrogen. DNA isolation was carried out according to the previously published
388 protocol (<https://www.protocols.io/view/plant-dna-extraction-and-preparation-for-ont-seque-bcvyiw7w>) with several modifications. Chloroform was used instead of dichloromethane
389 during the extraction step. Samples with added CTAB2 buffer were incubated at 65°C until the formation
390 of visible flakes. Dissolving of the precipitate in 1M NaCl was performed at 65°C for 15 minutes with
391 further cooling to the RT before adding isopropanol. Removal of contaminants residues and small-size
392 nucleotide fragments was performed by adding 2.5 volumes of 80% ethanol with further centrifugation
393 at 12000 g for 15 minutes. The obtained DNA pellet was washed with 70% ethanol and centrifuged at
394 12000 g for 5 minutes. Final DNA elution was carried out with 65 µL of nuclease-free water.

396 For whole-genome Nanopore sequencing and CANS, size-selection of the large DNA fragments was done
397 using SRE or XL Short Read Eliminator Kits (Circulomics, USA) according to the manufacture instructions.
398 The concentration and quality of isolated DNA were estimated by gel electrophoresis in 1% agarose gel,
399 NanoDrop One UV-Vis Spectrophotometer (Thermo Scientific, USA) and Quantus Fluorometer (Promega,
400 USA) using a DNA QuantiFluor ONE dsDNA System (Promega, USA). Only pure DNA (A260/A280 ~1.8 and
401 A260/A230 ~2.0 according to NanoDrop) with almost no differences in concentrations obtained by
402 Nanodrop and Quantus was used for sequencing.

403

404 gRNA design

405 To design gRNAs recognizing multiple copies of LTR retrotransposons we used a step-by-step strategy: (1)
406 using FlashFry 'discover' tool with the '-positionOutput -maxMismatch 0' arguments and reference
407 genome sequences (At 4.0 for *Ae. tauschii* and TAIR10 for *A.thaliana*) a pool of gRNAs and the
408 corresponding genomic sites were found followed by bed file generation; (2) reference sequences of
409 target LTR retrotransposons were blasted against the corresponding genomes followed by hit filtering
410 (minimum hit length 100bp and minimum e-value 1e-50) and bed file consisting blast hit positions was
411 generated; (3) the genomic location of predicted gRNA sites (step 1) were intersected with blast hit bed
412 file (step 2) using bedtools 'intersect' and OnOutTE ratio was calculated for each gRNAs:

413 OnOutTE = number of gRNA sites located on the blast hit from the corresponding TE / total number of
414 gRNA sites

415 (4) We selected gRNAs with OnOutTE ratio > 0.85.

416 For each retrotransposon a set of gRNAs were designed including one (*A. thaliana*) or two (*Ae. tauschii*)
417 gRNAs negative strand recognizing both LTRs and three gRNAs (one on negative and two on positive
418 strands) located between two LTR sequences. Sequences of gRNAs are provided in Supplemental Data Set
419 2.

420 In vitro transcription of sgRNAs

421 SgRNAs for CANS were produced by in vitro transcription from double-stranded DNA templates
422 containing T7 promoter that were assembled from two oligos, specific (gRNA1-5 in Supplemental Data
423 Set 2) and universal (CRISPR_R). All oligonucleotides were ordered in Evrogen (Moscow, Russia). To
424 synthesize double-stranded DNA templates, the following reagents were mixed per each reaction: 2µl of
425 unique sgRNA oligo (1 µM); 2µl of CRISPR_R (1 µM); 2µl of T7 forward (5'-
426 GGATCCTAATACGACTCACTATAG-3') and reverse (5'-AAAAAAGCACCGACTCGG-3') primers (100 µM of
427 each); 2µl of 50x dNTP (10 µM) (Evrogen, Moscow, Russia); 10µl of 10x Encyclo buffer (Evrogen, Moscow,
428 Russia); 1µl of Encyclo polymerase (Evrogen, Moscow, Russia); 79µl of RNase-free water. Oligos were
429 annealed and extended according to the PCR program: initial denaturation at 95°C (2 min); 30 cycles of
430 denaturation at 98°C (30 sec), annealing at 60°C (30 sec), elongation at 72°C (30 sec); final elongation at
431 72°C (1 min). Gel electrophoresis was performed after amplification and templates were column purified.

432 To transcribe the templates, the Highly Efficient RNA in vitro Synthesis Kit (Biolabmix, Novosibirsk,
433 Russia) was used. The subsequent reactions were set up and incubated at 37°C for 2h: 10µl of 5x T7
434 transcription buffer; 2µl of 25x DTT (250mM); 2µl of dNTP (25mM of each); 2µl of double-stranded DNA
435 template (150-200 ng); 1µl of T7 polymerase (150 units); 33µl RNase-free water. The sgRNAs were then
436 purified using RNA and miRNA Extraction Kit (Biolabmix, Novosibirsk, Russia) according to the
437 manufacturer instructions. The concentration and quality of prepared sgRNAs were estimated by
438 Nanodrop (Thermo Scientific, USA), Qubit (Thermo Scientific, USA) and gel electrophoresis in 2% agarose
439 gel.

440

441 Cas9/sgRNA ribonucleoprotein complexes (RNPs) assembly

442 For CANS, 50 ng of each sgRNA were used for RNP assembly. To obtain RNPs for sgRNA Mixture 2
443 the corresponding sgRNAs were pooled together at 1:1 molar ratio in 11 µl of MQ water. Before RNP
444 assembly all sgRNA mixtures were denaturated at 95°C for 5 min, then cooled on ice. The RNPs were
445 assembled by combining 11µl of sgRNA mix with 1µl of Cas9 nuclease in 3 µl of reaction buffer (Biolabmix,
446 Novosibirsk, Russia) at a final volume of 15µl followed by incubation for 30 min at room temperature and
447 kept on ice until usage.

448

449 Cas9 cleavage and library preparation

450 Before CANS DNA was dephosphorylated. For this, 2-8 µg of HMW DNA was diluted in 1x CutSmart Buffer
451 and 6 µl of Quick CIP enzyme (New England Biolabs, catalog no. M0508) was added followed by incubation
452 of the reaction at 37°C for 30 minutes. The reaction was stopped by heating (80°C, 2 minutes). DNA
453 cleavage by RNPs was carried out by mixing of 40 µl of dephosphorylated DNA, 15 µl of RNPs, 1.5 µl of
454 dATP (10 mM) and 1 µl Encyclo polymerase (Evrogen, Moscow, Russia). The mixture was incubated at
455 37°C for 30 minutes followed by 5 min at 72 °C. Nanopore sequencing adapters (AMX) were ligated to the
456 Cas9 cleaved DNA by mixing the following components: 25 µl of LNB buffer (Oxford Nanopore
457 Technologies, catalog no. SQK-LSK109), 5 µl of nuclease-free water, 12.5 µl of Quick T4 DNA Ligase (New
458 England Biolabs, NEBNext Companion Module for Oxford Nanopore Technologies Ligation Sequencing
459 catalog no. E7180S) and Adapter Mix (SQK-LSK109). The mixture was added to the Cas9 cleaved DNA and
460 incubated for 20 min at room temperature. The samples were purified by adding equal volume of TE
461 buffer (pH 8.0) and 0.3 volume of AMPure XP Beads (Beckman Coulter, catalog no. A63881) and washed
462 twice by LFB buffer (Oxford Nanopore Technologies, catalog no. SQK-LSK109). The DNA was eluted in 15 µl
463 of elution buffer (Oxford Nanopore Technologies, catalog no. SQK-LSK109). Then sequencing library was
464 prepared according to the Ligation Sequencing Kit SQK-LSK109 protocol (Oxford Nanopore Technologies,
465 catalog no. SQK-LSK109).

466 Library preparation for whole-genome sequencing of G-ddm1-3 was carried out using the Ligation
467 Sequencing Kit (Oxford Nanopore Technologies, catalog no. SQK-LSK109) with 0.5 - 1 µg of input HMW
468 DNA. The library for G-ddm1-1 and G-ddm1-2 WGS was prepared from 1 µg of both DNA using the Native
469 Barcoding Expansion 1-12 (Oxford Nanopore Technologies, catalog no. EXP-NBD104) and the Ligation
470 Sequencing Kit SQK-LSK109 (Oxford Nanopore Technologies).

471 Sequencing and basecalling

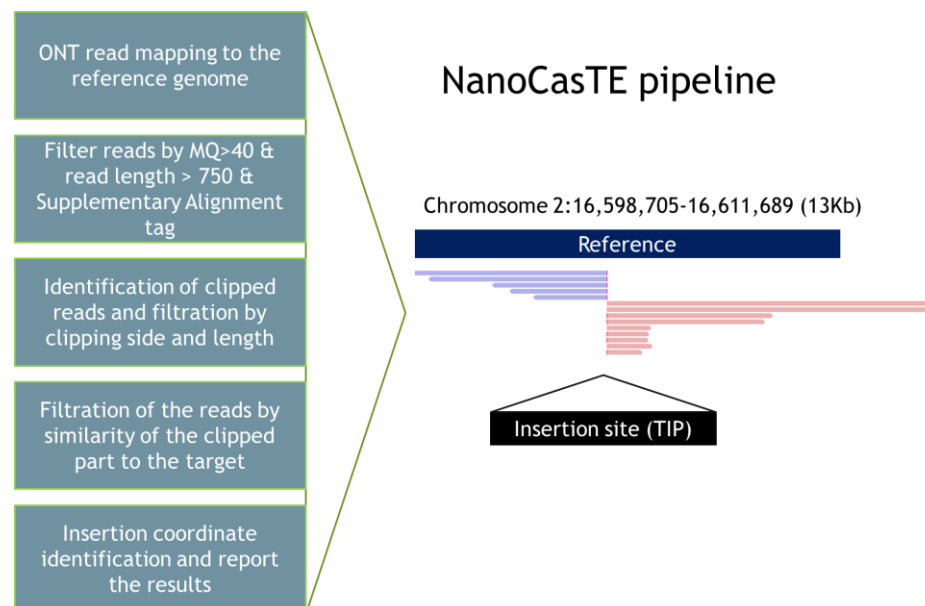
472 Sequencing was performed by MinION equipped by R9.4.1 or R10.3 flow cells. Sequencing process was
473 operated by MinKNOW software (v.19.12.5). The detailed information on Nanopore sequencing is
474 described in Supplemental Data Set 3.

475 Analysis

476 Base calling was done by Guppy (Version 3.2.10). The Aet_v4.0 genome sequence was downloaded from
477 the Gramene database (ftp://ftp.gramene.org/pub/gramene/release-63/fasta/aegilops_tauschii). The
478 TAIR10 genome was downloaded from NCBI database. The generated Nanopore reads were aligned to
479 the reference genomes (TAIR10 for *A. thaliana* and Aet_v4.0 for *Ae. tauschii*) using minimap2 software (Li,
480 2018) with the following parameters: -ax map-ont -t 100. The obtained sam file was converted to bam
481 format, sorted and indexed using SAMtools (Li et al., 2009). The obtained sorted bam files were used for
482 NanoCasTE (CANS data) and NanoWgsTE (WGS data) pipelines.

483 NanoCasTE pipeline for the identification of transposon insertions following CANS

484 NanoCasTE is written in python and can be run as a standalone tool or via Snakefile
485 (<https://github.com/Kirovez/NanoCasTE>). NanoCasTE involves several steps (Figure 7): (1) Mapping the
486 reads to the genome using minimap2 (Li, 2018), followed by the generation of sorted bam files using
487 SAMtools (Li et al., 2009); (2) Filtering the reads based on mapping quality (MQ > 40), read length (>750)
488 and SA tags (no supplementary alignments are allowed); (3) Selecting mapped reads with clipped heads
489 (reads mapped to the positive strand) or tails (reads mapped to the negative strand) with the length of
490 the clipped parts close to the distance from the gRNA positions to the TE end; (4) Filtering the reads based
491 on similarity of the clipped part to the target TE sequence; (5) Identifying the inserted sites and TE
492 orientation (+ or - strand) and outputting the results.



493

494 **Figure 7. Schematic view of the major steps of the NanoCasTE pipeline to identify transposon insertions**
495 **in the genome after CANS**

496

497 NanoCasTE uses a set of stringent criteria to specifically detect TE insertions in CANS data and to
498 distinguish them from noise signals. The pipeline reports the putative positions of the target TE insertion
499 as well as additional information that is useful for further analysis, including the number of selected reads
500 supporting the insertion, the total number of reads covering the insertion, the strand harboring the TE
501 insertion and the length of the clipped parts of the selected reads.

502 NanoCasTE was run for all samples using the following parameters: min_len_clipped: 0.6,
503 mapping_quality: 40, min_read_length: 750. The following lists of gRNAs were used: (1) for *A. thaliana*,
504 'TCTTGGTGATGAGAGTGAC, ACCCTGGATTTAAGGTGAGA, AGTTTAAGAGCTCTAGTATG,
505 CTACAAGGTCAATCGAAAGG, TCAACACATGAAAGTCCCGA'; (2) for *Ae. tauschii*,
506 'CCGGGTCGTCCCTTTCTATA, GTCCCTTTCTATAGGGAGGT, CTGAGCCGTTTCGATGAGAC,
507 TCCGGAGAATGACGTCACTC, TAAGCCGGAGATTTTTCTGT'. We used stricter filtration criteria for *Ae.*
508 *tauschii* insertion selection and kept only the insertions localized on chromosomes and possessing two or
509 more supporting reads.

510

511 NanoWgsTE pipeline

512 To detect TE or T-DNA insertions in the genome using Nanopore whole-genome sequencing data we
513 designed a pipeline called NanoWgsTE (<https://github.com/Kirovez/NanoWgsTE>). NanoWgsTE involves
514 several steps: (1) it converts fastq to fasta and performs BLAST search of the reads with the similarity to
515 target sequence followed by selection of the reads into a separate file; (2) the selected reads are mapped
516 to the genome by minimap2 (Li, 2018) followed by generation sorted bam file using SAMtools (Li et al.,
517 2009) and BAMtools (Barnett et al., 2011); (3) the reads with clipped ends are searched in the obtained
518 bam file assisted by pysam package (<https://github.com/pysam-developers/pysam>); (4) the genomic
519 positions where clipped ends of the reads were mapped are collected and output files (bed file with
520 putative insertion positions) and file with clipped read count data per putative insertion) are generated.
521 We run this pipeline with default parameters for EVD5 (AT5TE20395) and T-DNA (pBI121 sequence, NCBI
522 accession number: AF485783.1): -q 40 -mlc 500 -mbh 500. For detection of EVD5 insertions the fasta
523 sequence was used. All putative T-DNA and EVD5 insertions were manually checked by a local instance of
524 JBrowse (Buels et al., 2016).

525 Insertion validation by PCR

526 For validations of TEIs we used the combination of primers including one primer located on TEI flanking
527 regions and one primer located on TE. The sequences of primers are listed in Supplemental Data S2. The
528 PCR was performed in a reaction volume of 25 µl using Bio-Rad T100 Thermal Cycler (Bio-Rad, USA). The
529 reaction mixture contained 2.5 mM MgCl₂, 200 mM dNTPs (Dia-M, Moscow, Russia), 5 pmols of each
530 primer, 0.5 units of ColoredTaq polymerase (Sileks, Russia). The PCR reactions were carried out under the
531 following conditions: Touchdown-PCR was carried out according to the following cycling program: 94 °C
532 for 3 min, 94 °C for 30 s, 65 °C for 30 s and 72 °C for 2 min, followed by 6 cycles at decreasing annealing
533 temperatures in decrements of 1 °C per cycle, then 35 cycles of 30 s at 94 °C, 30 s at 55 °C, 2 min at 72 °C
534 and final extension at 72 °C for 5 min. For validation of *A.thaliana* TEIs the following PCR program was
535 used: 35 cycles of 94 °C for 30 s, 60 °C for 30 s and 72 °C for 2 min. Amplification products were visualized
536 by electrophoresis in 1.5% agarose gel and observed under UV light after ethidium bromide staining.
537 Primers used for TEI validation are listed in Supplemental Data Set 4.

538

539 *Aegilops tauschii* LTR retrotransposon identification and selection for CANS

540 Genomic sequence of *A. tauschii* (Aet_v4.0) was downloaded from the Gramene database
541 (ftp://ftp.gramene.org/pub/gramene/release-63/fast/aegilops_tauschii). LTR retrotransposons of
542 chromosome 1 of *A. tauschii* were predicted in the genome as previously described (Penin et al. 2021)
543 using LTRharvest (Ellinghaus et al., 2008) and LTRdigest (Steinbiss et al., 2009). The following arguments
544 were used for LTRdigest: -aaout yes -pptlen 10, 30 -pbsoffset 0, 3 -pdomevalcutoff 0.001. Hmm profiles
545 of transposon domains were downloaded from the GyDB database (Llorens et al., 2010). TESorter
546 software (Zhang et al., 2019) was used for LTR retrotransposon classification. The LTR sequences of all
547 predicted TEs were clustered using cd-hit with the following arguments: -T 100 -l 150 -d 100 -s 0.8 -aL 0.8
548 -aS 0.8 -A 80 -sc -sf -0.95 . Then we selected a single cluster containing complete TE (according to the
549 TESorter results) with identical LTRs pointing to its recent insertion time. This Ty3/Gypsy LTR
550 retrotransposon (Aty3_169) of the Retand family is located on chromosome 1D:169126894..169140327

551 of the reference Aet_v4.0 genome. Blast search of the similarity of Aty3_169 LTR sequence to all Aet_v4
552 genomic sequences was carried out on the Ensembl website and resulted in >300 hits of
553 Aty3_169 distributed along all chromosomes.

554

555 [Statistics and Data Visualization](#)

556 Statistical analysis was carried out in Rstudio Version 1.2.1335 (<http://www.rstudio.com/>) with R version
557 3.6.0. Visualization was carried out by ggplot2 (Wickham, 2011) and VennDiagram (Chen and Boutros,
558 2011) packages.

559

560 [Accession Numbers](#)

561 All Nanopore data generated in this study are deposited in the National Center for Biotechnology
562 Information SRA database (BioProject ID PRJNA736208).

563

564 [Acknowledgements](#)

565 We thank Ekaterina Badaeva (Laboratory of Genetic identification of plants, Vavilov Institute of General
566 Genetics, Russian Academy of Sciences, Moscow) and Vincent Colot (Institut de Biologie de l'Ecole
567 Normale Supérieure (IBENS), Paris, France) for providing the plant material.

568

569 [Author contribution statement](#)

570 I.K. and M.D. designed the research; I.K., P.M., S.G., R.K., M.O., M.D. and A.K. performed the
571 research. I.K., G.K., A.S. and M.D. analyzed data. I.K. wrote the article.

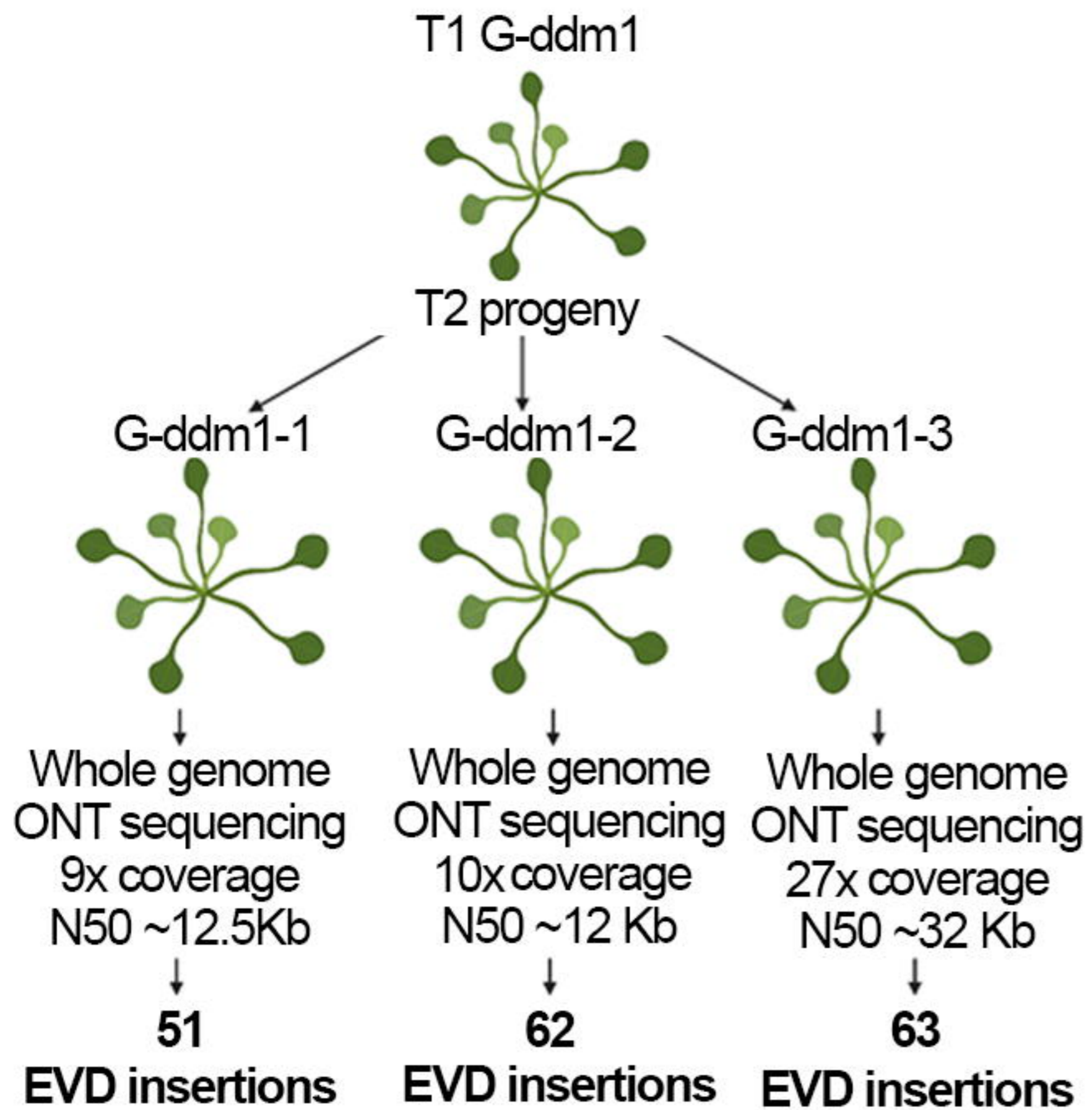
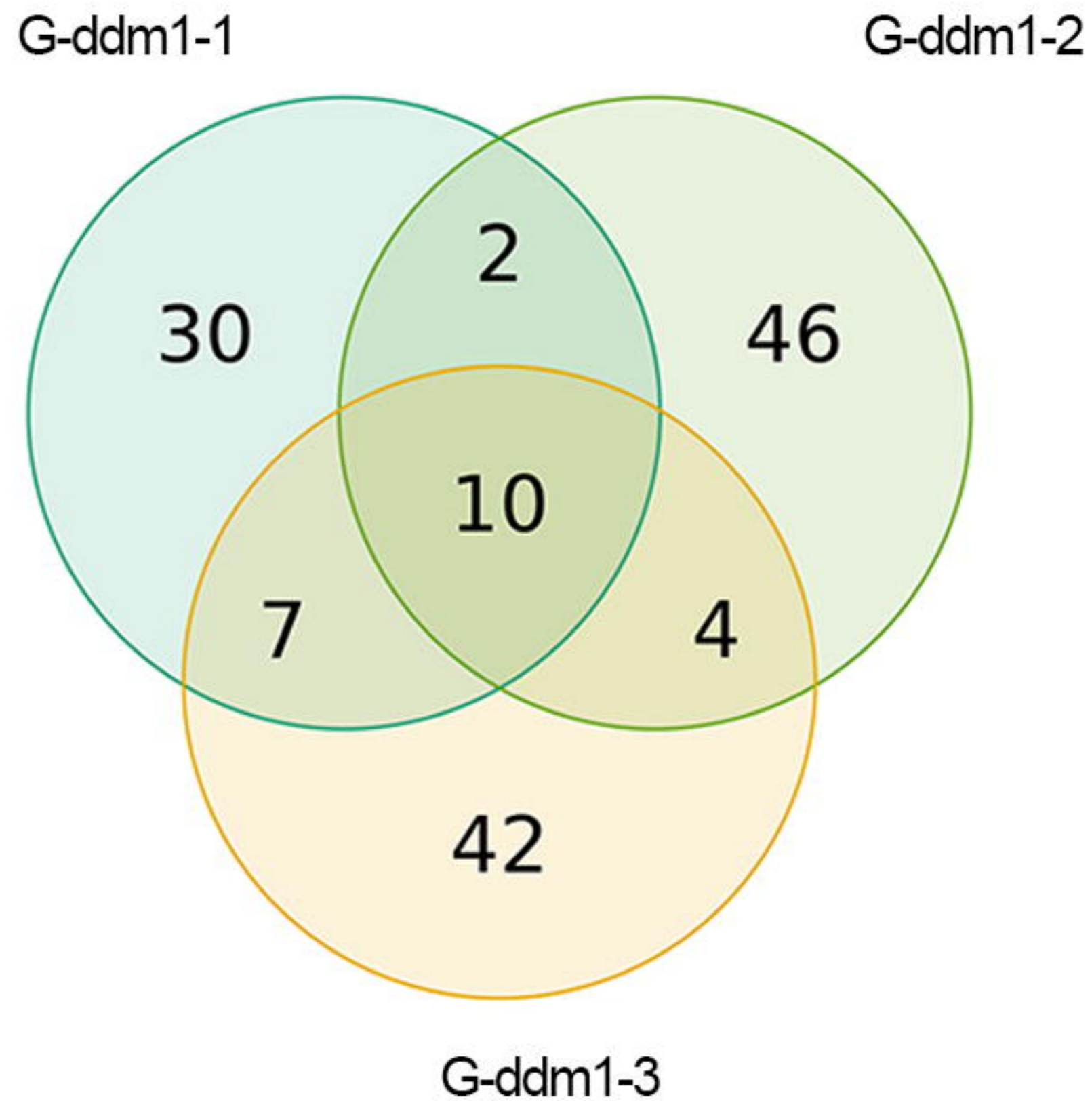
572

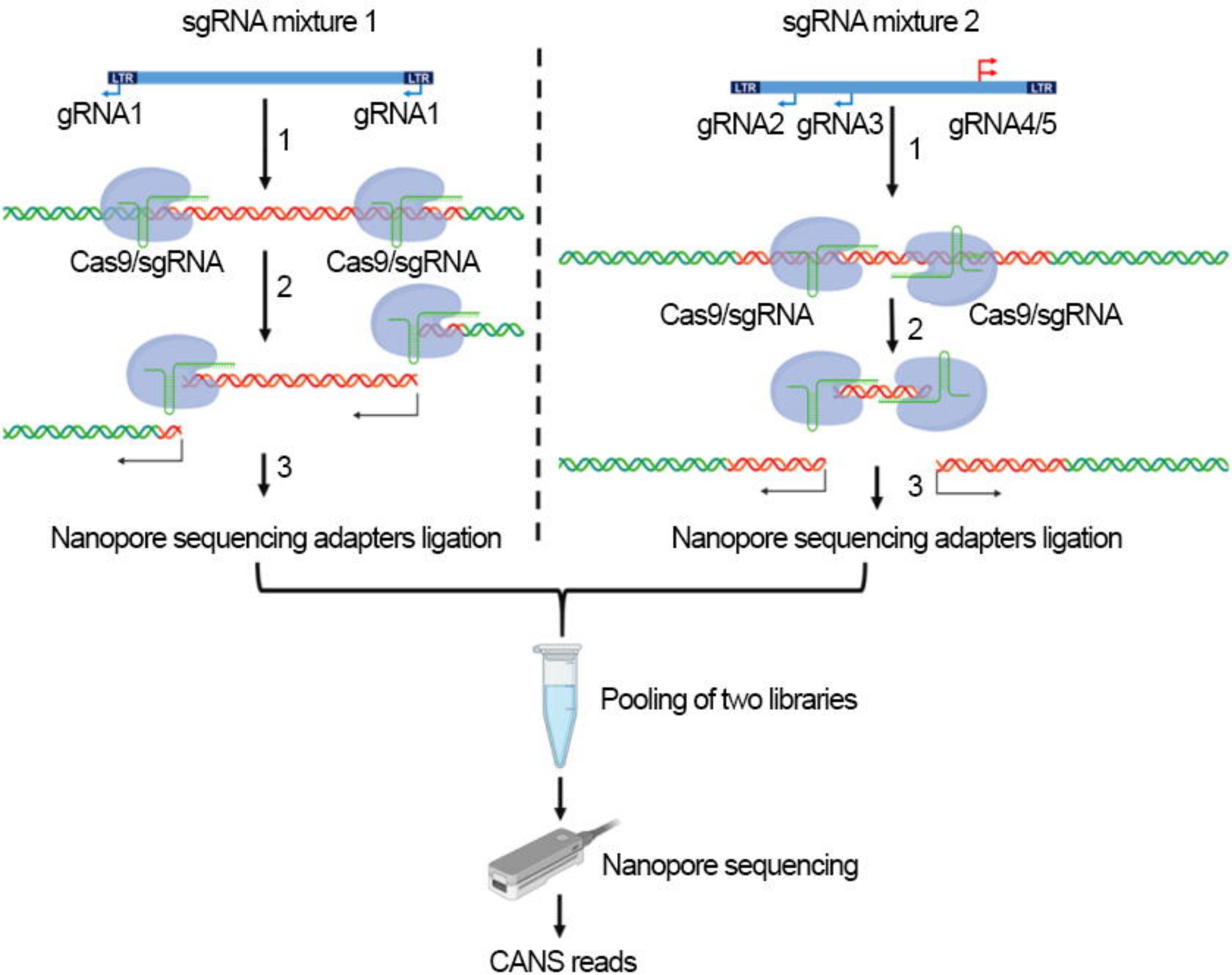
573 [References](#)

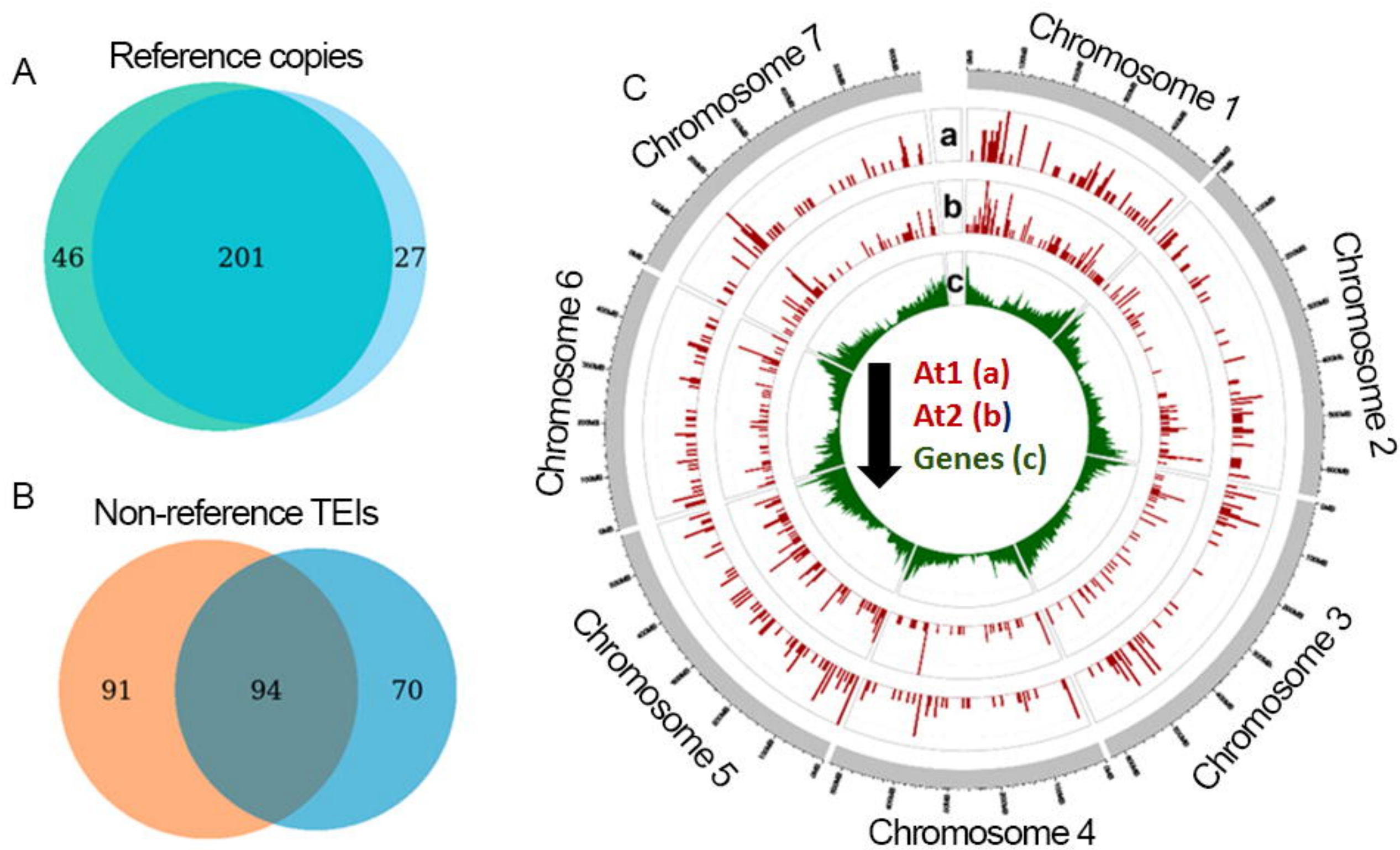
- 574 **Baduel, P., Leduque, B., Ignace, A., Gy, I., Gil, J., Loudet, O., Colot, V., and Quadrana, L.** (2021). Genetic
575 and environmental modulation of transposition shapes the evolutionary potential of *Arabidopsis*
576 *thaliana*. *Genome Biology* **22**, 138.
- 577 **Barnett, D.W., Garrison, E.K., Quinlan, A.R., Strömberg, M.P., and Marth, G.T.** (2011). BamTools: a C++
578 API and toolkit for analyzing and managing BAM files. *Bioinformatics* **27**, 1691-1692.
- 579 **Buels, R., Yao, E., Diesh, C.M., Hayes, R.D., Munoz-Torres, M., Helt, G., Goodstein, D.M., Elsik, C.G.,
580 Lewis, S.E., Stein, L., and Holmes, I.H.** (2016). JBrowse: a dynamic web platform for genome
581 visualization and analysis. *Genome Biology* **17**, 66.
- 582 **Butelli, E., Licciardello, C., Zhang, Y., Liu, J., Mackay, S., Bailey, P., Reforgiato-Recupero, G., and
583 Martin, C.** (2012). Retrotransposons control fruit-specific, cold-dependent accumulation of
584 anthocyanins in blood oranges. *Plant Cell* **24**, 1242-1255.
- 585 **Carpentier, M.-C., Manfroi, E., Wei, F.-J., Wu, H.-P., Lasserre, E., Llauro, C., Debladis, E., Akakpo, R.,
586 Hsing, Y.-I., and Panaud, O.** (2019). Retrotranspositional landscape of Asian rice revealed by
587 3000 genomes. *Nature Communications* **10**, 24.
- 588 **Chen, H., and Boutros, P.C.** (2011). VennDiagram: a package for the generation of highly-customizable
589 Venn and Euler diagrams in R. *BMC Bioinformatics* **12**, 35.
- 590 **Chen, J., Lu, L., Robb, S.M.C., Collin, M., Okumoto, Y., Stajich, J.E., and Wessler, S.R.** (2020). Genomic
591 diversity generated by a transposable element burst in a rice recombinant inbred population.
592 *Proceedings of the National Academy of Sciences of the United States of America* **117**, 26288-
593 26297.
- 594 **Chuong, E.B., Elde, N.C., and Feschotte, C.** (2017). Regulatory activities of transposable elements: from
595 conflicts to benefits. *Nature Reviews Genetics* **18**, 71-86.

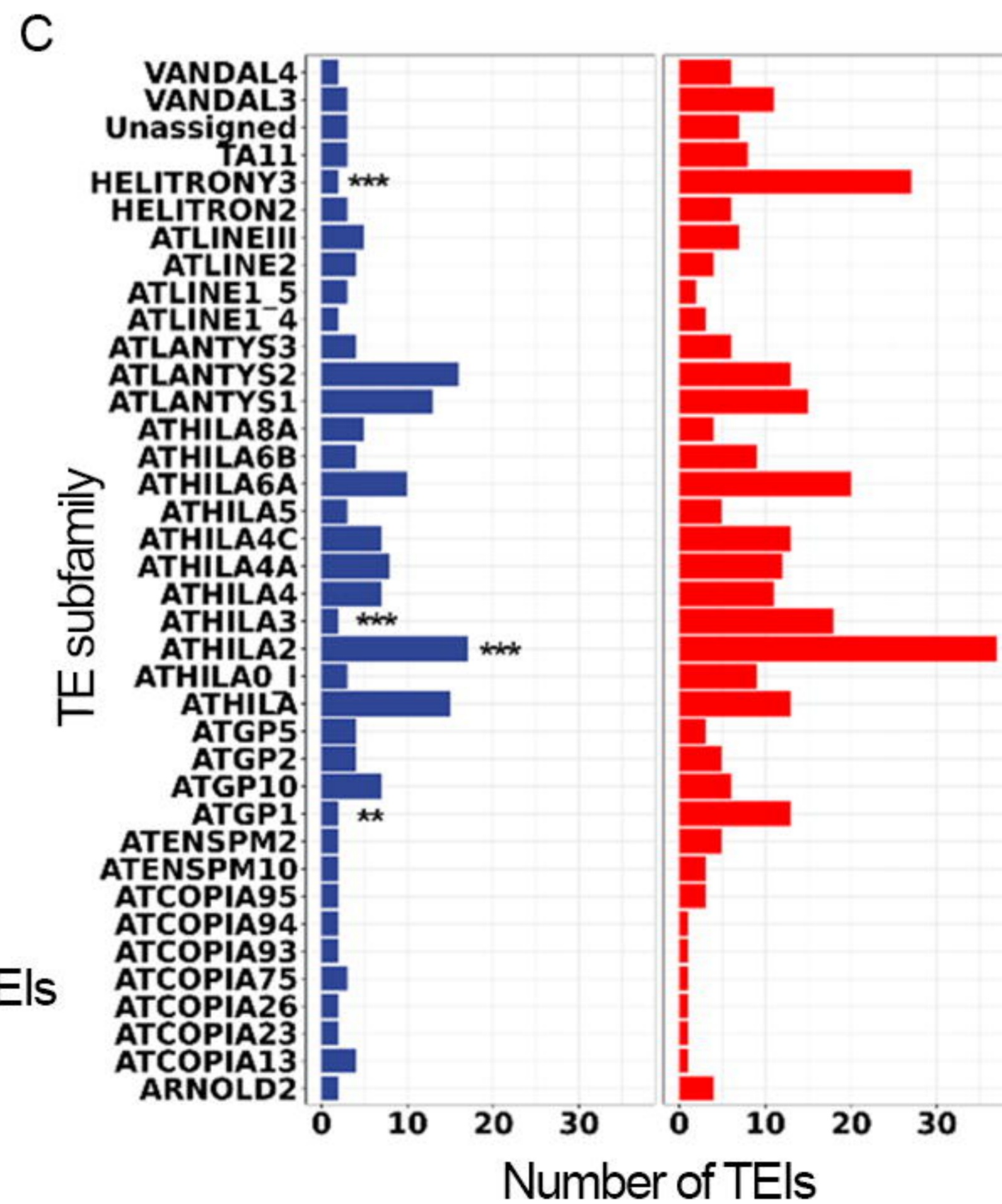
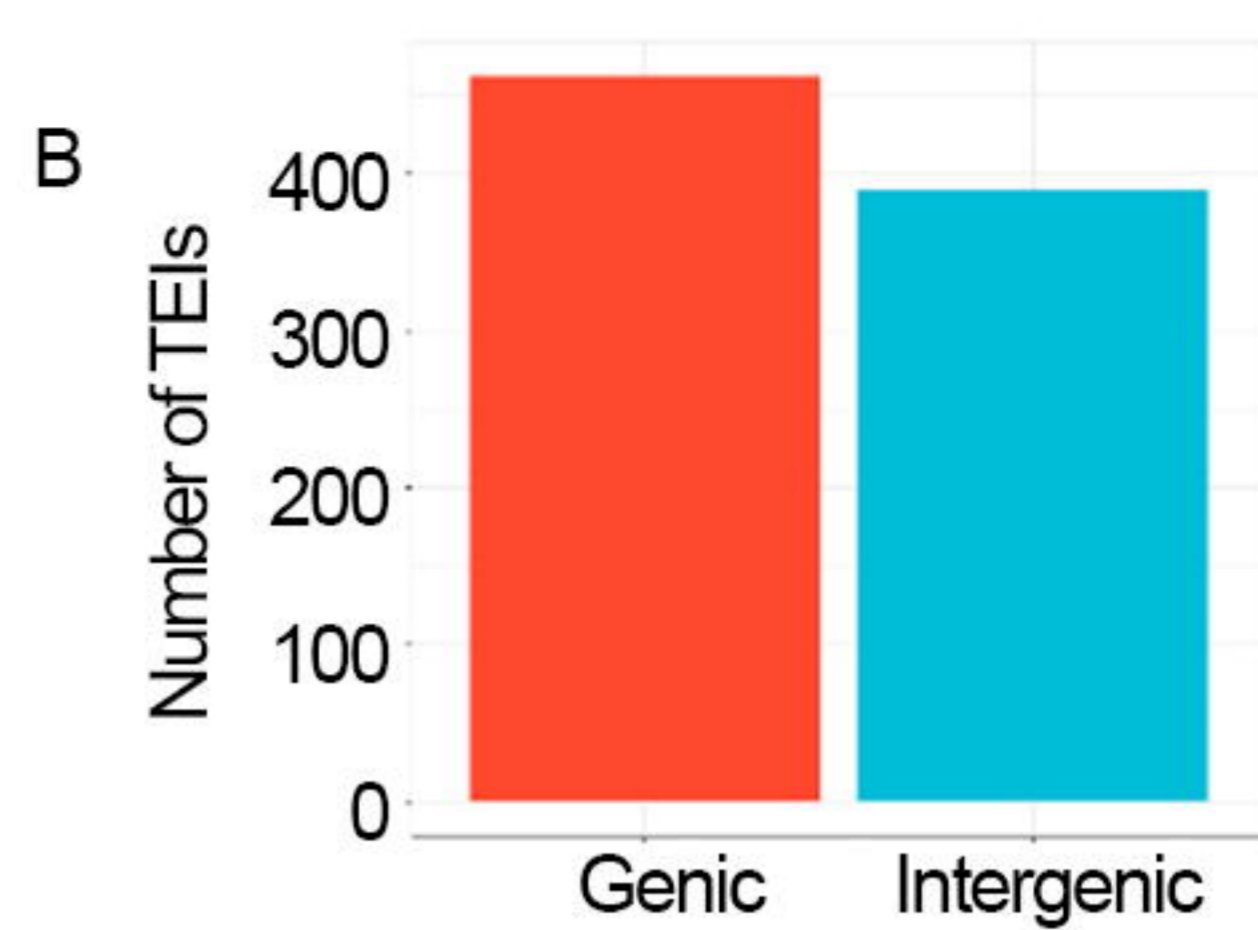
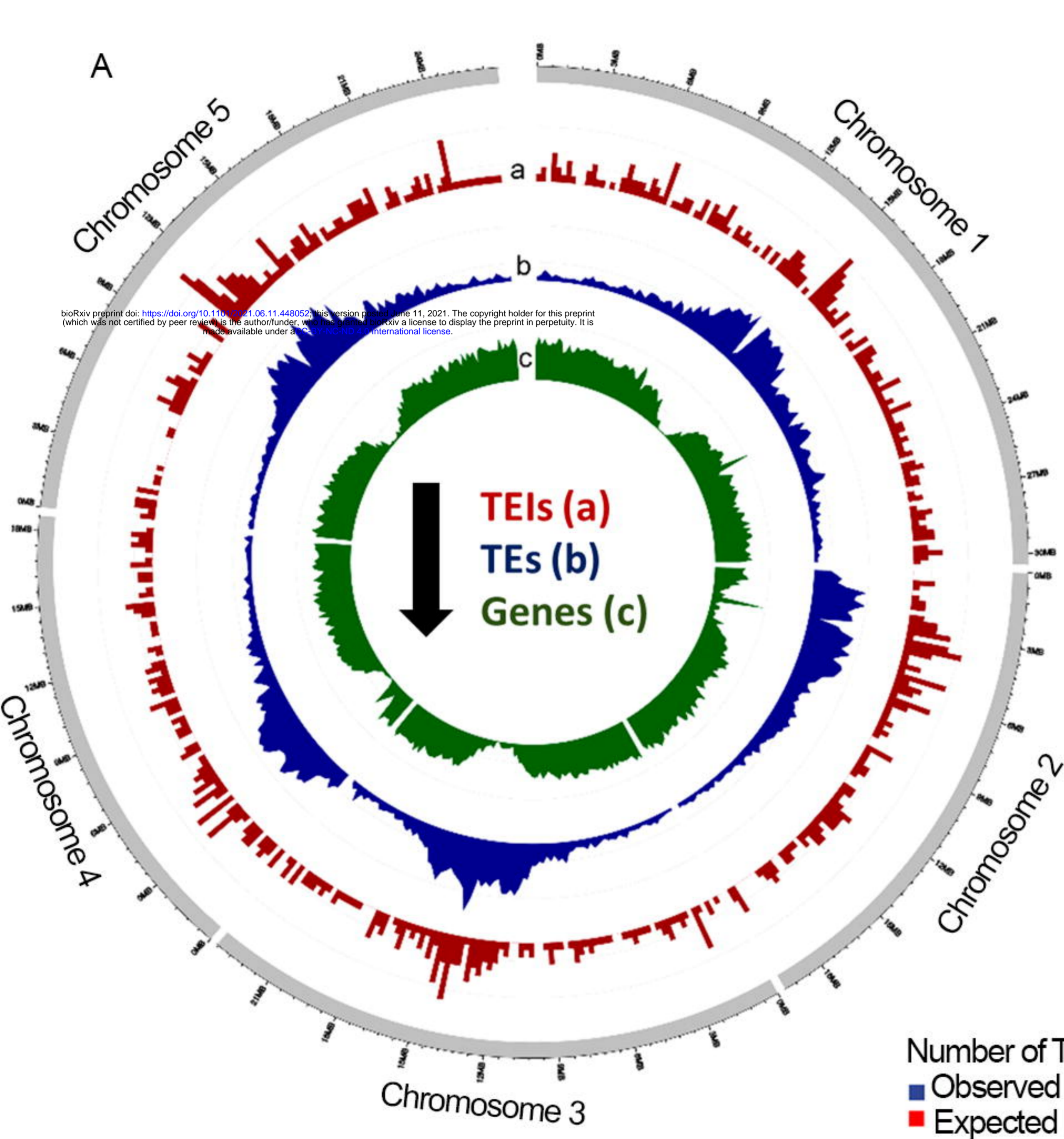
- 596 **Clough, S.J., and Bent, A.F.** (1998). Floral dip: a simplified method for *Agrobacterium*-mediated
597 transformation of *Arabidopsis thaliana*. *Plant J* **16**, 735-743.
- 598 **Dmitriev, A.A., Pushkova, E.N., Novakovskiy, R.O., Beniaminov, A.D., Rozhmina, T.A., Zhuchenko, A.A.,**
599 **Bolsheva, N.L., Muravenko, O.V., Povkhova, L.V., Dvorianinova, E.M., Kezimana, P., Snezhkina,**
600 **A.V., Kudryavtseva, A.V., Krasnov, G.S., and Melnikova, N.V.** (2021). Genome Sequencing of
601 Fiber Flax Cultivar Atlant Using Oxford Nanopore and Illumina Platforms. *Front Genet* **11**,
602 590282-590282.
- 603 **Domínguez, M., Dugas, E., Benchouaia, M., Leduque, B., Jiménez-Gómez, J.M., Colot, V., and**
604 **Quadrona, L.** (2020). The impact of transposable elements on tomato diversity. *Nature*
605 *Communications* **11**, 4058.
- 606 **Ellinghaus, D., Kurtz, S., and Willhoeft, U.** (2008). LTRharvest, an efficient and flexible software for de
607 novo detection of LTR retrotransposons. *BMC Bioinformatics* **9**, 18.
- 608 **Gabriel, T., Sharim, H., Fridman, D., Arbib, N., Michaeli, Y., and Ebenstein, Y.** (2018). Selective
609 nanopore sequencing of human BRCA1 by Cas9-assisted targeting of chromosome segments
610 (CATCH). *Nucleic Acids Research* **46**, e87-e87.
- 611 **Gilpatrick, T., Lee, I., Graham, J.E., Raimondeau, E., Bowen, R., Heron, A., Downs, B., Sukumar, S.,**
612 **Sedlazeck, F.J., and Timp, W.** (2020). Targeted nanopore sequencing with Cas9-guided adapter
613 ligation. *Nature Biotechnology* **38**, 433-438.
- 614 **Handsaker, R.E., Korn, J.M., Nemes, J., and McCarroll, S.A.** (2011). Discovery and genotyping of
615 genome structural polymorphism by sequencing on a population scale. *Nat Genet* **43**, 269-276.
- 616 **Jia, J., Zhao, S., Kong, X., Li, Y., Zhao, G., He, W., Appels, R., Pfeifer, M., Tao, Y., Zhang, X., Jing, R.,**
617 **Zhang, C., Ma, Y., Gao, L., Gao, C., Spannagl, M., Mayer, K.F.X., Li, D., Pan, S., Zheng, F., Hu, Q.,**
618 **Xia, X., Li, J., Liang, Q., Chen, J., Wicker, T., Gou, C., Kuang, H., He, G., Luo, Y., Keller, B., Xia, Q.,**
619 **Lu, P., Wang, J., Zou, H., Zhang, R., Xu, J., Gao, J., Middleton, C., Quan, Z., Liu, G., Wang, J.,**
620 **Yang, H., Liu, X., He, Z., Mao, L., Wang, J., and International Wheat Genome Sequencing, C.**
621 (2013). *Aegilops tauschii* draft genome sequence reveals a gene repertoire for wheat
622 adaptation. *Nature* **496**, 91-95.
- 623 **Jiang, N., Gao, D., Xiao, H., and van der Knaap, E.** (2009). Genome organization of the tomato sun locus
624 and characterization of the unusual retrotransposon Rider. *Plant J* **60**, 181-193.
- 625 **Jung, H., Winefield, C., Bombarely, A., Prentis, P., and Waterhouse, P.** (2019). Tools and Strategies for
626 Long-Read Sequencing and De Novo Assembly of Plant Genomes. *Trends in Plant Science* **24**,
627 700-724.
- 628 **Kirov, I., Omarov, M., Merkulov, P., Dudnikov, M., Gvaramiya, S., Kolganova, E., Komakhin, R., Karlov,**
629 **G., and Soloviev, A.** (2020a). Genomic and Transcriptomic Survey Provides New Insight into the
630 Organization and Transposition Activity of Highly Expressed LTR Retrotransposons of Sunflower
631 (*Helianthus annuus* L.). *Int J Mol Sci* **21**, 9331.
- 632 **Kirov, I., Dudnikov, M., Merkulov, P., Shingaliev, A., Omarov, M., Kolganova, E., Sigaeva, A., Karlov, G.,**
633 **and Soloviev, A.** (2020b). Nanopore RNA Sequencing Revealed Long Non-Coding and LTR
634 Retrotransposon-Related RNAs Expressed at Early Stages of Triticale SEED Development. *Plants*
635 **9**, 1794.
- 636 **Kobayashi, S., Goto-Yamamoto, N., and Hirochika, H.** (2004). Retrotransposon-induced mutations in
637 grape skin color. *Science* **304**, 982.
- 638 **Lee, Y.G., Choi, S.C., Kang, Y., Kim, K.M., Kang, C.-S., and Kim, C.** (2019). Constructing a Reference
639 Genome in a Single Lab: The Possibility to Use Oxford Nanopore Technology. *Plants* **8**, 270.
- 640 **Li, H.** (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094-3100.
- 641 **Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R.,**
642 **and Subgroup, G.P.D.P.** (2009). The Sequence Alignment/Map format and SAMtools.
643 *Bioinformatics* **25**, 2078-2079.
- 644 **Li, S., Jia, S., Hou, L., Nguyen, H., Sato, S., Holding, D., Cahoon, E., Zhang, C., Clemente, T., and Yu, B.**
645 (2019). Mapping of transgenic alleles in soybean using a nanopore-based sequencing strategy.
646 *Journal of Experimental Botany* **70**, 3825-3833.
- 647 **Lisch, D.J.N.R.G.** (2013). How important are transposons for plant evolution? *Nature Reviews Genetics*
648 **14**, 49-61.
- 649 **Llorens, C., Futami, R., Covelli, L., Domínguez-Escribá, L., Viu, J.M., Tamarit, D., Aguilar-Rodríguez, J.,**
650 **Vicente-Ripolles, M., Fuster, G., Bernet, G.P., Maumus, F., Munoz-Pomer, A., Sempere, J.M.,**

- 651 **Latorre, A., and Moya, A.** (2010). The Gypsy Database (GyDB) of mobile genetic elements:
652 release 2.0. *Nucleic Acids Research* **39**, D70-D74.
- 653 **López-Girona, E., Davy, M.W., Albert, N.W., Hilario, E., Smart, M.E.M., Kirk, C., Thomson, S.J., and**
654 **Chagné, D.** (2020). CRISPR-Cas9 enrichment and long read sequencing for fine mapping in
655 plants. *Plant Methods* **16**, 121.
- 656 **Madsen, E.B., Höijer, I., Kvist, T., Ameer, A., and Mikkelsen, M.J.** (2020). Xdrop: Targeted sequencing of
657 long DNA molecules from low input samples using droplet sorting. *Hum Mutat* **41**, 1671-1679.
- 658 **McDonald, T.L., Zhou, W., Castro, C., Mumm, C., Switzenberg, J.A., Mills, R.E., and Boyle, A.P.** (2021).
659 Cas9 targeted enrichment of mobile elements using nanopore sequencing,
660 2021.2002.2010.430605.
- 661 **Mirouze, M., Reinders, J., Bucher, E., Nishimura, T., Schneeberger, K., Ossowski, S., Cao, J., Weigel, D.,**
662 **Paszowski, J., and Mathieu, O.** (2009). Selective epigenetic control of retrotransposition in
663 *Arabidopsis*. *Nature* **461**, 427-430.
- 664 **Ni, P., Huang, N., Nie, F., Zhang, J., Zhang, Z., Wu, B., Bai, L., Liu, W., Xiao, C.-L., Luo, F., and Wang, J.**
665 (2021). Genome-wide Detection of Cytosine Methylations in Plant from Nanopore sequencing
666 data using Deep Learning, 2021.2002.2007.430077.
- 667 **Nuthikattu, S., McCue, A.D., Panda, K., Fultz, D., DeFraia, C., Thomas, E.N., and Slotkin, R.K.** (2013).
668 The Initiation of Epigenetic Silencing of Active Transposable Elements Is Triggered by RDR6 and
669 21-22 Nucleotide Small Interfering RNAs **162**, 116-131.
- 670 **Panda, K., and Slotkin, R.K.** (2020). Long-Read cDNA Sequencing Enables a “Gene-Like” Transcript
671 Annotation of Transposable Elements **32**, 2687-2698.
- 672 **Quadrana, L., Bortolini Silveira, A., Mayhew, G.F., LeBlanc, C., Martienssen, R.A., Jeddloh, J.A., and**
673 **Colot, V.** (2016). The *Arabidopsis thaliana* mobilome and its impact at the species level. *eLife* **5**,
674 e15716.
- 675 **Rabanus-Wallace, M.T., Hackauf, B., Mascher, M., Lux, T., Wicker, T., Gundlach, H., Baez, M., Houben,**
676 **A., Mayer, K.F.X., Guo, L., Poland, J., Pozniak, C.J., Walkowiak, S., Melonek, J., Praz, C.R.,**
677 **Schreiber, M., Budak, H., Heuberger, M., Steuernagel, B., Wulff, B., Börner, A., Byrns, B.,**
678 **Čížková, J., Fowler, D.B., Fritz, A., Himmelbach, A., Kaithakottil, G., Keilwagen, J., Keller, B.,**
679 **Konkin, D., Larsen, J., Li, Q., Myśków, B., Padmarasu, S., Rawat, N., Sesiz, U., Biyiklioglu-Kaya,**
680 **S., Sharpe, A., Šimková, H., Small, I., Swarbreck, D., Toegelová, H., Tsvetkova, N., Voylovkov,**
681 **A.V., Vrána, J., Bauer, E., Bolibok-Bragoszewska, H., Doležel, J., Hall, A., Jia, J., Korzun, V.,**
682 **Laroche, A., Ma, X.-F., Ordon, F., Özkan, H., Rakoczy-Trojanowska, M., Scholz, U., Schulman,**
683 **A.H., Siekmann, D., Stojatowski, S., Tiwari, V.K., Spannagl, M., and Stein, N.** (2021).
684 Chromosome-scale genome assembly provides insights into rye biology, evolution and
685 agronomic potential. *Nature Genetics* **53**, 564-573.
- 686 **Schalamun, M., Nagar, R., Kainer, D., Beavan, E., Eccles, D., Rathjen, J.P., Lanfear, R., and**
687 **Schwessinger, B.** (2019). Harnessing the MinION: An example of how to establish long-read
688 sequencing in a laboratory using challenging plant tissue from *Eucalyptus pauciflora*. *Mol Ecol*
689 *Resour* **19**, 77-89.
- 690 **Slotkin, R.K., and Martienssen, R.** (2007). Transposable elements and the epigenetic regulation of the
691 genome. *Nature Reviews Genetics* **8**, 272-285.
- 692 **Stangl, C., de Blank, S., Renkens, I., Westera, L., Verbeek, T., Valle-Inclan, J.E., González, R.C., Hensen,**
693 **A.G., van Roosmalen, M.J., Stam, R.W., Voest, E.E., Kloosterman, W.P., van Haften, G., and**
694 **Monroe, G.R.** (2020). Partner independent fusion gene detection by multiplexed CRISPR-Cas9
695 enrichment and long read nanopore sequencing. *Nature Communications* **11**, 2861.
- 696 **Steinbiss, S., Willhoeft, U., Gremme, G., and Kurtz, S.** (2009). Fine-grained annotation and classification
697 of de novo predicted LTR retrotransposons. *Nucleic Acids Research* **37**, 7002-7013.
- 698 **Sultana, T., Zamborlini, A., Cristofari, G., and Lesage, P.J.N.R.G.** (2017). Integration site selection by
699 retroviruses and transposable elements in eukaryotes. *Nature Reviews Genetics* **18**, 292-308.
- 700 **Tsukahara, S., Kobayashi, A., Kawabe, A., Mathieu, O., Miura, A., and Kakutani, T.** (2009). Bursts of
701 retrotransposition reproduced in *Arabidopsis*. *Nature* **461**, 423-426.
- 702 **Wickham, H.** (2011). ggplot2. *Wiley Interdisciplinary Reviews: Computational Statistics* **3.2** **3**, 180-185.
- 703 **Williams-Carrier, R., Stiffler, N., Belcher, S., Kroeger, T., Stern, D.B., Monde, R.-A., Coalter, R., and**
704 **Barkan, A.** (2010). Use of Illumina sequencing to identify transposon insertions underlying
705 mutant phenotypes in high-copy Mutator lines of maize. *Plant J* **63**, 167-177.

A**B**







ONT read mapping to the reference genome

Filter reads by MQ>40 & read length > 750 & Supplementary Alignment tag

Identification of clipped reads and filtration by clipping side and length

Filtration of the reads by similarity of the clipped part to the target

Insertion coordinate identification and report the results

Chromosome 2:16,598,705-16,611,689 (13Kb)

Reference



Insertion site

