

1 **T-lex2: genotyping, frequency estimation and re-annotation of transposable**
2 **elements using single or pooled next-generation sequencing data**

3

4 Anna-Sophie Fiston-Lavier^{1,2}, Maite G. Barrón^{3,4}, Dmitri A. Petrov¹ and Josefa González⁴

5

6 ¹Department of Biology, Stanford University, Stanford, CA 94305-5020, USA

7 ²Institut des Sciences de l'Evolution (ISEM), UMR5554 CNRS-Université Montpellier 2, France

8 ³Genomics, Bioinformatics and Evolution Group, Institut de Biotecnologia i de Biomedicina -

9 IBB/Department of Genetics and Microbiology, Campus Universitat Autònoma de Barcelona, Bellaterra,

10

08193, Spain.

11

⁴Institut de Biologia Evolutiva (CSIC-Universitat Pompeu Fabra), Barcelona, 08003, Spain.

12

13

14 Corresponding author:

15 Anna-Sophie Fiston-Lavier, PhD

16 Assistant professor, University Montpellier2,

17 Adaptation Genomic team, Institute of Sciences of Evolution (ISEM),

18 Bâtiment 22, 1er étage,

19 34095 Montpellier, FRANCE

20 +33 0467143289

21 asfiston@univ-montp2.fr

22

23

24

25 Abbreviations:

26 TE, transposable element

27 NGS, next-generation sequencing

28 LTR, long-terminal repeat

29 TSD, target site duplication

30 PTS, putative target site

31 PE, paired-end

32

33

34

35

36

37

38

39

40 **ABSTRACT**

41

42 Transposable elements (TEs) are the most active, diverse and ancient component
43 in a broad range of genomes. As such, a complete understanding of genome function
44 and evolution cannot be achieved without a thorough understanding of TE impact and
45 biology. However, in-depth analyses of TEs still represent a challenge due to the
46 repetitive nature of these genomic entities. In this work, we present a broadly applicable
47 and flexible tool: T-lex2. T-lex2 is the only available software that allows routine,
48 automatic, and accurate genotyping of individual TE insertions and estimation of their
49 population frequencies both using individual strain and pooled next-generation
50 sequencing (NGS) data. Furthermore, T-lex2 also assesses the quality of the calls
51 allowing the identification of miss-annotated TEs and providing the necessary
52 information to re-annotate them. Although we tested the fidelity of T-lex2 using the high
53 quality *Drosophila melanogaster* genome, the flexible and customizable design of T-
54 lex2 allows running it in any genome and for any type of TE insertion. Overall, T-lex2
55 represents a significant improvement in our ability to analyze the contribution of TEs to
56 genome function and evolution as well as learning about the biology of TEs. T-lex2 is
57 freely available online at <http://petrov.stanford.edu/cgi-bin/Tlex.html>.

58

59 INTRODUCTION

60

61 The key question in genomics is how genomes vary and evolve at both large and
62 fine scales. In order to answer such question, we need to be able to study genome
63 variation, *i.e.* identifying and analyzing functionally relevant SNPs (Single Nucleotide
64 Polymorphisms) and structural variants, both within and between populations. Next-
65 generation sequencing (NGS) technology has revolutionized this field by allowing one
66 to study variation of a large number of individuals and even cells within individuals (1-
67 3). Unfortunately, NGS technology generally provides fairly short, even if extremely
68 abundant, sequencing reads and thus is not perfectly designed for the study of structural
69 variants (4).

70 The study of one class of structural variants, transposable elements (TEs), is
71 particularly difficult to carry out using NGS data. TEs are repetitive, ubiquitous and
72 dynamic components of genomes that often vary in genomic location among members
73 of the same species. TEs are classified in three orders: DNA transposons, long-terminal
74 repeat (LTR) elements, and non-LTR elements, and each TE order is represented by
75 several TE families (5,6). This classification highlights the biological diversity in terms
76 of sequence, dynamics and evolution of these repetitive elements that makes their
77 analysis challenging.

78 TEs represent a large part of most of the eukaryotes genomes (7). A recent study
79 suggests that more than two-thirds of the human genome is composed of TEs (8), and in
80 plants, TEs may represent up to 90% (9). TEs are responsible for a large number of
81 mutations both in populations and somatically within individuals (10,11). Although
82 most of the TE-induced mutations are deleterious, evidence for an adaptive role of TE-
83 induced mutations is starting to accumulate (12,13). TEs can provide active promoter,

84 splice site, or terminator features that can affect the expression, structure and function of
85 nearby genes (14). TEs are also involved in the creation of transcriptional regulatory
86 networks, and in the generation of chromosomal rearrangements (12,15). Thus, knowing
87 the role of TEs on genome dynamics and evolution, it is crucial to identify and quantify
88 their impact (16,17). Given the increased number of sequenced genomes, a tool that
89 allows routine and automatic genotyping of individual TE insertions using a large
90 number of NGS samples is needed.

91 Recently several accurate computational approaches have been developed to
92 analyze TE insertions using NGS data (10,18-22). Designed for specific projects, all
93 these tools were built with a limited number of features, *e.g.* some tools were designed
94 to specifically call particular families of TEs (10), while others were designed
95 exclusively for pooled NGS data (22). Unfortunately, none of the currently available
96 tools allows a complete and in-depth analysis of individual annotated TE insertions.

97 With the idea of providing a tool to automatically call TE insertions in multiple
98 genomes and to estimate their population frequency, we recently designed an integrated,
99 flexible pipeline called T-lex (23). While most of the available approaches only detect
100 the presence of TE insertions, T-lex launches two complementary and distinct TE
101 detection modules: one to detect the “presence” and another to detect the “absence” of a
102 TE insertion. Both detection approaches are based on the analysis of the junction
103 sequences of a TE insertion and its flanking sequences as defined by genome annotation.
104 Thus, the accuracy of the TE calls mainly depends on the quality of the annotation and
105 more specifically on the delimitation of the individual TE insertions. Unfortunately, the
106 detection of the precise TE insertion sites is still one of the main challenges of the TE
107 annotation process (24). The veracity of the TE calls also depends on the genomic
108 environment of the TE insertion. For instance, a TE insertion located inside a

109 duplication with one of the copies lacking the TE insertion can be miscalled as
110 heterozygote (*i.e.*, both present and absent). In addition, a TE insertion flanked by other
111 repeats such as low-complexity regions or other TE sequences can also be miscalled.

112 Here, we present T-lex2, the version 2 of the original T-lex pipeline (23).
113 Besides improving the accuracy of the TE frequency estimate, T-lex2 now assesses the
114 quality of the TE calls. To achieve this, T-lex2 uses information from (i) a new module
115 specifically designed to identify target site duplications (TSDs), and (ii) the genomic
116 context of each insertion, to identify putatively miss-annotated TE insertions highly
117 likely to produce wrong calls. TSDs are short duplications, from two base pairs to 20 bp,
118 flanking the TE sequences as a result of the transposition mechanism of most TE
119 families. While current TSD detection approaches require knowing the biology of the
120 TEs and are limited to LTR or non-LTR elements, T-lex2 TSD detection module works
121 without *a priori* knowledge, and for all type of TEs. This new module allows to
122 precisely delimit the location of each TE insertion. In this new version, frequency of TE
123 insertions can be estimated using individual and pooled NGS data.

124 We provide evidence that the new features of T-lex2 help improving the
125 genotyping and frequency estimate of TE insertions, and allow re-annotating TE
126 locations using both individual and pooled *Drosophila melanogaster* genomic data.

127

128 MATERIAL AND METHODS

129

130 T-lex2 pipeline overview

131 T-lex2 is composed of five modules that can be run with individual strain or
132 pooled NGS data (Figure 1). (i) The TE-analysis module investigates the flanking
133 regions of the known TE insertions (*i.e.* annotated in a reference genome) to identify
134 those likely to return wrong calls. (ii) The TE-presence detection and (iii) the TE-
135 absence detection modules combine mapping and read depth coverage information to
136 identify reads providing evidence for the presence and for the absence, respectively, of
137 the known TE insertions. (iv) Then the TE-combine module combines the results of the
138 presence and absence detection modules to genotype TEs and/or to estimate their
139 frequencies in populations. Finally, (v) the TE-TSD detection module annotates TSDs in
140 an unbiased and accurate way. Each of these modules can be launched independently
141 and provides a large number of options allowing the user to carry out flexible and
142 customizable analyses. A detailed manual describing step-by-step how to run T-lex2 and
143 listing all Tlex-2 options is available on line at [http://petrov.stanford.edu/cgi-](http://petrov.stanford.edu/cgi-bin/Tlex2_manual.html)
144 [bin/Tlex2_manual.html](http://petrov.stanford.edu/cgi-bin/Tlex2_manual.html).

145

146 **TE-analysis module.** This module analyses the flanking sequences of each known TE
147 insertion to identify features that might interfere with the presence and absence detection
148 modules. The TE-analysis module (i) provides information about the presence of
149 repetitive elements using RepeatMasker (25), (ii) identifies miss-annotated poly-A tails,
150 and (iii) flags the TE insertions that are part of segmental duplications.

151 (i) To identify the presence of repetitive elements, the TE-analysis module extracts the
152 flanking sequences of each TE insertion (125 bp by default; see option “-f”) and

153 launches RepeatMasker to annotate the TEs, simple repeats and low-complexity regions
154 (25). The option “-s” allows specification of the name of the model organism and thus
155 specifies the repeat library that needs to be used by RepeatMasker (25). If at least one
156 flanking region shows a repeat density greater than the pre-specified value (50% by
157 default; option “-d”) the TE insertion is flagged as flanked by repetitive elements. By
158 default, these TE insertions are filtered out. This filtering step can be bypassed using the
159 option “-noFilterTE”.

160 (ii) A new feature of the TE-analysis module identifies putatively miss-annotated poly-A
161 tails by searching for stretches of A’s or T’s (by default more than five base pairs) at the
162 TE junctions, without *a priori* knowledge of the TE type. Non-LTR elements harbor
163 poly-A tails at their 3’ end that are known to be highly variable in length. Such
164 sequences are therefore very difficult to annotate automatically. As RepeatMasker only
165 annotates repeats longer than 20 bp, it cannot be used to identify short stretches of A’s
166 or T’s that may correspond to miss-annotated poly-A tail sequences (25).

167 (iii) When a TE insertion is part of duplication, other copies of the same duplication may
168 not contain the insertion. The analysis of the flanking sequences of such TE insertions
169 may lead to call this TE insertion as present and absent while it should only be called as
170 present. We added one new feature in the analysis of the TE flanking regions to identify
171 the TE insertions fully or partially part of segmental duplications. For each TE insertion,
172 the two extracted flanking sequences are blated against a reference genome (26) (BLAT
173 v. 34 default parameters; <http://www.soe.ucsc.edu/~kent>). When more than 50% of a
174 flanking sequence matches at more than one location on the reference genome with
175 more than 90% of sequence identity, the TE insertion is flagged as part of a duplication.
176 TE insertions partially part of duplication (notified as “sd_left” or “sd_right”) are
177 distinguished from the TE fully part of a duplication (notified as “sd”). When a TE is

178 fully part of a duplication while another copies do not contain the TE insertion, this TE
179 insertion is notified as “sd_noTE”.

180 The analysis of the TE flanking sequences including the RepeatMasker output, whether
181 the TE is flanked by a longer than annotated poly-A tail, and the segmental duplication
182 analysis results, are stored in a sub-directory called “Tanalyses”.

183

184 **TE-presence detection module.** The presence detection module detects sequencing
185 reads overlapping the flanking regions of the TE insertion (Figure 2A). The junction
186 sequences of each TE insertion are extracted and used as reference sequences. By
187 default, each TE junction sequence encompasses the terminal TE sequence (60 bp from
188 each end) and the flanking sequence (1 kilo base on each side). The lengths of the TE
189 sequence and the flanking sequence can be changed using the options “-b” and “-j”
190 respectively. While the original T-lex used a 20 bp region of the TE sequence, T-lex2
191 uses a 60 bp region at this step. The increased length improves sensitivity and can now
192 be used given the increased read lengths in the NGS data compared to even a few years
193 ago. The reads are mapped on these new extracted reference sequences (Figure 2A) and
194 are used to build a contig using Phrap (27) (version 1.090518; <http://phrap.org>). While
195 the original T-lex uses all the mapped reads at this step, T-lex2 increases the stringency
196 by using only the reads mapping with a minimum Phrep quality score (30 by default;
197 this can be changed using the option “-minQ”). The pipeline also reports now the
198 number of reads mapped on the TE junction at the minimum Phrep quality score (Figure
199 2A). Besides increasing the minimum length of the match inside and outside the TE
200 sequence that was required by the original version of the pipeline in order to return a
201 “presence” call, the new version of T-lex also requires a minimum read sequence
202 identity (Figure 2A). These two parameters are set at 15 bp and 95% by default and can

203 be modified using the options “-limp” and “-id”, respectively. By default, multiple
204 alignments of the contigs for each TE side are also returned. The analysis and TE calls
205 from this module are stored in a sub-directory called “Tpresence”.

206

207 **TE-absence detection module.** The absence module detects reads overlapping both
208 sides of the TE insertion, *i.e.* reads spanning across the TE insertion site (Figure 2B). It
209 starts by extracting the two flanking sequences for each TE insertion (125 bp each by
210 default; see option “-f”). The two sequences are then concatenated and the new
211 constructed sequence is used as the reference sequence for the absence detection (Figure
212 2B). This new constructed sequence approximates the ancestral sequence prior to the TE
213 insertion. However, note that traces of the TE insertion mechanism, called TSDs, may
214 also be encompassed in the new constructed sequence and located at the TE insertion
215 breakpoint. The length of the extracted sequences should be longer than the reads
216 themselves and total length of the new constructed sequence should be similar to the
217 insert size of the paired-end (PE) data in order to get reads and/or pairs spanning the TE
218 insertion breakpoints. T-lex2 uses now by default a length of 125 bp matches (instead of
219 100 bp in the original T-lex) for libraries of 100 bp read length and an insert size of 250
220 bp. The reads are then mapped on the new constructed sequence using SHRiMP2 (28)
221 (version 2.2.1, October 2011) as opposed to the original version of SHRiMP (29) used in
222 the original T-lex version. One major improvement of SHRiMP2 is that it can now
223 handle PE read sequencing data.

224 SHRiMP2, as SHRiMP, was specifically designed to handle long gaps and
225 polymorphisms (28,29). Such features allow the mapping of reads despite the presence
226 of the TSDs or despite miss-annotation. To handle long gaps, SHRiMP2 is launched
227 with a Smith-Waterman gap open score of -40 and a gap extension score of -1 for both

228 query and reference sequences. We also allow higher divergence with a Smith-
229 Waterman mismatch score of -20. SHRiMP2 maps the reads in every pair together (28)
230 (see option “-pairends”). If a pair does not map, SHRiMP2 attempts to map each read
231 individually (28). If the two reads of a pair map to the new constructed reference
232 sequence, only the pairs with reads mapping on opposite strands are selected (*i.e.*,
233 concordant pairs). When the two reads from the pair overlap in the mapping, the absence
234 module merges them and considers the pair as a single read. The absence module
235 specifically looks for “reads” spanning the TE junction with at least 15 bp (parameters
236 by default; see option “-v”) of overlap on each side. T-lex2 then runs RepeatMasker on
237 the selected reads (see option “-noit”) in order to test whether the mapping is due
238 primarily to simple repeats or low-complexity regions (25). After this step, T-lex2
239 selects the reads that have at least five non-repetitive (low-complexity or satellites) base
240 pairs at each end that both map to the flanking regions. Because of the stringency of our
241 approach, even a single read mapping to both flanking sides and thereby spanning the
242 insertion site of the TE is sufficient for the absence detection module to classify the TE
243 insertion as “absent” (Figure 2B). If no reads overlap the TE junction, two possibilities
244 are considered: (i) the TE is present in the strain (or fixed in the pool sample) or (ii) the
245 coverage of the data used is insufficient to detect reads providing evidence for the
246 absence. In order to distinguish between these two possibilities, the absence detection
247 module checks whether reads map on both flanking sides but fail to map over the
248 junction. If this is the case, given that the flanking regions are similar in length to the
249 reads themselves, the module concludes that the TE is present.
250 If reads do not map to both flanking sides and do not map over the junction site, the
251 module concludes that the coverage is insufficient and returns a “no call” as a result. By
252 default, the multiple alignments of the contigs for each TE breakpoint are also returned.

253 The analysis and TE calls from this module are stored in a sub-directory called
254 “Tabsence”.
255
256 **TE-combine module.** The results from the presence and absence detection approaches
257 can then be combined to return definitive TE detection calls. We defined five TE call
258 categories based on the evidence for presence and absence. Briefly, a TE is called
259 “present” or “absent” when the calls from both detection modules are congruent. A TE
260 is called as “polymorphic” when this TE is clearly detected as present and absent (*i.e.*,
261 the presence detection module detects it as present and the absence detection module
262 detects it as absent). A TE is called as “present/polymorphic” when the presence
263 detection module detects it as present while the absence detection module fails to return
264 a call. A TE is called as “absent/polymorphic” if the TE is detected as absent by the
265 absence detection approach, while the presence detection fails (see Table 1 in (23)).

266 T-lex2 can also combine the TE calls from several NGS data from the same
267 population and can estimate the population frequency of each annotated TE insertion
268 (see options “-combRes” and “-combine”). In T-lex2, the probability for a TE to be
269 present is associated with each of the above categories as follows: “present” 100%,
270 ”present/polymorphic” 75%, ”polymorphic” 50%, ”absent/polymorphic” 25%, ”absent”
271 zero percent, respectively.

272 Finally, T-lex2 pipeline can also estimate the TE frequency using pooled NGS
273 data by taking into account the number of reads supporting the presence and the number
274 of reads supporting the absence. In this case, the user just needs to add in the command
275 line the option “-pooled”. This approach is based on two assumptions: (i) individuals in
276 the pooled sample are equally represented, and (ii) the read-depth coverage is correlated
277 with the occurrence of a TE sequence in a population. Based on such assumptions, we

278 expect to observe a positive correlation between the frequency and the number of reads
279 providing evidence for the presence. We also expect to observe a negative correlation
280 between the frequency and the number of reads providing evidence for the absence. The
281 number of copies of TEs from the same family may also have an impact on the number
282 of reads mapping at the TE breakpoints suggesting that the number of reads supporting
283 the presence can be more biased than the number of reads supporting the absence.
284 Taking all these expectations into account, we designed and verified a metric based on
285 the local read depth coverage in the vicinity of the TE insertion to estimate TE
286 frequencies from pooled data: total number of reads supporting the presence (NP)
287 divided by the total number of reads supporting the presence (NP) and supporting the
288 absence (NA), *i.e.* TE frequency = NP/(NP+NA).

289 The file combining the TE calls from both TE detection modules is called “Tresults”.
290 The TE frequency estimates are stored in the “Tfreq” or Tfreq_pooled” if the option “-
291 pooled” is specified.

292

293 **TE-TSD detection module.** To detect the TSD of each annotated TE insertion, this
294 module analyzes the read alignments generated by the absence detection module (Figure
295 3). As all T-lex2 modules, this new module can be launched independently although it
296 does require the read alignments from the absence detection module (see option “-tsd”).
297 The TSD detection module looks for tandem duplications located on the absence
298 reference sequence spanning the TE breakpoint. It starts by assembling all the selected
299 reads supporting the absence for each TE using the Phrap program (27) (Figure 3).
300 Because Phrap requires a minimum of three sequences to build a contig, if fewer than
301 three reads are selected to support the absence call, the reads are considered
302 independently (Figure 3). Each contig (or read) is then re-aligned on the absence

303 reference sequence using BLAT v.34 program (26) (default parameters; Figure 3). Only
304 the absence calls with a gap larger than two base pairs and located close to the TE
305 breakpoint are selected and analyzed. If no gaps are observed, the TE-TSD detection
306 module returns “no Gap”. If a gap is observed, the motif present in front of the gap is
307 called the putative target site or PTS (Figure 3). Using the FastaGrep program (available
308 from bioinfo.ut.ee/download/dl.php?file=6 and executed with the default parameters), a
309 tool that looks for short and conserved motifs, the TSD detection module looks for the
310 copy of the PTS on the other side of the breakpoint. Based on the sensitivity of this
311 program and the short length of the TSDs, we only trust the matches with more than
312 80% sequence identity (Figure 3). T-lex2 then sorts the FastaGrep hits to identify the
313 sequence showing the closer and highly conserved motif. If FastaGrep fails to detect a
314 match, the TSD detection returns “noTSD” as a result, otherwise, the TSD detection
315 module returns “detected”. In the latter case, the PTS and its closest copy sequences are
316 also returned. Because several contigs can be generated for the same TE insertion,
317 several TSDs can also be returned. The results of the TSD detection process are stored
318 in a file called “TSDannot”.

319 TE-TSD detection module flags the TE insertions for which the TSDs cannot be clearly
320 defined. TSD detection may fail for old TE insertions for which the TSD is too
321 divergent, for truncated TEs, and when the boundaries of the TE insertions are not well
322 annotated (Figure 4). The latter TEs can be re-annotated by the user after manual
323 inspection.

324

325 **Re-annotation of TE insertions**

326 Because both TE detection approaches are based on the mapping of NGS reads to the
327 flanking TE sequences, the accuracy of the TE calls depends on having a correct TE

328 annotation. TE-analysis and TE-TSD detection modules help identifying putatively
329 miss-annotated TE insertions and thus putatively wrong T-lex2 calls.
330 A particular TE insertion can be miss-annotated because it is longer than annotated
331 (Figure 4A.II) or shorter than annotated (Figure 4A.III) in both cases the TE-TSD
332 detection will fail. If the TE insertion is present in a given genome, both the presence
333 detection and the absence detection modules will provide the correct call even though
334 the TE is longer or shorter than initially annotated (Figure 4B). However, if the TE
335 insertion is absent, and the TE is miss-annotated the presence and the absence modules
336 could return a wrong call (Figure 4B). For example, if the TE sequence is longer than
337 the official annotation and the TE is absent, the gap in the reads providing the absence
338 could be longer than the TSD size and the module could return present as a call (Figure
339 4B.k). On the other hand, when the TE is shorter than annotated and the TE is absent,
340 the gap will be present in the reference sequence and not in the reads. In such situation,
341 the module could return present as a call (Figure 4B.l).

342 The Tfreq_combine file of T-lex2 provides information from both the TE-
343 analyses and TE-TSD detection modules, *i.e.* polyA tail longer than the official
344 annotation and TSD detection failures, that allows classifying these TE insertions as
345 putatively miss-annotated TEs. This information can be use to perform manual curation
346 of these insertions and re-annotate them.

347

348 **Validation dataset**

349 Veracity of T-lex2 calls was assessed by calling a total of 755 well-studied TE
350 insertions in *Drosophila melanogaster* strains from the Drosophila Genome Reference
351 Panel (DGRP) project (30,31). Estimates of TE frequency based on pooled-PCR are
352 available for these 755 TEs (31,32). TE frequency estimates based on pooled-PCR

353 approach were categorized in four classes: “very rare” (*i.e.* TE frequency < 1.5%),
354 “rare” (*i.e.*, TE frequency ~2-15% but not absent), “common” (*i.e.*, TE frequency ~10-
355 98%) and “fixed” (*i.e.*, TE frequency > 98%) (31). The *D. melanogaster* reference
356 genome (release 5) was downloaded from Flybase website (<http://flybase.org>). The
357 annotations of the TE insertions were extracted from the release 5.43
358 (<http://flybase.org>). The resequencing data from 86 *D. melanogaster* DGRP single
359 strains from the freeze 2 were downloaded from the DGRP website (~20X;
360 <http://dgrp.gnets.ncsu.edu/>, Supplementary Table S1). We also called the same TE
361 insertions in the DGRP population using a high-depth sequencing data from a pool of 92
362 *D. melanogaster* DGRP strains (~60X) (33). 84 strains are common in the two datasets
363 (Supplementary Table S1). To investigate the effect of different read-depth coverage
364 and different number of flies in pooled samples, we also used T-lex2 to analyzed TE
365 frequencies in four other pools of lower depth coverage described in Zhu *et al.* 2012
366 (33). All the NGS data used in this work was generated using the *Illumina* technology.
367 Experimental validation of T-lex2 results was also done for 11 TE insertions known to
368 be polymorphic (32) in 24 single *D. melanogaster* North American DGRP strains using
369 a single PCR approach (Supplementary Table S1). The PCR conditions and primers
370 used were as described in Gonzalez *et al.* 2008 (32).

371

372 **Statistical analysis**

373 All the statistical analyses were performed using R (<http://cran.r-project.org/>).
374 Sensitivity and specificity were estimated using individual PCR frequencies as we have
375 the T-lex2 frequencies estimates using the same exact strains. Sensitivity and specificity
376 were computed such as: (i) Sensitivity = number of correct presence calls / total number
377 of presence calls and (ii) Specificity = number of correct absence calls / total number of

378 absence calls.

379

380

381 **RESULTS**

382

383 **T-lex2 is an improved and expanded version of T-lex**

384 T-lex was originally designed to genotype and estimate frequencies of the known
385 TE insertions in individual strain NGS data (23). The new version of T-lex incorporates
386 numerous new features in order to improve genotyping and frequency estimation of
387 known TE insertions (Table 1). T-lex2 now assesses the quality of the calls by analyzing
388 the genomic context of each TE insertion and detecting their TSDs that allows
389 identifying miss-annotated TE insertions. The new version of T-lex takes into account
390 the sequencing technology advances and can now handle PE sequencing and longer
391 reads. Unlike the previous version, T-lex2 can be run for individual strains and for
392 pooled NGS for the TE frequency estimates (Table 1).

393

394 **Accurate TE frequency estimates using NGS data from individual strains**

395 To test the quality of the T-lex2 calls, we compared T-lex2 TE frequency
396 estimates with (i) frequency estimates previously obtained using a pooled-PCR approach
397 (31,32,34) and (ii) frequency estimates obtained in this work by individual strain PCRs.

398 We first launched the pipeline to estimate the frequency of 755 TE insertions for
399 which we have an estimate of their population frequencies based on a pooled-PCR
400 approach (31,32,34). We run T-lex2 on 86 single DGRP (Supplementary Table S1; see
401 Material and Methods). The TE call distribution, *i.e.* proportion of TEs called as present,
402 absent, polymorphic, present/polymorphic, absent/polymorphic and no calls, is

403 consistent among the strains, excepted for 14 strains that differ drastically
404 (Supplementary Figure S1). For 11 out of the 14 strains, we observed more than 40% of
405 no calls, and for the other three strains (RAL-379, RAL-362 and RAL-313), more than
406 70% of TE insertions were detected as present. As we suspect these 14 strains to be
407 problematic (most likely a combination of poor sequence quality and extensive residual
408 heterozygosity), we decided to remove them from our analysis and analyze the calls in
409 the remaining 72 strains (Supplementary Figure S1).

410 Out of the 755 TE insertions, TE-analysis module detected six that were flanked
411 by other TE insertions, and three that were part of segmental duplications
412 (Supplementary Table S2). We removed these nine insertions because they are likely to
413 produce wrong calls (see Material and Methods). We then filtered out seven TEs for
414 which T-lex2 failed to return any calls, and 30 TE insertions for which T-lex2 fails to
415 return a call for more than 75% of the strains (Supplementary Table S2). All of these
416 removed TE insertions, 46 in total, are likely to be miss-annotated as we failed to call
417 them in all or most of the strains. We thus estimated the population frequencies of 709
418 TE insertions using 72 strains. Note that 65 out of the 709 TEs were identified in a first
419 T-lex2 run as miss-annotated based on results of the TE-TSD detection module (see
420 below). We manually re-annotated these 65 TE insertions, and we have re-run T-lex2 to
421 obtain accurate frequency estimates of these TEs.

422 T-lex2 frequency estimates of the 709 TEs in 72 strains are significantly and
423 positively correlated with previous estimates based on a pooled-PCR approach (Figure
424 5A: Spearman's $\rho = 0.87$, $P\text{-val} \ll 0.001$). The vast majority of the very rare and fixed
425 TE insertions are correctly detected by T-lex2 with 92.5% (357/386) of the very rare TE
426 insertions absent and 85% (91/107) of the fixed TE insertions fixed in the population
427 (Figure 5A). As the pooled-PCR and T-lex2 frequencies categories were not build the

428 same way, we expected to observe intermediate profiles for the rare and common TE
429 insertions. As expected, most of the rare TE insertions (114/146) are detected with a
430 zero or low frequency, while most of the common TE insertions (67/70) appear frequent
431 or fixed (Figure 5A).

432 Overall, we identified only 1.5% (11/709) of clear discrepant estimates between pooled-
433 PCR and T-lex2 frequency estimates. However, after manual curation, we suspected all
434 these TE insertions to be miss-classified by the experimental approach or miss-
435 annotated. For example, insertion *FBti0020042* is classified as present at very low
436 frequency by T-lex2 but is detected at high frequency using the pooled-PCR approach.
437 However, a single strain PCR approach confirms that this TE is present at a low
438 frequency as estimated using T-lex2 (32). The reclassification of this TE insertion,
439 allows to conclude that T-lex2 returns accurate TE frequencies with only ~1% of wrong
440 estimates.

441 To further test the accuracy of T-lex2 calls, we experimentally determined the
442 frequency of 11 well-studied TE insertions known to be polymorphic and thus providing
443 all types of calls (32). We then randomly selected 24 strains part of the bigger set of 86
444 DGRP strains (Supplementary Table S1) and performed single strain PCR for each of
445 the 11 analyzed TEs. We experimentally called the 11 TE insertions in the same 24
446 DGRP strains (Supplementary Table S3). We first compared the frequency estimates for
447 the 11 TE insertions obtained by T-lex2 using all the DGRP strains (86) and the subset
448 of 24 strains. This analysis confirms the robustness of the frequency estimates using the
449 subset of 24 strains (Figure 5B: Pearson's $\rho = 0.96$, $P\text{-val} \ll 0.001$). Out of the 264
450 calls (11 detections in 24 strains), we removed 53 for which we failed to get a call
451 experimentally or using T-lex2 and ended up analyzing 211 TE calls. We observed only
452 11 calls that were different between the experimental data and the T-lex2 calls

453 (Supplementary Table S3). This represents 5.2% of putatively wrong calls that is close
454 to the error rate previously estimated for the TE detection using PCR or the previous
455 version of T-lex, around five percent in both cases (23,31). We finally conservatively
456 estimated a high TE call accuracy with 99.14% of sensitivity and 89.58% of specificity
457 (see Material and Methods). Thus, we overall observed accurate TE frequency estimates
458 comparing the T-lex2 estimates with the pooled-PCR and the individual PCR
459 experimental estimates.

460

461 **Accurate TE population frequency estimates using pooled sequencing data**

462 We called the presence/absence of the 709 selected TE insertions using data
463 from a pool of 92 DGRP strains, containing the 86 single DGRP individual strains
464 previously used (33) (Supplementary Table S1). For 29 out of the 709 selected TE
465 insertions, T-lex2 failed to return a call using the pooled sample, most probably due to
466 low read coverage in these particular regions of the genome. We thus ended up
467 estimating the population frequency of 680 TE insertions (Supplementary Table S2). As
468 expected, we found that the number of reads providing evidence for the absence is
469 negatively correlated with the experimental TE frequency estimates (Figure 6A:
470 Spearman's $\rho = -0.60$, $P\text{-val} \ll 0.001$), while the number of reads providing evidence
471 for the presence is positively correlated (Figure 6B: Spearman's $\rho = 0.87$, $P\text{-val} \ll$
472 0.001). Interestingly, the number of reads supporting the presence does not seem to be
473 strongly biased by the number of copies of TEs from the same family indicating the
474 accuracy of our approach (see Material and Methods).

475 We observed significant and strong positive correlations between T-lex2 TE
476 frequency estimates using the pooled sequencing data and the pooled-PCR determined
477 frequencies (Figure 6C: Spearman's $\rho = 0.87$, $P\text{-val} \ll 0.001$) and between the T-lex2

478 TE frequency estimates using the pooled sequencing data and the T-lex2 estimates using
479 single strains (Figure 6D: Pearson's $\rho = 0.98$, $P\text{-val} \ll 0.001$). This finding supports
480 that the TE population frequency can be estimated with accuracy using pooled NGS
481 sample only based on the read depth coverage.

482 For 25% of the TEs, all of them polymorphic, we noted an overestimation of the TE
483 frequency using the pooled sample compared to the estimates using single strains
484 (Figure 6D). However, for the majority of TEs the difference between the two estimates
485 is $<10\%$, which can be explained by the sampling effect when the pool and the
486 individual strain libraries are constructed. To get accurate frequency estimates, each
487 selected individual should be equally represented in the pooled library. Unfortunately,
488 those conditions are hard to set up or control during library construction. For seven
489 percent of the TEs, the discrepancy between the two estimates was $>10\%$ and the pool
490 estimate was consistently higher than the individual strain estimate. However, we run T-
491 lex2 in a different set of individual and pooled strains and although five percent of the
492 TEs showed discrepant frequencies estimates between pooled and individual samples,
493 these discrepancies were both over- and under-estimations (Supplementary Figure S2).
494 These results suggest that T-lex2 frequency estimation is not biased and that the
495 observed overestimation for the DGRP strains is likely due to some particular feature of
496 this dataset such as sequencing coverage or number of strains used in pooled *vs*
497 individual strains.

498 To assess the effect of the sampling and coverage on the TE frequency estimates using
499 pooled data, we compared the TE frequency estimates using five different pooled
500 datasets (Table 2). The read depth coverage varies from 20X to ~ 60 X (33), the number
501 of strains from 42 to 92, and the number of flies sampled varies from 42 to 134. With
502 only 42 flies pooled and a minimum coverage of 20X, the correlation is strongly positive

503 (Table 2). The increase in the number of flies only does not seem to significantly
504 improve the TE frequency estimates (Table 2). However, the increase of the coverage
505 but not the number of strains improves significantly the TE frequency estimates (Table
506 2). Thus, the coverage matters more than the number of flies or strains pooled in the TE
507 population frequency estimation using pooled sequencing samples. Read-depth coverage
508 of 20X appears sufficient to obtain good estimates although the increase of the coverage
509 allows reducing the number of errors in the calling.

510

511 **T-lex2 detection of TSDs is unbiased and accurate**

512 The new TE-TSD detection module of T-lex2 allows the automatic annotation of TSDs
513 without *a priori* knowledge of the biology of the TEs. While only one copy of the TSD
514 is present in the genome without the TE insertion (*i.e.*, the real TE absence), two TSD
515 copies should be observed in the new constructed reference sequence used by T-lex2 for
516 the absence detection (*i.e.*, the TE absence reference sequence constructed by
517 computational removal of the TE sequence from the genome). Based on this
518 expectation, the new T-lex2 module post-processes the results from the absence module
519 to analyze the alignments of the reads spanning the TE insertion site to identify the
520 TSDs. This new module looks for the expected missing sequence in the alignment
521 flanking the TE insertion site (Figure 3). If the missing sequence corresponds strictly to
522 the TSD, another copy of this TSD should be detected in the vicinity of the missing
523 sequence. Contrary to most of the TSD detection approaches (35-37), no expected TSD
524 length is required for T-lex2 TSD detection.

525 Out of the 587 TEs with at least one read providing evidence for the absence, T-lex2
526 detects TSDs for 390 of them (Table 3). TSD length ranges from two base pairs to 19 bp
527 with an average and a median of five base pairs (Supplementary Figure S3). TSDs of the

528 LTR retroelements are consistently short: four base pairs long on average. TSD lengths
529 of DNA elements are also consistent but longer: seven base pairs on average. On the
530 contrary, TSDs of the non-LTR retroelements (LINE) are not conserved within families
531 and vary in length from four base pairs to 19 bp (Table 3). Overall, we identified for the
532 first time TSDs for 10 TE families: six DNA transposon families and four non-LTR
533 families (Table 3).

534 To assess the quality of the TSD annotation using T-lex2, we compared 148 TSD results
535 obtained by T-lex2 with the annotations obtained with *LTRharvest* software (35)
536 (Supplementary Table S4). For 92 out of the 148 TE insertions, the TSD annotations are
537 identical in length and motif. For seven TEs, TSDs have the same length but differ in
538 motif (Supplementary Table S4) and for the other 49 TEs TSD motifs are identical but
539 differ in length: seven TEs have shorter motifs and 42 TEs have longer motifs. All the
540 longer TSDs correspond to shifts of only one base pair except for one TSD that is three
541 base pairs longer than previously annotated. For 13 of these TEs, corresponding to
542 *Copia* elements, T-lex2 results are supported by experimental studies indicating that
543 *Copia* TSD are five base pairs long (38,39) and not four base pairs as detected with
544 *LTRHarvest* (Supplementary Table S4) (35).

545 TE-TSD detection module did not identify TSDs for 197 TEs out of the 587 with at least
546 one read providing evidence for the absence. TSD detection may fail for old TE
547 insertions for which the TSD is too divergent, for truncated TEs, and when the
548 boundaries of the TE insertions are not well annotated.

549

550 **T-lex2 allows reannotation of TEs**

551 Combining outputs from the TE-analysis module and the TE-TSD annotation, T-
552 lex2 allows manual curation and reannotation of TE insertions. To facilitate the

553 identification of the putatively miss-annotated TEs, information on TEs for which the
554 TSD detection failed and TEs with poly-A tail longer than annotated, are also given in
555 “Tfreq_combine” file. TSD detection fails for TEs that are longer (Figure 4A.II) and
556 shorter (Figure 4A.III) than annotated. We thus manually curated the 197 TEs for which
557 TSD could not be detected. For 65 out of the 197 TE insertions, manual curation
558 confirms the miss-annotation. All these 65 TEs are detected experimentally and by T-
559 lex2 at very low frequency (Supplementary Table S5). 64 of the 65 miss-annotated TE
560 insertions are longer than annotated and they all correspond to non-LTR TE insertions:
561 33 have a longer polyA tail than annotated, and the other 32 TEs have one or the two
562 junctions miss-annotated. The increase in length of these TEs ranges from five base
563 pairs to 40 bp (12 bp on average; Supplementary Table S5). The only TE insertion
564 detected as shorter than annotated is an LTR (*FBti0020118*; Supplementary Table S5).
565 Thus, T-lex2 allows identifying miss-annotated TEs and manual curation allows to
566 reannotate ~10% (65/709) of the analyzed TEs.

567

568

569 **DISCUSSION**

570 We present here the new version of the T-lex pipeline: T-lex2 (Figure 1). This integrated
571 Perl computational pipeline designed to analyze TE insertions in NGS data is freely
572 available and user-friendly since only one command line is necessary to run it
573 (<http://petrov.stanford.edu/cgi-bin/Tlex.html>). T-lex2 is also flexible. The pipeline is
574 composed of five distinct modules that can be launched independently and each one of
575 them includes a large number of options that allows performing flexible and
576 customizable analyses (see T-lex2 manual at [http://petrov.stanford.edu/cgi-](http://petrov.stanford.edu/cgi-bin/Tlex2_manual.html)
577 [bin/Tlex2_manual.html](http://petrov.stanford.edu/cgi-bin/Tlex2_manual.html)). T-lex2 runs with single-end and/or PE reads and can be

578 customized for NGS datasets with different read lengths, read quality and read depth.

579 Thus, T-lex is not only able to work with the different NGS datasets already available,

580 but also with upcoming datasets resulting from improvements in the NGS technology.

581

582 Briefly, T-lex2 has improved and/or expanded the four modules present in T-lex1 and

583 has incorporated a new module that allows detecting TSDs in an unbiased and accurate

584 way (Table 1). Besides individual strain NGS data, T-lex2 can now estimate TE

585 frequencies from pooled NGS data. This is an important new feature of T-lex2 since

586 sequencing pooled samples is an inexpensive and efficient means of obtaining

587 population genome data (22,33). Another important new feature of T-lex2 is the ability

588 to identify TEs likely to return wrong calls *i.e.* miss-annotated TEs, TEs part of

589 segmental duplications and TEs flanked by other TEs. This has been accomplished by

590 the expansion of the TE-filter module of T-lex1 and by the incorporation of the TSD

591 detection module (Figure 1). Being able to identify TEs likely to give wrong calls is

592 especially important for the analyses of TE population frequencies in genomes that do

593 not have a high-quality TE annotation. Note that even in an extremely high quality

594 genome as the *D. melanogaster* one, we re-annotated ~10% of the analyzed TEs. This

595 result also supports the opportunities offered by the NGS data to improve TE

596 annotations (40).

597 Although we tested the reliability of T-lex2 calls and frequency estimates using the fly

598 genome, the pipeline can be used for any species for which NGS data and TE annotation

599 is available. T-lex2 can be run for all type of TEs in any species because (i) it is a

600 flexible pipeline that allows to customize runs for NGS datasets of different qualities,

601 (ii) it is able to identify putatively wrong calls and re-annotate miss-annotate TEs, and

602 (iii) no information about the biology of the TE is needed to get accurate frequency

603 estimates.

604

605 There are other available tools that can be used to estimate TE frequencies (10,18-22).

606 However T-lex2 is to date the only tool that (i) combines presence and absence detection

607 analysis that allows the identification of heterozygotes and/or polymorphic TE

608 insertions; (ii) can be used for individual and pooled samples; (iii) provides the user with

609 an output file with the frequency of the analyzed TEs; (iv) can be easily customize to be

610 run with different NGS datasets in any organism; and (v) assess the quality of the calls

611 allowing to re-annotate miss-annotated TEs. Thus overall, T-lex2 is the most broadly

612 applicable and flexible tool available to date.

613

614 The flexibility of T-lex2, exemplified by the incorporation of the TE-TSD detection

615 module in this new version, will allow us to easily add new modules in a near future. We

616 are currently working on a new module designed to identify TEs not annotated in

617 reference genomes. However, T-lex2 can be used to genotype and estimate frequencies

618 of *de novo* TEs identified with already available software such as RetroSeq (19). The

619 ability of T-lex2 to re-annotate TEs would be specially useful to study this *de novo* TEs

620 as their junctions might not be as well annotated as the ones in the reference genomes.

621

622 Overall, we are providing a versatile tool that allows exploring the impact of TEs on

623 genome evolution as well as learning about TE biology. By analyzing the frequency of

624 TEs in different populations, we will be able to determine which proportion of the TE-

625 induced mutations are subject to strong purifying selection or are likely to be evolving

626 under positive selection. Additionally, accurate annotation of TE insertions will allow us

627 to study for example the target site preferences of different TE families. The availability

628 of a tool such as T-lex2 improves our ability to explore TE dynamics and TE biology in
629 an increasing number of species. This is important, as we cannot hope to understand
630 genome structure and function without a thorough understanding of its most active,
631 diverse, and ancient component.

632

633

634 **ACKNOWLEDGMENTS**

635 We would like to thank the ScaleGenomics (<http://scalegen.com/>) and Proclus
636 (Stanford University; <http://www.stanford.edu/>) computational platforms that we used to
637 investigate this study.

638

639 **FUNDING**

640 This work was supported by grants from the National Institutes of Health
641 [R01GM089926 to D.A.P.], from the European Commission [PCIG-GA-2011-293860
642 to J.G.], and from the Spanish Government [BFU-2011-24397 to J. G.]. MGB was
643 supported by a PIF fellowship from the Universitat Autònoma de Barcelona.

644

645

646

647 **Table 1.** Comparison of the original T-lex version and T-lex2 modules and features
 648
 649

T-lex2 Modules	Features	T-lex (23)	T-lex2
TE-analysis	Detection of low-complexity regions	✓	✓
	Detection of repeat-rich regions	✓	✓
	Detection of terminal poly-A stretches	✗	✓
	Detection of TEs part of segmental duplications	✗	✓
TE-presence detection	Analyzes inside the TE	✓ (20bp)	✓ (60bp)
	Read quality filtered	✗	✓ (phrep score ≥ 30)
	Read depth reported	✗	✓
	Min. match length	✓ (5bp)	✓ (15bp)
	Min. sequence identity	✗	✓ (95%)
TE-absence detection	Analyzes on each side of the TE	✓ (100bp)	✓ (125bp)
	PE data analysis	✗	✓
	Read quality filtered	✗	✓ (30)
	Read depth reported	✗	✓
TE-combine	Take into account the no calls	✗	✓
	Estimates using pooled sample	✗	✓
TE-TSD	TSD detection	✗	✓
	Flag of putative miss-annotated TEs	✗	✓

650
 651
 652

653 **Table 2.** Sampling effect on the population frequency estimates of 680 TE insertions.

654

655

Sample identifier	Read depth coverage	Nb. of strains pooled	Nb. of flies pooled	% (Nb.) of no calls	Pearson's ρ [IC95]
B1	20X	42	42	3.82(26)	0.968 [0.963-0.973]
B4	20X	50	50	4.70(32)	0.977 [0.973-0.980]
B3	40X	42	84	0.88(6)	0.980 [0.977-0.983]
B5	40X	92	92	0.88(6)	0.981 [0.978-0.983]
B3-B4	60X	92	134	0(0)	0.985 [0.982-0.987]

656

657

658 **Table 3.** TSD detection for the 53 TE families analyzed in this study. ND: Not
 659 determined.
 660
 661

TE order	TE superfamily	TE family	T-lex2				Known TSD length (bp)
			Nb. of TE insertions with TSD detected	TSD consensus motif ⁽¹⁾	Min. length (bp)	Max. length (bp)	
DNA	hAT	H	11	ACCCACAG	8	8	8 ⁽²⁾
DNA	Mariner/Tc1	FB	7	AAAACCGTC	9	9	ND
DNA	Mariner/Tc1	S2	1	AT	2	2	ND
DNA	P-element	1360	36	GTYYYYG	6	7	7 ⁽²⁾
DNA	Transib	hopper	7	CAATG	5	5	5 ⁽²⁾
DNA	Transib	transib2	5	CAVTG	5	5	5 ⁽²⁾
LINE	Jockey	BS	12	×	4	12	ND
LINE	Jockey	Doc	14	×	4	12	ND
LINE	Jockey	F-element	6	×	9	13	~10 ⁽³⁾
LINE	Jockey	G2	4	×	5	12	~10 ⁽³⁾
LINE	I-element	I-element	2	×	15	19	ND
LINE	Jockey	Ivk	1	CACATCATCTTAT	13	13	~10 ⁽⁴⁾
LINE	Jockey	jockey	7	×	7	15	ND
LINE	Jockey	Juan	1	ACGAAACAACATT A	14	14	ND
LINE	R1	Rt1a	2	×	6	14	ND
LINE	R1	Rt1b	1	GCTATC	6	6	ND
LINE	Jockey	X-element	1	TTTTGAAA	8	8	ND
LTR	copia	1731	1	TAAAT	5	5	4-6 ⁽⁵⁾
LTR	copia	copia	21	AAAAT	5	5	4-6 ⁽⁵⁾
LTR	Gypsy	17.6	3	ATAT	4	4	4-6 ⁽⁵⁾
LTR	Gypsy	297	20	ATAT	4	4	4 ⁽²⁾
LTR	Gypsy	412	16	AAAC	4	4	4 ⁽²⁾
LTR	Gypsy	blood	19	AACC	4	4	4-6 ⁽⁵⁾
LTR	Gypsy	Burdock	9	TRYA	4	4	4 ⁽²⁾
LTR	Gypsy	diver	6	AAGGG	5	5	4-6 ⁽⁵⁾
LTR	Gypsy	flea	9	TXHA	4	4	4-6 ⁽⁵⁾
LTR	Gypsy	gtwin	1	TACA	4	4	4 ⁽²⁾
LTR	Gypsy	gypsy	1	TACA	4	4	4 ⁽²⁾
LTR	Gypsy	gypsy2	1	ATA	3	3	4 ⁽²⁾
LTR	Gypsy	gypsy5	2	CGCG	4	4	4 ⁽²⁾
LTR	Gypsy	HMS-Beagle	3	TGYA	4	4	4 ⁽²⁾
LTR	Gypsy	HMS-Beagle2	2	×	4	7	4-6 ⁽⁵⁾
LTR	Gypsy	invader1	2	ATG	3	4	4-6 ⁽⁵⁾
LTR	Gypsy	invader3	1	AAAT	4	4	4 ⁽²⁾
LTR	Gypsy	invader4	1	TATA	4	4	4-6 ⁽⁵⁾

LTR	Gypsy	invader6	1	TATA	4	4	4-6 ⁽⁵⁾
LTR	Gypsy	McClintock	2	YATAT	5	5	4-6 ⁽⁵⁾
LTR	Gypsy	mdg1	10	AAGG	4	4	4 ⁽²⁾
LTR	Gypsy	mdg3	9	AGTA	4	4	4 ⁽²⁾
LTR	Gypsy	opus	3	TAA	3	4	4 ⁽²⁾
LTR	Gypsy	Quasimodo	5	ATAT	4	4	4 ⁽²⁾
LTR	Gypsy	rover	1	ATAT	4	4	4-6 ⁽⁵⁾
LTR	Gypsy	springer	1	GTA	3	3	4-6 ⁽⁵⁾
LTR	Gypsy	Stalker	4	CAAT	4	4	4-6 ⁽⁵⁾
LTR	Gypsy	Stalker2	4	AAAT	4	4	4 ⁽²⁾
LTR	Gypsy	Stalker4	2	ATAC	4	4	4-6 ⁽⁵⁾
LTR	Gypsy	Tabor	2	CCTG	4	4	4 ⁽²⁾
LTR	Gypsy	Tirant	15	ACG	2	4	2 ⁽²⁾
LTR	Gypsy	Transpac	4	ATAT	4	4	4 ⁽²⁾
LTR	Pao	3S18	4	GYRDT	5	5	4-6 ⁽⁵⁾
LTR	Pao	Max	2	CCTCA	5	5	5 ⁽²⁾
LTR	Pao	roo	84	AACAG	4	7	5 ⁽²⁾
LTR	Pao	rooA	1	ACAAT	5	5	4-6 ⁽⁵⁾

662 ⁽¹⁾TSD motifs represented in the standard IUB/IUPAC nucleic acid codes. The sign “*” indicates that no
663 consensus can be reported due to a lack of motifs or too variable motifs.

664 ⁽²⁾(41), ⁽³⁾(42), ⁽⁴⁾(43), ⁽⁵⁾(5).

665

666 **FIGURE LEGENDS**

667

668 **Figure 1. T-lex2 pipeline.**

669 Schematic representation of the five different T-lex modules (dark gray boxes) and of

670 the input files (white boxes) required to run T-lex2. SD: Segmental Duplication

671

672 **Figure 2. T-lex2 TE-presence detection and TE-absence detection modules.**

673 (A) TE-presence detection module is based on the mapping of the NGS reads on the TE

674 insertion junctions. (B) TE-absence detection module is based on the mapping of NGS

675 reads on the putative ancestral genomic sequence prior the TE insertion. Input files

676 required to run each module (white boxes), the different steps of the pipeline (light grey

677 boxes), and examples of reads providing evidence for the presence, absence, and no call

678 reads are depicted.

679

680 **Figure 3. T-lex2 TE-TSD detection module.**

681 Schematic representation of the procedure to identify TSDs. Input files required to run

682 this module and the different steps followed by the pipeline are also depicted.

683

684 **Figure 4. Detection of putatively mis-annotated TE insertions.**

685 (A) Schematic representation of a correctly annotated TE (I), a TE longer than annotated

686 (II) and a TE shorter than annotated (III). (B) Effect of miss-annotation on the TE calls

687 when the TE insertion is present and absent. When the TE insertion is absent, T-lex2

688 will give a correct call if the miss-annotation is short and an erroneous call in the miss-

689 annotation is long (B.h, B.i, B.k, and B.l). (1) If the miss-annotated region is longer than

690 the maximum read length, selected reads will not overlap the real TE junction(s). These

691 reads will be mapped inside the TE sequence. Consequently, the number of selected

692 reads will depend on the number of TE copies from the same family. (2) If the miss-
693 annotated region is longer than the maximum read length, selected reads will fully map
694 on the TE flanking region. The number of selected reads will be a proxy of the local
695 read depth coverage.

696

697 **Figure 5. TE frequency estimates using single strains.**

698 (A) Comparison of the frequency estimates of 709 TE insertions using T-lex2 and
699 pooled PCR approach. (B) Frequency estimates of 11 TE insertions using T-lex2 and
700 single PCR approach.

701

702 **Figure 6. TE frequency estimates using pooled samples.**

703 (A) Number of reads supporting the absence and (B) number of reads supporting the
704 presence for each of the TE frequency classes. (C) T-lex2 frequency estimates using
705 pooled sample *versus* TE frequency classes experimentally determined using a pooled-
706 PCR approach. (D) Comparison of T-lex2 frequency estimates using single and pooled
707 NGS data.

708

709 **REFERENCES**

710

- 711 1. Gilad, Y., Pritchard, J.K. and Thornton, K. (2009) Characterizing natural
712 variation using next-generation sequencing technologies. *Trends in genetics* :
713 *TIG*, **25**, 463-471.
- 714 2. Koboldt, D.C., Steinberg, K.M., Larson, D.E., Wilson, R.K. and Mardis, E.R.
715 (2013) The next-generation sequencing revolution and its impact on genomics.
716 *Cell*, **155**, 27-38.
- 717 3. Shapiro, E., Biezuner, T. and Linnarsson, S. (2013) Single-cell sequencing-based
718 technologies will revolutionize whole-organism science. *Nature reviews*.
719 *Genetics*, **14**, 618-630.
- 720 4. Treangen, T.J. and Salzberg, S.L. (2012) Repetitive DNA and next-generation
721 sequencing: computational challenges and solutions. *Nature reviews. Genetics*,
722 **13**, 36-46.
- 723 5. Wicker, T., Sabot, F., Hua-Van, A., Bennetzen, J.L., Capy, P., Chalhoub, B.,
724 Flavell, A., Leroy, P., Morgante, M., Panaud, O. *et al.* (2007) A unified
725 classification system for eukaryotic transposable elements. *Nature reviews*.
726 *Genetics*, **8**, 973-982.
- 727 6. Kapitonov, V.V. and Jurka, J. (2008) A universal classification of eukaryotic
728 transposable elements implemented in Repbase. *Nature reviews. Genetics*, **9**,
729 411-412; author reply 414.
- 730 7. Biemont, C. (2010) A brief history of the status of transposable elements: from
731 junk DNA to major players in evolution. *Genetics*, **186**, 1085-1093.
- 732 8. de Koning, A.P., Gu, W., Castoe, T.A., Batzer, M.A. and Pollock, D.D. (2011)
733 Repetitive elements may comprise over two-thirds of the human genome. *PLoS*
734 *genetics*, **7**, e1002384.
- 735 9. Devos, K.M., Beales, J., Ogihara, Y. and Doust, A.N. (2005) Comparative
736 sequence analysis of the phytochrome C gene and its upstream region in
737 allohexaploid wheat reveals new data on the evolution of its three constituent
738 genomes. *Plant molecular biology*, **58**, 625-641.
- 739 10. Ewing, A.D. and Kazazian, H.H., Jr. (2011) Whole-genome resequencing allows
740 detection of many rare LINE-1 insertion alleles in humans. *Genome research*,
741 **21**, 985-990.
- 742 11. Iskow, R.C., McCabe, M.T., Mills, R.E., Torene, S., Pittard, W.S., Neuwald,
743 A.F., Van Meir, E.G., Vertino, P.M. and Devine, S.E. (2010) Natural
744 mutagenesis of human genomes by endogenous retrotransposons. *Cell*, **141**,
745 1253-1261.
- 746 12. Cowley, M. and Oakey, R.J. (2013) Transposable elements re-wire and fine-tune
747 the transcriptome. *PLoS genetics*, **9**, e1003234.
- 748 13. Casacuberta, E. and Gonzalez, J. (2013) The impact of transposable elements in
749 environmental adaptation. *Molecular ecology*, **22**, 1503-1517.
- 750 14. Akagi, K., Li, J. and Symer, D.E. (2013) How do mammalian transposons induce
751 genetic variation? A conceptual framework: the age, structure, allele frequency,
752 and genome context of transposable elements may define their wide-ranging
753 biological impacts. *BioEssays : news and reviews in molecular, cellular and*
754 *developmental biology*, **35**, 397-407.
- 755 15. Chenais, B., Caruso, A., Hiard, S. and Casse, N. (2012) The impact of
756 transposable elements on eukaryotic genomes: from genome size increase to

- 757 genetic adaptation to stressful environments. *Gene*, **509**, 7-15.
- 758 16. Hua-Van, A., Le Rouzic, A., Boutin, T.S., Filee, J. and Capy, P. (2011) The
759 struggle for life of the genome's selfish architects. *Biology direct*, **6**, 19.
- 760 17. Makalowski, W., Pande, A., Gotea, V. and Makalowska, I. (2012) Transposable
761 elements and their identification. *Methods in molecular biology*, **855**, 337-359.
- 762 18. Hormozdiari, F., Hajirasouliha, I., Dao, P., Hach, F., Yorukoglu, D., Alkan, C.,
763 Eichler, E.E. and Sahinalp, S.C. (2010) Next-generation VariationHunter:
764 combinatorial algorithms for transposon insertion discovery. *Bioinformatics*, **26**,
765 i350-357.
- 766 19. Keane, T.M., Wong, K. and Adams, D.J. (2013) RetroSeq: transposable element
767 discovery from next-generation sequencing data. *Bioinformatics*, **29**, 389-390.
- 768 20. Lee, E., Iskow, R., Yang, L., Gokcumen, O., Haseley, P., Luquette, L.J., 3rd,
769 Lohr, J.G., Harris, C.C., Ding, L., Wilson, R.K. *et al.* (2012) Landscape of
770 somatic retrotransposition in human cancers. *Science*, **337**, 967-971.
- 771 21. Robb, S.M., Lu, L., Valencia, E., Burnette, J.M., 3rd, Okumoto, Y., Wessler,
772 S.R. and Stajich, J.E. (2013) The Use of RelocaTE and Unassembled Short
773 Reads to Produce High-Resolution Snapshots of Transposable Element
774 Generated Diversity in Rice. *G3 (Bethesda)*, **3**, 949-957.
- 775 22. Kofler, R., Betancourt, A.J. and Schlotterer, C. (2012) Sequencing of pooled
776 DNA samples (Pool-Seq) uncovers complex dynamics of transposable element
777 insertions in *Drosophila melanogaster*. *PLoS genetics*, **8**, e1002487.
- 778 23. Fiston-Lavier, A.S., Carrigan, M., Petrov, D.A. and Gonzalez, J. (2011) T-lex: a
779 program for fast and accurate assessment of transposable element presence using
780 next-generation sequencing data. *Nucleic acids research*, **39**, e36.
- 781 24. Lerat, E. (2010) Identifying repeats and transposable elements in sequenced
782 genomes: how to find your way through the dense forest of programs. *Heredity*,
783 **104**, 520-533.
- 784 25. Smit, A., Hubley, R & Green, P. . (1996-2010) RepeatMasker Open-3.0.
785 **<Erreur ! Référence de lien hypertexte non valide.. .**
- 786 26. Kent, W.J. (2002) BLAT--the BLAST-like alignment tool. *Genome research*,
787 **12**, 656-664.
- 788 27. Green, P. (2009) 2009. Phrap, version 1.090518. <http://phrap.org>. .
- 789 28. David, M., Dzamba, M., Lister, D., Ilie, L. and Brudno, M. (2011) SHRiMP2:
790 sensitive yet practical SHort Read Mapping. *Bioinformatics*, **27**, 1011-1012.
- 791 29. Rumble, S.M., Lacroute, P., Dalca, A.V., Fiume, M., Sidow, A. and Brudno, M.
792 (2009) SHRiMP: accurate mapping of short color-space reads. *PLoS Comput*
793 *Biol*, **5**, e1000386.
- 794 30. Mackay, T.F., Richards, S., Stone, E.A., Barbadilla, A., Ayroles, J.F., Zhu, D.,
795 Casillas, S., Han, Y., Magwire, M.M., Cridland, J.M. *et al.* (2012) The
796 *Drosophila melanogaster* Genetic Reference Panel. *Nature*, **482**, 173-178.
- 797 31. Petrov, D.A., Fiston-Lavier, A.S., Lipatov, M., Lenkov, K. and Gonzalez, J.
798 (2011) Population genomics of transposable elements in *Drosophila*
799 *melanogaster*. *Mol Biol Evol*, **28**, 1633-1644.
- 800 32. Gonzalez, J., Lenkov, K., Lipatov, M., Macpherson, J.M. and Petrov, D.A.
801 (2008) High rate of recent transposable element-induced adaptation in
802 *Drosophila melanogaster*. *PLoS Biol*, **6**, e251.
- 803 33. Zhu, Y., Bergland, A.O., Gonzalez, J. and Petrov, D.A. (2012) Empirical
804 validation of pooled whole genome population re-sequencing in *Drosophila*
805 *melanogaster*. *PloS one*, **7**, e41901.
- 806 34. Gonzalez, J., Karasov, T.L., Messer, P.W. and Petrov, D.A. (2010) Genome-

- 807 wide patterns of adaptation to temperate environments associated with
808 transposable elements in *Drosophila*. *PLoS genetics*, **6**, e1000905.
- 809 35. Ellinghaus, D., Kurtz, S. and Willhoeft, U. (2008) LTRharvest, an efficient and
810 flexible software for de novo detection of LTR retrotransposons. *BMC*
811 *Bioinformatics*, **9**, 18.
- 812 36. Szak, S.T., Pickeral, O.K., Makalowski, W., Boguski, M.S., Landsman, D. and
813 Boeke, J.D. (2002) Molecular archeology of L1 insertions in the human genome.
814 *Genome Biol*, **3**, research0052.
- 815 37. Rawal, K. and Ramaswamy, R. (2011) Genome-wide analysis of mobile genetic
816 element insertion sites. *Nucleic acids research*, **39**, 6864-6878.
- 817 38. Dunsmuir, P., Brorein, W.J., Jr., Simon, M.A. and Rubin, G.M. (1980) Insertion
818 of the *Drosophila* transposable element copia generates a 5 base pair duplication.
819 *Cell*, **21**, 575-579.
- 820 39. Rubin, G.M., Brorein, W.J., Jr., Dunsmuir, P., Flavell, A.J., Levis, R., Strobel,
821 E., Toole, J.J. and Young, E. (1981) Copia-like transposable elements in the
822 *Drosophila* genome. *Cold Spring Harb Symp Quant Biol*, **45 Pt 2**, 619-628.
- 823 40. Chaparro, C. and Sabot, F. (2012) Methods and software in NGS for TE
824 analysis. *Methods in molecular biology*, **859**, 105-114.
- 825 41. Linheiro, R.S. and Bergman, C.M. (2012) Whole genome resequencing reveals
826 natural target site preferences of transposable elements in *Drosophila*
827 *melanogaster*. *PloS one*, **7**, e30008.
- 828 42. Kapitonov, V.V. and Jurka, J. (2003) Molecular paleontology of transposable
829 elements in the *Drosophila melanogaster* genome. *Proc Natl Acad Sci U S A*,
830 **100**, 6569-6574.
- 831 43. Berezikov, E., Bucheton, A. and Busseau, I. (2000) A search for reverse
832 transcriptase-coding sequences reveals new non-LTR retrotransposons in the
833 genome of *Drosophila melanogaster*. *Genome Biol*, **1**, RESEARCH0012.
- 834

Figure 1.

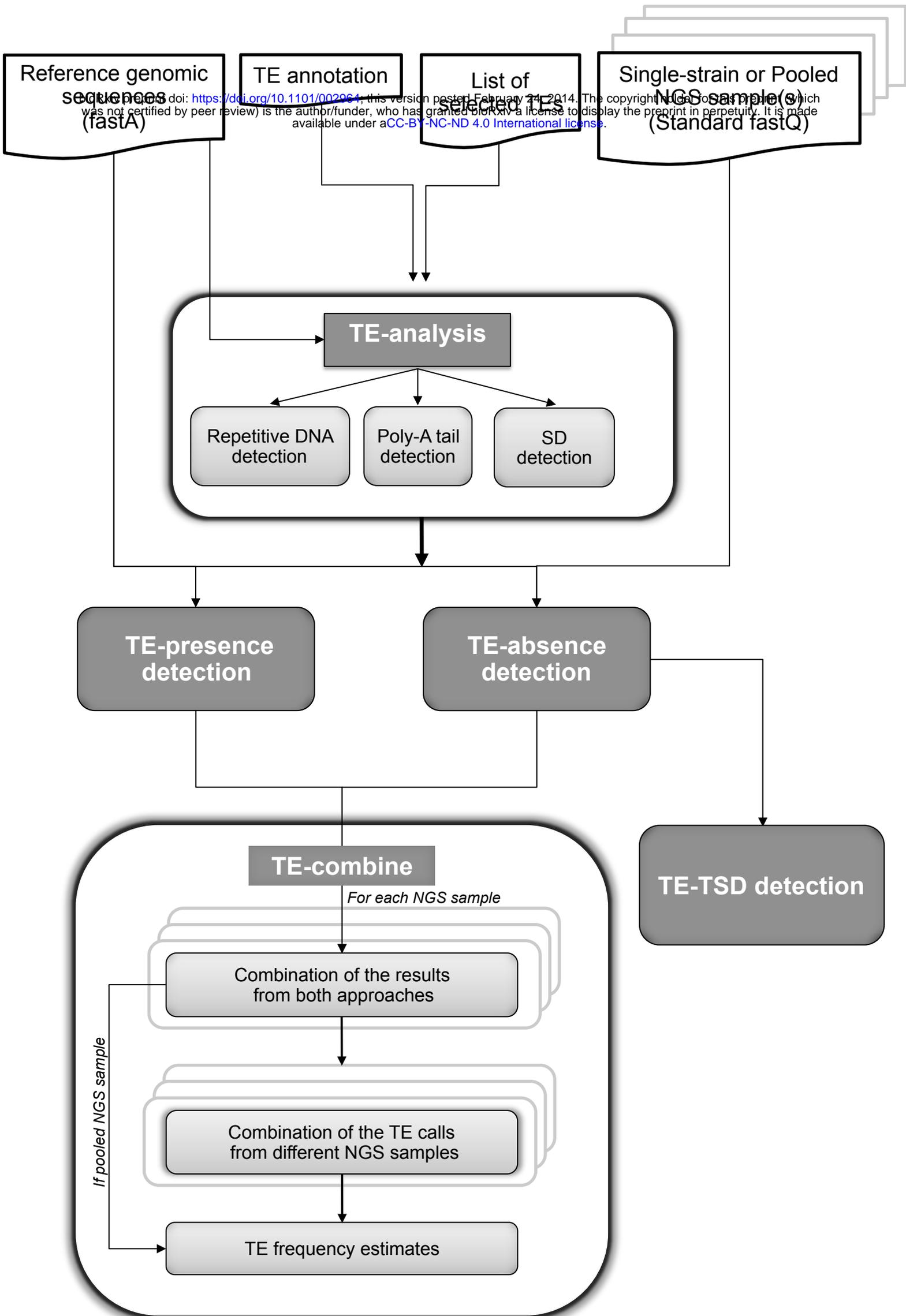
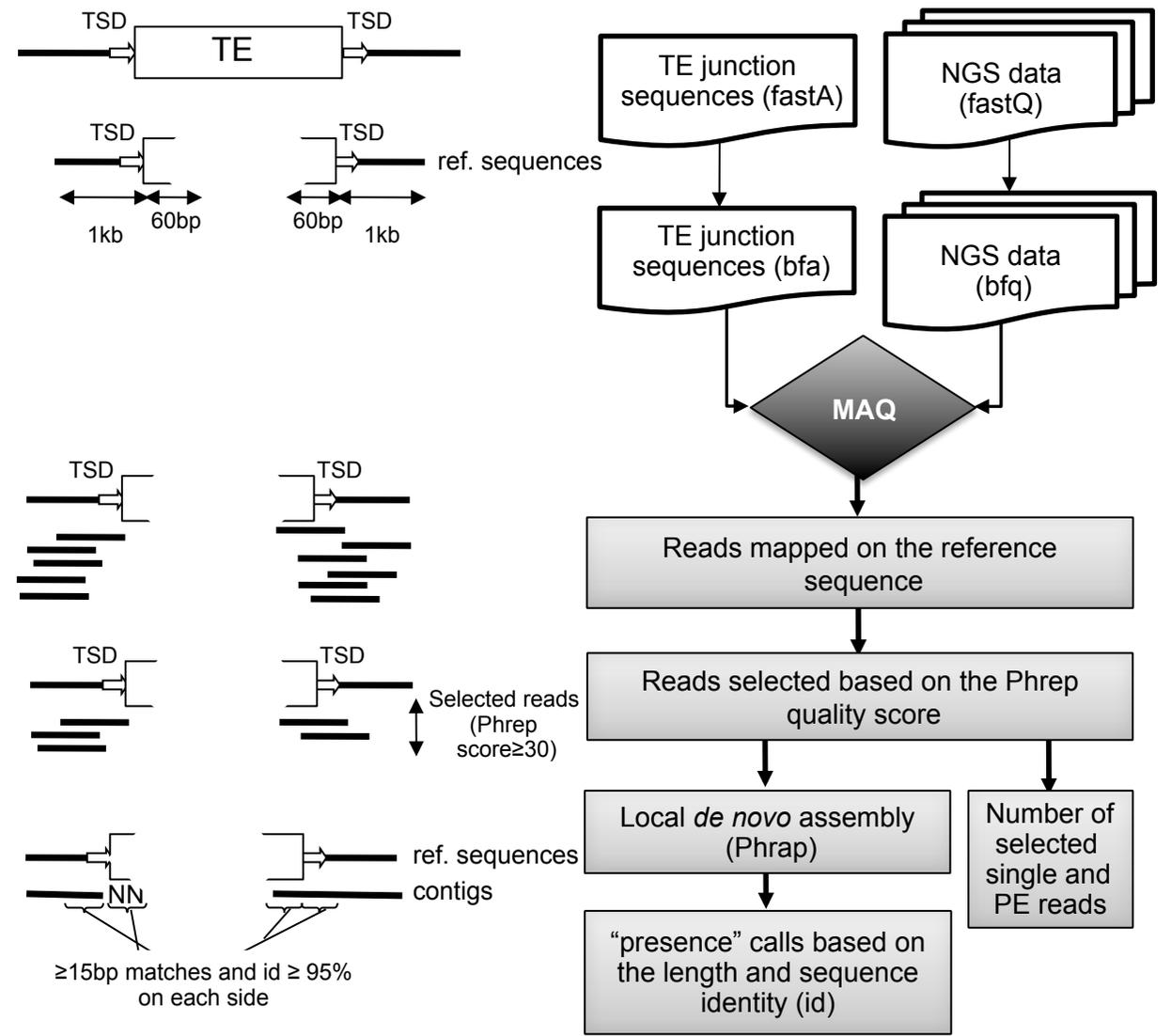
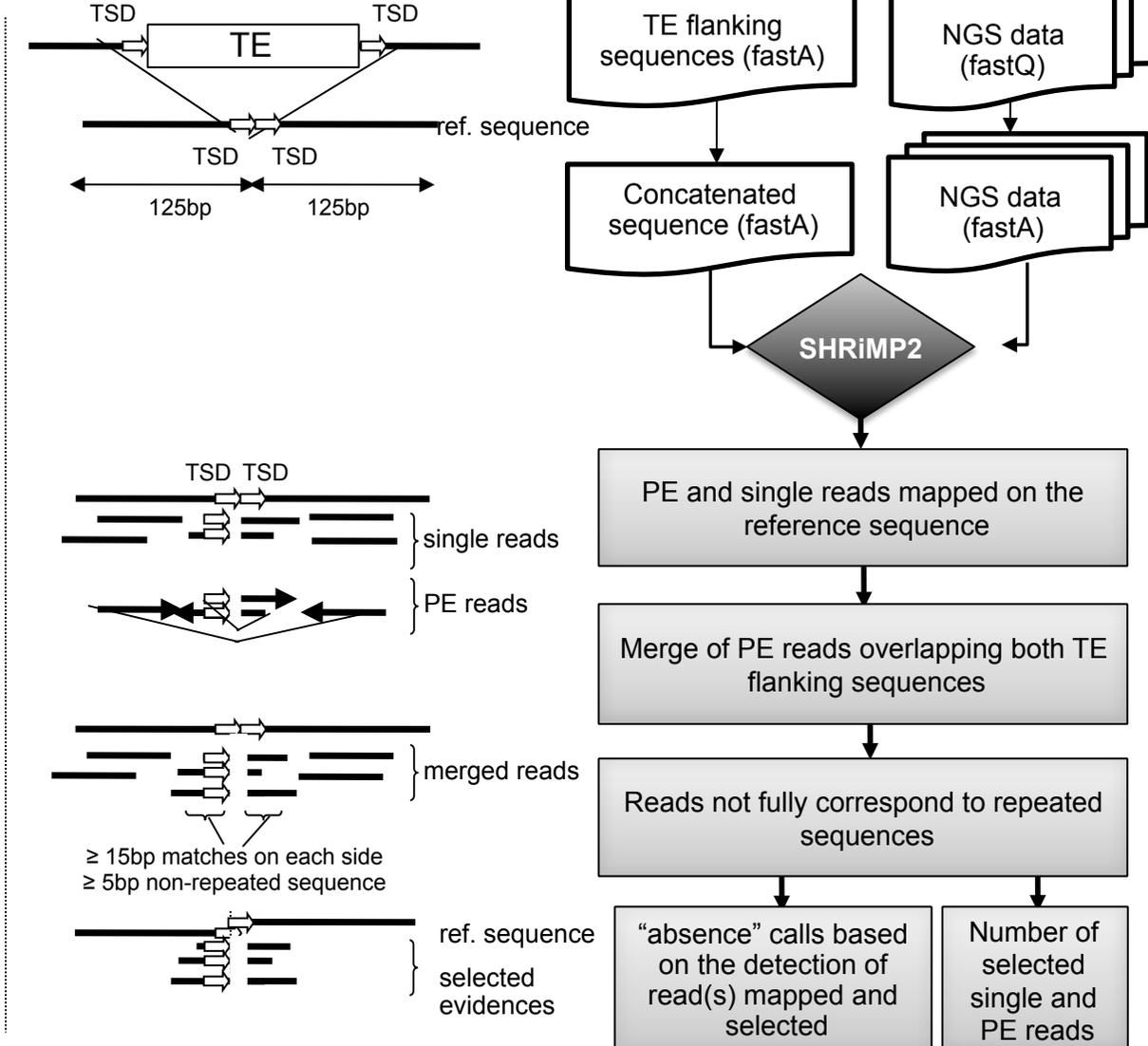


Figure 2.

A. TE-presence detection module



B. TE-absence detection module



TE present	ref. sequence	TTAACACTTCATACAGTTGCGGTAACAATAAA	TTATTTTAACCGCACCTGTATGTTAACTAACT
	contig	TTAACACTTCATACAGTTGCGGTAACAANNNN	TTATTTTAACCGCACCTGTATGTTAACTAACT
TE absent	ref. sequence	TATAGTCTAGCCAAAGACACTAGAATAACAAG	TTGTTATTCTAGTGCTTTGGTCTAGCGGACAT
	contig	TATAGTCTANNNNNNNNNNNNNNNNNNNNNNNNN	NNNNNNNNNNNNNNNNNNNNATAGTCTAGCGGACAT
no call	ref. sequence	TGTGTATATATGTAGTGTATCTACCCTCACT	ATTACTTATTACTTATTACAATATATATATAA
	contig	NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN	NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNTATAA

←----- TE sequence -----→

TE present	ref. sequence	CCTAAGTTGTAAACGTTGGTACAGCCCGGTACAGCC	TAGTATAGTATGGTAAATGCAA
	read1	CCTAAGTTGTAAACGTTGGTACAGCC	
	read2		GGTACAGCC TAGTATAGTATGGTAAATGCAA
TE absent	ref. sequence	TTGAATATTTTAATCTCTGAAAACCTGAAA	AATAGCTTTTAGATATCCTCTTATCAAC
	read	TTGAATATTTTAATC-----CTGAAAA	AATAGCTTTTAGATATCCTCTTATCAAC
no call	ref. sequence	GACTTGCTAAGGGTGTAAAGCTCAGTAAATAACCGTAAATA	TAAACAGGTTACCCTGA
	No read		

Figure 3.

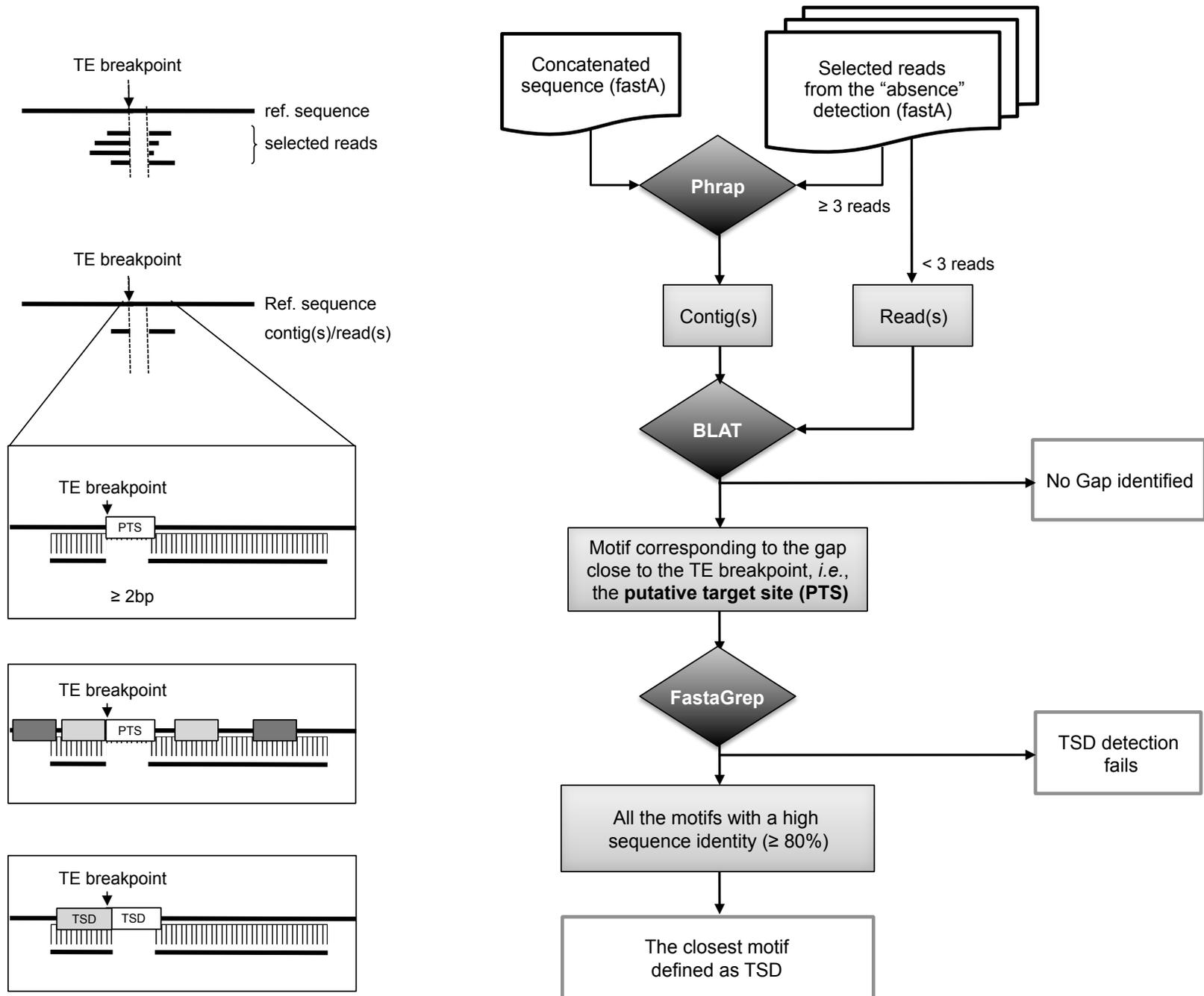
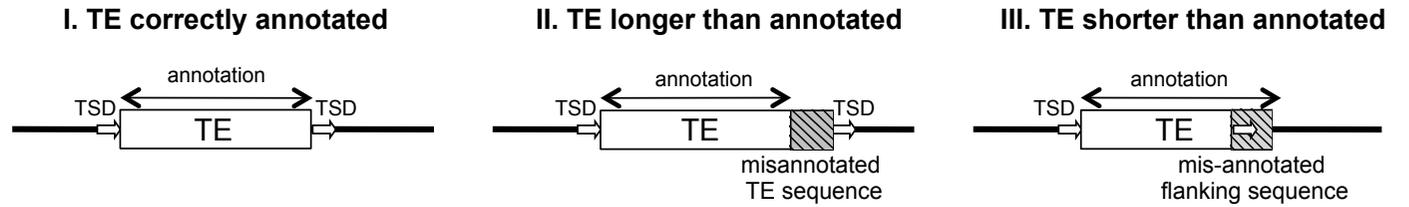


Figure 4.

A. TE annotation and misannotation



B. Effect of misannotation on TE calls

		TE insertion is PRESENT	TE insertion is ABSENT
TE-presence detection	I. TE correctly annotated	a) ref. sequences present	g) ref. sequences absent
	II. TE longer than annotated	b) present ⁽¹⁾	h) absent or no call
	III. TE shorter than annotated	c) present ⁽²⁾	i) absent or present
TE-absence detection	I. TE correctly annotated	d) ref. sequence present	j) ref. sequence absent
	II. TE longer than annotated	e) present	k) absent or present
	III. TE shorter than annotated	f) present	l) absent or present

(1) If the mis-annotated region is longer than the maximum read length, selected reads will not overlap the real TE junction(s). These reads will be mapped inside the TE sequence. Consequently, the number of selected reads will depend on the number of TE copies from the same family.

(2) If the mis-annotated region is longer than the maximum read length, selected reads will fully map on the TE flanking region. The number of selected reads will be a proxy of the local read depth coverage.

Figure 5.

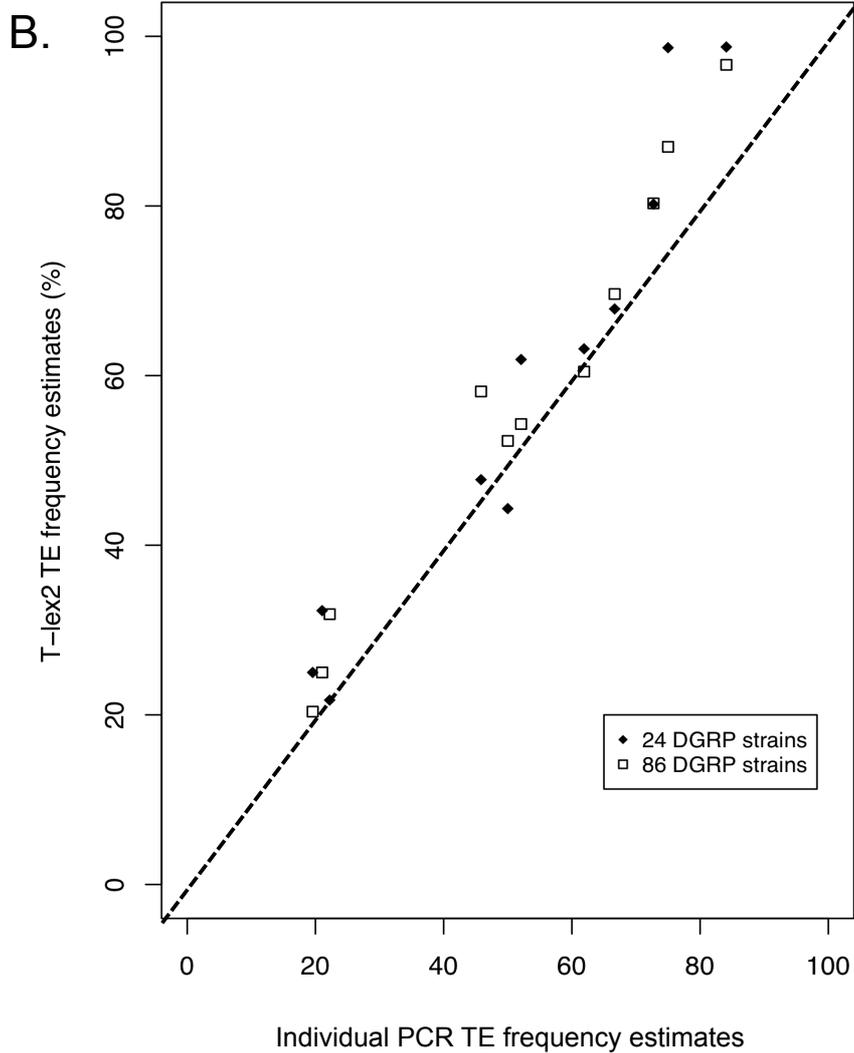
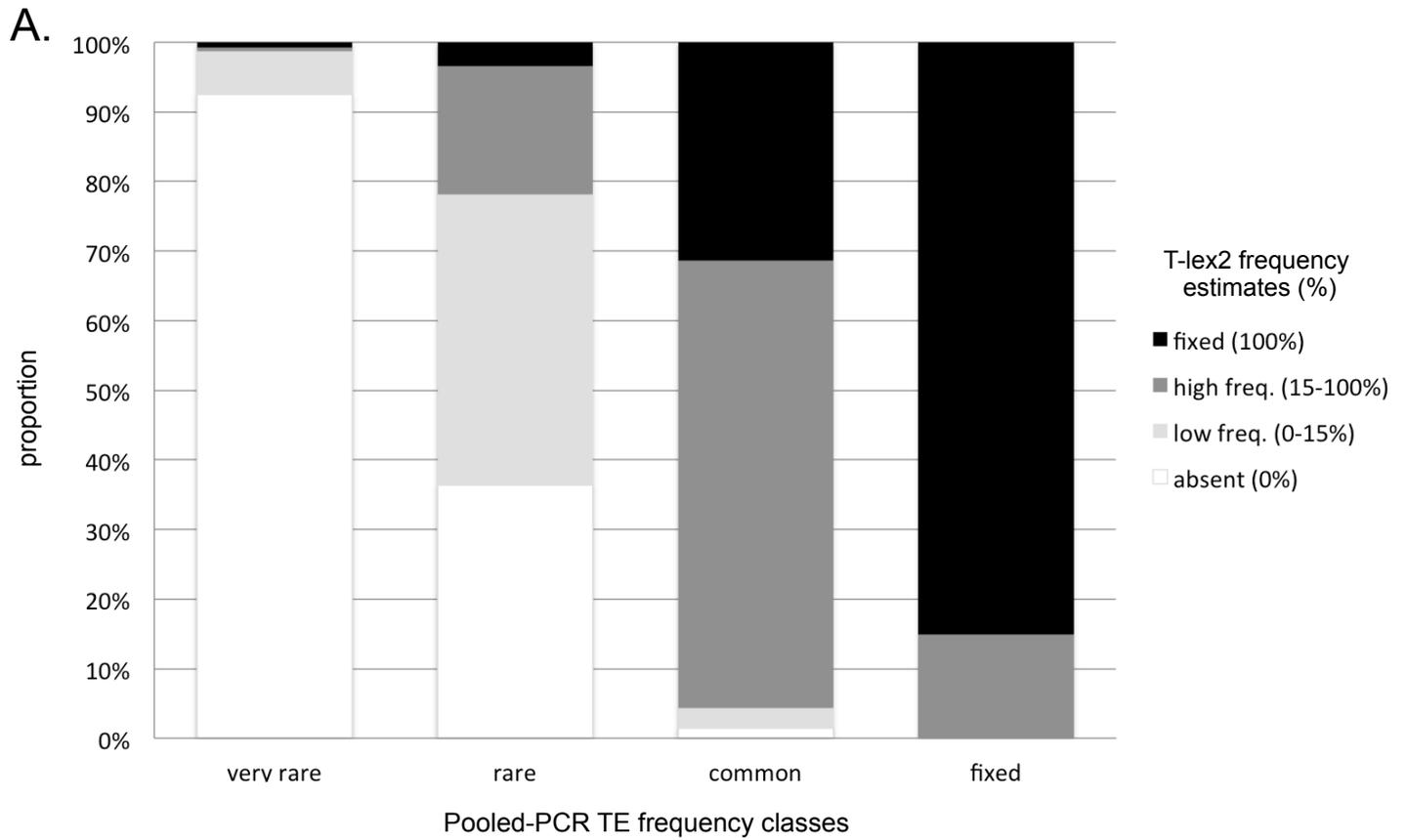


Figure 6.

