

1 **The *Arabidopsis thaliana* mobilome and its impact at the species level**

2

3 Leandro Quadrana¹, Amanda Bortolini Silveira¹, George F. Mayhew², Chantal LeBlanc⁴, Robert A.
4 Martienssen³, Jeffrey A. Jeddloh², and Vincent Colot^{1*}

5 ¹Institut de Biologie de l'Ecole Normale Supérieure (IBENS), Centre National de la Recherche Scientifique
6 (CNRS), Institut National de la Santé et de la Recherche Médicale (INSERM), Ecole Normale Supérieure, F-
7 75005 Paris, France

8 ²Roche NimbleGen, Inc, Madison, WI, 53719, USA

9 ³Howard Hughes Medical Institute-Gordon and Betty Moore Foundation, Watson School of Biological
10 Sciences, Cold Spring Harbor Laboratory, 1 Bungtown Road, Cold Spring Harbor, NY 11724, USA

11 ⁴Department of Molecular, Cellular & Developmental Biology. Yale University. 352A Osborn Memorial
12 Laboratories. 165 Prospect St, New Haven, CT, 06511, USA

13

14

15

16 * Corresponding author: colot@biologie.ens.fr

17

18 Transposable elements (TEs) are powerful motors of genome evolution yet a comprehensive assessment
19 of recent transposition activity at the species level is lacking for most organisms. Here, using genome
20 sequencing data for 211 *Arabidopsis thaliana* accessions taken from across the globe, we identify
21 thousands of recent transposition events involving half of the 326 TE families annotated in this plant
22 species. We further show that the composition and activity of the “mobilome” vary extensively between
23 accessions in relation to climate and genetic factors. Moreover, TEs insert equally throughout the genome
24 and are rapidly purged by natural selection from gene-rich regions because they frequently affect genes,
25 in multiple ways. Remarkably, loci controlling adaptive responses to the environment are the most
26 frequent transposition targets observed. These findings demonstrate the pervasive, species-wide impact
27 that a rich mobilome can have and the importance of transposition as a recurrent generator of large-
28 effect alleles.

29

30 INTRODUCTION

31 Transposable elements (TEs) are sequences that move and replicate around the genome. Depending on whether
32 their mobilization relies on a RNA or DNA intermediate, they are classified as retrotransposons (class I) or
33 DNA transposons (class II), respectively (Slotkin and Martienssen 2007). TEs are further subdivided into
34 distinct families, the prevalence of which differs between organisms because of a complex array of factors,
35 including variable transposition activity and diverse selection pressures (Barrón et al. 2014). Given their mobile
36 nature, TEs pose multiple threats to the physical and functional integrity of genomes. In particular, TEs can
37 disrupt genes through insertion and also through excision in the case of DNA transposons. Thus, TE
38 mobilization is a source of both germline and somatic mutations (Richardson et al. 2015; Barrón et al. 2014;
39 Lisch 2013). Although TEs are endogenous mutagens with potentially catastrophic effects, their mobilization
40 might sometimes be beneficial. In fact, soon after their discovery, Barbara McClintock named TEs “controlling
41 elements” to emphasize their role in the control of gene action (McClintock 1956). In mammals, transposition
42 of some *LINE 1* retrotransposons occurs extensively during embryogenesis as well as in the adult brain, again
43 suggesting functional relevance of somatic TE mobilization (Richardson, Morell, and Faulkner 2014).
44 Nonetheless, TEs are under tight control to limit their mutational impact both within and across generations. In
45 plants and mammals, a major control is through epigenetic silencing mechanisms, including DNA methylation
46 and these mechanisms can in turn have “epimutagenic” effects on adjacent genes (Richards, 2006; Slotkin and
47 Martienssen 2007; Weigel and Colot 2012; Heard and Martienssen 2014; Quadrana and Colot 2016).
48 Despite the many documented short-term as well as evolutionary consequences of TE mobilization (Rebollo,
49 Romanish, and Mager 2012; Trono 2016; Kazazian 2004; Elbarbary, Lucas, and Maquat 2016; Bennetzen and
50 Wang 2014; Lisch 2013), TEs are among the least investigated components of genomes, mainly because they
51 are present in multiple, often degenerated copies, which complicate analysis. Thus, a species-wide view of the
52 mobilome - i.e. of the set of TE families with transposition activity - is lacking for most organisms.
53 Studies in humans suggest that although at least half of the 3Gb genome is made up of TE sequences, mainly
54 belonging to *LINE 1* and *SINE* retrotransposon families, a few these only and none of the other TE families
55 contain mobile copies (Richardson et al. 2015). In contrast, the number of TE families that have retained
56 transposition activity is much larger in the mouse and these families include several so-called endogenous
57 retroviruses (ERVs) in addition to *LINEs* and *SINEs* (Richardson et al. 2015). In *Drosophila*, which has a much
58 smaller genome (~120 Mb) characterized by a large repertoire of class I and class II TE families, the situation

59 is again different with most of these TE families likely mobile (Mackay et al. 2012; Cridland et al. 2013; Kofler,
60 Nolte, and Schlötterer 2015; Rahman et al. 2015). However, in this species and even more so in mammals, the
61 population genetics of the mobilome remains poorly characterized.

62 The flowering plant *A. thaliana* is particularly attractive for conducting a systematic survey of the mobilome
63 and of its molecular as well as phenotypic impact at the species level. First, like *Drosophila*, *A. thaliana* has a
64 compact genome and a large repertoire of class I and class II TE families (The Arabidopsis Genome Initiative
65 2000; Ahmed et al. 2011; Joly-Lopez and Bureau 2014). Thus, most TE families are of relatively small size,
66 which facilitates their study. Second, *A. thaliana* occupies a wide range of habitats across the globe and
67 representative accessions have been extensively characterized both genetically and phenotypically (Weigel and
68 Nordborg 2015). Third, whole genome sequencing has been performed for >1000 *A. thaliana* accessions and
69 DNA methylome as well as transcriptome data are also available for hundreds of these (Schmitz et al. 2013;
70 Long et al. 2013; Dubin et al. 2015; Cao et al. 2011). Finally, genome-wide association studies (GWASs) are
71 straightforward in this species (Weigel and Nordborg 2015).

72 Here we present a comprehensive assessment of the *A. thaliana* mobilome, which provides important novel
73 insights into the population genetics of TE mobilization and radically changes the prevailing view of limited
74 transposition potential in this species (Hu et al. 2011; Maumus and Quesneville 2014). Specifically, we show
75 that the *A. thaliana* mobilome is composed of a very large number of class I and class II TE families overall,
76 but differs markedly among accessions. We further show that TE mobilization is a complex trait and we have
77 identified environmental as well as genetic factors that influence transposition in nature. These factors include
78 the annual temperature range, the TE themselves and multiple gene loci, notably *MET2a*, which encodes a
79 poorly characterized DNA methyltransferase. In addition, we present compelling evidence that TEs insert
80 throughout the genome with no overt bias and that the mobilome has a pervasive impact on the expression and
81 DNA methylation status of adjacent genes. These and other observations indicate that purifying selection is
82 most probably the main factor responsible for the differential accumulation of TE sequences along the *A.*
83 *thaliana* genome and notably their clustering in pericentromeric regions. Finally, we reveal the importance of
84 the mobilome as a generator of large-effect alleles at loci underlying adaptive traits. Collectively, our
85 approaches and findings provide a unique framework for detailed studies of the dynamics and impact of
86 transposition in nature.

87

88 RESULTS

89 Composition of the *A. thaliana* mobilome

90 The reference genome sequence of *A.thaliana* is 125 Mb long (TAIR 10) and contains ~32000, mostly
91 degenerate TE copies that belong to 326 distinct families (The Arabidopsis Genome Initiative 2000; Ahmed et
92 al. 2011). So far, transposition activity has been documented experimentally for nine TE families, mainly based
93 on studies carried out in the reference accession Col-0 (Ito and Kakutani 2014; Tsay et al. 1993). To assess
94 species-wide the composition of the *A.thaliana* mobilome, we used publically available Illumina short genome
95 sequence reads (Schmitz et al. 2013; Schneeberger et al. 2011). First, we looked for TE copy number variation
96 (CNV) between the reference accession Columbia (Col-0) and 211 accessions taken from across the globe. To
97 limit the problem posed by the presence of TEs in multiple copies across the genome, with varying degrees of
98 similarity to each other, we performed an aggregated CNV analysis based on the 11,851 annotated Col-0 TE
99 sequences longer than 300bp (see ‘Materials and methods’). CNVs were detected for 263 TE families (Figure
100 1A and B; Figure 1-Source data 1; see ‘Materials and methods’), in keeping with the results of a previous study
101 indicating that the vast majority of the TE sequences annotated in the Col-0 reference genome are absent from
102 that of at least one of 80 accessions analyzed (Cao et al. 2011).

103 Since CNVs could reflect either recent TE mobilization or the gain or loss of TE copies through other types of
104 chromosomal rearrangements, we then looked among the unmapped Illumina short reads for so-called “split-
105 reads” that contain TE extremities. Crucially, because most TE families generate short target site duplications
106 (TSDs) of fixed size upon insertion, TSDs can serve as signatures of *bona fide* transposition events. We
107 therefore developed a pipeline for the systematic identification of split-reads covering TE junctions that are
108 absent from the reference genome and that produce, when mapped to the insertion site, a sequence overlap of
109 the size of TSDs (3 to 15 bp, depending on the TE family, Figure 1C and D; see ‘Materials and methods’). Our
110 pipeline differs in that respect from that used in another study to detect the presence/absence of reference and
111 non-reference TE insertions in the same set of accessions (Stuart et al. 2016). Results produced by our pipeline
112 for the 292 annotated TE families that create TSDs upon transposition were verified visually to eliminate false
113 positives (Figure 1-figure Supplement 2; see ‘Materials and methods’). Following this approach, non-reference
114 TE insertions with TSDs were identified for 131 TE families in total (Figure 1-Source data 2), which all also
115 show CNV (Figure 1-Source data 1).

116 Most (86%) non-reference TE insertions with TSDs are private or shared only by a few accessions and thus

117 they typically correspond to recently derived alleles, as expected (Figure 1E). Moreover, recent transposition
118 activity is only detected for between four and 66 TE families in any given accession, thus indicating large
119 variations in the composition of the mobilome among accessions (Figure 1F). However, we have probably
120 identified most of the annotated TE families that compose the mobilome at the species level, since the number
121 of TE families defined as mobile by the split-reads approach reaches a plateau after examining just 74
122 accessions (Figure 1G). The 53 class I *COPIA* families and the 40 class II Mutator-like (*MuDR*) families are the
123 most mobile, as they account for 1408 and 729 of the 2835 non-reference TE insertions with TSDs identified in
124 total, respectively (Figure 1H and Figure 1-Source data 2 and 3). However, the number of non-reference
125 insertions per accession is always small (<16) for any given family (Figure 1-Source data 2), thus suggesting a
126 lack of recent transposition bursts.

127 The ability to detect non-reference TE insertions with TSDs using split-reads is strongly dependent on read
128 depth as well as sequence composition at the insertion site (Figure 1-figure supplement 1A and B; Hénaff et al.
129 2015). To assess the extent of this limitation, we used the assembled Ler-1 genome sequence recently obtained
130 using PacBio long reads (see ‘Materials and methods’). Although not annotated, this sequence assembly can
131 nonetheless serve to identify by whole genome comparison the Col-0 TEs flanked by TSDs that are absent from
132 the corresponding position in Ler-1 (see ‘Materials and Methods’). A total of 142 TEs belonging to 80 distinct
133 families were identified in this way (Figure 1-Source data 4), two numbers that are consistent with estimates
134 obtained using other approaches (Ziolkowski et al. 2009; Hénaff et al. 2015). In contrast, we could detect only
135 78 Col-0-specific TEs with TSDs belonging to 49 TE families when using the split-reads pipeline to map Col-0
136 short reads onto the assembled Ler-1 genome. These results indicate therefore that the split-reads approach has
137 a low sensitivity (Figure 1-figure supplement 1C; 45% false negatives and 10% false positives; False Discovery
138 Rate: 15.3%).

139 To obtain an independent estimation of the composition of the mobilome, we also performed TE sequence
140 capture (TE-capture; Baillie et al. 2011). Briefly, probes were designed to cover the 5’ and 3’ extremities of 310
141 TE elements belonging to 181 distinct families, including 117 of the 131 TE families identified as mobile using
142 the split-reads approach (see ‘Materials and methods’). Using genomic DNA extracted from 12 randomly
143 chosen accessions (Figure 1-figure supplement 1D and E), we could validate by TE-capture most (87%) of the
144 non-reference TE insertions with TSDs that were detected by the split-reads approach (Figure 1-figure
145 supplement 1F; see ‘Materials and methods’). As expected, TE-capture also uncovered many additional non-

146 reference TE insertions with TSDs for these TE families (Figure 1-figure supplement 1F). However, no such
147 insertions were detected for the other TE families that could be captured but which were not identified as
148 mobile by the split-reads approach in any of the 12 accessions. These results confirm that despite the low
149 sensitivity of the split-reads approach, we have probably identified most of the TE families with TSDs that
150 compose the *A. thaliana* mobilome at the species level. Finally, non-reference insertions were also identified for
151 30 TE families (including 15 *HELITRON* families) that could not be analyzed using our split-reads pipeline
152 because they do not produce TSDs or have insertion sites located in low complexity regions (Figure 1-figure
153 supplement 1B). Since most of the non-reference insertions for these 30 TE families are present in only one or
154 two of the 12 accessions examined (Figure 1-figure supplement 1G), they likely reflect recent transposition
155 events. Thus, there are altogether at least 165 TE families with recent transposition activity at the species level.
156 Moreover, based on the TE-capture data, we can estimate that depending on their divergence time, any two
157 accessions have accumulated between ~200 and ~300 newly transposed TE copies (Figure 1-figure supplement
158 1H).

159

160 **TE mobilization as a complex trait**

161 The observation that the composition of the mobilome differs extensively between accessions (Figure 1F)
162 suggests that it is influenced by both environmental and genetic factors. To try to identify such factors, we first
163 established that copy number (CN) correlates positively with the number of TE insertions with TSDs that are
164 detected by TE-capture (Figure 2-figure supplement 1; see ‘Materials and methods’). Thus, CNV provides a
165 reliable and quantitative estimator of differential TE mobilization between accessions, which we used to
166 analyze the 113 TE families that were defined as mobile based both on the split-reads approach and TE-capture
167 (Figure 2-Source data 1).

168 Controlling for population stratification and considering thirteen geo-climatic variables (Hancock et al. 2011),
169 we uncovered robust correlations with CN for 15 class I and class II TE families. Among these, *ATCOPIA2* and
170 *ATCOPIA78* share the highest number of geo-climatic variables correlated with CN (Figure 2-figure
171 supplement 2). Moreover, the strongest correlation is between temperature annual range and CN for
172 *ATCOPIA78* (Figure 2A and B). Given that at least one member of this TE family is transcriptionally induced
173 by heat shock in the Col-0 accession (Ito et al. 2011; Cavrak et al. 2014), *ATCOPIA78* provides a compelling
174 case of a causal link between climate and TE mobilization.

175 We next explored the possibility of using GWASs to identify genetic variants influencing TE mobilization (see
176 ‘Materials and methods’). For 33 TE families, a disproportionately large number of SNPs are associated with
177 CNV, preventing further analysis. For the remaining 80 TE families, SNPs in linkage disequilibrium with each
178 other and associated with CNVs delineate 230 loci. Moreover, 34% of these loci are also identified by GWAS
179 using whole genome sequencing data obtained for another 180 accessions taken from Sweden (Long et al. 2013)
180 (Figure 2-figure supplement 3A). This substantial overlap suggests a similar genetic architecture for the *A.*
181 *thaliana* mobilome at both the local and global scales, which prompted us to perform a joined GWAS using all
182 391 accessions in order to increase both sensitivity and specificity. Depending on the TE family, GWAS
183 identified between 0 and 33 loci and collectively, associations explain between 2% and 67% of the variance in
184 CN (Figure 2C).

185 Among the 334 loci detected in total, 130 encompass sites with reference or non-reference TE insertions
186 (Figure 2C and Supplementary file 1). Furthermore, each of these local (*cis*) genetic variants explains on
187 average 5.2% of the total variance compared with 2.2% for distal (*trans*) genetic variants. The higher
188 explanatory power of *cis* variants is of course to be expected, as the TEs themselves are the primary
189 determinants of the transposition process. Indeed, almost all *cis* SNPs that map to TE sequences in the reference
190 genome are likely causal as they affect sequences involved in transposition, for example the long terminal
191 repeats (LTRs) and the various open reading frames of LTR retrotransposons (Figure 2-figure supplement 3B
192 and C).

193 While *cis* loci collectively explain most CN variance for class I TE families, this is not the case for class II TE
194 families (Figure 2-figure supplement 3D). Given that class II TEs move by a cut and paste mechanism, some
195 *trans* loci could in fact correspond to sites of excision. However, we could not find evidence of excision
196 footprints, such as small insertions or deletions. Alternatively, the larger fraction of CN variance explained by
197 *trans* loci for class II TE families may in part result from many of these families being non-autonomous, i.e.
198 requiring factor(s) encoded by other TEs for their mobilization. Consistent with this possibility, the proportion
199 of CN variance associated with *trans* loci as well as the number of TE annotations overlapping *trans* loci are
200 higher for non-autonomous than autonomous class II TE families (Figure 2-figure supplement 3E). Although
201 we did not investigate *trans* mobilization in depth, we readily identified one presumed case, involving the non-
202 autonomous and autonomous MuDR families *ATDNAIT9A* and *VANDALI6*, respectively (Figure 2-figure
203 supplement 3F). Of note, CN do not co-vary between these two families, which could indicate that their

204 transposition is differentially controlled. Finally, “false” *trans* loci could also be caused by non-reference TE
205 insertions that are sufficiently frequent to be in linkage disequilibrium with SNPs used for the GWASs but that
206 we have failed to detect. However, such *trans* loci are not expected to be more prevalent for class II than class I
207 TE families and should be very rare in any case given that the probability of missing moderately frequent (>5%)
208 non-reference TE insertions by both the split-reads approach and TE-capture is low (Figure 2-figure supplement
209 3G).

210 Since transposition is controlled by multiple protein activities in addition to those encoded by TEs, we also
211 examined genes located within or adjacent to *trans* loci. Overall, these genes do not appear to be enriched for
212 any particular function and most of them are specific to a single TE family (Figure 2D and Figure 2-figure
213 supplement 3H). These observations indicate either a complex genetic architecture of mobilome variation
214 and/or spurious *trans* associations such as those considered above. Nonetheless, 22 *trans* loci stand out as they
215 show association with CNV for two or more TE families (Figure 2D) and a causal link is evident in two cases.
216 Indeed, the locus associated with CNV for respectively the retrotransposon and DNA transposon families
217 *ATGP2* and *ATENSPM2* encodes the transcription factor ARF23, which recognizes motifs that are enriched in
218 the sequence of these TEs. The second locus is associated with CNV for the largest number of TE families (four
219 class I and three class II families, Figure 2E) and encodes the MET2a protein, a poorly characterized homolog
220 of the main DNA maintenance methyltransferase MET1. Moreover, one of the *MET2a* SNPs is presumably
221 causal as it leads to a non-synonymous amino-acid substitution (G519E) in a conserved domain of the protein
222 (Figure 2F) that in the mammalian homolog Dnmt1 is required for the targeting to replication foci (Klein et al.
223 2011). A role for *MET2a* is also supported by the observation that *met2a* mutant plants (Stroud et al. 2013) lose
224 some DNA methylation exclusively over mobile TE families. Furthermore, loss of methylation is more
225 pronounced when only considering the seven TE families that show a *MET2a* association (Figure 2G).
226 Intriguingly, CHG sites (where H=A, T or C), which are poor substrates for MET1 or Dnmt1 compared to CG
227 sites (Law and Jacobsen 2010), are the most affected in the *met2a* mutant. Whether or not this observation
228 reflects an atypical recognition specificity for MET2a remains to be determined. Finally, we note that GWASs
229 failed to detect any association with genes known to be involved in the epigenetic silencing of TEs (Ito and
230 Kakutani 2014) such as *MET1* and *DDMI*, presumably because of their essential function.

231

232 **Genome localization of newly inserted TEs**

233 In *A. thaliana* as in many other eukaryotes, TE sequences tend to cluster in pericentromeric regions (The
234 Arabidopsis Genome Initiative 2000). Mechanistically, such clustering may result from insertion bias, selective
235 constraints or differential elimination of TE copies through ectopic homologous recombination (Barrón et al.
236 2014). To distinguish between these possibilities, we looked at the genomic location of the 2835 distinct non-
237 reference TE insertions with TSDs detected with the split-reads approach and found that they are distributed
238 almost evenly along chromosomes (Figure 3A). Since a similar distribution is observed for the non-reference
239 TE insertions with TSDs detected specifically using TE capture (Figure 3-figure supplement 1A), we can rule
240 out an ascertainment bias of the split-reads approach towards non-reference TE insertions located along the
241 chromosome arms. However, there is a clear trend towards a more pericentromeric localization when only
242 considering non-reference TE insertions with TSDs that are shared by two or more accessions and that are thus
243 presumably more ancestral (Figure 3B and Figure 3-figure supplement 1B). Moreover, the density of non-
244 reference TE insertions with TSDs positively correlates with the recombination rate but negatively with gene
245 density (Figure 3-figure supplement 1C, D and E). Finally, except for the *COPIA* families, non-reference TE
246 insertions with TSDs are globally under-represented within genes, where they are expected to be most
247 detrimental (Figure 3C and Figure 3-figure supplement 1F and Figure 3-figure supplement 2A). Collectively,
248 these observations provide strong evidence that TEs are preferentially purged over time from the chromosome
249 arms because of their deleterious effects on adjacent genes rather than as a consequence of ectopic homologous
250 recombination.

251 Although most TE families show not overt insertion bias at the genome scale, there are clear local insertion
252 preferences. In agreement with previous observations (Fu et al. 2013; Miyao et al. 2003), private non-reference
253 *COPIA* and *MuDR* insertions with TSDs are enriched at coding sequences and transcriptional start sites (TSS),
254 respectively (Figure 3-figure supplement 2A). In addition, insertion sites for most TE superfamilies are
255 enriched in specific DNA sequence motifs or exhibit biased sequence composition (Figure 3-figure supplement
256 2B and C; Supplementary file 2). For example, *LINEs* tend to insert within poly(A) tracks, as expected for this
257 superfamily of non-LTR retrotransposons, which integrate into the genome via poly(A)-dependent, target site-
258 primed reverse transcription (TPRT; Beck et al. 2011).

259

260 **Impact of newly inserted TEs on the expression of adjacent genes**

261 Transcriptome analyses in the reference accession Col-0 have revealed that *A.thaliana* genes nearest to TE

262 sequences are expressed at lower levels compared with the genome-wide distribution of gene expression,
263 suggesting that TE insertions tend to reduce the expression of neighboring genes (Hollister and Gaut 2009). To
264 investigate more directly the impact of TEs on the genes within or near which they insert, we examined RNA-
265 seq data available for 144 accessions (Schmitz et al. 2013). Specifically, we considered all non-reference TE
266 insertions with TSDs and calculated for each gene located within 1kb of them (1616 genes in total), the ratio
267 between the expression level in the accession(s) harboring the insertion and the median expression level in the
268 accessions devoid of the insertion. Expression ratios expected under the null hypothesis (no effect of the TE
269 insertions) were calculated by taking 10^6 randomly chosen sets of 1616 genes and assigning for each set the TE
270 insertion ‘presence/absence’ label randomly among the 144 accessions (see ‘Materials and Methods’).
271 Comparison of the distribution of the observed and expected expression ratios indicates that for a large fraction
272 of genes, expression is indeed significantly altered when TEs insert within or near them (Figure 3D, $P < 1.9 \times 10^{-5}$).
273 These alterations are most pronounced for the *COPIA* insertions, which are overrepresented in genes and less
274 pronounced for the *MuDR* insertions, despite the latter being overrepresented around the TSS of genes (Figure
275 3E and Figure 3-figure supplement 2A). Although other TE superfamilies show similar trends, we could not
276 draw firm conclusions in these cases because of insufficient statistical power (Figure 3-figure supplement 3).
277 This notwithstanding, it is clear that TE insertions induce both increases and decreases in gene expression with
278 equal frequency (Figure 3D and E). Thus our findings contradict the prevailing view of a dampening effect of
279 TE insertions on the expression of adjacent genes (Hollister and Gaut 2009) and suggest instead a stronger
280 selection against TE insertions when they occur close to highly expressed genes.

281 To complement the re-analysis of transcriptome data, we also measured by RT-qPCR the expression level of 19
282 genes with recent *COPIA* or *MuDR* insertions, using nine different accessions grown under control conditions
283 or subjected to heat shock. *COPIA* insertions were found to have more dramatic and systematic effects on gene
284 expression in stressed plants (Figure 3-figure supplement 4) which in the case of *ATCOPIA78* can be related to
285 its transcriptional sensitivity to heat shock (Ito et al. 2011; Cavrak et al. 2014). On the other hand, we could not
286 detect any effect of the *MuDR* insertions under the two conditions tested (Figure 3-figure supplement 4). These
287 findings are in agreement with those of the transcriptome analysis and indicate in addition that the effect of TE
288 insertions on the expression of adjacent genes can vary substantially in relation to the environment.

289

290 **Impact of newly inserted TEs on the DNA methylation status of adjacent sequences**

291 TE sequences are typically targeted by the RNA-directed DNA methylation (RdDM) machinery in *A. thaliana*
292 (Lippman et al. 2004; Lister et al. 2008; Cokus et al. 2008) and we have previously provided genome-wide
293 evidence that DNA methylation can spread from RdDM targets to flanking sequences, with consequences on
294 gene expression (Ahmed et al. 2011). To investigate the effect of new TE insertions on the DNA methylation
295 status of adjacent sequences, we used MethylC-Seq data available for 140 accessions (Schmitz et al. 2013).
296 Analysis of this data set first indicated that mobile TE families have on average higher CG, CHG and CHH
297 methylation than non-mobile TE families (Figure 4A). Furthermore, DNA methylation is also higher for most
298 mobile TE families in the accessions that have experienced recent transposition activity (Figure 4-figure
299 supplement 1A). These observations prompted us to examine in addition methylome data obtained for several
300 mutation accumulation (MA) lines (Becker et al. 2011; Schmitz et al. 2011). We found that mobile TE families
301 suffer less sporadic DNA methylation loss than non-mobile families (Figure 4B). These findings are entirely
302 consistent with DNA methylation playing an important role in the control of TE mobility and suggest in turn
303 that most of the recent TE insertions we have identified are present in the methylated state. Moreover, given
304 that DNA methylation is likely established over newly inserted TE copies by RdDM in a progressive manner
305 across multiple generations (Teixeira et al. 2009; Mari-Ordóñez et al. 2013), unmethylated non-reference TE
306 insertions should be mainly private and reflect very recent transposition events.

307 Based on these considerations, we next analyzed the DNA methylation status of 1543 TE insertion sites for
308 which reliable data could be extracted across all 140 accessions (Figure 4C). Approximately 10% of insertion
309 sites are methylated in most accessions, including systematically in the one(s) containing the TE. As expected,
310 these sites are preferentially located within TE-rich, pericentromeric regions. In contrast, another 40% of sites
311 are devoid of methylation in the accession(s) containing the TE insertion as well as in most of the other
312 accessions. This absence of adjacent DNA methylation could indicate either that the TE insertions themselves
313 are unmethylated or else that DNA methylation does not spread from them. Finally, 50% of sites are methylated
314 exclusively or almost exclusively in the accession(s) with the TE insertion (Figure 4-figure supplement 1B),
315 thus suggesting that at these sites TEs are methylated and that DNA methylation spread into adjacent sequences.

316 Why some sites may be refractory to DNA methylation spreading when others are not is unclear, as we did not
317 identify any feature that could distinguish them, such as the identity of the TE or the sequence composition at
318 the insertion site.

319 Further analysis of DNA methylation associated with TE insertions indicates that it involves the three sequence

320 contexts and that it generally extends for up to 300 bp on both sides of the insertion (Figure 4D and E; Figure 4-
321 figure supplement 1B), a distance that closely matches that previously reported for the spreading of DNA
322 methylation from RdDM targets (Ahmed et al. 2011). For 243 insertion sites however, DNA methylation
323 extends over much longer distances (up to 3.5 kb) on one or the other side of the insertion (Figure 4D and E).
324 While most of these sites lie within or close to genes, the TE insertions are not preferentially orientated with
325 respect to gene transcription (Figure 4F), which rules out sense-antisense transcription as the likely trigger for
326 this long-distance DNA methylation. Proximity to pericentromeric heterochromatin can also be ruled out,
327 because most of the genes with long-distance DNA methylation are located on the chromosome arms (Figure 4-
328 figure supplement 1C). To explore potential mechanisms further, we made use of the wealth of epigenomic
329 data available in Col-0 to examine the 142 Col-0 TE insertions with TSDs that are absent from the assembled
330 Ler-1 genome sequence (Figure 1-Source data 4). Methylome data (Stroud et al. 2013) indicate that 121 of
331 these 142 Col-0 TE copies are methylated and that DNA methylation tends to extend into flanking sequences,
332 predominantly over short distances, but occasionally over much longer distances (<300pb: 63 TE
333 insertions; >1kb: 36 TE insertions; Figure 1-Source data 4). These results confirm those obtained for the non-
334 reference TE insertions. In addition, analysis of Col-0 small RNA-seq data (Fahlgren et al. 2009) indicates that
335 in contrast to short-distance DNA methylation, long-distance DNA methylation aligns with 24-nt siRNAs
336 (Figure 4-figure supplement 1D). Thus, genes affected by the latter type of DNA methylation have presumably
337 become secondary targets of RdDM, as was shown for a transgene (Kanno et al. 2008; Daxinger et al. 2009).
338 Moreover, genes affected by long-distance DNA methylation in accessions other than Col-0 tend in this
339 accession, where they are by definition in the ancestral state, to be most highly expressed in both pollen and
340 seeds and more highly expressed in these two organs than genes affected by short-distance methylation (Figure
341 4G). Given that RdDM activity is also maximal in these organs (Teixeira and Colot 2010), our observations
342 suggest that secondary RdDM results from the concomitance of strong transcription and strong RdDM at target
343 loci.

344 Finally, our analysis of TE-associated DNA methylation indicates that it accounts for at least 7% of the so-
345 called gene C-DMRs (i.e. regions of differential methylation at CG, CHG and CHH sites) identified in nature,
346 which are typically low frequency gain of DNA methylation variants (Schmitz et al. 2013). These and similar
347 findings reported recently (Stuart et al. 2016) confirm and extend previous results that first indicated that many
348 natural gene C-DMRs are not *bona fide* epialleles but rather new alleles caused by TE insertions (Schmitz et al.

349 2013). However, close examination of one TE-insertion allele shared among 13 accessions indicates that it is
350 subject to epigenetic variation in nature, as it is present in the unmethylated state in one accession (Figure 4-
351 figure supplement 1B).

352

353 **TE insertions as motors of adaptive changes**

354 Although TEs tend to insert with no overt bias at the genome scale (Figure 3A), we detected nineteen 10kb
355 windows with a high load of non-reference TE insertions (Figure 5A). Such enrichment could result from
356 insertion preference or else an absence of strong negative selection. In fact, three of these 10kb windows span
357 genes encoding nucleotide-binding domain and leucine-rich repeat containing (NLR) proteins, which function
358 as immune receptors in plants and are known to be under strong diversifying selection (Chae et al. 2014).
359 Moreover, a fourth 10kb window spans the gene *FLC*, which encodes a key repressor of flowering and is one of
360 the main genetic factors causing natural variation in the onset of flowering (Ietswaart, Wu, and Dean 2012).
361 Remarkably, the *FLC* locus has the highest number of non-reference TE insertions (seven in total) across the
362 genome. These insertions belong to several COPIA families and affect four distinct *FLC* haplotypes in total
363 (Figure 5A and Figure 5-figure supplement 1A). Moreover, five insertions are located within the first intron
364 (Figure 5B), which plays an important role in the epigenetic regulation of *FLC* in response to cold (Ietswaart,
365 Wu, and Dean 2012). Although four of these insertions as well as another intronic insertion were previously
366 identified among early flowering accessions (Liu et al. 2004; Lempe et al. 2005), causality could not be
367 established unequivocally because of numerous other sequence polymorphisms in complete linkage
368 disequilibrium. To obtain direct proof of a causal role for the seven TE insertions we identified, we used
369 publically available transcriptomic (Schmitz et al. 2013) as well as phenotypic data (Li et al. 2010; Lempe et al.
370 2005) and compared *FLC* expression as well as flowering time among accessions that have the same *FLC*
371 haplotype but differ by the presence or absence of a TE insertion (see ‘Material and Methods’). Results of these
372 comparisons indicate that the TE-containing accessions have systematically much reduced *FLC* expression and
373 flower much earlier than their TE-free counterparts (Figure 5C and D; Figure 5-figure supplement 1B and C).
374 Thus, we can conclude that TEs are recurrent generators of major effect *FLC* alleles, which in turn suggests that
375 they contribute significantly to the high level of allelic heterogeneity observed at this locus (Li et al. 2014)

376

377 **DISCUSSION**

378 We have shown that the *A. thaliana* mobilome is particularly rich at the species level, being composed of at
379 least half of the 326 TE families that are annotated in the reference genome. This finding is at odds with the
380 prevailing view that most TE families are mere molecular fossils in *A. thaliana*, since they contain a much
381 lower proportion of “young”, i.e. non-degenerated TE copies than in the close relative *A. lyrata* (Hu et al. 2011;
382 Maumus and Quesneville 2014). Furthermore, we provide definitive evidence that TEs insert throughout the
383 genome, with no overt bias towards pericentromeric regions, which contrasts with the observed clustering of
384 annotated TE sequences around centromeres. However, these discrepancies are easily resolved, since we have
385 also shown that despite the richness of the *A. thaliana* mobilome, most TEs tend to be rapidly purged by natural
386 selection in this species when they insert in the chromosome arms, which are gene-rich. Indeed, our systematic
387 survey indicates that TEs have pervasive effects on the expression and DNA methylation status of genes near or
388 within which they insert. Incidentally, the deleterious effects associated with most transposition events in *A.*
389 *thaliana* may also explain in part the fact that we did not detect any recent transposition bursts, as these should
390 be strongly counter-selected. Furthermore, because *A. thaliana* is predominantly self-fertilizing, the purging of
391 deleterious TE insertions should be accelerated in this species compared to *A. lyrata*, which is an obligated out-
392 crosser. Given this difference in mating systems, the TE population dynamics of *A. thaliana* and *A. lyrata* are
393 expected to differ significantly (Lockton and Gaut 2010). Thus, homologous recombination could play a more
394 prominent role in the elimination of TE insertions in *A. lyrata*, as is thought to be the case in *D. melanogaster*
395 (Barrón et al. 2014). However, comprehensive studies similar to those presented here remain to be performed
396 for *A. lyrata* in order to identify conclusively the forces that shape the TE landscape in this species.

397 We have also shown that the composition and activity of the mobilome vary greatly between accessions.
398 GWASs revealed that this variation is caused in part by sequence polymorphisms within the TEs themselves
399 (*cis* variation), which is in agreement with empirical data and theoretical models indicating that TE families
400 contain only one or a few active, autonomous (i.e. master) TE copies at any one time (Beck et al. 2011). The
401 fact that we could readily detect such *cis* variants may again be linked to the mating system of *A. thaliana*,
402 which on the one hand should increase the probability that disabling mutations accumulate within the few
403 active TE copies that are present within a given lineage, before these copies could transpose further; and on the
404 other hand should decrease the probability of acquiring new active copies through crosses.

405 Another important result of the GWASs is that natural variation at the *MET2a* locus, which encodes a poorly
406 characterized DNA methyltransferase, has a significant impact on mobilome composition and activity across

407 accessions, being associated with differential transposition activity for seven distinct class I and class II TE
408 families. While the role of MET2a in transposition control remains to be determined experimentally, it is
409 noteworthy that none of the epigenetic repressors of TE activity identified through genetic screens, such as
410 MET1 or DDM1, are associated with natural variation of the mobilome, presumably because of their essential
411 function. Altogether, these findings illustrate the power of GWASs in identifying the genetic factors affecting
412 transposition in nature.

413 Being a complex trait, TE mobilization is also modulated by environmental factors, and we have identified
414 temperature annual range as a clear contributor to the variation in *ATCOPIA78* mobilization across accessions.
415 Remarkably, this TE family has generated several rare alleles with large effects at the *FLC* locus, which is a
416 major genetic determinant of the onset of flowering in nature. It is therefore tempting to speculate that
417 *ATCOPIA78* may endow *A. thaliana* with a unique ability to adapt to global warming and the associated
418 increase in droughts by facilitating the creation of early flowering *FLC* alleles. Additionally, these observations
419 may provide insights into how *A. thaliana* has been able to colonize efficiently the entire northern hemisphere
420 from a few glacial refugia located in Southern Europe (François et al. 2008).

421 In summary, our findings have far reaching implications, as they indicate that part of the missing heritability
422 that plagues many GWASs may be accounted for by recent and thus rare TE insertion alleles with large effects
423 (Vinkhuyzen et al. 2013; Brachi, Morris, and Borevitz 2011). More generally, our study highlights the need for
424 similar species-wide explorations of the mobilome in a variety of organisms in order to assess the true
425 mutational and epimutational impact of transposition as well as its contribution to natural phenotypic variation.
426 In this respect, it can be anticipated that the advent of long read sequencing technologies will greatly facilitate
427 such studies, especially in organisms with large, repeat-rich genomes.

428

429 **MATERIALS AND METHODS**

430 **Data sources.** Whole genome sequencing data (Illumina short paired-end reads) for *A. thaliana* accessions
431 were obtained from the NCBI SRA archive for 210 accessions collected worldwide (Schmitz et al. 2013)
432 SRA012474 as well as for the 180 Swedish accessions (Long et al. 2013; SRA052536). Illumina short paired-
433 end reads from Ler-1 (Schneeberger et al. 2011) and the assembled Ler-1 genome obtained using PacBio long
434 reads (unpublished) were retrieved from
435 <http://1001genomes.org/data/MPI/MPISchneeberger2011/releases/current/Ler-1/Reads/> and <https://github.com/>

436 PacificBiosciences/DevNet/wiki/Arabidopsis-P5C3, respectively. High quality SNP imputations for these
437 accessions were obtained from the 1001 Genomes Project (<http://1001genomes.org/>). Processed MethylC-seq
438 and RNA-seq data from accessions collected worldwide (Schmitz et al. 2013) were obtained from NCBI GEO
439 (GSE43857 and GSE43858, respectively). Processed BS-seq data for the *met2a* mutant (Stroud et al. 2013)
440 were obtained from NCBI GEO (GSE39901). Flowering time data were retrieved from the phenotypic
441 information deposited at <https://www.arabidopsis.org/servlets/TairObject?type=germplasm&id=6530472136> as
442 well as (Y. Li et al. 2010; Lempe et al. 2005).

443 Local climate data and coordinates for the 211 accessions analyzed using CNV were retrieved from (Hancock
444 et al. 2011) (http://bergelson.uchicago.edu/wp-content/uploads/2015/04/allvars948_notnormd_011311.txt.zip).
445 TE-capture sequencing data have been deposited in the ENA short read archive under project number
446 PRJEB11706 and accessions codes ERR1121179 to ERR1121190.

447

448 **CNV analysis.** We performed an aggregated CNV analysis for each TE family by considering only annotated
449 TE sequences longer than 300bp as a joined pseudo-annotation. Families composed in the Col-0 reference
450 genome exclusively of TE sequences shorter than 350bp (seven families, *ARI2*, *AR3*, *ATCLUST1*, *ATDNATA1*,
451 *ATMSAT1*, *ATSINE1* and *ATSINE3*), that contain too many (>1000) copies or that overlap with genomic regions
452 with abnormal coverage (four families, *ATREP15*, *ATREP10D*, *ATREP3* and *HELITRONY3*; see next section)
453 were not considered. Illumina pair-end short reads were mapped onto the TAIR10 reference genome using
454 Bowtie2 (using the arguments `--mp 13 --rdg 8,5 --rfg 8,5 --very-sensitive`) and PCR-duplicates were removed
455 using Picard. Read coverage (RC) was computed in non-overlapping windows of 100bp spanning the joined
456 pseudo-annotation for each TE family. RC over each bin was corrected for GC content (Yoon et al. 2009) and
457 normalized to the genome-wide RD. Normalized RC across joined pseudo-annotations were compared between
458 each accession and two independent re-sequenced Col-0 reference genomes (Schmitz et al. 2013; Long et al.
459 2013). CNVs were detected by performing a distribution-independent permutation test for each of the two
460 references. In order to be as stringent as possible, the maximum *P*-value from these two comparisons were
461 considered. These *P*-values were compared with an empirical null-distribution of *P*-values constructed by
462 randomizing a million times the TE annotation labels over the 100bp windows. We defined the false discovery
463 rate (FDR) of an observed *P*-value as the fraction of significant (that is, below the observed *P*-value in question)
464 hits found in the randomized set. We considered statistically significant TE-CNVs when FDR was below 10^{-5} .

465

466 **Filtering out genomic regions with aberrant coverage or low sequence complexity.** Read mapping to the
467 reference *A. thaliana* genome sequence revealed several regions with very high coverage, which correspond
468 mainly to the 180bp centromeric repeat unit, the 45S and 5S rDNA repeat units, ATHILA TE sequences,
469 telomeres and plastid DNA-like sequences. These as well as the few regions not suitable for sequence
470 alignment because of low complexity were identified by mapping genome sequencing reads obtained for two
471 independent lines of the Col-0 accession (Schmitz et al. 2013; Long et al. 2013) onto the TAIR10 genome
472 sequence. RC was calculated on consecutive non-overlapping windows of 100bp across the entire genome.
473 After correction for GC content (Yoon et al. 2009), consecutive windows (allowing one window gap) with a RC
474 greater or lower than three median absolute deviations from the median RC signal were merged to define larger
475 segments. These sequences spanned 1,125,487 bp (0.94% of the TAIR10 genome sequence) and were excluded
476 from all analyses.

477

478 **Identification of non-reference TE insertions with TSDs using split-reads.** The split-read analysis was
479 performed in four steps: (i) extraction of reads not mapping to the reference genome; (ii) forced mapping to a
480 collection of TE sequence extremities and soft-clipping of mapped reads; (iii) mapping of clipped reads to the
481 reference genome; (iv) identification of clipped reads that reveal target site duplications (TSDs). Briefly, for
482 each accession we retrieved reads that do not map to the TAIR10 reference genome (containing the SAM flag
483 4), but we did not make use of information about discordant paired-end reads, unlike in the Jitterbug approach
484 (Hénaff et al. 2015). Unmapped reads were then aligned (using Bowtie2 in --local mode to allow for soft clip
485 alignments) to a collection of 5' and 3' TE sequence extremities (300bp) obtained from the Col-0 for each TE
486 family or from Repbase Update (Jurka et al. 2005) in the case of the ARNOLDY2, ATCOPIA62, ATCOPIA95,
487 TA12, TAG1 families, which do not contain copies with intact extremities in the Col-0 genome. Next, we
488 selected all reads with one end (≥ 20 nt) mapping to a TE extremity (by locating reads where the CIGAR string
489 contained only one 'S' character with a value equal or greater to 20). These reads were recursively soft clipped
490 by 1nt and mapped to the TAIR10 reference genome using Bowtie2 (using the arguments --end-to-end --mp 13
491 --rdg 8,5 --rfg 8,5 --very-sensitive) until the soft-clipped read length reached 20nt. Read clusters composed of
492 five or more reads clipped from the same extremity and overlapping with read clusters composed of reads
493 clipped from the other extremity were taken to indicate the presence of a bona fide TE insertion only if the size

494 of the overlap was equal or less than 2-fold larger than that reported for TSDs for the corresponding TE family.
495 Putative non-reference TE insertions overlapping genomic regions with aberrant coverage (as defined above) or
496 located within inner pericentromeres (Clark et al. 2007) or spanning the corresponding donor TE sequence were
497 filtered out. Presence or absence of these putative non-reference TE insertions was verified in other accessions
498 by relaxing the parameters used to detect them in the first place. Specifically, we asked for the presence of a
499 minimum of two rather than ten soft-clipped reads spanning the corresponding TSD coordinates. This improved
500 the discovery of putative TE-insertions that are shared by more than one accession.

501

502 **Estimation of sensitivity.** The rate of false negatives was estimated indirectly by first searching for TE
503 insertions with TSDs that are present in the Col-0 reference genome but absent from the Ler-1 assembled
504 genome sequence recently obtained using PacBio long-reads (unpublished). Specifically, annotated TE
505 sequences together with 1kb of upstream and downstream sequences were aligned using BLAT to the Ler-1
506 genome sequence. Heavily truncated TE sequences (covering less than 50% of full-length copies) were not
507 considered, as most of these are unlikely to have been mobilized recently. Out of the 7931 TE sequences
508 analyzed, 200 are not present in the Ler-1 genome sequence and among these, 142 have TSDs (Table S2). This
509 value was then compared to that obtained by mapping Col-0 Illumina short split-reads to the Ler-1 genome
510 sequence. A total of 64 TE insertions with TSDs could not be detected using that approach, thus giving a false
511 negative rate of 45% (Figure 1-figure supplement 1C). Furthermore, 15 TE insertions were called that are in
512 fact not present in the Col-0 genome sequence, thus giving a false positive rate of 10%; (i.e. FDR of 15.3%;
513 Figure 1-figure supplement 1C).

514

515 **Manual filtering of residual false positives.** We analyzed two re-sequenced Col-0 genomes (Schmitz et al.
516 2013; Long et al. 2013) using our split-read pipeline. Non-reference TE insertions with TSDs detected in other
517 accessions and also identified in at least one of the two Col-0 genomes were filtered out, as they are most likely
518 errors. Additionally, after manual inspection using the Integrative Genomics Viewer (IGV) (Robinson et al.
519 2011), we eliminated all putative non-reference TE insertions with TSDs for which read-clusters comprise a
520 disproportionately high number of reads (more than three MADs over the median whole genome coverage) as
521 well as all putative non-reference TE insertions with TSDs that are present in 50% of accessions, as these are
522 frequently the result of sequencing errors or mapping artifacts. Indeed, the pattern of 5' and 3' split-reads at the

523 insertion site differs radically between false and *bona fide* non-reference TE insertions, which enable their
524 discrimination by visual inspection (Figure 1-figure supplement 2). We also observed a high rate of non-
525 reference TE insertions with TSDs within TE sequences, mostly in the case of the ATHILA and GYPSY
526 families. Visual inspection of all these insertions indicates that they most likely result from mapping artifacts
527 and they were therefore excluded manually. Finally, we examined 12 non-reference TE insertions with TSDs by
528 PCR and all were validated (Supplementary file 3).

529

530 **TE sequence capture.**

531 A complete description of our TE-capture design and protocol will be published elsewhere. Briefly 2.1M
532 biotinylated capture probes (Roche-NimbleGen) were designed to cover 200bp at each end of 310 potentially
533 mobile TEs (Supplementary file 4), belonging to 181 of the 326 TE families annotated in TAIR10. These 181
534 TE families include 117 of the 131 mobile TE families identified using the split-read approach as well as most
535 other TE families for which non-degenerate and thus potentially mobile copies are present in the Col-0 genome.
536 Mobile or potentially mobile TE families with skewed sequence composition or very high copy number
537 (including all ATHILA families) were excluded from the design in order to avoid overt biases in the capture of
538 TE sequences. Twelve DNA sequence libraries for the accessions Pi-0, Bl-1, Sp-0, Ta-0, Lip-0, Bla-1, Stw-0,
539 Cvi-0, Kondara, Jm-0, Gre-0 and Rubezhnoe-1 were prepared following the Illumina TruSeq paired-end kit,
540 with some modifications. One μg of genomic DNA was sheared to a median fragment size of 400-450bp.
541 Fragments were subjected to end repair, A-tailing and index adapter ligation following the manufacturer's
542 instructions. Libraries were then amplified through 7 cycles of ligation-mediated PCR using the KAPA HiFi
543 Hot Start Ready Mix and primers AATGATACGGCGACCACCGAGA and
544 CAAGCAGAAGACGGCATAACGAG at a final concentration of $2\mu\text{M}$. PCR products were cleaned-up using
545 the Qiaquick PCR Purification kit following the manufacturer's instructions, except that DNA was eluted in
546 PCR grade water. Amplified DNA libraries were then pooled and one μg of the multiplex pool was used in the
547 first hybridization step to capture probes (72 hours). Captured DNA was recovered using Streptavidin
548 Dynabeads, washed and PCR amplified (5 cycles) as above. Amplified captured DNA was then subjected to a
549 second round of hybridization (12 hours), recovery and amplification (14 cycles). Hybridization, captured
550 DNA recovery and washing were performed as described in the NimbleGen SeqCapture EZ protocol. Pre-
551 capture and post-second capture PCR products were run on an Agilent Bioanalyser DNA 1000 chip.

552 Enrichment for capture TE sequences was confirmed by qPCR and estimated at 50-100 fold depending on the
553 TE. Pair-end sequencing was performed with Illumina HiSeq 2000 and 100bp reads. Between 10.7 and 16.1
554 million pairs were sequenced per library. Pairs were mapped to the TAIR10 reference genome using Bowtie2
555 with the arguments --mp 13 --rdg 8,5 --rfg 8,5 --very-sensitive -X 1000. Enrichment factor was ~250 fold on
556 average (Figure 1-figure supplement 1E). Discordant pairs were remapped using Bowtie2 with the parameters -
557 -mp 13 --rdg 8,5 --rfg 8,5 -D 20 -R 10 -N 0 -L 15 -i S,1,0.50 -k 100. TE insertions were detected in each
558 accession using discordantly mapped reads and the algorithm Hydra (Quinlan et al. 2010) using the arguments:
559 pairDiscordants: -d 50000 -z 1000 -x0 -r 100; dedupDiscordants: -s3; hydra: -ms 50 -li -use all -mld (10*mad)
560 -mno (median+(20*mad). Comparing the results of TE-capture with those obtained using our split-read
561 pipeline gives respectively false negative and positive rates of 56% and 7% (i.e. FDR of 14%; Figure 1-figure
562 supplement 1E). These values are similar to those obtained using the assembled genome sequence of Ler-1 and
563 confirm the low sensitivity of the split-read pipeline.

564

565 **GWAS of CNVs.** Unlike for the initial CNV analysis, we considered this time only the annotated TE sequences
566 that are at least half the size of the corresponding reference TE sequence (obtained from the *A. thaliana*
567 RepeatMasker repeat library, <http://repeatmasker.org>) as this improved significantly the correlation between
568 CNV and the number of non-reference TE insertions. For each of the 113 families identified as mobile by both
569 the split-reads approach and TE-capture, aggregated CNV was obtained for each accession by summing the
570 CNVs estimated for each annotated copy. GWAS on this CNV values was carried out with imputed SNPs
571 (MAF > 5%) from the 211 accessions collected worldwide (Schneeberger et al. 2011; Schmitz et al. 2013), the
572 180 Swedish accessions (Long et al. 2013) and the joined dataset using a linear mixed model (LMM) (Kang et
573 al. 2010). Kinship matrix was included in the model as a random effect to control population structure. Out of
574 the 113 TE families, 33 displayed a genome inflation factor (GIF) greater than 1.10 and were excluded from
575 subsequent analysis. A conservative threshold value (P -value $< 1 \times 10^{-8}$) was set to call statistically associated
576 SNPs. Proximal associated SNPs in linkage disequilibrium ($r^2 > 0.2$) were identified using Plink (Purcell et al.
577 2007) and combined in blocks to build statistical associated intervals, which were expanded by 1 or 5kb on
578 either side. Variance explained by the leading SNP within each locus was calculated using the following
579 equation: $\sigma_{SNP}^2 = 2(MAF)(1 - MAF) \frac{\beta^2}{\sigma_y^2}$, in which MAF is the minor allele frequency, β is the SNP effect

580 estimated by the LMM and σ_y^2 is the variance of the phenotype Y. Total variance explained by *cis* and *trans* loci
581 was computed as the sum of the single-locus explained variances, under the assumption of additive
582 contributions. The probability P(f) of missing a non-reference TE insertion with TSD underlying a “false” *trans*
583 locus through both the split reads and TE-capture approaches was calculated using the following equation:
584 $P(f) = FN \times \binom{n}{12-k} \times f^{12-k} \times (1-f)^k$ where f is the frequency of the minor allele for the *trans* locus under
585 consideration (reported in the Supplementary File 1), FN is the false negative rate for the split-reads
586 bioinformatic pipeline (0.56), n is the number of accessions analyzed by TE-capture (12) and k is the number of
587 accessions without the non-reference TE insertion with TSD among the 12 accessions analyzed by TE-capture
588 (k ranges from 0 to 12 and calculations were all performed using k=12). .

589

590 **Climate analysis.** We selected 12 geo-climatic variables representing different ecological layers: Aridity,
591 number of frosty days, number of consecutive frost-free days, day length in the spring, maximum temperature
592 in the warmest month, minimum temperature in the coldest month, temperature annual range,
593 photosynthetically active radiation, precipitation in the wettest month, precipitation in the driest month, relative
594 humidity in the summer, landscape slope and thermal (Hancock et al. 2011). CNs for the 113 families
595 confirmed as mobile by the split-read approach and TE-capture were used to calculate a partial Mantel
596 correlation with the 13 geo-climatic variables. Population structure kinship was included in the test to control
597 population stratification. Partial Mantel tests were conducted using the “ecodist” package in R. A threshold
598 of $P < 0.01$, corrected for multiple testing (8.33×10^{-4}), was set to call statistically correlated variables. In
599 addition to using the partial Mantel test, we also applied linear models to regress CN as a function of the
600 climatic variables. Although this method does not control for population structure, it largely confirmed the
601 associations found by the partial Mantel test.

602

603 **Characterization of *cis* and *trans* associations.** Protein coding genes, miRNA genes, ncRNAs and TE
604 annotations were retrieved from
605 (ftp://ftp.arabidopsis.org/home/tair/Genes/TAIR10_genome_release/TAIR10_gff3/TAIR10_GFF3_genes.gff).

606 GWAS intervals were defined as *cis* associations if they overlap with a TE annotation or a non-reference TE
607 insertion with TSDs of the same family. All other GWAS intervals were defined as *trans* associations. The *cis*
608 associations are over-represented 34 times ($P\text{-value} < 1 \times 10^{-16}$) when compared to randomly chosen genomic

609 intervals of the same size, which is consistent with TE activity being primarily determined by the TE sequence
610 itself. All genomic annotations overlapping with *trans* intervals (within 1 or 5kb) were considered as being
611 putatively causal.

612

613 **Analysis of the localization of non-reference TE insertions with TSDs along chromosomes.** To assess if
614 non-reference TE insertions with TSDs are enriched in pericentromeric regions, their number within these
615 regions was compared with that expected from a random distribution. Insertions with different allele counts
616 (private, shared by 2-10 accessions, shared by >10 accessions) were considered separately. TE distribution in
617 the reference genome (TAIR10) was obtained by counting the number of TE sequences located within
618 pericentromeres (minus genomic regions showing aberrant coverage and inner pericentromeres). The expected
619 distribution for the 2835 non-reference TE insertions with TSDs was calculated by randomizing 10^6 times their
620 position across the chromosomes (genomic regions showing coverage deviation, the inner pericentromeres, or
621 coordinates spanning the corresponding donor TE sequence were excluded). This set of random positions was
622 used as a control for all subsequent analyses. Insertion distribution over genes and neighboring sequences was
623 performed using a meta-gene. Briefly, protein coding gene features were extracted from the TAIR10 annotation
624 and coordinates of non-reference TE insertions with TSDs were crossed with the set of genic features according
625 to the following stepwise hierarchy: 5' UTR > 3' UTR > exon > intron > intergenic regions. For insertions that
626 do not overlap protein-coding genes, the distance to the closest gene was calculated and reported as negative or
627 positive distance according to the gene orientation. Expected insertion distribution under the null hypothesis
628 was retrieved by applying this procedure for each of the 10^6 randomized sets of insertions. To assess if non-
629 reference TE insertions with TSDs are enriched within small clusters, we divided the genome into 10kb non-
630 overlapping windows and we counted the number of insertions events within them. The observed and expected
631 (random) densities of non-reference TE insertions with TSDs per window were compared and significant
632 enrichment was declared when the number of insertions found within a window was in the upper 0.005% tail of
633 the random distribution.

634

635 **Reconstruction of the historical recombination landscape.** Historical recombination was estimated using
636 LDhat (McVean et al. 2004) as described before (Choi et al. 2013). Briefly, biallelic SNPs ($MAF \geq 0.1$) from
637 the 210 accessions collected worldwide were selected and split into blocks of 5,000 SNPs with overlap of 500

638 SNPs. SNPs located in inner centromeres were excluded from the analysis. Blocks of SNPs were formatted
639 using the ‘convert’ program. A likelihood lookup table was generated for 210 individuals with program
640 ‘complete’ using the following parameters: -n 210 -rhomax 100 -n_pts 100 -theta 0.001. Population-scaled
641 recombination rate (ρ /kb, $\rho = 4Ner$, where N_e is the effective population size and r is the per-generation
642 recombination rate) was calculated using the ‘interval’ program with the following parameters: -its 60000000 -
643 bpen 5 -sam 40000. Recombination rates for contiguous blocks were joined at overlap position 250.
644 Population-scaled recombination rate map is provided in Supplementary file 5.

645

646 **Identification of DNA sequence motifs overrepresented at non-reference TE insertion sites with TSDs.**

647 Sequence spanning non-reference insertion sites were analyzed using Bioprospector Release 2 (Liu, Brutlag,
648 and Liu 2001) using the following parameters: -r 1 -n 200 -a 1 -W ‘TSD size’ . Background sequence
649 distribution for the reference genome was obtained using the “genomebg” program. Sequence logo was
650 produced using the seqLogo package version 1.36.0.

651

652 **Assessing the impact of non-reference TE insertions with TSDs on gene expression.** All genes with a non-
653 reference TE insertion with TSDs within 1kb were retrieved and their expression analyzed using transcriptome
654 data available for 144 accessions (Schmitz et al. 2013). For each gene, we calculated the ratio between the
655 median gene expression level for the accessions harboring the TE insertion and the median gene expression
656 level for accessions lacking that insertion. Distribution plots of observed gene expression ratios were compared
657 to the expected distribution under the null hypothesis (random effect). This expected distribution was obtained
658 by calculating the gene expression ratio for 10^6 randomly chosen sets of genes for which the TE insertion
659 presence/absence “label” was randomly assigned between the accessions. Statistical significance of differences
660 between the observed and expected distributions was determined using the Kolmogorov-Smirnov test.

661

662 **Expression levels of genes affected by non-reference insertions with TSDs.** RNA was extracted using the
663 RNeasy plant mini kit (Qiagen from plants grown under normal conditions (10 days old seedlings grown in
664 liquid medium) or subjected to heat shock treatment (Ito et al. 2011). RT-qPCR was performed as described
665 previously (Silveira et al. 2013). Primers details are given in Supplementary file 3. RT-qPCR results (one
666 biological replicate only) are indicated relative to those obtained for a gene (*AT5G13440*) that shows invariant

667 expression under multiple conditions.

668

669 **Gene expression profile of genes affected by long-distance DNA methylation.** Expression data in the Col-0

670 accession were obtained from <http://www.weigelworld.org/resources/>

671 microarray/AtGenExpress/AtGE_dev_gcRMA.txt.zip/at_download/file for the 224 and 162 genes affected

672 respectively by short- and long-distance DNA methylation in accessions with non-reference TE insertions with

673 TSDs (by definition, these TE insertions are absent in Col-0). Triplicate data for each developmental time point

674 was averaged and then normalized across the developmental time-point series. Average expression level was

675 then calculated for each time point for all genes affected by short-distance DNA methylation and compared to

676 the average calculated for all genes affected by long-distance DNA methylation. Statistical significance of

677 differences between these two averages was calculated using the non-parametric Mann-Whitney U test.

678

679 **Determination of *FLC* haplotypes.** Haplotype analysis was performed as described previously (Li et al. 2014).

680 Briefly, SNPs within 100 kb of *FLC* were retrieved for the 211 worldwide accessions and used as input into

681 fastPHASE version 1.4.0 (Scheet and Stephens 2006). Default parameters were kept, except for the -Pzp option.

682 For each SNP, haplotype membership with the highest likelihood was assigned.

683

684 **Code availability.** Source code for the split-read pipeline can be accessed at

685 <https://github.com/LeanQ/SPLITREADER>

686

687

688 REFERENCES AND NOTES

689

690 Ahmed, Ikhlaq, Alexis Sarazin, Chris Bowler, Vincent Colot, and Hadi Quesneville. 2011. "Genome-Wide
691 Evidence for Local DNA Methylation Spreading from Small RNA-Targeted Sequences in Arabidopsis."
692 *Nucleic Acids Research* 39 (16): 6919–31. doi:10.1093/nar/gkr324.

693 Baillie, J Kenneth, Mark W Barnett, Kyle R Upton, Daniel J Gerhardt, Todd A. Richmond, Fioravante De Sapio,
694 Paul M. Brennan, et al. 2011. "Somatic Retrotransposition Alters the Genetic Landscape of the Human
695 Brain." *Nature* 479 (7374): 534–37. doi:10.1038/nature10531.

696 Barrón, Maite G., Anna-Sophie Fiston-Lavier, Dmitri a. Petrov, and Josefa González. 2014. "Population
697 Genomics of Transposable Elements in Drosophila." *Annual Review of Genetics* 48 (1): 561–81.
698 doi:10.1146/annurev-genet-120213-092359.

699 Beck, Christine R, José Luis Garcia-Perez, Richard M Badge, and John V Moran. 2011. "LINE-1 Elements in
700 Structural Variation and Disease." *Annu Rev Genomics Hum Genet* 12: 187–215. doi:10.1146/annurev-

- 701 genom-082509-141802.
- 702 Becker, Claude, Jörg Hagmann, Jonas Müller, Daniel Koenig, Oliver Stegle, Karsten Borgwardt, and Detlef
703 Weigel. 2011. "Spontaneous Epigenetic Variation in the *Arabidopsis thaliana* Methylome." *Nature* 480
704 (480): 245–52. doi:10.1038/nature10555.
- 705 Bennetzen, Jeffrey L, and Hao Wang. 2014. "The Contributions of Transposable Elements to the Structure,
706 Function, and Evolution of Plant Genomes." *Annual Review of Plant Biology* 65. Annual Reviews: 505–
707 30. doi:10.1146/annurev-arplant-050213-035811.
- 708 Brachi, Benjamin, Geoffrey P Morris, and Justin O Borevitz. 2011. "Genome-Wide Association Studies in
709 Plants: The Missing Heritability Is in the Field." *Genome Biology* 12 (10): 232. doi:10.1186/gb-2011-12-
710 10-232.
- 711 Cao, Jun, Korbinian Schneeberger, Stephan Ossowski, Torsten Günther, Sebastian Bender, Joffrey Fitz, Daniel
712 Koenig, et al. 2011. "Whole-Genome Sequencing of Multiple *Arabidopsis thaliana* Populations." *Nature*
713 *Genetics* 43 (10): 956–63. doi:10.1038/ng.911.
- 714 Cavrak, Vladimir V, Nicole Lettner, Suraj Jamge, Agata Kosarewicz, Laura Maria Bayer, and Ortrun Mittelsten
715 Scheid. 2014. "How a Retrotransposon Exploits the Plant's Heat Stress Response for Its Activation."
716 *PLoS Genet* 10 (1): e1004115. doi:10.1371/journal.pgen.1004115.
- 717 Chae, Eunyong, Kirsten Bomblies, Sang-Tae Kim, Darya Karelina, Maricris Zaidem, Stephan Ossowski,
718 Carmen Martín-Pizarro, et al. 2014. "Species-Wide Genetic Incompatibility Analysis Identifies Immune
719 Genes as Hot Spots of Deleterious Epistasis." *Cell* 159: 1341–51. doi:10.1016/j.cell.2014.10.049.
- 720 Choi, Kyuha, Xiaohui Zhao, Krystyna a Kelly, Oliver Venn, James D Higgins, Nataliya E Yelina, Thomas J
721 Hardcastle, et al. 2013. "Arabidopsis Meiotic Crossover Hot Spots Overlap with H2A.Z Nucleosomes at
722 Gene Promoters." *Nature Genetics* 45 (11): 1327–36. doi:10.1038/ng.2766.
- 723 Clark, Richard M., Gabriele Schweikert, Christopher Toomajian, Stephan Ossowski, Georg Zeller, Paul Shinn,
724 Norman Warthmann, et al. 2007. "Common Sequence Polymorphisms Shaping Genetic Diversity in
725 *Arabidopsis thaliana*." *Science* 317: 338–42. doi:10.1126/science.1138632.
- 726 Cokus, Shawn J, Suhua Feng, Xiaoyu Zhang, Zugen Chen, Barry Merriman, Christian D Haudenschild,
727 Sriharsa Pradhan, Stanley F Nelson, Matteo Pellegrini, and Steven E Jacobsen. 2008. "Shotgun Bisulphite
728 Sequencing of the Arabidopsis Genome Reveals DNA Methylation Patterning." *Nature* 452 (7184): 215–
729 19. doi:10.1038/nature06745.
- 730 Cridland, Julie M., Stuart J. Macdonald, Anthony D. Long, and Kevin R. Thornton. 2013. "Abundance and
731 Distribution of Transposable Elements in Two Drosophila QTL Mapping Resources." *Molecular Biology*
732 *and Evolution* 30 (10): 2311–27. doi:10.1093/molbev/mst129.
- 733 Daxinger, Lucia, Tatsuo Kanno, Etienne Bucher, Johannes van der Winden, Ulf Naumann, Antonius J M
734 Matzke, and Marjori Matzke. 2009. "A Stepwise Pathway for Biogenesis of 24-Nt Secondary siRNAs and
735 Spreading of DNA Methylation." *The EMBO Journal* 28 (1): 48–57. doi:10.1038/emboj.2008.260.
- 736 Dubin, Manu J, Pei Zhang, Dazhe Meng, Marie-Stanislas Remigereau, Edward J Osborne, Francesco Paolo
737 Casale, Philipp Drewe, et al. 2015. "DNA Methylation in Arabidopsis Has a Genetic Basis and Shows
738 Evidence of Local Adaptation." *eLife* 4: e05255. doi:10.7554/eLife.05255.
- 739 Elbarbary, Reyad A., Bronwyn A. Lucas, and Lynne E. Maquat. 2016. "Retrotransposons as Regulators of Gene
740 Expression." *Science* 351 (6274): aac7247. doi:10.1126/science.aac7247.
- 741 Fahlgren, Noah, Christopher M Sullivan, Kristin D Kasschau, Elisabeth J Chapman, Tyler W H Backman, Scott
742 A Givan, and James C Carrington. 2009. "Computational and Analytical Framework for Small RNA
743 Profiling by High-Throughput Sequencing." *RNA*, no. 15: 992–1002. doi:10.1261/rna.1473809.
- 744 François, Olivier, Michael G. B. Blum, Mattias Jakobsson, and Noah A. Rosenberg. 2008. "Demographic
745 History of European Populations of *Arabidopsis thaliana*." *PLoS Genetics* 4 (5): e1000075.
746 doi:10.1371/journal.pgen.1000075.
- 747 Fu, Yu, Akira Kawabe, Mathilde Etcheverry, Tasuku Ito, Atsushi Toyoda, Asao Fujiyama, Vincent Colot,
748 Yoshiaki Tarutani, and Tetsuji Kakutani. 2013. "Mobilization of a Plant Transposon by Expression of the

- 749 Transposon-Encoded Anti-Silencing Factor.” *The EMBO Journal* 32 (17): 2407–17.
750 doi:10.1038/emboj.2013.169.
- 751 Hancock, Angela M, Benjamin Brachi, Nathalie Faure, Matthew W Horton, Lucien B Jarymowycz, F Gianluca
752 Sperone, Chris Toomajian, Fabrice Roux, and Joy Bergelson. 2011. “Adaptation to Climate Across the
753 *Arabidopsis thaliana* Genome.” *Science* 334 (6052): 83–86. doi:10.1126/science.1209244.
- 754 Heard, Edith, and Robert A. Martienssen. 2014. “Transgenerational Epigenetic Inheritance: Myths and
755 Mechanisms.” *Cell* 157 (1): 95–109. doi:10.1016/j.cell.2014.02.045.
- 756 Hénaff, Elizabeth, Luís Zapata, Josep M. Casacuberta, and Stephan Ossowski. 2015. “Jitterbug: Somatic and
757 Germline Transposon Insertion Detection at Single-Nucleotide Resolution.” *BMC Genomics* 16 (1): 768.
758 doi:10.1186/s12864-015-1975-5.
- 759 Hill, W G, and Alan Robertson. 1966. “The Effect of Linkage on Limits to Artificial Selection.” *Genetics
760 Research* 8 (03): 269–94. doi:10.1017/S0016672300010156.
- 761 Hollister, Jesse D, and Brandon S Gaut. 2009. “Epigenetic Silencing of Transposable Elements: A Trade-off
762 between Reduced Transposition and Deleterious Effects on Neighboring Gene Expression.” *Genome
763 Research* 19 (8): 1419–28. doi:10.1101/gr.091678.109.
- 764 Hu, Tina T, Pedro Pattyn, Erica G Bakker, Jun Cao, Jan-Fang Cheng, Richard M Clark, Noah Fahlgren, et al.
765 2011. “The *Arabidopsis lyrata* Genome Sequence and the Basis of Rapid Genome Size Change.” *Nature
766 Genetics* 43 (5): 476–81. doi:10.1038/ng.807.
- 767 Ietswaart, Robert, Zhe Wu, and Caroline Dean. 2012. “Flowering Time Control: Another Window to the
768 Connection between Antisense RNA and Chromatin.” *Trends in Genetics* 28 (9): 445–53.
769 doi:10.1016/j.tig.2012.06.002.
- 770 Ito, Hidetaka, Etienne Bucher, Marie Mirouze, Isabelle Vaillant, and Jerzy Paszkowski. 2011. “An siRNA
771 Pathway Prevents Transgenerational Retrotransposition in Plants Subjected to Stress.” *Nature* 472 (7341):
772 115–19. doi:10.1038/nature09861.
- 773 Ito, Hidetaka, and Tetsuji Kakutani. 2014. “Control of Transposable Elements in *Arabidopsis thaliana*.”
774 *Chromosome Research* 22 (2): 217–23. doi:10.1007/s10577-014-9417-9.
- 775 Joly-Lopez, Zoé, and Thomas E. Bureau. 2014. “Diversity and Evolution of Transposable Elements in
776 *Arabidopsis*.” *Chromosome Research* 22 (2): 203–16. doi:10.1007/s10577-014-9418-8.
- 777 Jurka, Jerzy, Vladimir V Kapitonov, A Pavlicek, P Klonowski, O Kohany, and J Walichiewicz. 2005. “Repbse
778 Update, a Database of Eukaryotic Repetitive Elements.” *Cytogenetic and Genome Research* 110 (1-4):
779 462–67.
- 780 Kang, Hyun Min, Jae Hoon Sul, Susan K. Service, Noah A. Zaitlen, Sit-yeek Kong, Nelson B. Freimer, Chiara
781 Sabatti, and Eleazar Eskin. 2010. “Variance Component Model to Account for Sample Structure in
782 Genome-Wide Association Studies.” *Nature Genetics* 42 (4): 348–54. doi:10.1038/ng.548.
- 783 Kanno, Tatsuo, Etienne Bucher, Lucia Daxinger, Bruno Huettel, Gudrun Böhmendorfer, Wolfgang Gregor, David
784 P Kreil, Marjori Matzke, and Antonius J M Matzke. 2008. “A Structural-Maintenance-of-Chromosomes
785 Hinge Domain-containing Protein Is Required for RNA-Directed DNA Methylation.” *Nature Genetics* 40
786 (5): 670–75. doi:10.1038/ng.119.
- 787 Kazazian, Haig H. 2004. “Mobile Elements: Drivers of Genome Evolution.” *Science* 303 (5664): 1626–32.
788 doi:10.1126/science.1089670.
- 789 Klein, Christopher J, Maria-Victoria Botuyan, Yanhong Wu, Christopher J Ward, Garth a Nicholson, Simon
790 Hammans, Kaori Hojo, et al. 2011. “Mutations in DNMT1 Cause Hereditary Sensory Neuropathy with
791 Dementia and Hearing Loss.” *Nature Genetics* 43 (6): 595–600. doi:10.1038/ng.830.
- 792 Kofler, Robert, Viola Nolte, and Christian Schlötterer. 2015. “Tempo and Mode of Transposable Element
793 Activity in *Drosophila*.” *PLoS Genetics* 11 (7): e1005406. doi:10.1371/journal.pgen.1005406.
- 794 Law, Julie A., and Steven E. Jacobsen. 2010. “Establishing , Maintaining and Modifying DNA Methylation
795 Patterns in Plants and Animals.” *Nature Reviews Genetics* 11 (3): 204–20. doi:10.1038/nrg2719.

- 796 Lempe, Janne, Sureshkumar Balasubramanian, Sridevi Sureshkumar, Anandita Singh, Markus Schmid, and
797 Detlef Weigel. 2005. "Diversity of Flowering Responses in Wild *Arabidopsis thaliana* Strains." *PLoS*
798 *Genetics* 1 (1): e6. doi:10.1371/journal.pgen.0010006.
- 799 Li, Peijin, Daniele Filiault, MS Mathew S Box, Envel Kerdaffrec, Cock van Oosterhout, Amity M Wilczek,
800 Johanna Schmitt, et al. 2014. "Multiple FLC Haplotypes Defined by Independent Cis-Regulatory
801 Variation Underpin Life History Diversity in *Arabidopsis thaliana*." *Genes & Development* 28 (15):
802 1635–40. doi:10.1101/gad.245993.114.
- 803 Li, Y., Y. Huang, J. Bergelson, M. Nordborg, and J. O. Borevitz. 2010. "Association Mapping of Local Climate-
804 Sensitive Quantitative Trait Loci in *Arabidopsis thaliana*." *Proceedings of the National Academy of*
805 *Sciences* 107 (49): 21199–204. doi:10.1073/pnas.1007431107.
- 806 Lippman, Zachary, Anne-Valerie Gendrel, Michael Black, Matthew W Vaughn, Neilay Dedhia, W Richard
807 McCombie, Kimberly Lavine, et al. 2004. "Role of Transposable Elements in Heterochromatin and
808 Epigenetic Control." *Nature* 430 (6998): 471–76. doi:10.1038/nature02651.
- 809 Lisch, Damon. 2013. "How Important Are Transposons for Plant Evolution?" *Nature Reviews Genetics* 14 (1):
810 49–61. doi:10.1038/nrg3374.
- 811 Lister, Ryan, Ronan C O Malley, Julian Tonti-filippini, Brian D Gregory, Charles C Berry, A Harvey Millar, and
812 Joseph R Ecker. 2008. "Highly Integrated Single-Base Resolution Maps of the Epigenome in *Arabidopsis*." *Cell*
813 133 (3): 523–36. doi:10.1016/j.cell.2008.03.029.
- 814 Liu, Jun, Yuehui He, Richard Amasino, and Xuemei Chen. 2004. "siRNAs Targeting an Intronic Transposon in
815 the Regulation of Natural Flowering Behavior in *Arabidopsis*." *Genes & Development*, no. 732: 2873–78.
816 doi:10.1101/gad.1217304.
- 817 Liu, X, D L Brutlag, and J S Liu. 2001. "Bioprospector: Discovering Conserved DNA Motifs in Upstream
818 Regulatory Regions of Co-Expressed Genes." *Pacific Symposium on Biocomputing* 6: 127–38.
- 819 Lockton, Steven, and Brandon S Gaut. 2010. "The Evolution of Transposable Elements in Natural Populations
820 of Self-Fertilizing *Arabidopsis thaliana* and Its Outcrossing Relative *Arabidopsis lyrata*." *BMC*
821 *Evolutionary Biology* 10: 10. doi:10.1186/1471-2148-10-10.
- 822 Long, Quan, Fernando a Rabanal, Dazhe Meng, Christian D Huber, Ashley Farlow, Alexander Platzer, Qingrun
823 Zhang, et al. 2013. "Massive Genomic Variation and Strong Selection in *Arabidopsis thaliana* Lines from
824 Sweden." *Nature Genetics* 45 (8): 884–90. doi:10.1038/ng.2678.
- 825 Mackay, Trudy F. C., Stephen Richards, Eric a. Stone, Antonio Barbadilla, Julien F. Ayroles, Dianhui Zhu,
826 Sònia Casillas, et al. 2012. "The *Drosophila melanogaster* Genetic Reference Panel." *Nature* 482 (7384):
827 173–78. doi:10.1038/nature10811.
- 828 Marí-Ordóñez, Arturo, Antonin Marchais, Mathilde Etcheverry, Antoine Martin, Vincent Colot, and Olivier
829 Voinnet. 2013. "Reconstructing de Novo Silencing of an Active Plant Retrotransposon." *Nature Genetics*
830 45 (9): 1029–39. doi:10.1038/ng.2703.
- 831 Maumus, Florian, and Hadi Quesneville. 2014. "Ancestral Repeats Have Shaped Epigenome and Genome
832 Composition for Millions of Years in *Arabidopsis thaliana*." *Nature Communications* 5 (May): 4104.
833 doi:10.1038/ncomms5104.
- 834 McClintock, Barbara. 1956. "Intranuclear Systems Controlling Gene Action and Mutation." In *Brookhaven*
835 *Symposia in Biology*, 58–74.
- 836 McVean, Gilean A T, Simon R Myers, Sarah Hunt, Panos Deloukas, David R Bentley, and Peter Donnelly. 2004.
837 "The Fine-Scale Structure of Recombination Rate Variation in the Human Genome." *Science* 304 (5670):
838 581–84. doi:DOI: 10.1126/science.1092500.
- 839 Miyao, Akio, Katsuyuki Tanaka, Kazumasa Murata, Hiromichi Sawaki, Shin Takeda, Kiyomi Abe, Yoriko
840 Shinozuka, Katsura Onosato, and Hirohiko Hirochika. 2003. "Target Site Specificity of the Tos17
841 Retrotransposon Shows a Preference for Insertion Withn Genes and against Insertion in Retrotransposon-
842 Rich Regions of the Genome." *The Plant Cell* 15 (August): 1771–80. doi:10.1105/tpc.012559.ements.
- 843 Purcell, Shaun, Benjamin Neale, Kathe Todd-Brown, Lori Thomas, Manuel A R Ferreira, David Bender, Julian

- 844 Maller, Pamela Sklar, Paul I W De Bakker, and Mark J Daly. 2007. "PLINK: A Tool Set for Whole-
845 Genome Association and Population-Based Linkage Analyses." *The American Journal of Human Genetics*
846 81 (3): 559–75. doi:10.1086/519795.
- 847 Quadrana, Leandro, and Vincent Colot. 2016. "Plant Transgenerational Epigenetics." *Annual Review of*
848 *Genetics* 50 (1): null. doi:10.1146/annurev-genet-120215-035254.
- 849 Quinlan, Aaron R, Royden a Clark, Svetlana Sokolova, Mitchell L Leibowitz, Yujun Zhang, Matthew E Hurles,
850 Joshua C Mell, and Ira M Hall. 2010. "Genome-Wide Mapping and Assembly of Structural Variant
851 Breakpoints in the Mouse Genome." *Genome Research* 20 (5): 623–35. doi:10.1101/gr.102970.109.
- 852 Rahman, Reazur, Gung-Wei Chirn, Abhay Kanodia, Yuliya A. Sytnikova, Björn Brembs, Casey M. Bergman,
853 and Nelson C. Lau. 2015. "Unique Transposon Landscapes Are Pervasive across *Drosophila*
854 *melanogaster* Genomes." *Nucleic Acids Research* 43 (22): 10655–72. doi:10.1093/nar/gkv1193.
- 855 Rebollo, Rita, Mark T Romanish, and Dixie L Mager. 2012. "Transposable Elements: An Abundant and Natural
856 Source of Regulatory Sequences for Host Genes." *Annual Review of Genetics* 46 (January): 21–42.
857 doi:10.1146/annurev-genet-110711-155621.
- 858 Richardson, Sandra R, Santiago Morell, and Geoffrey J Faulkner. 2014. "L1 Retrotransposons and Somatic
859 Mosaicism in the Brain." *Annual Review of Genetics* 48 (1): 1–27. doi:10.1146/annurev-genet-120213-
860 092412.
- 861 Richardson, Sandra R., Aurélien J. Doucet, Huira C. Kopera, John B. Moldovan, José Luis Garcia-Perez, and
862 John V. Moran. 2015. "The Influence of LINE-1 and SINE Retrotransposons on Mammalian Genomes." *Microbiology Spectrum* 3 (2): 1–26. doi:10.1128/microbiolspec.MDNA3-0061-2014.
- 864 Robinson, James T, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S Lander, Gad Getz, and
865 Jill P Mesirov. 2011. "Integrative Genomics Viewer." *Nature Biotechnology* 29 (1): 24–26.
866 doi:10.1038/nbt0111-24.
- 867 Scheet, Paul, and Matthew Stephens. 2006. "A Fast and Flexible Statistical Model for Large-Scale Population
868 Genotype Data: Applications to Inferring Missing Genotypes and Haplotypic Phase." *The American*
869 *Journal of Human Genetics* 78 (4): 629–44. doi:10.1086/502802.
- 870 Schmitz, Robert J., Matthew D. Schultz, Mathew G. Lewsey, Ronan C. O'Malley, Mark A. Urich, Ondrej
871 Libiger, Nicholas J. Schork, and Joseph R. Ecker. 2011. "Transgenerational Epigenetic Instability Is a
872 Source of Novel Methylation Variants." *Science* 334 (6054). American Association for the Advancement
873 of Science: 369–73. doi:10.1126/science.1212959.
- 874 Schmitz, Robert J., Matthew D. Schultz, Mark A. Urich, Joseph R. Nery, Mattia Pelizzola, Ondrej Libiger,
875 Andrew Alix, et al. 2013. "Patterns of Population Epigenomic Diversity." *Nature* 495 (7440): 193–98.
876 doi:10.1038/nature11968.
- 877 Schneeberger, Korbinian, Stephan Ossowski, Felix Ott, Juliane D Klein, Xi Wang, Christa Lanz, Lisa M Smith,
878 et al. 2011. "Reference-Guided Assembly of Four Diverse *Arabidopsis thaliana* Genomes." *Proceedings*
879 *of the National Academy of Sciences* 108 (25): 10249–54. doi:10.1073/pnas.1107739108.
- 880 Silveira, Amanda Bortolini, Charlotte Trontin, Sandra Cortijo, Joan Barau, Luiz Eduardo Vieira Del Bem,
881 Olivier Loudet, Vincent Colot, and Michel Vincentz. 2013. "Extensive Natural Epigenetic Variation at a
882 de Novo Originated Gene." *PLoS Genetics* 9 (4): e1003437. doi:10.1371/journal.pgen.1003437.
- 883 Slotkin, R Keith, and Robert Martienssen. 2007. "Transposable Elements and the Epigenetic Regulation of the
884 Genome." *Nature Reviews. Genetics* 8 (4): 272–85. doi:10.1038/nrg2072.
- 885 Stroud, Hume, Maxim V.C. Greenberg, Suhua Feng, Yana V. Bernatavichute, and Steven E. Jacobsen. 2013.
886 "Comprehensive Analysis of Silencing Mutants Reveals Complex Regulation of the Arabidopsis
887 Methyloome." *Cell* 152 (1-2): 352–64. doi:10.1016/j.cell.2012.10.054.
- 888 Stuart, Tim, Steven R Eichten, Jonathan Cahn, Justin Borevitz, and Ryan Lister. 2016. "Population Scale
889 Mapping of Transposable Element Diversity Reveal Link to Gene Regulation and Epigenetic Variation."
890 *bioRxiv* 0897: 0–3. doi:10.1101/039511.
- 891 Teixeira, Felipe K, and Vincent Colot. 2010. "Repeat Elements and the Arabidopsis DNA Methylation

- 892 Landscape.” *Heredity* 105 (1). Nature Publishing Group: 14–23. doi:10.1038/hdy.2010.52.
- 893 Teixeira, Felipe K, Fabiana Heredia, Alexis Sarazin, François Roudier, Martine Boccara, Constance Ciaudo,
894 Corinne Cruaud, et al. 2009. “A Role for RNAi in the Selective Correction of DNA Methylation Defects.”
895 *Science* 323 (5921): 1600–1604. doi:10.1126/science.1165313.
- 896 The Arabidopsis Genome Initiative. 2000. “Analysis of the Genome Sequence of the Flowering Plant
897 *Arabidopsis thaliana*.” *Nature* 408 (6814): 796–815. doi:10.1038/35048692.
- 898 Trono, Didier. 2016. “Transposable Elements, Polydactyl Proteins, and the Genesis of Human-Specific
899 Transcription Networks.” *Cold Spring Harbor Symposia on Quantitative Biology* LXXX (ii).
900 doi:10.1101/sqb.2015.80.027573.
- 901 Tsay, Yi Fang, Mary J. Frank, Tania Page, Caroline Dean, and Nigel M. Crawford. 1993. “Identification of a
902 Mobile Endogenous Transposon in *Arabidopsis thaliana*.” *Science* 260 (5106): 342–44.
903 doi:10.1126/science.8385803.
- 904 Vinkhuyzen, Anna a.E., Naomi R. Wray, Jian Yang, Michael E. Goddard, and Peter M. Visscher. 2013.
905 “Estimation and Partition of Heritability in Human Populations Using Whole-Genome Analysis Methods.”
906 *Annual Review of Genetics* 47 (1): 75–95. doi:10.1146/annurev-genet-111212-133258.
- 907 Weigel, Detlef, and Vincent Colot. 2012. “Epialleles in Plant Evolution.” *Genome Biology* 13 (10): 249.
908 doi:10.1186/gb-2012-13-10-249.
- 909 Weigel, Detlef, and Magnus Nordborg. 2015. “Population Genomics for Understanding Adaptation in Wild
910 Plant Species.” *Annual Review of Genetics* 49 (1): 315–38. doi:10.1146/annurev-genet-120213-092110.
- 911 Yoon, Seungtae, Zhenyu Xuan, Vladimir Makarov, Kenny Ye, and Jonathan Sebat. 2009. “Sensitive and
912 Accurate Detection of Copy Number Variants Using Read Depth of Coverage.” *Genome Research* 19 (9):
913 1586–92. doi:10.1101/gr.092981.109.
- 914 Ziolkowski, Piotr A., Grzegorz Koczyk, Lukasz Galganski, and Jan Sadowski. 2009. “Genome Sequence
915 Comparison of Col and Ler Lines Reveals the Dynamic Nature of Arabidopsis Chromosomes.” *Nucleic
916 Acids Research* 37 (10): 3189–3201. doi:10.1093/nar/gkp183.

917
918

919 **Acknowledgements.** The authors declare no competing financial interest, except GFM and JAJ who
920 declare a competing interest as employees of Roche NimbleGen Inc. We thank members of the Colot
921 lab and especially Mathilde Etcheverry for discussions. We thank Edith Heard and Pierre Capy for
922 critical reading of an earlier version of the manuscript. This work was supported by the European
923 Union Seventh Framework Programme Network of Excellence EpiGeneSys (Award 257082, to VC),
924 the Investissements d’Avenir ANR-10-LABX-54 MEMO LIFE, ANR-11-IDEX-0001-02 PSL*
925 Research University and ANR-12-ADAP-0020-01 (to VC) and the Chaire Blaise Pascal (to RAM).
926 LQ was the recipient of postdoctoral fellowships from the ANR-10-LABX-54 MEMO LIFE and
927 ANR-11-IDEX-0001-02 PSL* Research University. ABS was the recipient of postdoctoral
928 fellowships from the ANR-12-ADAP-0020-01 and from the Brazilian National Council for Scientific
929 and Technological Development (CNPq). LQ and VC conceived and designed the study. LQ

930 performed all of the bioinformatic analyses. ABS, GFM and JAJ designed the TE-sequence capture,
931 ABS performed the experiment and CL and RAM performed the sequencing of captured fragments.
932 ABS performed all of the other experimental analyses. LQ and VC wrote the paper, with additional
933 input from all authors. Correspondence and requests for materials should be addressed to
934 colot@biologie.ens.fr.

935

936 FIGURES LEGENDS

937 **Figure 1. Overview of the *A. thaliana* mobilome.** (A) Genome browser tracks showing normalized
938 sequencing coverage over the two full-length *ATCOPIA31* elements annotated in the reference genome (Col-0).
939 CNV is detected as increased or decreased coverage in other accessions. Number of copies is indicated on the
940 right. (B) Heat map representing CNVs (log₂ ratio) for 317 TE families and 211 *A. thaliana* accessions. TE
941 families with statistically significant CNV in at least one accession are indicated. Figure 1-Source data 1
942 contains absolute copy number estimation of TE sequences. (C) Schematic representation of the bioinformatics
943 pipeline to identify non-reference TE insertions with TSD using split-reads. 1- Reads are mapped on a
944 collection of TE extremities from annotated TE sequences and reference sequences (Repbase update). 2- Reads
945 aligning partially over TE extremities are extracted and clipped. 3- The unmapped portion of these split-reads
946 are re-mapped on the Arabidopsis reference genome. 4- Non-reference TE insertions with TSDs are identified
947 by searching for overlapping clusters of 5' and 3' split-reads. (D) Genome browser tracks showing split-reads
948 for two non-reference *ATCOPIA31* insertions and TSD reconstruction. Figure 1-Source data 2 contains the
949 coordinates of all non-reference TE insertions with TSDs. (E) Distribution frequency of allele counts for non-
950 reference TE insertions with TSDs. (F) Number of mobile TE families per accession identified using split-read
951 and TE-sequence capture. (G) Cumulative plot of the number of mobile TE families detected with increasing
952 numbers of accessions. (H) The total number of non-reference TE insertions with TSDs is indicated in relation
953 to the number of accessions with such insertions, for each of the 131 mobile TE families. Asterisks indicate the
954 nine TE families with experimental evidence of transposition (Ito and Kakutani 2014; Tsay et al. 1993). Figure
955 1-Source data 3 contains the total number of distinct non-reference TE insertions with TSD for each TE family
956 and super-family. Figure 1-figure supplement 1 shows TE-capture results. Figure 1- figure supplement 2
957 contains IGV screenshots showing the pattern of split-reads characteristic of true- and false-positive non-

958 reference TE insertions with TSDs.

959

960 **Figure 2. Environmental and genetic factors associated with differential mobilome activity.** (A) Copy
961 number (CN, red circles) of *ATCOPIA78* in accessions distributed across the globe. Annual temperature range
962 is also shown. (B) Partial Mantel correlation between *ATCOPIA78* CN and annual temperature range. (C)
963 Fraction of CNV variance explained by SNPs (*cis*, and *trans*) and partial Mantel correlation with geo-climatic
964 variables. (D) Distribution of *cis* and *trans* loci in the joined analysis (391 accessions) and number of TE
965 families associated with a given *trans* locus. A complete list of the GWAS results is provided in Supplementary
966 file 1. (E) Manhattan plots displaying GWAS results for the seven TE families with a *MET2a* association. The
967 leading SNP within each interval is indicated as a red diamond. Colors indicate the extent of linkage
968 disequilibrium (r^2) with the leading SNP. (F) Schematic view of the MET2a protein (TD: targeting domain;
969 BAH: bromo adjacent homology domain) and sequence alignment of the TD. The amino acid substitution
970 (G519E) that is present in some accessions is indicated (red arrow). (G) Average DNA methylation level over
971 non-mobile, mobile and *MET2a*-associated TE families in WT and *met2a* Col-0 seedlings (Stroud et al. 2013).
972 Statistically significant differences are indicated (MWU test). Figure 2-figure supplement 1 shows the positive
973 correlation between CN and number of non-reference TE insertions with TSDs. Figure 2-figure supplement 2
974 shows climate association to CNVs. Figure 2-figure supplement 3 shows GWAS results for CNVs.

975

976 **Figure 3. Genomic localization of non-reference TE insertions.** (A) Density of non-reference TE insertions
977 with TSDs (blue) and of annotated TE sequences (red) along the reference sequence of chromosome 1. Inner
978 pericentromeric regions are masked. (B) Fraction of private and shared non-reference TE insertions with TSDs
979 and of annotated TE sequences in outer pericentromeric regions. Statistically significant differences are
980 indicated (chi square test). (C) Observed/expected ratio (O/E) of private non-reference TE insertions with TSDs
981 in and around genes. Errors bars are defined as 95% confidence intervals. (D) Cumulative distribution of gene
982 expression ratios between alleles harboring and lacking non-reference TE insertions. Statistically significant
983 differences were calculated using the KS test. (E) As D, but only for COPIA (green) or MuDR (red) non-
984 reference TE insertions with TSDs. Figure 3-figure supplement 1 shows detailed analysis of the distribution of
985 non-reference TE insertions with TSDs along the genome. Figure 3-figure supplement 2 shows local TE
986 insertion preferences. Figure 3-figure supplement 3 shows global expression levels of gene affected by non-

987 reference TE insertions. Figure 3-figure supplement 4 shows expression levels of genes affected in some
988 accessions by a non-reference insertion with TSD in plants grown under control conditions or subjected to heat
989 stress.

990

991 **Figure 4. DNA methylation of non-reference TE insertion sites.** (A) Boxplot representation of average DNA
992 methylation level for mobile and non-mobile TE families across all accessions. (B) O/E ratio of spontaneous
993 DMRs identified in mutation accumulation lines (Becker et al. 2011; Schmitz et al. 2011) for non-mobile and
994 mobile TE families. Statistically significant differences were calculated using a chi square test. (C) Average
995 DNA methylation level in 50bp windows upstream and downstream of 1543 insertions sites for accessions
996 lacking or containing a given non-reference TE insertion with TSD. (D) Genome browser tracks showing
997 examples of insertion sites respectively associated with short- and long-distance DNA methylation. (E) Meta-
998 analysis of DNA methylation around non-reference TE insertions sites. (F) Distribution of non-reference TE
999 insertions associated with short- or long-distance DNA methylation according to their position relative to genes
1000 (stacked bar plot) and proportion of insertions in the two possible orientations relative to the closest gene (pie
1001 charts). (G) Average expression level in different organs and at different developmental time points (in Col-0)
1002 of genes with non-reference TE insertions with TSDs and affected by short- (blue) or long-distance (red) DNA
1003 methylation. Error bars are s.e.m. Statistical significance of differences was calculated using a MWU test.
1004 Figure 4-figure supplement 1 shows DNA methylation of TE families and impact on sequences flanking non-
1005 reference TE insertions with TSDs.

1006

1007 **Figure 5. Local enrichment of non-reference TE insertions with TSDs.** (A) Density of non-reference TE
1008 insertions with TSDs in 10kb windows. The 19 regions statistically enriched in such insertions are indicated by
1009 red bars. (B) Position and identity of the seven non-reference TE insertions with TSDs spanning the *FLC* locus.
1010 (D) and (E) Level of *FLC* expression (D) and flowering time (E) for accessions of same *FLC* haplotype but
1011 differing by the presence or absence of the relevant TE insertion. Errors bars are s.e.m. Figure 5-figure
1012 supplement 1 shows the reconstruction of the *FLC* haplotypes and additional analyses of the effect on flowering
1013 time of non-reference TE insertions with TSDs at the locus.

1014

1015

1016 **Figure 1-figure supplement 1. Validation of the *A. thaliana* mobilome by TE-capture.** (A) The number of
1017 non-reference TE insertions with TSDs identified by the split-read pipeline is plotted against the corresponding
1018 genome sequencing coverage for each accession. Accessions analyzed by TE-capture are highlighted in red. (B)
1019 Genome browser tracks showing examples of non-reference TE insertions identified by TE-capture only. (C)
1020 Overlap between TE insertions with TSDs identified specifically in Col-0 using the Ler-1 genome assembly as
1021 a reference and either whole genome sequence alignment or the split-reads pipeline. The percentage of false
1022 positives (FP), true positives (TP) and false negatives (FN) as well as the false discovery rate (FDR) are
1023 indicated. (D) Description of the TE-capture design and workflow. (E) TE-capture enrichment of target
1024 sequences. (F) Overlap between non-reference insertions with TSDs identified by split-read analysis and TE-
1025 capture. The percentage of FP, TP and FN as well as the FDR are indicated. (G) Distribution frequency of allele
1026 counts for non-reference TE insertions identified using the split-read approach and TE-capture among the 12
1027 accessions analysed. (H) Number of SNPs plotted against the number of non-reference TE insertions identified
1028 by TE-capture between any two accessions.

1029

1030 **Figure 1-figure supplement 2. Visual inspection of true- and false-positive non-reference TE insertions**
1031 **with TSDs.** IGV screenshots showing split-reads for non-reference TE insertions with TSDs that are validated
1032 or not by TE-capture (true- and false-positives, respectively). Split-reads are shown for 12 different accessions.
1033 Accessions containing the non-reference TE insertion with TSD are indicated in red.

1034

1035 **Figure 2-figure supplement 1. Pearson correlation between TE CN and number of TE sequences identified**
1036 **by TE capture.**

1037

1038 **Figure 2-figure supplement 2. Climate association to TE CNV.** Heat map representing partial Mantel
1039 correlation coefficient between TE CN and geo-climatic variables. TE families with statistically significant
1040 correlations ($P < 8.33 \times 10^{-4}$) are indicated.

1041

1042 **Figure 2-figure supplement 3. GWAS of CNVs.** (A) Overlap between GWAS results obtained using CNVs
1043 and SNPs from world-wide accessions and from Swedish accessions. (B) Manhattan plot of the GWAS results
1044 for *ATCOPIA69* CNVs. The leading SNP (red diamond) is located within the TE itself. Colors indicate the

1045 extent of linkage disequilibrium (r^2) to the leading SNP. Distribution of CN values associated with the leading
1046 SNP and the common allele are shown on the right. (C) Schematic overview of *ATCOPIA69*. LTR, long
1047 terminal repeat; gag, nucleocapsid protein; pro, protease. The sequence alignment on the right indicates the
1048 position of the amino acid change (red arrow) caused by the leading SNP in the transposase protein. Conserved
1049 amino acids are highlighted in blue. (D) Average fraction of CN variance explained by *cis* and *trans* loci for
1050 Class I and Class II TE families. (E) Average fraction of CN variance explained by *cis* and *trans* loci for
1051 autonomous and non-autonomous class II TE families. Ratio of the observed over expected (O/E) number of
1052 TE annotations overlapping *trans* loci for autonomous and non-autonomous class II TE families. Statistically
1053 significant differences were calculated by resampling 10,000 times the coordinates of the *trans* loci. (F)
1054 Manhattan plot of the GWAS results for *ATDNAIT9A*. Distribution of CN values associated with the leading
1055 SNP and the common allele are shown on the right. Note that both *ATDNAIT9A* and *VANDALI6* show similar
1056 insertion preference towards the TSS of genes. (G) Probability of missing a non-reference TE insertion with
1057 TSD as a function of the allele frequency of the *trans* locus identified by GWAS. (H) Gene ontology of genes
1058 overlapping or close to *trans* loci.

1059

1060 **Figure 3-figure supplement 1. Distribution of non-reference TE insertions with TSDs along the genome.**

1061 (A) Density of non-reference TE insertions with TSDs detected by split-reads (black) or TE-sequence capture
1062 (red) across the five chromosomes. (B) Fraction of non-reference TE insertions with TSDs detected in
1063 pericentromeric regions using TE-sequence capture as a function of allele frequency (f). (C) Historical
1064 recombination landscape estimated using genome sequencing data for 211 accessions. (D) Density of non-
1065 reference TE insertions with TSDs as a function of the density of coding sequences. (E) Density of non-
1066 reference TE insertions with TSDs as a function of the recombination rate. The correlation between these two
1067 variables still holds after correcting for the partial correlation between gene density and recombination rate
1068 ($r=0.32$, $P < 4^{-23}$). (F) Distribution frequency of allele counts for non-reference TE insertions with TSDs located
1069 either within or close to genes or away from genes.

1070

1071 **Figure 3-figure supplement 2. Local TE insertion preferences.** (A) Metagene analysis of the distribution of
1072 private non-reference TE insertions with TSDs for four Class II and two Class I TE superfamilies. UTR,
1073 untranslated transcribed region. (B) Sequence motifs for non-reference insertion sites. (C) GC-content for non-

1074 reference insertion sites (including 50bp upstream and downstream). Blue bars represent the GC content for
1075 100bp-long sequences randomly chosen from the reference genome sequence. Bars represent average GC-
1076 content \pm SD. Statistically significant differences were calculated using a permutation test

1077

1078 **Figure 3-figure supplement 3. Global expression levels of genes affected by non-reference TE insertions**
1079 **of different TE superfamilies.** Cumulative distribution of gene expression ratios between alleles harboring and
1080 lacking non-reference TE insertions for the different TE superfamilies. The number of genes analyzed is
1081 indicates in each case. Statistically significant differences between the observed and expected distributions of
1082 expression ratios were calculated using a KS test.

1083

1084 **Figure 3-figure supplement 4. Expression levels of selected genes affected by non-reference TE insertions**
1085 **with TSDs.** Accession(s) with a non-reference TE insertion with TSD are indicated in red for each gene. Plants
1086 were grown under standard conditions (Ctrl.) or subjected to a heat shock (HS). RT-qPCR results (three
1087 technical replicates) are indicated relative to those obtained for a gene that shows invariant expression under
1088 multiple conditions (see ‘Material and Methods’).

1089

1090 **Figure 4-figure supplement 1. DNA methylation of TE families and impact on sequences flanking non-**
1091 **reference TE insertions with TSDs.** (A) Boxplot representation of average DNA methylation levels for non-
1092 mobile TE families (across all accessions) and for mobile TE families (separately for accessions with or without
1093 evidence of mobility). (B) Genome browser tracks for one insertion site. DNA methylation and split-reads are
1094 indicated whenever present. Two accessions (Zdr-1 and Knox-18) have DNA methylation (red arrows) but lack
1095 split reads supportive of the presence of the non-reference TE insertion (black crosses). Another accession
1096 (Anholt-1) has no DNA methylation (red cross) yet contains the non-reference TE insertion (black arrow).
1097 Another 12 accessions have the non-reference TE insertion and all have DNA methylation at the insertion site.
1098 (C) Expected density distribution of the fraction of non-reference TE insertions with TSDs that are located in
1099 the chromosome arms and are associated with long-distance DNA methylation. The observed fraction is
1100 indicated by the vertical line. (D) Genome browser tracks showing the density of 24-nt siRNAs over two Col-0
1101 TE insertions with TSDs that are associated with long- and short-distance DNA methylation, respectively (top
1102 and bottom panels). Brackets indicate the absence of the insertions in the other accessions.

1103

1104 **Figure 5-figure supplement 1. Reconstruction of the *FLC* haplotypes and additional analyses of the effect**
1105 **on flowering time of non-reference TE insertions with TSDs at the locus. (A)** SNPs identify 20 distinct
1106 haplotypes at the *FLC* locus (+/- 50Kb). The seven non-reference TE insertions with TSDs (indicated A to G, as
1107 in Figure 5B) located within *FLC* affect four distinct haplotypes, as shown on the right. **(B)** and **(C)** Flowering
1108 time (Y. Li et al. 2010; Lempe et al. 2005) associated with *FLC* alleles belonging to the same haplotype but
1109 differing by the presence or absence of a non-reference TE insertion with TSD.

1110

1111 **Figure 1-Source data 1. Copy number estimation of TE sequences. (A)** Copy number estimation based on
1112 read coverage for the 317 TE families analyzed across 211 *A. thaliana* accessions collected worldwide. Column
1113 descriptions are provided in **(B)**.

1114

1115 **Figure 1-Source data 2. Coordinates of non-reference TE insertions with TSDs. (A)** Coordinates and
1116 presence or absence call (1 and 0, respectively) across the 211 *A. thaliana* accessions. Description of columns is
1117 provided in **(B)**.

1118

1119 **Figure 1-Source data 3. Number of distinct non-reference TE insertions with TSDs identified by the split-**
1120 **reads approach for each TE family and super-family.**

1121

1122 **Figure 1-Source data 4. TE insertions with TSDs present in Col-0 but absent in Ler-1. (A)** Genomic
1123 coordinates of the insertion in Col-0 and of the corresponding empty site in Ler-1. Description of columns is
1124 provided in **(B)**.

1125

1126 **Figure 2-Source data 1. Copy number estimation used for the geo-climatic associations and GWASs. (A)**
1127 Copy number estimation based on read coverage for the 131 mobile TE families analyzed across 211 *A.*
1128 *thaliana* accessions collected worldwide. **(B)** Copy number estimation based on read coverage for the same 131
1129 mobile TE families across 180 *A. thaliana* accessions from Sweden.

1130

1131 **Supplementary file 1. Summary of GWAS results for CNV.** Distribution of CN values, Manhattan plot and

1132 QQ-plot across the joined data set (391 accessions) for the indicated TE families. Summary statistics of
1133 associations are indicated below. MAF indicates Minor Allele Frequency in the joined dataset. Genes within
1134 GWAS intervals are indicated (*MET2a* is in bold).

1135

1136 **Supplementary file 2. DNA sequence motifs at insertions sites.** Sequence logo of the overrepresented DNA
1137 sequence motifs at insertion sites (\pm 30bp) is shown for the 79 mobile TE families with at least 10 non-
1138 reference TE insertions with TSDs. The number of sequences used in each case is indicated.

1139

1140 **Supplementary file 3. PCR validation of non-reference TE insertions with TSDs and list of primer**
1141 **sequences used in this study.**

1142

1143 **Supplementary file 4. List of TE-capture targets.**

1144

1145 **Supplementary file 5. Historical population-scaled recombination rate map for *A. thaliana***







