1  **A large-scale genome and transcriptome sequencing analysis reveals the mutation**

2  **landscapes induced by high-activity adenine base editors in plants**

3  Shaofang Li[1, #, *], Lang Liu[1, 2, 3, #], Wenxian Sun[3], Xueping Zhou[1, 4], Huanbin Zhou[1, 2,*]

4

5  [1]State Key Laboratory for Biology of Plant Diseases and Insect Pests, Institute of Plant

6  Protection, Chinese Academy of Agricultural Sciences, Beijing 100193, China

7  [2]Scientific Observing and Experimental Station of Crop Pests in Guilin, Ministry of

8  Agriculture and Rural Affairs, Guilin 541399, China

9  [3]Department of Plant Pathology, China Agricultural University, Beijing, 100193, China

10  [4]State Key Laboratory of Rice Biology, Institute of Biotechnology, Zhejiang University,

11  Hangzhou, Zhejiang, China

12

13  [#]These authors contributed equally to this work

14  [*]To whom correspondence should be addressed: shaofangli2021@hotmail.com,

15  hbzhou@ippcaas.cn

16

17  **Running title:** ABE-induced DNA and RNA mutations

18  **Data deposition:**

19  Reviewer links to deposited GEO data with token mzafkeomtjobvuv at the following

20  link: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE185497

21

**Abstract**

**Background**: The high-activity adenine base editors (ABEs), engineered with the recently-developed tRNA adenosine deaminases (TadA8e and TadA9), show robust base editing activity but raise concerns about off-target effects.

**Results**: In this study, we performed a comprehensive evaluation of ABE8e- and ABE9-induced DNA and RNA mutations in *Oryza sativa*. Whole-genome sequencing analysis of plants transformed with four ABEs, including SpCas9n-TadA8e, SpCas9n-TadA9, SpCas9n-NG-TadA8e, and SpCas9n-NG-TadA9, revealed that ABEs harboring TadA9 lead to a higher number of off-target A-to-G (A>G) single-nucleotide variants (SNVs), and that those harboring the CRISPR/SpCas9n-NG lead to a higher total number of off-target SNVs in the rice genome. An analysis of the T-DNAs carrying the ABEs indicated that the on-target mutations could be introduced before and/or after T-DNA integration into plant genomes, with more off-target A>G SNVs forming after the ABEs had integrated into the plant genome. Furthermore, we detected off-target A>G RNA mutations in plants with high expression of ABEs but not in plants with low expression of ABEs. The off-target A>G RNA mutations tended to cluster, while off-target A>G DNA mutations rarely clustered.

**Conclusion**: Our findings that Cas proteins, TadA variants, temporal expression of ABEs, and expression levels of ABEs contribute to ABE specificity in rice provide insight into the specificity of ABEs and suggest alternative ways to increase ABE specificity besides engineering TadA variants.


**Keywords:** adenine base editor, TadA variants, single nucleotide variant, T-DNA insertion, off-target, *Oryza sativa* L.

## Background

48

49  Single-nucleotide variants (SNVs), a universal feature of plant, animal, and human

50  genomes, have been widely identified in association with agronomic traits and human

51  diseases[1-3]. Various clustered regularly interspaced short palindromic repeats

52  (CRISPR)/CRISPR-associated protein (Cas)-mediated base editing tools (e.g., ABEs

53  and cytosine base editors), which efficiently produce desired point mutations in

54  genomic DNA without causing double-stranded DNA breaks [4], have been used

55  widely in laboratory research, crop and animal breeding, as well as human gene therapy

56  [5-7]. Since the mutation of G•C base pairs to A•T base pairs is the primary form of *de*

57  *novo* mutations [8], ABEs that catalyze the conversion of A•T base pairs to G•C base

58  pairs have great potential to correct human pathogenic point mutations [9]. However,

59  potential DNA and RNA off-target mutations remain a serious concern and threaten to

60  limit the application of ABEs.

61  The pioneer ABE7s, which are composed of a tRNA adenosine deaminase

62  (TadA7.10) and CRISPR/Cas systems, perform remarkably clean and efficient A•T to

63  G•C conversions in the genomes of a variety of species, including human, mouse, and

64  rice, without inducing obvious genome-wide off-target DNA mutations [10-14].

65  However, the editing efficiency of ABE7s varies in a locus-dependent manner [11, 13].

66  Subsequently, high-activity ABEs, such as those containing TadA8.17, TadA8.20,

67  TadA8e, and TadA9, have been developed and used in different organisms [15-17], but

68  a whole-genome assessment of the off-target DNA mutations induced by TadA8e and

69  TadA9 has not yet been investigated.

70  The tRNA adenosine deaminase TadA, a key component of ABEs, induces site-

71  specific inosine formation on RNAs [18]. Recently, it was reported that TadAs, ABE7s,

72  and ABE8es induced a significantly higher number or higher mutation ratio of RNA A-

73  to-G (A>G) SNVs when compared to Cas proteins or GFP [9, 19, 20], and that ABE8.17

74  and ABE8.20 induced very low levels of adenosine deamination in mRNAs if ABEs

75  were delivered as messenger RNAs in mammalian cells [17]. Thus, several labs have

3

76    developed improved TadA variants with reduced RNA activity [19, 21]. However,

77    RNA A>G mutations induced by ABEs are complicated due to the large genomes in

78    the heterogeneous mammalian cells as well as the conversion of adenosines into

79    inosines mediated by endogenous adenosine deaminase RNA specific (ADAR) family.

80    In addition, ABE-induced RNA mutations have never been reported in plant yet.

81       The relatively small genome (~0.4 Gb) of self-pollinated rice and the absence of

82    endogenous ADAR family make rice an ideal model organism to examine the DNA

83    and RNA specificity of gene editing tools. Here, we investigated the off-target DNA

84    and RNA mutations induced by ABE8es and ABE9s in rice through whole-genome

85    sequencing (WGS) and transcriptome sequencing.

86

87    **Results**

88    **ABEs induced sgRNA-independent heterozygous DNA mutations**

89    To assess the off-target activity effects of high-activity tRNA adenosine deaminases

90    (TadA8e and TadA9), we chose four ABEs that are composed of different variants of

91    tRNA adenosine deaminase and CRISPR/Cas systems with different PAM

92    compatibility: rBE46b (SpCas9n-TadA8e), rBE49b (SpCas9n-TadA9), rBE50

93    (SpCas9n-NG-TadA8e), and rBE53 (SpCas9n-NG-TadA9) (Fig. 1a). For each ABE,

94    three constructs with one or two sgRNAs and one construct without sgRNAs were

95    generated. After *Agrobacterium tumefaciens*-mediated transformation, we obtained

96    three independent transgenic plants for each construct except 46bM and 49bM, which

97    had three plants from two independent transformation events (Additional file 2: Table

98    S1). We examined the on-target mutations in the 36 plants carrying the ABE plus

99    sgRNA(s) through Sanger sequencing and identified the desired mutations in 35 plants

100   (Additional file 1: Figure S1 and S2). To assess the effects of tissue culture and

101   *Agrobacterium* infection, three independently regenerated plants subjected to tissue

102   culture and six independently regenerated plants subjected to *Agrobacterium* infection

103   without vectors were selected for WGS. We also sequenced 10 wild-type rice variety

4

104 Kitaake plants to filter out background mutations (Fig. 1b). To ensure high confidence

105 in base calling, we sequenced all 71 plants at an average coverage of 41× (Additional

106 file 2: Table S2). SNVs in each plant were identified using three independent variant-

107 calling software systems: GATK, Strelka2, and Lofreq [14, 22-24]. Small insertions or

108 deletions (indels) were called independently by GATK and Strelka2. SNVs identified

109 by all three methods and indels identified by two methods were kept for later analysis

110 (Fig. 1b). All the SNVs and indels identified in the 10 Kitaake plants were considered

111 background mutations and removed from the analysis. The sgRNA-guided on-target

112 and off-target loci were located by Criflash [25] (Additional file 2: Table S3).

113 Consistent with Sanger sequencing results, A>G on-target mutations were observed in

114 35 out of 36 plants (Additional file 1: Figure S3) and removed in the following off-

115 target analysis. No mutations were detected at 33 predicted sgRNA-guided off-target

116 sites with 2-3 nt mismatches. For plants that had undergone tissue culture (Control

117 group 1: C1) and *Agrobacterium* infection (Control group 2: C2), we identified around

118 200-400 SNVs and around 250-350 indels from each plant (Additional file 1: Figure

119 S4a, b). For plants carrying an ABE, we identified around 200-800 SNVs and 200-500

120 indels (Additional file 1: Figure S4a, b). Six types of SNVs were identified separately

121 in control plants and in those carrying ABEs. We discovered that A>G/T>C SNVs

122 constituted a higher proportion in plants with ABEs (Additional file 1: Figure S4c, d).

123 For simplicity, we referred to the number of A>G SNVs as the total number of A>G

124 and T>C SNVs, and we referred to the percentage of A>G SNVs as the percentage of

125 the total number of A>G and T>C SNVs versus the total number of all six types of

126 SNVs throughout the manuscript. Consistently, the number and percentage of A>G

127 SNVs in plants with ABEs were higher when compared to both control groups,

128 indicating that ABEs induce the genomic mutations of the A•T base pairs to G•C base

129 pairs (Additional file 1: Figure S4c-f).

130  A few homozygous SNVs and indels were detected in all sequenced plants

131 (Additional file 1: Figure S5a, b). We counted the number of plants with the same

132  mutation sites and found that the homozygous mutations tended to be present in more

133  than one plant, while the heterozygous mutations tended to be present in a single plant

134  (Additional file 1: Figure S5c, d). These homozygous mutations could be the remaining

135  background mutations or mutations induced by tissue culture, *Agrobacterium* infection,

136  or ABEs. The induced mutations in the two alleles are two independent events

137  following binomial distribution, so the probability of the homozygous mutations is $p^2$,

138  the probability of being wild type (WT) is $(1-p)^2$, and the probability of the

139  heterozygous mutations is $2 * p * (1-p)$, assuming that the induced mutation ratio for

140  each allele was p and the ratio of the WT allele was 1-p. A binomial test for all loci of

141  homozygous SNVs or indels revealed that these loci did not follow a binomial

142  distribution (Additional file 1: Figure S5e and Additional file 2: Tables S4 and S5),

143  indicating that these homozygous mutations remain background SNVs and indels.

144  These data suggests that ABEs induce sgRNA-independent heterozygous DNA

145  mutations.

146

147  **Genome-wide analysis of ABE-induced single-nucleotide mutations**

148  After background homozygous mutations were removed, we recalculated the number

149  of SNVs and indels in the plants (Additional file 1: Figure S6a, b and Additional file 2:

150  Table S6). We did not observe any significant differences of SNVs or indels induced

151  by tissue culture or *Agrobacterium* infection (Fig. 2a, b). Therefore, we used the plants

152  that had been infected with *Agrobacterium* as the control group in the following

153  analysis. Consistent with the finding that ABEs do not cause double-strand DNA breaks,

154  plants with ABEs did not show a higher number of indels (Fig. 2a). The number of total

155  SNVs and A>G SNVs was significantly higher in plants harboring rBE50 and rBE53

156  than in the control groups, while the number and the percentage of A>G SNVs were

157  significantly higher in plants harboring rBE49b and rBE53 than in plants in control

158  groups (Fig. 2b). We did not observe a significantly higher number of SNVs or a higher

159  percentage of A>G SNVs in plants harboring rBE46b (Fig. 2b).

160    We next examined whether Cas proteins or TadA variants play distinct roles in

161    inducing off-target DNA mutations by comparing plants harboring rBE46b with those

162    harboring rBE49b as well as rBE50 to rBE53 to characterize TadA8e and TadA9, and

163    compared plants with rBE46b to rBE50 and rBE49b to rBE53 to characterize the role

164    of SpCas9n and SpCas9n-NG in off-target effects. Although there was no significant

165    difference between TadA8e and TadA9 when the total number of SNVs was considered,

166    plants harboring TadA9 had a higher number and a higher percentage of A>G SNVs

167    (Fig. 2c), indicating that TadA9-based ABEs lead to a higher number of A>G SNVs.

168    Plants harboring SpCas9n-NG had a higher number of SNVs as well as a higher number

169    of A>G SNVs, but not a higher percentage of A>G SNVs (Fig. 2d), indicating that

170    SpCas9n-NG-based ABEs lead to a higher number of SNVs.

171    We classified all SNVs into six types and calculated the percentage of each type of

172    SNV versus the total number of SNVs. We observed a higher percentage of C>A/G>T

173    SNVs in plants harboring TadA8e (Additional file 1: Figure S7). We further mapped

174    all SNVs and A>G SNVs to different genic and intergenic regions and calculated the

175    ratio of SNVs in given regions versus in the whole genome. As a result, the number of

176    A>G SNVs and the total number of SNVs were higher at all genic and intergenic

177    regions in plants for all four types of ABEs, while A>G SNVs were enriched in genic

178    regions and depleted in intergenic regions (Fig. 2e and Additional file 1: Figure S8). In

179    addition, we mapped total SNVs as well as A>G SNVs to the 12 rice chromosomes and

180    established that they were distributed throughout the rice genome (Additional file 1:

181    Figure S9).

182

183    **T-DNA insertion influences the single-nucleotide mutations**

184    We detected genome-wide off-target SNVs induced by tissue culture from three

185    plants, those induced by *Agrobacterium* infection without vectors in six plants, and

186    those in 48 plants transformed by *Agrobacterium* infection with ABEs. We compared

187    SNVs from the individual plants to those identified in all other plants to examine the

188    overlapping SNVs. Among 1,596 comparisons, we found none of the common SNVs

189    in 1,567 comparisons, and 1-7 overlapping SNVs in 27 comparisons (Additional file 2:

190    Table S7), which indicates the randomness of off-target DNA mutations induced by

191    tissue culture, *Agrobacterium* infection, and ABEs. In addition, we detected 147

192    overlapping SNVs in the comparison of lines 46bM_s2 and 46bM_s3, and 85

193    overlapping SNVs in the comparison of lines 49bM_s2 and 49bM_s3. Notably,

194    46bM_s2 and 46bM_s3 as well as 49bM_s2 and 49bM_s3 are plants regenerated from

195    the same resistant calli (Additional file 2: Table S1). The T-DNA insertion sites in the

196    genomes of three plants transformed with 46bM and three plants transformed with

197    49bM were located by T-LOC (Li *et al.* in preparation). We determined that lines

198    46bM_s2 and 46bM_s3 were derived from the same T-DNA integration event, whereas

199    line 46bM_s1 was from a different T-DNA integration event (Fig. 3a). Similarly, lines

200    49bM_s2 and 49bM_s3, but not line 49bM_s1, harbored the same T-DNA insertion site

201    (Fig. 3a). Surprisingly, plants carrying the same T-DNA insertion event did not always

202    have the same sgRNA-guided on-target mutations (Additional file 1: Figure S3). To

203    validate this phenomenon, we also sequenced line 49bAG_s4, which was regenerated

204    from the same resistant callus as line 49bAG_s3. We established that lines 49bAG_s3

205    and 49bAG_s4 had the same T-DNA insertion site, which differed from that of lines

206    49bAG_s1 and 49bAG_s2 (Fig. 3a), and that lines 49bAG_s3 and 49bAG_s4 had

207    different on-target editing events (Additional file 1: Figure S10a). We further

208    characterized the off-target SNVs in these plants and found that different plants with

209    the same T-DNA insertion had both unique SNVs and common SNVs (Fig. 3b and

210    Additional file 1: Figure S11b). We defined three sequential stages in *Agrobacterium*-

211    transformed callus: stage 1, the period after the T-DNA plasmid has entered the callus

212    cell and before it has integrated into the genome; stage 2, the period after the T-DNA

213    has integrated into the genome and before the callus cell has divided; and stage 3, the

214    period after the callus cell has divided. Since the off-target mutation happens randomly,

215    the unique SNVs should occur at stage 3, while common SNVs should occur at both

8

216     stage 1 and stage 2. A higher percentage of A>G SNVs was observed among the unique

217     SNVs when compared to common SNVs in plants transformed with 49bM and 49bAG

218     (Fig. 3b), indicating that ABEs integrated into the reference genome are more prone to

219     cause A>G SNVs.

220        We next examined the integrity of T-DNA regions containing both a complete left

221     border (LB) and right border (RB) and identified four plants with a partial T-DNA

222     insertion characterized by the missing TadA8e, TadA9, or SpCas9n-NG fragment

223     (Additional file 1: Figure S11a). However, desired on-target mutations were detected

224     in three out of four plants (Additional file 1: Figure S3), suggesting that sgRNA-

225     dependent on-target A>G editing could occur before T-DNA integration into the rice

226     genome. We further checked the off-target SNVs between plants with or without

227     complete T-DNA insertion and found that plants with a complete T-DNA insertion had

228     a higher number of total SNVs, a higher number of A>G SNVs, and a higher percentage

229     of A>G SNVs when compared to those with partial T-DNA insertion (Fig. 3c).

230

231     **ABEs induce transcriptome-wide A>G RNA mutations**

232     To examine whether ABEs induce RNA off-target mutations, we profiled the

233     transcriptomes of three plants subjected to *Agrobacterium* infection without vectors,

234     three transformed plants carrying functional SpCas9 only, three plants carrying rBE46b

235     without sgRNAs, nine plants carrying rBE46b with one or two sgRNAs, three plants

236     carrying rBE49b without sgRNAs, and nine plants carrying rBE49b with one or two

237     sgRNAs (Fig. 1b). SNVs were called independently by GATK, Strelka2, and Lofreq

238     from each transcriptome and the corresponding genome data. We kept SNVs called by

239     all three methods in transcriptome data but not in genome sequencing data. In addition,

240     SNVs detected from plants in the *Agrobacterium* infection group were removed as

241     background mutations. Overall, the number of SNVs, the number of A>G SNVs, and

242     the percentage of A>G SNVs were not significantly higher in plants harboring rBE46b

243     and rBE49b than in those harboring SpCas9 nuclease (Additional file 1: Figure S12a

244    and Additional file 2: Table S9); however, A>G SNVs constituted a higher proportion

245    in plants harboring rBE46b and rBE49b than in plants harboring SpCas9 nuclease only

246    (Additional file 1: Figure S12b). When SNVs were counted separately for each plant,

247    we found that transcriptomes R49AG_s2 and R49AG_s3 had more than 100 A>G

248    SNVs and that A>G SNVs were barely detected in plants harboring SpCas9 only (Fig.

249    4a). In contrast to the randomness of DNA off-target SNVs, the ratio of ABE-induced

250    RNA off-target SNVs in transcriptomes R49AG_s2 and R49AG_s3 A>G correlated

251    with each other (Fig. 4b), indicating that ABEs might have preferred RNA editing

252    sequence content. As expected, we identified a conserved YAN-enriched (Y = T, C and

253    N = A, T, C, G) motif at ABE-edited RNA loci (Fig. 4c). We combined the SNV loci

254    detected in all transcriptomes as ABE-targeting RNA loci, computed the A>G editing

255    ratio in each transcriptome with sufficient read coverage (read number higher than 10),

256    and performed a Wilcoxon test that compared the A>G editing ratio of each plant

257    containing ABEs versus the A>G editing ratio in three plants that contained SpCas9.

258    Although the number of ABE-targeted RNA loci with sufficient reads were comparable

259    in all sequenced transcriptomes (Additional file 1: Figure S12c), transcriptomes from

260    eight plants harboring rBE46b and rBE49b, including R46AG_s1, R46AG_s3,

261    R46GG_s1, R49AG_s2, R49AG_s3, R49bg_s1, R49bg_s2, and R49bg_s3, had

262    significantly higher A>G editing ratios (Fig. 4d). Since these eight plants also had

263    detectable numbers of A>G RNA SNVs, we concluded that ABEs (rBE46b and rBE49b)

264    induced RNA editing in these eight plants but not in the remaining 16 plants. We

265    examined the ABE-induced DNA off-target mutations, but found no differences

266    between the plants with RNA mutations and those without RNA mutations (Additional

267    file 1: Figure S12f). When the reads per million (RPM) value of ABEs

268    (SpCas9n/SpCas9n-NG and TadA8e/TadA9) was calculated, we found that the

269    transcript levels of ABEs were significantly higher in the eight plants with RNA

270    mutations than in the 16 plants without RNA mutations (Fig. 4e). Given the high

271    concordance between ABE transcript abundance and the A>G editing ratio, we

272   wondered whether RNA A>G editing would cease after the T-DNA insertion

273   segregated out in the next generation. Two transgenic and two transgene-free plants

274   were selected in the $T_1$ population of line 49AG_s2 and subjected to transcriptome

275   analysis. As expected, A>G RNA editing was eliminated in the two $T_1$ plants that lacked

276   the ABE transgene but remained active in the two plants with transgenes (Fig. 4f and

277   Additional file 1: Figure S12d).

278

279   **ABEs induce clustered off-target editing**

280   Given that ABEs lead to multiple A>G editing events at the sgRNA-dependent on-

281   target window, we wondered whether they function the same way at the A>G off-target

282   editing loci. We examined the A>G mutations located within the 5' and 3' 30-bp

283   flanking region of every ABE-induced A>G off-target locus in the transcriptome data.

284   After counting A>G SNVs for which the A>G conversion rate was higher than 0.05

285   and also counting A sites in cases where the read coverage was higher than 10, we

286   determined the ratios of A>G SNVs at every flanking position. In eight transcriptomes

287   with RNA off-target editing, A>G SNVs were consistently distributed in the flanking

288   regions (Fig 5a). We refer to SNVs with flanking SNVs as clustered SNVs. By contrast,

289   no flanking A>G editing occurred in plants lacking RNA off-target SNVs or in plants

290   harboring SpCas9 nuclease (Fig. 5b and Additional file 1: Figure S13 and S14). Of

291   these A>G off-target RNA editing events, there were SNVs with a high number of

292   flanking A>G mutations and high occurrence in many transcriptomes, and there were

293   also SNVs with a low number of flanking A>G mutations and occurrence in a few

294   transcriptomes (Fig. 5a and Additional file 1: Figure S15).

295       We performed similar studies on DNA off-target SNVs but did not observe general

296   patterns of flanking A>G editing. However, we did identify 25 loci with more than one

297   A>G SNV from 12 plants (Additional file 2: Table S10); some loci contained 5-10 A>G

298   SNVs, and others contained 2-3 A>G SNVs (Fig. 5d and Additional file 1: Figure S16).

299   Overall, 45% of these SNVs were located in the genic region, which is higher than the

300  30% observed for all A>G SNVs in the genic region, consistent with the tendency of

301  off-target A>G SNVs to occur in the genic region (Fig. 5e). We classified these 12

302  plants into group 1, and the remaining 36 plants carrying ABEs into group 2. The

303  number of SNVs and A>G SNVs and the percentage of A>G SNVs were significantly

304  higher for plants in group 1 compared to plants in group 2 (Fig. 5f).

305

306  **Discussion**

307  ABE8s and ABE9s have been developed by several groups to overcome the limitation

308  of ABE7s [15-17]. Their robust editing efficiency raised another question: How is the

309  specificity of those high-activity ABEs engineered with TadA8e and TadA9

310  deaminases? Compared to mouse and human genomes (each ~3 Gb), the rice genome

311  (~0.4 Gb) is small, making WGS of individuals more feasible. In addition, rice is self-

312  pollinating, circumventing the challenges of population heterogeneity of human cells,

313  and lacks innate A-to-I RNA editing, facilitating analyses of ABE-induced RNA editing.

314  Therefore, we performed a comprehensive evaluation of ABE8- and ABE9-induced

315  genetic mutations through WGS and transcriptome sequencing in rice.

316      Cas proteins and TadA variants play different roles in ABE-induced DNA off-

317  target mutations: ABEs harboring SpCas9n-NG, an engineered SpCas9 protein

318  recognizing a flexible protospacer adjacent motif (PAM) [26, 27], result in a higher

319  number of total SNVs; those harboring TadA9, a TadA variant with robust activity [16],

320  lead to a higher number of specific A>G SNVs. Plants transformed with the ABE

321  rBE46b (SpCas9n-TadA8e) did not have more SNVs or a higher percentage of A>G

322  SNVs than plants subjected to *Agrobacterium* infection, suggesting that selection of

323  SpCas9n and TadA8e eliminates most sgRNA-independent DNA mutations induced by

324  ABEs. Given that no sgRNA-dependent off-target mutations were observed, we

325  conclude that optimization of sgRNA design is an efficient way of eliminating sgRNA-

326  dependent off-target mutations.

12

327    Using deeply sequenced genomes and transcriptomes, we systematically studied

328    ABE-induced RNA mutations. ABEs induce RNA A>G mutations in one-third of

329    plants with high ABE expression but do not induce mutations in two-thirds of plants

330    with low ABE expression. When ABEs segregated out, RNA mutations diminished. In

331    addition, T-DNA integration analysis suggested that stable ABEs induce more off-

332    target SNVs than those whose T-DNA has not been integrated into the genome.

333    Together, these data highlight the importance of controlling the expression of ABEs in

334    future applications, such as using inducible or photoactivatable transcription systems,

335    ribonucleoprotein-based delivery in clinic gene therapy [28, 29], and transgene-free

336    gene-edited plants in crop breeding.

337    Without the noise from A-to-I mutations mediated by ADAR proteins, we were able

338    to obtain a clean set of ABE-induced RNA mutations and discovered that ABEs induced

339    clustered A>G mutations, which provided useful information for defining and

340    characterizing true ABE RNA targets. Furthermore, given the existence of common and

341    unique mutations in plants regenerated from the same callus, we provide robust

342    experimental evidence that plants with different on-target editing could be derived from

343    the same T-DNA insertion event with a shared set of off-target SNVs. Therefore, we

344    highly recommend using two independent transgenic lines from separated calli (with

345    two different T-DNA insertion sites and two sets of non-overlapping SNVs) in gene

346    function studies.

347

348    **Conclusions**

349    The properties of the small genome, self-pollination and the absence of ADAR proteins

350    make rice a model organism to employ large-scale sequencing approaches to evaluate

351    ABEs' off-target activity. The pioneering comprehensive analysis of ABE-induced

352    DNA and RNA mutations using whole-genome and transcriptome sequencing in rice

353    sheds light on defining and characterizing ABEs' specificity. The discovery that Cas

354    proteins, TadA variants, transient expression, and the expression level of ABEs

13

355     contribute to ABEs' specificity in rice points out alternative ways improving ABEs'

356     specificity including combinatorial optimization of Cas/deaminase (SpCas9n-TadA8e),

357     temporal control of ABEs' expression besides the traditional protein engineering of

358     deaminases.

359

360     **Materials and Methods**

361     **Plasmid construction**

362     In this study, five rice (*Oryza sativa*) genomic loci (*OsACC*, *OsGS1*, *OsMPK13*,

363     *OsGSK3*, and *OsGSK4*) and four rice genomic loci (*OsACC*, *OsGS1*, *OsMPK13*, and

364     *OsTms9*) were targeted by rBE46b and rBE49b, respectively. Three genes (*OsSERK2*,

365     *OsDEP2*, and *OsGSK4*) were targeted by both rBE50 and rBE53. Plant IDs and their

366     corresponding information are described in Additional file 2: Table S1. The rBE46b,

367     rBE49b, rBE50, and rBE53 expression plasmids were constructed as previously

368     reported [16]. The empty entry vector without any spacer was cloned into pUbi:rBE46b,

369     pUbi:rBE49b, pUbi:rBE50, and pUbi:rBE53 using Gateway technology to yield ABEs

370     without sgRNAs (Additional file 2: Table S1).

371

372     *Agrobacterium*-**mediated rice transformation and plant growth**

373     The genome editing constructs were individually introduced into the *Agrobacterium*

374     *tumefaciens* strain EHA105 via the freeze-thaw transformation method, and 2-week-

375     old calli derived from immature seeds of the Geng rice variety Kitaake were infected

376     by each *Agrobacterium* strain. After 4 weeks of culture on MSD medium supplemented

377     with 50 mg/L hygromycin (Roche, Germany), the resistant callus lines were transferred

378     onto RM plates to generate transgenic rice seedlings. All information on target gene

379     mutations of each seedling examined in this study is given in Additional file 2: Table

380     S1.

381        As controls, seedlings were regenerated from calli infected with *Agrobacterium*

382     harboring *rBE* genes only and the empty EH105 strain. To eliminate WGS artifacts

383    caused by *Agrobacterium* infection, plants were also obtained directly from rice tissue

384    that had not been co-cultured with *Agrobacterium* cells. All plants were grown in the

385    greenhouse under a 16-h-light/8-h-dark photoperiod, 28/25°C temperature cycle, and

386    75% humidity.

387

388    **DNA and RNA extractions**

389    Genomic DNA of 4-week-old rice plants was extracted using the CTAB method (Li et

390    al., 2016). Approximately 200 mg of fresh rice leaves was collected in a 2-ml centrifuge

391    tube containing disposable metal balls. After being quickly frozen in liquid nitrogen,

392    samples were ground to a fine powder using a tissue grinding apparatus (Jingxin, China).

393    Following chloroform extraction, isopropanol precipitation, and 70% EtOH washing,

394    genomic DNAs were eluted with 50 μl of double-distilled water supplemented with 1

395    μl of 10 U/μL RNase I (Thermo Fisher Scientific, USA) and stored at −80°C for later

396    experiments.

397    RNA was extracted with TRIzol reagent (Takara, Japan) according to the

398    manufacturer's instructions. Briefly, 100 mg of fresh rice leaves was sampled, quickly

399    frozen in liquid nitrogen, and ground to a powder with a tissue grinding apparatus. Then,

400    1 ml of TRIzol reagent was added to the sample followed by chloroform and

401    isopropanol treatment. Finally, RNA pellets were dissolved in 50 μl of RNase-free

402    water (0.1% DEPC-treated) and stored at −80°C for later experiments.

403

404    **On-target mutation detection**

405    The on-target genomic regions were amplified using Phanta Max Super-Fidelity DNA

406    Polymerase (Vazyme, China) and locus-specific primers (Additional file 2: Table S1)

407    with genomic DNAs used as the template. PCR amplicons were subjected to Sanger

408    sequencing, and Bioedit software was used for sequence data analysis.

409

410    **Whole-genome analysis of genetic mutations**

15

411    RNA-free genomic DNAs (0.2 µg) from each sample were used to construct the DNA

412    libraries using a NEBNext Ultra DNA Library Prep Kit for Illumina (NEB, USA)

413    following the manufacturer's instructions. DNA libraries were sequenced on the

414    Illumina platform in the 150-nt paired-end mode with an average coverage depth of 40×

415    (Additional file 2: Table S2).

416         The clean reads were mapped to the Kitaake genome V3 from Phytozome

417    (https://data.jgi.doe.gov/refine-download/phytozome) via BWA [30] and sorted using

418    samtools (v1.9) [31]. The Genome Analysis Toolkit (GATK v4.2) was used to mark

419    duplicated reads and recalibrate base qualities [24]. To identify high-quality genetic

420    changes at the genomic scale, we applied three independent germline variant-calling

421    methods: GATK, LoFreq [22], and Strelka2 [22]. We documented SNVs identified by

422    all three methods and indels identified by GATK and Strelka. All genetic changes

423    identified by the three methods in the 10 Kitaake plants were combined and used as

424    background mutations. The genetic mutation ratios were calculated using an in-house

425    R program and 'AC' value from GATK's results. Both background mutations and

426    homozygous mutations were removed from the SNVs as well as indels. The IGV

427    browser was used to demonstrate sgRNA-directed on-target mutations [32]. Then, the

428    on-target mutations were removed for off-target analysis. sgRNA-dependent off-target

429    mutations were discovered using Crisflash [25], and the genetic on-target mutations

430    were    assessed    using    the    IGV    browser.    A    gene    annotation    file

431    (OsativaKitaake_499_v3.1.gene_exons.gtf) from the Phytozome website was used to

432    define different genomic regions, such as gene regions, exon regions, and intergenic

433    regions. The ggpubr, ggbio, and VennDiagram R libraries were used to draw the graphs.

434

435    **Analysis of T-DNA insertion sites and ABE transcripts**

436    The clean reads were mapped to T-DNA sequences using BWA and sorted using

437    samtools. The T-DNA insertion sites were located through T-LOC (Li *et al*. in

438    preparation). The coverage of T-DNAs between the left border (LB) and right border

439    (RB) was assessed using the R library ShortRead. The expression of ABEs was

440    quantified as the average raw read number of Cas proteins and TadA variants

441    normalized by the total read number in millions.

442

443    **Analysis of ABE-induced RNA mutations**

444    DNA-free RNAs (0.2 μg) were used to construct the RNA-seq libraries using a NEB

445    Next Ultra RNA Library Prep Kit for Illumina (NEB, USA) following the

446    manufacturer's instructions. RNA-seq libraries were sequenced on the Illumina

447    platform in the 150-nt paired-end mode (Additional file 2: Table S8).

448    The clean reads were mapped to the Kitaake V3 genome and annotation from

449    Phytozome via STAR aligner with a maximum of eight mismatches per paired-end read

450    [33]. GATK was used to mark duplicate reads and split reads that contained Ns in their

451    cigar string and to recalibrate base qualities. SNVs were called by GATK, LoFreq, and

452    Strelka2 for each transcriptome dataset and corresponding genome dataset. The SNVs

453    identified by three methods in the transcriptome data but not in the genome data were

454    kept for later analysis. All the genetic changes identified by the three methods in three

455    *Agrobacterium*-infected plants were combined and used as background mutations and

456    were removed from the SNVs identified in plants transformed with SpCas9, rBE46b,

457    and rBE49b. The A>G mutation ratios of off-target RNA loci were calculated through

458    in-house Python programs. The 30- and 3-bp flanking sequences of the off-target RNA

459    SNVs were extracted from the Kitaake reference genome and subjected to motif

460    prediction using WebLogo3 (http://weblogo. threeplusone.com/) [34].

461

462    **Calculation of flanking A>G mutations in genome and transcriptome data**

463    We combined all A>G off-target SNVs obtained from plants with RNA off-target

464    activities. For each A>G SNV, we calculated the number of reads with nucleotide A, T,

465    G, and C separately in the 5′ and 3′ 30-bp region with a read coverage larger than 10.

466    The genetic change ratio was calculated as the number of Gs divided by the total number

17

467 of As and Gs if the reference is A. The genetic change ratio was calculated as the

468 number of Cs divided by the total number of Cs and Ts if the reference is T. Positions

469 with an A>G mutation ratio of higher than 0.05 were used as the numerator, while

470 positions of A/T with a read coverage larger than 10 were used as the denominator.

471 Similarly, we combined all A>G off-target SNVs obtained from plants through WGS

472 and calculated the percentage of A>G mutations at the 5′ and 3′ 30-bp flanking regions.

473

474 **Parameters of boxplots used in this study**

475 The horizontal line in the box represents the median value, and the bottom and top of

476 the box are the lower (Q1) and upper quartiles (Q3), respectively. The upper whisker is

477 min(max($x$), Q3 + 1.5 * IQR), and the lower whisker is max(min($x$), Q1 - 1.5 * IQR).

478 IQR (interquartile range) = Q3 - Q1. Black dots located outsides the whiskers are

479 outliers.

480

481 **Acknowledgements**

482 We thank Sujie Zhang and Yongjie Kuang for assistance with RNA manipulation.

483

484 **Author' contributions**

485 S.L., W.S., X.Z., and H.Z. designed and guided the research. S.L. performed

486 bioinformatic analysis and L.L. performed experiments. S.L. and H.Z. wrote the

487 manuscript. All authors read and approved the final manuscript.

488

489 **Funding**

490 This work was supported by grants from the National Natural Science Foundation of

491 China (31871948) and the Central Public-interest Scientific Institution Basal Research

492 Fund (Y2020PT26) to H.Z.

493

494 **Ethics approval and consent to participate**

495    Not applicable.

496

**Competing Interests statement**

498    The authors declare that they have no competing financial interests.

499

**Availability of data and materials**

501    All data in the study has been included in the manuscript and additional files. All

502    sequencing genome and transcriptome data have been deposited in the NCBI database

503    with the accession number GSE185497.

504

**Author details**

506    [1]State Key Laboratory for Biology of Plant Diseases and Insect Pests, Institute of Plant

507    Protection, Chinese Academy of Agricultural Sciences, Beijing 100193, China

508    [2]Scientific Observing and Experimental Station of Crop Pests in Guilin, Ministry of

509    Agriculture and Rural Affairs, Guilin 541399, China

510    [3]Department of Plant Pathology, China Agricultural University, Beijing, 100193, China

511    [4]State Key Laboratory of Rice Biology, Institute of Biotechnology, Zhejiang University,

512    Hangzhou, Zhejiang, China

513

**References**

515    1.    Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, Gu B, Hart

516          J, Hoffman D, Hoover J, et al: **ClinVar: public archive of interpretations of**

517          **clinically relevant variants.** *Nucleic Acids Res* 2016, **44:**D862-868.

518    2.    Henikoff S, Comai L: **Single-nucleotide mutations for plant functional**

519          **genomics.** *Annu Rev Plant Biol* 2003, **54:**375-401.

520    3.    Wang W, Mauleon R, Hu Z, Chebotarov D, Tai S, Wu Z, Li M, Zheng T,

521          Fuentes RR, Zhang F, et al: **Genomic variation in 3,010 diverse accessions of**

522          **Asian cultivated rice.** *Nature* 2018, **557:**43-49.

523    4.    Gaudelli NM, Komor AC, Rees HA, Packer MS, Badran AH, Bryson DI, Liu

524          DR: **Programmable base editing of A·T to G·C in genomic DNA without**

525          **DNA cleavage.** *Nature* 2017, **551:**464-471.

526    5.    Liu L, Kuang Y, Yan F, Li S, Ren B, Gosavi G, Spetz C, Li X, Wang X, Zhou

527          X, Zhou H: **Developing a novel artificial rice germplasm for dinitroaniline**

528          **herbicide resistance by base editing of *OsTubA2*.** *Plant Biotechnol J* 2021,

529          **19:**5-7.

530    6.    Kuang Y, Li S, Ren B, Yan F, Spetz C, Li X, Zhou X, Zhou H: **Base-editing-**

531          **mediated artificial evolution of *OsALS1 in planta* to develop novel**

532          **herbicide-tolerant rice germplasms.** *Mol Plant* 2020, **13:**565-572.

533    7.    Ryu SM, Koo T, Kim K, Lim K, Baek G, Kim ST, Kim HS, Kim DE, Lee H,

534          Chung E, Kim JS: **Adenine base editing in mouse embryos and an adult**

535          **mouse model of Duchenne muscular dystrophy.** *Nat Biotechnol* 2018,

536          **36:**536-539.

537    8.    Acuna-Hidalgo R, Veltman JA, Hoischen A: **New insights into the generation**

538          **and role of *de novo* mutations in health and disease.** *Genome Biol* 2016,

539          **17:**241.

540    9.    Rees HA, Wilson C, Doman JL, Liu DR: **Analysis and minimization of**

541          **cellular RNA editing by DNA adenine base editors.** *Sci Adv* 2019,

542          **5:**eaax5717.

543    10.   Rees HA, Liu DR: **Base editing: precision chemistry on the genome and**

544          **transcriptome of living cells.** *Nat Rev Genet* 2018, **19:**770-788.

545    11.   Yan F, Kuang Y, Ren B, Wang J, Zhang D, Lin H, Yang B, Zhou X, Zhou H:

546          **Highly efficient A·T to G·C base editing by Cas9n-guided tRNA adenosine**

547          **deaminase in rice.** *Mol Plant* 2018, **11:**631-634.

548    12.   Zeng Y, Li J, Li G, Huang S, Yu W, Zhang Y, Chen D, Chen J, Liu J, Huang X:

549          **Correction of the marfan syndrome pathogenic FBN1 mutation by base**

550 **editing in human cells and heterozygous embryos.** *Mol Ther* 2018, **26:**2631-
551 2637.

552 13. Huang TP, Zhao KT, Miller SM, Gaudelli NM, Oakes BL, Fellmann C, Savage
553 DF, Liu DR: **Circularly permuted and PAM-modified Cas9 variants**
554 **broaden the targeting scope of base editors.** *Nat Biotechnol* 2019, **37:**626-
555 631.

556 14. Jin S, Zong Y, Gao Q, Zhu Z, Wang Y, Qin P, Liang C, Wang D, Qiu JL, Zhang
557 F, Gao C: **Cytosine, but not adenine, base editors induce genome-wide off-**
558 **target mutations in rice.** *Science* 2019, **364:**292-295.

559 15. Richter MF, Zhao KT, Eton E, Lapinaite A, Newby GA, Thuronyi BW, Wilson
560 C, Koblan LW, Zeng J, Bauer DE, et al: **Phage-assisted evolution of an**
561 **adenine base editor with improved Cas domain compatibility and activity.**
562 *Nat Biotechnol* 2020, **38:**883-891.

563 16. Yan D, Ren B, Liu L, Yan F, Li S, Wang G, Sun W, Zhou X, Zhou H: **High-**
564 **efficiency and multiplex adenine base editing in plants using new TadA**
565 **variants.** *Mol Plant* 2021, **14:**722-731.

566 17. Gaudelli NM, Lam DK, Rees HA, Sola-Esteves NM, Barrera LA, Born DA,
567 Edwards A, Gehrke JM, Lee SJ, Liquori AJ, et al: **Directed evolution of**
568 **adenine base editors with increased activity and therapeutic application.**
569 *Nat Biotechnol* 2020, **38:**892-900.

570 18. Wolf J, Gerber AP, Keller W: **TadA, an essential tRNA-specific adenosine**
571 **deaminase from *Escherichia coli*.** *EMBO J* 2002, **21:**3841-3851.

572 19. Zhou C, Sun Y, Yan R, Liu Y, Zuo E, Gu C, Han L, Wei Y, Hu X, Zeng R, et
573 al: **Off-target RNA mutation induced by DNA base editing and its**
574 **elimination by mutagenesis.** *Nature* 2019, **571:**275-278.

575 20. Grunewald J, Zhou R, Iyer S, Lareau CA, Garcia SP, Aryee MJ, Joung JK:
576 **CRISPR DNA base editors with reduced RNA off-target and self-editing**
577 **activities.** *Nat Biotechnol* 2019, **37:**1041-1048.

578    21.    Li J, Yu W, Huang S, Wu S, Li L, Zhou J, Cao Y, Huang X, Qiao Y: **Structure-**
579           **guided engineering of adenine base editor with minimized RNA off-**
580           **targeting activity.** *Nat Commun* 2021, **12:**2287.

581    22.    Kim S, Scheffler K, Halpern AL, Bekritsky MA, Noh E, Kallberg M, Chen X,
582           Kim Y, Beyter D, Krusche P, Saunders CT: **Strelka2: fast and accurate calling**
583           **of germline and somatic variants.** *Nat Methods* 2018, **15:**591-594.

584    23.    Wilm A, Aw PP, Bertrand D, Yeo GH, Ong SH, Wong CH, Khor CC, Petric R,
585           Hibberd ML, Nagarajan N: **LoFreq: a sequence-quality aware, ultra-**
586           **sensitive variant caller for uncovering cell-population heterogeneity from**
587           **high-throughput sequencing datasets.** *Nucleic Acids Res* 2012, **40:**11189-
588           11201.

589    24.    DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C,
590           Philippakis AA, del Angel G, Rivas MA, Hanna M, et al: **A framework for**
591           **variation discovery and genotyping using next-generation DNA sequencing**
592           **data.** *Nat Genet* 2011, **43:**491-498.

593    25.    Jacquin ALS, Odom DT, Lukk M: **Crisflash: open-source software to**
594           **generate CRISPR guide RNAs against genomes annotated with individual**
595           **variation.** *Bioinformatics* 2019, **35:**3146-3147.

596    26.    Nishimasu H, Shi X, Ishiguro S, Gao L, Hirano S, Okazaki S, Noda T,
597           Abudayyeh OO, Gootenberg JS, Mori H, et al: **Engineered CRISPR-Cas9**
598           **nuclease with expanded targeting space.** *Science* 2018, **361:**1259-1262.

599    27.    Ren B, Liu L, Li S, Kuang Y, Wang J, Zhang D, Zhou X, Lin H, Zhou H: **Cas9-**
600           **NG greatly expands the targeting scope of the genome-editing toolkit by**
601           **recognizing NG and other atypical PAMs in rice.** *Mol Plant* 2019, **12:**1015-
602           1026.

603    28.    Dow LE, Fisher J, O'Rourke KP, Muley A, Kastenhuber ER, Livshits G,
604           Tschaharganeh DF, Socci ND, Lowe SW: **Inducible *in vivo* genome editing**
605           **with CRISPR-Cas9.** *Nat Biotechnol* 2015, **33:**390-394.

606   29.   Nihongaki Y, Yamamoto S, Kawano F, Suzuki H, Sato M: **CRISPR-Cas9-**
607         **based photoactivatable transcription system.** *Chem Biol* 2015, **22:**169-174.

608   30.   Li H, Durbin R: **Fast and accurate long-read alignment with Burrows-**
609         **Wheeler transform.** *Bioinformatics* 2010, **26:**589-595.

610   31.   Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham
611         A, Keane T, McCarthy SA, Davies RM, Li H: **Twelve years of SAMtools and**
612         **BCFtools.** *Gigascience* 2021, **10**.

613   32.   Thorvaldsdottir H, Robinson JT, Mesirov JP: **Integrative genomics viewer**
614         **(IGV): high-performance genomics data visualization and exploration.**
615         *Brief Bioinform* 2013, **14:**178-192.

616   33.   Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P,
617         Chaisson M, Gingeras TR: **STAR: ultrafast universal RNA-seq aligner.**
618         *Bioinformatics* 2013, **29:**15-21.

619   34.   Crooks GE, Hon G, Chandonia JM, Brenner SE: **WebLogo: a sequence logo**
620         **generator.** *Genome Res* 2004, **14:**1188-1190.

621

622

623

624

625

626    **Figure Legends**

627    **Fig. 1 Profiling of off-target effects caused by ABE-mediated base editing in rice.**

628    **a** The gene architecture of four base editors: rBE46b, rBE49b, rBE50, and rBE53. Ubi-

629    P, maize ubiquitin 1 promoter, *NLS*, nuclear localization sequence; NOS, nopaline

630    synthase terminator. **b** Diagram of the experimental design. For plants in pink

631    rectangles, both genomes and transcriptomes were sequenced. For plants in blue

632    rectangles, only genomes were sequenced.

633

634    **Fig. 2 Characterization of ABE-induced genomic mutations.**

635    **a, b** Number of indels, SNVs, and A>G SNVs, and percentage of A>G SNVs identified

636    for plants that had undergone tissue culture (C1) or *Agrobacterium* infection (C2) and

637    plants harboring SpCas9n-TadA8e (rBE46b), SpCas9n-TadA9 (rBE49b), SpCas9n-

638    NG-TadA8e (rBE50), and SpCas9n-NG-TadA9 (rBE53). In each plot, each dot

639    represents the number of indels, SNVs, and A>G SNVs, and the percentage of A>G

640    SNVs from an individual plant; each middle line represents the median value; and each

641    upper line and lower line represent the standard errors. **c** Number of SNVs and A>G

642    SNVs, and percentage of A>G SNVs were compared for ABE-edited plants harboring

643    TadA8e or TadA9: rBE46b versus rBE49b, and rBE50 versus rBE53. **d** Number of

644    SNVs and A>G SNVs, and percentage of A>G SNVs were compared for ABE-edited

645    plants harboring SpCas9n or SpCas9n-NG: rBE46b versus rBE50, and rBE49b versus

646    rBE53. **e** Percentage of A>G SNVs at given regions for plants in control groups or

647    carrying one of the four ABEs. Each bar represents the mean value, and each error bar

648    represents the standard error. (ns) denotes $p$-value > 0.1, (*) denotes $p$-value < 0.1, (**)

649    denotes $p$-value < 0.01, and (***) denotes $p$-value < 0.001 (one-tailed Wilcoxon test).

650

651    **Fig. 3 ABE-induced DNA mutations in different T-DNA insertion events.**

652    **a** IGV browser views showing the read coverages at T-DNA insertion sites. Lines

653    46bM_s2 and 46bM_s3, 49bM_s2 and 49bM_s3, and 49bAG_s3 and 49bAG_s4 were

24

654 germinated from the same calli. Regions in red rectangles are the T-DNA insertion sites.

655 **b** Number of SNVs and A>G SNVs, and percentage of A>G SNVs. Set 1 represents

656 the unique SNVs only in 46bM_s2, 49bM_s2, and 49bAG_s3. Set 2 represents the

657 unique SNVs only in 46bM_s3, 49bM_s3, and 49bAG_s4. Overlap represents the

658 overlapping SNVs in 46bM_s2 and 46bM_s3, 49bM_s2 and 49bM_s3, and 49bAG_s3

659 and 49bAG_s4. **c** Number of SNVs and A>G SNVs, and percentage of A>G SNVs in

660 plants with partial or whole T-DNA insertions of rBE50 or rBE53. Each bar represents

661 the mean value, each error bar represents the standard error, and each dot represents the

662 number of SNVs, the number of A>G SNVs, and percentage of A>G SNVs of each

663 plant. (ns) denotes $p$-value > 0.1, (*) denotes $p$-value < 0.1 (one-tailed Wilcoxon test).

664

665 **Fig. 4 Transcriptome-wide ABE-induced off-target mutations**

666 **a** Number of SNVs and A>G SNVs, and percentage of A>G SNVs in plants harboring

667 SpCas9 (Cas), SpCas9n-TadA8e (rBE46b), and SpCas9n-TadA9 (rBE49b). **b** Ratios of

668 A>G mutations were calculated for A>G SNV loci detected in lines R49bAG_s2 and

669 R49bAG_s3 and shown in the scatterplot. The Pearson correlation coefficient ($r$) was

670 also calculated, and the red line is the diagonal line. **c** A sequence logo derived from

671 edited adenines from all RNA-seq data. Bits account for how much each column is

672 conserved and how much the nucleotide frequencies obtained in the profile differ from

673 those that would have been obtained by aligning oligonucleotides chosen at random. **d**

674 Boxplot showing ratios of A>G mutations at all RNA A>G SNV loci for plants

675 harboring SpCas9, rBE46b, and rBE49b. A Wilcoxon test was conducted between

676 every plant harboring ABEs versus plants harboring Cas only, and the -log10 $p$-value

677 is shown. **e** Bar plot showing the average RPM values of ABEs for plants without RNA

678 mutations and plants with RNA mutations. Each bar represents the mean value, each

679 error bar represents the standard error, and each dot represents the ABE RPM value of

680 each plant. (***) denotes $p$-value < 0.001 (one-tailed Wilcoxon test). **f** Ratios of A>G

681 mutations of all A>G RNA SNV loci were calculated for one 49bAG_s2 $T_0$ plant and

25

682    four 49bAG_s2 $T_1$ plants (left). -log10 *p*-value of Wilcoxon test on A>G ratios between

683    five 49bAG_s2 plants versus plants harboring SpCas9 (middle). RPMs of ABEs are

684    shown in the bar plot (right). N1 and N2 are $T_1$ 49bAG_s2 plants with a T-DNA

685    insertion, while N3 and N4 are $T_1$ 49bAG_s2 plants without a T-DNA insertion.

686

687    **Fig. 5 ABE-induced clustered RNA and DNA A>G SNVs**

688    **a** An IGV genome browser view showing representative loci with clustered A>G SNVs

689    in transcriptomes. **b** Ratios of A>G mutations were calculated in flanking 5′ and 3′ 30-

690    bp regions centered at A>G RNA SNV loci. Lines R49bAG_s2 and R49bAG_s3 with

691    RNA mutations and line RCas_s1 with SpCas9 only are shown. **c** Boxplot showing

692    number of A>G SNVs in the flanking 5′ and 3′ 30-bp regions separately for RNA SNVs

693    in many (3-8) or few (1-2) plants. **d** IGV genome browser views showing representative

694    SNV loci with flanking A>G SNVs in whole-genome sequencing. **e** Ratios of clustered

695    SNVs located in genic regions. **f** Plants with ABEs were classified into two groups:

696    group 1 with clustered SNVs and group 2 without clustered SNVs. Number of SNVs

697    and A>G SNVs, and percentage of A>G SNVs are shown separately for plants in group

698    1 and plants in group 2. (**) denotes *p*-value < 0.01, and (***) denotes *p*-value < 0.001

699    (one-tailed Wilcoxon test). In IGV genome browser views, the grey bar represents a

700    sequenced nucleotide that is the same as the reference genome, while bars in other

701    colors represent sequenced nucleotides that are partially or totally different from the

702    reference genome: red represents nucleotide A, green represents nucleotide T, orange

703    represents nucleotide G, and blue represents nucleotide C. The height of each color bar

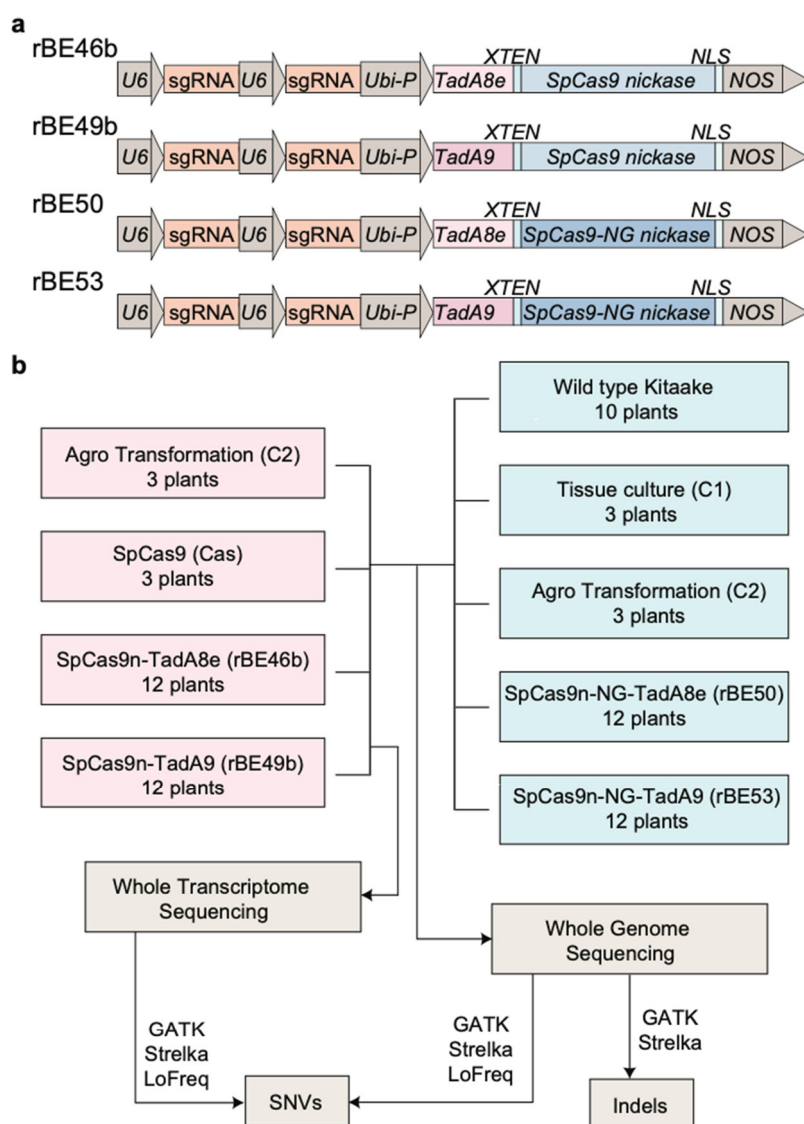704    represents the relative composition of each nucleotide.

705

706

**Fig. 1 Profiling of off-target effects caused by ABE-mediated base editing in rice.**

**a** The gene architecture of four base editors: rBE46b, rBE49b, rBE50, and rBE53. Ubi-P, maize ubiquitin 1 promoter, *NLS*, nuclear localization sequence; NOS, nopaline synthase terminator. **b** Diagram of the experimental design. For plants in pink rectangles, both genomes and transcriptomes were sequenced. For plants in blue rectangles, only genomes were sequenced.
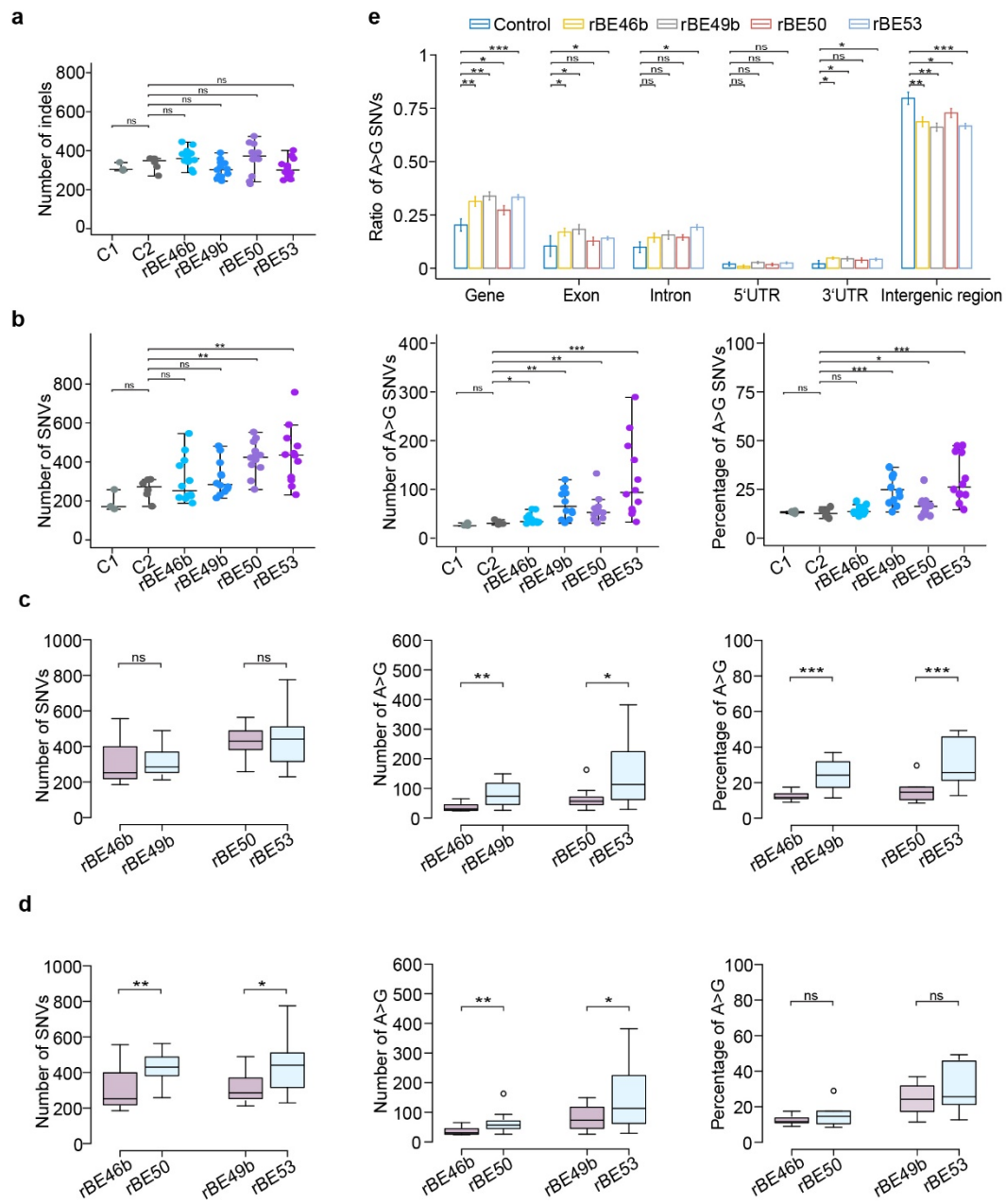
713

**Fig. 2 Characterization of ABE-induced genomic mutations.**

**a, b** Number of indels, SNVs, and A>G SNVs, and percentage of A>G SNVs identified for plants that had undergone tissue culture (C1) or *Agrobacterium* infection (C2) and plants harboring SpCas9n-TadA8e (rBE46b), SpCas9n-TadA9 (rBE49b), SpCas9n-NG-TadA8e (rBE50), and SpCas9n-NG-TadA9 (rBE53). In each plot, each dot represents the number of indels, SNVs, and A>G SNVs, and the percentage of A>G SNVs from an individual plant; each middle line represents the median value; and each upper line and lower line represent the standard errors. **c** Number of SNVs and A>G

742    SNVs, and percentage of A>G SNVs were compared for ABE-edited plants harboring

743    TadA8e or TadA9: rBE46b versus rBE49b, and rBE50 versus rBE53. **d** Number of

744    SNVs and A>G SNVs, and percentage of A>G SNVs were compared for ABE-edited

745    plants harboring SpCas9n or SpCas9n-NG: rBE46b versus rBE50, and rBE49b versus

746    rBE53. **e** Percentage of A>G SNVs at given regions for plants in control groups or

747    carrying one of the four ABEs. Each bar represents the mean value, and each error bar

748    represents the standard error. (ns) denotes $p$-value > 0.1, (*) denotes $p$-value < 0.1, (**)

749    denotes $p$-value < 0.01, and (***) denotes $p$-value < 0.001 (one-tailed Wilcoxon test).

750

751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
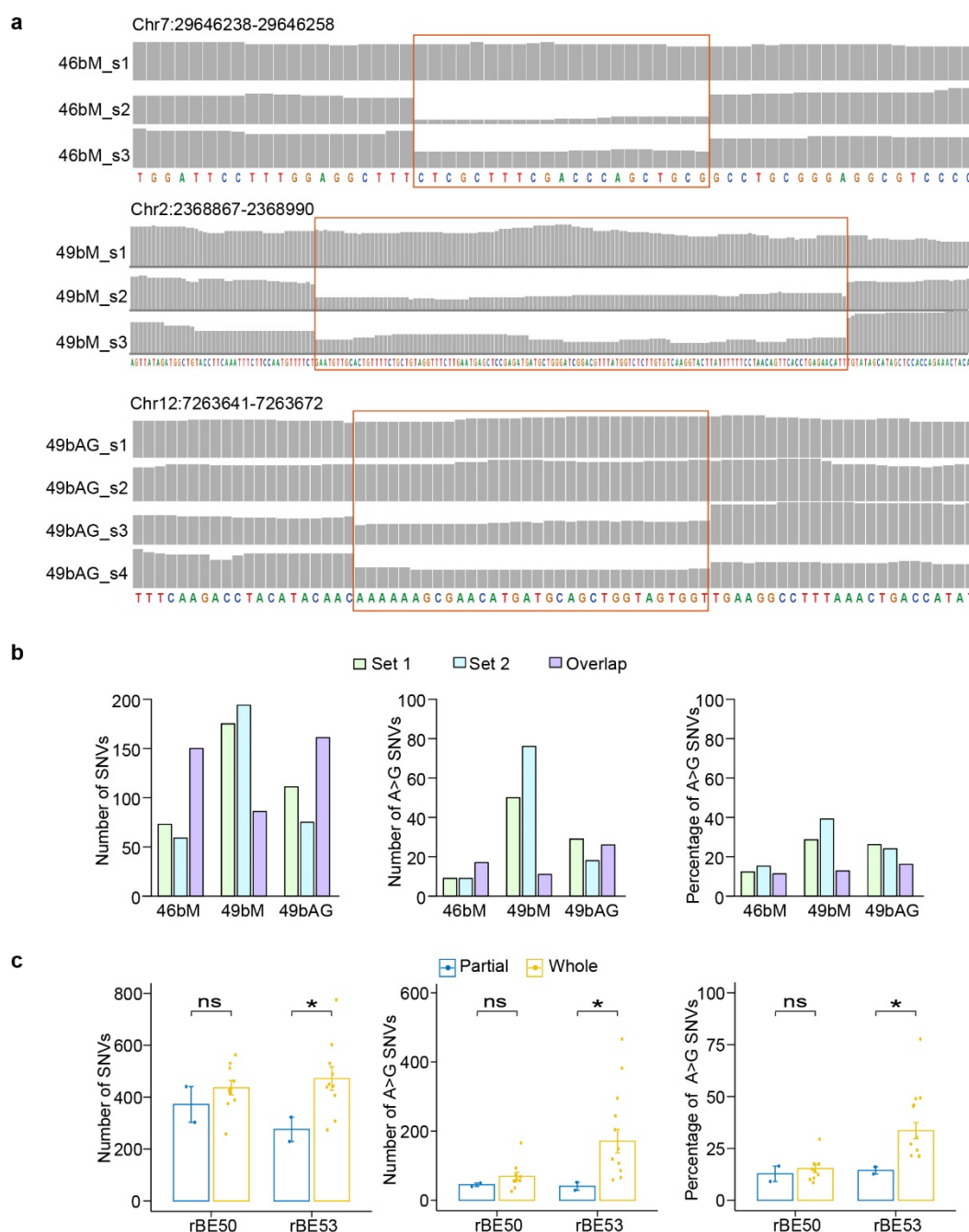767
768
769
770
771



772 **Fig. 3 ABE-induced DNA mutations in different T-DNA insertion events.**

773 **a** IGV browser views showing the read coverages at T-DNA insertion sites. Lines

774 46bM_s2 and 46bM_s3, 49bM_s2 and 49bM_s3, and 49bAG_s3 and 49bAG_s4 were

775 germinated from the same calli. Regions in red rectangles are the T-DNA insertion sites.

776 **b** Number of SNVs and A>G SNVs, and percentage of A>G SNVs. Set 1 represents

777 the unique SNVs only in 46bM_s2, 49bM_s2, and 49bAG_s3. Set 2 represents the

778 unique SNVs only in 46bM_s3, 49bM_s3, and 49bAG_s4. Overlap represents the

30

779    overlapping SNVs in 46bM_s2 and 46bM_s3, 49bM_s2 and 49bM_s3, and 49bAG_s3

780    and 49bAG_s4. **c** Number of SNVs and A>G SNVs, and percentage of A>G SNVs in

781    plants with partial or whole T-DNA insertions of rBE50 or rBE53. Each bar represents

782    the mean value, each error bar represents the standard error, and each dot represents the

783    number of SNVs, the number of A>G SNVs, and percentage of A>G SNVs of each

784    plant. (ns) denotes $p$-value > 0.1, (*) denotes $p$-value < 0.1 (one-tailed Wilcoxon test).
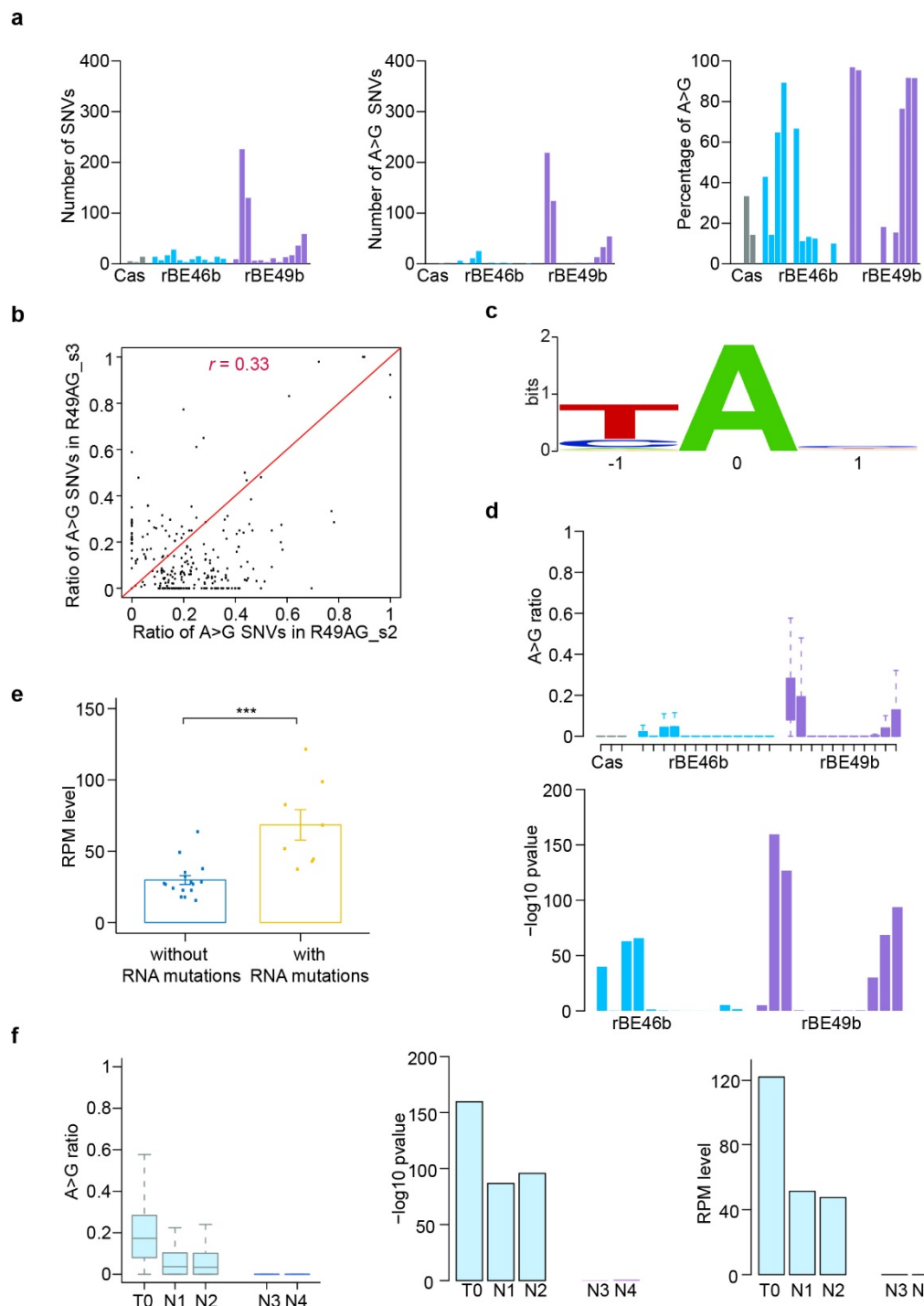
785

786

**Fig. 4 Transcriptome-wide ABE-induced off-target mutations**

**a** Number of SNVs and A>G SNVs, and percentage of A>G SNVs in plants harboring

SpCas9 (Cas), SpCas9n-TadA8e (rBE46b), and SpCas9n-TadA9 (rBE49b). **b** Ratios of

A>G mutations were calculated for A>G SNV loci detected in lines R49bAG_s2 and

R49bAG_s3 and shown in the scatterplot. The Pearson correlation coefficient (*r*) was

815    also calculated, and the red line is the diagonal line. **c** A sequence logo derived from

816    edited adenines from all RNA-seq data. Bits account for how much each column is

817    conserved and how much the nucleotide frequencies obtained in the profile differ from

818    those that would have been obtained by aligning oligonucleotides chosen at random. **d**

819    Boxplot showing ratios of A>G mutations at all RNA A>G SNV loci for plants

820    harboring SpCas9, rBE46b, and rBE49b. A Wilcoxon test was conducted between

821    every plant harboring ABEs versus plants harboring Cas only, and the -log10 $p$-value

822    is shown. **e** Bar plot showing the average RPM values of ABEs for plants without RNA

823    mutations and plants with RNA mutations. Each bar represents the mean value, each

824    error bar represents the standard error, and each dot represents the ABE RPM value of

825    each plant. (***) denotes $p$-value < 0.001 (one-tailed Wilcoxon test). **f** Ratios of A>G

826    mutations of all A>G RNA SNV loci were calculated for one 49bAG_s2 $T_0$ plant and

827    four 49bAG_s2 $T_1$ plants (left). -log10 $p$-value of Wilcoxon test on A>G ratios between

828    five 49bAG_s2 plants versus plants harboring SpCas9 (middle). RPMs of ABEs are

829    shown in the bar plot (right). N1 and N2 are $T_1$ 49bAG_s2 plants with a T-DNA

830    insertion, while N3 and N4 are $T_1$ 49bAG_s2 plants without a T-DNA insertion.
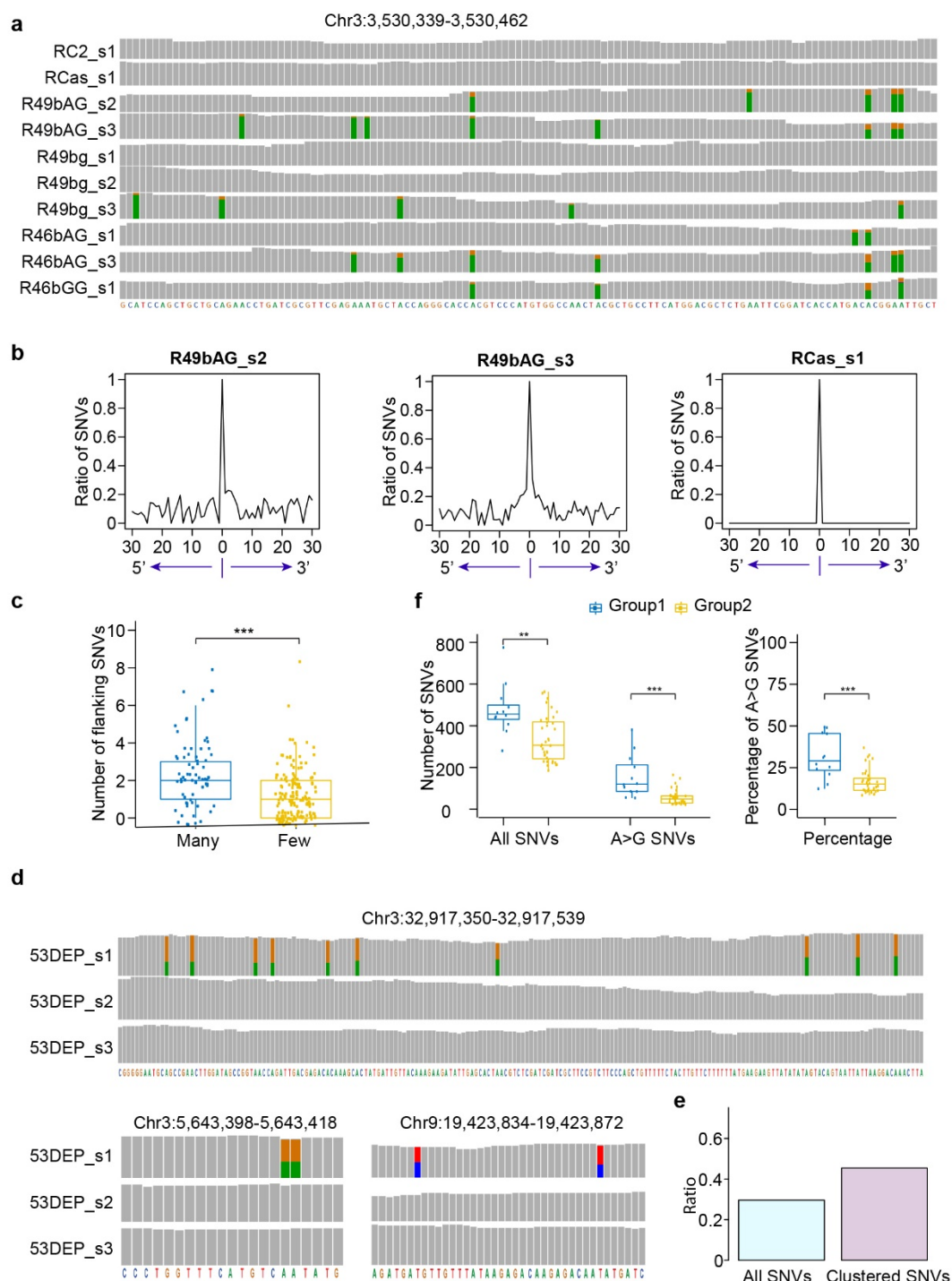
831

832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854



855 **Fig. 5 ABE-induced clustered RNA and DNA A>G SNVs**

856 **a** An IGV genome browser view showing representative loci with clustered A>G SNVs

857 in transcriptomes. **b** Ratios of A>G mutations were calculated in flanking 5′ and 3′ 30-

858 bp regions centered at A>G RNA SNV loci. Lines R49bAG_s2 and R49bAG_s3 with

859 RNA mutations and line RCas_s1 with SpCas9 only are shown. **c** Boxplot showing

860    number of A>G SNVs in the flanking 5′ and 3′ 30-bp regions separately for RNA SNVs

861    in many (3-8) or few (1-2) plants. **d** IGV genome browser views showing representative

862    SNV loci with flanking A>G SNVs in whole-genome sequencing. **e** Ratios of clustered

863    SNVs located in genic regions. **f** Plants with ABEs were classified into two groups:

864    group 1 with clustered SNVs and group 2 without clustered SNVs. Number of SNVs

865    and A>G SNVs, and percentage of A>G SNVs are shown separately for plants in group

866    1 and plants in group 2. (**) denotes $p$-value < 0.01, and (***) denotes $p$-value < 0.001

867    (one-tailed Wilcoxon test). In IGV genome browser views, the grey bar represents a

868    sequenced nucleotide that is the same as the reference genome, while bars in other

869    colors represent sequenced nucleotides that are partially or totally different from the

870    reference genome: red represents nucleotide A, green represents nucleotide T, orange

871    represents nucleotide G, and blue represents nucleotide C. The height of each color bar

872    represents the relative composition of each nucleotide.