# Transcription-coupled repair in *Drosophila melanogaster* is independent of the mismatch repair pathway

Lauri Törmä[1,2], Claire Burny[1,2], Viola Nolte[1], Kirsten-André Senti[1], Christian Schlötterer[1,#]

1) Institut für Populationsgenetik, Vetmeduni Vienna, Vienna, Austria

2) Vienna Graduate School of Population Genetics, Vetmeduni Vienna, Vienna, Austria

#) corresponding author.

Keywords: Transcription-coupled repair, *Drosophila melanogaster*, mutational strand bias, mismatch-repair, *spellchecker1*

Correspondence:

Christian Schlötterer

Institut für Populationsgenetik, Vetmeduni Vienna, 1210 Wien, Austria

Email: christian.schloetterer@vetmeduni.ac.at

1    **Abstract**

2    Transcription-coupled repair (TCR) removes base damage on the transcribed strand of a gene to

3    ensure a quick resumption of transcription. Based on the absence of key enzymes for TCR and

4    empirical evidence, TCR was thought to be missing in *Drosophila melanogaster*. The recent

5    demonstration of TCR in S2 cells raises the question about the involved genes. Since the

6    mismatch repair (MMR) pathway serves a central role in TCR, at least in *Escherichia coli*, we

7    studied the mutational signatures in flies with a deletion of the MMR gene *spellchecker1* (*spel1*),

8    a MutS homolog. Whole-genome sequencing of mutation accumulation (MA) lines obtained 7,345

9    new single nucleotide variants (SNVs) and 5,672 short indel mutations, the largest data set from

10   an MA study in *D. melanogaster.* Based on the observed mutational strand-asymmetries, we

11   conclude that TCR is still active without *spel1*. The operation of TCR is further confirmed by a

12   negative association between mutation rate and gene expression. Surprisingly, the TCR

13   signatures are detected for introns, but not for exons. We propose that an additional exon-specific

14   repair pathway is masking the signature of TCR. This study presents the first step towards

15   understanding the molecular basis of TCR in *Drosophila melanogaster*.

**Background**

DNA continuously undergoes a large number of spontaneous chemical modifications leading to DNA damage (1, 2). The damaged bases can cause mutations, block DNA replication, and interfere with transcription (3). To repair some of these adducts, nucleotide excision repair (NER) removes the damaged strand at short distances from both sides of the lesion and a new strand is synthesized to fill the gap (4). NER has two pathways to recognize these lesions: the global genomic repair (GGR) and the transcription-coupled repair (TCR) (5). GGR scans the whole genome and the DNA damage is recognized by helix-distorting lesions (4). TCR detects the base damage from RNA polymerase stalling in actively transcribed DNA (4, 5) and leads to a mutational asymmetry between the strands (6). Transcription inhibition can be dangerous to a cell or even an organism (7, 8) so a quick resumption of transcription is vital. TCR is found in most bacterial species and many eukaryotes (9) and a defective pathway causes strong disease phenotypes, such as xeroderma pigmentosum and Cockayne's syndrome in humans (10).

*Drosophila melanogaster* presents an interesting case where the GGR pathway for NER is present but TCR was thought to be missing (11, 12). The lack of fly homologs for genes required for TCR in other organisms — CSA/ERCC8 and CSB/ERCC6 — suggested that the pathway was lost during evolution (11). Furthermore, biochemical studies failed to detect TCR after UV-induced damage in *D. melanogaster* cell cultures (13, 14). In addition to the indirect evidence for TCR which is based on positive correlation of the compositional skews in introns (15) with expression (16), a recent study showed that TCR is operating in *Drosophila* S2 cells (17). This result raises the important question of how *Drosophila* is able to perform TCR when the key genes CSA and CSB are absent.

In *E. coli,* mismatch repair genes MutS and MutL are required for TCR (18). In yeast, genes required for NER interact with MMR genes (19) but MMR deficient cells are still performing TCR

42    (20). In humans, MMR repair also interacts with NER (21) and evidence suggests that the pathway

43    is involved in TCR of UV and oxidative damage (22). However, the issue remains controversial

44    (23). Given the uncertainty about the functional basis of TCR in *Drosophila*, we determined the

45    influence of the MutS homolog *spellchecker1* (*spel1*) on TCR in flies using MMR deficient

46    mutation accumulation (MA) lines. The lack of MMR in MA lines is expected to result in a high

47    number of mispaired bases. Such bases do not only lead to mutations but also cause local

48    structural and dynamic distortions in the DNA structure (24) and are hotspots for DNA damage

49    due to the higher susceptibility of unpaired bases to chemical modifications (25). For example,

50    the loss of *msh2* in mice, *Trypanosoma brucei*, and *T. cruzi* increases oxidative damage of

51    guanine by reactive oxygen species (26–28), which is repaired by TCR in murine cells (29).

52

53    Using a mismatch repair-deficient background, we find mutational asymmetries that are

54    negatively associated with germline expression intensities, demonstrating functional TCR without

55    *spel1.* Based on the absence of the TCR signatures in exons, we propose that an additional,

56    exon-specific repair mechanism is operating.

57

58    **Results**

59    We generated a *spel1* null mutant using CRISPR with guide RNAs targeting the 5' and 3' ends of

60    the gene. Our mutant contained a double insertion of the template plasmid with the backbone of

61    the vector (Supplementary Figure 1), a frequent event arising from the recombination of two

62    plasmids into the locus (30). We propagated seven independent lines for 10 generations by

63    brother-sister mating and identified 7,345 new single nucleotide variants (SNV) and 5,672 indels

64    in females from these mutation accumulation lines. With 73.5% of the non-synonymous

65    substitutions on the autosomes and 77.0% on the X chromosome, our data did not significantly

66    deviate from the 75% expected under neutrality (31) (Fisher's exact test (FET), p=0.4143 for the

67    autosomes; FET, p=0.6573 for X).

68

69     The presence of TCR can be detected by mutational asymmetries between the transcribed and

70     non-transcribed strands (6). The identification of mutational asymmetries is critically dependent

71     on the correct null hypothesis. *Drosophila* introns have a skewed base composition, which

72     depends on transcription levels (16). We confirmed that the fraction of thymines and cytosines on

73     the transcribed strand is significantly negatively associated with expression in both ovaries (Wald

74     test, OR=0.9977, p<2.2e-16 for thymines; Wald test, OR=0.9997, p=1.16e-13 for cytosines) and

75     testes (Wald test, OR=0.9982, p<2.2e-16 for thymines; Wald test, OR=0.9994, p<2.2e-16 for

76     cytosines) (Figure 1. a,b). We accounted for this by including the bias into the formulation of a

77     null hypothesis for the expected number of mutations on the transcribed and on the non-

78     transcribed strands. We calculated the expected bias with two different approaches: from the

79     mutated genes and from a sample of genes with a similar expression as the mutated genes (see

80     Methods). Both approaches produced highly consistent results.

81

82     5,071 SNVs located in genes were used to test the ratio of bias-adjusted mutation rates on the

83     transcribed and non-transcribed strands for every mutation type. Without TCR, a rate ratio (RR)

84     of 1 is expected and the statistical significance can be determined with a Poisson test. After

85     multiple testing correction, C>A mutations occurred less often (Poisson test, RR=0.82; 95% CI:

86     0.70-0.95, adjusted p-value=0.038), (Figure 1. c; Supplementary Table 1) and T>C mutations

87     more often on the transcribed strand (Poisson test, RR=1.12, 95% CI: 1.02-1.22, adjusted p-

88     value=0.039) (Figure 1. c; Supplementary Table 1). Similar results were obtained using the gene

89     expression sampling scheme (see Methods, Supplementary Figure 2). Assuming that TCR is

90     causing this bias, this implies that cytosine and adenine are more likely to experience base

91     damage than other bases in MMR deficient flies.

92

93   DNA repair and damage processes can differ between exons and introns (32, 33). We therefore

94   analyzed exons and introns separately. After excluding SNVs which overlapped both exon and

95   intron annotations, C>A mutations occurred less often on the transcribed strand (RR=0.732, 95%

96   CI: 0.597-0.895, adjusted p-value=0.014), but exonic C>A mutations did not (RR=0.995, 95% CI:

97   0.783-1.263, adjusted p-value=1) (Figure 1. d; Supplementary Table 2). Despite intronic T>C

98   mutations occurring slightly more often on the transcribed strand, this was not significant (Poisson

99   test, RR=1.113, 95% CI: 0.996-1.243, adjusted p-value=0.155). However, looking for the effect of

100  the 5' and 3' bases flanking the mutation, we observed that the A[T>C]N context is exhibiting a

101  significant strand bias in introns with a rate ratio of 1.472 (Poisson test, 95% CI: 1.138-1.910,

102  adjusted p-value=0.01) but not in exons (Poisson test, RR=0.894, 95% CI: 0.627-1.280, adjusted

103  p-value=0.866). No other contexts exhibited strand bias (Figure 1. d; Supplementary Table 2).

104  Since the null hypothesis was not adjusted for triplet composition, we updated our null hypothesis

105  to take into account the 5' and 3' flanking bases by performing a permutation test (see Methods)

106  and obtained similar results. Intronic A[T>C]N mutations still exhibited a significant strand bias

107  (permutation test, p=0.001) while exonic mutations did not (permutation test, p=0.215)

108  (Supplementary                                  Figure                                  3).

109

110  To confirm that the strand bias is caused by TCR, we tested for expression differences in genes

111  containing C>A or A[T>C]N mutations. In the case of an active TCR, a correlation between strand

112  asymmetry and gene expression is expected, because DNA damage on the transcribed strand is

113  more likely to be detected in highly expressed genes. Thus, mutations arising from DNA damage

114  on the transcribed strand should be found in lowly expressed genes. We used the FlyAtlas2 (34)

115  expression data set from ovaries and testes as a proxy for the expression environment where the

116  mutations occurred. Consistent with these predictions, we found that the genes with intronic C>A

117  mutations on the transcribed strand have on average lower expression in both ovaries (one-sided

118  Wilcoxon rank-sum test, adjusted p-value=0.022) and testes (one-sided Wilcoxon rank-sum test,

119   adjusted p-value=0.022) than genes with intronic C>A mutations on the non-transcribed strand

120   (Figure 2. a). As expected from the lack of strand bias, the expression level of genes with exonic

121   C>A mutations were not different (Figure 2. a). Genes with context-dependent A[T>C]N mutations

122   were not differentially expressed (Figure 2. b). This could be due to either a lack of power or

123   because the expression data used does not reflect the expression environment where the base

124   damage occurred.

125

126   While the expression analysis suggests that TCR is responsible for the strand bias for C>A

127   mutations, it is important to rule out the alternative explanation of a mutagenic effect of

128   transcription on the non-transcribed strand. We used a randomization procedure (see Methods)

129   to test if C>A mutations occur more frequently on the non-transcribed strand of highly expressed

130   genes. Consistent with previous observations (16, 31), we found no evidence that transcription is

131   mutagenic neither in testes (randomization test, p=0.611) nor in ovaries (randomization test,

132   p=0.403) (Supplementary Figure 4) ruling it out as the source of the strand bias.

133

134   Based on the combined evidence, we conclude that TCR is operating in *Drosophila* biasing the

135   C>A mutations and *spel1* is not required. Nevertheless, it is not clear why TCR signatures are

136   only detected for introns, but not for exons. Two different explanations can account for the lack of

137   mutational strand bias in exons for the C>A mutations: i) TCR requires *spel1* in exons or ii) an

138   additional DNA repair mechanism is operating on exons, which erases the signal of TCR. The

139   two explanations can be distinguished based on their different predictions for the relative mutation

140   rates. MMR dependence for exons predicts an increased mutation rate for exons on the

141   transcribed strand while the latter predicts a reduced exonic mutation rate for the non-transcribed

142   strand. To test these hypotheses, we performed a permutation test while controlling for the triplet

143   context in exons and introns to test for relative mutation rate differences. We found no evidence

144   of elevated exonic mutation rate on the transcribed strand (permutation test, p=0.4762)

145 (Supplementary Figure 5) showing that the lack of *spel1* does not cause the missing strand bias.

146 We found signs - although nonsignificant - of reduced exonic mutation rate on the non-transcribed

147 strand for C>A mutations (Supplementary Figure 5) (permutation test, p=0.067) suggesting that

148 the lack of exonic strand bias may be caused by a favorable repair. The A[T>C]N did not show

149 differences in the relative mutation rates on the transcribed (permutation test, p=0.475) or the

150 non-transcribed strand (permutation test, p=0.126) (Supplementary Figure 6).

151

152 **Discussion**

153 We demonstrated that TCR is independent of MMR in flies by uncovering TCR-induced mutational

154 asymmetries in intronic C>A mutations in MMR deficient *D. melanogaster* mutation accumulation

155 lines. Because UV-light was not used during the experiment, we are able to demonstrate that

156 TCR in flies is not only limited to UV-induced damage, as previously seen (17), but can also repair

157 other types of DNA damage. The C>A mutations can arise from mismatches with oxidatively

158 damaged DNA (35). An important finding is that TCR does not cause mutational asymmetry in

159 exons. We ruled out that this is caused by the MMR deficiency and found support for a pathway

160 that protects exons over introns thus masking the signatures of TCR. A similar finding was made

161 in human cells where less oxidative DNA damage accumulates in exons than in introns — possibly

162 due to a favorable repair (33). If a similar process is occurring in flies, as our data suggest, we

163 propose that the global repair pathway of nucleotide excision repair is favoring exons over introns.

164 The global repair does not discriminate between the transcribed and non-transcribed strands and

165 detects the same lesions as TCR thus explaining the lack of strand bias, the gene expression

166 difference, and the signs of reduced exonic mutation rate on the non-transcribed strand for C>A

167 mutations.

168

169 In summary, generating the largest de novo mutation data set from an MA study in *D.*

170 *melanogaster*, we demonstrated that TCR operates against DNA damage in the germline

171   independent of the MMR pathway. We uncovered differences in mutational processes of exons

172   and introns and attribute this to an additional repair operating on exons. We anticipate the use of

173   *spel1* mutations will become a widely used approach to study mutation patterns in a broad range

174   of species.

175

176   **Materials and methods**

177   **Generating the *spel1* deletion and mutation accumulation.**

178   The *spel1* null mutant was generated from an isogenized *Oregon-R* strain using the CRISPR-

179   Cas9 genome engineering tool. The 2nd and the 3rd chromosomes were isogenized with balancers

180   and the variation on the X chromosome was reduced by 5 generations of full-sib mating. Two

181   gRNAs targeting the second and the last exon of *spel1* were cloned with the Gibson Assembly®

182   Cloning Kit (New England Biolabs) into a BbsI (10,000 units/ml, NEB, R0539) digested pCDF4

183   (50) (Addgene plasmid # 49411; http://n2t.net/addgene:49411; RRID:Addgene_49411)

184   expression vector. The ligation product was transformed into SURE2 cells and the construct was

185   verified by Sanger sequencing.

186

187   A template for homology-directed repair was generated by Golden gate cloning. 1 kb homology

188   arms were amplified from genomic DNA with primers LT41-LT44. Purified amplicons were mixed

189   (30 ng each) with 50 ng pJET1.2-STOP-dsRed (51) (Addgene plasmid # 60944 ;

190   http://n2t.net/addgene:60944 ; RRID:Addgene_60944), 50 ng pBS-GGAC-ATGC (51) (Addgene

191   plasmid # 60949 ; http://n2t.net/addgene:60949 ; RRID:Addgene_60949), 1.5 µl 10x T4 ligation

192   buffer, 1 µl BsmBI (10,000 units/ml, NEB, R0580), and water was added to 14 µl. After incubation

193   for min at 55°C 1 µl T4 ligase (400,000 units/ml, NEB, M0202) was added. Ligation was performed

194   by cycling the reaction between 5 min in 42°C and 5 min in 16°C overnight. Final digestion was

195   performed for 30 min in 55°C followed by 10 min at 80°C to inactivate the enzyme.  The ligation

196   product was treated with Plasmid-safe nuclease (10,000 units/ml, Epicentre, E3101K) and

197  transformed into SURE2 cells. Positive colonies were identified with colony PCR and recovered

198  plasmids were verified by sequencing.

199

200  The germline transformation was achieved by microinjecting a mixture of the template (500 ng/µl),

201  the gRNA expression vector (100 ng/µl), and pHsp70-Cas9 (52) (250 ng/µl) (Addgene plasmid #

202  60944 ; http://n2t.net/addgene:60944 ; RRID:Addgene_60944) into dechorionated fly embryos.

203  F1 progeny were screened for the 3XP3::DsRed marker and a correct targeting of *spel1* was

204  confirmed with PCR and sequencing. A PCR was performed to detect the double integration

205  where the template plasmid integrates twice into the locus with the backbone. All the primers used

206  in this study are listed in Supplementary Table 3.

207

208  We performed 10 generations of mutation accumulation with full-sib mating and sequenced

209  individual females from 7 surviving lines.

210

211  **Library preparation and sequencing**

212  Genomic DNA was extracted from a single female fly of each MA line using a standard high salt

213  extraction method (36) with RNase A treatment. From each female, 70 ng genomic DNA was

214  used to prepare paired-end libraries with the NEBNext Ultra II FS DNA Library Prep Kit (New

215  England Biolabs, Ipswich, MA) using only 10% of the reagents recommended in the original

216  protocol of the supplier. After double sided size selection targeting an insert size of 300 bp,

217  libraries were amplified with dual-index primers using 5 PCR cycles. After purification with

218  AMPureXP beads (Beckman Coulter, Brea, CA), the 7 libraries were quantified using the Qubit

219  dsDNA HS Kit (Invitrogen, Carlsbad, CA), combined in equimolar amounts with additional 4

220  libraries from another experiment and sequenced on one lane of a HiSeq2500 using a 2x125bp

221  protocol.

222

223 **QC and reads mapping**

224 Libraries were first demultiplexed using ReadTools (37) (version 1.5.2;

225 AssignReadGroupByBarcode --splitSample, --maximumMismatches 1, providing the

226 corresponding barcodes). The raw reads were assessed for their quality using FastQC software

227 (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/). Low quality tails at 3´ end were

228 trimmed using ReadTools (--mottQualityThreshold 20, --minReadLength 50, --disable5pTrim true)

229 and BAM files were converted to compressed FASTQ files using ReadTools (ReadsToFastq --

230 interleavedInput true --barcodeInReadName true --outputFormat GZIP). As FastQC detected

231 residual levels of adapter contamination, adapter cleaning was performed with the BBTools suit

232 (38) using BBDuk (version 38.32; ktrim=r k=23 mink=11 hdist=1 tbo).

233

234 Processed paired-end reads were mapped to the *D. melanogaster* reference genome release

235 6.24 indexed with the bwa index command using BWA-MEM (39) (version 0.7.17; bwamem) on

236 a Hadoop cluster using DistMap (40) (version 2.7.5).

237

238 PCR duplicates were removed using PICARD (http://broadinstitute.github.io/picard/)

239 MarkDuplicates tool (version 2.21.3; REMOVE_DUPLICATES=true

240 VALIDATION_STRINGENCY=SILENT). We kept mapped reads with each segment properly

241 aligned and removed reads that mapped equally well to multiple positions or have a low mapping

242 quality using SAMtools (41) (version 1.9; -b -q 20 -f 0x002 -F 0x004 -F 0x008). We clipped the

243 processed overlapping paired-end reads using the BamUtil suit (42) (version 1.0.13; bam

244 clipOverlap --in --out --stats).

245

246 **Variants inventory**

247 The fasta reference was indexed using SAMtools faidx command. Processed BAM files were

248 sorted and indexed with SAMtools for each chromosome arm (2L, 2R, 3L, 3R, X) separately using

249 SAMtools view command. We then added a unique read group tag per sample using the PICARD

250 AddOrReplaceReadGroups command. We increased the accuracy of variant calling by using two

251 different tools; Freebayes (45, cloned from https://github.com/ekg/freebayes) (version v0.9.10-3-

252 g47a713e) and GATK HaplotypeCaller (46) (version 4.0.12.0) and kept only variants that were

253 identified with both tools. To use the parallel version of Freebayes, we split the reference into 1Mb

254 regions with Freebayes fasta_generate_regions.py script (python version 2.7.17). For each

255 chromosome arm, we used the freebayes-parallel executable (-C 1 –F 0.01 --min-base-quality

256 20, all other options set to default), providing the 1 Mb regions file and individual BAM files.

257 Second, we followed (31) and used the GATK HaplotypeCaller with --heterozygosity 0.01 option

258 with default settings.

259

260 We obtained two raw lists of variants per chromosome arm in a VCF format (43). Two different

261 filtering procedures were applied for each variant caller.

262 For Freebayes, each raw list of variants was filtered as follows:

263     i. to remove variants based on depth at the variant position using BCFtools (44)

264        (version 1.8; filter –i SAF>0 && SAR>0 && (SAF+SAR+SRF+SRR)>5),

265     ii. to suppress variants within 5-bp of an INDEL using BCFtools (filter -g 5),

266     iii. to keep variants with at most 2 alleles denoted by reference and alternate alleles

267        using VCFtools (43) (version 0.1.15; --vcf --min-alleles 1 --max-alleles 2 --recode-

268        INFO-all --recode),

269     iv. to simplify multi-nucleotide polymorphisms into SNPs using vt cloned from

270        https://github.com/atks/vt (45) (version 0.57721; decompose_blocksub, normalize

271        commands successively),

272     v. to filter for QUAL>40 using VCFtools (--vcf –minQ 40 --recode-INFO-all --recode).

273 For GATK, we used GATK VariantFiltration with the options  --filter-expression "QD < 2.0" --filter-

274 name "QD" --filter-expression "FS > 60.0" --filter-name "FS" --filter-expression "MQ < 40.0" --filter-

275    name "MQ" --filter-expression "MQRankSum < -12.5" --filter-name "MQRankSum" --filter-

276    expression "ReadPosRankSum < -8.0" --filter-name "ReadPosRankSum".

277

278    We intersected the two filtered VCF files retaining only variants with the same position using

279    BEDtools (46) (version 2.27.1; intersect –u –a -b -wa –header). We then extracted private SNPs

280    using BEDtools (intersect –v –a -b –header) providing all bgzipped and tabix-indexed (47) (version

281    1.8; -p vcf) VCF per line. Finally, we subtracted the variants lists with the variants called from 10

282    individual *spel1* null flies, which did not go through MA, as a quality control for residual ancestral

283    alternative alleles after having applied a similar pipeline; we masked the X region

284    6240639:6686943 from line 5 using  BEDtools (intersect –v –a -b –header) where some residual

285    variants were observed. We obtained a final set of 7,345 SNPs and 5,672 INDELs.

286

287    For our analyses we relied on the genome annotation from flybase Dmel-all-filtered-r6.30.gff

288    (downloaded                                                                                    from:

289    ftp://ftp.flybase.net/genomes/Drosophila_melanogaster/dmel_r6.30_FB2019_05/gff/    in    May

290    2019).

291

292    **Statistical analyses**

293    All statistical analyses were done with R (48) (version 3.5.0).

294

295    Fraction of non-synonymous mutations compared to neutral expectation

296    The SnpEff software (49) (version 4.3) was used to distinguish synonymous and nonsynonymous

297    mutations in the longest transcript of each gene. We performed a Fisher's exact test to compare

298    the observed and expected number of synonymous and non-synonymous mutations. Following

299    (31), we used odds of 1:3 for synonymous and non-synonymous mutations as a neutral

300    expectation.

301    Skew of intronic base composition

302    Gene expression data from ovaries and testes tissues were obtained from FlyAtlas2 (34),

303    representing 16,781 genes. FPKM gene expression values were grouped into 40 bins, separately

304    for ovaries and testes with the mltools::bin_data (50) (version 0.3.5; binType="quantile") R

305    function. Since alternative splicing may generate ambiguous signals, 7 bases from the 5' end and

306    35 bases from the 3' end were removed from introns to exclude genomic regions containing

307    splicing sequences as recommended in (15). AT (CG) skews were then calculated as the number

308    of T (C) on the transcribed strand over the total number of A and T (C and G) bases. For each

309    tissue  and  type  of  skew,  we  fitted  a  Generalized  Linear  Model  (51)  using  the

310    glm(cbind(#transcribed, #total-#transcribed), family="binomial") R function, and reported the Wald

311    test     p-values     corresponding     to     the     binned     gene     expression     covariate.

312

313    Mutational strand bias

314    We  restricted  our  analysis  to  unambiguous  exons  and  introns  and  excluded  annotations

315    overlapping  with  other  genes  located  on  a  different  strand  using  the  BEDtools  intersect  -s

316    command  (a  GTF  with  the  final  annotation  can  be  found  in  the  Dryad  repository).

317    We used the Bioconductor MutationalPatterns package (52) (version 1.12.0) to count the different

318    mutation types on the transcribed and non-transcribed strands. Our first approach was to estimate

319    the expected mutation rate from the base composition on the transcribed and non-transcribed

320    strand of genes with at least one mutation. Since without strand bias a ratio of 1 is expected, we

321    calculated its significance and 95% confidence intervals using the poisson.test R function.

322    In the second approach, we accounted for the impact of gene expression intensity on base

323    composition. For each of the 40 expression bins, we randomly sampled the same number of

324    genes as observed being mutated in our SNPs set and calculated the expected strand bias from

325    the sample.

326    We repeated the approaches for the expected intronic and exonic biases, using exclusively either

327    intronic or exonic sequences. The p-values were corrected for multiple testing using the

328    Benjamini-Hochberg procedure.

329    In order to take the 5' and 3' flanking bases of the A[T>C]N mutations into account in the null

330    hypothesis, we adapted a permutation procedure from (32) to test for strand bias in exons and

331    introns (Supplementary Figure 3). Briefly, we obtained the frequency of mutations for each of the

332    4 A[T>C]N contexts (triplets) genome-wide and rescaled the frequencies to sum up to 1. In

333    parallel, we used the GATK tool CallableLoci (53) to obtain the callable sites per line and the

334    BEDtools suit (maskfasta and getfasta commands) to mask the reference for non-callable sites.

335    For both strands, we then retained as a sampling pool the number of callable triplets in the

336    mutated genes for exons, introns, summed over each line, and multiplied it with the rescaled

337    frequency to weight the sampling according to the genome-wide prevalence of triplets. Finally, we

338    redistributed the observed number of mutations on the transcribed and non-transcribed strand

339    separately 10,000 times to get the expected number of mutations on the transcribed strand in

340    introns and exons. The p-values were calculated as the number of times the sampled value was

341    higher than the observed one divided by 10,000.

342

343    <u>Gene expression analysis for C>A and A[T>C]N mutations</u>

344    Gene expression differences between genes containing C>A and A[T>C]N mutations on different

345    strands were tested with either one-sided (intron) or two-sided (exon) Wilcoxon rank-sum test on

346    the FPKM scale. We used a one-sided test for intronic sequences because the strand bias

347    predicts the direction of gene expression difference. The p-values were corrected for multiple

348    testing using the Benjamini-Hochberg procedure.

349

350    <u>Mutagenic effect of transcription</u>

351    To test if transcription is mutagenic, we performed a randomization test similarly as (54)

352    (Supplementary Figure 4). We randomly picked 221 genes, corresponding to the number of C>A

353    intronic mutations overlapping with the FlyAtlas2 data (over 236) on the non-transcribed strand,

354    and computed the mean expression in ovaries and testes separately. The sampling was weighted

355    by the length of the introns. This was done 10,000 times. For each tissue, a p-value was calculated

356    as the number of times the randomly sampled mean values exceeded the observed mean divided

357    by 10,000.

358

359    Decreased exonic mutation rate for C>A and A[T>C]N mutations

360    We used a similar permutation procedure as described above in the mutational strand bias

361    subsection to test for reduced exonic mutation rates (Supplementary Figures 5, 6). We modified

362    the sampling pool of callable triplets to include the genome-wide exons and introns with strands

363    separated.

364

365    **Code and data availability**

366    The code (R and bash scripts) will be accessible in the following github repository: ***, available

367    upon publication.

368    The final set of SNPs and INDELs as well as the updated annotation and intermediate files can

369    be found from the following dryad repository: ***, available upon publication.

370    Raw reads will be available in the following SRA project: ***, available upon publication.

371

372    **Authors contribution**

373    L. T. performed experiments, V. N. performed sequencing,  L. T., C. B. analyzed the data,  L. T.,

374    C. B., V. N., C. S. wrote the paper, K.S. supervised the project and provided feedback,  L. T., C.

375    S. designed the study.

376 **Acknowledgments**

381

382

383 **References**

384 1. T. Lindahl, Instability and decay of the primary structure of DNA. *Nature* **362**, 709–715 (1993).

385 2. J. Nakamura, *et al.*, Highly sensitive apurinic/apyrimidinic site assay can detect spontaneous
386 and chemically induced depurination under physiological conditions. *Cancer Res.* **58**, 222–
387 225 (1998).

388 3. W. Wang, C. Walmacq, J. Chong, M. Kashlev, D. Wang, Structural basis of transcriptional
389 stalling and bypass of abasic DNA lesion by RNA polymerase II. *Proc. Natl. Acad. Sci. U. S.*
390 *A.* **115**, E2538–E2545 (2018).

391 4. J. A. Marteijn, H. Lans, W. Vermeulen, J. H. J. Hoeijmakers, Understanding nucleotide
392 excision repair and its roles in cancer and ageing. *Nat. Rev. Mol. Cell Biol.* **15**, 465–481
393 (2014).

394 5. P. C. Hanawalt, G. Spivak, Transcription-coupled DNA repair: two decades of progress and
395 surprises. *Nat. Rev. Mol. Cell Biol.* **9**, 958–970 (2008).

396 6. P. Green, *et al.*, Transcription-associated mutational asymmetry in mammalian evolution.
397 *Nat. Genet.* **33**, 514–517 (2003).

398 7. L. E. Kerry, *et al.*, Selective inhibition of RNA polymerase I transcription as a potential
399 approach to treat African trypanosomiasis. *PLoS Negl. Trop. Dis.* **11**, e0005432 (2017).

400 8. A. Zheleva, D. Michelot, Z. D. Zhelev, Sensitivity of alpha-amanitin to oxidation by a
401 lactoperoxidase-hydrogen peroxide system. *Toxicon* **38**, 1055–1063 (2000).

402 9. J. A. Eisen, P. C. Hanawalt, A phylogenomic study of DNA repair genes, proteins, and
403 processes. *Mutat. Res.* **435**, 171–213 (1999).

404 10. A. R. Lehmann, The xeroderma pigmentosum group D (XPD) gene: one gene, two functions,
405 three diseases. *Genes Dev.* **15**, 15–23 (2001).

406 11. J. J. Sekelsky, M. H. Brodsky, K. C. Burtis, DNA repair in Drosophila: insights from the
407 Drosophila genome sequence. *J. Cell Biol.* **150**, F31–6 (2000).

408   12. J. Sekelsky, DNA Repair in Drosophila: Mutagens, Models, and Missing Genes. *Genetics*
409        **205**, 471–490 (2017).

410   13. J. G. de Cock, *et al.*, Repair of UV-induced (6-4)photoproducts measured in individual genes
411        in the Drosophila embryonic Kc cell line. *Nucleic Acids Res.* **20**, 4789–4793 (1992).

412   14. P. J. van der Helm, E. C. Klink, P. H. Lohman, J. C. Eeken, The repair of UV-induced
413        cyclobutane pyrimidine dimers in the individual genes Gart, Notch and white from isolated
414        brain tissue of Drosophila melanogaster. *Mutat. Res.* **383**, 113–124 (1997).

415   15. M. Touchon, A. Arneodo, Y. d'Aubenton-Carafa, C. Thermes, Transcription-coupled and
416        splicing-coupled strand asymmetries in eukaryotic genomes. *Nucleic Acids Res.* **32**, 4969–
417        4978 (2004).

418   16. J. Bergman, A. J. Betancourt, C. Vogl, Transcription-Associated Compositional Skews in
419        Drosophila Genes. *Genome Biol. Evol.* **10**, 269–275 (2018).

420   17. N. Deger, Y. Yang, L. A. Lindsey-Boltz, A. Sancar, C. P. Selby, Drosophila, which lacks
421        canonical transcription-coupled repair proteins, performs transcription-coupled repair. *J. Biol.*
422        *Chem.* **294**, 18092–18098 (2019).

423   18. I. Mellon, G. N. Champe, Products of DNA mismatch repair genes mutS and mutL are
424        required for transcription-coupled nucleotide-excision repair of the lactose operon in
425        Escherichia coli. *Proc. Natl. Acad. Sci. U. S. A.* **93**, 1292–1297 (1996).

426   19. P. Bertrand, D. X. Tishkoff, N. Filosi, R. Dasgupta, R. D. Kolodner, Physical interaction
427        between components of DNA mismatch repair and nucleotide excision repair. *Proc. Natl.*
428        *Acad. Sci. U. S. A.* **95**, 14278–14283 (1998).

429   20. K. S. Sweder, *et al.*, Mismatch repair mutants in yeast are not defective in transcription-
430        coupled DNA repair of UV-induced DNA damage. *Genetics* **143**, 1127–1135 (1996).

431   21. J. Zhao, A. Jain, R. R. Iyer, P. L. Modrich, K. M. Vasquez, Mismatch repair and nucleotide
432        excision repair proteins cooperate in the recognition of DNA interstrand crosslinks. *Nucleic*
433        *Acids Res.* **37**, 4420–4429 (2009).

434   22. A. Bellacosa, Functional interactions and signaling properties of mammalian DNA mismatch
435        repair proteins. *Cell Death Differ.* **8**, 1076–1092 (2001).

436   23. K. Kobayashi, P. Karran, S. Oda, K. Yanaga, The involvement of mismatch repair in
437        transcription coupled nucleotide excision repair. *Hum. Cell* **18**, 103–115 (2005).

438   24. G. Rossetti, *et al.*, The structural impact of DNA mismatches. *Nucleic Acids Res.* **43**, 4309–
439        4321 (2015).

440   25. D. Mu, *et al.*, Recognition and repair of compound DNA lesions (base damage and mismatch)
441        by human mismatch repair and excision repair systems. *Mol. Cell. Biol.* **17**, 760–769 (1997).

442   26. T. L. DeWeese, *et al.*, Mouse embryonic stem cells carrying one or two defective Msh2 alleles
443        respond abnormally to oxidative stress inflicted by low-level radiation. *Proc. Natl. Acad. Sci.*
444        *U. S. A.* **95**, 11915–11920 (1998).

445   27. C. Colussi, *et al.*, The mammalian mismatch repair pathway removes DNA 8-oxodGMP

incorporated from the oxidized dNTP pool. *Curr. Biol.* **12**, 912–918 (2002).

28. V. Grazielle-Silva, *et al.*, Distinct Phenotypes Caused by Mutation of MSH2 in Trypanosome Insect and Mammalian Life Cycle Forms Are Associated with Parasite Adaptation to Oxidative Stress. *PLoS Negl. Trop. Dis.* **9**, e0003870 (2015).

29. F. Le Page, A. Klungland, D. E. Barnes, A. Sarasin, S. Boiteux, Transcription coupled repair of 8-oxoguanine in murine cells: the ogg1 protein is required for repair in nontranscribed sequences but not in transcribed sequences. *Proc. Natl. Acad. Sci. U. S. A.* **97**, 8397–8402 (2000).

30. Y. Ding, A. Berrocal, T. Morita, K. D. Longden, D. L. Stern, Natural courtship song variation caused by an intronic retroelement in an ion channel gene. *Nature* **536**, 329–332 (2016).

31. Z. J. Assaf, S. Tilk, J. Park, M. L. Siegal, D. A. Petrov, Deep sequencing of natural and experimental populations of Drosophila melanogaster reveals biases in the spectrum of new mutations. *Genome Res.* **27**, 1988–2000 (2017).

32. J. Frigola, *et al.*, Reduced mutation rate in exons due to differential mismatch repair. *Nat. Genet.* **49**, 1684–1692 (2017).

33. A. R. Poetsch, S. J. Boulton, N. M. Luscombe, Genomic landscape of oxidative DNA damage and repair reveals regioselective protection from mutagenesis. *Genome Biol.* **19**, 215 (2018).

34. D. P. Leader, S. A. Krause, A. Pandit, S. A. Davies, J. A. T. Dow, FlyAtlas 2: a new version of the Drosophila melanogaster expression atlas with RNA-Seq, miRNA-Seq and sex-specific data. *Nucleic Acids Res.* **46**, D809–D815 (2018).

35. S. A. Lujan, *et al.*, Heterogeneous polymerase fidelity and mismatch repair bias genome variation and composition. *Genome Res.* **24**, 1751–1764 (2014).

36. S. A. Miller, D. D. Dykes, H. F. Polesky, A simple salting out procedure for extracting DNA from human nucleated cells. *Nucleic Acids Res.* **16**, 1215 (1988).

37. D. Gómez-Sánchez, C. Schlötterer, ReadTools: A universal toolkit for handling sequence data from different sequencing platforms. *Mol. Ecol. Resour.* **18**, 676–680 (2018).

38. B. Bushnell, BBMap. *sourceforge*.

39. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).

40. R. V. Pandey, C. Schlötterer, DistMap: a toolkit for distributed short read mapping on a Hadoop cluster. *PLoS One* **8**, e72614 (2013).

41. H. Li, *et al.*, The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).

42. M. R. Breese, Y. Liu, NGSUtils: a software suite for analyzing and manipulating next-generation sequencing datasets. *Bioinformatics* **29**, 494–496 (2013).

43. P. Danecek, *et al.*, The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).

483   44. H. Li, A statistical framework for SNP calling, mutation discovery, association mapping and
484       population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–
485       2993 (2011).

486   45. A. Tan, G. R. Abecasis, H. M. Kang, Unified representation of genetic variants. *Bioinformatics*
487       **31**, 2202–2204 (2015).

488   46. A. R. Quinlan, I. M. Hall, BEDTools: a flexible suite of utilities for comparing genomic features.
489       *Bioinformatics* **26**, 841–842 (2010).

490   47. H. Li, Tabix: fast retrieval of sequence features from generic TAB-delimited files.
491       *Bioinformatics* **27**, 718–719 (2011).

492   48. R Core Team, R: A Language and Environment for Statistical Computing (2018).

493   49. P. Cingolani, *et al.*, A program for annotating and predicting the effects of single nucleotide
494       polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster strain w1118; iso-
495       2; iso-3. *Fly* **6**, 80–92 (2012).

496   50. B. Gorman, mltools: Machine Learning Tools.

497   51. P. McCullagh, *Generalized linear models* (Routledge, 2019).

498   52. F. Blokzijl, R. Janssen, R. van Boxtel, E. Cuppen, MutationalPatterns: comprehensive
499       genome-wide analysis of mutational processes. *Genome Med.* **10**, 33 (2018).

500   53. A. McKenna, *et al.*, The Genome Analysis Toolkit: a MapReduce framework for analyzing
501       next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).

502   54. C. Park, W. Qian, J. Zhang, Genomic evidence for elevated mutation rates in highly
503       expressed genes. *EMBO Rep.* **13**, 1123–1129 (2012).

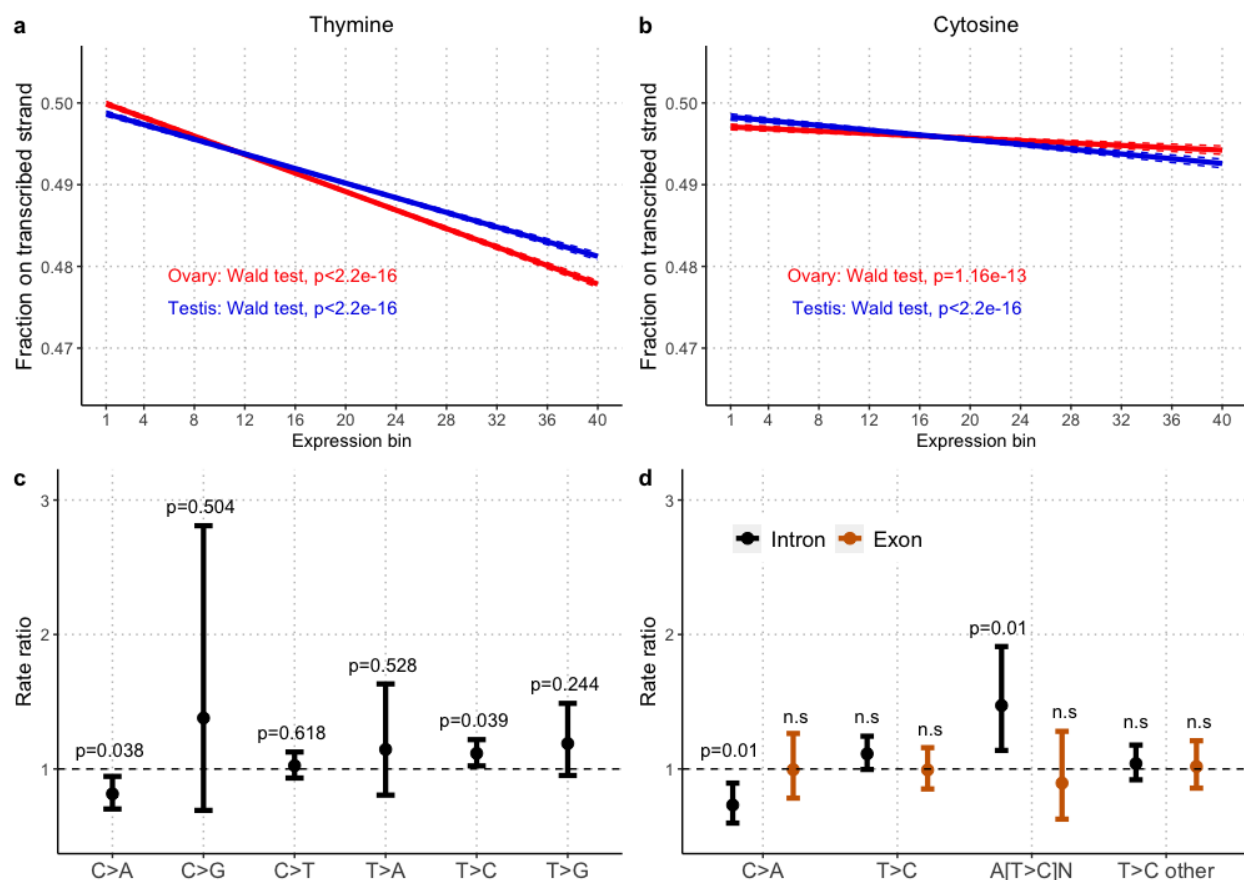504   55. J. T. Robinson, *et al.*, Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).

Figure 1. The interplay between transcription-associated base skews and mutational bias within genes. *Top*. Correlation between binned gene expression of 16,781 genes in ovaries (red) and testes (blue) and the fraction of a) thymines and b) cytosines on the transcribed strand. The regression line and its confidence interval are in plain and dotted lines respectively using the Generalized Linear Modeling framework. The p-values (Wald tests) correspond to the binned gene expression covariate. The intercept line of 1 indicates the absence of differences in rates of mutations on the transcribed and non-transcribed strands. The significance threshold is set to 5%. *Bottom*. Estimated rate ratios (RR) for different substitution types on the transcribed and non-transcribed strands in c) genes and partitioned for d) introns (black) and exons (orange). Adjusted p-values using the Benjamini-Hochberg procedure are reported and 95% Poisson confidence intervals are indicated by whiskers.
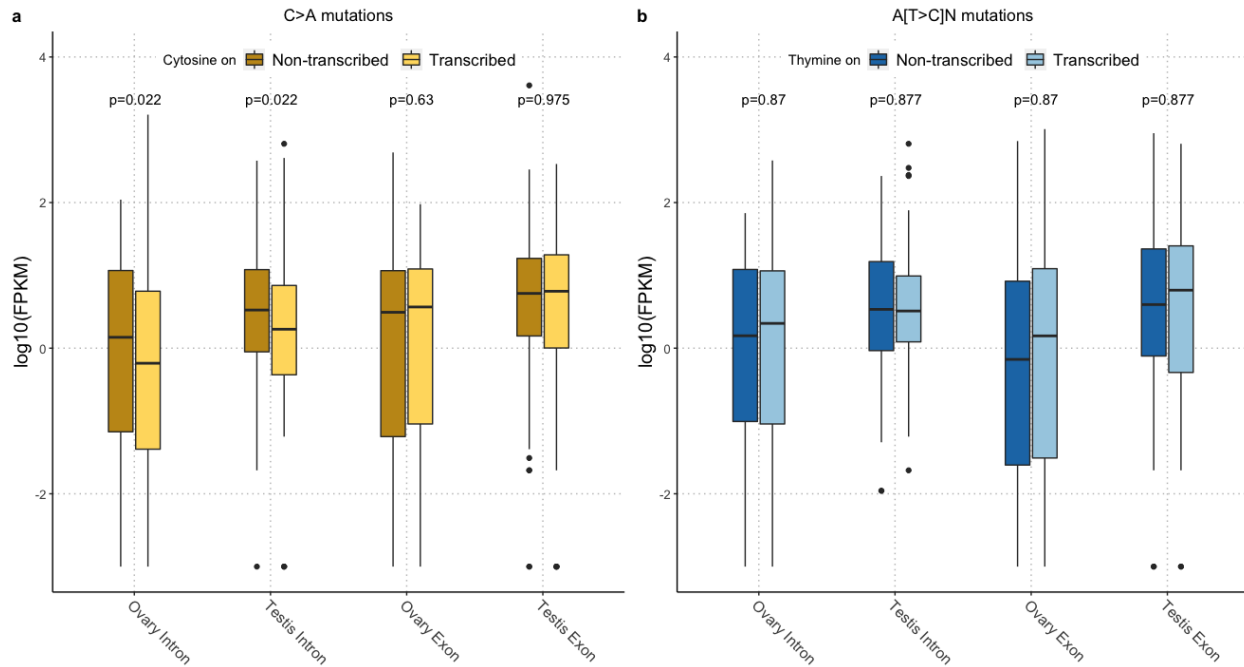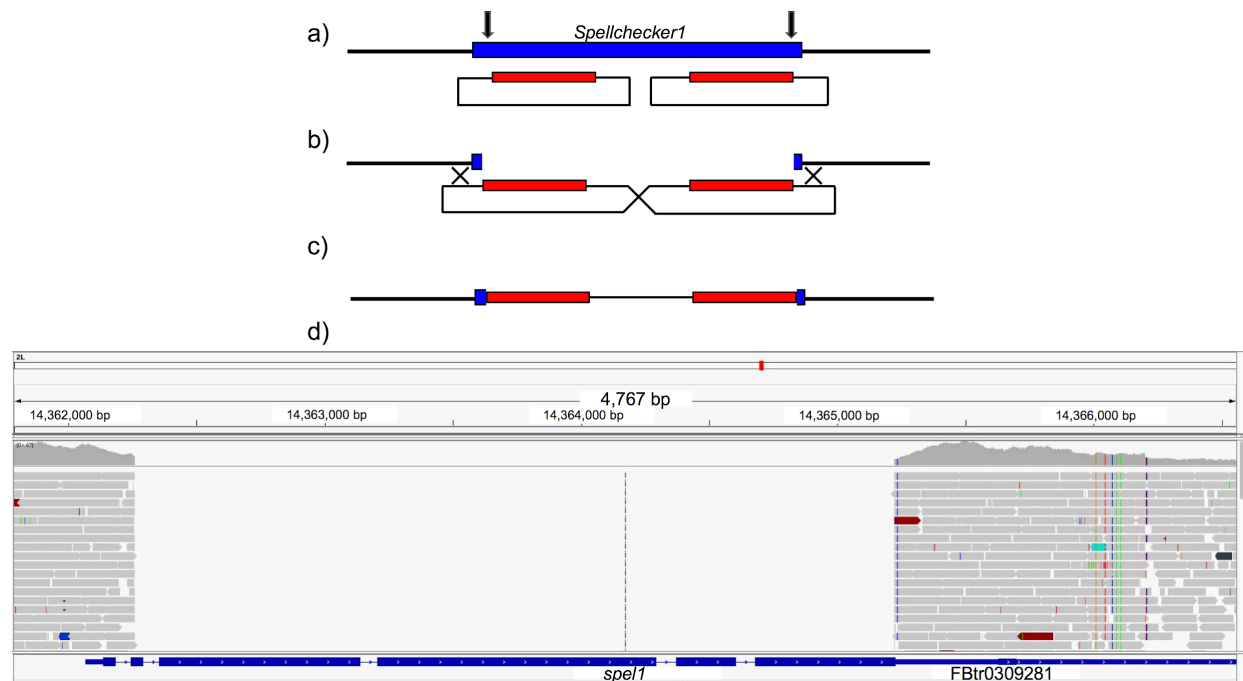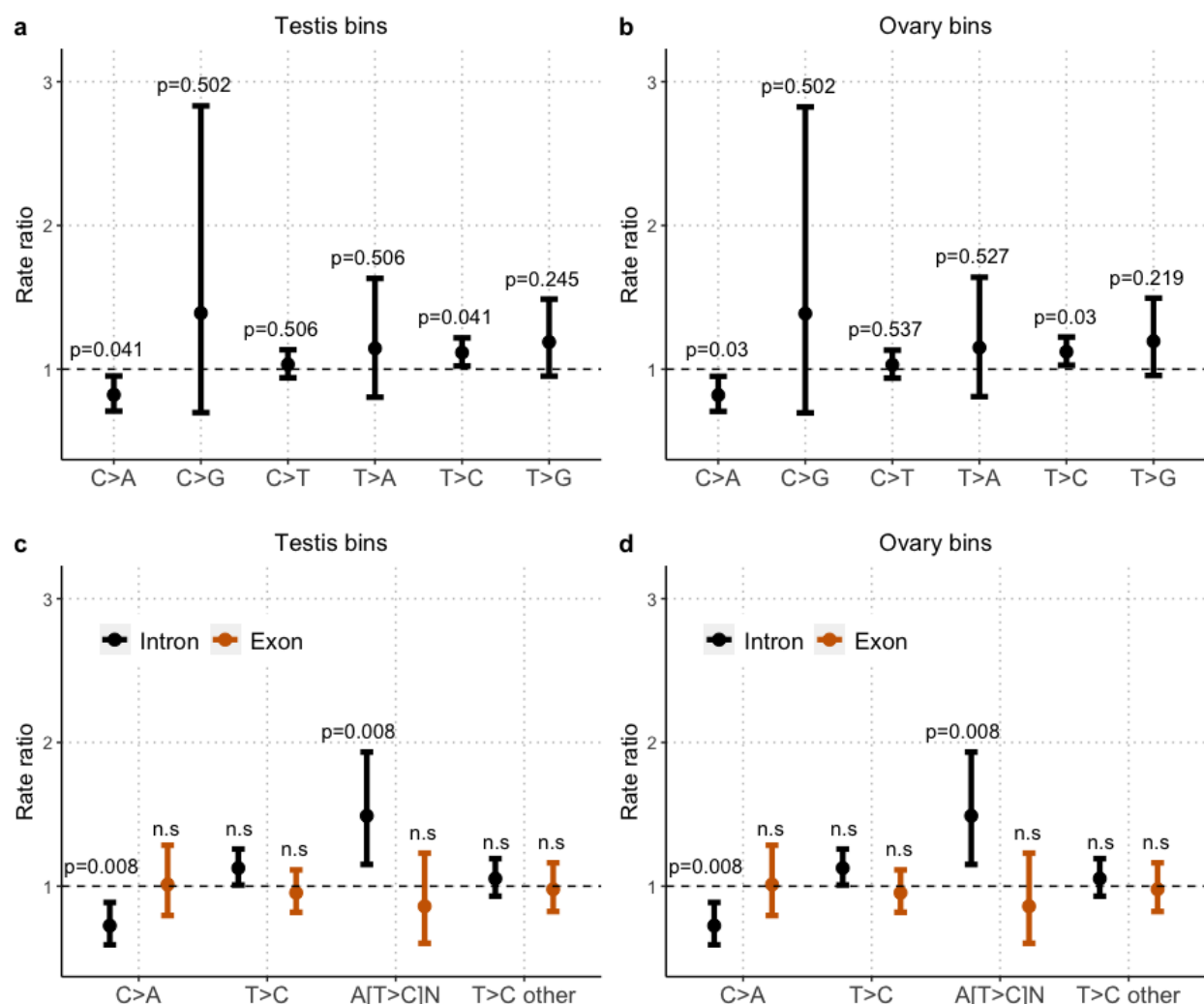
Figure 2. Introns with C>A mutations on the transcribed strand have lower expression levels. Boxplots of gene expression levels (log-10 transformed FPKM + $c$, with $c = 0.001$ to include genes not being expressed) of expressed mutated genes in ovaries and testes with a) C>A mutations and b) A[T>C]N mutations on the transcribed (light) and non-transcribed (dark) strands in exons and introns. The p-values are from one-sided (introns) or two-sided (exons) Wilcoxon rank-sum tests done on all genes and were adjusted per mutation type with the Benjamini-Hochberg procedure.
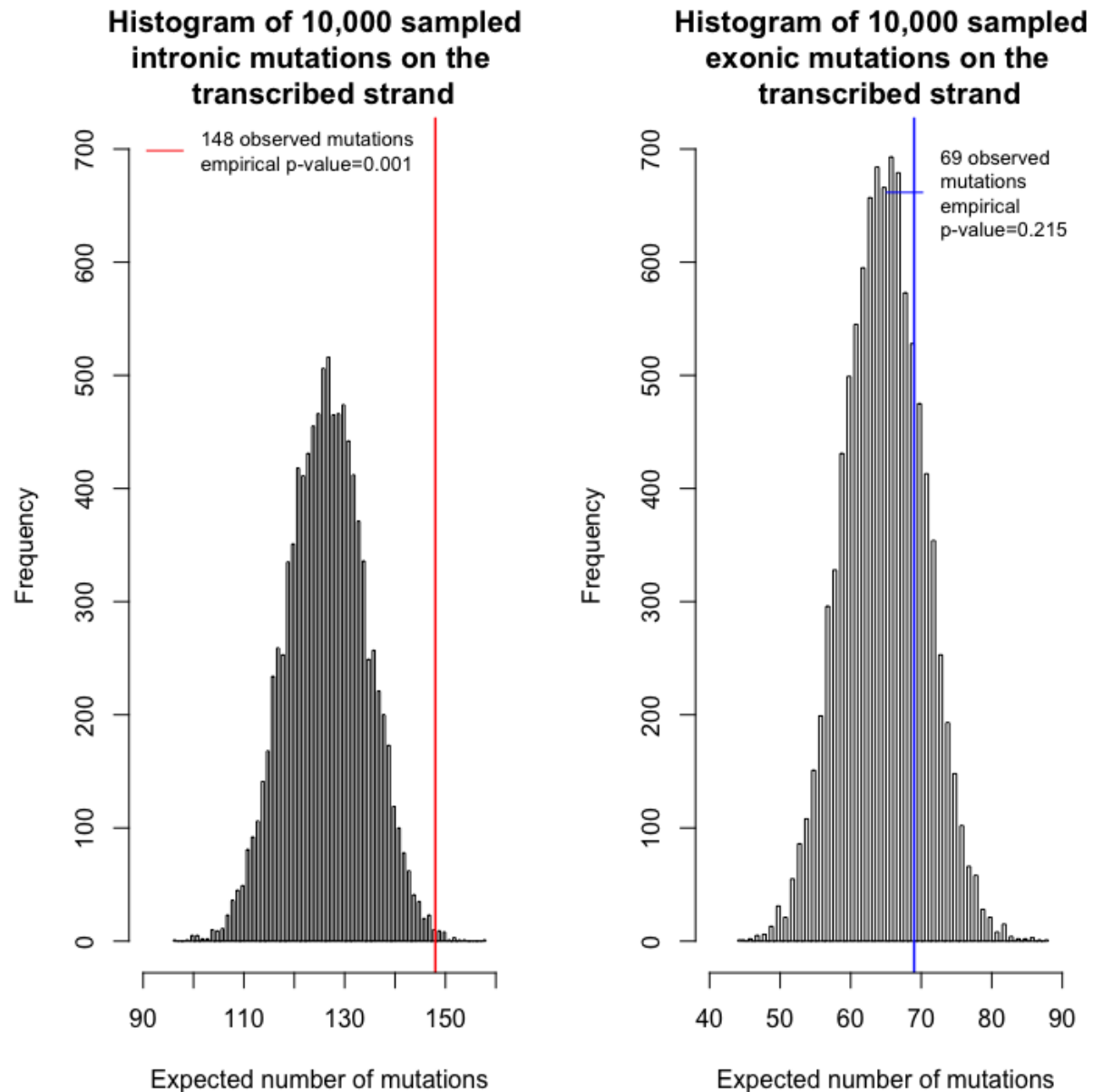
Supplementary Figure 1. A schematic overview showing the double insertion of the template plasmid and the short Illumina-read based confirmation of the *spel1* deletion. a) Two gRNAs targeting the gene are indicated by black arrows b) the resulting double-stranded break is repaired by two template plasmids which recombine with each other c) the resulting allele contains the backbone of the plasmid flanked by dsRed cassettes d) a screenshot of the short read coverage at the *spel1* locus visualized by IGV (55).
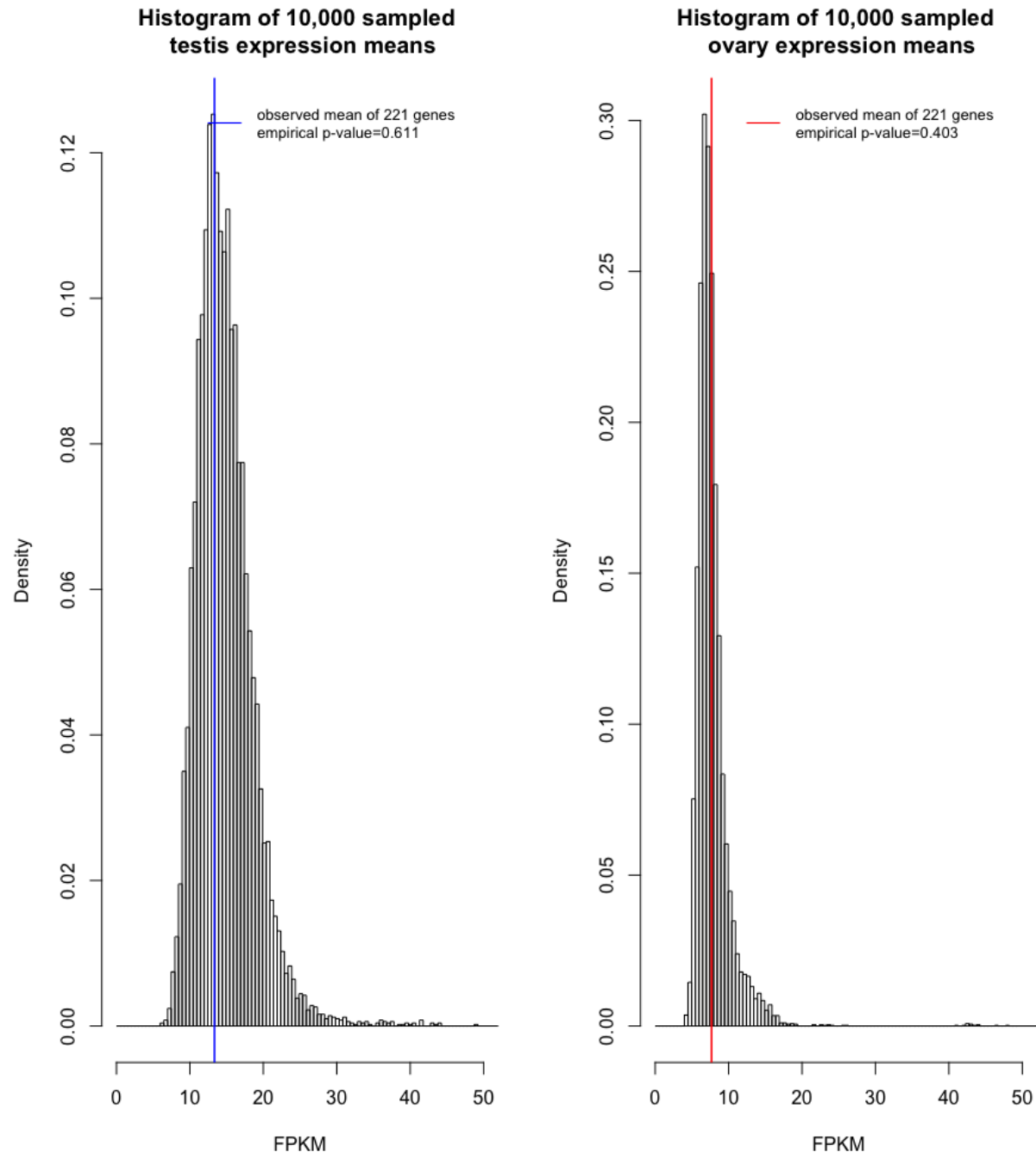
532
533　Supplementary Figure 2. Mutational bias within genes tested using the gene expression sampling
534　scheme. Estimated rate ratios (RR) on the transcribed strand by the non-transcribed strands (dot)
535　in genes and partitioned in d) introns (black) and exons (orange) based on testis (a, c) and b)
536　ovary (b, d) expression bins. Adjusted p-values using the Benjamini-Hochberg procedure are
537　reported and 95% Poisson confidence intervals are represented by segments. The intercept line
538　of 1 indicates the absence of differences in rates of mutations on the transcribed and non-
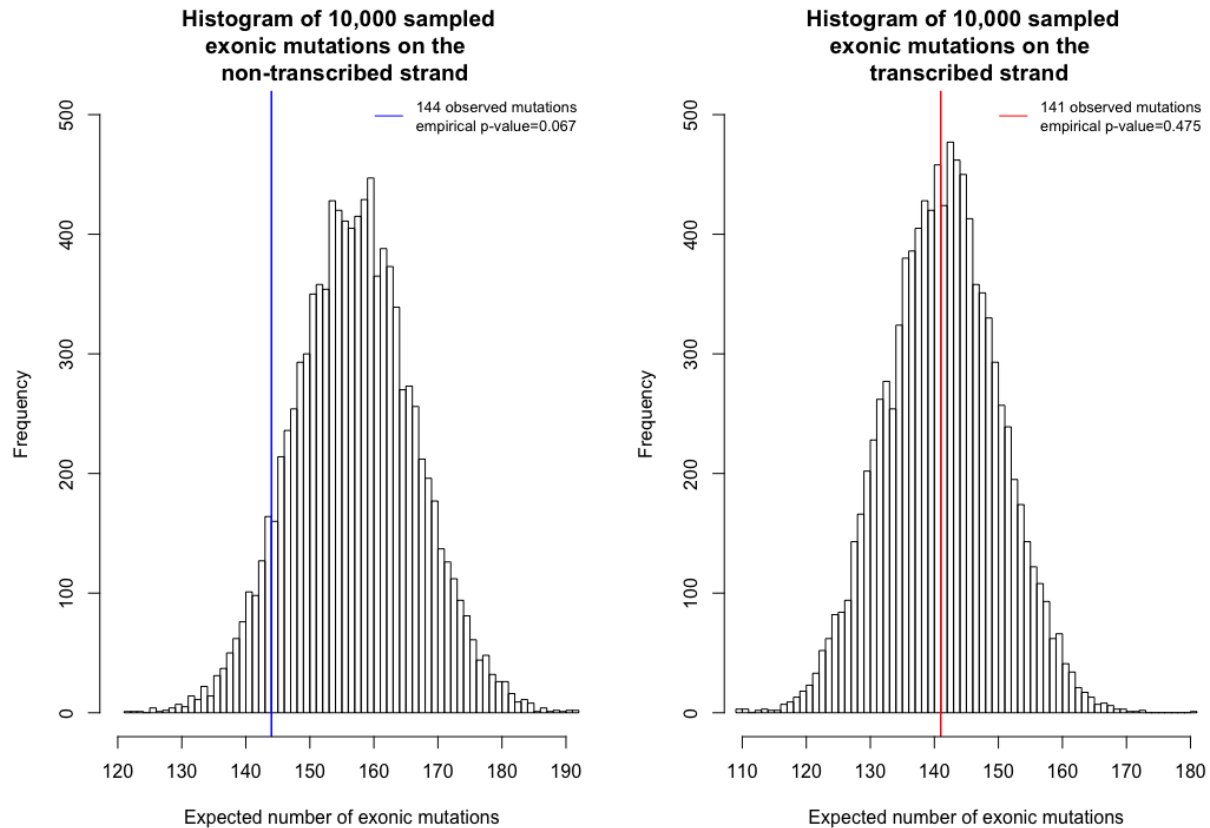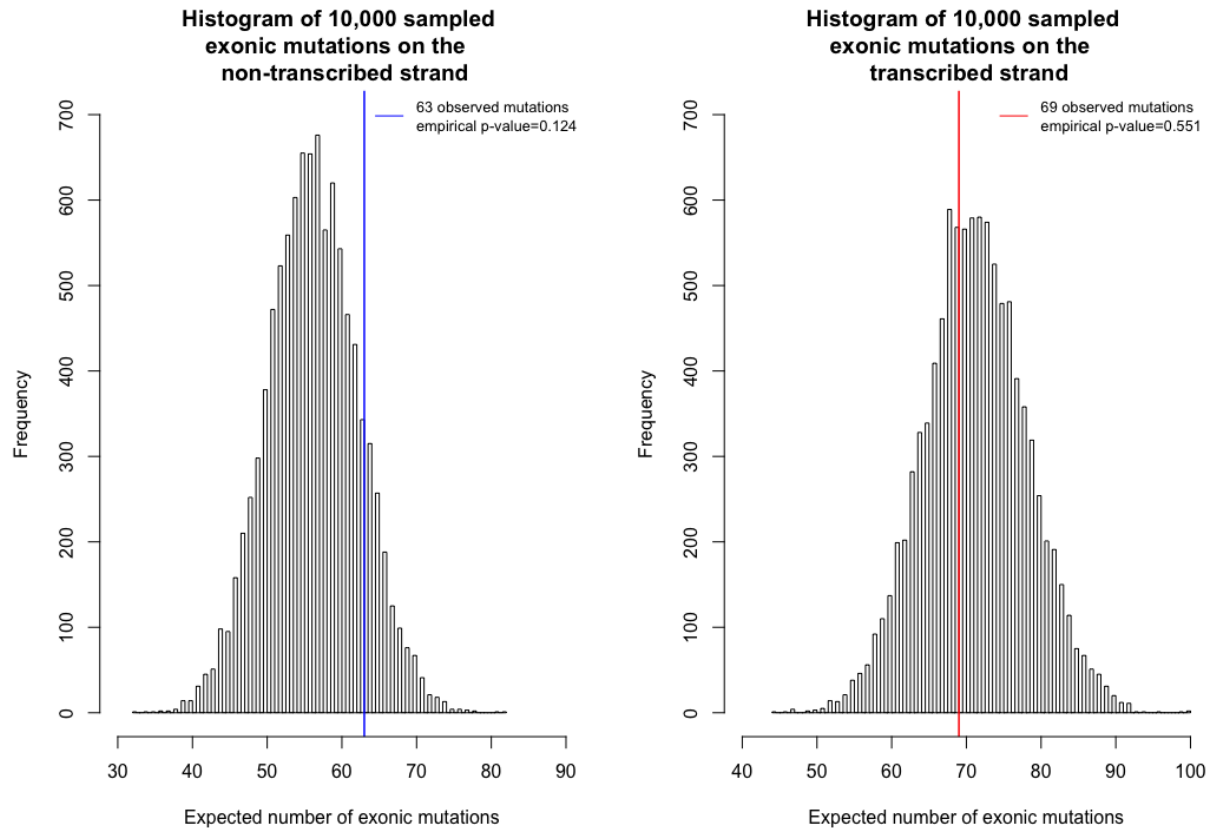539　transcribed strands. The significance threshold is set to 5%.

540

Supplementary Figure 3. Permutation test to test strand bias for A[T>C]N mutations while taking the composition of the 5′ and 3′ flanking bases into account. The histograms show the expected number of A[T>C]N on the transcribed strand in introns (left) and exons (right). The red line shows the observed mutations on the transcribed strand in introns and the blue line exons.

545
546  Supplementary Figure 4. Randomization test to rule out the mutagenic effect of transcription.
547  Histograms show the frequency distribution of the mean expression level of 221 randomly
548  sampled genes in testes (left) and ovaries (right) from 10,000 trials. Vertical lines show the
549  observed mean expression of genes with C>A mutations located on the non-transcribed strand.
550  The binning is the same between both histograms.

551
552 Supplementary Figure 5. Permutation test to test for reduced exonic C>A mutation rate.
553 Histograms show the sampled number of exonic mutations when the expectation was calculated
554 based on the non-transcribed strand (left) and transcribed strand (right). Vertical lines show the
555 observed number for C>A mutations when the cytosine was located on the non-transcribed strand
556 (blue) and on the transcribed strand (red).

557
558    Supplementary Figure 6. Permutation test to test for reduced exonic A[T>C]N mutation rate.
559    Histograms show the sampled number of exonic mutations when the expectation was calculated
560    based on the non-transcribed strand (left) and transcribed strand (right). Vertical lines show the
561    observed number for A[T>C]N mutations when the thymine was located on the non-transcribed
562    strand (blue) and on the transcribed strand (red).

563    Supplementary Table 1. Observed SNV counts in genes. Raw p-values were obtained using a
564    Poisson test and corrected with a Benjamini-Hochberg procedure.
565

| Type | Transcribed | Non-transcribed | Total | Rate ratio (95% CI) | Adjusted p-value |
|------|-------------|-----------------|-------|---------------------|------------------|
| C>A | 329 | 405 | 734 | 0.815 (0.702-0.945) | 0.0376 |
| C>G | 22 | 16 | 38 | 1.379 (0.692-2.809) | 0.5041 |
| C>T | 896 | 877 | 1,773 | 1.024 (0.933-1.126) | 0.6180 |
| T>A | 72 | 63 | 135 | 1.145 (0.805-1.632) | 0.5279 |
| T>C | 1087 | 976 | 2,063 | 1.116 (1.023-1.218) | 0.0385 |
| T>G | 178 | 150 | 328 | 1.189 (0.951-1.488) | 0.2444 |

566

567 Supplementary Table 2. Observed SNV counts in introns and exons. Raw p-values were obtained
568 using a Poisson test and corrected using the Benjamini-Hochberg procedure.
569

| Type | Transcribed | Non-transcribed | Total | Rate ratio (95% CI) | Adjusted p-value | Annotation |
|---|---|---|---|---|---|---|
| C>A | 171 | 236 | 407 | 0.732 (0.597-0.895) | 0.00982 | Intron |
| T>C | 667 | 620 | 1287 | 1.112 (0.996-1.243) | 0.15454 | Intron |
| A[T>C]N | 148 | 104 | 252 | 1.472 (1.138-1.910) | 0.00982 | Intron |
| T>C Other | 519 | 516 | 1,035 | 1.040 (0.919-1.178) | 0.86566 | Intron |
| C>A | 141 | 144 | 285 | 0.994 (0.783-1.264 | 1 | Exon |
| T>C | 373 | 307 | 680 | 0.993 (0.851-1.158) | 1 | Exon |
| A[T>C]N | 69 | 63 | 132 | 0.894 (0.627-1.280) | 0.86566 | Exon |
| T>C Other | 304 | 244 | 548 | 1.018 (0.857-1.210) | 1 | Exon |

570

571     Supplementary Table 3. Primers used to generate the *sple1* null mutant.

572

| Name | Sequence | Description |
|------|----------|-------------|
| LT1 | CTGAATATGGGTAAGCTGATAAGC | Left gRNA cut site sequencing For |
| LT2 | GTGTGATCAAAGAACCTCACTGTAGT | Left gRNA cut site sequencing Rev |
| LT3 | CTGATGACTTTACTCTGCTGTATCAAG | Right gRNA cut site sequencing For |
| LT4 | GATAAGTACGTAGAACAACTGCCTCTT | Right gRNA cut site sequencing Rev |
| LT39 | TATATAGGAAAGATATCCGGGTGAACTTCGTTCCATGCGAAACTCGGCGGTTTTAGAGCTAGAAATAGCAAG | Left gRNA primer |
| LT40 | ATTTTAACTTGCTATTTCTAGCTCTAAAACCCATTTCAAGCCTCACTAGCGACGTTAAATTGAAAATAGGT | Right gRNA primer |
| LT22 | GACACAGCGCGTACGTCCTTCG | Sequencing primer for pCFD4 |
| LT41 | CACACCACGTCTCAGGACCCTCATCAGTCTGGATCTGTGCTC | Left homology arm For |
| LT42 | CACACCACGTCTCACTGGACGCGCATTTGTGTCTGCAAAC | Left homology arm Rev |
| LT43 | CACACCACGTCTCATGTTGGCTTGAAATGGACGTAGGGTC | Right homology arm For |
| LT44 | CACACCACGTCTCAGCATTCTGAAGAGCGTGATGTGGAATGTC | Right homology arm Rev |
| LT58 | ACGGAGAAGGCGGAAATTGTG | Targeting check left For |
| LT26 | GGATGGGACAAGTCGCCATG | Targeting check left Rev/Template colony PCR left Rev |
| LT25 | CGATTAAGTTGGGTAACGCCAGG | Template colony PCR left For |
| LT27 | TGTGTGGAATTGTGAGCGGATAAC | Targeting check right For/Template colony PCR right For |
| LT59 | TTTGGATGCTGTTAAGCGTTGC | Targeting check right Rev |
| LT28 | TGTGTGGAATTGTGAGCGGATAAC | Template colony PCR right Rev |
| LT101 | CTCTTGATCCGGCAAACAAACC | Backbone check For |
| LT102 | GGGAGTCAGGCAACTATGGATG | Backbone check Rev |

573