

# RelocaTE2: a high resolution transposable element polymorphism mapping tool for population resequencing

Jinfeng Chen<sup>1,2,3</sup>, Travis R Wrightsman<sup>3</sup>, Susan R Wessler<sup>2,3</sup>, Jason E. Stajich<sup>Corresp. 1,2</sup>

<sup>1</sup> Department of Plant Pathology & Microbiology, University of California, Riverside, Riverside, CA, United State

<sup>2</sup> Institute for Integrative Genome Biology, University of California, Riverside, Riverside, CA, United States

<sup>3</sup> Department of Botany and Plant Sciences, University of California, Riverside, Riverside, CA, United States

Corresponding Author: Jason E. Stajich

Email address: jason.stajich@ucr.edu

**Background** Transposable element (TE) polymorphisms are important components of population genetic variation. The functional impacts of TEs in gene regulation and generating genetic diversity have been observed in multiple species, but the frequency and magnitude of TE variation is under appreciated. Inexpensive and deep sequencing technology has made it affordable to apply population genetic methods to whole genomes with methods that identify single nucleotide and insertion/deletion polymorphisms. However, identifying TE transposition events or polymorphisms can be challenging due to the repetitive nature of these sequences, which hamper both the sensitivity and specificity of analysis tools.

**Methods** We have developed the tool RelocaTE2 (<http://github.com/stajichlab/RelocaTE2>) for identification of TE polymorphisms at high sensitivity and specificity. RelocaTE2 searches for known TE sequences in whole genome sequencing reads from second generation sequencing platforms such as Illumina. These sequence reads are used as seeds to pinpoint chromosome locations where TEs have transposed. RelocaTE2 detects target site duplication (TSD) of TE insertions allowing it to report TE polymorphism loci with single base pair precision.

**Results and Discussion** The performance of RelocaTE2 is evaluated using both simulated and real sequence data. RelocaTE2 demonstrates a higher level of sensitivity and specificity when compared to other tools. Even in highly repetitive regions, such as those tested on rice chromosome 4, RelocaTE2 is able to report up to 95% of simulated TE insertions with less than 0.1% false positive rate using 10-fold genome coverage resequencing data. RelocaTE2 provides a robust solution to identify TE polymorphisms and can be incorporated into analysis workflows in support of describing the complete genotype from light coverage genome sequencing.

1 **RelocaTE2: a high resolution transposable element polymorphism mapping tool for**  
2 **population resequencing**

3 Jinfeng Chen<sup>1,2,3</sup>, Travis R. Wrightsman<sup>2</sup>, Susan R. Wessler<sup>2,3</sup> and Jason E. Stajich<sup>1,3,\*</sup>

4 <sup>1</sup> Department of Plant Pathology and Microbiology, University of California-Riverside,  
5 Riverside CA 92521, USA.

6 <sup>2</sup> Department of Botany and Plant Sciences, University of California-Riverside, Riverside CA  
7 92521, USA.

8 <sup>3</sup> Institute for Integrative Genome Biology, University of California-Riverside, Riverside, CA  
9 92521, USA.

10 \*To whom correspondence should be addressed. E-mail: [jason.stajich@ucr.edu](mailto:jason.stajich@ucr.edu)

**11 Abstract****12 Background**

13 Transposable element (TE) polymorphisms are important components of population genetic  
14 variation. The functional impacts of TEs in gene regulation and generating genetic diversity have  
15 been observed in multiple species, but the frequency and magnitude of TE variation is under  
16 appreciated. Inexpensive and deep sequencing technology has made it affordable to apply  
17 population genetic methods to whole genomes with methods that identify single nucleotide and  
18 insertion/deletion polymorphisms. However, identifying TE transposition events or  
19 polymorphisms can be challenging due to the repetitive nature of these sequences, which hamper  
20 both the sensitivity and specificity of analysis tools.

**21 Methods**

22 We have developed the tool RelocaTE2 (<http://github.com/stajichlab/RelocaTE2>) for  
23 identification of TE polymorphisms at high sensitivity and specificity. RelocaTE2 searches for  
24 known TE sequences in whole genome sequencing reads from second generation sequencing  
25 platforms such as Illumina. These sequence reads are used as seeds to pinpoint chromosome  
26 locations where TEs have transposed. RelocaTE2 detects target site duplication (TSD) of TE  
27 insertions allowing it to report TE polymorphism loci with single base pair precision.

**28 Results and Discussion**

29 The performance of RelocaTE2 is evaluated using both simulated and real sequence data.  
30 RelocaTE2 demonstrates a higher level of sensitivity and specificity when compared to other  
31 tools. Even in highly repetitive regions, such as those tested on rice chromosome 4, RelocaTE2 is  
32 able to report up to 95% of simulated TE insertions with less than 0.1% false positive rate using  
33 10-fold genome coverage resequencing data. RelocaTE2 provides a robust solution to identify  
34 TE polymorphisms and can be incorporated into analysis workflows in support of describing the  
35 complete genotype from light coverage genome sequencing.

## 36 **Introduction**

37 Transposable elements (TE), mobile DNA of the genome, are drivers of genomic innovation  
38 (Bennetzen & Wang 2014; Cordaux & Batzer 2009). They can act as mutagens to disrupt gene  
39 functions or induce novel gene functions by providing enhancers or promoters that alter host  
40 gene expression (Feschotte 2008; Lisch 2013). In plants, TEs have been shown to contribute to  
41 several key trait innovations in crop domestication (Lisch 2013). Systematic analysis of TE  
42 insertions and gene expression also suggests widespread roles of TEs in altering gene regulation  
43 (Kunarso et al. 2010; Lynch et al. 2011; Sundaram et al. 2014). It was found that 600-2000  
44 genetic variants between individuals in the human population and 200-300 variants between  
45 *Arabidopsis* accessions could be attributed to TE polymorphism (Quadrana et al. 2016; Stewart  
46 et al. 2011). Although the magnitude of these polymorphisms is small compared to SNPs or  
47 other insertion/deletions, some TE polymorphisms are proximal to protein coding genes and can  
48 have large impacts on gene function or gene regulation (Cowley & Oakey 2013; Quadrana et al.  
49 2016; Stewart et al. 2011).

50 Two categories of bioinformatics tools have been developed to identify TE polymorphisms from  
51 population resequencing data. One type employs a strategy similar to that used to discover  
52 structural variations. These tools identify discordant pairs of sequence reads based on the  
53 chromosomal position of read alignments to indicate genomic inversions, insertions, deletions or  
54 other complex rearrangements (Campbell et al. 2008; Korbel et al. 2007). Software for TE  
55 mapping scrutinize genomic loci with discordant read pairs to see if known TE sequences are can  
56 be implicated near the rearrangement site. These tools, such as Retroseq (Keane et al. 2013) and  
57 TEMP (Zhuang et al. 2014), are generally highly sensitive and can locate insertion sites to a 10-  
58 50 bp resolution. A second category of tools first identify by similarity, any sequence reads  
59 containing known TE sequences. The tools excise the TE sequence from the reads and search the  
60 remaining 5' or 3' flanking sequence against the host organism genome sequence to find the  
61 element's genomic location. These tools, including RelocaTE (Robb et al. 2013), T-lex2 (Fiston-  
62 Lavier et al. 2015), and ITIS (Jiang et al. 2015), are able to detect the exact location of insertion  
63 sites and TSDs characteristic of TEs. This second category of tools is ideal for identifying new  
64 insertions from population resequencing data because it can accurately detect an insertion  
65 location and identify the sequence of TSD. However, most of these tools are designed to search  
66 with only a single TE at a time, which sacrifices speed for increased sensitivity and specificity.

67 The extended runtime limits the feasibility of applying these tools when searching thousands of  
68 TEs in hundreds or thousands of individuals.

69 In RelocaTE2, an improved version of RelocaTE, we implement a junction-based approach that  
70 can search multiple template TEs in the same pass through the sequencing data, streamlining the  
71 computational approach. Using simulated datasets, we show that RelocaTE2 is highly sensitive  
72 even in low coverage resequencing data or on chromosomes with high repetitive sequence  
73 content has a specificity of greater than 99%. RelocaTE2 performed as the most sensitive and  
74 specific tool in our tests profiling human and rice population genomics data and can be widely  
75 used for analyzing population dynamics of TEs.

## 76 **Materials & Methods**

### 77 **RelocaTE2 Workflow**

78 RelocaTE2 is based on the previous algorithm implemented in RelocaTE (Robb et al. 2013),  
79 which uses junction reads to find insertion sites of TEs. In RelocaTE2, we re-implement the  
80 search strategy to enable identification of multiple TEs in a single search, greatly increasing the  
81 speed and enabling searches for hundreds or thousands of candidate TE families in a genome  
82 (Fig.1). We also implement new features in the algorithm to automatically identify TSDs and  
83 remove false junction reads (Fig.1).

84 Briefly, the workflow initiates by matching a library of known repeat elements against short  
85 sequence reads generated by next generation sequencing, typically Illumina paired-end reads,  
86 using BLAT with the sensitive setting "`-minScore=10 -tileSize=7`" (Kent 2002; Robb  
87 et al. 2013). Every read with similarity to repeat elements is denoted as an informative read and  
88 will contain a partial or complete copy of a TE. Informative reads that contain partial matches at  
89 the boundaries of the repeat elements are trimmed to remove the TE region so that only the  
90 regions flanking the element remain in either one or both of the paired-end reads (denoted as  
91 junction reads). Untrimmed versions of each junction read and its pair (denoted as full reads) are  
92 used as controls to filter false positive junction reads.

93 Sequence reads comprised entirely of repeat elements are ignored, but their read pair is kept  
94 (denoted as supporting reads). These junction, full, and supporting reads, are all aligned to the  
95 reference genome using BWA (v0.6.2) with the default setting "`-l 32 -k 2`" (Li & Durbin  
96 2009). Mapped reads are sorted by chromosome order and windows of 2,000 bp are evaluated to

97 define insertion clusters. In each insertion cluster, additional subclusters are further refined based  
98 on the mapping position of junction reads to address the possible scenario of multiple insertion  
99 sites within a window. TSD position and sequence are identified if the subcluster is supported by  
100 junction reads from both upstream and downstream of the TE insertion site.

101 Next, a series of cleaning steps are used to filter low quality candidate insertion sites: i.) remove  
102 insertion sites that are only supported by low quality junction reads (map quality < 29); ii.)  
103 remove insertion site only supported by less than 3 junction reads total on the left & right flank  
104 when there are additional insertion sites which pass these filters in the same window. iii.) remove  
105 insertion sites only supported by junction reads and located within 10 bp range of an annotated  
106 TE in the reference genome. RelocaTE2 output reports the number of junction reads and  
107 supporting reads from both upstream and downstream of candidate TE insertion sites. Only  
108 confident insertions, defined as having at least one supporting junction read flanking the  
109 upstream or downstream of insertion sites and at least one junction read or one supporting read  
110 supporting the other end of TE insertion, are provided in the default output:

111 "ALL.all\_nonref\_insert.gff". Additional information about all candidate sites are provided in  
112 alternative output file: "ALL.all\_nonref\_insert.raw.gff".

### 113 **Simulated data for evaluation of TE insertion tools**

114 Simulated datasets were created by randomly inserting TEs into sequence records of  
115 chromosomes 3 (OsChr3) and 4 (OsChr4) of rice (*Oryza sativa japonica*). OsChr3 is primarily  
116 made up of euchromatic regions, whereas OsChr4 has the sequence complexity consistent with  
117 heterochromatic regions and is a typical feature of many plant genomes (Zhao et al. 2002).  
118 Fourteen TE families found in rice genomes comprised of 7 DNA Transposons: *mPing*, *nDart*,  
119 *Gaijin*, *spmlike*, *Truncator*, *mGing*, *nDarz* and 7 RNA Retrotransposons: *Bajie*, *Dasheng*,  
120 *Retro1*, *RIRE2*, *RIRE3*, *Copia2*, *karma*, were used. The insertion simulations were performed by  
121 choosing 200 random insertion sites on each chromosome in three independent replicates. Each  
122 simulated insertion site was generated by selecting one random chromosome position and then  
123 one random TE and TSDs of the expected size was generated for each TE. After generating 200  
124 random sites and TE assignments, a new genome sequence was generated with the TEs inserted  
125 at corresponding locations. A GFF3 file with the insertion locations recorded was produced to  
126 support the evaluation of the performance of each tool on the dataset. Paired-end reads of all  
127 simulated chromosomes were simulated by pIRS (pirs simulate -l 100 -x coverage -m 500 -v

128 100) (Hu et al. 2012). For each dataset, simulate sequence reads at sequence depths of 1, 2, 3, 4,  
129 5, 6, 7, 8, 9, 10, 15, 20, 40-fold coverage were generated.

### 130 **Real sequence data for evaluation of TE insertion tools**

131 Three sets of data, an individual human genome, HuRef, an individual rice genome, IR64, and  
132 population resequencing data of 50 rice and wild rice genomes (Levy et al. 2007; Schatz et al.  
133 2014; Xu et al. 2012), were used to evaluate the performance of RelocaTE2 and TEMP. High  
134 quality reference genome assemblies of HuRef and IR64 were used to evaluate the accuracy of  
135 TE genotyping tools by comparing the assembled sequences to the reference genome. The HuRef  
136 (also known as Venter genome) has been extensively studied for TE insertions (Xing et al.  
137 2009). Previous work identified 574 *Alu* elements that have been experimentally verified and can  
138 be treated as a gold standard data set for evaluation (Hormozdiari et al. 2010; Xing et al. 2009).  
139 Paired-end sequence reads of 10-fold depth were simulated from HuRef as test dataset by pIRS  
140 (`pirs simulate -l 100 -x coverage -m 500 -v 100`) (Hu et al. 2012).  
141 RelocaTE2 and TEMP were tested and their results compared to the Human Genome Reference  
142 Consortium genome (GRCh36 or hg18). A second dataset, the finished reference genome  
143 assembly of rice strain IR64, was explored utilizing available Illumina sequencing reads (Schatz  
144 et al. 2014). RelocaTE2 and TEMP were tested on libraries of 100 bp paired-end Illumina short  
145 reads (SRA accession: SRR546439) aligned to the rice reference genome (MSU7). A third  
146 dataset of resequencing data from 50 strains of rice and wild rice population with an average  
147 sequencing depth of 17-fold. RelocaTE2 and TEMP were tested on the sequencing libraries from  
148 each of these 50 strains to assess their consistency across datasets with varying sequence depth  
149 and genetic diversity. RelocaTE and ITIS were not included in the biological data testing  
150 because of the prohibitively long run times on these large datasets and their poor performance on  
151 simulated datasets.

### 152 **Detection of TE insertions using RelocaTE2, RelocaTE, TEMP and ITIS**

153 RelocaTE2, RelocaTE, TEMP and ITIS were run with default parameter settings on simulated  
154 data. The results were filtered to evaluate the best performing parameters for each tool.  
155 RelocaTE2 was tested with parameters "`--len_cut_match 10 --len_cut_trim 10`  
156 `--mismatch 2 --aligner blat`", which uses BLAT as the search engine (`--aligner`  
157 `blat`), allows for 2 mismatches (`--mismatch 2`) in matched sequence between reads and  
158 repeat elements (`--len_cut_match 10`), and only keeps sequence fragments that have at

159 least 10 bp after trimming repeat elements from reads (`--len_cut_trim 10`). RelocaTE  
160 was tested using parameters "`--len_cutoff 10 --mismatch 0`", which uses BLAT as  
161 search engine by default and allowed 0 bp mismatch (`--mismatch 0`) for matched sequence  
162 between reads and repeat elements (`--len_cutoff 10`). It should be noted that the mismatch  
163 setting in RelocaTE is the ratio of base pairs in the alignment between reads and repeat elements  
164 that can be mismatched, not an integer number of allowed mismatches, as used in RelocaTE2.  
165 Singleton calls from RelocaTE's results, which are sites supported by only one read, were  
166 removed. TEMP was tested with parameters "`-m 3`", which allow for three mismatches  
167 between reads and repeat elements. Singleton calls from TEMP's results were removed when  
168 testing on simulated data to achieve a balance between sensitivity and specificity. ITIS was  
169 tested with default parameters, which filtered TE calls with at least one read supporting from  
170 both ends of TE insertions. For analysis of the HuRef genome, the IR64 genome and the 50 rice  
171 and wild rice strains, RelocaTE2 and TEMP were run with default parameter settings as  
172 described above. The TEMP results were filtered to keep only TE calls with supporting and/or  
173 junction reads from both ends of TE insertions to achieve a comparable balance between  
174 sensitivity and specificity.

## 175 **Results and Discussions**

### 176 **Performance of RelocaTE2, RelocaTE, TEMP and ITIS on simulated data**

177 RelocaTE2 was first compared to RelocaTE, TEMP and ITIS using the simulated datasets. Each  
178 dataset of simulated rice chromosomes, OsChr3 and OsChr4, was virtually sheared to simulated  
179 paired-end short reads at a coverage ranging from 1X to 40X. At high sequencing coverage  
180 ( $\geq 10X$ ), RelocaTE2, TEMP and ITIS were able to identify  $>99\%$  of simulated insertions on  
181 OsChr3, whereas the performance of RelocaTE was much lower (85%) (Fig.2A). At lower  
182 sequencing coverage, e.g. 3X, only RelocaTE2 and TEMP were able to achieve  $\geq 95\%$  sensitivity  
183 on OsChr3 (Fig.2A). Furthermore, TEMP was able to identify 83% and 93% of simulated  
184 insertions on OsChr3 at very low sequence coverage of 1X and 2X, respectively (Fig.2A).  
185 RelocaTE2 had a sensitivity of 53% and 83% on OsChr3 for the 1X and 2X coverage due to the  
186 removal of TE insertions supported by only one read (singleton) or supported by reads from only  
187 one end of TE insertions (insufficient insertions), which can result in many false positives  
188 (Fig.2A).

189 RelocaTE2, RelocaTE and TEMP showed >99% specificity on OsChr3 at multiple levels of  
190 sequence coverage (Fig.2B). In contrast, the specificity of ITIS was much lower (<90%), even  
191 when run on the high sequence coverage dataset on OsChr3 (Fig.2B). In comparing recall rates  
192 of TSDs, RelocaTE2 and ITIS had similar performance and achieved the highest recall rate of  
193 98% and 91% respectively, on OsChr3 at ~10X coverage (Fig.2C). The recall rate of TSDs for  
194 both TEMP and RelocaTE depended on sequence depth and achieved only 37% and 60%,  
195 respectively, at 10X coverage (Fig.2C). All the tools performed worse on OsChr4 as compared to  
196 OsChr3 (Fig.2D-F). RelocaTE2 demonstrated a lower average sensitivity (92%) on OsChr4  
197 when compared OsChr3 (96%) (Fig.2A,D). Similarly, TEMP had a slightly lower sensitivity  
198 (95%) on OsChr4 than on OsChr3 (97%) (Fig.2A,D). However, RelocaTE2 and RelocaTE  
199 demonstrated high level of the specificity (>99%) while TEMP performed at a slightly lower  
200 specificity (98%) on OsChr4 compared to >99% on OsChr3 (Fig.2B,E). In comparing TSD  
201 accuracy on OsChr4, on average 81% of RelocaTE2 calls correctly identified the TSD, whereas  
202 only 31% of TEMP calls were correct (Fig.2C,F).

### 203 **Evaluation of RelocaTE2 and TEMP on biological datasets**

204 We evaluated TE identifying tools in the HuRef genome and benchmark the sensitivity and  
205 specificity of these tools using 574 experimental verified *Alu* insertions in HuRef genome and  
206 genomic comparison between HuRef genome and GRCh36. RelocaTE2 and TEMP reported  
207 similar results and identified 83% (479/574) and 76% (438/574) of standard insertion sites  
208 (Fig.3A). Comparing the HuRef genome with GRCh36 suggested that 89% and 95% of  
209 insertions identified by RelocaTE2 and TEMP, respectively, were real insertions (Fig.3A). In  
210 addition, RelocaTE2 predicted TE insertion sites with higher precision ( $9 \pm 6$  bp) compared to  
211 TEMP ( $366 \pm 170$  bp).

212 RelocaTE2 and TEMP were used to analyze data from the rice strain IR64 and the results were  
213 evaluated by comparing the genome assembly of IR64 with MSU7. RelocaTE2 identified 648  
214 insertion sites while the genome comparison revealed that 93% of insertions were true positives  
215 (Fig.3A). TEMP identified 362 insertions, of which 50% (183/362) overlapped with RelocaTE2  
216 (Fig.3A). The specificity of TEMP was estimated to be 86%, slightly lower than RelocaTE2  
217 (93%) (Fig.3A). However, TEMP was found to be less sensitive than RelocaTE2 in the rice  
218 genome, only calling 362 sites as compared to 648 by RelocaTE2 (38% vs. 90%, Fig.3A).

219 Moreover, RelocaTE2 predicted TE insertion junctions of  $3 \pm 1$  bp, which was much smaller  
220 than TEMP ( $393 \pm 199$  bp).

221 RelocaTE2 and TEMP were used to identify TE polymorphisms in 50 resequenced rice and wild  
222 rice strains, which contain substantial sequence diversity and population structure (Xu et al.  
223 2012). The results from these two tools were well correlated ( $R^2 = 0.96$ ,  $P$  value =  $2.2e-16$ ) and  
224 predicted more TE insertions in the diverged population of wild rice, *O. nivara* and *O. rufipogon*,  
225 and even in the *indica* population than *japonica* rice which close to the reference genome  
226 (Fig.3B). On average 72% of the sites predicted in these 50 rice and wild rice strain by  
227 RelocaTE2 and TEMP overlapped. Many insertion sites from TEMP were predicted with only  
228 supporting read flanking one end of an insertion, which produced large variations in predicted  
229 junctions of TE insertion sites ( $118 \pm 151$  bp). In contrast, RelocaTE2 reported most of TE  
230 insertions supported by junction reads or supporting reads on both ends, which resulted in  
231 accurate insertion junction predictions ( $3 \pm 2$  bp).

### 232 **Runtime performance**

233 We implemented the searching process for TE insertion to run on multiple processors in Python.  
234 The process is relatively memory efficient. When searching TEs in the rice genome for example,  
235 one process generally uses less than 1 Gb memory. The running time of RelocaTE2 depends on  
236 number of processors used. Searching 3000 templates of transposable elements with 20X  
237 genome coverage sequencing data of the rice genome takes 3-4 hours for RelocaTE2 using 32  
238 CPUs including the alignment steps. TEMP identifies transposable element insertions from a  
239 BAM file. It takes ~1 hours for TEMP for the same project using single process. RelocaTE  
240 (version 1) and ITIS take at least days for the same rice datasets and can be prohibitively difficult  
241 to run on large datasets with multiple templates due to the serial searching approach of their  
242 implementation.

### 243 **Conclusions**

244 We present RelocaTE2 as a new tool for mapping TE polymorphisms to base-pair resolution  
245 from resequencing data. RelocaTE2 identifies multiple TE families in a single search with high  
246 sensitivity and specificity. The evaluation of these tools on simulated and biological datasets  
247 support the use of RelocaTE2 for analysis of genomes of plants and animals and indicate it can  
248 generate very high quality genotyping of TE insertions from resequencing datasets of modest

249 sequencing depths. The high resolution mapping of TE insertions sites will enable detailed  
250 analysis of the interaction of TEs and genes and as structural variations that vary in populations.

251

### 252 **Competing Interests**

253 The authors declare that there are no competing interests.

### 254 **Author Contributions**

255 Jinfeng Chen conceived and designed the study, wrote the code, analyzed the data, wrote the  
256 paper, and prepared figures and/or tables.

257 Travis R. Wrightsman tested the code, wrote the manual, and reviewed drafts of the papers.

258 Susan R. Wessler and Jason E. Stajich conceived and designed the study, wrote the paper, and  
259 reviewed drafts of the papers.

### 260 **Data Availability**

261 The source code in Python, manual, and sample data of RelocaTE2 are available for download at  
262 <https://github.com/stajichlab/RelocaTE2>.

263

264

### 265 **Acknowledgements**

### 266 **Funding**

267 This research is supported by US National Science Foundation grant IOS-1027542 and a grant  
268 from the W. M. Keck Foundation to SR Wessler and JE Stajich and funds to SR Wessler from  
269 the Howard Hughes Medical Institute 52008110 and US Department of Agriculture Grant 2013-  
270 38422-20955.

### 271 **References**

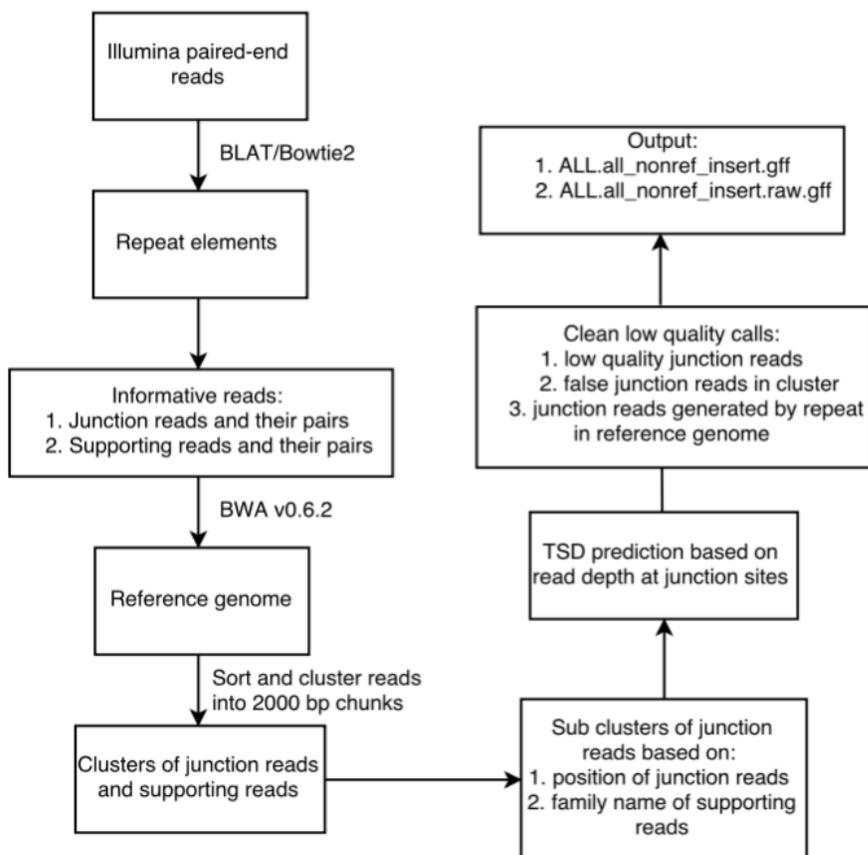
272 Bennetzen JL, and Wang H. 2014. The contributions of transposable elements to the  
273 structure, function, and evolution of plant genomes. *Annu Rev Plant Biol* 65:505-530.  
274 10.1146/annurev-arplant-050213-035811

275 Campbell PJ, Stephens PJ, Pleasance ED, O'Meara S, Li H, Santarius T, Stebbings LA, Leroy C,  
276 Edkins S, Hardy C, Teague JW, Menzies A, Goodhead I, Turner DJ, Clee CM, Quail MA,  
277 Cox A, Brown C, Durbin R, Hurler ME, Edwards PA, Bignell GR, Stratton MR, and  
278 Futreal PA. 2008. Identification of somatically acquired rearrangements in cancer

- 279 using genome-wide massively parallel paired-end sequencing. *Nat Genet* 40:722-  
280 729. 10.1038/ng.128
- 281 Cordaux R, and Batzer MA. 2009. The impact of retrotransposons on human genome  
282 evolution. *Nat Rev Genet* 10:691-703. 10.1038/nrg2640
- 283 Cowley M, and Oakey RJ. 2013. Transposable elements re-wire and fine-tune the  
284 transcriptome. *PLoS Genet* 9:e1003234. 10.1371/journal.pgen.1003234
- 285 Feschotte C. 2008. Transposable elements and the evolution of regulatory networks. *Nat*  
286 *Rev Genet* 9:397-405. 10.1038/nrg2337
- 287 Fiston-Lavier AS, Barron MG, Petrov DA, and Gonzalez J. 2015. T-lex2: genotyping,  
288 frequency estimation and re-annotation of transposable elements using single or  
289 pooled next-generation sequencing data. *Nucleic Acids Res* 43:e22.  
290 10.1093/nar/gku1250
- 291 Hormozdiari F, Hajirasouliha I, Dao P, Hach F, Yorukoglu D, Alkan C, Eichler EE, and  
292 Sahinalp SC. 2010. Next-generation VariationHunter: combinatorial algorithms for  
293 transposon insertion discovery. *Bioinformatics* 26:i350-357.  
294 10.1093/bioinformatics/btq216
- 295 Hu X, Yuan J, Shi Y, Lu J, Liu B, Li Z, Chen Y, Mu D, Zhang H, Li N, Yue Z, Bai F, Li H, and Fan  
296 W. 2012. pIRS: Profile-based Illumina pair-end reads simulator. *Bioinformatics*  
297 28:1533-1535. 10.1093/bioinformatics/bts187
- 298 Jiang C, Chen C, Huang Z, Liu R, and Verdier J. 2015. ITIS, a bioinformatics tool for accurate  
299 identification of transposon insertion sites using next-generation sequencing data.  
300 *BMC Bioinformatics* 16:72. 10.1186/s12859-015-0507-2
- 301 Keane TM, Wong K, and Adams DJ. 2013. RetroSeq: transposable element discovery from  
302 next-generation sequencing data. *Bioinformatics* 29:389-390.  
303 10.1093/bioinformatics/bts697
- 304 Kent WJ. 2002. BLAT--the BLAST-like alignment tool. *Genome Res* 12:656-664.  
305 10.1101/gr.229202. Article published online before March 2002
- 306 Korbelt JO, Urban AE, Affourtit JP, Godwin B, Grubert F, Simons JF, Kim PM, Palejev D,  
307 Carriero NJ, Du L, Taillon BE, Chen Z, Tanzer A, Saunders AC, Chi J, Yang F, Carter NP,  
308 Hurles ME, Weissman SM, Harkins TT, Gerstein MB, Egholm M, and Snyder M. 2007.

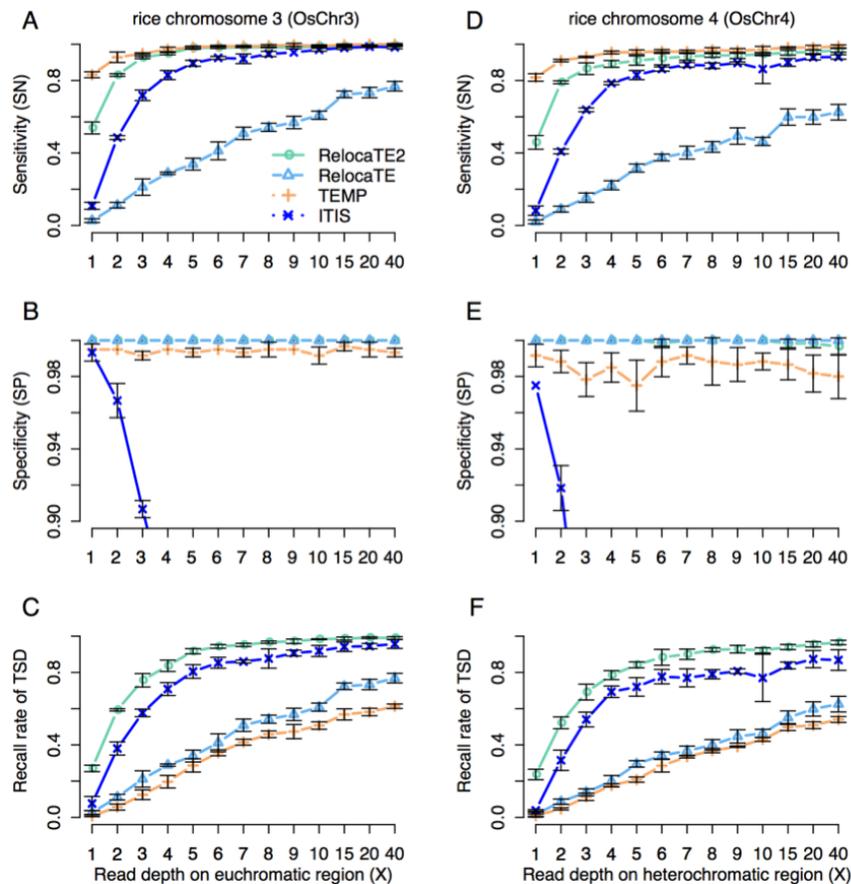
- 309 Paired-end mapping reveals extensive structural variation in the human genome.  
310 *Science* 318:420-426. 10.1126/science.1149504
- 311 Kunarso G, Chia NY, Jeyakani J, Hwang C, Lu X, Chan YS, Ng HH, and Bourque G. 2010.  
312 Transposable elements have rewired the core regulatory network of human  
313 embryonic stem cells. *Nat Genet* 42:631-634. 10.1038/ng.600
- 314 Lee E, Iskow R, Yang L, Gokcumen O, Haseley P, Luquette LJ, 3rd, Lohr JG, Harris CC, Ding L,  
315 Wilson RK, Wheeler DA, Gibbs RA, Kucherlapati R, Lee C, Kharchenko PV, Park PJ,  
316 and Cancer Genome Atlas Research N. 2012. Landscape of somatic  
317 retrotransposition in human cancers. *Science* 337:967-971.  
318 10.1126/science.1222077
- 319 Levy S, Sutton G, Ng PC, Feuk L, Halpern AL, Walenz BP, Axelrod N, Huang J, Kirkness EF,  
320 Denisov G, Lin Y, MacDonald JR, Pang AW, Shago M, Stockwell TB, Tsiamouri A,  
321 Bafna V, Bansal V, Kravitz SA, Busam DA, Beeson KY, McIntosh TC, Remington KA,  
322 Abril JF, Gill J, Borman J, Rogers YH, Frazier ME, Scherer SW, Strausberg RL, and  
323 Venter JC. 2007. The diploid genome sequence of an individual human. *PLoS Biol*  
324 5:e254. 10.1371/journal.pbio.0050254
- 325 Li H, and Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler  
326 transform. *Bioinformatics* 25:1754-1760. 10.1093/bioinformatics/btp324
- 327 Lisch D. 2013. How important are transposons for plant evolution? *Nat Rev Genet* 14:49-61.  
328 10.1038/nrg3374
- 329 Lynch VJ, Leclerc RD, May G, and Wagner GP. 2011. Transposon-mediated rewiring of gene  
330 regulatory networks contributed to the evolution of pregnancy in mammals. *Nat*  
331 *Genet* 43:1154-1159. 10.1038/ng.917
- 332 Quadrana L, Bortolini Silveira A, Mayhew GF, LeBlanc C, Martienssen RA, Jeddloh JA, and  
333 Colot V. 2016. The *Arabidopsis thaliana* mobilome and its impact at the species level.  
334 *Elife* 5. 10.7554/eLife.15716
- 335 Robb SM, Lu L, Valencia E, Burnette JM, 3rd, Okumoto Y, Wessler SR, and Stajich JE. 2013.  
336 The use of RelocaTE and unassembled short reads to produce high-resolution  
337 snapshots of transposable element generated diversity in rice. *G3 (Bethesda)* 3:949-  
338 957. 10.1534/g3.112.005348

- 339 Schatz MC, Maron LG, Stein JC, Hernandez Wences A, Gurtowski J, Biggers E, Lee H, Kramer  
340 M, Antoniou E, Ghiban E, Wright MH, Chia JM, Ware D, McCouch SR, and McCombie  
341 WR. 2014. Whole genome de novo assemblies of three divergent strains of rice,  
342 *Oryza sativa*, document novel gene space of aus and indica. *Genome Biol* 15:506.  
343 10.1186/PREACCEPT-2784872521277375
- 344 Stewart C, Kural D, Stromberg MP, Walker JA, Konkel MK, Stutz AM, Urban AE, Grubert F,  
345 Lam HY, Lee WP, Busby M, Indap AR, Garrison E, Huff C, Xing J, Snyder MP, Jorde LB,  
346 Batzer MA, Korbel JO, Marth GT, and Genomes P. 2011. A comprehensive map of  
347 mobile element insertion polymorphisms in humans. *PLoS Genet* 7:e1002236.  
348 10.1371/journal.pgen.1002236
- 349 Sundaram V, Cheng Y, Ma Z, Li D, Xing X, Edge P, Snyder MP, and Wang T. 2014. Widespread  
350 contribution of transposable elements to the innovation of gene regulatory  
351 networks. *Genome Res* 24:1963-1976. 10.1101/gr.168872.113
- 352 Xing J, Zhang Y, Han K, Salem AH, Sen SK, Huff CD, Zhou Q, Kirkness EF, Levy S, Batzer MA,  
353 and Jorde LB. 2009. Mobile elements create structural variation: analysis of a  
354 complete human genome. *Genome Res* 19:1516-1526. 10.1101/gr.091827.109
- 355 Xu X, Liu X, Ge S, Jensen JD, Hu F, Li X, Dong Y, Gutenkunst RN, Fang L, Huang L, Li J, He W,  
356 Zhang G, Zheng X, Zhang F, Li Y, Yu C, Kristiansen K, Zhang X, Wang J, Wright M,  
357 McCouch S, Nielsen R, Wang J, and Wang W. 2012. Resequencing 50 accessions of  
358 cultivated and wild rice yields markers for identifying agronomically important  
359 genes. *Nat Biotechnol* 30:105-111. 10.1038/nbt.2050
- 360 Zhao Q, Zhang Y, Cheng Z, Chen M, Wang S, Feng Q, Huang Y, Li Y, Tang Y, Zhou B, Chen Z, Yu  
361 S, Zhu J, Hu X, Mu J, Ying K, Hao P, Zhang L, Lu Y, Zhang LS, Liu Y, Yu Z, Fan D, Weng  
362 Q, Chen L, Lu T, Liu X, Jia P, Sun T, Wu Y, Zhang Y, Lu Y, Li C, Wang R, Lei H, Li T, Hu  
363 H, Wu M, Zhang R, Guan J, Zhu J, Fu G, Gu M, Hong G, Xue Y, Wing R, Jiang J, and Han  
364 B. 2002. A fine physical map of the rice chromosome 4. *Genome Res* 12:817-823.  
365 10.1101/gr.48902
- 366 Zhuang J, Wang J, Theurkauf W, and Weng Z. 2014. TEMP: a computational method for  
367 analyzing transposable element polymorphism in populations. *Nucleic Acids Res*  
368 42:6826-6838. 10.1093/nar/gku323



369

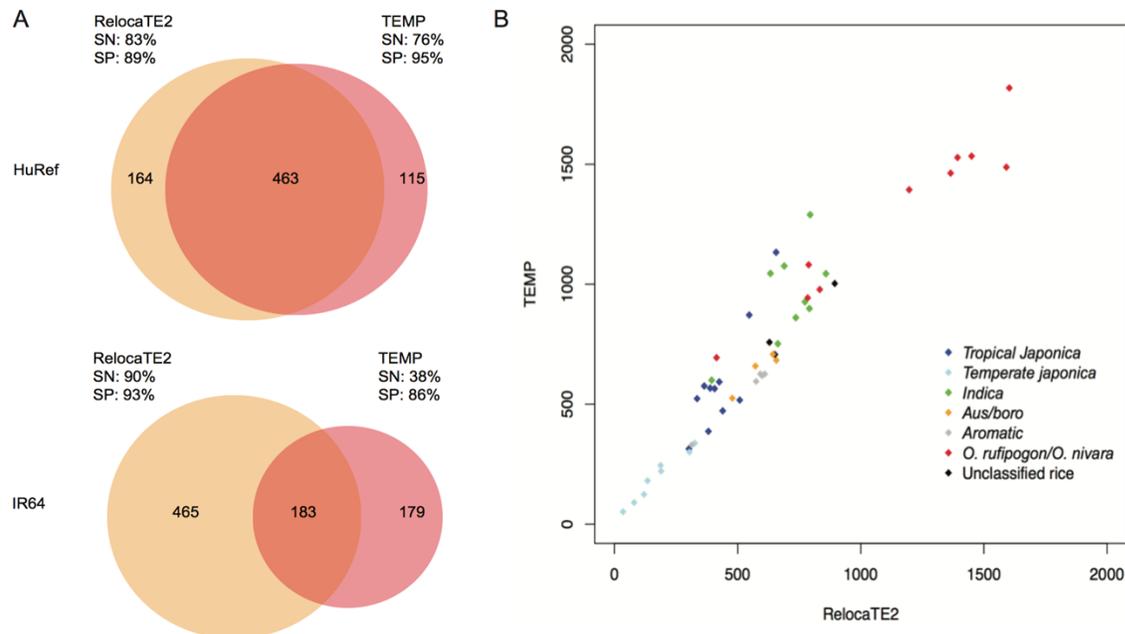
370 **Figure 1. Workflow for identification of transposable element insertions in population**371 **resequencing data using Illumina paired-end reads.**



372

373 **Figure 2. Performance of RelocaTE2, RelocaTE, TEMP and ITIS on simulated rice data.**

374 Simulations of 200 random transposable element (TE) insertions were generated for rice  
 375 chromosome 3 (OsChr3) and rice chromosome 4 (OsChr4) with three replicates. A series of  
 376 datasets with different sequence depths (from 1X to 40X) were generated for each simulation  
 377 dataset. Sensitivity (SN), Specificity (SP) and Recall rate of target site duplication (TSD) of each  
 378 tool were estimated for each of these datasets and plotted against sequence depth. The error bars  
 379 show the standard deviation of three replicates with different sets of 200 random TE insertions.  
 380 SN was defined as the percentage of calls within 100 base pairs of 200 random TE insertions. SP  
 381 was defined as the percentage of calls not within 100 base pairs of 200 random TE insertions.  
 382 Recall rate of TSD was defined as the percentage of true positive calls that correctly matched the  
 383 simulated TSD of TE insertions. The results illustrate how tools perform on chromosomes which  
 384 are primarily euchromatic or heterochromatic using OsChr3 and OsChr4 respectively.



385  
 386 **Figure 3. Performance of RelocaTE2 and TEMP on biological dataset in HuRef genome,**  
 387 **IR64 genome, and 50 rice and wild rice strains.** **A.** Venn diagram of the overlap in non-  
 388 reference TE insertions identified in the HuRef genome and the rice IR64 genome using  
 389 RelocaTE2 and TEMP. Sensitivity (SN) and Specificity (SP) were assessed by comparing the  
 390 assembled HuRef genome to the GRCh36 reference genome and the assembled IR64 genome to  
 391 the MSU7 reference genome. SN was defined as the percentage of validated calls out of all  
 392 validated calls by either RelocaTE2 or TEMP. SP was defined as the percentage of validated  
 393 calls out of all calls by each tool. **B.** Comparison of the number of non-reference TE insertions of  
 394 14 TE families in 50 rice and wild rice strains identified by RelocaTE2 and TEMP. Strains are  
 395 color-coded based on subpopulation classification.