# Sensitive detection of pre-integration intermediates of long terminal repeat retrotransposons in crop plants

Cho, Jungnam; Benoit, Matthias; Catoni, Marco; Drost, Hajk-Georg; Brestovitsky, Anna; Oosterbeek, Matthijs; Paszkowski, Jerzy

[Link to publication on Research at Birmingham portal](#)

1  **Title**

2  Sensitive detection of pre-integration intermediates of LTR retrotransposons in crop plants

3

4  **Authors**

5  Jungnam Cho[1,2,3,]*, Matthias Benoit[1], Marco Catoni[1,4], Hajk-Georg Drost[1], Anna Brestovitsky[1],

6  Matthijs Oosterbeek[5,6] and Jerzy Paszkowski[1,]*

7  [1]The Sainsbury Laboratory, University of Cambridge, Cambridge CB2 1LR, UK

8  [2]National Key Laboratory of Plant Molecular Genetics (NKLPMG), CAS Center for Excellence

9  in Molecular Plant Sciences, Institute of Plant Physiology and Ecology (SIPPE), 200032

10  Shanghai, P. R. China

11  [3]CAS-JIC Centre of Excellence for Plant and Microbial Science (CEPAMS), Chinese Academy

12  of Sciences, 200032 Shanghai, P. R. China

13  [4]School of Biosciences, University of Birmingham, Birmingham BI5 2TT, UK

14  [5]Laboratory of Molecular Biology, Wageningen University, Wageningen 6708PB, The

15  Netherlands

16  [6]Current address: Laboratory of Nematology, Wageningen University, Wageningen 6708PB,

17  The Netherlands

18  *Corresponding author

19

20  **Correspondence:**

21  Jungnam Cho jungnam.cho@slcu.cam.ac.uk and Jerzy Paszkowski

22  jerzy.paszkowski@slcu.cam.ac.uk

23

**Abstract**

Retrotransposons have played an important role in the evolution of host genomes[1,2]. Their impact is mainly deduced from the composition of DNA sequences that have been fixed over evolutionary time[2]. Such studies provide important "snapshots" reflecting the historical activities of transposons but do not predict current transposition potential. We previously reported Sequence-Independent Retrotransposon Trapping (SIRT) as a method that, by identification of extrachromosomal linear DNA (eclDNA), revealed the presence of active LTR retrotransposons in *Arabidopsis*[3]. However, SIRT cannot be applied to large and transposon-rich genomes, as found in crop plants. We have developed an alternative approach named ALE-seq (*a*mplification of *L*TR of *e*clDNAs followed by *seq*uencing) for such situations. ALE-seq reveals sequences of 5' LTRs of eclDNAs after two-step amplification: *in vitro* transcription and subsequent reverse transcription. Using ALE-seq in rice, we detected eclDNAs for a novel *Copia* family LTR retrotransposon, *Go-on*, which is activated by heat stress. Sequencing of rice accessions revealed that *Go-on* has preferentially accumulated in *indica* rice grown at higher temperatures. Furthermore, ALE-seq applied to tomato fruits identified a developmentally regulated *Gypsy* family of retrotransposons. A bioinformatic pipeline adapted for ALE-seq data analyses is used for the direct and reference-free annotation of new, active retroelements. This pipeline allows assessment of LTR retrotransposon activities in organisms for which genomic sequences and/or reference genomes are either unavailable or of low quality.

**Main**

Chromosomal copies of activated retrotransposons containing long terminal repeats (LTRs) are transcribed by RNA polymerase II, followed by reverse transcription of transcripts to extrachromosomal linear DNAs (eclDNA); these integrate back into host chromosomes[3]. Because of the two obligatory template switches during reverse transcription, the newly synthetized eclDNA is flanked by LTRs of identical sequence. Their subsequent divergence due to the accumulation of mutations correlates well with length of time since the last transposition, and thus transposon age[4]. However, the age of LTR retrotransposons cannot

54    be used to predict their current transpositional potential. Moreover, predictions are further

55    complicated by recombination events that occur with high frequency between young and

56    old members of a retrotransposon family[5]; thus old family members also contribute to the

57    formation of novel recombinant elements that insert into new chromosomal positions[5].

58    Although, retrotransposon activities can be relatively easily measured at the transcriptional

59    level[6], the presence of transcripts is a poor predictor of transpositional potential due to

60    posttranscriptional control of this process[7,8]. In addition, direct detection of transposition by

61    genome-wide sequencing to identify new insertions is too expensive and time-consuming to

62    be applied as a screening method. Clearly, the development of an expeditious approach to

63    identify active retrotransposons that predict their transposition potential would be

64    welcomed. We previously described the SIRT strategy for *Arabidopsis* that led to the

65    identification of eclDNA of a novel retroelement and subsequent detection of new

66    insertions[3]. Thus, the presence of eclDNAs, the last pre-integration intermediate, was shown

67    to be a good predictor of retrotransposition potential.

68

69    **Results**

70    <u>Development of ALE-seq</u>

71        Retrotransposons include a conserved sequence known as the primer binding site

72    (PBS), where binding of the 3' end of cognate tRNA initiates the reverse transcription

73    reaction[3]. Met-iCAT (Methionine tRNA-CAT anticodon) PBS was chosen for SIRT as it is the

74    site present in the majority of annotated *Arabidopsis* retrotransposons[3]. To examine

75    whether Met-iCAT PBS sequences are also predominant in LTR retrotransposons of other

76    plants, we used the custom-made software *LTRpred* for *de novo* annotation of LTR

77    retrotransposons in rice and tomato genomes (see Materials and methods). Young

78    retroelements were selected by filtering for at least 95% identity between the two LTRs and

79    subsequently examined for their cognate tRNAs (Supplementary Figure 1). As in *Arabidopsis*,

80    around 80% of LTR retrotransposons in the tomato genome contained Met-iCAT PBS

81    (Supplementary Figure 1). In contrast, only 30% harboured Met-iCAT PBS in rice, and Arg-

82    CCT (Arginine tRNA-CCT anticodon) PBS was found in 60% of young LTR retrotransposons

83    (Supplementary Figure 1). Nonetheless, we used Met-iCAT PBS in our initial experiments

84  because most retrotransposons known to be active in rice callus (e.g. *Tos17* and *Tos19*)

85  contain Met-iCAT PBS. Initially, SIRT was performed on DNA extracted from rice leaves and

86  calli; however, we did not detect eclDNAs for *Tos17* and *Tos19* in rice tissues by this method

87  (Supplementary Figure 2). We reasoned that the short stretch of PBS used for primer design

88  in SIRT may have impaired PCR efficiency due to the many PBS-related sequences present in

89  larger genomes containing a high number of retroelements, as is the case in rice.

90      To counter this problem, we developed an alternative method, named ALE-seq, with

91  significantly improved selectivity and sensitivity of eclDNA detection. A crucial difference to

92  SIRT is that ALE-seq amplification of eclDNA is separated into two reactions: *in vitro*

93  transcription and reverse transcription (Figure 1a). This decoupling of the use of the two

94  priming sequences followed by the digestion of non-templated DNA and RNA is significantly

95  more selective and efficient than the single PCR amplification in SIRT.

96      ALE-seq starts with ligation to the ends of eclDNA of an adapter containing a T7

97  promoter sequence at its 5' end and subsequent *in vitro* transcription with T7 RNA

98  polymerase. The synthesized RNA is then reverse transcribed using the primer that binds

99  the transcripts at the PBS site. The adapter and the oligonucleotides priming reverse

100  transcription are anchored with partial Illumina adapter sequences (Supplementary Table 1),

101  which allows the amplified products to be directly deep-sequenced in a strand-specific

102  manner. The ALE-seq-sequences derived from retrotransposon eclDNAs are predicted to

103  contain the intact 5' LTR up to the PBS site, flanked by Illumina paired-end sequencing

104  adapters. We used the Illumina MiSeq platform for sequencing because its long reads of 300

105  bp from both ends cover the entire LTR lengths of most potentially active elements. It is

106  worth noting that the Illumina adapters were tagged to the intact LTR DNA without

107  fragmentation of the amplicons. This together with the long reads of MiSeq allowed us to

108  reconstitute the complete LTR sequences, even in the absence of the reference genome

109  sequence. The reconstituted LTRs were analysed using the alignment-based approach that

110  complements the mapping-based approach when the reference genome is incomplete

111  (Figure 1b).

112      First, we tested ALE-seq on *Arabidopsis* by examining  heat-stressed Col-0

113  *Arabidopsis* plants[9], *met1-1* mutant[3] and epi12[8], a *met1*-derived epigenetic recombinant

114 inbred line. ALE-seq cleanly and precisely recovered sequences of complete LTRs for *Onsen*,

115 *Copia21* and *Evade* in samples containing their respective eclDNA (Supplementary Figure

116 3)[3,8,9]. Due to priming of the reverse transcription reaction at PBS, the reads were explicitly

117 mapped to the 5' but not to the 3' LTR, although the two LTRs have identical sequences. The

118 ALE-seq reads have well-defined extremities, starting at the position marking the start of

119 LTRs and finishing at the PBS, which is consistent with their eclDNA origin. The ends of LTRs

120 can also be inspected for conserved sequences that would further confirm their eclDNA

121 origin (Supplementary Figure 4). This reduced ambiguity of read mapping in ALE-seq analysis,

122 combined with the clear-cut detection of LTR ends, allows for explicit and precise

123 assignment of ALE-seq results to active LTR retrotransposons.

124 Since SIRT failed to detect eclDNAs of rice retrotransposons known to be activated in

125 rice callus, we examined whether ALE-seq would identify their eclDNAs. As shown in Figure

126 1c to f, ALE-seq unambiguously detected eclDNAs of *Tos17* and *Tos19* in rice callus, but not

127 in leaf samples. To test whether detection of 5' LTR sequences requires the entire ALE-seq

128 procedure, we performed control experiments with depleted ALE-seq reactions, for example,

129 in the absence of enzymes for either ligation, *in vitro* transcription, or reverse transcription.

130 All incomplete procedures failed to produce sequences containing 5' LTRs derived from

131 eclDNAs (Figure 1e and f). Taken together, the data show that ALE-seq can detect eclDNAs

132 of LTR retrotransposons in *Arabidopsis* as well as in rice with considerably greater efficiency

133 than the SIRT method.

134 To examine the suitability of ALE-seq for quantitative determination of eclDNA levels,

135 we carried out a reconstruction experiment spiking 100 ng of genomic DNA from rice callus

136 with differing amounts of PCR-amplified full-length *Onsen* DNA from 1 ng to 100 fg (Figure

137 2a to d). The results in Figure 2a and b show that the readouts of ALE-seq for *Onsen*

138 correlate well with the input amounts ($R^2$=0.99). The initial ALE-seq steps of ligation and *in*

139 *vitro* transcription impinged proportionally on the input DNA, resulting in unbiased

140 quantification of the eclDNA and minimal quantitative distortion of the final ALE-seq data.

141 Noticeably, the levels of *Tos17* were similar in all the spiked samples, indicating that

142 addition of *Onsen* DNA did not influence the detection sensitivity of *Tos17*, at least for the

143 amounts tested (Figure 2c and d). Thus, ALE-seq can be used to accurately determine

144 eclDNA levels.

145     Most rice retrotransposons harbour Arg-CCT PBS (Supplementary Figure 1). We

146     tested whether the reverse transcription reaction can be multiplexed to capture both types

147     of retrotransposons (containing Arg-CCT or Met-iCAT PBS) and whether multiplexing of the

148     reverse transcription primers compromises the sensitivity of the procedure. ALE-seq was

149     performed on DNA from rice callus, testing each of the reverse transcription primers

150     separately or as a mixture of both primers in a single reaction. As shown in Figure 2e and

151     Supplementary Figure 5, the levels of *Tos17* recorded in the samples with both primers were

152     similar to the Met-iCAT primer alone. Importantly, we also detected the eclDNAs of the

153     *RIRE2* element containing Arg-CCT PBS (Figure 2f), which was known to be transpositionally

154     active in rice callus[7].


155


156     <u>Identification of *Go-on* retrotransposon using ALE-seq</u>

157     We next used ALE-seq to search for novel active rice retrotransposons. Since many

158     plant retrotransposons are transcriptionally activated by abiotic stresses[9,10], we subjected

159     rice plants to heat stress before subjecting them to ALE-seq. In this way we identified a

160     *Copia*-type retrotransposon able to synthetize eclDNA in the heat-stressed plants (Figure 3a

161     to c) and named this element *Go-on* (the Korean for 'high temperature'). The three

162     retrotransposons with the highest eclDNA levels in heat-stress conditions all belong to the

163     *Go-on* family (Figure 3b and Supplementary Figure 6). Although, eclDNAs were detected for

164     all three copies, *Go-on3* seems to be the youngest and, thus, possibly the most active family

165     member, containing identical LTRs and a complete ORF (Supplementary Figure 6). As

166     depicted in Supplementary Figure 6, the 5' LTR sequences of the three *Go-on* copies are

167     identical; thus the ALE-seq reads derived from *Go-on3* LTR were also cross-mapped to other

168     copies that are possibly inactive or have reduced activities. To further determine whether

169     sequences of *Go-on* LTRs recovered by ALE-seq are indeed derived from *Go-on3* or also from

170     other family members, we performed an ALE-seq experiment using RT primers located

171     further downstream of the PBS, including sequences specific for each *Go-on* family member

172     (Supplementary Figure 6). The amplified ALE-seq products revealed that the eclDNAs

173     produced in heat-stressed rice originated only from *Go-on3* (Supplementary Figure 6). We

174   validated the production of eclDNAs of *Go-on3* by sequencing the junction of the adapter

175   and the 5' end of LTR (Supplementary Figure 6) and by qPCR (Supplementary Figure 7).

176        Next, we examined whether *Go-on3* is transcriptionally activated in rice subjected to

177   heat stress. RNA-seq and the RT-qPCR data clearly showed that *Go-on* is strongly activated

178   in heat-stress conditions (Figure 3d and Supplementary Figure 7). Similar to many other

179   retrotransposons including *ONSEN* of *Arabidopsis*[9,11,12], the LTR sequence of *Go-on3*

180   contains *cis*-acting regulatory element such as the heat shock transcription factor HSFC1-

181   binding sequence motif (Supplementary Figure 7), which is suggestive to its heat stress-

182   mediated transcriptional activation (Figure 3d). To determine whether *Go-on* is also

183   activated in *indica* rice, we heat-stressed plants of *IR64* for three days and examined *Go-on*

184   RNA and DNA levels. Similar to *japonica* rice, *Go-on* RNA and DNA accumulated markedly

185   under heat stress (Supplementary Figure 8), suggesting that the trigger for *Go-on* activation

186   is conserved in both of these evolutionarily distant rice genotypes. Analysis of the RNA-seq

187   data from the heat-stressed rice plants revealed a poor correlation between the mRNA and

188   eclDNA levels of retrotransposons (Supplementary Figure 9). Given that eclDNAs captured

189   by ALE-seq in *Arabidopsis* and rice (Figure 1c to f and Supplementary Figure 3) are all known

190   well for their transposition competence, this possibly agrees with the notion that the

191   eclDNA level is a better predictor of retrotransposition than the RNA level.

192        To possibly relate accumulation of *Go-on* copies in plant populations grown in

193   different temperatures, we analysed the historical retrotransposition of *Go-on* using the

194   genome resequencing data of rice accessions from the 3,000 Rice Genome Project[13]. First,

195   we retrieved the raw sequencing data for all 388 *japonica* rice accessions and the same

196   number of randomly selected sequences of *indica* rice accessions. Using the Transposon

197   Insertion Finder (TIF) tool[14], *japonica* and *indica* sequences were analysed for the number of

198   *Go-on* copies and their genome-wide distribution. Only non-reference insertions that were

199   absent in the reference genome were scored and the cumulative number of new insertions

200   was plotted (Figure 3e to g). Figure 3e shows that the *indica* rice population grown in a

201   warmer climate[15] accumulated significantly more *Go-on* copies than the *japonica* population.

202   As controls, we also examined the accumulation of *Tos17* and *Tos19*, which were not

203   activated by heat stress in our ALE-seq profile (Figure 3a and b). Both retrotransposons

204   showed more transposition events in *japonica* than in *indica* rice (Figure 3f, g and

7

205 Supplementary Figure 10). Therefore, the copy number of *Go-on* in rice accessions

206 correlated with their growth temperatures, which could possibly be related to occasional

207 *Go-on* activation in elevated ambient temperatures.

208

209 <u>Identification of *FIRE* retrotransposon using ALE-seq</u>

210          It was reported previously that the tomato genome (*Solanum lycopersicum*)

211 experiences a significant loss of DNA methylation in fruits during their maturation, which

212 leads to transcriptional activation of retrotransposons[16]. However, it was not known

213 whether these transcriptionally activated tomato transposons synthesise eclDNA. It was

214 questionable whether the ALE-seq strategy is sensitive enough to detect eclDNA in the ~950

215 Mb tomato genome, which is almost three times as large as ~400 Mb of rice[17]. To address

216 these questions, ALE-seq was carried out on DNA samples from fruits at 52 days post

217 anthesis (DPA), when the loss of DNA methylation is most pronounced[16], and from leaves as

218 a control. It is important to note that we used tomato cultivar (cv.) M82 for these

219 experiments, as it is commonly used for genetic studies[18,19], and that the sequence of the

220 current tomato reference genome is based on cv. Heinz 1706[17]. Since retrotransposon

221 sequences and their chromosomal distributions differ largely between genomes of different

222 varieties within the same plant species[20–22], we could not use the standard mapping-based

223 annotation of the ALE-seq results. As a consequence, we developed a reference-free and

224 alignment-based approach that adopts the clustering of reads based on their sequence

225 similarities (Figure 1b). Briefly, the reads from both samples were pooled and then clustered

226 by sequence homology (See Materials and methods). The consensus of each cluster was

227 determined and used as the reference in paired-end mapping. Subsequently, the consensus

228 sequences were used for a BLAST search against the reference genome for the closest

229 homologues. In this way, the BLAST search was able to map the clustered ALE-seq output to

230 reference genome annotated retrotransposons, which are most similar to the ALE-seq

231 recovered sequences. Applying this strategy, we identified a retroelement belonging to a

232 *Gypsy* family (*FIRE, Fruit-Induced RetroElement*) that produces significant amounts of

233 eclDNA at 52 DPA during fruit ripening (Figure 4a and b). We also determined the transcript

234 levels of the *FIRE* element in leaves and 52 DPA fruit samples. As shown in Figure 4c, fruit

235    RNA levels were enhanced twofold compared to leaves, where *FIRE* eclDNA was barely

236    detectable (Figure 4a). Finally, we found that the DNA methylation status of the *FIRE*

237    element was lower in fruits than leaves in all three sequence contexts (Figure 4d and f). In

238    contrast, the DNA methylation levels of sequences directly flanking *FIRE* were similar in

239    leaves and fruits (Figure 4e to g).

240

241    **Discussion**

242          Recently, a novel active retrotransposon was identified in rice by sequencing

243    extrachromosomal circular DNA (eccDNA) produced as a by-product of retrotransposition or

244    by nuclear recombination reactions of eclDNAs[23,24]. Although the method of eccDNA

245    sequencing has certain advantages over SIRT, such as increased sensitivity and the recovery

246    of sequences of the entire element, it also has certain limitations. For example, the method

247    requires relatively large amounts of starting material but still shows serious limits in

248    efficiency and indicative power for retrotransposition. The method did not detect the

249    eccDNA of *Tos19* in rice callus, where this transposon is known to move[23], however, direct

250    comparison of both methods on the same biological samples was not performed. More

251    importantly, eccDNAs may also be the result of genomic DNA recombination[25] and these

252    background products may be misleading when extrapolating to the transpositional potential

253    of a previously unknown element. In this respect, ALE-seq is a significantly improved tool

254    that largely overcomes the above-mentioned limitations of previous methods and requires

255    only 100 ng of plant DNA.

256          The heat-responsiveness of *Go-on*, the novel heat-activated *Copia* family

257    retrotransposon of rice detected using ALE-seq, seems to be conferred by *cis*-acting DNA

258    elements embedded in the LTR, which are similar to the heat-activated *Onsen*

259    retrotransposon in *Arabidopsis*[11,12]. Although heat stress can induce production of mRNA

260    and eclDNA of *Onsen*, its retrotransposition is tightly controlled by the small interfering RNA

261    pathway[9]. Given that real-time transposition of rice retrotransposons has only been

262    detected in epigenetic mutants[26,27] and triggered by tissue culture conditions causing vast

263    alterations in the epigenome[7],or as a result of interspecific hybridization[28], an altered

264    epigenomic status seems to be an important prerequisite for retrotransposition. In fact, we

265 failed to detect transposed copies of *Go-on* in the progeny of heat-stressed rice plants. Thus,

266 although *Go-on* produces eclDNAs after heat stress, it may be mobilized only at low

267 frequency in wild type rice due to epigenetic restriction of retrotransposition. Nevertheless,

268 on an evolutionary scale, the higher number of new insertions of *Go-on* in *indica* rice

269 populations grown at elevated temperatures might suggest its potential mobility.

270 Many retrotransposons are transcriptionally reactivated during specific

271 developmental stages or in particular cell types[29,30]. In tomato, fruit pericarp exhibits a

272 reduction in DNA methylation during ripening[16]. This is largely attributed to higher

273 transcription of the *DEMETER-LIKE2* DNA glycosylase gene[31]. Despite massive transcriptional

274 reactivation of retrotransposons in tomato fruits, it has been difficult to determine whether

275 further steps toward transposition also take place. Using ALE-seq, we identified eclDNA that

276 we annotated using a reference-free and alignment-based approach to a novel *FIRE* element.

277 *FIRE* has 164 copies in the reference tomato genome and in a conventional mapping-based

278 approach the ALE-seq reads of *FIRE* cross-mapped to multiple copies, making it difficult to

279 assign eclDNA levels to particular family members (Supplementary Figure 11). Therefore,

280 our strategy can be used in situations where sequence of the reference genome is

281 unavailable or the mapping of reads is hindered by the high complexity and multiplicity of

282 the retrotransposon population.

283 ALE-seq could also be applied to non-plant systems. For example, numerous studies

284 in various eukaryotes, including mammals, found that retrotransposons are transcriptionally

285 activated by certain diseases or at particular stages during embryo development[32,33]. It was

286 also suggested that retrotransposition might be an important component of disease

287 progression[34]. Given that the direct detection of retrotransposition is challenging, it would

288 be interesting to use ALE-seq to determine whether such temporal relaxations of epigenetic

289 transposon silencing also result in the production of the eclDNAs, as the direct precursor of

290 the chromosomal integration of a retrotransposon.

291

292

**Materials and methods**

Plant materials

Seeds of *Oryza sativa ssp. japonica cv. Nipponbare* and *Oryza sativa ssp. indica cv. IR64* were surface-sterilized in 20% bleach for 15 min, rinsed three times with sterile water and germinated on ½-MS media. Rice plants were grown in 10 h light / 14 h dark at 28°C and 26°C, respectively. For heat-stress experiments, 1-week-old rice plants were transferred to a growth chamber at 44°C and 28°C in light and dark, respectively. Rice callus was induced by the method used for rice transformation as previously described[35].

Tomato plants (*Solanum lycopersicum cv. M82*) were grown under standard greenhouse conditions (16 h supplemental lighting of 88 w/m$^2$ at 25°C and 8 h at 15°C). Tomato leaf tissue samples were taken from 2-month-old plants. Tomato fruit pericarp tissues were harvested at 52 days post anthesis (DPA).


Annotation of LTR retrotransposons

Functional *de novo* annotation of LTR retrotransposons for the genomes of TAIR10 (Arabidopsis), MSU7 (rice) and SL2.50 (tomato) was achieved by the *LTRpred* pipeline (https://github.com/HajkD/LTRpred) using the parameter configuration: minlenltr = 100, maxlenltr = 5000, mindistltr = 4000, maxdisltr = 30000, mintsd = 3, maxtsd = 20, vic = 80, overlaps = "no", xdrop = 7, motifmis = 1 , pbsradius = 60, pbsalilen = c(8,40), pbsoffset = c(0,10), quality.filter = TRUE, n.orf = 0. The plant-specific tRNAs used to screen for primer binding sites (PBS) were retrieved from GtRNAdb[36] and plantRNA[37] and combined in a custom fasta file. The hidden Markov model files for gag and pol protein conservation screening were retrieved from Pfam[38] using the protein domains RdRP_1 (PF00680), RdRP_2 (PF00978), RdRP_3 (PF00998), RdRP_4 (PF02123), RVT_1 (PF00078), RVT_2 (PF07727), Integrase DNA binding domain (PF00552), Integrase zinc binding domain (PF02022), Retrotrans_gag (PF03732), RNase H (PF00075) and Integrase core domain (PF00665). Computationally reproducible scripts for generating annotations can be found at http://github.com/HajkD/ALE.

321

## ALE-seq library preparation

323 Genomic DNA was extracted using a DNeasy Plant Mini Kit (Qiagen) following the

324 manufacturer's instruction. Genomic DNA (100 ng) was used for adapter ligation with 4 µl of

325 50 µM adapter DNA. After an overnight ligation reaction at 4°C, the adapter-ligated DNA

326 was purified by AMPure XP beads (Beckman Coulter) at a 1:0.5 ratio. *In vitro* transcription

327 reactions were performed using a MEGAscript RNAi kit (Thermo Fisher) with minor

328 modifications. Briefly, the reaction was carried out for 4 h at 37°C and the template DNA

329 was digested prior to RNA purification. Purified RNA (3 µg) was subjected to reverse

330 transcription (RT) using a Transcriptor First Strand cDNA Synthesis Kit (Roche). Transcriptor

331 First Strand cDNA Synthesis Kit was chosen because the RTase of the kit is thermostable.

332 This allowed the RT reaction at higher temperature (55°C) that reduces the RT-inhibiting

333 RNA secondary structure formation. The custom RT primers were added as indicated for

334 each experiment. After the RT reaction, 1 µl of RNase A/T1 (Thermo Fisher) was added to

335 digest non-templated RNA and the reaction mixture was incubated at 37°C for at least 30

336 min. Single-stranded first strand cDNA was PCR-amplified by 25 cycles using Illumina TruSeq

337 HT dual adapter primers and the PCR product was purified by AMPure XP beads (Beckman

338 Coulter) at a 1:1 ratio. After purification, the eluted DNA was quantified using a KAPA

339 Library Quantification Kit (KAPA Biosystems) and run on the MiSeq v3 2 X 300 bp platform in

340 the Department of Pathology of the University of Cambridge. Due to the nature of ALE-seq

341 that specifically amplifies ecDNAs, some ecDNA-free samples did not produce enough

342 library DNAs which, although suboptimal loading, were nevertheless sequenced. It is

343 advisable to spike in PCR-amplified retrotransposon DNA as described below. The

344 oligonucleotide sequences are provided in Supplementary Table 1.

345

## Preparation of full-length *Onsen* DNA

347 The full-length *Onsen* copy (AT1TE12295) was amplified using Phusion High-Fidelity DNA

348 polymerase (New England Biolabs). PCR products were run on 1% agarose gels. The full-

349 length fragment was then purified by QIAquick Gel Extraction (Qiagen) and its concentration

350     was measured using the Qubit Fluorometric Quantitation system (Thermo Fisher). Primers

351     used for amplification are listed in Supplementary Table 1.

352

353     <u>RT-qPCR analyses</u>

354     Samples were ground in liquid nitrogen using mortar and pestle. An RNeasy Plant Mini Kit

355     (Qiagen) was used to extract total RNA following the manufacturer's instructions. The

356     amount of extracted RNA was estimated using the Qubit Fluorometric Quantitation system

357     (Thermo Fisher). cDNAs were synthesized using a SuperScript VILO cDNA Synthesis Kit

358     (Invitrogen). Real-time quantitative PCR was performed in the LightCycler 480 system

359     (Roche) using primers listed in Supplementary Table 1. LightCycler 480 SYBR green I master

360     premix (Roche) was used to prepare the reaction mixture in a volume of 10 µl. The results

361     were analysed by the ΔΔCt method.

362

363     <u>RNA-seq library construction</u>

364     Total RNA was prepared as described above. An Illumina TruSeq Stranded mRNA Library

365     Prep kit (Illumina) was used according to the manufacturer's instructions. The resulting

366     library was run on an Illumina NextSeq 500 machine (Illumina) in the Sainsbury Laboratory

367     at the University of Cambridge.

368

369     <u>Analysis of next-generation sequencing data</u>

370     For RNA-seq data analysis, the adapter and the low-quality sequences were removed by

371     Trimmomatic software[39]. The cleaned reads were mapped to the MSU7 version of the rice

372     reference genome (http://rice.plantbiology.msu.edu) using TopHat2[40]. The resulting

373     mapping files were processed to the Cufflinks/Cuffquant/Cuffnorm pipeline[41] guided by the

374     annotation file which includes the MSU7 reference gene annotation

375     (http://rice.plantbiology.msu.edu/pub/data/Eukaryotic_Projects/o_sativa/annotation_dbs/

376     pseudomolecules/version_7.0/all.dir/) and our custom retrotransposon annotation.

377    Visualization of sequencing data was performed using an Integrative Genomics Viewer

378    (IGV)[42].

379    For ALE-seq data analysis, the adapter sequence was removed from the raw reads using

380    Trimmomatic software. For the mapping-based approach, paired-end reads were mapped to

381    the reference genomes (Arabidopsis, TAIR10; rice, MSU7; tomato, SL2.50) using Bowtie2[43]

382    with minor optimization. In most short-read sequencing platforms, it is often difficult to

383    assign the multi-mapped reads of TEs to precise genomic location. However, as MiSeq

384    outputs relatively longer reads, we presumed that ALE-seq reads have less ambiguity than

385    other sequencing platforms and set the parameters dealing with multi-mappers to default.

386    It is only the maximum fragment length option which is set to 500 by default that was

387    manipulated to 3000 (-X 3000). The numbers of reads mapped throughout each

388    retrotransposon were counted by the featureCounts tool of the SubRead package[44] using

389    the custom annotation file created by *LTRpred*. Since featureCounts recognizes multi-

390    mappers by SAM file's NH tag that bowtie2 does not generate, multi-mapped reads are

391    counted as one read aligned to a single genomic location, which reduces quantitation bias

392    that often happens to multi-mappers. IGV was used to visualize the sequencing data. For the

393    alignment-based approach, the forward and reverse reads were merged to yield the full-

394    length fragment sequences and converted to fasta files using the BBTools

395    (https://jgi.doe.gov/data-and-tools/bbtools/). The fasta files created for all the samples

396    were concatenated to get a master fasta file that is later inputted to CD-HIT software[45] to

397    cluster the reads by sequence similarity with the following options: -c 0.95, -ap 1, -g 1. CD-

398    HIT outputs a fasta file of representative reads for each cluster. The resulting fasta file was

399    used as reference for paired-end mapping of initial fastq files. The mapped reads were

400    counted with the featureCounts tool. Those clusters that significantly differed in the number

401    of mapped reads in different samples were further analysed for their identities using BLAST

402    search.

403    For Bisulfite sequencing analysis, raw sequenced reads derived from tomato fruits (52 DPA)

404    and leaves were downloaded from the public repository (SRP008329)[16] and re-analysed as

405    previously described[46], with minor modifications. Briefly, high-quality sequenced reads were

406    mapped with Bismark[47] on the cv. Heinz 1706 reference genome (https://solgenomics.net),

407    including a chloroplast sequence obtained from GenBank database (NC_007898.3) to

14

408     estimate the conversion rate. After methylation call and correction for unconverted

409     cytosines, the methylation proportions at each cytosine position with a coverage of at least

410     3 reads were used to generate a bedGraph file for each cytosine context, using the R

411     Bioconductor packages DMRCaller[48] and Rtracklayer[49]. The IGV browser was used to

412     visualize the methylation profiles.

413

414     <u>Detection of retrotransposon insertions</u>

415     The insertions of selected retrotransposons were detected from the genome resequencing

416     data of *japonica* and *indica* rice accessions downloaded from the 3,000 rice genome project

417     (PRJEB6180). The Transposon Insertion Finder (TIF) program[14] was used to identify the split

418     reads in the fastq files and detect newly integrated copies. We used MSU7

419     (http://rice.plantbiology.msu.edu) and ShuHui498 (http://www.mbkbase.org) for the

420     reference of *japonica* and *indica* rice, respectively. Only non-reference insertions were

421     considered and common insertions found in multiple accessions were counted as a single

422     retrotransposition event.

423

424     <u>Data availability</u>

425     The next generation sequencing data that support the findings of this study are available in

426     the Sequence Read Archive (SRA) repository with the identifier SRP155920.

427

428     <u>Code availability</u>

429     The custom scripts used in this study are available in http://github.com/HajkD/ALE.

430

**References**

432    1.     Lisch, D. How important are transposons for plant evolution? *Nat Rev Genet* **14,** 49–
433           61 (2012).

434    2.     Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory activities of transposable
435           elements: from conflicts to benefits. *Nat. Rev. Genet.* **18,** 71–86 (2017).

436    3.     Griffiths, J., Catoni, M., Iwasaki, M. & Paszkowski, J. Sequence-Independent
437           Identification of Active LTR Retrotransposons in Arabidopsis. *Mol. Plant* **11,** 508–511
438           (2017).

439    4.     Ma, J. & Bennetzen, J. L. Rapid recent growth and divergence of rice nuclear genomes.
440           *Proc. Natl. Acad. Sci. U. S. A.* **101,** 12404–12410 (2004).

441    5.     Sanchez, D. H., Gaubert, H., Drost, H., Zabet, N. R. & Paszkowski, J. High-frequency
442           recombination between members of an LTR retrotransposon family during
443           transposition bursts. *Nat. Commun.* **8,** 1–6 (2017).

444    6.     Picault, N. *et al.* Identification of an active LTR retrotransposon in rice. *Plant J.* **58,**
445           754–765 (2009).

446    7.     Sabot, F. *et al.* Transpositional landscape of the rice genome revealed by paired-end
447           mapping of high-throughput re-sequencing data. *Plant J.* **66,** 241–246 (2011).

448    8.     Mirouze, M. *et al.* Selective epigenetic control of retrotransposition in Arabidopsis.
449           *Nature* **461,** 427–430 (2009).

450    9.     Ito, H. *et al.* An siRNA pathway prevents transgenerational retrotransposition in
451           plants subjected to stress. *Nature* **472,** 115–119 (2011).

452    10.    Paszkowski, J. Controlled activation of retrotransposition for plant breeding. *Curr.*
453           *Opin. Biotechnol.* **32,** 200–206 (2015).

454    11.    Cavrak, V. V. *et al.* How a Retrotransposon Exploits the Plant's Heat Stress Response
455           for Its Activation. *PLoS Genet.* **10,** 1–12 (2014).

456    12.    Pietzenuk, B. *et al.* Recurrent evolution of heat-responsiveness in Brassicaceae COPIA

457          elements. *Genome Biol.* 1–15 (2016). doi:10.1186/s13059-016-1072-3

458   13.   The 3, 000 rice genomes project. The 3 , 000 rice genomes project. *Gigascience* **3,** 1–6

459          (2014).

460   14.   Nakagome, M. *et al.* Transposon Insertion Finder (TIF): a novel program for detection

461          of de novo transpositions of transposable elements. *BMC Bioinformatics* **15,** 1–9

462          (2014).

463   15.   Xiong, Z. Y. *et al.* Latitudinal Distribution and Differentiation of Rice Germplasm : Its

464          Implications in Breeding. *Crop Sci.* **51,** 1050–1058 (2011).

465   16.   Zhong, S. *et al.* Single-base resolution methylomes of tomato fruit development

466          reveal epigenome modifications associated with ripening. *Nat. Biotechnol.* **31,** 154–

467          159 (2013).

468   17.   Consortium, T. tomato genome. The tomato genome sequence provides insights into

469          fleshy fruit evolution. *Nature* **485,** 635–641 (2012).

470   18.   Eshed, Y. & Zamir, D. An Introgression Line Population of Lycopersicon pennellii in the

471          Cultivated Tomato Enables the Identification and Fine Mapping of Yield-Associated

472          QTL. *Genetics* **141,** 1147–1162 (1995).

473   19.   Eshed, Y. & Zamir, D. Less-Than-Additive Epistatic Interactions of Quantitative Trait

474          Loci in Tomato. *Genetics* **143,** 1807–1817 (1996).

475   20.   Quadrana, L. *et al.* The Arabidopsis thaliana mobilome and its impact at the species

476          level. *Elife* **5,** 1–25 (2016).

477   21.   Stuart, T. *et al.* Population scale mapping of transposable element diversity reveals

478          links to gene regulation and epigenomic variation. *Elife* **5,** 1–27 (2016).

479   22.   Wei, B. *et al.* Genome-wide characterization of non-reference transposons in crops

480          suggests non-random insertion. *BMC Genomics* **17,** 1–13 (2016).

481   23.   Lanciano, S. *et al. Sequencing the extrachromosomal circular mobilome reveals*

482          *retrotransposon activity in plants*. *PLOS Genetics* **13,** (2017).

483    24.    Møller, H. D. *et al.* Formation of Extrachromosomal Circular DNA from Long Terminal
484           Repeats of Retrotransposons in Saccharomyces cerevisiae. *G3 (Bethesda).* **6,** 453–462
485           (2015).

486    25.    Møller, H. D., Parsons, L., Jørgensen, T. S., Botstein, D. & Regenberg, B.
487           Extrachromosomal circular DNA is common in yeast. *Proc. Natl. Acad. Sci. U. S. A.* **112,**
488           3114–3122 (2015).

489    26.    Cheng, C. *et al.* Loss of function mutations in the rice chromomethylase OsCMT3a
490           cause a burst of transposition. *Plant J.* **83,** 1069–1081 (2015).

491    27.    Cui, X. *et al* Control of transposon activity by a histone H3K4 demethylase in rice.
492           *Proc. Natl. Acad. Sci. U. S. A.* **110,** 1953–1958 (2013).

493    28.    Wang, Z. H. *et al.* Genomewide Variation in an Introgression Line of Rice-Zizania
494           Revealed by Whole-Genome re-Sequencing. *PLoS One* **8,** 1–12 (2013).

495    29.    Li, H., Freeling, M. & Lisch, D. Epigenetic reprogramming during vegetative phase
496           change in maize. *Proc. Natl. Acad. Sci. U. S. A.* **107,** 22184–22189 (2010).

497    30.    Slotkin, R. K. *et al.* Epigenetic Reprogramming and Small RNA Silencing of
498           Transposable Elements in Pollen. *Cell* **136,** 461–472 (2009).

499    31.    Liu, R. *et al* A DEMETER-like DNA demethylase governs tomato fruit ripening. *Proc.*
500           *Natl. Acad. Sci.* **112,** 10804–10809 (2015).

501    32.    Goodier, J. L. Retrotransposition in tumors and brains. *Mobile DNA* **5,** 1–6 (2014).

502    33.    Baillie, J. K. *et al.* Somatic retrotransposition alters the genetic landscape of the
503           human brain. *Nature* **479,** 534–537 (2011).

504    34.    Mullins, C. S. & Linnebacher, M. Human endogenous retroviruses and cancer :
505           Causality and therapeutic possibilities. *World J. Gastroenterol.* **18,** 6027–6035 (2012).

506    35.    Cho, J. & Paszkowski, J. Regulation of rice root development by a retrotransposon
507           acting as a microRNA sponge. *Elife* **6,** 1–21 (2017).

508    36.    Chan, P. P. & Lowe, T. M. GtRNAdb 2 . 0 : an expanded database of transfer RNA

509           genes identified in complete and draft genomes. *Nucleic Acids Res.* **44,** 184–189

510           (2016).

511    37.   Daujat, M. *et al.* PlantRNA , a database for tRNAs of photosynthetic eukaryotes.

512           *Nucleic Acids Res.* **41,** 273–279 (2012).

513    38.   Finn, R. D. *et al.* Pfam : the protein families database. *Nucleic Acids Res.* **42,** 222–230

514           (2014).

515    39.   Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic : a flexible trimmer for Illumina

516           sequence data. *Bioinformatics* **30,** 2114–2120 (2014).

517    40.   Kim, D. *et al.* TopHat2 : accurate alignment of transcriptomes in the presence of

518           insertions , deletions and gene fusions. *Genome Biol.* **14,** 1–13 (2013).

519    41.   Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals

520           unannotated transcripts and isoform switching during cell differentiation. *Nat.*

521           *Biotechnol.* **28,** 516–520 (2010).

522    42.   Robinson, J. T. *et al.* Integrative Genomics Viewer. *Nat. Biotechnol.* **29,** 24–26 (2011).

523    43.   Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat.*

524           *Methods* **9,** 357–359 (2013).

525    44.   Liao, Y., Smyth, G. K. & Shi, W. The Subread aligner : fast , accurate and scalable read

526           mapping by seed-and-vote. *Nucleic Acids Res.* **41,** 1–17 (2013).

527    45.   Li, W. & Godzik, A. Cd-hit : a fast program for clustering and comparing large sets of

528           protein or nucleotide sequences. *Bioinformatics* **22,** 1658–1659 (2006).

529    46.   Catoni, M. *et al.* DNA sequence properties that predict susceptibility to epiallelic

530           switching. *EMBO* **36,** 617–628 (2017).

531    47.   Krueger, F. & Andrews, S. R. Bismark : a flexible aligner and methylation caller for

532           Bisulfite-Seq applications. *Bioinformatics* **27,** 1571–1572 (2011).

533    48.   Catoni, M., Tsang, J. M. F., Greco, A. P. & Zabet, N. R. DMRcaller : a versatile R /

534           Bioconductor package for detection and visualization of differentially methylated

535    regions in CpG and non-CpG contexts. *Nucleic Acids Res.* 1–11 (2018).

536  49.  Lawrence, M., Gentleman, R. & Carey, V. rtracklayer : an R package for interfacing

537    with genome browsers. *Bioinformatics* **25,** 1841–1842 (2009).

538

**Author Contributions**

J.C. and J.P. conceived the research. J.C., M.B., M.C., H.-G.D., A.B. and M.O. performed experiment. J.C., M.B., M.C. and H.-G.D. analysed data. J.C. and J.P. wrote and revised the manuscript.

**Competing Interests**

The authors declare that no competing interests exist.

**Corresponding authors**

Correspondence to Jungnam Cho (jungnam.cho@slcu.cam.ac.uk) and Jerzy Paszkowski (jerzy.paszkowski@slcu.cam.ac.uk).

**Figure Legends**

557     Figure 1. Detection of eclDNA by ALE-seq

558     **a**, The workflow of ALE-seq. The colour code is indicated in a box. **b**, Analysis pipeline of ALE-

559     seq results. The sequenced reads can be mapped to the reference genome or aligned to

560     each other to obtain a cluster consensus. **c** and **d**, Genome-wide plots of rice ALE-seq results

561     from leaf (**c**) and callus (**d**). The levels are shown as number of reads mapped to each

562     retrotransposon. Dots represent annotated retrotransposons; those corresponding to *Tos17*

563     and *Tos19* are indicated. **e** and **f**, Read coverage plots mapped to *Tos17* (**e**) and *Tos19* (**f**).

564     The black bars represent retrotransposons and white arrowheads indicate LTRs.

565

566     Figure 2. Sensitivity and specificity of eclDNA detection by ALE-seq

567     **a-d**, ALE-seq reconstruction experiment with varying amounts of PCR-amplified *Onsen* DNA

568     added to rice callus DNA. Genome browser image with the read coverage (**a** and **c**) and

569     quantitated read counts (**b** and **d**) for *Onsen* (**a** and **b**) and *Tos17* (**c** and **d**) loci. The amounts

570     of *Onsen* DNA added were 1 ng, 100 pg, 10 pg, 1 pg or 100 fg; 100 ng of rice callus DNA was

571     used. Note that read coverage values are $log_{10}$-converted in **a**. For **b** and **d**, values are shown

572     as $log_{10}$-converted counts per million sequenced reads. **e** and **f**, Read coverage plots for the

573     ALE-seq of rice callus using different RT primers. *Tos17* and *RIRE2* transposons are depicted

574     below the plots as in Figure 1.

575

576     Figure 3. Identification of a novel heat-activated retrotransposon in rice

577     **a** and **b**, Genome-wide plots of rice ALE-seq results as in Figure 1. Control (**a**) and heat-

578     stressed (**b**) rice plants were used. One-week-old seedlings were subjected to heat stress

579     (44°C) for 3 days. Met-iCAT PBS primer was used in RT. The levels are shown as the number

580     of reads mapped to retroelements. Three *Go-on* copies are indicated in **b**. **c**, Read coverage

581     plot for *Go-on3*. **d**, RNA-seq data showing *Go-on3* and a neighbouring gene. RNA-seq data

582     were generated using the same plant materials as in **a** and **b**. The experiment was repeated

583     independently two times with similar results. **e-g**, Cumulative plots for the number of non-

584 reference insertions of *Go-on* (**e**), *Tos17* (**f**), and *Tos19* (**g**) in the genomes of 388 *japonica*

585 and *indica* rice accessions. The statistical difference was determined by iterating random

586 selection of 200 accessions out of 388 and performing the two-tailed Wilcoxon test. ** *P*

587 $=2.2e^{-16}$.

588

589 Figure 4. Identification of a tomato retrotransposon activated in fruit pericarp

590 **a**, Read coverage plot for the *FIRE* retrotransposon identified in tomato fruit pericarp by

591 ALE-seq. Met-iCAT PBS primer was used in RT. **b** and **c**, The DNA (**b**) and RNA (**c**) levels of

592 *FIRE* in leaves and fruits determined by qPCR. The levels are means of two biological

593 replicates. Normalization was done against *SlGAPDH* (Solyc03g111010) and *SlCAC*

594 (Solyc08g006960) for DNA and RNA analyses, respectively. **d**, Genome browser image for

595 the DNA methylation levels at *FIRE* element in leaves and fruits of tomato. The levels are

596 shown as percent methylation of each cytosine. **e**-**g**, Violin plots for DNA methylation levels

597 at the upstream (**e**), *FIRE* (**f**) and downstream (**g**) regions. Only cytosines supported by at

598 least three reads in both samples were considered. In *FIRE* locus, for example, 4,032 out of

599 4,078 cytosines in both strands were analysed. The upstream and downstream regions are

600 immediate flanking sequences taken for the same length as *FIRE* of 9.362 kb. P-values were

601 determined by a two-sided Fisher's t-test using 558 CG and 717 CHG sites at *FIRE* locus.

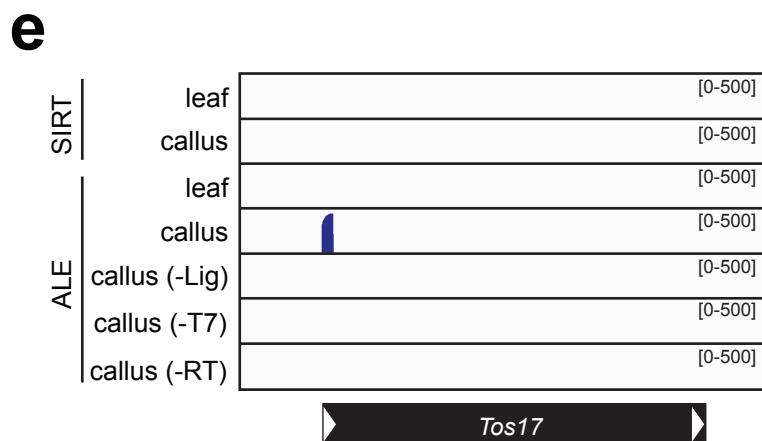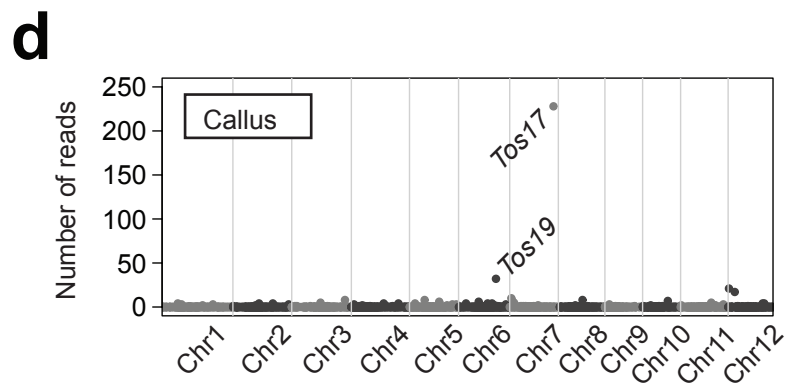602 Other samples with insignificant statistical difference are not shown for the p-values.
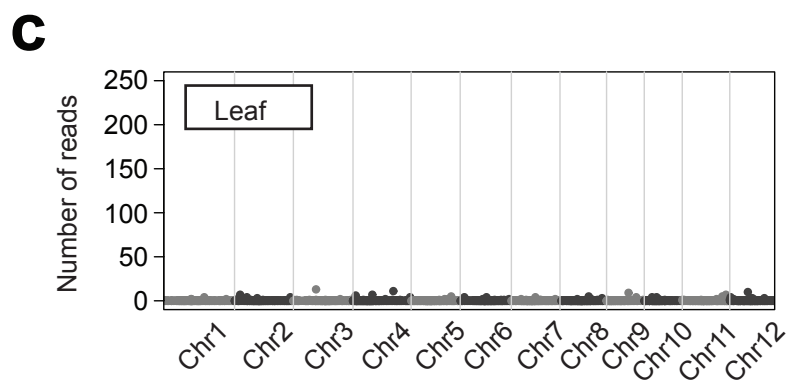
603

**Figure 1.**

# Figure 2.

**a**

[0-7]

[0-7]

[0-7]

[0-7]

[0-7]

*Onsen* DNA

*Onsen* (AT1TE12295)

**b**

Counts per million (Log10)

*Onsen* DNA (ng, Log10)

**c**

[0-500]

[0-500]

[0-500]

[0-500]

[0-500]

*Onsen* DNA

*Tos17* (Chr7:26694787-26698920)

**d**

Counts per million (Log10)

*Onsen* DNA (ng, Log10)

**e**

Met-iCAT [0-500]

Arg-CCT [0-500]

Met-iCAT + Arg-CCT [0-500]

*Tos17* (Chr7:26694787-26698920)

**f**

Met-iCAT [0-100]

Arg-CCT [0-100]

Met-iCAT + Arg-CCT [0-100]

*RIRE2* (Chr8:20710296-20720336)

**Figure 3.**

Figure 4.

1

2 **Supplementary Figure 1.** PBS sequences of LTR retrotransposons in *Arabidopsis*, rice and

3 tomato

4 The frequency of tRNAs used for targeting PBS. LTR retrotransposons were annotated by

5 *LTRpred* (http://github.com/HajkD/ALE) and selected for young elements by filtering LTR

6 similarities higher than 95%. The total numbers of retrotransposons analysed in each

7 species are shown below the plots.

8

**a**



**b**



9

10    **Supplementary Figure 2.** SIRT results from leaves and calli of rice

11    **a** and **b,** Genome-wide plots for SIRT performed in leaves (**a**) and in calli (**b**) of rice.

12

**Supplementary Figure 3.** ALE-seq detection of eclDNAs of *Arabidopsis* retrotransposons

Genome-wide plots (**a**, **b**, **d** and **f**) and read coverage plots (**c**, **e** and **g**) for ALE-seq profiles of *Arabidopsis* Col-0 wt (**a**), heat-stressed Col-0 (**b** and **c**), *met1-1* (**d** and **e**), and epi12 (**f** and **g**).

16

5' end of LTR          3' end of LTR

*Arabidopsis thaliana*

*Oryza sativa*

*Solanum lycopersicum*

17

18    **Supplementary Figure 4.** Conservation of end sequences of LTR

19    The conserved sequences of 5' and 3' ends of LTR. The first and last five nucleotides of LTRs

20    are displayed. The images were generated by the WebLogo tool

21    (http://weblogo.berkeley.edu/logo.cgi).

22

23

**Supplementary Figure 5.** ALE-seq detection of eclDNAs of rice retrotransposons using

multiplexed PBS primers

Genome-wide plot for ALE-seq profiles of rice callus using pooled PBS primers of Met-iCAT

and Arg-CCT.

**a**

| | | LTR similarity |
|---|---|---|
| *Go-on1* (Chr4:26160202 -26164923) | Δ7bp | 99.18% |
| *Go-on2* (Chr4:31588773 -31593491) | T>C | 100% |
| *Go-on3* (Chr9:11858139 -11862873) | 1039 1525 1731 1777 2109 3145 3879 4123 4533 4596 — GAG INT RTase RNaseH | 100% |

100%   a                                                            b  c

**b**

```
Go-on1  123  TGGTATCAGAGCCAATCGGCTGGTGGCTGGCGACGG-------CTAAACCCTAGCCTCGCCGGAG
Go-on2  123  TGGTATCAGAGCCAATCGGCTGGCGGCTGGCGACGGGCGACGGCTAAACCCTAGCCTCGCCGGAG
Go-on3  123  TGGTATCAGAGCCAATCGGCTGGTGGCTGGCGACGGGCGACGGCTAAACCCTAGCCTCGCCGGAG
clone1  123  TGGTATCAGAGCCAATCGGCTGGTGGCTGGCGACGGGCGACGGCTAAACCCTAGCCTCGCCGGAG
clone2  123  TGGTATCAGAGCCAATCGGCTGGTGGCTGGCGACGGGCGACGGCTAAACCCTAGCCTCGCCGGAG
clone3  123  TGGTATCAGAGCCAATCGGCTGGTGGCTGGCGACGGGCGACGGCTAAACCCTAGCCTCGCCGGAG
clone4  123  TGGTATCAGAGCCAATCGGCTGGTGGCTGGCGACGGGCGACGGCTAAACCCTAGCCTCGCCGGAG
clone5  123  TGGTATCAGAGCCAATCGGCTGGTGGCTGGCGACGGGCGACGGCTAAACCCTAGCCTCGCCGGAG
```

**c**

ACACGACGCTCTTCCGATCT TGTTGAGTTATGTATGTGTT

29

30   **Supplementary Figure 6.** *Go-on* retrotransposon family

31   **a**, Schematic structure of *Go-on* retrotransposons. The genomic coordinates and LTR

32   similarities of each copy are shown at the left and right, respectively. Red boxes, ORFs;

33   green boxes, regions encoding protein domains; blue boxes, PBS; white arrowheads, LTRs.

34   Note that the sequences of the upstream LTRs through the PBS are identical in all three

35   copies. The sequence variation specific for each element is indicated. Protein domains were

36   predicted by NCBI BLASTP tool (https://blast.ncbi.nlm.nih.gov/Blast.cgi). Nucleotide

37   positions indicating the start and end of ORF and protein domains are provided. Primers

38   used for sequencing and qPCR analyses are shown as arrows. **b**, Multiple sequence

39   alignment of the genomic sequences of three *Go-on* copies and the sequenced ALE clones.

40   ALE-seq was performed using the RT primer specific to *Go-on3* indicated as "a" in **a**. The

41    resulting single-stranded first strand cDNA was PCR-amplified, cloned to the pGEM T-easy

42    vector, and sequenced. Multiple sequence alignment was performed by ClustalW

43    (http://www.genome.jp/tools-bin/clustalw) and visualized by boxshade tools

44    (https://www.ch.embnet.org/software/BOX_form.html). **c**, Sequencing of the ALE-seq

45    product of *Go-on3* showing the junction region of the adapter and LTR. Sequences in red

46    and black are the adapter and *Go-on* LTR, respectively.

47

**a**  Relative DNA levels vs Days of heat stress (0, 1, 3, 5, 5+3r). Control (black), Heat (red). n.s., p=0.0062, p=0.0008, p=0.0013

**b**  Relative RNA levels vs Days of heat stress (0, 1, 3, 5, 5+3r). Control (black), Heat (red). n.s., p=0.0005, p=0.0005, p=0.0135

**c**

```
5'-TGTTGAGTTATGTATGTGTTGGCCCATGAGGCCCATATACTAC
    |         |         |         |         |
    1        10        20        30        40

TCATATGTACATGTATATAGCAGAGTTAGAGAAATGAAAAAAGTAG
          |         |         |         |
         50        60        70        80

TGAAGCTTCTAGAGAAAAATTCCCAAAACTTCATGGTATCAGAGC-3'
  |         |         |         |         |    |
 90        100       110       120       130  134
```

48

49  **Supplementary Figure 7.** Heat stress-triggered transcriptional activation of *Go-on*

50  **a** and **b**, The relative levels of DNA (**a**) and RNA (**b**) of *Go-on3* determined by qPCR. Heat

51  treatment (44°C) was applied to 1-week-old rice seedlings for the periods indicated; +3r

52  means 3 days of recovery in normal growth conditions after heat stress. The levels are

53  means ± sd of three biological replicates. For DNA analysis, Day 0 levels are set to 3,

54  reflecting three genomic copies of *Go-on* in *japonica* rice. Normalization was done against

55  *eEF1α*. P-values were calculated by a two-tailed Student's t-test; n.s., not significant. **c**, The

56  sequence of the left LTR and PBS of *Go-on3*. The sequence in red is the heat-related HSFC1-

57  binding sequence motif predicted by PlantPan 2.0 tool

58  (http://plantpan2.itps.ncku.edu.tw/index.html) with statistically significant enrichment of

59  *P*= 4.28e$^{-10}$ as determined by Fisher's exact test. The enrichment of the sequence motif was

60    calculated by comparing the ten most similar sequences of *Go-on* found in rice genome with

61    1,000 random genomic loci of 150bp. The PBS is shown in blue.

62

63

**Supplementary Figure 8.** Heat stress-triggered activation of *Go-on* in *indica* rice

**a** and **b**, The qPCR analyses for DNA (**a**) and RNA (**b**) levels of *Go-on* in *indica* rice. The levels are means ± sd of three biological replications. The levels of control sample are set to 2 (**a**) reflecting 2 genomic copies of *Go-on* in *indica* rice. P-values are determined by two-sided Student's t-test.

69

## a



## b



70

71  **Supplementary Figure 9.** Comparison of mRNA and ecIDNA levels

72  **a**, Scatter plot for log2-fold changes (FCs) in RNA-seq and ALE-seq profiles in the control and

73  heat-stressed rice plants used in Figure 3. FCs were calculated by dividing heat samples

74  values by control samples values of CPM (counts per million reads) and FPKM (fragments

75  per kb per million reads) for ALE-seq and RNA-seq data, respectively. Each dot represents an

76    individual retroelement and the dashed lines mark log2-FC one. Th retrotransposon in red

77    has log2-FC higher than one in both ALE-seq and RNA-seq. **b**, Read coverage plot for a

78    selected retrotransposon showing evidence of transcriptional activation upon heat stress

79    not followed by synthesis of eclDNAs.

80

81

**Supplementary Figure 10.** Retrotransposon insertions in *japonica* and *indica* rice

**a-c**, Density plots for number of non-reference insertions in randomly selected 200 accessions out of 388 iterated by 1,000 times.

85

86

**Supplementary Figure 11.** ALE-seq profile of tomato leaves and fruits

**a** and **b**, Genome-wide plots for ALE-seq profiles performed in tomato leaves (**a**) and fruits 52 DPA (**b**). Each dot represents an individual retrotransposon.

90

91    **Supplementary Table 1.** Sequences of oligonucleotides used in this study

92    T7 promoter sequence is underlined and in bold is partial Illumina adapter sequence.

| Primer | Sequence (5' → 3') |
|---|---|
| ALE adapter top strand | AGAGAG<u>TAATACGACTCACTATAGGG</u>**ACACGACGCTCTTCCGATCT** |
| ALE adapter bottom strand | **AGATCGGAAGAGCGTCGTGT**<u>CCCTATAGTGAGTCGTATTA</u>CTCTCT |
| ALE RT Met-iCAT-R | AGACGTGTGCTCTTCCGATCTGCTCTGATACCA |
| ALE RT Arg-CCT-R | AGACGTGTGCTCTTCCGATCTCCTGGCGCGCCA |
| ONSEN full length-F | TGTTGAAAGTTAAACTTGATTTTG |
| ONSEN full length-R | TGTTAGAGTAAAATTCTTTTAG |
| Go-on-F (b of Figure S5) | GGCAGAATACAGGGCAATGTC |
| Go-on-R (c of Figure S5) | GCCGACTTATTGTCACACCAC |
| Go-on RT-R (a of Figure S5) | TCTCTGCACGCCTCGACAAG |
| eEF1α-F | GCACGCTCTTCTTGCTTTCACTCT |
| eEF1α-R | AAAGGTCACCACCATACCAGGCTT |
| FIRE RT-F | GAGTTGGCTACGTATCGTTTGC |
| FIRE RT-R | AGCCTCCACAAATTCATCCCAT |
| FIRE copy number-F | GGTGTTCTCGTTGTGGTAAGT |
| FIRE copy number-R | TAAGGTGACACTCCCTCATAGT |
| SlCAC-F | CCTCCGTTGTGATGTAACTGG |
| SlCAC-R | ATTGGTGGAAAGTAACATCATCG |
| SlGAPDH-F | ATGCTCCCATGTTTGTTGTGGGTG |
| SlGAPDH-R | TTAGCCAAAGGTGCAAGGCAGTTC |

93

94    **Supplementary Table 2.** Summary of ALE-seq libraries

95    Numbers of reads sequenced and mapped are summarised. Both unique-mappers and

96    multi-mappers are considered. "% mapped to LTRs" refers to all reads mapped throughout

97    retrotransposon.

| Samples | Reads sequenced | % mapped to genome | % mapped to LTRs | % not mapped to LTRs | Accession number |
|---|---|---|---|---|---|
| Arabidopsis Col-0 | 56,057 | 93.20 | 17.77 | 75.43 | SAMN09748167 |
| Arabidopsis heat-stressed | 45,554 | 90.33 | 15.85 | 74.48 | SAMN09748168 |
| Arabidopsis *met1-1* | 58,029 | 94.79 | 16.03 | 78.76 | SAMN09748169 |
| Arabidopsis epi12 | 45,545 | 96.19 | 31.18 | 65.01 | SAMN09748170 |
| Rice leaf | 27,063 | 97.54 | 12.33 | 85.21 | SAMN09748171 |
| Rice callus | 25,183 | 95.67 | 13.58 | 82.09 | SAMN09748172 |
| Rice callus -Lig | 37,610 | 83.48 | 11.75 | 71.73 | SAMN09748173 |
| Rice callus -T7 | 870 | 2.82 | 0.12 | 2.70 | SAMN09748174 |
| Rice callus -RT | 516 | 1.26 | 0 | 1.26 | SAMN09748175 |
| Rice callus pooled PBS | 22,939 | 90.06 | 14.16 | 75.90 | SAMN09748176 |
| Rice non-stressed | 31,819 | 90.39 | 12.29 | 78.10 | SAMN09748177 |
| Rice heat-stressed | 31,525 | 97.63 | 13.54 | 84.09 | SAMN09748178 |
| Tomato leaf | 46,421 | 96.65 | 13.24 | 83.41 | SAMN09748179 |
| Tomato fruit 52 DPA | 73,067 | 96.97 | 28.53 | 68.44 | SAMN09748180 |

98

99   **Supplementary Table 3.** Non-reference insertions of *Go-on*.

100  Neo-insertions of *Go-on* detected by TIF. The positions are provided as coordinates of target
101  site duplication.

| Chromosome | Start | End | Accession | Category |
|---|---|---|---|---|
| Chr4 | 31968290 | 31968294 | ERS467756 | *Indica* |
| Chr4 | 31968290 | 31968294 | ERS467757 | *Indica* |
| Chr1 | 43029985 | 43029989 | ERS467791 | *Indica* |
| Chr4 | 31968290 | 31968294 | ERS467794 | *Indica* |
| Chr4 | 31968290 | 31968294 | ERS467830 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467831 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS467831 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS467843 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467872 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS467876 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467877 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467878 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS467878 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS467880 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS467883 | *indica* |
| Chr11 | 1485949 | 1485953 | ERS467910 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467915 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS467924 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS467925 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS467926 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS467927 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS467934 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467938 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS467943 | *indica* |
| Chr8 | 3790507 | 3790511 | ERS467943 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS467944 | *indica* |
| Chr4 | 31960603 | 31960607 | ERS467952 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS467952 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS467959 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS467960 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467961 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS467962 | *indica* |
| Chr4 | 13737528 | 13737532 | ERS467962 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS467962 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467966 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467969 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS467969 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467979 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS467980 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467986 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS467995 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS467995 | *indica* |
| Chr12 | 16846383 | 16846387 | ERS467996 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS467996 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS467996 | *indica* |
| Chr5 | 287377 | 287381 | ERS467998 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS467999 | *indica* |
| Chr8 | 3790507 | 3790511 | ERS468001 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS468004 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468004 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS468006 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468006 | *indica* |
| Chr4 | 31960603 | 31960607 | ERS468008 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS468011 | *indica* |

| Chr8 | 5518300 | 5518306 | ERS468011 | *indica* |
|---|---|---|---|---|
| Chr1 | 43029985 | 43029989 | ERS468014 | *indica* |
| Chr12 | 16846383 | 16846387 | ERS468016 | *indica* |
| Chr4 | 31960603 | 31960607 | ERS468018 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468023 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468025 | *indica* |
| Chr8 | 5518299 | 5518306 | ERS468028 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468029 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468042 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS468048 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468049 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468050 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468052 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468053 | *indica* |
| Chr4 | 13737528 | 13737532 | ERS468055 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468055 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS468059 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS468060 | *indica* |
| Chr6 | 31319589 | 31319593 | ERS468065 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468066 | *indica* |
| Chr4 | 13737528 | 13737532 | ERS468068 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468071 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS468072 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS468073 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468074 | *indica* |
| Chr4 | 31960603 | 31960607 | ERS468075 | *indica* |
| Chr12 | 16846383 | 16846387 | ERS468077 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS468078 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468084 | *indica* |
| Chr11 | 29783248 | 29783252 | ERS468086 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS468087 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS468088 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468088 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468089 | *indica* |
| Chr12 | 22672020 | 22672024 | ERS468095 | *indica* |
| Chr4 | 31960603 | 31960607 | ERS468101 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468102 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468104 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468106 | *indica* |
| Chr12 | 16846383 | 16846387 | ERS468111 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468112 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468115 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468121 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468126 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468131 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468133 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS468134 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468136 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468138 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468139 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468142 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468154 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468157 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468160 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468161 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468163 | *indica* |
| Chr1 | 20049396 | 20049400 | ERS468166 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468169 | *indica* |
| Chr5 | 19524953 | 19524957 | ERS468170 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468174 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468184 | *indica* |

| Chr8 | 5518300 | 5518306 | ERS468186 | *indica* |
|------|---------|---------|-----------|----------|
| Chr4 | 31968290 | 31968294 | ERS468187 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468187 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468191 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468192 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468193 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468195 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468202 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS468202 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468204 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468205 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468207 | *indica* |
| Chr8 | 5518300 | 5518306 | ERS468209 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS468210 | *indica* |
| Chr11 | 30168035 | 30168039 | ERS468212 | *indica* |
| Chr5 | 287377 | 287381 | ERS468212 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468215 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468222 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468230 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468232 | *indica* |
| Chr4 | 31968290 | 31968294 | ERS468234 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468237 | *indica* |
| Chr4 | 13737528 | 13737532 | ERS468240 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS468249 | *indica* |
| Chr1 | 43029985 | 43029989 | ERS468250 | *indica* |
| Chr3 | 431859 | 431863 | ERS468252 | *indica* |
| Chr7 | 29937784 | 29937788 | ERS468252 | *indica* |
| Chr11 | 29783248 | 29783252 | ERS468255 | *indica* |
| Chr1 | 19257271 | 19257275 | ERS467801 | *japonica* |
| Chr5 | 23638163 | 23638167 | ERS467889 | *japonica* |
| Chr11 | 1514174 | 1514178 | ERS467893 | *japonica* |
| Chr8 | 5635769 | 5635774 | ERS467904 | *japonica* |
| Chr11 | 1514174 | 1514178 | ERS468026 | *japonica* |
| Chr5 | 258622 | 258626 | ERS468308 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468310 | *japonica* |
| Chr1 | 41968356 | 41968360 | ERS468380 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468383 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468384 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468387 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468402 | *japonica* |
| Chr6 | 24954457 | 24954461 | ERS468442 | *japonica* |
| Chr6 | 22413483 | 22413487 | ERS468446 | *japonica* |
| Chr7 | 29379081 | 29379085 | ERS468449 | *japonica* |
| Chr8 | 5635769 | 5635774 | ERS468449 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468456 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468458 | *japonica* |
| Chr8 | 5635769 | 5635774 | ERS468595 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468596 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468604 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468613 | *japonica* |
| Chr8 | 5635768 | 5635775 | ERS468617 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468620 | *japonica* |
| Chr6 | 22413483 | 22413487 | ERS468649 | *japonica* |
| Chr5 | 258622 | 258626 | ERS468684 | *japonica* |
| Chr7 | 29379081 | 29379085 | ERS468704 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS468705 | *japonica* |
| Chr5 | 258622 | 258626 | ERS468721 | *japonica* |
| Chr5 | 23638163 | 23638167 | ERS468734 | *japonica* |
| Chr1 | 41968356 | 41968360 | ERS468902 | *japonica* |
| Chr1 | 41968356 | 41968360 | ERS468917 | *japonica* |
| Chr8 | 5635769 | 5635774 | ERS468993 | *japonica* |

| Chr2 | 1659944 | 1659948 | ERS469049 | *japonica* |
|------|---------|---------|-----------|-----------|
| Chr8 | 5635768 | 5635774 | ERS469069 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469132 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469177 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469199 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469215 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469302 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469307 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469556 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469602 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469604 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469605 | *japonica* |
| Chr8 | 5635768 | 5635775 | ERS469637 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469650 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469668 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469669 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469689 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469694 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469696 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469699 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469746 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469758 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469845 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469880 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469978 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS469985 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS470129 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS470132 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS470188 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS470344 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS470439 | *japonica* |
| Chr8 | 5635768 | 5635774 | ERS470516 | *japonica* |