

 Open access • Journal Article • DOI:10.1007/S00299-020-02554-8

Genomic re-assessment of the transposable element landscape of the potato genome.

— [Source link](#) 

Diego Zavallo, Juan Manuel Crescente, Juan Manuel Crescente, Magdalena Gantuz ...+5 more authors

Institutions: Spanish National Research Council, National Scientific and Technical Research Council, International Trademark Association

Published on: 20 May 2020 - Plant Cell Reports (Springer Science and Business Media LLC)

Topics: Genome evolution, Genome and Retrotransposon

Related papers:

- [Bioinformatics and genomic analysis of transposable elements in eukaryotic genomes.](#)
- [Discovering and detecting transposable elements in genome sequences](#)
- [The genomic ecosystem of transposable elements in maize.](#)
- [Transposable elements contribute to dynamic genome content in maize](#)
- [Evolutionary history of mammalian transposons determined by genome-wide defragmentation](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/genomic-re-assessment-of-the-transposable-element-landscape-1k3d5t2u1r>

1 **Genomic re-assessment of the transposable element landscape of the potato**
2 **genome**

3
4
5 Diego Zavallo*¹, Juan Manuel Crescente*^{2,4}, Magdalena Gantuz^{3,4}, Melisa Leone^{1,5},
6 Leonardo Sebastian Vanzetti^{2,4}, Ricardo Williams Masuelli^{3,4}, and Sebastian Asurmendi¹
7

8 ¹Instituto de Agrobiotecnología y Biología Molecular (IABIMO), Instituto Nacional de
9 Tecnología Agropecuaria (INTA), Consejo Nacional de investigaciones Científicas y
10 Tecnológicas (CONICET), Los Reseros y Nicolas Repeto, Hurlingham, Argentina

11 ²Grupo Biotecnología y Recursos Genéticos, EEA INTA Marcos Juárez, Ruta 12 km 3,
12 2580 Marcos Juárez, Argentina

13 ³Instituto de Biología Agrícola de Mendoza (IBAM), Facultad de Ciencias Agrarias
14 (FCA), CONICET-UNCuyo, Almirante Brown 500, M5528AHB, Chacras de Coria,
15 Mendoza, Argentina

16 ⁴Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET), Argentina

17 ⁵Agencia Nacional de Promoción Científica y Tecnológica (ANPCyT), Argentina.
18

19 * Both authors contributed equally to this work

20
21 Corresponding author:

22 Diego Zavallo: E-mail: zavallo.diego@inta.gob.ar; Phone: (5411) 4621 447/1676, int.
23 3526. Orcid N°: 0000-0002-9021-2175

24 Sebastian Asurmendi: E-mail: asurmendi.sebastian@inta.gob.ar; Phone: (5411) 4621
25 447/1676, int. 3639. Orcid N°: 0000-0001-9516-5948

26
27
28
29 **Abstract**

30
31 Transposable elements (TEs) are DNA sequences with the ability to auto-replicate
32 and move throughout the host genome. TEs are major drivers in stress response and
33 genome evolution. Given their significance, the development of clear and efficient TE
34 annotation pipelines has become essential for many species. The latest *de novo* TE
35 discovery tools, along with available TEs from Repbase and sRNA-seq data, allowed us
36 to perform a reliable potato TEs detection, classification and annotation through an open
37 -source and freely available pipeline (https://github.com/DiegoZavallo/TE_Discovery).
38 Using a variety of tools, approaches and rules, our pipeline revealed that ca. 16% of the
39 potato genome can be clearly annotated as TEs. Additionally, we described the
40 distribution of the different types of TEs across the genome, where LTRs and MITEs
41 present a clear clustering pattern in pericentromeric and subtelomeric/telomeric regions
42 respectively. Finally, we analyzed the insertion age and distribution of LTR
43 retrotransposon families which display a distinct pattern between the two major

44 superfamilies. While older Gypsy elements concentrated around heterochromatic
45 regions, younger Copia elements located predominantly on euchromatic regions.
46 Overall, we delivered not only a reliable, ready-to-use potato TE annotation files, but
47 also all the necessary steps to perform *de novo* detection for other species.

48

49 **Keywords:** Transposable elements; *Solanum tuberosum*; potato; TEs annotation;
50 Retrotransposons; DNA transposons

51

52 **Key Message**

53

54 We provide a comprehensive and reliable potato TE landscape, based on a wide
55 variety of identification tools and integrative approaches, producing clear and ready-to-
56 use outputs for the scientific community.

57

58 **Introduction**

59

60 Transposable elements (TEs) are DNA sequences with the ability to auto-replicate
61 and move throughout the genome. In plants, TEs can occupy a large proportion of the
62 genome, representing more than half of the total genomic DNA in some cases. For
63 example, they comprise about 85% and 88% of wheat and maize genomes, respectively
64 (Schnable et al., 2009; Appels et al., 2018) which may indicate the relevance of these
65 elements to genome architecture and size (Roessler et al., 2018).

66 Furthermore, over the past few years an increasing number of studies have shed
67 light on their importance in gene regulation (Hirsch and Springer, 2017; Judd and
68 Feschotte, 2018) stress response and genome evolution (Hosaka and Kakutani, 2018).

69 TEs are usually divided into two major classes based on their mechanism of
70 transposition. Class I elements (or retrotransposons) propagate via a reverse-
71 transcribed RNA intermediate, whereas class II elements (or DNA transposons) move
72 through a “cut and paste” mechanism. Another type of classification is by their ability to
73 transpose on their own (i.e. autonomous), characteristic shared by some TEs from both
74 classes. The non-autonomous elements can move but rely on autonomous TEs for their
75 mobility (Wicker et al., 2007).

76 Given their significance, the identification, classification and annotation of TEs has
77 emerged as a new field of great interest in science, which involves both wet-lab biology
78 and bioinformatics. As Hoen et al. (2015) describe in their review on TE annotation
79 benchmarking, precise detection and annotation of TEs is a difficult task due to their
80 great diversity, both within and among genomes. TEs differ across multiple attributes,
81 including transposition mechanism, sequence, length, repetitiveness, and chromosomal

82 distribution. In addition, whereas recently inserted TEs have a relatively low variability
83 within the family, over time they accumulate mutations and diverge, making them harder
84 to detect (Hoen et al., 2015).

85 There are two main strategies for TE annotation: homology-based and *de novo*
86 identification, which can also be referred to as library-based and signature-based,
87 respectively. The homology-based strategy uses libraries of known TEs such as the
88 Repbase repository (Jurka et al., 2005) to screen genomes in order to identify similar
89 sequences, most commonly by using RepeatMasker (Smit et al., 1996). On the other
90 hand, *de novo* approaches use characteristic structural features, such as LTRs (Long
91 Terminal Repeats) for retrotransposons and TIRs (Terminal Inverted Repeats) for DNA
92 transposons, to identify new elements. Moreover, autonomous TEs have conserved
93 structures like RT (reverse transcriptase) or TR (transposase) that can also be used for
94 accurate TE identification (Wicker et al., 2007). Several tools based on structural
95 features, such as LTRharvest (Ellinghaus et al., 2008) and TIRvish command, which are
96 part of the GenomeTools suite (Gremme et al., 2013), are available. Other tools based
97 on this criteria are specifically designed to discover different types of TE families (e.g.
98 SINE Scan (Mao and Wang, 2016), HelitronScanner (Xiong et al., 2014) and MITE
99 Tracker (Crescente et al., 2018)). Another *de novo* based strategy relies on the most
100 important biological mechanism that silences TEs - RNA directed DNA Methylation
101 (RdDM) - in which double-stranded RNAs (dsRNAs) are processed into 21-24 nt small
102 interfering RNAs (siRNAs) and guide methylation on homologous DNA loci. For
103 instance, TASR (for Transposon Annotation using Small RNAs) tool uses sRNAs
104 Illumina data as guide for TE annotation/identification (El Baidouri et al., 2015).
105 However, complete and accurate TEs annotation will likely require a combination of both
106 homology-based and *de novo* methods together with an additional manual curation step
107 (Platt et al., 2016).

108 Nevertheless, all of these tools often give a representative sequence for each family
109 (usually a full-length TE sequence), but fail presenting their copies across the genome,
110 not only for other potentially autonomous copies, but also for members that are partially
111 or entirely deficient in one or more domains. Furthermore, in the case of non-
112 autonomous TEs (e.g. SINEs, TRIMs and MITEs), the amount of copies across the
113 genome is high due to their repetitive content and their short length; however this
114 information is usually not presented.

115 For the scientific community, especially in genomics, the need for reliable annotation
116 is becoming fundamental. Generally, for each genome sequencing project, annotations
117 consisting mainly of protein-coding genes (structural and functional) and miRNAs genes
118 are made, whereas TEs remain poorly annotated. As an example, the Gramene

119 Database (<http://www.gramene.org>) is a curated, open-source database for comparative
120 functional genomics in crops and model plant species with information on more than two
121 million genes from 67 plant genomes. By contrast, the PGSB Repeat Element Database
122 (Nussbaumer et al., 2012) which compiles publicly available repeat sequences from
123 TREP (Wicker et al., 2002), TIGR repeats (Ouyang and Buell, 2004), PlantSat (Macas
124 et al., 2002), Genbank and *de novo* detected LTR retrotransposon sequences from
125 LTR_STRUC (McCarthy and McDonald, 2003) only comprises 62.000 sequences.
126 Nonetheless, in the latest version of sunflower genome, Badouin et al. (2017) performed
127 a comprehensive search for repeat elements by developing a new tool called Tephra
128 (<https://github.com/sestaton/tephra>) (Staton, 2018), which discovers and annotates all
129 types of transposons. Tephra combines existing specific transposon discovery tools for
130 all types of TEs, classifies and annotates them, but still lacks information on copy
131 numbers across the genome. A recent study addressed this issue by applying a method
132 called “Russian doll” due to its nesting strategy. This method builds nested libraries
133 establishing different search rules for each one of them. The first one includes only
134 “potentially autonomous TEs”, the second one contains “total TEs”, including non-
135 autonomous and a third one that also includes uncategorized “repeated elements”
136 (Bertheliet et al., 2018).

137 *Solanum tuberosum*, the cultivated potato, is the third most important food crop after
138 rice and wheat, and the main horticultural crop (Devaux et al., 2014). The sequencing of
139 the *S. tuberosum* genome resulted in an assembly of 727 Mb of 810.6 Mb sequenced.
140 Because most potato cultivars are autotetraploid ($2n=4x=48$) and highly heterozygous,
141 sequencing was performed on a homozygous doubled-monoploid potato clone. The
142 latest potato genome assembly (4.03) contains 39.031 annotated genes and 62.2% of
143 the genome corresponds to repetitive elements at scaffold level (Consortium et al.,
144 2011).

145 Many attempts have been made to discover repetitive elements in Solanaceae
146 families. Most of these studies were mainly focused on tandem repetitive elements,
147 whereas studies of complex repetitive elements were mostly performed on limited
148 groups of TEs (Mehra et al., 2015). The researchers assessed the complex repetitive
149 elements in potato (*S. tuberosum*) and tomato (*S. lycopersicum*) genomes, identifying
150 629,713 and 589,561 repetitive elements, respectively. Mehra et al. used
151 RepeatModeler (<http://repeatmasker.org/RepeatModeler.html>), which employs a
152 repetitiveness-based strategy, and enriched the amount of repeat families previously
153 identified in the 4.03 version of the potato genome with RepeatMasker (Smit et al.,
154 1996).

155 In this study we present an optimized pipeline of transposable elements detection
156 and annotation from *S. tuberosum*. Our strategy relies on the combination of the latest
157 *de novo* TE identification tools, available TEs from Repbase and Illumina sRNA-seq
158 data to obtain TEs. We then find copies and applied a series of filters depending on the
159 TE family to obtain a comprehensive and curated whole-genome atlas of potato
160 transposable elements. Furthermore, we provide to the research community our
161 pipeline, annotation results and files, which are publicly available at
162 https://github.com/DiegoZavallo/TE_Discovery to encourage reproducibility and the
163 eventual implementation of our framework in diverse organisms.

164

165 **Materials and Methods**

166

167 **Input data**

168

169 The latest assembly version of *Solanum tuberosum* genome sequence was
170 downloaded from the Potato Genome Sequencing Consortium (PGSC v4.03 of the
171 doubled monoploid *S. tuberosum* Group Phureja DM1-3). Potato TEs sequences from
172 Repbase Giri (Jurka et al., 2005) were downloaded prior registration. Note that LTRs
173 transposons are divided in LTR (Long Terminal Repeat) and I (Internal) sequences,
174 hence must be concatenated. We gathered 126 LTRs, 18 LINEs, 2 SINEs and 42 TIRs
175 family sequences
(<https://www.girinst.org/repbase/update/search.php?query=tuberosum&querytype=Taxonomy>).

178 Illumina sRNA-seq data for TASR run was generated by our lab (data unpublished);
179 however, any available data from public repositories such as SRA
180 (<https://www.ncbi.nlm.nih.gov/sra>) could be used as input.

181

182 **Transposable elements identification**

183

184 TEs obtained from different sources were merged together according to each
185 element classification. Tephra, which uses several structure-based tools, was applied to
186 harvest different kind of TEs in the potato genome. *Tephra all* command (which runs all
187 subcommands) was executed using *tephra config.yml* file with default configuration
188 parameters with the exception of *S. tuberosum* TEs from Repbase for the *repeatdb*
189 parameter. A total of 1,325 Helitrons, 7,694 LTRs, 2,994 MITEs, 2,414 TIRs and 7,011
190 TRIMs families were found with this program. No non-LTR TEs (LINEs and SINEs) were
191 found by Tephra. TASR (El Baidouri et al., 2015), a tool for *de novo* discovery of TEs
192 using small RNA data, was also used in this work. 21, 22 and 24 nt sRNAs were parsed

193 from files belonging to all treatments and replicates from our sRNA-Seq data and
194 subsequently concatenated to be used as input. TASR.v.1.1.pl perl script was run with
195 default parameters except for: *-cpu 14, -nsirna 10, and -cnumber 5*. A total of 1,916
196 families were found and presented as multifasta files comprising all the elements for
197 each family. A perl script provided by TASR developer was run in order to generate a
198 consensus sequence for each family. Then we created a single multifasta file with the
199 entire consensus. The PASTEC tool (Hoede et al., 2014) was used for this purpose,
200 since TASR does not classify TEs into categories. From the 1,916 families discovered, a
201 total of 891 LTRs, 49 LINEs, 15 SINEs, 9 TRIMs, 2 LARDs, 84 TIRs, 35 MITEs and 5
202 Helitrons were classified. For MITEs discovery, MITE Tracker (Crescente et al., 2018)
203 was employed using default parameters. A total of 1,045 MITEs elements were
204 detected. SINE_Scan (Mao and Wang, 2016), an efficient de novo tool to discover
205 SINEs was also used in this work with default parameters. A total of 13 SINEs families
206 were detected. As a result, we obtained 8,711 LTRs elements from Tephra, TASR and
207 Repbase, 67 LINEs elements from TASR and Repbase, 30 SINEs elements from TASR,
208 Repbase and SINE_Scan, 540 TIRs elements from Tephra, TASR and Repbase, 4,074
209 MITEs elements from Tephra, TASR and MITE Tracker and 1,330 Helitrons elements
210 from Tephra and TASR.

211

212 **Pipeline description**

213

214 To detect, filter and annotate TEs copies across the potato genome the obtained TEs
215 list was subjected to an in house pipeline. The pipeline was developed using bash
216 scripts and Jupyter notebooks.

217 *1.1 add_annotation.ipynb* uses TEs sequences retrieved from each program and
218 merge them together into one multi-fasta file per studied TE type (LTRs, LINEs, SINEs,
219 TRIMs, LARDs, TIRs, MITEs and Helitrons). This script adds to each sequence a
220 unique identifier containing an auto-incremented number, TE classification and the
221 program source from which it was obtained.

222 *2.1 vsearch.sh* uses VSearch program to cluster similar sequences that share 80%
223 identity, according to the 80/80/80 rule in the study of Wicker et al. (2007). It is executed
224 once per TE type, thus obtaining as a result TEs clustered by type.

225 *2.2 vsearch_merge.ipynb* script uses vsearch outputs to create fasta files containing
226 one TE per family. It also adds a family description indicating which program the
227 members came from.

228 *3.1 blast.sh* performs a genome-wide BLASTn search using the files from the
229 previous step and searches for TEs in the potato genome. For this task, the script uses

230 the following parameters: *-perc identity 80, -evalue 10e-3 and -task blastn* (except for
231 LTRs). *-qcov hsp perc 80* was used for SINEs, TRIMs and MITEs, whereas *-qcov hsp*
232 *perc 50* was set for the rest.

233 *3.2 blast filter.ipynb* filters BLASTn results by using parameters according to Table 1.
234 First, each file is filtered by a length range defined by min len and max len. Afterwards, a
235 length threshold range is calculated by multiplying the query length by a min and max
236 subject length subject length. The subject sequence has to be inside this range to be
237 considered valid. Later, minimum identity percentage and query coverage are required
238 for the sequences to be valid. Finally, duplicated hits are removed by searching those
239 whose start and end positions overlaps within a margin of plus or minus 5 nt.

240 *4.1 annotate.ipynb* transforms the BLAST tab-delimited results to a *gff3* format file,
241 adding a detailed description for each TE. The description includes TE id (a numeric
242 identifier of the element after clustering) source name (original id name of the element
243 before clustering) type (family type of the element), source (program or tool from which
244 the TEs were detected) and unique id (unique identifier for copy element).

245

246 **LTR age**

247

248 To determine tentative LTRs insertion age we used *tephra ltrage* command
249 implemented by the Tephra package. This command uses the Tephra-discovered full
250 length LTRs, aligns the LTR sequences and generates a neighbor-joining guide tree with
251 MUSCLE (Edgar 2004). The alignment and guide tree are used to generate an
252 alignment in PHYLIP format. A likelihood divergence estimate was calculated with
253 baseml from PAML (Yang 2007) by specifying the K80 substitution model. This
254 divergence value (hereafter d) was used to calculate LTR-RT age with the formula $T =$
255 $d/2r$, where $r = 1e8$ is the default substitution rate.

256

257 **Data resource**

258

259 Scripts from this work, including all pipeline steps as well as circos ideogram,
260 distance histogram and LTR age plot scripts are available at
261 ([https://github.com/DiegoZavallo/TE Discovery](https://github.com/DiegoZavallo/TE_Discovery)). Annotation and fasta files are available
262 as supporting information.

263

264 **Results and Discussion**

265

266 The goal of the present work was to assemble a comprehensive TE repertoire of the
267 potato genome and provide legible, ready-to-use files for the scientific community. To

268 address this issue, we used a combination of different approaches to identify TEs such
269 as similarity-based, structure-based and mapping-based strategies. Moreover, we
270 introduce a set of scripts to gather, detect, filter and ultimately annotate TE copies from
271 all classes across the potato genome and present *gff3* files of TE features. Additionally,
272 we display the accumulation and distribution across the genome of the different types of
273 TEs and a table summarizing data and metrics, such as distances to nearest gene and
274 LTRs insertion ages.

275

276 **TEs in the potato genome**

277

278 Public *S. tuberosum* Rebase library (Jurka et al., 2005), the potato reference
279 genome and Illumina small RNA-seq data were used as **input data**. **TEs family**
280 **detection** was executed by applying two “All TEs” tools: Tephra and TASR, which
281 discover *de novo* TEs with structure-based and map-based approaches, respectively. To
282 complement the search, we applied “Specific TEs” tools: MITE Tracker and SINE_Scan.
283 The obtained sequences were merged into multi-fasta files, one for each type of TE, and
284 headers were renamed. A **clustering** step was carried out with Vsearch to reduce
285 redundancy (Online Resource 1). Next, **detection of copies** in the genome of the
286 different TE families was achieved by conducting a BLAST search with specific
287 parameters according to the type of TE (see Materials and Methods section). A **copy**
288 **filter** step was implemented by establishing several rules with very stringent criteria to
289 detect “potential autonomous TEs” and a second set of rules for “All TEs” with more
290 relaxed parameters to account for autonomous and non-autonomous TEs of all types
291 (Table 1). Finally, an **annotation** step was performed to generate *gff3* files containing a
292 description for each element that includes TE type and the detection tool that identified
293 it. Figure 1 shows an overview of the pipeline used to detect and re-annotate the potato
294 TEs.

295 TE content is highly variable in plants and usually displays a positive correlation with
296 genome size. For instance, as much as 85 % of maize genome or 70 % of Norway
297 spruce genome (Nystedt et al., 2013) has been annotated as transposons including
298 unclassified ones, whereas in the more compact *Arabidopsis thaliana* genome TE
299 content is only 21 % (Ahmed et al., 2011). In potato, the data presented by Mehra et al.
300 (2015) comprised an annotation file of 1,061,377 repetitive elements, including rRNA,
301 tRNA, simple repeats and low complexity elements which represents almost 50% of the
302 genome. When only the most complex elements (i.e. transposons) were taking into
303 account, the coverage percentage dropped to nearly 34%.

304 However, our pipeline revealed a TE content of ~16% (excluding the unanchored
305 ChrUn), representing half of the genomic coverage according to the data presented by
306 Mehra et al. (2015).

307 Of those ~16%, LTRs comprised around 13% of the potato genome, which
308 corresponds to over 80% of the total TEs. The most abundant superfamily was Gypsy,
309 whereas the other types of TEs barely made up 1% of the genome coverage. For
310 instance, each DNA TE (TIRs, MITEs and Helitrons) covered almost the same
311 percentage of the genome with 0.51, 0.72 and 0.72 %, respectively. These coverage
312 ratio patterns are in agreement with results from most plant genomes that have TE
313 identification projects (Du et al., 2010; Andorf et al., 2016; Badouin et al., 2017; Alaux et
314 al., 2018). Table 2 summarizes the amount and diversity of all identified TEs, filtered
315 copies and proportion in the genome of all TE families.

316 To understand more deeply the discrepancy between the coverage of the TE genome
317 presented here and the one reported by Mehra et al. (2015), we should observe other
318 differences between both works. For instance, even though we describe all types of TE
319 families reported by Wicker et al. (2007), including LARDs, TRIMs and MITEs, which
320 were absent in Mehra study, we annotated less than half of the sequences. We applied
321 a clustering method to decrease redundancy and established a set of filters selected
322 specifically for each type of TEs, which, is aimed to reduce false signal. Moreover, as
323 already mentioned, Mehra et al. used a unique tool that relies on repetitiveness-based
324 strategy, leaving aside a wide variety of detection methods that in this work we have
325 combined, which is indicated to improve TE detection efficiency (Kamoun et al., 2013;
326 Hoen et al., 2015; Arensburger et al., 2016).

327 We scored a total of 243,010 elements compared to 629,713 complex elements
328 previously found by Mehra et. al. and when compared, 198,025 (81%) of our sequences
329 overlapped with their set, which evidences the effectiveness of both annotation
330 pipelines.

331 To test the effectiveness of the pipeline, we run it on a well annotated genome such
332 as soybean. SoyTEdb represents an example of a thoroughly annotated TE database in
333 which a combination of structure-based and homology-based approaches was used to
334 structurally annotate and clearly categorize TEs in the soybean genome (Du et al.,
335 2010). The authors reported over 38,000 TEs representing ~17% of the genome.
336 However, when they informed the genome coverage they included fragments defined by
337 RepeatMasker (i.e. low complex repeats) rising up to 58% of the soybean genome.
338 These data may indicate the existence of a large set of repetitive sequences in the
339 genome that cannot be annotated as TE with current knowledge about the structure of
340 TE. One hypothesis that could arise from this is that the structural patterns that

341 represent these non-annotated TEs are not yet thoroughly described, or they are just
342 repetitive DNA that simply cannot be assigned as TEs. We used the 38,000 annotated
343 TEs to run our pipeline for copy elements discovery and filtering steps and we came out
344 with similar genomic coverage (~23%) and more than 75% of the TEs sequences
345 overlapped. Even though this was only a test, since the full run of our pipeline on
346 another species would require a more comprehensive approach which exceeds this
347 work, it validates that this pipeline does not underestimate the occurrence of TEs in the
348 genome.

349 Finally, we deliver ready-to-use annotation files in *gff3* format of all TEs annotated
350 with our pipeline with detailed descriptions in the ninth column including *TE id*, *source*
351 *name*, *type*, *source* and *unique id* (Online Resource 2).

352

353 **Distribution of TEs across the potato genome**

354

355 It is well know that the distribution, amount and genome coverage of TEs vary greatly,
356 particularly between plants and animals, where LTRs and non-LTRs (LINEs and SINEs)
357 are the predominant type of TE, respectively (Chalopin et al. 2015). Moreover, TE
358 distribution is highly dependent on the family type. Some TEs are more prone to
359 concentrate in regions near protein coding genes while others are more equally
360 distributed along the genome.

361 To assess the type-specific landscape of the diverse TE categories, we performed
362 circos ideograms for each TE type separated by class, as well as gene density as
363 reference. Each concentric circle represents the coverage percentage of one type of TE.
364 Since they have different coverage ranges, each type has its own color pallet to
365 appreciate the distribution across the chromosomes (Figure 2).

366 Left panel of Figure 2 shows the class I elements, where LTRs stand out, not only for
367 their clear pattern of clustering around centromeric and pericentromeric regions, but also
368 for their high coverage in some areas of the chromosomes. For instance, each dark red
369 line of the second circle represent up to 36% of LTRs coverage per Mb, which explains
370 the 13% genome-wide coverage for this kind of TE (Table 2). Furthermore, LTRs
371 distribution is virtually opposite to protein coding gene distribution (Figure 2, left panel).
372 This behavior has already been reported for other plants (Baucom et al., 2009; Paterson
373 et al., 2009; Badouin et al., 2017).

374 Conversely, SINEs, which are the least represented type (besides LARDs which are
375 not shown since we discovered only 18 elements) seem to have an even distribution
376 pattern with a slight tendency towards telomeric and subtelomeric regions in some
377 chromosomes. LINEs and TRIMs also display homogeneous distribution patterns with

378 some coverage hotspots, mainly near telomeric regions. Due to their length and filtering
379 parameters established, only 248 LINES were found in the genome using our
380 methodology, contrasting with more than 50,000 elements found by Mehra et al. (2015).
381 Given that Heitkam et al. (2014) extracted 59,390 intact LINE sequences from 23 plant
382 genomes in order to classify them into families, it is somehow unlikely that potato alone
383 could have 50,000 LINES. A look at the data presented by Mehra et al. (2015) shows
384 that some elements annotated as LINES are small fragments with some identity to
385 LINES given by the RepeatMasker tool which in our case were removed by the filtering
386 process.

387 In a recent work, Gao et al. (2016) performed a comprehensive analysis of TRIMs in
388 48 plant genomes, including *S. tuberosum*. They observed that TRIMs are generally
389 enriched in genic regions and likely play a role in gene evolution (Gao et al., 2016).
390 They discovered 12,473 copies in the potato genome representing 0.46% of the
391 genome, which is consistent with our results (Table 2).

392 Right panel of Figure 2 shows class II TEs distribution in which MITEs display a clear
393 concentration pattern around gene-rich subtelomeric regions. According to previous
394 works, MITEs are often found close to or within genes, where they affect gene
395 expression (Bureau and Wessler, 1994b). Indeed, MITEs may affect gene regulation via
396 small RNA pathways, in addition to play the canonical role of TEs in the evolution by
397 altering gene structure (Kuang et al., 2009; Gagliardi et al., 2019).

398 TIRs and Helitrons displayed an unbiased distribution across the chromosomes
399 (Figure 2, right panel). Helitron chromosome distribution seems to vary by species.
400 While in *Arabidopsis* Helitrons are enriched in gene-poor pericentromeric region, in
401 maize they are more abundant in gene-rich regions (Yang and Bennetzen, 2009). Rice,
402 on the other hand, exhibited a more erratic pattern of Helitron distribution (Yang and
403 Bennetzen, 2009), more similar to our results.

404 In sum, this kind of analysis allowed us to have a holistic view of the different
405 distribution patterns of TEs across the genome and to elucidate their potential role in
406 transcriptional gene regulation.

407

408 **TEs and genes**

409

410 The importance of TEs accumulation near genes, where these elements could
411 influence gene expression, has been extensively reported (Bureau and Wessler
412 1994b,a; Wang et al., 2013). As we described above, we found a strong positive
413 correlation between LTRs and pericentromeric regions as well as a positive correlation
414 between MITEs and rich-gene subtelomeric regions. Hence, we determined the

415 distances from the different types of elements to the closest gene in the genome in
416 order to assess how many TEs are likely within or close to genes. For this purpose, we
417 computed metrics such as median distance to the closest gene and percentage of TEs
418 overlapping gene transcripts or near coding regions. LTRs were found to be generally
419 far from genes with a median distance of 10.25 kb; 9.01% overlap with coding sequence
420 genes and only 5.11% are in the immediate vicinity (up to 1 kb from the nearest gene),
421 totaling a 32.23% within the first 5 kb, including elements within transcripts (Figure 3 and
422 Online Resource 3). TRIMs exhibit a similar behavior, with a median distance of 8.12kb,
423 6.43% elements located within a gene and 6.19% in the range of 0-1 kb to the nearest
424 transcript, comprising a 36.62% in the first 5 kb. In contrast, LINEs and SINEs are close
425 to genes (median distance of 4.16 kb and 3.86 kb, respectively) with more than 20%
426 elements within a transcript and rising up to >50% in the 0-5 Kb range (52.44% and
427 56.20% respectively).

428 Class II TEs also have more than 50% of their elements within the first 5 kb.
429 However, TIRs and MITEs distribution differs from that of Helitrons. TIRs and MITEs are
430 the only types of TEs that appear to have less elements within a gene than in the 0-1 kb
431 range (TIRs: 10.05% against 15.5% and a median distance of 3.32 kb; MITEs: 7.07%
432 against 14.21% and a median distance of 3.53 kb). On the other hand, Helitrons display
433 a similar pattern to LINEs and SINEs with a median of 4.56 kb and 23.65% elements
434 within a transcript. These differences can be observed in the less pronounced curve of
435 TIR and MITE histograms demonstrating that a significant proportion of elements are
436 indeed in the proximity of genes but not necessarily within them. (Figure 3, right panel
437 and Online Resource 3).

438 Overall, DNA TEs and nonLTRs retrotransposons have more than 50% of the
439 elements inserted within transcripts or in a range of 0-5 kb from the nearest gene.
440 Several studies have reported examples of transcriptional impact due to TEs insertion
441 near genes in tomato (Xiao et al., 2008; Quadrana et al., 2014), potato (Momose et al.,
442 2010; Kloosterman et al., 2013), melon (Martin et al., 2009) and orange (Butelli et al.,
443 2012) among others. Several mechanisms including disruption of promoter or reduction
444 of transcription through the spread of epigenetic silencing often suppress expression.
445 However, TE can also introduce new sequences in the promoter, leading to up-
446 regulation of proximal gene (Yan et al., 2004; Cowley and Oakey, 2013; Dubin et al.,
447 2018). Insertion of a TE into the coding sequence can disrupt gene function, generally
448 resulting in loss-of-function mutations, particularly if located in an exon. Intronic TEs can
449 also be harmful, for instance by altering splicing patterns (Saze et al., 2013; Ong-
450 Abdullah et al., 2015).

451 In contrast, LTRs and TRIMs located near genes barely exceed 30% and, with a few
452 exceptions, they appear in intergenic, heterochromatic and gene-poor-regions (Kumar
453 and Bennetzen, 1999). In *Arabidopsis*, Wang et al. (2013) reported that gene expression
454 is positively correlated with the distance of the gene to the nearest TE, and negatively
455 correlated with the number of proximal TEs. Whether LTR-like TEs specifically target
456 these regions or if they are simply not selected against and accumulated in regions
457 nearby genes remains unclear (Sigman and Slotkin, 2016).

458

459 **LTR elements age and evolution**

460

461 LTRs are the most represented TEs in the majority of plant genomes, encompassing
462 more than 75% of the nuclear genome of some species (Kumar and Bennetzen, 1999;
463 Paz et al., 2017). For this reason, we decided to explore the evolutionary history of
464 LTRs during potato genome evolution.

465 We analyzed the age and distribution of the elements discovered by Tephra (Staton,
466 2018) by means of *tephra ltrage* command that allowed the characterization of
467 phylogenetic substructure within families of LTR retrotransposons. A total of 8,034 full-
468 length LTRs were analyzed in terms of chromosomal distribution and insertion age by
469 plotting all LTR elements together (Figure 4, upper panel), or grouped into superfamilies
470 (Figure 4, bottom panel).

471 Younger LTRs are enriched in euchromatic subtelomeric regions and correspond
472 mainly to Copia (RLC) family of TEs (average insertion age of 2.54 mya), whereas older
473 LTRs are more abundant in heterochromatic pericentromeric regions where Gypsy
474 (RLG) elements are mostly located (average insertion age of 4.14 mya). These
475 distributions are in agreement with previous findings in maize (Sun et al., 2018), wheat
476 (Luo et al., 2017) and tomato (Paz et al., 2017).

477 Unclassified LTRs (RLX) display a more even distribution, although slightly towards
478 to pericentromeric regions, as well as an intermediate insertion age (average insertion
479 age of 3.78 mya) (Figure 4, bottom panel).

480 To determine whether these aging differences were a family trait or a genomic region-
481 dependent mutation/substitution rate, we divided each chromosome in euchromatic and
482 heterochromatic regions. In order to do so, we compared the pachytene karyotype
483 previously published (Consortium et al., 2011) with the chromosomal ideograms we
484 produced, and plotted the frequency of each LTR family by insertion age according to
485 the karyotype determined region (Figure 5, upper panel). Gypsy family (RLG) display a
486 Gaussian distribution centered around four millions years both in euchromatic and
487 heterochromatic regions, whereas Copia family (RLC) have a chi-square distribution

488 with a peak between two and three millions years on both regions. These results
489 suggest that the insertion age of LTRs retrotransposons depends on the superfamily
490 and not on differential mutation/substitution rate by region.

491 Furthermore, we plotted ten random independent families encompassing different
492 member sizes from RLG, RLC and RLX to assess the age distribution of individual
493 elements by region (Figure 5, bottom panel). Only four out of the ten evaluated Gypsy
494 families showed significant differences (t-test $p < 0.05$) in aging by region. In all the
495 cases, heterochromatic elements were older than euchromatic elements, which reflect a
496 slight age difference within elements of the same family owing to their insertion sites.
497 This suggests that insertions into heterochromatic regions are more likely to persist for
498 longer periods of time.

499 Overall, these variances may have a component based on differential
500 mutation/substitution rate by region but more importantly it is strongly affected by
501 superfamily traits. A recent study by Quadrana et al. (2019) described the mechanisms
502 for which several TEs of the Copia superfamily preferentially integrate within genes by
503 association with H2A.Z-containing nucleosomes. Moreover, they suggest that the role of
504 H2A.Z in the integration of Copia retrotransposons has been evolutionary conserved
505 since the last common ancestor of plants and fungi (Quadrana et al., 2019).

506 On the other hand, Gypsy superfamily harbors a chromodomain that interacts with
507 repressive histone marks such as H3K9m2 which targets to heterochromatin (Sultana
508 et al., 2017).

509

510 **Conclusion**

511

512 TEs have been historically neglected in genome assembly projects, partially due to
513 their repetitive nature but also because their heterogeneity in sequence, size, number of
514 copies, distribution and mutation rates both inter and intra species, make them very
515 difficult to detect accurately (Bourque et al., 2018).

516 However, in the past years, TEs discovery and annotation for the main crops have
517 emerged with diverse results of coverage, complexity and accuracy. For instance, by
518 using the CLARITE software on the IWGSC RefSeq v1.0 genome assembly, Alaux et al.
519 (2018) found over 5 million elements from all types of TEs in wheat.

520 The Rice TE Database collects repeat sequences and TEs of several species of
521 *Oryza* (rice) genus. All sequences have been characterized adopting Wicker's
522 classification code and extending it by encoding new TE superfamilies and non-TE
523 repeats. Particular emphasis was given to the proper classification of sequences and to
524 the removal of nested insertions (Copetti et al., 2015).

525 As we described above, potato has some TE sequences annotated on the RepBase,
526 essentially, a list based on RepeatMasker provided by the Potato Genome Sequencing
527 Consortium and the work from Mehra et al. (2015), which is based on RepeatModeler.
528 However, this worldwide important crop still lacked a comprehensive, multiple-based
529 discovery approach focused on transposon elements annotation.

530 This work arises as a need of a reliable potato TE atlas for ongoing projects that
531 involve epigenetic and transcriptional regulation for different *Solanum tuberosum* in
532 contrasting environments.

533 Plants are known for their phenotypic plasticity and good adaptation to environmental
534 changes due to their sessile condition. There is an increasing evidence of the impact
535 that TEs have on the transcriptome on response to stress. For example, TEs may have
536 direct effects on genes regulation, by providing them new coding or regulatory
537 sequences, changes on the epigenetic status of the chromatin close to genes, and more
538 subtle effects by imposing diverse evolutionary constraints to different chromosomal
539 regions (Makarevitch et al., 2015; Vicient and Casacuberta, 2017; Cambiagno et al.,
540 2018). Finding common patterns or specific alterations on these features (TEs) and
541 linking them with sRNAs profiles and DNA methylation patterns could help researchers
542 to elucidate the underlying mechanisms of transcriptional changes during different
543 stress conditions.

544 As mentioned above, the current data regarding TEs on potato genome is either
545 scarce or imprecise for transcriptional analysis. Given the importance of transposon
546 elements, we provide here a comprehensive potato TE landscape, based on a wide
547 variety of identification tools and approaches, clustering methods, copies detection,
548 filtering rules and clear outputs, that the scientific community will likely use for metadata
549 analysis.

550

551 **Acknowledgements**

552

553 The authors would like to thank Dr. Soledad Lucero and Dr. Julia Sabio y Garcia for
554 the assistance with English-language editing, and Humberto Julio Debat for his critical
555 comments on the manuscript. This work was supported by the Instituto Nacional de
556 Tecnología Agropecuaria (INTA) and by ANPCyT PICT 2015-1532, and PICT 2016-
557 0429. The funders had no role in this study design, data collection and analysis,
558 decision to publish or preparation of the manuscript.

559

560 **Author contributions**

561

562 **DZ:** Conceptualization, Data curation, Formal analysis, Investigation, Methodology,
563 Visualization, Writing - original draft, Writing - review & editing. **JMC:** Formal analysis,
564 Data curation, Investigation, Methodology, Software, Writing – review & editing.
565 **MG:** Investigation, Writing – review & editing. **ML:** Investigation, Writing – review &
566 editing. **LSV:** Formal analysis, Investigation, Writing – review & editing. **RWM:** Funding
567 acquisition, Investigation, Writing – review & editing. **SA:** Conceptualization, Funding
568 acquisition, Investigation, Project administration, Resources, Supervision, Writing -
569 review & editing.

570

571 **Conflict of interest**

572

573 The authors have no conflicts of interest to declare

574

575 **Electronic Supplementary Material**

576

577 **Online Resource 1.** Fasta files containing family sequences for each type transposon
578 elements.

579

580 **Online Resource 2.** gff3 file containing all TEs copy coordinates including Chr00.

581

582 **Online Resource 3.** Table resuming metrics of TEs to the nearest genes.

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601 **References**

602

603

604 Ahmed, I., Sarazin, A., Bowler, C., Colot, V., and Quesneville, H. (2011). Genome-wide
605 evidence for local DNA methylation spreading from small RNA-targeted sequences in
606 *Arabidopsis*. *Nucleic Acids Research*, 39(16):6919–6931.

607

608 Alaux, M., Rogers, J., Letellier, T., Flores, R., Alfama, F., Pommier, C., Mohellibi, N.,
609 Durand, S., Kimmel, E., Michotey, C., et al. (2018). Linking the International Wheat
610 Genome Sequencing Consortium bread wheat reference genome sequence to wheat
611 genetic and phenomic data. *Genome biology*, 19(1):111.

612

613 Andorf, C. M., Cannon, E. K., Portwood, J. L., Gardiner, J. M., Harper, L. C., Schaeffer, M.
614 L., Braun, B. L., Campbell, D. A., Vinnakota, A. G., Sribalusu, V. V., et al. (2016).
615 MaizeGDB update: new tools, data and interface for the maize model organism
616 database. *Nucleic Acids Research*, 44(D1):D1195–D1201.

617

618 Appels, R., Eversole, K., Feuillet, C., Keller, B., Rogers, J., Stein, N., Pozniak, C. J.,
619 Choulet, F., Distelfeld, A., Poland, J., et al. (2018). Shifting the limits in wheat research
620 and breeding using a fully annotated reference genome. *Science*, 361(6403):eaar7191.

621

622 Arensburger, P., Piégu, B., and Bigot, Y. (2016). The future of transposable element
623 annotation and their classification in the light of functional genomics - what we can learn
624 from the fables of Jean de la Fontaine? *Mobile Genetic Elements*, 6(6):e1256852.

625

626 Badouin, H., Gouzy, J., Grassa, C. J., Murat, F., Staton, S. E., Cottret, L., Lelandais-
627 Brière, C., Owens, G. L., Carrère, S., Mayjonade, B., et al. (2017). The sunflower
628 genome provides insights into oil metabolism, flowering and Asterid evolution. *Nature*,
629 546(7656):148.

630

631 Baucom, R. S., Estill, J. C., Chaparro, C., Upshaw, N., Jogi, A., Deragon, J.-M.,
632 Westerman, R. P., SanMiguel, P. J., and Bennetzen, J. L. (2009). Exceptional diversity,
633 non-random distribution, and rapid evolution of retroelements in the B73 maize genome.
634 *PLoS Genetics*, 5(11):e1000732.

635

636 Berthelie, J., Casse, N., Daccord, N., Jamilloux, V., Saint-Jean, B., and Carrier, G. (2018).
637 A transposable element annotation pipeline and expression analysis reveal potentially
638 active elements in the microalga *Tisochrysis lutea*. *BMC Genomics*, 19(1):378.

639

640 Bourque, G., Burns, K. H., Gehring, M., Gorbunova, V., Seluanov, A., Hammell, M.,
641 Imbeault, M., Izsvak, Z., Levin, H. L., Macfarlan, T. S., et al. (2018). Ten things you
642 should know about transposable elements. *Genome Biology*, 19(1):199.

643

644 Bureau, T. E. and Wessler, S. R. (1994a). Mobile inverted-repeat elements of the tourist
645 family are associated with the genes of many cereal grasses. *Proceedings of the*
646 *National Academy of Sciences*, 91(4):1411–1415.

647

648 Bureau, T. E. and Wessler, S. R. (1994b). Stowaway: a new family of inverted repeat
649 elements associated with the genes of both monocotyledonous and dicotyledonous
650 plants. *The Plant Cell*, 6(6):907–916.

651

- 652 Butelli, E., Licciardello, C., Zhang, Y., Liu, J., Mackay, S., Bailey, P., Reforgiato-Recupero,
653 G., and Martin, C. (2012). Retrotransposons control fruit-specific, cold-dependent
654 accumulation of anthocyanins in blood oranges. *The Plant Cell*, 24(3):1242–1255.
655
- 656 Cambiagno, D. A., Nota, F., Zavallo, D., Rius, S., Casati, P., Asurmendi, S., and Alvarez,
657 M. E. (2018). Immune receptor genes and pericentromeric transposons as targets of
658 common epigenetic regulatory elements. *The Plant Journal*, 96(6):1178–1190.
659
- 660 Chalopin, D., Naville, M., Plard, F., Galiana, D., and Voff, J.-N. (2015). Comparative
661 analysis of transposable elements highlights mobilome diversity and evolution in
662 vertebrates. *Genome Biology and Evolution*, 7(2):567–580.
663
- 664 Consortium, P. G. S. et al. (2011). Genome sequence and analysis of the tuber crop
665 potato. *Nature*, 475(7355):189.
666
- 667 Copetti, D., Zhang, J., El Baidouri, M., Gao, D., Wang, J., Barghini, E., Cossu, R.M.,
668 Angelova, A., Maldonado, L. CE., Roffler, S., Ohyanagi, H., Wicker, T., Fan, C., Zuccolo,
669 A., Chen, M., Costa de Oliveria, A., Han, B., Henry, R., Hsing, Y.-I., Kurata, N., Wang,
670 W., Jackson, S.A., Panaud, O., Wing, R.A. (2015). RiTE database: a resource database
671 for genus-wide rice genomics and evolutionary biology. *BMC Genomics*, 16:538
672
- 673 Cowley, M. and Oakey, R. J. (2013). Transposable elements re-wire and fine-tune the
674 transcriptome. *PLoS Genetics*, 9(1):e1003234.
675
- 676 Crescente, J. M., Zavallo, D., Helguera, M., and Vanzetti, L. S. (2018). Mite tracker: an
677 accurate approach to identify miniature inverted-repeat transposable elements in large
678 genomes. *BMC Bioinformatics*, 19(1):348.
679
- 680 Devaux, A., Kromann, P., and Ortiz, O. (2014). Potatoes for sustainable global food
681 security. *Potato Research*, 57(3-4):185–199.
682
- 683 Du, J., Grant, D., Tian, Z., Nelson, R. T., Zhu, L., Shoemaker, R. C., and Ma, J. (2010).
684 SoyTEdb: a comprehensive database of transposable elements in the soybean genome.
685 *BMC Genomics*, 11(1):113.
686
- 687 Dubin, M. J., Scheid, O. M., and Becker, C. (2018). Transposons: a blessing curse.
688 *Current Opinion in Plant Biology*, 42:23–29.
689
- 690 Edgar, R. C. (2004). MUSCLE: multiple sequence alignment with high accuracy and high
691 throughput. *Nucleic Acids Research*, 32(5):1792–1797.
692
- 693 El Baidouri, M., Kim, K. D., Abernathy, B., Arikat, S., Maumus, F., Panaud, O., Meyers, B.
694 C., and Jackson, S. A. (2015). A new approach for annotation of transposable elements
695 using small rna mapping. *Nucleic Acids Research*, 43(13):e84–e84.
696
- 697 Ellinghaus, D., Kurtz, S., and Willhoeft, U. (2008). LTRharvest, an efficient and flexible
698 software for de novo detection of LTR retrotransposons. *BMC Bioinformatics*, 9(1):18.
699
- 700 Gagliardi, D., Cambiagno, D. A., Arce, A. L., Tomassi, A. H., Giacomelli, J. I., Ariel, F. D.,
701 Manavella, P. A. (2019). Dynamic regulation of chromatin topology and transcription by
702 inverted repeat-derived small RNAs in sunflower. *Proceedings of the National Academy
703 of Sciences*, 116 (35): 17578-1758.
704

- 705 Gao, D., Li, Y., Do Kim, K., Abernathy, B., and Jackson, S. A. (2016). Landscape and
706 evolutionary dynamics of terminal repeat retrotransposons in miniature in plant
707 genomes. *Genome Biology*, 17(1):7.
708
- 709 Gremme, G., Steinbiss, S., and Kurtz, S. (2013). Genometools: a comprehensive software
710 library for efficient processing of structured genome annotations. *IEEE/ACM*
711 *Transactions on Computational Biology and Bioinformatics (TCBB)*, 10(3):645–656.
712
- 713 Heitkam, T., Holtgrawe, D., Dohm, J. C., Minoche, A. E., Himmelbauer, H., Weisshaar, B.,
714 and Schmidt, T. (2014). Profiling of extensively diversified plant lines reveals distinct
715 plant-specific subclades. *The Plant Journal*, 79(3):385–397.
716
- 717 Hirsch, C. D. and Springer, N. M. (2017). Transposable element influences on gene
718 expression in plants. *Biochimica et Biophysica Acta (BBA)-Gene Regulatory*
719 *Mechanisms*, 1860(1):157–165.
720
- 721 Hoede, C., Arnoux, S., Moisset, M., Chaumier, T., Inizan, O., Jamilloux, V., and
722 Quesneville, H. (2014). PASTEC: an automatic transposable element classification tool.
723 *PloS One*, 9(5):e91929.
724
- 725 Hoen, D. R., Hickey, G., Bourque, G., Casacuberta, J., Cordaux, R., Feschotte, C., Fiston-
726 Lavier, A.-S., Hua-Van, A., Hubley, R., Kapusta, A., et al. (2015). A call for benchmarking
727 transposable element annotation methods. *Mobile DNA*, 6(1):13.
728
- 729 Hosaka, A. and Kakutani, T. (2018). Transposable elements, genome evolution and
730 transgenerational epigenetic variation. *Current Opinion in Genetics & Development*,
731 49:43–48.
732
- 733 Judd, J. and Feschotte, C. (2018). Gene expression: Transposons take remote control.
734 *eLife*, 7:e40921.
735
- 736 Jurka, J., Kapitonov, V. V., Pavlicek, A., Klonowski, P., Kohany, O., and Walichiewicz, J.
737 (2005). Repbase update, a database of eukaryotic repetitive elements. *Cytogenetic and*
738 *Genome Research*, 110(1-4):462–467.
739
- 740 Kamoun, C., Payen, T., Hua-Van, A., and Filée, J. (2013). Improving prokaryotic
741 transposable elements identification using a combination of de novo and profile hmm
742 methods. *BMC Genomics*, 14(1):700.
743
- 744 Kloosterman, B., Abelenda, J. A., Gomez, M. d. M. C., Oortwijn, M., de Boer, J. M.,
745 Kowitzanich, K., Horvath, B. M., van Eck, H. J., Smaczniak, C., Prat, S., et al. (2013).
746 Naturally occurring allele diversity allows potato cultivation in northern latitudes. *Nature*,
747 495(7440):246.
748
- 749 Kuang, H., Padmanabhan, C., Li, F., Kamei, A., Bhaskar, P. B., Ouyang, S., Jiang, J.,
750 Buell, C. R., and Baker, B. (2009). Identification of miniature inverted-repeat
751 transposable elements (MITEs) and biogenesis of their siRNAs in the Solanaceae: new
752 functional implications for MITes. *Genome Research*, 19(1):42–56.
753
- 754 Kumar, A. and Bennetzen, J. L. (1999). Plant retrotransposons. *Annual Review of*
755 *Genetics*, 33(1):479–532.
756
- 757 Luo, M.-C., Gu, Y. Q., Puiu, D., Wang, H., Twardziok, S. O., Deal, K. R., Huo, N., Zhu, T.,
758 Wang, L., Wang, Y., et al. (2017). Genome sequence of the progenitor of the wheat D
759 genome *aegilops tauschii*. *Nature*, 551(7681):498.

760
761 Macas, J., Meszaros, T., and Nouzova, M. (2002). PlantSat: a specialized database for
762 plant satellite repeats. *Bioinformatics*, 18(1):28–35.
763
764 Makarevitch, I., Waters, A. J., West, P. T., Stitzer, M., Hirsch, C. N., Ross-Ibarra, J., and
765 Springer, N. M. (2015). Transposable elements contribute to activation of maize genes
766 in response to abiotic stress. *PLoS Genetics*, 11(1):e1004915.
767
768 Mao, H. and Wang, H. (2016). Sine_scan: an efficient tool to discover short interspersed
769 nuclear elements (SINEs) in large-scale genomic datasets. *Bioinformatics*, 33(5):743–
770 745.
771
772 Martin, A., Troadec, C., Boualem, A., Rajab, M., Fernandez, R., Morin, H., Pitrat, M.,
773 Dogimont, C., and Bendahmane, A. (2009). A transposon-induced epigenetic change
774 leads to sex determination in melon. *Nature*, 461(7267):1135.
775
776 McCarthy, E. M. and McDonald, J. F. (2003). LTR_STRUC: a novel search and
777 identification program for LTR retrotransposons. *Bioinformatics*, 19(3):362–367.
778
779 Mehra, M., Gangwar, I., and Shankar, R. (2015). A deluge of complex repeats: the
780 *Solanum* genome. *PloS One*, 10(8):e0133962.
781
782 Momose, M., Abe, Y., and Ozeki, Y. (2010). Miniature inverted-repeat transposable
783 elements of stowaway are active in potato. *Genetics*, 186(1):59–66.
784
785 Nussbaumer, T., Martis, M. M., Roessner, S. K., Pfeifer, M., Bader, K. C., Sharma, S.,
786 Gundlach, H., and Spannagl, M. (2012). MIPS PlantsDB: a database framework for
787 comparative plant genome research. *Nucleic Acids Research*, 41(D1):D1144–D1151.
788
789 Nystedt, B., Street, N. R., Wetterbom, A., Zuccolo, A., Lin, Y.-C., Scofield, D. G., Vezzi, F.,
790 Delhomme, N., Giacomello, S., Alexeyenko, A., et al. (2013). The Norway spruce
791 genome sequence and conifer genome evolution. *Nature*, 497(7451):579.
792
793 Ong-Abdullah, M., Ordway, J. M., Jiang, N., Ooi, S.-E., Kok, S.-Y., Sarpan, N., Azimi, N.,
794 Hashim, A. T., Ishak, Z., Rosli, S. K., et al. (2015). Loss of *Karma* transposon
795 methylation underlies the mantled somaclonal variant of oil palm. *Nature*,
796 525(7570):533.
797
798 Ouyang, S. and Buell, C. R. (2004). The TIGR Plant Repeat Databases: a collective
799 resource for the identification of repetitive sequences in plants. *Nucleic Acids Research*,
800 32(suppl 1):D360–D363.
801
802 Paterson, A. H., Bowers, J. E., Bruggmann, R., Dubchak, I., Grimwood, J., Gundlach, H.,
803 Haberer, G., Hellsten, U., Mitros, T., Poliakov, A., et al. (2009). The *Sorghum bicolor*
804 genome and the diversification of grasses. *Nature*, 457(7229):551.
805
806 Paz, R. C., Kozaczek, M. E., Rosli, H. G., Andino, N. P., and Sanchez-Puerta, M. V.
807 (2017). Diversity, distribution and dynamics of full-length Copia and Gypsy LTR
808 retroelements in *Solanum lycopersicum*. *Genetica*, 145(4-5):417–430.
809
810 Platt, R. N., Blanco-Berdugo, L., and Ray, D. A. (2016). Accurate transposable element
811 annotation is vital when analyzing new genome assemblies. *Genome Biology and*
812 *Evolution*, 8(2):403–410.
813

- 814 Quadrana, L., Almeida, J., Asís, R., Duffy, T., Dominguez, P. G., Bermúdez, L., Conti, G.,
815 Da Silva, J. V. C., Peralta, I. E., Colot, V., et al. (2014). Natural occurring epialleles
816 determine vitamin E accumulation in tomato fruits. *Nature Communications*, 5:4027.
817
- 818 Quadrana, L., Etcheverry, M., Gilly, A., Caillieux, E., Madoui, M. A., Guy, J., & Aury, J. M.
819 (2019). Transposition favors the generation of large effect mutations that may facilitate
820 rapid adaptation. *Nature Communications*, 10(1), 3421.
821
- 822 Roessler, K., Bousios, A., Meca, E., and Gaut, B. S. (2018). Modeling interactions
823 between transposable elements and the plant epigenetic response: a surprising reliance
824 on element retention. *Genome Biology and Evolution*, 10(3):803–815.
825
- 826 Saze, H., Kitayama, J., Takashima, K., Miura, S., Harukawa, Y., Ito, T., and Kakutani, T.
827 (2013). Mechanism for full-length RNA processing of *Arabidopsis* genes containing
828 intragenic heterochromatin. *Nature Communications*, 4:2301.
829
- 830 Schnable, P. S., Ware, D., Fulton, R. S., Stein, J. C., Wei, F., Pasternak, S., Liang, C.,
831 Zhang, J., Fulton, L., Graves, T. A., et al. (2009). The B73 maize genome: complexity,
832 diversity, and dynamics. *Science*, 326(5956):1112–1115.
833
- 834 Sigman, M. J. and Slotkin, R. K. (2016). The first rule of plant transposable element
835 silencing: location, location, location. *The Plant Cell*, 28(2):304–313.
836
- 837 Smit, A. F., Hubley, R., and Green, P. (1996). Repeatmasker.
838
- 839 Staton, E. (2018). Tephra a tool for discovering transposable elements and describing
840 patterns of genome evolution.
841
- 842 Sultana, T., Zamborlini, A., Cristofari, G., & Lesage, P. (2017). Integration site selection by
843 retroviruses and transposable elements in eukaryotes. *Nature Reviews Genetics*, 18(5),
844 292.
845
- 846 Sun, S., Zhou, Y., Chen, J., Shi, J., Zhao, H., Zhao, H., Song, W., Zhang, M., Cui, Y.,
847 Dong, X., et al. (2018). Extensive intraspecific gene order and gene structural variations
848 between Mo17 and other maize genomes. *Nature Genetics*, 50(9):1289.
849
- 850 Vicent, C. M. and Casacuberta, J. M. (2017). Impact of transposable elements on
851 polyploid plant genomes. *Annals of Botany*, 120(2):195–207.
852
- 853 Wang, X., Weigel, D., and Smith, L. M. (2013). Transposon variants and their effects on
854 gene expression in *Arabidopsis*. *PLoS Genetics*, 9(2):e1003255.
855
- 856 Wicker, T., Matthews, D. E., and Keller, B. (2002). TREP: a database for Triticeae
857 repetitive elements. *Trends in Plant Science*, 7(12):561–562.
858
- 859 Wicker, T., Sabot, F., Hua-Van, A., Bennetzen, J. L., Capy, P., Chalhoub, B., Flavell, A.,
860 Leroy, P., Morgante, M., Panaud, O., et al. (2007). A unified classification system for
861 eukaryotic transposable elements. *Nature Reviews Genetics*, 8(12):973.
862
- 863 Xiao, H., Jiang, N., Schaffner, E., Stockinger, E. J., and Van Der Knaap, E. (2008). A
864 retrotransposon-mediated gene duplication underlies morphological variation of tomato
865 fruit. *Science*, 319(5869):1527–1530.
866

867 Xiong, W., He, L., Lai, J., Dooner, H. K., and Du, C. (2014). HelitronScanner uncovers a
 868 large overlooked cache of Helitron transposons in many plant genomes. Proceedings of
 869 the National Academy of Sciences, page 201410068.

870
 871 Yang, L. and Bennetzen, J. L. (2009). Distribution, diversity, evolution, and survival of
 872 Helitrons in the maize genome. Proceedings of the National Academy of Sciences,
 873 106(47):19922–19927.

874
 875 Yang, Z. (2007). PAML 4: phylogenetic analysis by maximum likelihood. Molecular Biology
 876 and Evolution, 24(8):1586–1591.

877

878

879

880

881

882

883 **Table 1:** TEs filtering parameters. Parameters established to filter copies element for
 884 “Autonomous TEs” and “All TEs” which includes autonomous and non-autonomous
 885 elements. *min-len*: minimum element length; *max-len*: maximum element length; *min-*
 886 *pid*: minimum query identity percentage; *min-threshold*: minimum threshold length
 887 between query and subject; *max-threshold*: maximum threshold length between query
 888 and subject; *min-cov*: minimum query coverage percentage; *overlap*: margin of plus or
 889 minus nucleotides overlap between query and subject to be considered as duplicates.

Family	min-len (nt)	max- len (nt)	min-pid (%)	min- threshol d (ratio)	max- threshol d (ratio)	min- cov (%)	overlap (nt)
Autonomous TEs							
LTR	650	∞	95	0.9	1.1	95	5
LINE	1500	∞	95	0.9	1.1	95	5
TIR	700	∞	95	0.9	1.1	95	5
Helitron	2000	∞	95	0.9	1.1	95	5
All TEs							
LTR	650	∞	80	0.5	1.5	50	5
LINE	1500	∞	80	0.5	1.5	50	5
SINE	150	800	80	0.8	1.2	80	5
TRIM	600	∞	90	0.9	1.1	90	5
LARD	4000	∞	80	0.5	1.5	50	5
TIR	700	∞	80	0.5	1.5	50	5
MITE	50	800	80	0.85	1.15	90	5
Helitron	2000	∞	80	0.5	1.5	50	5

890

891

892

893

894

895

896

897

898

899

900 **Table 2:** Quantity and diversity of all identified TEs. Number of TE families, number of
 901 TE copies and TEs genome coverage (%) of all types of TEs.

Class	Order/Family	TEs identified	TEs copies	Genome coverage (%)
I	LTR/Copia	2,541	13,044	2.46
I	LTR/Gypsy	5,736	170,380	10.68
I	LTR/Unclassified	29	604	0.06
	Total LTRs	8,306	184,028	13.2
I	LINE	65	248	0.1
I	SINE	24	2,477	0.07
I	TRIM	6,908	13,491	0.67
I	LARD	2	18	0.002
Total Class I		15,305	200,262	14.04
II	TIR/hAT	112	382	0.05
II	TIR/Mariner	1,466	1,591	0.29
II	TIR/Harbinger	16	219	0.02
II	TIR/Mutator	575	941	0.13
II	TIR/CACTA	131	117	0.02
II	TIR/Unclassified	3	110	0.01
	Total TIRs	2,303	3,380	0.51
II	MITEs	3,515	38,205	0.72
II	Helitrons	1,322	1,163	0.72
Total Class II		7,140	42,748	1.99
Total TEs		22,445	243,010	16.03

902

903

904

905 **Figure legends**

906

907 **Figure 1.** Overview of the TE Discovery pipeline. *1_Input_data:* last genome assembly,
 908 available TE sequences from Repbase, Gini and sRNA-seq Illumina data are used as
 909 input data. *2_TEs family detection:* the detection of putative TEs is performed by using
 910 four tools combining two detection approaches. The resulted sequences are merged
 911 into multi-fasta files for each type of TEs. *3_TEs clustering:* Each TE type sequences
 912 are clustered with Vsearch to reduce redundancy. *4_TEs copy detection:* blastn is
 913 performed to the clustered sequences with specific parameters according to each type
 914 of TE to detect copies across the genome. *5_TEs copy filter:* the detected copies are
 915 subjected to filtering steps to detect “potential autonomous TEs” and/or “All TEs”
 916 including non- autonomous TEs. *6_TEs annotation:* TEs annotation *gff3* files are
 917 generated for each type of TEs with detailed descriptions and merged into one single
 918 file.

919

920 **Figure 2.** Comprehensive circos ideograms of TEs of the potato genome. Left panel:
 921 Retrotransposon TEs (Class I). Right panel: DNA TEs (Class II). Each concentric circle
 922 represents a different type of TE with their own color pallet and range of coverage so the
 923 distribution across the chromosomes can be appreciated. Each line within the circles

924 represents coverage percentage per Mb. (1) Gene distribution. (a) LTR distribution. (b)
925 SINE distribution. (c) LINE distribution. (d) TRIM distribution. (e) TIR distribution. (f)
926 MITE distribution. (g) Helitron distribution.

927

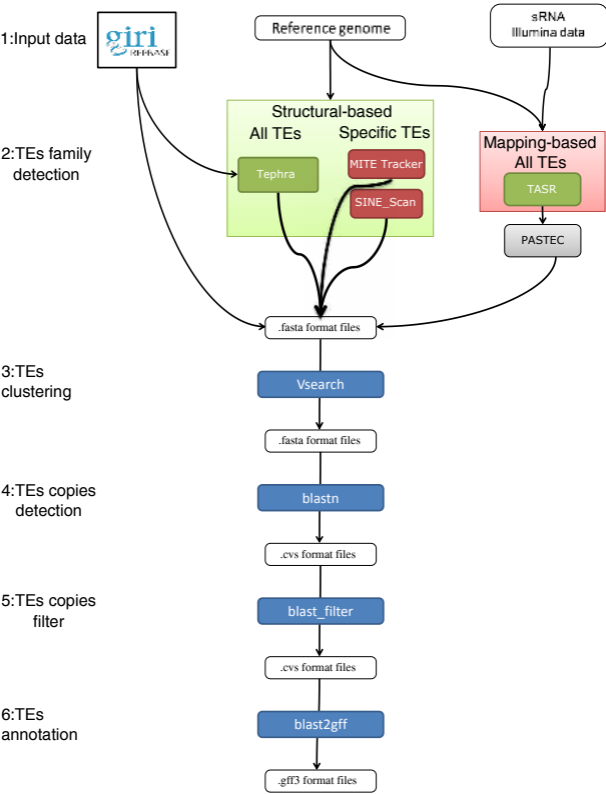
928 **Figure 3.** Frequency histograms of TE distance to the nearest gene. Left panel:
929 Retrotransposon TEs (Class I). Right panel: DNA TEs (Class II). Bars on the left side of
930 the red line represent the TE distance within the first 5kb to the nearest gene.

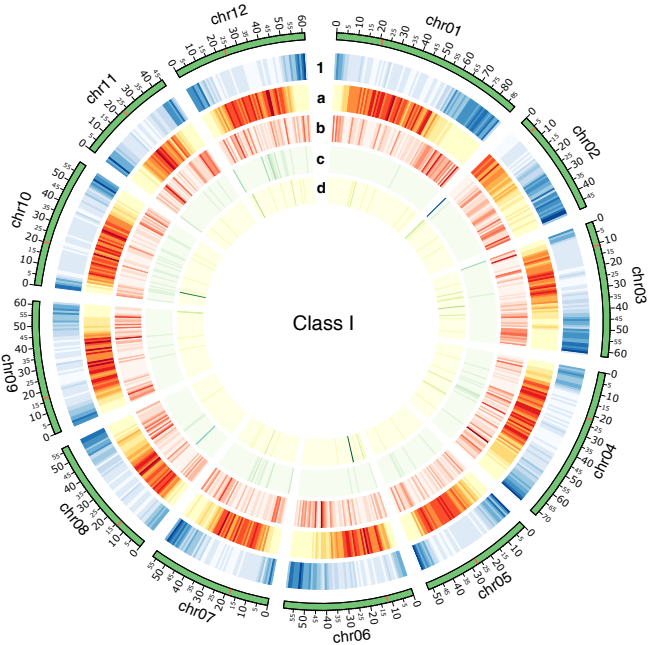
931

932 **Figure 4.** Chromosomal ideograms of LTRs age per family. Upper panel: chromosomal
933 distribution of insertion age of all superfamilies of full-length LTR retrotransposons.
934 Lower panel: chromosomal distribution of insertion age of Gypsy, Copia and
935 Unclassified LTR retrotransposons separately.

936

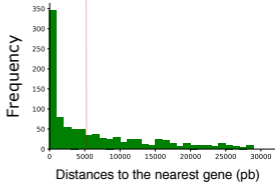
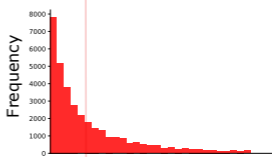
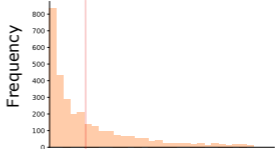
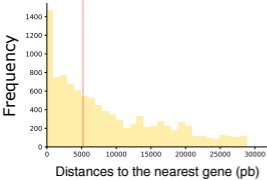
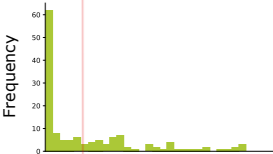
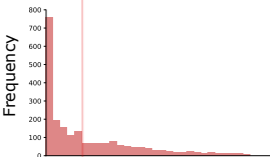
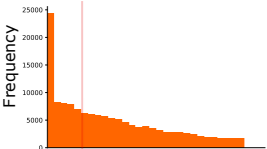
937 **Figure 5.** Number of LTR family by insertion age according to genome chromatin state.
938 Upper panel: frequency histograms of Gypsy, Copia and Unclassified LTR families by
939 their insertion age (in millions of years) separated by heterochromatin regions (blue) and
940 euchromatic regions (red). Lower panel: Scatter plots of ten random independent
941 families from Gypsy, Copia and Unclassified LTRs assessing the age distribution of
942 individual elements by their heterochromatin (blue) and euchromatin (red) state. * shows
943 significant differences between families (t-test $p < 0.05$).

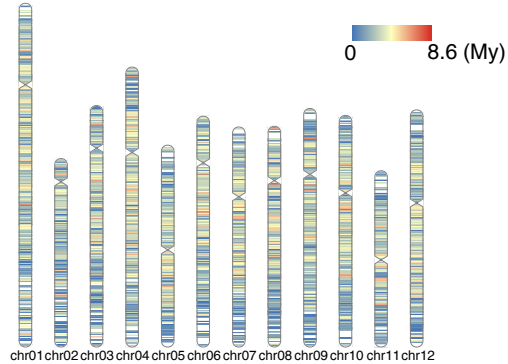




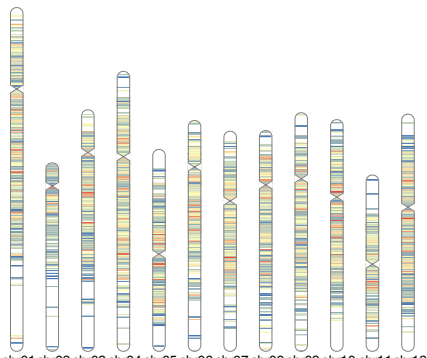
1 Gene distribution 0.5 39.8
 a LTR distribuion 0 36.1
 b SINE distribuion 0 0.36
 c LINE distribuion 0 2.89
 d TRIM distribuion 0 4.29
 coverage (% per Mb)

1 Gene distribution 0.5 39.8
 e TIR distribuion 0 2.34
 f MITE distribuion 0 2.09
 g Helitron distribuion 0 5.11
 coverage (% per Mb)

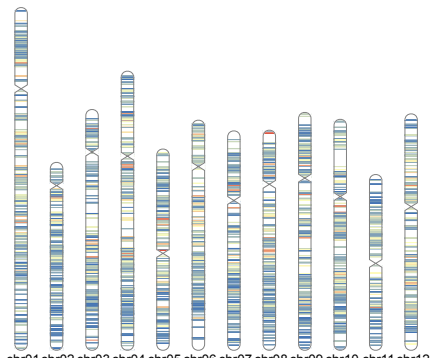




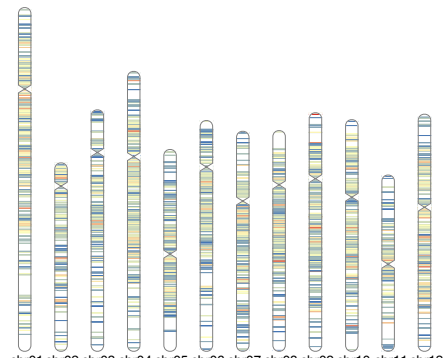
all LTRs



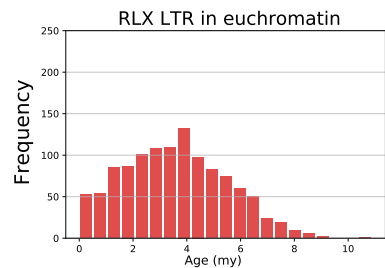
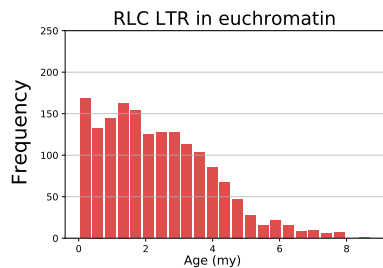
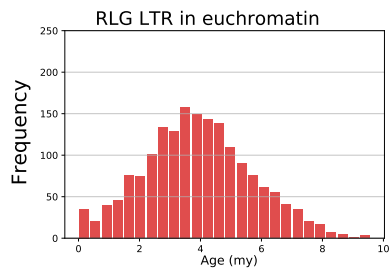
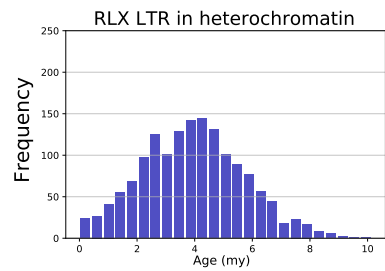
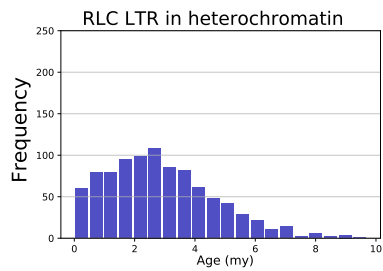
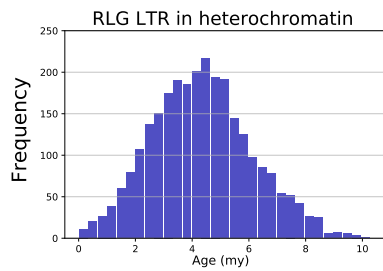
LTR Gypsy (RLG)



LTR Copia (RLC)



LTR Unclassified (RLX)



■ heterochromatin
■ euchromatin

