

# Animal, fungi, and plant genome sequences harbour different non-canonical splice sites

**Katharina Frey<sup>1</sup>, Boas Pucker<sup>1,\*</sup>**

<sup>1</sup>Genetics and Genomics of Plants, Center for Biotechnology (CeBiTec),  
Bielefeld University, Bielefeld, Germany

\*Corresponding author

Email addresses:

BP: [bpucker@cebitec.uni-bielefeld.de](mailto:bpucker@cebitec.uni-bielefeld.de)

KF: [katharina.frey@uni-bielefeld.de](mailto:katharina.frey@uni-bielefeld.de)

ORCIDs:

<sup>1</sup> BP: <https://orcid.org/0000-0002-3321-7471>

<sup>2</sup> KF: <https://orcid.org/0000-0002-4022-8531>

<sup>3</sup>

## 4 **Abstract**

5 Most protein encoding genes in eukaryotes contain introns which are inter-  
6 woven with exons. After transcription, introns need to be removed in order  
7 to generate the final mRNA which can be translated into an amino acid  
8 sequence by the ribosome. Precise excision of introns by the spliceosome  
9 requires conserved dinucleotides which mark the splice sites. However,  
10 there are variations of the highly conserved combination of GT at the 5'  
11 end and AG at the 3' end of an intron in the genome. GC-AG and AT-AC  
12 are two major non-canonical splice site combinations which are known for  
13 many years. During the last few years, various minor non-canonical splice  
14 site combinations were detected with all possible dinucleotide permuta-  
15 tions. Here we expand systematic investigations of non-canonical splice  
16 site combinations in plant genomes to all eukaryotes by analysing fungal  
17 and animal genome sequences. Comparisons of splice site combinations  
18 between these three kingdoms revealed several differences such as a sub-  
19 stantially increased CT-AC frequency in fungal genomes. In addition, high  
20 numbers of GA-AG splice site combinations were observed in two animal  
21 species. In depth investigation of splice site usage based on RNA-Seq  
22 read mappings indicates a generally higher flexibility of the 3' splice site  
23 compared to the 5' splice site.

## 24 Introduction

25 Splicing, the removal of introns after transcription, is an essential step dur-  
26 ing the generation of mature mRNAs in eukaryotes. This process allows  
27 variation which provides the basis for quick adaptation to changing con-  
28 ditions [1,2]. Alternative splicing, e.g. skipping exons, results in an enor-  
29 mous diversity of synthesized proteins and therefore substantially expands  
30 the diversity of products encoded in eukaryotic genomes [3–6]. The full  
31 range of functions as well as the evolutionary relevance of introns are still  
32 under discussion [7]. However, introns are energetically expensive for the  
33 cell to maintain as the transcription of introns costs time and energy and  
34 the removal of introns has to be exactly regulated [8]. Dinucleotides at  
35 both intron/exon borders mark the splice sites and are therefore highly  
36 conserved [9]. GT at the 5' end and AG at the 3' end of an intron form the  
37 canonical splice site combination on DNA level. More complexity arises  
38 through non-canonical splice site combinations, which deviate from the  
39 highly conserved canonical one. Besides the major non-canonical splice  
40 site combinations GC-AG and AT-AC, several minor non-canonical splice  
41 site combinations have been detected before [9, 10].

42  
43 Furthermore, the position of introns in homologous genes across organ-  
44 isms, which diverged 500-1500 million years ago, are not conserved [11].  
45 In addition, many intron sequences mutate at a higher rate due to hav-  
46 ing much less of an impact on an organism's reproductive fitness com-  
47 pared to a mutation located within an exon [12]. These factors, along with  
48 the existence of several non-cannonical splice sites, make the complete  
49 prediction of introns, even in non-complex organisms like yeast, almost  
50 impossible [13, 14]. Moreover, most introns which can be predicted com-  
51 putationally still lack experimental support [15].

52  
53 Splice sites are recognised during the splicing process by a complex of  
54 snRNAs and proteins, the spliceosome [16]. U2-spliceosome and U12-  
55 spliceosome are two subtypes of this complex which comprise slightly dif-

56 ferent proteins with equivalent functions [17–19]. Although the terminal  
57 dinucleotides are important for the splicing process, these splice sites are  
58 not sufficient to determine which spliceosome is processing the enclosed  
59 intron [20]. This demonstrates the complexity of the splicing process which  
60 involves additional signals present in the DNA. Even though multiple mech-  
61 anisms could explain the splicing process, the exact mechanism of non-  
62 canonical splicing is still not completely resolved [5].

63  
64 Branching reaction and exon ligation are the two major steps of splic-  
65 ing [21,22]. In the branching reaction, the 2'-hydroxyl group of the branch-  
66 point adenosine initiates an attack on the 5'-phosphate of the donor splice  
67 site [23,24]. This process leads to the formation of a lariat structure. Next,  
68 the exons are ligated and the intron is released through activity of the 3'-  
69 hydroxyl group of the 5'exon at the acceptor splice site [21].

70  
71 Previous in-depth analyses of non-canonical splice sites in fungi and an-  
72 imals were often focused on a single or a small number of species [9,  
73 25,26]. Several studies focused on canonical GT-AG splice sites but ne-  
74 glected non-canonical splice sites [27,28]. Our understanding of splice  
75 site combinations is more developed in plants compared to other king-  
76 doms [10,29–33]. Previous works reported 98 % GT-AG splice site com-  
77 binations in fungi [25], 98.7 % in plants [10] and 98.71 % in animals [9].  
78 Consequently, the proportion of non-canonical splice sites is around or be-  
79 low 2 % [9,10,25]. To the best of our knowledge, it is not known if the value  
80 reported for mammals is representative for all animals. The combined pro-  
81 portion of minor non-canonical splice sites is even lower e.g. 0.09 % in  
82 plants, but still exceeding the frequency of the major non-canonical AT-  
83 AC splice sites [10]. Despite this apparently low frequency, non-canonical  
84 splice site combinations have a substantial impact on gene products, es-  
85 pecially on exon-rich genes [10]. About 40 % of genes with 40 exons are  
86 affected (AdditionalFile 11).

87  
88 Consideration of non-canonical splice sites is important for gene predic-

89 tion approaches, because these sites cannot be identified *ab initio* [29].  
90 Moreover, as many human pathogenic mutations occur at the donor splice  
91 site [34], it is of great interest to understand the occurrence and usage of  
92 non-canonical splice sites. Therefore, several non-canonical splice sites  
93 containing AG at the acceptor site were investigated in human fibro-  
94 blasts [34]. Alongside this, fungi are interesting due to pathogenic proper-  
95 ties and importance in the food industry [35]. Since splicing leads to high  
96 protein diversity [3–6], the analysis of splicing in fungi is important with re-  
97 spect to biotechnological applications e.g. development of new products.

98  
99 In this study, a collection of annotated genome sequences from 130 fungi  
100 and 489 animal species was screened for canonical and non-canonical  
101 splice site combinations. RNA-Seq data sets were harnessed to identify  
102 biologically relevant and actually used splice sites. Non-canonical splice  
103 site combinations, which appeared at substantially higher frequency in a  
104 certain kingdom or species, were analysed in detail. As knowledge about  
105 splice sites in plants was available from previous investigations [10, 29],  
106 a comparison between splice sites in fungi, animals and plants was per-  
107 formed.

108

## 109 **Results and Discussion**

### 110 **Analysis of non-canonical splice sites**

111 In total, 64,756,412 and 2,302,340 splice site combinations in animals  
112 and fungi, respectively, were investigated based on annotated genome se-  
113 quences (AdditionalFile 1 and 2). The average frequency of the canonical  
114 splice site combination GT-AG is 98.3 % in animals and 98.7 % in fungi,  
115 respectively. These values exceed the 97.9 % previously reported for  
116 plants [10], thus indicating a generally higher frequency of non-canonical  
117 splice site combinations in plants. As previously speculated [10], a gen-

118 erally more complex splicing system in plants could be an adaptation to  
119 changing environments. Since most plants are not able to change their  
120 geographic location, the tolerance for unfavourable conditions should be  
121 stronger than in animals. The lower proportion of non-canonical splice  
122 sites in fungi compared to animals seems to contradict this hypothesis.  
123 However, the genome size and complexity needs to be taken into account  
124 here. The average animal genome is significantly larger than the average  
125 fungal genome (Mann-Whitney U-Test;  $p=5.64e-68$ ) (AdditionalFile 3).  
126 Average percentages of the most important splice site combinations were  
127 summarized per kingdom and over all analysed genomes (Table 1). The  
128 number of canonical and non-canonical splice site combinations per species  
129 was also summarized (AdditionalFile 4 and 5). A higher percentage of  
130 non-canonical splice sites was observed in animals in comparison to fungi.  
131 Several species strongly exceeded the average values for major and minor  
132 non-canonical splice sites. The fungal species *Meyerozyma guilliermondi*  
133 shows approximately 6.67 % major and 13.33 % minor non-canonical  
134 splice sites. *Eurytemora affinis* and *Oikopleura dioica* reveal approximately  
135 10 % minor non-canonical splice sites. In summary, the observed frequen-  
136 cies of canonical and major non-canonical splice site combinations are  
137 similar to the pattern previously reported for plants [10], but some essen-  
138 tial differences and exceptions were found in animals and fungi.

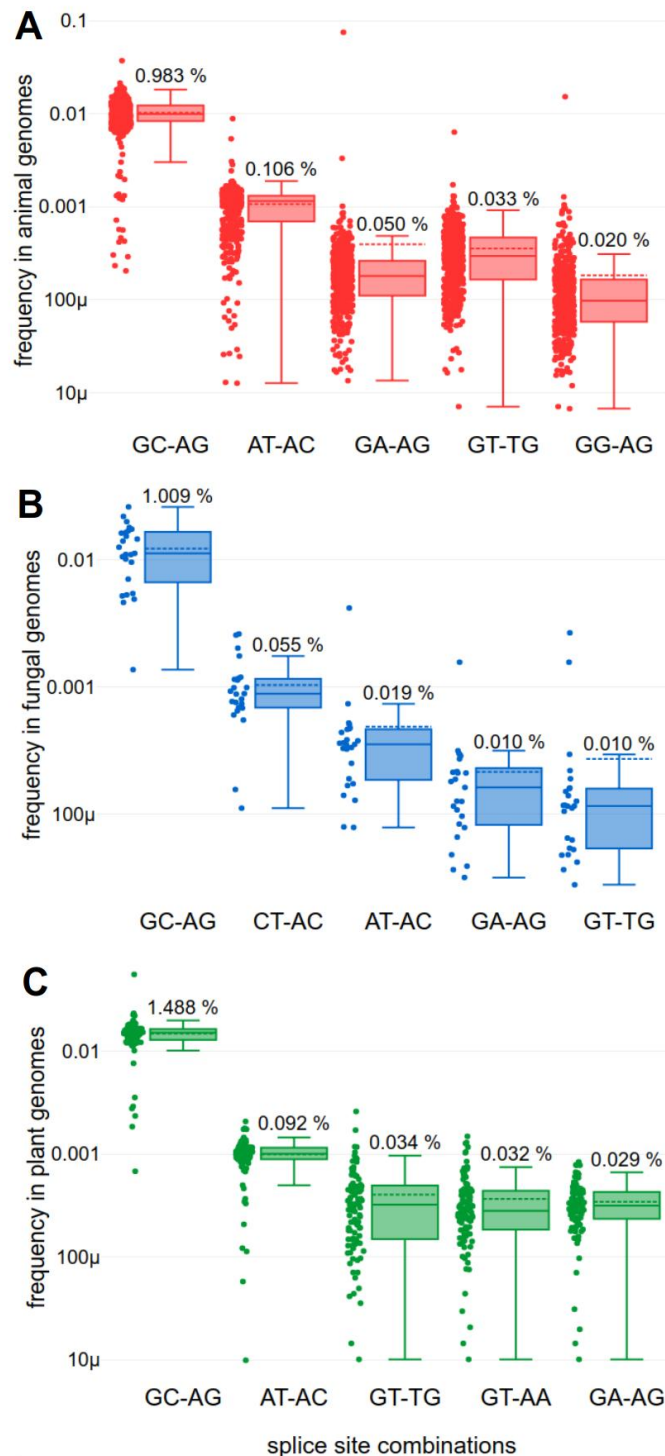
**Table 1: Splice site combination frequencies in animals, fungi, and plants.** Only the most frequent combinations are displayed here and all minor non-canonical splice site combinations are summarized as one group ("others"). A full list of all splice site combinations is available (AdditionalFile 6 and 7).

	<b>GT-AG</b>	<b>GC-AG</b>	<b>AT-AC</b>	<b>others</b>
<b>animals</b>	98.334 %	0.983 %	0.106 %	0.577 %
<b>fungi</b>	98.715 %	1.009 %	0.019 %	0.257 %
<b>plants</b>	97.886 %	1.488 %	0.092 %	0.534 %
<b>all</b>	98.265 %	1.074 %	0.101 %	0.560 %

139 Different properties of the genomes of all investigated species were anal-  
140 ysed to identify potential explanations for the splice site differences (Ad-  
141 ditionalFile 8 and 9). In fungi, the average number of introns per gene  
142 is 1.49 and the average GC content is 47.1 % ( $\pm 7.39$ ). In animals, each  
143 gene contains on average 6.95 introns and the average GC content is 39.4  
144 % ( $\pm 3.87$ ). This difference in the GC content could be associated with the  
145 much lower frequency of AT-AC splice site combinations and the higher fre-  
146 quency of CT-AC splice site combinations in fungi (Figure 1). CT-AC has a  
147 higher GC content than the AT rich AT-AC splice site combination. A gen-  
148 erally higher GC content could result in the higher GC content within splice  
149 site combinations due to the overall mutations rates in these species.

150 A comparison of the genome-wide GC content to the GC content of all  
151 splice sites revealed a weak correlation in the analysed fungi ( $r \approx 0.236$ ,  
152  $p \approx 0.008$ ). Species with a high genomic GC content tend to show a high  
153 GC content in the splice site combinations in the respective species. A  
154 similar correlation ( $r \approx 0.4$ ,  $p < 0.001$ ) was found in plant and animal species  
155 as well (AdditionalFile 10). Additionally, the GC content in fungal genomes  
156 is substantially exceeding the average GC content of plant and animal  
157 genomes.

158 The most frequent non-canonical splice site combinations show differ-  
159 ences between animals, fungi, and plants (Figure 1). In fungal species,  
160 the splice site CT-AC is more frequent than the splice site combination AT-  
161 AC. Regarding the splice site combination GA-AG in animals, two outliers  
162 are clearly visible: *Eurytemora affinis* and *Oikopleura dioica* show more  
163 GA-AG splice site combinations than GC-AG splice site combinations.



**Figure 1: Frequencies of non-canonical splice site combinations in animals, fungi, and plants.** The frequency of non-canonical splice site combinations across the 489 animal (red), 130 fungal (blue) and 121 plant (green) genomes is shown. Normalization of the absolute number of each splice site combination was performed per species based on the total number of splice sites. The frequency of the respective splice site combination of each species is shown on the left hand side and the percentage of the respective splice site combination on top of each box plot.



164 Despite overall similarity in the pattern of non-canonical splice site combi-  
165 nations between kingdoms, specific minor non-canonical splice sites were  
166 identified at much higher frequency in some fungal and animal species.  
167 First, RNA-Seq data was harnessed to validate these unexpected splice  
168 site combinations. Next, the frequencies of selected splice site combina-  
169 tions across all species of the respective kingdom were calculated. The  
170 correlation between the size of the incorporated RNA-Seq data sets and  
171 the number of supported splice sites was examined as well (AdditionalFile  
172 11). In animals, there is a correlation ( $r \approx 0.417$ ,  $p \approx 0.022$ ) between num-  
173 ber of supported splice sites and total number of sequenced nucleotides  
174 in RNA-Seq data. For fungi, no correlation between number of splice sites  
175 and size of the RNA-Seq data sets could be observed. It is important  
176 to note that the the number of available RNA-Seq data sets from fungi  
177 was substantially lower. Further, analysis of introns with canonical and  
178 non-canonical splice site combinations, respectively, revealed that a higher  
179 number of introns is associated with a higher proportion of non-canonical  
180 splice sites (AdditionalFile 12).

## 181 **High diversity of non-canonical splice sites in animals**

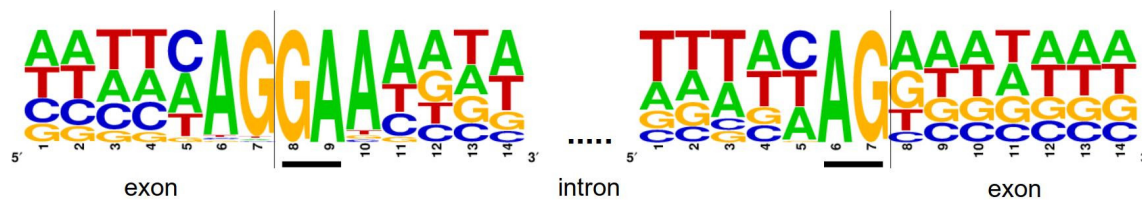
182 Kupfer *et al.* suggested that splicing may differ between fungi and ver-  
183 tebrates [25]. Our results indicate substantial differences in the diver-  
184 sity of splice site combinations other than GT-AG and GC-AG in fungi  
185 ( $H' \approx 0.0277$ ) and animals ( $H' \approx 0.0637$ ) (Kruskal-Wallis:  $p \approx 0.00000$ ). Be-  
186 sides the overall high proportion of minor non-canonical splice sites (Table  
187 1), differences between species are high (Figure 1). The slightly higher in-  
188 terquartile range of splice site combination frequencies in animal species  
189 and especially in plant species (Figure 1A and C), together with the rel-  
190 atively high frequency of "other" splice sites in animals and plants (Table  
191 1) suggest more variation of splice sites in the kingdoms of animals and  
192 plants compared to the investigated fungal species. Thus, the high di-  
193 versity of splice sites could be associated with the higher complexity of  
194 animal and plant genomes. In addition, the difference in prevalence be-

195 between the major non-canonical splice site combination GC-AG and minor  
196 non-canonical splice site combinations is smaller in animals compared to  
197 fungi and plants (Figure 1).

198

199 GA-AG is a frequent non-canonical splice site combination in some an-  
200 imal species. Two species, namely *Eurytemora affinis* and *Oikopleura*  
201 *dioica*, showed a much higher abundance of GA-AG splice site combi-  
202 nations compared to the other investigated species (Figure 1A). RNA-Seq  
203 reads support 5,795 (28.68 %) of all GA-AG splice site combinations of  
204 these species. In both species, the number of the GA-AG splice site com-  
205 bination exceeds the number of the major non-canonical splice site com-  
206 bination GC-AG.

207 For *Eurytemora affinis*, the high frequency of the GA-AG splice site combi-  
208 nations was described previously for 36 introns [36]. We quantified the pro-  
209 portion of GA-AG splice site combinations to 3.2 % (5,345) of all 166,392  
210 supported splice site combinations in this species. The donor splice site  
211 GA is flanked by highly conserved upstream AG and a downstream A (Fig-  
212 ure 2).



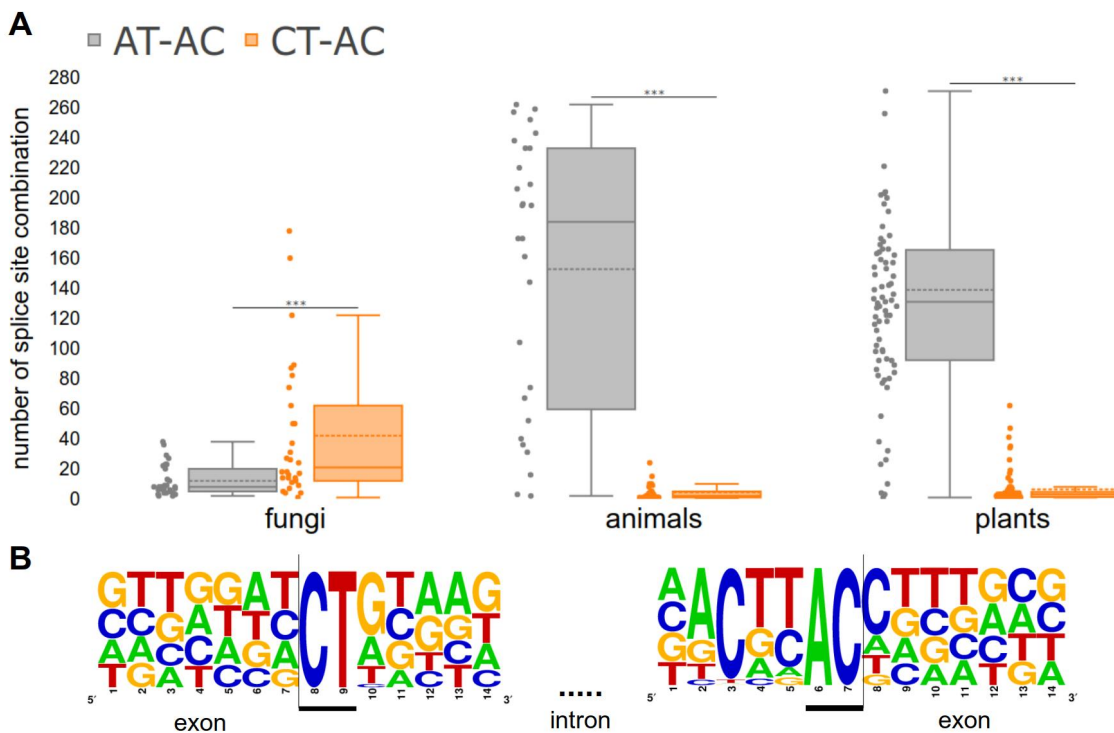
**Figure 2: Flanking positions of GA-AG splice site combinations in *Eurytemora affinis* and *Oikopleura dioica*.** All 5,795 supported splice site combinations of these two species were investigated. Seven exonic and seven intronic positions are displayed at the donor and acceptor splice sites. Underlined bases represent the terminal dinucleotides of the intron i.e. the donor and acceptor splice site.

213 Efficient splicing of the splice site combination GA-AG was detected in hu-  
214 man fibroblast growth factor receptor genes [37]. Further, it was suggested  
215 that this splicing event is, among other sequence properties, dependent on  
216 a canonical splice site six nucleotides upstream [37], which does not exist

217 in the species investigated here (Figure 2). An analysis of all five potential  
218 U1 snRNAs in this species did reveal one single nucleotide polymorphism  
219 in the binding site of the 5' splice site from C to T in one of these U1  
220 snRNAs. This could result in the binding of AG/GGAAGT or AGG/GAAGT  
221 instead of AG/GTAAGT. Although this would imply an elegant way for the  
222 splicing of GA-AG splice sites, the same variation was also detected in  
223 putative human U1 snRNAs. Therefore, another mechanism seems to be  
224 responsible for splicing of introns containing the GA-AG splice site combi-  
225 nation.

## 226 **CT-AC is a frequent splice site combination in fungi**

227 Although the general frequency pattern of fungal splice site combinations  
228 is similar to plants and animals, several fungal species displayed a high  
229 frequency of minor non-canonical CT-AC splice site combinations. This  
230 co-occurs with a lower frequency of AT-AC splice site combinations.  
231 Non-canonical splice sites in fungi were, so far, only described in stud-  
232 ies which focussed on a single or a few species. An analysis in the  
233 oomycota species *Phytophthora sojae*, which is a fungus-like microorgan-  
234 ism [38, 39], revealed 3.4 % non-canonical splice site combinations GC-  
235 AG and CT-AC [40]. Our findings indicate, that the minor non-canonical  
236 splice site combination CT-AC occurs with a significantly (Mann-Whitney  
237 U-Test;  $p \approx 0.00035$ ) higher frequency than the major non-canonical splice  
238 site combination AT-AC. In contrast, the frequency of AT-AC in animals  
239 and plants exceeds the CT-AC frequency significantly ( $p < 0.001$ ) (Figure  
240 3A). For the splice site combination CT-AC a sequence logo, which shows  
241 the conservation of this splice site in four selected species, was designed  
242 (Figure 3B). In summary, we conclude that CT-AC is a major non-canonical  
243 splice site combination in fungi, while AT-AC is not.



**Figure 3: CT-AC frequency exceeds AT-AC frequency in fungi.** A) Number of the minor non-canonical splice site combination CT-AC in comparison to the major non-canonical splice site combination AT-AC in each kingdom ( $p < 0.001$ ). B) Sequence logo for the splice site combination CT-AC in four selected fungal species (*Alternaria alternata*, *Aspergillus brasiliensis*, *Fomitopsis pinicola* and *Zymoseptoria tritici*). In total, 67 supported splice sites with this combination were used to generate the sequence logo.

244 The highest frequencies of the splice site combination CT-AC, supported  
 245 by RNA-Seq reads, were observed in *Alternaria alternata*, *Aspergillus brasiliensis*,  
 246 *Fomitopsis pinicola* and *Zymoseptoria tritici* (approx. 0.08 - 0.09 %).  
 247 As AT-AC was described as major non-canonical splice site, these findings  
 248 indicate a different splice site pattern in fungi compared to animals and  
 249 plants (Figure 3).

## 250 Intron size analysis

251 In total, 8,060,924, 737,783 and 2,785,484 transcripts across animals,  
 252 fungi and plants, respectively, were selected to check whether the intron  
 253 lengths are multiples of three. Introns with this property could be kept in  
 254 the final transcript without causing a shift in the reading frame. There is  
 255 no significant difference between introns with different splice site combina-  
 256 tions (Table 2). The ratio of introns with a length divisible by 3 is very close  
 257 to 33.3 % which would be expected based on an equal distribution. The  
 258 only exception are minor non-canonical splice site combinations in fungi  
 259 which are slightly less likely to occur in introns with a length divisible by 3.

**Table 2: Proportion of introns with length divisible by 3. The results of intron length analysis for selected splice site combinations for animals, fungi and plants are shown.**

	splice site combination	frequency of introns divisible by 3	total number of introns divisible by 3
animals	GT-AG	0.333862150381	n=63677347
	AT-AC	0.325106284189	n=68919
	GC-AG	0.330352389911	n=636823
	others	0.327633755094	n=496411
fungi	GT-AG	0.33932356858	n=2273756
	AT-AC	0.331775700935	n=428
	GC-AG	0.333577333793	n=23224
	others	0.3125	n=6240
plants	GT-AG	0.332967299596	n=14227286
	AT-AC	0.326150175229	n=13411
	GC-AG	0.329271562364	n=216326
	others	0.323971037399	n=93638

## 260 **Conservation of non-canonical splice site combinations** 261 **across species**

262 In total, *A. thaliana* transcripts containing 1,073 GC-AG, 64 AT-AC and 19  
263 minor non-canonical splice sites were aligned to transcripts of all plant  
264 species. Homologous intron positions were checked for non-canonical  
265 splice sites. GC-AG splice site combinations were conserved in 9,830  
266 sequences, matched with other non-canonical splice site combinations in  
267 121 cases, and aligned to GT-AG in 13,045 sequences. Given that the  
268 dominance of GT-AG splice sites was around 98 %, the number observed  
269 here indicates a strong conservation of GC-AG splice site combinations.  
270 AT-AC splice site combinations were conserved in 967 other sequences,  
271 matched with other non-canonical splice site combinations in 93 cases,  
272 and aligned to GT-AG in 157 sequences. These numbers indicate a con-  
273 servation of AT-AC splice site combinations, which exceeds the conserva-  
274 tion of GC-AG splice site combinations substantially. Minor non-canonical  
275 splice sites were conserved in 48 other sequences, matched with other  
276 non-canonical splice site combinations in 64 cases, and were aligned to  
277 a canonical GT-AG splice site in 213 cases. This pattern suggests that  
278 most non-canonical splice site combinations are either (A) mutations of  
279 the canonical ones or (B) mutated towards GT-AG splice site combina-  
280 tions.

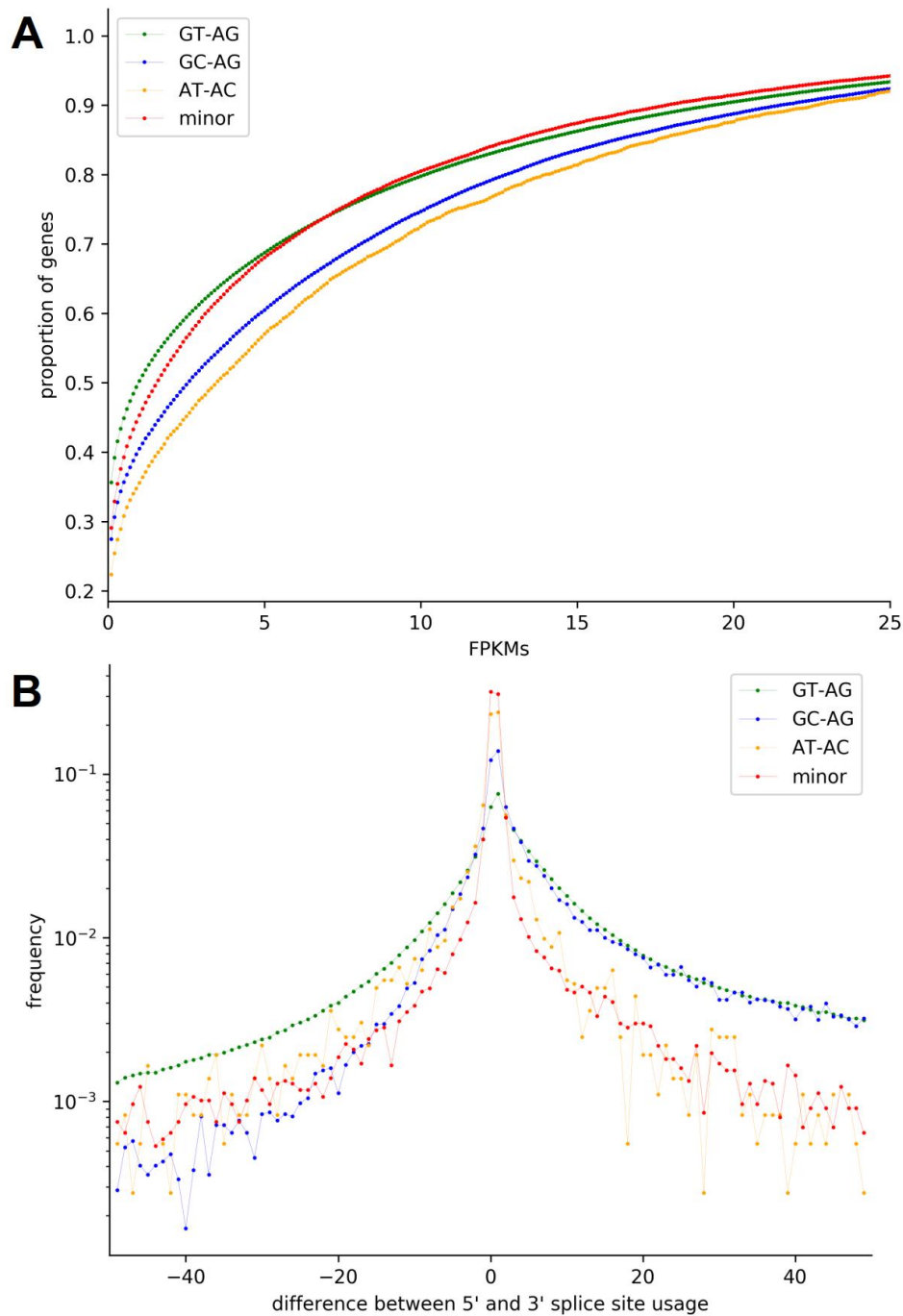
281 The power of this analysis is currently limited by the quality of the align-  
282 ment. Although splice site combinations should be aligned properly in  
283 most cases, small differences in the number could be caused by ambigu-  
284 ous situations. It is likely that both hypothesis stated above are partly valid.  
285 To assign each splice site combination to A or B, a manual inspection of  
286 the observed phylogenetic pattern would be required.

## 287 **Usage of non-canonical splice sites**

288 Non-canonical splice site combinations were described to have regula-  
289 tory roles by slowing down the splicing process [41]. Previous reports

290 also indicated that non-canonical splice site combinations might appear in  
291 pseudogenes [9, 10]. To analyse a possible correlation of non-canonical  
292 splice sites with low transcriptional activity, we compared the transcript  
293 abundance of genes with non-canonical splice site combinations to genes  
294 with only canonical GT-AG splice site combinations (Figure 4A). Genes  
295 with at least one non-canonical splice site combination are generally less  
296 likely to be lowly expressed than genes with only canonical splice sites.  
297 While this trend holds true for all analysed non-canonical splice site com-  
298 bination groups, GC-AG and AT-AC containing genes display especially  
299 low proportions of genes with low FPKMs. We speculate that a stronger  
300 transcriptional activity of genes with non-canonical splice sites compen-  
301 sates for lower turnover rates in the splicing process. The regulation of the  
302 genes might be shifted from the transcriptional to the post-transcriptional  
303 level. This trend is similar for animals and plants (AdditionalFile 13). In  
304 fungi, genes with minor non-canonical splice sites display relatively high  
305 proportions of genes with low FPKMs.  
306 Moreover, a higher number of non-canonical splice sites per gene is as-  
307 sociated with a lower expression. This leads to the suggestion, that non-  
308 canonical splice sites occur more often within pseudogenes.

309



**Figure 4: Usage of non-canonical splice sites in plant species.** A) Comparison of the transcript abundance (FPKMs) of genes with non-canonical splice site combinations to genes with only canonical GT-AG splice site combinations. GC-AG and AT-AC containing genes display especially low proportions of genes with low FPKMs. This leads to a higher transcript abundance of genes with low FPKMs. B) Comparison of the usage of 5' and 3' splice sites. On the x-axis, the difference between the 5' splice site usage and the usage of the 3' splice site is shown. A fast drop of values when going to the negative side of the x-axis indicates that the 3' splice site is probably more flexible than the 5' splice site.



310 Introns are mostly defined by phylogenetically conserved splice sites, but  
311 nevertheless some variation of these splice sites is possible [9, 10, 25, 26,  
312 40]. To understand the amount of flexibility in respect to different terminal  
313 dinucleotides, we compared the usage of donor and acceptor splice sites  
314 over 4,141,196 introns in plants, 3,915,559 introns in animals and 340,619  
315 introns in fungi (Figure 4B). The plot shows that the 3' splice site seems  
316 to be more flexible than the 5' splice site which was observed in all three  
317 kingdoms. Our observations align well with previous findings of a higher  
318 flexibility at the 3' splice site compared to the 5' splice site. A mutated 5'  
319 splice site represses the removal of the upstream intron [10, 42, 43]. Fur-  
320 ther, for plants and animals, the difference between the usage of the 5'  
321 splice site and the 3' splice site is notably higher for introns with the splice  
322 site combination GC-AG.

323

324 Although *bona fide* non-canonical splice site combinations are present in  
325 many plant transcripts [10], additional isoforms of the genes might exist.  
326 To evaluate the relevance of such alternative isoforms, we assessed the  
327 contribution of isoforms to the overall abundance of transcripts of a gene.  
328 Therefore, the usage of splice sites flanking an intron was compared to  
329 the average usage of splice sites. This reveals how often a certain intron  
330 is removed by splicing. Introns with low usage values might only be in-  
331 volved in minor transcript isoforms. While most introns display no or very  
332 small differences, GT-AG introns deviate from this trend. This indicates  
333 that non-canonical splice site combinations are frequently part of the dom-  
334 inant isoform. Again, these findings were similar for all of the investigated  
335 kingdoms.

336

## 337 **Conclusion**

338 Our investigation of non-canonical splice sites in animals, fungi and plants  
339 revealed kingdom specific differences. Animal species with a high propor-  
340 tion of GA-AG splice site combinations were examined. Further, properties

341 of introns and splice sites were analysed. One aspect of this analysis is,  
342 that the 3' splice site seems to be more flexible than the 5' splice site,  
343 which was observed in all three kingdoms. In fungi, the splice site com-  
344 bination CT-AC is more frequent than the splice site combination AT-AC.  
345 This makes CT-AC a major non-canonical splice site combination in fungal  
346 species, while AT-AC should be considered a minor non-canonical splice  
347 site in fungi. Overall, our findings demonstrate the importance of con-  
348 sidering non-canonical splice sites despite their low relative frequency in  
349 comparison to the canonical splice site combination GT-AG. RNA-Seq data  
350 confirmed the existence and usage of numerous non-canonical splice site  
351 combinations. By neglecting non-canonical splice sites, *bona fide* genes  
352 might be excluded or at least structurally altered.

## 353 **Methods**

### 354 **Analysis and validation of splice site combinations**

355 Genome sequences (FASTA) and corresponding annotations (GFF3) of  
356 130 fungal species and 489 animal species were retrieved from the  
357 NCBI. Representative transcript and peptide sequences were extracted  
358 as described before [10]. General statistics were calculated using a  
359 Python script [10]. The completeness of all data sets was assessed with  
360 BUSCO v3 [44] using the reference data sets 'fungi odb9' and 'meta-  
361 zoa odb9', respectively [45] (AdditionalFile 14 and 15). To validate the  
362 detected splice site combinations, paired-end RNA-Seq data sets were  
363 retrieved from the Sequence Read Archive [46] (AdditionalFile 16 and  
364 17). The following validation approach [10] utilized STAR v2.5.1b [47]  
365 for the read mapping and Python scripts for downstream processing  
366 (<https://doi.org/10.5281/zenodo.2586989>). An overview of the RNA-Seq  
367 read coverage depth of splice sites in animals [48] and fungi [49] is avail-  
368 able. RNA-Seq read mappings with STAR and HiSat2 were compared  
369 based on a gold standard generated by exonerate, because a previ-

370 ous report [50] indicated a superiority of STAR. All transcripts with non-  
371 canonical splice sites in *A. thaliana* and *Oryza sativa* were considered.  
372 When investigating the alignment of RNA-Seq reads over non-canonical  
373 splice sites, we observed a high accuracy for both mappers without a  
374 clear difference between them. Previously described scripts [10] were  
375 adjusted for this analysis and updated versions are available on github  
376 (<https://doi.org/10.5281/zenodo.2586989>). The distribution of genome  
377 sizes was analysed using the Python package dabest [51]. Sequence  
378 logos for the analysed splice sites were designed at [http://weblogo.](http://weblogo.berkeley.edu/logo.cgi)  
379 [berkeley.edu/logo.cgi](http://weblogo.berkeley.edu/logo.cgi) [52].

### 380 **Calculation of the splice site diversity**

381 A custom Python script was applied to calculate the Shannon diversity in-  
382 dex ( $H'$ ) [53] of all splice site combinations in fungi, animals and plants  
383 (<https://doi.org/10.5281/zenodo.2586989>). To determine the significance  
384 of the obtained results, a Kruskal-Wallis test [54] was calculated using the  
385 Python package scipy [55]. Further, the interquartile range of all distribu-  
386 tions was examined.

### 387 **Investigation of a common non-canonical splice site in** 388 **fungi**

389 A Mann-Whitney U Test implemented in the Python package scipy was  
390 performed to analyse differences in the number of minor non-canonical  
391 splice site combinations. The observed distributions were visualized in  
392 a boxplot (<https://doi.org/10.5281/zenodo.2586989>) constructed with the  
393 Python package plotly [56].

### 394 **Detection of potential U1 snRNAs**

395 A potential U1 snRNA of *Pan troglodytes* (obtained from the NCBI) was  
396 subjected to BLASTn [57] against the genome sequences of selected

397 species. Hits with a score above 100, with at least 80 % similarity and  
398 with the conserved sequence at the 5' end of the snRNA [58] were in-  
399 vestigated, as these sequences are potential U1 snRNAs. The obtained  
400 sequences were compared and small nucleotide variants were detected.

## 401 **Correlation between the GC content of the genome and** 402 **the GC content of the splice sites**

403 The Pearson correlation coefficient between the GC content of the genome  
404 sequence of each species and the GC content of the respective splice site  
405 combination was calculated using the Python package `scipy`. Splice site  
406 combinations were weighted with the number of occurrences while calcu-  
407 lating the GC content. Finally, the correlation coefficient and the p-value  
408 were determined. For better visualization, a scatter plot was constructed  
409 with the Python package `plotly` [56].

## 410 **Phylogeny of non-canonical splice sites**

411 All *A. thaliana* transcripts with non-canonical splice sites were subjected  
412 to BLASTn searches against the transcript sequences of all other plant  
413 species previously studied [10]. The best hit per species was selected for  
414 an alignment against the respective genomic region with `exonerate` [59].  
415 Next, splice site combinations were extracted and aligned. This align-  
416 ment utilized MAFFT v7 [60] by representing different splice site com-  
417 binations as amino acids. Finally, splice site combinations aligned with  
418 the non-canonical splice site combinations of *A. thaliana* were analysed  
419 (<https://doi.org/10.5281/zenodo.2586989>).

## 420 **Usage of non-canonical splice sites**

421 Genes were classified based on the presence/absence of non-canonical  
422 splice combinations into four groups: GT-AG, GC-AG, AT-AC, and minor  
423 non-canonical splice site genes. When having different non-canonical

424 splice sites, genes were assigned into multiple groups. Next, the tran-  
425 scription of these genes was quantified based on RNA-Seq using feature-  
426 Counts [61] based on the RNA-Seq read mapping generated with STAR.  
427 Binning of the genes was performed based on the fragments per kilobase  
428 transcript length per million assigned reads (FPKMs). Despite various  
429 shortcomings [62], we consider FPKMs to be acceptable for this analysis.  
430 Outlier genes with extremely high values were excluded from this analysis  
431 and the visualization. Next, a cumulative sum of the relative bin sizes was  
432 calculated. The aim was to compare the transcriptional activity of genes  
433 with different splice site combinations i.e. to test whether non-canonical  
434 splice site combinations are enriched in lowly transcribed genes.

435  
436 Usage of splice sites was calculated per intron as previously described  
437 [10]. The difference between both ends of an intron was calculated. The  
438 distribution of these differences per splice site type were analysed. In-  
439 trons were grouped by their splice site combination. The average of both  
440 coverage values of the directly flanking exon positions was calculated as  
441 estimate of the local expression around a splice site combination. Next,  
442 the sequencing coverage of a transcript was estimated by multiplying 200  
443 bp (assuming 2x100 nt reads) with the number of read counts per gene  
444 and normalization to the transcript length. The difference between both  
445 values was calculated for each intron to assess its presence in the major  
446 isoform.

447

## 448 **Acknowledgments**

449 We thank members of Genetics and Genomics of Plants for discussion of  
450 preliminary results. We are very grateful to Hanna Schilbert, Janik Siele-  
451 mann, and Iain Place for helpful comments on the manuscript.

## 452 References

- 453 [1] Moore, Melissa J and Sharp, Phillip A, "Site-specific modification of pre-mRNA: the  
454 2'-hydroxyl groups at the splice sites," *Science*, vol. 256, no. 5059, pp. 992–997,  
455 1992.
- 456 [2] Barbosa-Morais, Nuno L and Irimia, Manuel and Pan, Qun and Xiong, Hui Y and  
457 Gueroussov, Serge and Lee, Leo J and Slobodeniuc, Valentina and Kutter, Claudia  
458 and Watt, Stephen and Çolak, Recep and others, "The evolutionary landscape of  
459 alternative splicing in vertebrate species," *Science*, vol. 338, no. 6114, pp. 1587–  
460 1593, 2012.
- 461 [3] Ben-Dov, Claudia and Hartmann, Britta and Lundgren, Josefin and Valcárcel, Juan,  
462 "Genome-wide analysis of alternative pre-mRNA splicing," *Journal of Biological*  
463 *Chemistry*, vol. 283, no. 3, pp. 1229–1233, 2008.
- 464 [4] Matlin, Arianne J and Clark, Francis and Smith, Christopher WJ, "Understanding  
465 alternative splicing: towards a cellular code," *Nature Reviews Molecular Cell Biology*,  
466 vol. 6, no. 5, p. 386, 2005.
- 467 [5] Sibley, Christopher R and Blazquez, Lorea and Ule, Jernej, "Lessons from non-  
468 canonical splicing," *Nature Reviews Genetics*, vol. 17, no. 7, p. 407, 2016.
- 469 [6] Maniatis, Tom and Tasic, Bosiljka, "Alternative pre-mRNA splicing and proteome  
470 expansion in metazoans," *Nature*, vol. 418, no. 6894, p. 236, 2002.
- 471 [7] Xue, Min and Chen, Bing and Ye, Qingqing and Shao, Jingru and Lyu, Zhangxia and  
472 Wen, Jianfan, "Sense-antisense gene overlap causes evolutionary retention of the  
473 few introns in *Giardia* genome and the implications," *bioRxiv*, 2018. doi: 10.1101/  
474 333310.
- 475 [8] Chorev, Michal and Carmel, Liran, "The function of introns," *Frontiers in Genetics*,  
476 vol. 3, 2012.
- 477 [9] Burset, M and Seledtsov, IA and Solovyev, VV, "Analysis of canonical and non-  
478 canonical splice sites in mammalian genomes," *Nucleic Acids Research*, vol. 28,  
479 no. 21, pp. 4364–4375, 2000.
- 480 [10] Pucker, Boas and Brockington, Samuel F, "Genome-wide analyses supported by  
481 RNA-Seq reveal non-canonical splice sites in plant genomes," *BMC Genomics*,  
482 vol. 19, no. 1, p. 980, 2018. doi: <https://doi.org/10.1186/s12864-018-5360-z>.
- 483 [11] Bon, Elisabeth and Casaregola, Serge and Blandin, Gaëlle and Llorente, Bertrand  
484 and Neuvéglise, Cécile and Munsterkotter, Martin and Guldener, Ulrich and Mewes,  
485 Hans-Werner and Helden, Jacques Van and Dujon, Bernard and others, "Molecular  
486 evolution of eukaryotic genomes: hemiascomycetous yeast spliceosomal introns,"  
487 *Nucleic Acids Research*, vol. 31, no. 4, pp. 1121–1135, 2003.

- 488 [12] Logsdon, John M, “The recent origins of spliceosomal introns revisited,” *Current*  
489 *Opinion in Genetics & Development*, vol. 8, no. 6, pp. 637–648, 1998.
- 490 [13] Burge, Chris and Karlin, Samuel, “Prediction of complete gene structures in human  
491 genomic DNA1,” *Journal of Molecular Biology*, vol. 268, no. 1, pp. 78–94, 1997.
- 492 [14] Stanke, Mario and Waack, Stephan, “Gene prediction with a hidden Markov model  
493 and a new intron submodel,” *Bioinformatics*, vol. 19, no. suppl\_2, pp. ii215–ii225,  
494 2003.
- 495 [15] Davis, Carrie A and Grate, Leslie and Spingola, Marc and Ares Jr, Manuel, “Test of  
496 intron predictions reveals novel splice sites, alternatively spliced mRNAs and new  
497 introns in meiotically regulated genes of yeast,” *Nucleic Acids Research*, vol. 28,  
498 no. 8, pp. 1700–1706, 2000.
- 499 [16] Wahl, Markus C and Will, Cindy L and Lührmann, Reinhard, “The spliceosome:  
500 design principles of a dynamic RNP machine,” *Cell*, vol. 136, no. 4, pp. 701–718,  
501 2009.
- 502 [17] Sharp, Phillip A and Burge, Christopher B, “Classification of introns: U2-type or  
503 U12-type,” *Cell*, vol. 91, no. 7, pp. 875–879, 1997.
- 504 [18] Hall, Stephen L and Padgett, Richard A, “Requirement of U12 snRNA for in vivo  
505 splicing of a minor class of eukaryotic nuclear pre-mRNA introns,” *Science*, vol. 271,  
506 no. 5256, pp. 1716–1718, 1996.
- 507 [19] Turunen, Janne J and Niemelä, Elina H and Verma, Bhupendra and Frilander, Mikko  
508 J, “The significant other: splicing by the minor spliceosome,” *Wiley Interdisciplinary*  
509 *Reviews: RNA*, vol. 4, no. 1, pp. 61–76, 2013.
- 510 [20] Dietrich, Rosemary C and Incorvaia, Robert and Padgett, Richard A, “Terminal in-  
511 tron dinucleotide sequences do not distinguish between U2-and U12-dependent in-  
512 trons,” *Molecular Cell*, vol. 1, no. 1, pp. 151–160, 1997.
- 513 [21] Wilkinson, Max E and Fica, Sebastian M and Galej, Wojciech P and Norman, Chris-  
514 tine M and Newman, Andrew J and Nagai, Kiyoshi, “Postcatalytic spliceosome struc-  
515 ture reveals mechanism of 3’-splice site selection,” *Science*, vol. 358, no. 6368,  
516 pp. 1283–1288, 2017.
- 517 [22] Burge, Christopher B and Tuschl, Thomas and Sharp, Phillip A, “Splicing of pre-  
518 cursors to mRNAs by the spliceosomes,” *Cold Spring Harbor Monograph Series*,  
519 vol. 37, pp. 525–560, 1999.
- 520 [23] Roca, Xavier and Krainer, Adrian R and Eperon, Ian C, “Pick one, but be quick: 5’  
521 splice sites and the problems of too many choices,” *Genes & Development*, vol. 27,  
522 no. 2, pp. 129–144, 2013.

- 523 [24] Shi, Yigong, “The spliceosome: a protein-directed metalloribozyme,” *Journal of*  
524 *Molecular Biology*, vol. 429, no. 17, pp. 2640–2653, 2017.
- 525 [25] Kupfer, Doris M and Drabenstot, Scott D and Buchanan, Kent L and Lai, Hongsh-  
526 ing and Zhu, Hua and Dyer, David W and Roe, Bruce A and Murphy, Juneann W,  
527 “Introns and splicing elements of five diverse fungi,” *Eukaryotic Cell*, vol. 3, no. 5,  
528 pp. 1088–1100, 2004.
- 529 [26] Kitamura–Abe, Sumie and Itoh, Hitomi and Washio, Takanori and Tsutsumi, Akihiro  
530 and Tomita, Masaru, “Characterization of the splice sites in GT–AG and GC–AG  
531 introns in higher eukaryotes using full-length cDNAs,” *Journal of Bioinformatics and*  
532 *Computational Biology*, vol. 2, no. 02, pp. 309–331, 2004.
- 533 [27] Michael, Deutsch and Manyuan, Long, “Intron—exon structures of eukaryotic model  
534 organisms,” *Nucleic Acids Research*, vol. 27, no. 15, pp. 3219–3228, 1999.
- 535 [28] Modrek, Barmak and Resch, Alissa and Grasso, Catherine and Lee, Christopher,  
536 “Genome-wide detection of alternative splicing in expressed sequences of human  
537 genes,” *Nucleic Acids Research*, vol. 29, no. 13, pp. 2850–2859, 2001.
- 538 [29] Pucker, Boas and Holtgräwe, Daniela and Weisshaar, Bernd, “Consideration of  
539 non-canonical splice sites improves gene prediction on the Arabidopsis thaliana  
540 Niederzenz-1 genome sequence,” *BMC Research Notes*, vol. 10, no. 1, p. 667, 2017.  
541 doi: <https://doi.org/10.1186/s13104-017-2985-y>.
- 542 [30] Sparks, Michael E and Brendel, Volker, “Incorporation of splice site probability mod-  
543 els for non-canonical introns improves gene structure prediction in plants,” *Bioinfor-*  
544 *matics*, vol. 21, no. Suppl\_3, pp. iii20–iii30, 2005.
- 545 [31] Dubrovina, AS and Kiselev, KV and Zhuravlev, Yu N, “The role of canonical and  
546 noncanonical pre-mRNA splicing in plant stress responses,” *BioMed Research In-*  
547 *ternational*, vol. 2013, 2013.
- 548 [32] Alexandrov, Nikolai N and Troukhan, Maxim E and Brover, Vyacheslav V and Tatari-  
549 nova, Tatiana and Flavell, Richard B and Feldmann, Kenneth A, “Features of Ara-  
550 bidopsis genes and genome discovered using full-length cDNAs,” *Plant Molecular*  
551 *Biology*, vol. 60, no. 1, pp. 69–85, 2006.
- 552 [33] Niu, Xiangli and Luo, Di and Gao, Shaopei and Ren, Guangjun and Chang, Lijuan  
553 and Zhou, Yuke and Luo, Xiaoli and Li, Yuxiang and Hou, Pei and Tang, Wei and oth-  
554 ers, “A conserved unusual posttranscriptional processing mediated by short, direct  
555 repeated (SDR) sequences in plants,” *Journal of Genetics and Genomics*, vol. 37,  
556 no. 1, pp. 85–99, 2010.
- 557 [34] Erkelenz, Steffen and Theiss, Stephan and Kaisers, Wolfgang and Ptok, Johannes  
558 and Walotka, Lara and Müller, Lisa and Hillebrand, Frank and Brillen, Anna-Lena



- 559 and Sladek, Michael and Schaal, Heiner, "Ranking noncanonical 5' splice site us-  
560 age by genome-wide RNA-seq analysis and splicing reporter assays," *Genome Re-*  
561 *search*, vol. 28, no. 12, pp. 1826–1840, 2018.
- 562 [35] Grützmann, Konrad and Szafranski, Karol and Pohl, Martin and Voigt, Kerstin and  
563 Petzold, Andreas and Schuster, Stefan, "Fungal alternative splicing is associated  
564 with multicellular complexity and virulence: a genome-wide multi-species study,"  
565 *DNA Research*, vol. 21, no. 1, pp. 27–39, 2013.
- 566 [36] Robertson, Hugh M, "Non-canonical GA and GG 5'Intron Donor Splice Sites Are  
567 Common in the Copepod *Eurytemora affinis*," *G3: Genes, Genomes, Genetics*,  
568 pp. g3–300189, 2017.
- 569 [37] Brackenridge, Simon and Wilkie, Andrew OM and Sreaton, Gavin R, "Efficient  
570 use of a 'dead-end'GA 5' splice site in the human fibroblast growth factor recep-  
571 tor genes," *The EMBO Journal*, vol. 22, no. 7, pp. 1620–1631, 2003.
- 572 [38] Tyler, Brett M, "Phytophthora sojae: root rot pathogen of soybean and model  
573 oomycete," *Molecular Plant Pathology*, vol. 8, no. 1, pp. 1–8, 2007.
- 574 [39] Förster, Helga and Coffey, Michael D and Elwood, Hille and Sogin, Mitchell L, "Se-  
575 quence analysis of the small subunit ribosomal RNAs of three zoosporic fungi and  
576 implications for fungal evolution," *Mycologia*, pp. 306–312, 1990.
- 577 [40] Shen, Danyu and Ye, Wenwu and Dong, Suomeng and Wang, Yuanchao and Dou,  
578 Daolong, "Characterization of intronic structures and alternative splicing in *Phytoph-*  
579 *thora sojae* by comparative analysis of expressed sequence tags and genomic se-  
580 quences," *Canadian journal of Microbiology*, vol. 57, no. 2, pp. 84–90, 2011.
- 581 [41] Aebi, M and Hornig, H and Padgett, RA and Reiser, J and Weissmann, C, "Se-  
582 quence requirements for splicing of higher eukaryotic nuclear pre-mRNA," *Cell*,  
583 vol. 47, no. 4, pp. 555–565, 1986.
- 584 [42] Talerico, MELISSA and Berget, SUSAN M, "Effect of 5'splice site mutations on  
585 splicing of the preceding intron.," *Molecular and Cellular Biology*, vol. 10, no. 12,  
586 pp. 6299–6305, 1990.
- 587 [43] Berget, Susan M, "Exon recognition in vertebrate splicing," *Journal of biological*  
588 *Chemistry*, vol. 270, no. 6, pp. 2411–2414, 1995.
- 589 [44] Simão, Felipe A and Waterhouse, Robert M and Ioannidis, Panagiotis and Krivent-  
590 seva, Evgenia V and Zdobnov, Evgeny M, "BUSCO: assessing genome assembly  
591 and annotation completeness with single-copy orthologs," *Bioinformatics*, vol. 31,  
592 no. 19, pp. 3210–3212, 2015.

- 593 [45] Kriventseva, Evgenia V and Tegenfeldt, Fredrik and Petty, Tom J and Waterhouse,  
594 Robert M and Simao, Felipe A and Pozdnyakov, Igor A and Ioannidis, Panagiotis  
595 and Zdobnov, Evgeny M, “OrthoDB v8: update of the hierarchical catalog of or-  
596 thologs and the underlying free software,” *Nucleic Acids Research*, vol. 43, no. D1,  
597 pp. D250–D256, 2014.
- 598 [46] Leinonen, Rasko and Sugawara, Hideaki and Shumway, Martin and International  
599 Nucleotide Sequence Database Collaboration, “The sequence read archive,” *Nu-  
600 cleic Acids Research*, vol. 39, no. suppl\_1, pp. D19–D21, 2010.
- 601 [47] Dobin, Alexander and Davis, Carrie A and Schlesinger, Felix and Drenkow, Jorg  
602 and Zaleski, Chris and Jha, Sonali and Batut, Philippe and Chaisson, Mark and  
603 Gingeras, Thomas R, “STAR: ultrafast universal RNA-seq aligner,” *Bioinformatics*,  
604 vol. 29, no. 1, pp. 15–21, 2013.
- 605 [48] Pucker, Boas and Frey, Katharina, “RNA-Seq read coverage depth of splice sites in  
606 animals,” 2019. doi: [10.4119/unibi/2934226](https://doi.org/10.4119/unibi/2934226).
- 607 [49] Pucker, Boas and Frey, Katharina, “RNA-Seq read coverage depth of splice sites in  
608 fungi,” 2019. doi: [10.4119/unibi/2934220](https://doi.org/10.4119/unibi/2934220).
- 609 [50] Dobin, Alexander and Gingeras, Thomas R., “Comment on “TopHat2: accurate  
610 alignment of transcriptomes in the presence of insertions, deletions and gene fu-  
611 sions” by Kim et al.,” 2013. doi: <https://doi.org/10.1101/000851>.
- 612 [51] Ho, Joses and Tumkaya, Tayfun and Aryal, Sameer and Choi, Hyungwon and  
613 Claridge-Chang, Adam, “Moving beyond P values: Everyday data analysis with es-  
614 timation plots,” *bioRxiv*, p. 377978, 2018. doi: <https://doi.org/10.1101/377978>.
- 615 [52] Crooks, Gavin E and Hon, Gary and Chandonia, John-Marc and Brenner, Steven E,  
616 “WebLogo: a sequence logo generator,” *Genome Research*, vol. 14, no. 6, pp. 1188–  
617 1190, 2004.
- 618 [53] Heip, Carlo, “A new index measuring evenness,” *Journal of the Marine Biological  
619 Association of the United Kingdom*, vol. 54, no. 3, pp. 555–557, 1974.
- 620 [54] Breslow, Norman, “A generalized Kruskal-Wallis test for comparing K samples sub-  
621 ject to unequal patterns of censorship,” *Biometrika*, vol. 57, no. 3, pp. 579–594,  
622 1970.
- 623 [55] Eric Jones and Travis Oliphant and Pearu Peterson and others, “SciPy: Open source  
624 scientific tools for Python,” 2001. url: <http://www.scipy.org/>.
- 625 [56] Plotly Technologies Inc., “Collaborative data science,” 2015.
- 626 [57] Altschul, Stephen F and Gish, Warren and Miller, Webb and Myers, Eugene W and  
627 Lipman, David J, “Basic local alignment search tool,” *Journal of Molecular Biology*,  
628 vol. 215, no. 3, pp. 403–410, 1990.

- 629 [58] Stark, Holger and Dube, Prakash and Lührmann, Reinhard and Kastner, Berthold,  
630 “Arrangement of RNA and proteins in the spliceosomal U1 small nuclear ribonucleo-  
631 protein particle,” *Nature*, vol. 409, no. 6819, p. 539, 2001.
- 632 [59] Slater, Guy St C and Birney, Ewan, “Automated generation of heuristics for biological  
633 sequence comparison,” *BMC Bioinformatics*, vol. 6, no. 1, p. 31, 2005.
- 634 [60] Katoh, Kazutaka and Standley, Daron M, “MAFFT multiple sequence alignment soft-  
635 ware version 7: improvements in performance and usability,” *Molecular Biology and*  
636 *Evolution*, vol. 30, no. 4, pp. 772–780, 2013.
- 637 [61] Liao, Yang and Smyth, Gordon K and Shi, Wei, “featureCounts: an efficient general  
638 purpose program for assigning sequence reads to genomic features,” *Bioinformatics*,  
639 vol. 30, no. 7, pp. 923–930, 2013.
- 640 [62] Conesa, Ana and Madrigal, Pedro and Tarazona, Sonia and Gomez-Cabrero, David  
641 and Cervera, Alejandra and McPherson, Andrew and Szczesniak, Michał Wojciech  
642 and Gaffney, Daniel J and Elo, Laura L and Zhang, Xuegong and others, “A survey  
643 of best practices for RNA-seq data analysis,” *Genome Biology*, vol. 17, no. 1, p. 13,  
644 2016.

## 645 **Additional Files**

646 **AdditionalFile 1.** List of genome sequences and annotations of the investigated animal  
647 species.

648  
649 **AdditionalFile 2.** List of genome sequences and annotations of the investigated fungal  
650 species.

651  
652 **AdditionalFile 3.** Distribution of genome sizes of all species.

653  
654 **AdditionalFile 4.** Distribution of canonical and non-canonical splice sites per species in  
655 the animal kingdom.

656  
657 **AdditionalFile 5.** Distribution of canonical and non-canonical splice sites per species in  
658 the fungal kingdom.

659  
660 **AdditionalFile 6.** List of all possible splice site combinations in animal species.

661  
662 **AdditionalFile 7.** List of all possible splice site combinations in fungal species.

663  
664 **AdditionalFile 8.** Genome statistics concerning each analysed animal species.

665  
666 **AdditionalFile 9.** Genome statistics concerning each analysed fungal species.

667  
668 **AdditionalFile 10.** Correlation between the GC content of the genome and the GC con-  
669 tent of the splice sites per kingdom.

670  
671 **AdditionalFile 11.** Correlation between the size of the used RNA-Seq data sets and the  
672 number of supported splice sites.

673  
674 **AdditionalFile 12.** Proportion of genes with non-canonical splice sites in dependence of  
675 the number of introns.

676  
677 **AdditionalFile 13.** Usage of non-canonical splice sites in animals and fungi.

678  
679 **AdditionalFile 14.** Non-canonical splice sites in BUSCOs and in all genes were as-  
680 sessed per species in the animal kingdom.

681  
682 **AdditionalFile 15.** Non-canonical splice sites in BUSCOs and in all genes were as-  
683 sessed per species in the fungal kingdom.

684  
685 **AdditionalFile 16.** List of Sequence Read Archive accession numbers of the investigated  
686 animal RNA-Seq data sets.

687  
688 **AdditionalFile 17.** List of Sequence Read Archive accession numbers of the investigated  
689 fungal RNA-Seq data sets.

690