

1 **Evolutionarily Conserved Alternative Splicing**

2 **Across Monocots**

3 **Wenbin Mei***, **Lucas Boatwright***, **Guanqiao Feng†**, **James C. Schnable‡**, **W.**

4 **Brad Barbazuk***, §,**

5 *Department of Biology, University of Florida, Gainesville, USA

6 †Graduate Program in Plant Molecular and Cellular Biology, University of Florida,

7 Gainesville, FL, USA

8 ‡Agronomy & Horticulture, University of Nebraska–Lincoln, Beadle Center E207,

9 Lincoln, Nebraska, USA

10 §Genetics Institute, University of Florida, Gainesville, USA

11 **Running title: Evolutionarily conserved alternative splicing**

12 **Keywords:** conserved alternative splicing events, monocot, grass family, RS and

13 RS2Z subfamilies, R2R3-MYB

14 **** Correspondence author:** W. Brad Barbazuk

15 Address: Department of Biology, University of Florida, Gainesville, USA

16 Telephone: 1(352) 273-8624

17 Email: bbarbazuk@ufl.edu

18 **Abstract**

19 One difficulty when identifying and analyzing alternative splicing (AS)
20 events in plants is distinguishing functional AS from splicing noise. One way to
21 add confidence to the validity of a splice isoform is to observe that it is conserved
22 across evolutionarily related species. We use a high throughput method to
23 identify junction based conserved AS events from RNA-Seq data across nine
24 plant species including: five grass monocots (maize, sorghum, rice, *Brachypodium*
25 and foxtail millet), plus two non-grass monocots (banana and African oil palm),
26 the eudicot *Arabidopsis* and the basal angiosperm *Amborella*. In total, 9,804
27 conserved AS events within 19,235 genes were identified conserved between 2
28 or more species studied. In grasses containing large regions of conserved
29 synteny, the frequency of conserved AS events is twice that observed for genes
30 outside of conserved synteny blocks. In plant-specific RS and RS2Z subfamilies,
31 we observe both conservation and divergence of AS events after the whole
32 genome duplication in maize. In addition, plant-specific RS and RS2Z subfamilies
33 are highly connected with R2R3-MYB in splicing networks. Furthermore, we
34 discovered that the network based on genes harboring conserved AS events is
35 enriched for phosphatases, kinases and ubiquitylation genes, which suggests that
36 AS may participate in regulating signaling pathways. These data lay the
37 foundation for identifying and studying conserved AS events in the monocots,
38 particularly across grass species, and this conserved AS resource identifies an
39 additional layer between genotype to phenotype that may impact future crop
40 improvement efforts.

41 **Introduction**

42 Alternative splicing (AS) is an important post-transcriptional process that
43 can produce two or more transcript isoforms from a pre-mRNA. AS occurs in the
44 spliceosome by removing introns and joining exons through the selective use of
45 splice sites (Lee and Rio 2015) and is governed by *cis*- regulatory elements such
46 as splicing enhancers/silencers, and *trans*- regulatory elements including SR
47 proteins and hnRNP proteins (Busch and Hertel 2012, Wang and Burge 2008,
48 Reddy et al. 2012, Reddy et al. 2013). AS isoforms that are translated can
49 influence proteome diversity, while others that are not translated are
50 purposefully non-functional and can act to post-transcriptionally modulate gene
51 product levels (Fu et al. 2009, Hammond, Wachter, and Breaker 2009). AS
52 participates in many important processes during the lifecycle of plants (Staiger
53 and Brown 2013) and AS occurs in response to many abiotic stresses
54 (Mastrangelo et al. 2012), for example red/blue light (Shikata et al. 2014, Wu et
55 al. 2014), salt stress (Feng et al. 2015), drought (Thatcher et al. 2016), flooding
56 (Syed et al. 2015) and temperature (Streitner et al. 2013, James et al. 2012).

57 RNA-Seq has become a standard tool to investigate transcriptomes and AS
58 isoforms. While many studies have identified AS in individual species (Li, Xiao,
59 and Zhu 2014, Thatcher et al. 2014, Shen et al. 2014, Filichkin et al. 2010,
60 Mandadi and Scholthof 2015), few report investigating conserved AS across
61 multiple plant genomes. Most published studies have focused on identifying
62 conserved AS events between only two species. For example Severing et al.
63 (2009) identified 56 protein-coding, conserved AS events between *Arabidopsis*
64 and rice orthologous gene sets, and 30 conserved AS events leading to non-
65 functional isoforms subjected to nonsense-mediated decay. Similarly, another
66 study detected 537 AS events conserved between *Arabidopsis* and *Brassica*

67 (Darracq and Adams 2013) and 71 conserved AS events were identified between
68 *Populus* and *Eucalyptus* (Xu et al. 2014). With respect to comparing more than
69 two species, one study found 16 AS transcripts conserved between
70 *Brachypodium*, rice, and *Arabidopsis* by performing all-vs-all BLAST between EST
71 sequences (Walters et al. 2013). Recently, discovery and identification of AS
72 events conserved broadly across eudicots has been reported (Chamala et al.
73 2015). Chuang et al. (2015) used available Sanger EST sets to characterize AS in
74 grass species but did not focus on conserved AS events. There are no
75 corresponding multi-species studies that leverage the deep NGS resources to
76 characterize and identify conserved AS events within monocots and the grass
77 family.

78 Many monocot species are economically important and grass species
79 (Poaceae) are a particularly important source of calories. Within the grass family
80 there are two major clades. The BEP clade is made up of the Bambusoideae,
81 Ehrhartoideae, and Pooideae and contains rice and *Brachypodium* (Zhao et al.
82 2013); the PACMAD clade includes Panicoideae, Arundinoideae, Chloridoideae,
83 Micrairoideae, Aristidoideae, and Danthonioideae, and includes maize, sorghum
84 and foxtail millet (Cotton et al. 2015). In addition to these two major clades,
85 there are three clades of basal grasses lineage Anomochlooideae, Pharoideae,
86 and Puelioideae (Cotton et al. 2015). Published whole genome sequences exist
87 for many monocot species including the grass species maize (Schnable et al.
88 2009), sorghum (Paterson et al. 2009), rice (Goff et al. 2002), *Brachypodium*
89 (Vogel et al. 2010) and foxtail millet (Bennetzen et al. 2012), the non-grass
90 species banana (D'Hont et al. 2012) and African oil palm (Singh et al. 2013).

91 To characterize AS and identify conserved events due to species shared
92 ancestral a computational method was employed that identifies high confidence
93 isoforms in any species with RNA-Seq data and a reference genome sequence,
94 and then uses a junction based approach to identify conserved AS events across
95 multiple species. In this analysis, we focused on seven monocot species that have
96 available reference genome sequences and deep RNA-Seq data sets. Five of these
97 are grasses (maize, rice, sorghum, foxtail millet and *Brachypodium*) that also
98 have substantial retained gene collinearity which eases the identification of
99 synteny and aids orthologue calls (Schnable, Freeling, and Lyons 2012). This
100 analysis also includes banana and African palm oil, which extends the discovery
101 of conserved events into non-grass monocot species. Although there are other
102 non-grass monocot species that have available reference genomes, such as
103 duckweed and orchid, there was not sufficient RNA-Seq data available for these
104 species to be included in this analysis. Likewise, species representing basal grass
105 lineages were excluded due to lack of reference genome sequences. Finally, the
106 eudicot *Arabidopsis thaliana* and the basal angiosperm *Amborella* were included
107 to provide outgroups for comparison. Therefore, the species studied allowed us
108 to identify conserved AS events at specific reference points across the grasses
109 and extends into non-grass monocots, eudicots and the sister species to all
110 angiosperms.

111 Examination of conserved events uncovers likely functionally important
112 AS isoforms of SR protein genes, stress-response transcription factors, and
113 members of known protein-protein interaction networks. In addition the
114 conserved AS event collection provides an opportunity to test the functional
115 sharing hypothesis (Su et al. 2006) between single-copy vs. non-single-copy

116 genes, something that was not examined in previous studies of conserved AS
117 events. This hypothesis predicts more instances of conserved alternative splicing
118 events will be found in single copy genes than in non-single-copy genes. To
119 measure to extent of selection pressure on AS, we examined Ka/Ks values on
120 exons that flank conserved AS events. We also present network analysis based
121 on the *Arabidopsis* protein-protein interaction data that provides evidence for a
122 close association between splicing factors and stress responsive transcription
123 factors that undergo conserved AS, and between splicing factors and R2R3-MYB
124 across species. Finally, we discuss several noteworthy cases of conserved AS
125 events across angiosperms including SPA3 possessing a conserved alternative
126 donor site, a KH RNA-binding domain containing protein with a widely
127 conserved exon skipping event, and the stress response NAC transcription factor
128 RD26 that harbors an intron retention event conserved across the angiosperms.

129 **Materials and Methods**

130 ***Genome Assemblies and Annotations***

131 All analyses were performed on publicly available sequence and
132 annotation collections for the following species: foxtail millet (*Setaria italica*),
133 rice (*Oryza sativa L. ssp. Japonica*), *Brachypodium* (*Brachypodium distachyon*), the
134 non-grass monocots banana (*Musa acuminata*) and African oil palm (*Elaeis*
135 *guineensis*), the eudicot *Arabidopsis thaliana* and the basal angiosperm *Amborella*
136 (*Amborella trichopoda*), which is a sister taxon to all angiosperms.

137 The Maize RefGen_v2 genome sequence and working gene set annotation
138 were retrieved from MaizeGDB (Lawrence et al. 2004). Reference genome
139 sequence and gene annotation for oil palm were retrieved from the Malaysian

140 Palm Oil Board (genomsawit.mpob.gov.my). Reference genome and annotation
141 for Banana were retrieved from the banana genome hub (Droc et al. 2013). All
142 remaining reference genome assemblies and annotations were retrieved from
143 Phytozome v10 (Goodstein et al. 2012) (Table S1).

144 ***Transcriptome Data Collection***

145 AS in maize (*Zea mays*) and sorghum (*Sorghum bicolor*) was previously
146 identified based on public RNA-Seq data and described in (Mei et al. In Review),
147 and available from figshare (<http://10.6084/m9.figshare.4205079>). Transcript
148 sequences and annotations of *Arabidopsis thaliana* were retrieved from
149 Araport11 PreRelease 20151202 (Cheng et al. 2016), which presents the
150 updated TAIR10 annotation by inclusion of 113 RNA-Seq data sets. For the
151 remaining 6 species, Illumina RNA-Seq, 454 and Sanger EST transcript sequences
152 were downloaded from NCBI GeneBank. The sequence collections are described
153 in Table S2.

154 ***Method to Detect Alternative Splicing***

155 The process and software parameters used to assemble transcripts from
156 Illumina short reads was described in full details (Mei et al. In Review). Briefly,
157 three approaches were taken to assemble transcripts from Illumina short reads
158 to maximize isoform detection. The assembly software used was annotation
159 guided Cufflinks 2.2.1 (Trapnell et al. 2012), genome-guided Trinity 2.0.4
160 (Grabherr et al. 2011) and annotation guided StringTie 1.0.0 (Pertea et al. 2015).
161 Isoforms built from cufflinks were required to have minimal expression FPKM
162 value of 0.1 and required full reads support for the whole isoforms.

163 We masked vector and contaminant sequence in the Sanger ESTs using
164 SeqClean (<http://seqclean.sourceforge.net/>). 454 transcripts were constructed
165 using Newbler v2.8 (Roche 2010) with parameters “-cdna -urt”. Cufflinks, Trinity
166 and StringTie assembled isoforms were merged with the Newbler 454
167 assemblies and the cleaned ESTs. Each species specific sequence collection was
168 aligned back to the appropriate reference genome sequence followed by
169 additional assembly and clustering to remove redundancy using the PASA 2.0
170 package (Haas et al. 2003). Isoform quality control steps described in Mei et al.
171 (In Review) were used on the PASA outputted transcript set to remove isoforms
172 that were poorly supported and/or potentially poor quality/poorly assembled.
173 Briefly, each splice junction in a final retained isoform was required to have a
174 minimum entropy score of 2 (Sturgill et al. 2013). Retained introns that define
175 potential intron-retention isoforms were expected over a minimum of 90% of
176 their length with a minimal median coverage depth of 10, and a minimum intron
177 retention isoform ratio of 10%. The intron retention rate was calculated as
178 described in Marquez et al. (2012), where the median coverage of the retained
179 intron was divided by the number of reads supporting the splice junction. Lastly,
180 we required each isoform to have a minimum FPKM of 1 and compose at least
181 5% of the total isoform abundance for that gene in at least one tissue examined.

182 We performed the above methods to identify and define AS events in rice,
183 foxtail millet, *Brachypodium*, banana, African oil palm, and *Amborella*. The
184 transcripts from Araport11 PreRelease 20151202 for *Arabidopsis* were realigned
185 to the *Arabidopsis* reference genome using PASA to identify and characterize the
186 set of *Arabidopsis* AS events. AS isoforms defined previously within maize and

187 sorghum Mei et al. (In Review) were used to complete the AS isoforms sets
188 analyzed in this study.

189 ***Orthofinder Clustering Across Nine Species***

190 Orthologous sequences were defined and grouped in Orthogroups across
191 the nine species by OrthoFinder with default setting (Emms and Kelly 2015).
192 OrthoFinder was designed for plant genome orthogroup identification and
193 accounts for gene length bias.

194 ***Identify Sets of Syntenic Genes Across Five Grass Species***

195 Syntenic orthologs across five grass species (*Zea mays*, *Sorghum bicolor*,
196 *Setaria italica*, *Oryza sativa*, *Brachypodium distachyon*) were defined using the
197 methodology employed by (Schnable, Freeling, and Lyons 2012). Gene models
198 were updated according to new annotations from Phytozome 10 (Goodstein et al.
199 2012). 11,995 syntenic gene lists across five species (sorghum, rice, foxtail millet,
200 *Brachypodium*, either maize1/maize2 or both) were used for downstream
201 analysis.

202 ***Identification of Conserved Alternative Splicing***

203 AS was independently identified in each species (see above) and all AS
204 isoform sequences of the same event type (intron retention, exon skipping etc.)
205 for a given species were pooled. The computational methodology to detect
206 conserved AS relies on the creation of splice anchor sequence tags (SAST) for
207 each AS event and then comparing the SASTs for a given event type with the
208 SASTs defining the same event type between species. SASTs were generated for

209 each AS event by extracting up to 300bp of sequence from both sides of the
210 splice junction that defines the AS event types described in Figure 1. The tags
211 from a particular AS event (i.e. exon skipping SASTs) from each species were
212 compared to one another using tBLASTx from WU-BLAST (Gish, W. (1996-2003)
213 <http://blast.wustl.edu>) with an E-value cutoff of 1e-5. If the length of two
214 matching tags were both less than 30bp, this matching was not considered. In
215 the cases tags being compared were less than 100 bp, the default tBLASTx E-
216 value cutoff value was used. To classify conserved AS, SASTs flanking a AS event
217 splice junction are required to be conserved between genes across species, and
218 those genes must belong to the same orthogroup. Conserved AS events identified
219 by this process were clustered and each cluster may contain conserved AS events
220 identified between orthologous genes across species (1:1) as well as between
221 paralogs. In an ideal setting, one conserved AS cluster represents one ancient
222 conserved AS event across multiple species. However, if a gene within sorghum
223 is present in two copies (homeologues) in maize, and all three genes harbor a
224 conserved AS event, all three events will be assigned to the same event cluster.
225 Therefore, with respect to the gene in sorghum and the two homeologous copies
226 in maize, this cluster identifies 1 conserved event defined by 2 homeologous
227 events in maize and 1 event in sorghum. Therefore, the total number of events
228 making up a cluster may be larger than the number of species harboring the
229 event. Because we identify these individual events within each species as
230 instances of a conserved event we generalize the AS event cluster as a
231 “conserved AS event”.

232 ***Measure Selection Pressure on Alternative Splicing***

233 We looked for evidence of selection pressure on AS events that were
234 conserved by examining the difference between pairwise non-synonymous
235 substitution rates (K_a) vs. synonymous substitution rates (K_s) for residues
236 defined by exon sequences flanking the conserved AS event in grass species.
237 These are the same regions that were used to construct SASTs (Figure 1). Both
238 upstream and downstream alignments required minimal alignment lengths of 40
239 codons. Selection pressure on AS events at the amino acid levels was measured
240 by determining the K_a/K_s ratios using KaKs_Calculator2.0 through model
241 averaging (Wang et al. 2010). To minimize errors introduced by alignments, or
242 due to paralogs that may result when comparing loci across large phylogenetic
243 distances, the K_a/K_s evaluation of conserved AS events was limited to those
244 events originating within orthologous grass genome loci that also exhibit
245 conserved synteny (defined above).

246 ***Testing the Functional Sharing Hypothesis: Whether Single Copy Genes Have*** 247 ***More Conserved AS Events***

248 The functional sharing hypothesis predicts that conserved AS events
249 would have a higher incidence in single copy genes compared to non-single copy
250 genes. We used a set composed of 2,717 multi-exon genes identified as “strictly”
251 and “mostly” single copy identified across 20 flowering plant genomes by De
252 Smet et al. (2013). We compared the number of single-copy genes exhibiting
253 conserved AS events to the number of non-single copy genes with conserved AS
254 events.

255 ***Conserved AS in SR and hnRNP Proteins***

256 Maize SR protein genes were retrieved from Rauch et al. (2014); Sorghum
257 SR protein genes were retrieved from Richardson et al. (2011). Maize hnRNP
258 protein encoding genes were retrieved from PlantGDB (Duvick et al. 2008). The
259 primary protein sequences of each SR protein gene in maize and sorghum were
260 used for phylogeny analysis. Protein sequences were aligned with MUSCLE
261 (version 3.8.31) using default parameters (Edgar 2004). The alignments were
262 used to construct a maximum likelihood phylogenetic tree with RAxML (version
263 8.1.12) software using a gamma distribution and LG4X model (Stamatakis 2014,
264 Le, Dang, and Gascuel 2012). A manual search along the phylogenetic tree was
265 performed to determine conserved AS events in SR families between maize and
266 sorghum based on conserved isoform structure. The exon-intron gene structure
267 was visualized with Fancygene version 1.4 (Rambaldi and Ciccarelli
268 2009) with AS events further labeled.

269 ***Network Analysis in Stress Responsive Transcription Factors and Conserved*** 270 ***Splicing Genes***

271 3,150 *Arabidopsis* transcription factors (TFs) responsive to 14 diverse
272 stresses were retrieved from STIFDB V2.0 (Naika et al. 2013) and we were
273 focused on five major stresses (ABA, drought, cold, NaCl and light). Protein-
274 protein interaction networks based on transcription factors (TFs) in *Arabidopsis*
275 with conserved AS in response to five major stresses were built with using
276 STRING v10 (Szklarczyk et al. 2014). In addition, a network based on *Arabidopsis*
277 genes demonstrating conserved AS in *Amborella*, *Arabidopsis* and at least one
278 monocot species was also constructed. The cluster of networks constructed was
279 identified and analyzed in Cytoscape v3.3.0 (Shannon et al. 2003).

280 **Go Term Annotation**

281 Genes harboring conserved AS events between *Amborella*, *Arabidopsis*,
282 and at least one monocot species were further evaluated for functional
283 annotation and GO term enrichment based on the *Arabidopsis* gene name using
284 the agriGO analysis toolkit (Du et al. 2010). We use the Fisher's exact test in
285 singular enrichment analysis, followed by multiple test corrections with the
286 Hochberg FDR test. Only associations with a minimum corrected p-value of 0.05
287 were considered.

288 **Data Availability**

289 AS events table and GFF files for the AS isoforms in seven species (*Setaria*
290 *italica*, *Oryza sativa* L. ssp. *Japonica*, *Brachypodium distachyon*, *Musa acuminata*,
291 and *Elaeis guineensis*, the eudicot *Arabidopsis thaliana* and the basal angiosperm
292 *Amborella*) are made available to the community via figshare (figureshare
293 address will made immediately available once the manuscript has been
294 accepted). In addition, a table detailing the AS events identified as conserved
295 during this study will be made available. The *Zea mays* and *Sorghum bicolor* are
296 already available community via figshare
297 (<http://10.6084/m9.figshare.4205079>).

298 **Results**

299 ***Up to 54.6% of Expressed Multi-Exon Genes Exhibit Evidence of AS Among***
300 ***Nine Species***

301 Five types of AS events were considered during this study: Intron
302 retention (IntronR), Alternative acceptor (AltA), Alternative donor (AltD), Exon
303 skipping (ExonS) and Alternate terminal exon (AltTE) (Figure 1). The percentage
304 of expressed multi-exon genes that exhibit AS was determined for each species.
305 The nine species evaluated exhibit AS across a range of 28.7% to 54.6% of multi-
306 exon genes (Table S3). The largest proportion of expressed multi-exon genes
307 with evidence of AS was in *Amborella* (54.6%). Over 60% of the intron retention
308 events in *Amborella* were removed during the filtering process (described in
309 Materials and Methods). Although the number of *Amborella* genes that exhibit AS
310 is similar to that of other species examined during this analysis, *Amborella* has a
311 lower proportion of multi-exon genes in the genome relative to other species.
312 Therefore, despite the similar number of multi-exon genes that undergo AS,
313 *Amborella* exhibits a high fraction of multi-exon genes that produce AS
314 transcripts compared to the other species considered here. Intriguingly,
315 *Amborella* genes that undergo AS produce the highest average number of
316 isoforms/gene compared to other species (Table S3). Whether this feature
317 reflects the evolutionary position of *Amborella* in angiosperm, or is the result of
318 the absence of lineage specific whole genome duplication events – whereas
319 species with additional lineage specific WGD events have lost or partitioned AS
320 isoforms among paralogues through sub-functionalization (Jiang et al. 2013) –
321 remains to be investigated. Maize has the largest collection of RNA-Seq data
322 among the nine species examined and also has the largest number of AS events,
323 but the percentage of multi-exon genes that undergo AS is similar to other
324 species. While the proportions of genes that undergo AS and the total number of
325 AS events is higher than previously reported for many plant species (*Amborella*

326 Genome Project 2013, Campbell et al. 2006, Panahi et al. 2014, Abdel-Ghany et al.
327 2016, Walters et al. 2013), a substantially lower fraction of plant multi-exon
328 genes undergo AS relative to human genes where a reported 95% of multi-exon
329 genes undergo AS (Pan et al. 2008). Similar to previous genome-wide AS studies
330 (Thatcher et al. 2014, Li, Xiao, and Zhu 2014, Mandadi and Scholthof 2015,
331 Chamala et al. 2015), intron retention was the most commonly observed AS
332 event type in monocots. Sorghum and *Amborella* were found to have the lowest
333 percentage of intron retention events, 35.6% and 37.2%, respectively, while
334 African oil palm and *Brachypodium* have the highest proportion of intron
335 retention events, 61.4% and 57.6%, respectively (Table S3).

336 ***AS Events Conserved Between Monocots, Arabidopsis and Amborella***

337 We identified 9,804 conserved AS event clusters in nine species (Table 1).
338 59.0% of conserved AS events represent intron retention. Approximately 68.7%,
339 18.6%, 7.3% and 5.4% conserved AS events are conserved across two species,
340 three species, four species and at least five species, respectively (Table 1). 1,816
341 conserved AS event clusters were found to include *Amborella*. Of these, 1,015
342 involved conserved intron retention events, 450 alternative acceptor events, 205
343 alternative donor events, 114 exon skipping events and 32 alternative terminal
344 exon events (Table 2). In addition, we identified 80 conserved AS events
345 conserved only between *Amborella* and *Arabidopsis* and absent from seven
346 monocot species examined. These 80 AS events likely represent AS events lost
347 after the divergence of eudicots and monocots including 34 intron retention, 31
348 alternative acceptor, 11 alternative donor, and 4 exon skipping events (Table 2).
349 Analysis of nine species revealed 9,804 conserved AS event belong to 19,235

350 genes (Figure 2), with an apparent increased number of genes with conserved AS
351 events within grass species (Figure 2); *Arabidopsis* had the least (944) number of
352 genes harboring conserved AS events while maize had the most (3,687).

353 GO term enrichment analysis was performed on *Arabidopsis* genes with
354 AS events that are conserved within *Amborella*, *Arabidopsis* and monocots to
355 identify whether genes with conserved AS events participated in particular
356 biological processes. 129 AS events meet this criteria. These 129 conserved AS
357 events represent 64 intron retention events, 43 alternative acceptor events, 11
358 alternative donor events, 10 exon skipping events and 1 alternative terminal
359 exon events (Table 2). GO term enrichment based on the *Arabidopsis* genes
360 producing these 129 conserved AS events suggest an over-representation of
361 kinase and phosphorylation activity (Figure S1). This result is consistent with
362 the proposal that AS impacts protein kinase mediated signaling pathways
363 (Reddy et al. 2013), which may modulate the transfer of developmental or
364 environmental cues.

365 TFIIIA possesses an exon skipping event conserved across land plants (Fu
366 et al. 2009, Barbazuk 2010). This event was also identified in all species
367 examined in our study, although our stringent depth requirements initially
368 filtered it out of the banana analysis. In addition, a widely conserved AS event
369 was discovered within SPA3, a member of the SPA protein family, that plays a
370 pivotal role in light signaling as a suppressor of photomorphogenesis (Laubinger
371 and Hoecker 2003). Initial investigation by Shikata et al. (2014) identified an
372 intron retention event and an alternative donor event within *Arabidopsis* SPA3;
373 both splice isoforms harbor premature terminal codons that result in truncated
374 protein products. These truncated proteins retain interaction with

375 CONSTITUTIVE PHOTOMORPHOGENIC 1 (COPI) but lose the ability to bind to
376 DAMAGED DNA-BINDING PROTEIN 1 (DDB1). Evidence suggests the
377 phytochrome could mediate the production of AS transcripts of SPA3 in response
378 to some light conditions, thus promoting plant photomorphogenesis (Shikata et
379 al. 2014). Both the intron retention and alternative donor AS events were
380 identified in our analysis: the alternative donor event is conserved across all nine
381 species examined and the intron retention event is conserved in four species:
382 maize, foxtail millet, African oil palm and *Arabidopsis*. The absence of the intron
383 retention event in the other species examined could be due to limited data depth
384 or tissue sampling. Further investigation is required to determine whether or not
385 this intron retention is conserved broadly across the angiosperms.

386 Alternative donor and acceptor events were reported to be significantly
387 enriched among events conserved between *Brassica* and *Arabidopsis* (Darracq
388 and Adams 2013). We found significant enrichment of alternative donor and
389 alternative acceptor events conserved across at least three species relative to
390 events conserved in only two species ($P < 0.0001$, Fisher's exact test, Table 1),
391 similar enrichment patterns were observed for exon skipping events conserved
392 across at least four species relative to clusters conserved in only two species
393 ($P < 0.05$, Fisher's exact test, Table 1). Expectedly, the overall numbers of AS
394 events that are conserved decreases as the number of species they are conserved
395 within increases. In addition, when considering AS conserved across multiple
396 species, as the number of species that harbor the conserved AS increase, the
397 proportion of conserved AS that are intron retention decrease relative to the
398 other event types. This suggests that requiring an AS event to be conserved
399 between species selects for biologically relevant and important events, and that

400 the trend we observe with intron retention might imply that intron retention
401 suffers a higher proportion of ‘noisy’ or non-relevant splicing relative to other AS
402 types.

403 Examining the lengths of conserved retained introns relative to the non-
404 conserved retained introns across 9 species does not reveal a universal trend. In
405 *Arabidopsis*, *Brachypodium*, and foxtail millet, the length of conserved retained
406 introns is significantly shorter than that of non-conserved (i.e. unique to only one
407 species) retained intron (Table 3). However, in *Amborella* and African oil palm,
408 the length of conserved retained introns is significantly longer than that of non-
409 conserved retained introns. In banana, maize, sorghum and rice, there is no
410 significant difference in retained intron length between conserved and non-
411 conserved retained introns. However in the related grass species maize, rice and
412 sorghum that exhibit large segments of genome co-linearity (Gale and Devos
413 1998) the mean length of conserved retained intron is shorter than non-
414 conserved retained introns, while the median lengths of conserved retained
415 introns are similar or larger than those of non-conserved retained introns (Table
416 3). A previous report suggested that the lengths of retained introns within
417 *Arabidopsis* that are conserved in *Brassica* are shorter than non-conserved
418 retained introns (Darracq and Adams 2013). We observed the same length trend
419 in *Arabidopsis*, but this is not consistent in other species.

420 *Arabidopsis* TFs responsive to 14 diverse stress signals were retrieved
421 from STIFDB V2.0 (Naika et al. 2013). TFs responsive to oxidative stress,
422 dehydration, ABA, NaCl, drought, light, and cold have more AS relative to TFs
423 responsive to the other 7 stresses (Table S4). Since the oxidative stress and
424 dehydration responsive TFs only have 1 and 8 TFs with conserved AS,

425 respectively, we focus on the five major stresses where AS appears to play a
426 significant role: ABA, NaCl, drought, light, and cold. The absolute number and
427 percentage of AS in these TFs responsive to five major stresses are described in
428 Table S4 together with the number and the percentage of conserved AS in these
429 TFs. Two genes are particularly interesting, AT4G27410 (NAC transcription
430 factor RD26) and AT1G78070 (Transducin/WD40 repeat-like superfamily
431 protein). These two genes have conserved AS events and are responsive to all
432 five stresses based on the data (Figure S2). RD26 has a conserved intron
433 retention event across 7 of the 9 species we examined (absent in sorghum and
434 *Brachypodium*). AT1G78070 has a conserved alternative terminal exon between
435 sorghum, banana and *Arabidopsis*. These results indicate that there are AS events
436 in stress-response TFs conserved among species across large phylogenetic
437 distances, which suggests that AS may play an important role in the activity of
438 some TFs during stress response.

439 ***AS Conserved Between Monocot Grass and Non-Grass Species***

440 There are 239 AS events conserved between banana, African oil palm, and
441 at least one of the five grass species examined (Table 2). 34 out of 239 conserved
442 AS events are conserved across banana, African oil palm and five grass species
443 (Table 2). 204 conserved AS events exist within the banana and African oil palm
444 lineages but are absent from the grass species studied. 224 conserved AS events
445 are conserved across the five grass species examined, and these are comprised of
446 60 intron retention events, 106 alternative acceptor events, 37 alternative donor
447 events and 21 exon skipping events (Table 2). Within the PACMAD clade, 1,544,
448 1,488 and 1,323 AS events are conserved between maize vs. sorghum, maize vs.

449 foxtail millet, and sorghum vs. foxtail millet, respectively. In the BEP clade, 1,658
450 conserved AS events are conserved between rice and *Brachypodium*. 4,333
451 conserved AS events are conserved in at least one species in the BEP clade and
452 one species in the PACMAD clade (Table 2). Maize clock genes (i.e circadian
453 rhythm) GRMZM2G033962 (pseudoresponse regulator protein 37, *PRR37*) and
454 GRMZM2G095727 (pseudoresponse regulator protein 73, *PRR73*) share a
455 conserved intron retention that is also conserved across rice, *Brachypodium*, and
456 foxtail millet. Our data suggests that this intron retention event is likely specific
457 to members of the grass family; the role of intron retention in *PRR37* and *PRR73*
458 controlling photoperiodic flowering needs further experimental investigation.

459 Since syntenic genes were more likely to transcribe transcripts and
460 translate proteins than non-syntenic regions (Walley et al. 2016), we tested the
461 hypothesis that AS and conserved AS is enriched in syntenic regions across grass
462 species compared to non-syntenic regions. There are 11,996 syntenic gene
463 clusters across five grass species including maize, sorghum, foxtail millet, rice,
464 and *Brachypodium*. Since maize has undergone an additional whole genome
465 duplication relative to the others, syntenic genes could be maize1 gene, maize2
466 gene or both. In each of five grass species, we compared the proportion of genes
467 within syntenic relationships that undergo AS vs. non-syntenic genes, and the
468 proportion genes within syntenic relationships that have conserved AS events vs.
469 the proportion of non-syntenic genes that have conserved AS events. The
470 proportion of genes with syntenic relationships that exhibit AS is approximately
471 twice that of genes in non-syntenic relationships (Figure 3). This trend is
472 consistent across five grass species. Similarly, the proportion of genes that have
473 conserved AS events that reside in syntenic regions is approximately twice that

474 of genes with conserved AS events that do not reside within syntenic regions
475 (Figure 3). This suggests that genes within conserved syntenic blocks between
476 present within the five grass species studied are enriched in both AS, and
477 conserved AS events.

478 We examined selection pressure at the amino acid level on AS events
479 conserved within grasses by performing pairwise comparisons of Ka/Ks ratios
480 across the exonic regions flanking each conserved AS event (Figure 4). In total
481 4,083 flanking exonic regions were examined, 1,999 from upstream exons and
482 2,084 from downstream exons. We separated the Ka/Ks values for the flanking
483 upstream and downstream portions of each alternative splicing type and plotted
484 these values (Figure 4). In general, Ka/Ks values of exonic regions flanking
485 conserved AS events are smaller than 1 (Figure 4), which suggest these regions
486 are under purifying selection. There are 161 pairwise flanking region pairs with
487 Ka/Ks value above 1, indicative potential candidate for positive selection. These
488 include 67 from intron retention, 41 from alternative acceptor, 40 from
489 alternative donor, 8 from exon skipping and 5 from terminal alternate exon. We
490 did not detect a statistically significant difference in Ka/Ks between exons
491 upstream and downstream of conserved intron retention, alternative acceptor,
492 and exon skipping events. However, a significant difference in Ka/Ks was
493 observed between upstream and downstream exons flanking conserved
494 alternative donor events ($P < 0.001$, Wilcoxon Rank-Sum Test) and alternate
495 terminal exon ($P < 0.001$, Wilcoxon Rank-Sum Test). In general, these results
496 suggest that the regions flanking conserved AS events in grasses are undergoing
497 purifying selection, which may be maintaining *cis*- signals required for AS and
498 somewhat responsible for maintaining sequence conservation across species.

499 ***Conserved AS is More Common Among Single Copy Genes***

500 We tested the functional sharing hypothesis between single copy genes
501 vs. non-single copy genes by measuring the number of genes within each
502 category that exhibit conserved AS events. 2,985 single copy genes including
503 strict and mostly single copy were described in *Arabidopsis* (De Smet et al. 2013).
504 Han et al. (2014) suggested single copy genes had increased levels of AS relative
505 to genes from large gene families and a decreased proportion of AS compared to
506 genes from small gene families. Similarly, Su and Gu (2012) suggested that single
507 copy genes had more isoforms per gene compared to the genes from large gene
508 families but similar average isoform counts per gene compared to the genes from
509 small gene families.

510 We compared the number of conserved AS genes between single copy
511 genes vs. non-single copy genes irrespective of gene family size. 2,717 out of
512 2,985 single copy genes are multi-exon genes. 1,113 out of the 2,717 multi-exon
513 genes have evidence of AS in *Arabidopsis* based on the Araport11 transcripts. In
514 the 19,700 non-single copy multi-exon genes, there are 6,751 that exhibit AS.
515 Among multi-exon genes, the single copy genes have significantly more AS
516 compared to the rest of multi-exon genes ($P < 0.0001$, Chi-square with Yates'
517 correction). Among the 944 *Arabidopsis* genes with conserved AS events in our
518 analysis, 173 were from 1,113 single copy AS genes, while the remaining 771
519 were from the 6,751 non-single copy AS genes. Single copy genes in *Arabidopsis*
520 showed significantly enriched in conserved AS events ($P < 0.0001$, Chi-square
521 with Yates' correction). These results suggest that among multi-exon genes in

522 *Arabidopsis*, a greater proportion of single copy genes both undergo AS and
523 harbor conserved AS events than do non-single copy genes.

524 ***Conserved AS is More Common Among SR Protein Genes Relative to hnRNP***
525 ***Proteins***

526 Similar to previous observations (Kalyna et al. 2006, Rauch et al. 2014,
527 Chamala et al. 2015), AS events were detected in SR proteins. 19 out of 21 SR
528 proteins in maize demonstrate AS, and 13 out of 19 alternatively spliced SR
529 protein genes in maize showed 24 conserved AS events with other species
530 including 13 intron retention events, 4 alternative acceptor events, 4 alternative
531 donor events and 3 exon skipping events (Table S5). 10 out of 24 conserved AS
532 events are present in at least five species. Richardson et al. (2011) identified two
533 SR proteins (GRMZM2G110143 and GRMZM2G170365) in maize with evidence
534 of positive selection (Supplemental Table5). These two genes also have
535 conserved AS events across 2 and 4 species, respectively. Zm-SCL30
536 (GRMZM2G065066), previously identified in maize, responds to both cold and
537 heat response (Mei et al. In Review). SCL30 has alternative acceptor, alternative
538 donor, and exon skipping events conserved across grass species and *Amborella*,
539 as well as two intron retention events only conserved within the grass lineage.
540 Altogether, there are 4 maize SR proteins (2 in the RS subfamily, 1 in the SC
541 subfamily, and 1 in the SCL subfamily) that demonstrate 6 conserved AS events
542 with *Amborella* (Table S5), which supports the existence of a deeply conserved
543 splicing mechanism in SR proteins. 10 out of 40 hnRNP proteins in maize exhibit
544 evidence for 17 conserved AS events, 10 of which are intron retention events, 4
545 alternative acceptor events, 2 alternative donor events, and 1 exon skipping

546 events (Supplemental Table5). None of these 17 conserved AS events are
547 conserved across more than four species. 3 of the 17 conserved AS events from 3
548 hnRNP genes in maize are conserved with *Amborella* (Table S5), while 6 out of
549 24 AS events from 4 SR proteins in maize are conserved with *Amborella*. These
550 results suggest that the signals for AS within SR protein genes across species are
551 better conserved than those within the hnRNP protein encoding genes.

552 ***Conservation and Subfunctionalization of SR Proteins After Whole Genome*** 553 ***Duplication***

554 We used a phylogenetic approach to identify conserved AS across the SR
555 protein subfamilies in maize and sorghum. An ancestral exon skipping event is
556 conserved across all members of this subfamily in maize and sorghum (Figure
557 5A). Intriguingly, all three maize SR subfamily genes were reduced to one copy
558 after whole genome duplication; one gene lost a homeologous copy from maize
559 subgenome1, and two genes lost their homeologous copies from subgenome2.
560 The plant-specific RS subfamily also has a shared exon skipping in the second
561 long intron region coding RRM domain both in maize and sorghum (Figure 5A).
562 The exon skipping isoform generates a complete RRM domain. Kalyna et al.
563 (2006) suggested this long intron is conserved from green algae to angiosperms,
564 and the splice site is retained between monocot and eudicot lineages. Our
565 analysis suggests that this exon skipping event is preserved in maize, sorghum,
566 rice, *Brachypodium*, foxtail millet and *Amborella*. In the plant-specific RSZZ
567 subfamily, maize maintained both copies of the genes after whole genome
568 duplication. Remarkably, two sorghum genes Sobic.009G022100 and
569 Sobic.009G022200 are next to each other but in opposite directions on the

570 chromosome, which suggests one copy may be the result of a tandem
571 duplication. Based on phylogenetic evidence Sobic.009G022200 might be the
572 older copy (Richardson et al. 2011). Not only did Sobic.009G022200 and two
573 duplicate genes in maize retain an exon skipping event, but they also contain two
574 intron retention events surrounding the skipped exon (Figure 5B).
575 Sobic.009G022100 and the two duplicate genes in maize (GRMZM2G099317 and
576 GRMZM2G474658) apparently diverged their AS isoforms (Figure 5B).
577 Sobic.009G022100 only preserved the exon skipping isoform, which would
578 generate a complete RRM domain. GRMZM2G099317 retained the intron
579 retention event on the right side, and the exon skipping event. For the second
580 copy, GRMZM2G474658, the intron retention event on the left side is retained as
581 well as an alternative acceptor event. There is one additional sorghum gene
582 Sobic.003G064400 without conserved synteny in maize, which suggests maize
583 might have lost both syntenic copies that had represented orthologues of this
584 sorghum gene. The AS pattern of this gene is similar to AS pattern of
585 GRMZM2G474658 (Figure 5B). The alternative long intron in the RS2Z subfamily
586 is also conserved from mosses to angiosperms (Kalyna et al. 2006).

587 ***R2R3-MYB is Tightly Connected with a Plant-Specific SR Protein Subfamily***

588 We used a network building approach to identify the degree of
589 interactivity of TFs with conserved AS in response to five stresses (ABA, drought,
590 cold, NaCl and light) based on protein-protein interactions in STRING v10
591 (Szklarczyk et al. 2014). Several interesting relationships showed up in the
592 network (Figure S3). Ubiquitin genes such as UBQ10, UBQ11, and UBQ14 are
593 centralized in the network, which suggests AS might play a role in the

594 susceptibility of translated isoforms to ubiquitin mediated degradation.
595 Ubiquitin genes are closely associated with splicing related genes (such as
596 RS2Z33, RS41, RS2Z32, RS40, U2AF65A, SR45, SCL33) via *Glycine Rich Protein 7*
597 (GRP7). GRP7 is part of the circadian clock and negatively auto-regulates its own
598 protein abundance by producing a non-productive AS isoform that is subject to
599 the nonsense-mediated decay (NMD) pathway (Schöning et al. 2007, Schöning et
600 al. 2008). In addition, we identified the MYB gene network is associated with
601 splicing network genes through the proline-rich spliceosome-associated protein
602 (AT4G21660) (Figure S3). The MYB genes identified in this network belong to
603 the R2R3-MYB family. The R2R3-MYB family has undergone expansion in plants
604 (Du et al. 2014) and plays many important plant-specific roles. Furthermore,
605 there are many kinase and phosphatase related genes in close association with
606 genes responsible for phosphorylation and dephosphorylation of splicing
607 regulators, such as serine-threonine protein kinase (CIPK3), calcium-dependent
608 protein kinase 29 (CPK29), CBL-interacting protein kinase 9 (CIPK9), and
609 calcineurin b-like protein (CBL1). CBL1 has been identified as a salt-tolerance
610 gene and connects with components of the spliceosome (Feng et al. 2015). CBL2
611 and CBL3 together with CIPK3 and CIPK9 can form a calcium network to regulate
612 magnesium levels, while CBL1 is a calcium sensor (Tang et al. 2015). Here, AS of
613 CBL1 and CIPK3, CIPK9 could potentially participate in regulating magnesium
614 levels through a feedback loop sensitive to calcium concentration.

615 We also examined the networks containing genes in *Arabidopsis* with
616 conserved AS events with *Amborella* and at least one species of monocot. Most of
617 the genes present in the network were from the plant-specific SR protein
618 subfamily, such as RS2Z, RS, and SR45a. In addition, TFIIIA is connected to the

619 splicing protein network via RNA-binding protein (AT4G35785). TFIIA has an
620 exon skipping event that is conserved across the land plants (Barbazuk 2010, Fu
621 et al. 2009). Similar to the network based on TFs in response to stress, we
622 identified a close connection between the R2R3-MYB class with a splicing
623 protein network via TFIID-1 (AT3G13445) and proline-rich spliceosome-
624 associated family protein (AT4G21660) (Figure 6). These R2R3-MYBs are
625 involved in many important functions, such as glucosinolate biosynthesis,
626 phenylpropanoid pathway, conical epidermal cell outgrowth, drought and
627 pathogens ABA- and JA- mediated, and cuticular wax biosynthesis (Table S6).

628 **Discussion**

629 In this study, we utilized a bioinformatics scheme to detect AS events and
630 identify those which are conserved across 9 plant species by mining publicly
631 available RNA-Seq data. In total, across all 9 species, 9,804 AS events conserved
632 between 2 or more species are defined within 19,235 genes. New sequenced
633 plant genomes and vast amounts of RNA-Seq are increasingly deposited into
634 NCBI's Sequence Read Archive (SRA). These are treasure troves of data waiting
635 to be mined for new discoveries. Many are already tapping into this wealth, as
636 exemplified by (Nellore et al. 2016) who examined splice junction variants in the
637 human genome, and Chamala et al. (2015) who examined conserved AS events
638 across several species of eudicots. The method we describe here can be used to
639 mine to identify conserved AS events within any species with an available
640 genome sequence and deep RNA-Seq data.

641

642 ***AS Shows Widely Conserved Across Species but Likely Undergone***

643 ***Complicated Patterns of Gain and Loss***

644 Each species examined has one to three thousand genes with evidence of
645 conserved AS. *Arabidopsis* has least number of genes with conserved AS (944),
646 while maize has the greatest number of genes (3,687) (Figure 2). The smallest
647 percentage of conserved AS genes among genes that undergo AS is *Arabidopsis*
648 (12.0%), while foxtail millet has the greatest percentage of conserved AS genes
649 among genes with AS (36.2%). We identified 1,413 genes with conserved AS
650 events in *Amborella*, which account for 18.3% of the genes in *Amborella* with
651 evidence of AS. One conserved exon skipping event occurs within the gene
652 encoding an RNA-binding KH domain-containing protein. This conserved AS
653 event was detected in 8 of the 9 species RNA-Seq data examined, but was
654 apparently absent in *Arabidopsis*. However, this isoform is reported within the
655 AtRTD *Arabidopsis* transcriptome data (Zhang et al. 2015) indicating that it is
656 broadly conserved across all species investigated here. RNA-binding KH domain
657 proteins are vital for heat stress-responsive gene regulation (Guan et al. 2013).
658 The alternatively spliced exon within this RNA-binding KH domain-containing
659 protein has high sequence similarity across nine species studied (Figure 7) and
660 this AS event is likely conserved broadly across angiosperms. Similar to a
661 previous study that identified that alternative donor and alternative acceptor
662 events were found to be significantly enriched between *Arabidopsis* and *Brassica*
663 (Darracq and Adams 2013), we identified a significant enrichment of alternative
664 donor and alternative acceptor events that were conserved in at least three
665 species relative to those conserved in only two species, and a significant
666 enrichment of exon skipping events conserved in at least four species relative to

667 those conserved in only two species. We observed that *Arabidopsis*,
668 *Brachypodium* and foxtail millet retained introns that are conserved have
669 significantly shorter lengths compared to non-conserved retained introns, which
670 is consistent with a previous report in *Arabidopsis* (Table 3) (Darracq and Adams
671 2013). However, this trend is not observed across all the species examined.
672 Conserved retained introns in *Amborella* and African oil palm are longer than
673 those retained introns that were not conserved, and there are no significant
674 differences in the lengths of conserved retained introns vs. non-conserved
675 introns in banana, maize, rice, and sorghum.

676 The plant specific RS and RS2Z subfamilies associate with the R2R3-MYB
677 class within the network built from *Arabidopsis* genes demonstrating conserved
678 AS with *Amborella* and at least one species of monocots. Many R2R3-MYB genes
679 are represented in this network and possess plant-specific functions (Table S6).
680 Li et al. (2006) identified conserved AS events within MYB59 and MYB48. Both
681 genes are represented within our network and associate with SR proteins via a
682 proline-rich spliceosome-associated protein AT4G21660 (Figure 6). Together,
683 these results suggest a conserved connection between splicing factors and R2R3-
684 MYB transcription factors. These two protein classes may cooperate during
685 developmental processes and stress response in plants.

686 ***Grass Lineage Enriched for Conserved Alternative Splicing***

687 Within monocots we identified lineage-specific conserved AS events.
688 There are 204 conserved AS events between banana and African oil palm that
689 are absent in the grass species examined, likely lost in the grass family, and 239
690 conserved AS events conserved in banana, African oil palm and at least one grass

691 species. The number of conserved AS between two species in the PACMAD clade
692 (maize, sorghum, and foxtail millet) is similar to the number of conserved AS
693 events between two species within the BEP clade (rice and *Brachypodium*). 224
694 conserved AS events were found conserved across five grass species examined
695 (maize, sorghum, foxtail millet, rice, and *Brachypodium*). In addition, the
696 proportion of genes within syntenic blocks with AS events in the grasses is
697 approximately twice that compared to those outside of syntenic blocks, and this
698 trend is also mirrored across those genes with conserved AS events (Figure 3).

699 In general, the average values of K_a/K_s across different AS events are less
700 than 1, which suggests that the flanking exonic regions are undergoing purifying
701 selection for conserved AS events. In contrast, previous studies have suggested
702 only moderate selection pressure on the alternative vs. the constitutive regions
703 (Xing and Lee 2006, Chen et al. 2006, Xing and Lee 2005). However, the selection
704 on the flanking region of conserved AS is less explored, particularly in plants. In
705 intron retention, alternative acceptor and exon skipping, there is no significant
706 difference in K_a/K_s between flanking upstream and downstream exonic regions
707 (Figure 4); however, we did observe significantly higher K_a/K_s in the upstream
708 compared to downstream flanking regions in alternative donor and alternative
709 terminal exon. These results indicate that the flanking downstream regions have
710 stronger purifying selection compared to the upstream regions in AltD and ATE
711 events, which suggest that a proportion of flanking upstream and downstream
712 exonic regions of conserved AS events might harbor splicing regulatory
713 elements.

714 We observed both conserved and subfunctionalization patterns of AS
715 within SR protein encoding genes after the whole genome duplication event in

716 the maize lineage. The plant-specific RS subfamily in maize and sorghum harbor
717 conserved exon skipping although the maize orthologue has been reduced to a
718 single copy after whole genome duplication (Figure 5A). In the RS2Z family,
719 Maize has retained both homeologous copies of the RS2Z subfamily genes. In one
720 case, it has maintained the conserved AS events in both maize homeologues
721 including exon skipping and intron retention events (Figure 5B), however, in
722 another case, the AS events have diverged between the sorghum gene and its two
723 maize homoelogs (Figure 5B).

724 ***Important Functional Conserved AS in SR Protein in Response to***
725 ***Environment Stress***

726 Many stress responsive TFs have conserved AS. We focused on identifying
727 conserved AS in stress responsive TFs during five major stresses: ABA, salt
728 (NaCl), drought, light, and cold. Two genes AT1G78070 (Transducin/WD40
729 repeat-like superfamily protein) and AT4G27410 (RD26) are expressed during
730 all five stresses. RD26 encodes a NAC transcription factor induced in response to
731 desiccation. In our data, RD26 has conserved AS in 7 out of the 9 species we
732 examined (except sorghum and *Brachypodium*), however we don't have data to
733 test whether this conserved AS is stress responsive. As illustrated in Figure S4,
734 the exonic flanking sequence of this intron retention event has the highest
735 sequence conservation across species than any other region of the gene,
736 although there remains some similarity between maize and sorghum in the
737 conserved intron region. This conserved intron retention event is likely broadly
738 conserved across angiosperms.

739 Within the protein-protein interaction network based on *Arabidopsis*
740 stress response TFs with conserved AS events identified in our data sets, we
741 detect phosphatases, kinases and several ubiquitin genes in the center of the
742 network (Figure 6). This result is in line with the recent discovery that
743 alternative exons and exons (exon-like introns) have more phosphorylation
744 sites and ubiquitination sites compared to constitutive exons (Marquez et al.
745 2015). In these conserved AS genes across the angiosperms, AS may be involved
746 in signaling pathways responsible for de/phosphorylation and protein
747 degradation pathways, which suggests an important regulatory role in plants.
748 Furthermore, we detected an association between plant-specific RS2Z and RS
749 subfamily splicing proteins and R2R3-MYB proteins via proline-rich
750 spliceosome-associated proteins, suggesting the R2R3-MYB family and splice
751 regulators may interact during abiotic stress responses (Figure 6). Another
752 interesting finding in the network is SPA3 (Figure 6). In moss, red and blue light
753 photoreceptors were shown to regulate AS and intron retention was mis-
754 regulated in moss mutants defective in the red light sensing phytochromes (Wu
755 et al. 2014), which suggests an ancient connection between light regulation and
756 conserved AS in land plants. We found both conserved alternative donor and
757 intron retention events in SPA3, which acts as a negative regulator during light
758 signaling via suppression of photomorphogenesis (Laubinger and Hoecker
759 2003). The truncated protein produced by two AS events in SPA3 will properly
760 interact with COPI but loses the ability to bind to DDB1 (Shikata et al. 2014).
761 SPA3 is connected to 9 additional kinase genes in the network, all of which have
762 conserved AS across the angiosperms (Figure 6). This suggests that AS of kinases
763 might play an important role in the regulation of the light signal, perhaps by

764 affecting a kinase mediated signal cascade linked to SPA3. Inclusive, our study
765 suggests the conserved AS landscape in plants is complicated and needs further
766 functional study to link the phenotype to conserved AS events and identify the
767 regulatory code for functional relevant AS.

768 **Authors' Contributions**

769 WM, WBB designed the work. WM, WBB, LB, GF, JCS analyzed the data. WM and
770 WBB wrote the manuscript with input from LB, GF and JCS.

771 **Acknowledgements**

772 We thank Daniel Gates and Emily Josephs provide helpful comments for the
773 manuscript. This work was supported by Department of Biological Sciences at
774 University of Florida, Florida Genetics Institute, Graduate Student Fellowship
775 and College of Liberal Arts and Sciences Dissertation Fellowship from University
776 of Florida awarded to WM and National Science Foundation grants IOS-0922742
777 & IOS-1547787 (WBB).

778 **Competing Interests**

779 The author(s) declare(s) that they have no competing interests.

Literature Cited

- Abdel-Ghany, Salah E, Michael Hamilton, Jennifer L Jacobi, Peter Ngam, Nicholas Devitt, Faye Schilkey, Asa Ben-Hur, and Anireddy SN Reddy. 2016. "A survey of the sorghum transcriptome using single-molecule long reads." *Nature communications* 7.
- Amborella, Genome Project. 2013. "The Amborella genome and the evolution of flowering plants." *Science* 342 (6165):1241089. doi: 10.1126/science.1241089.
- Barbazuk, W Brad. 2010. "A conserved alternative splicing event in plants reveals an ancient exonization of 5S rRNA that regulates TFIIA." *RNA biology* 7 (4):397-402.
- Bennetzen, Jeffrey L, Jeremy Schmutz, Hao Wang, Ryan Percifield, Jennifer Hawkins, Ana C Pontaroli, Matt Estep, Liang Feng, Justin N Vaughn, and Jane Grimwood. 2012. "Reference genome sequence of the model plant *Setaria*." *Nature biotechnology* 30 (6):555-561.
- Busch, Anke, and Klemens J Hertel. 2012. "Evolution of SR protein and hnRNP splicing regulatory factors." *Wiley Interdisciplinary Reviews: RNA* 3 (1):1-12.
- Campbell, M. A., B. J. Haas, J. P. Hamilton, S. M. Mount, and C. R. Buell. 2006. "Comprehensive analysis of alternative splicing in rice and comparative analyses with *Arabidopsis*." *Bmc Genomics* 7. doi: 10.1186/1471-2164-7-327.
- Chamala, Srikar, Guanqiao Feng, Carolina Chavarro, and W Brad Barbazuk. 2015. "Genome-wide identification of evolutionarily conserved alternative

splicing events in flowering plants." *Frontiers in bioengineering and biotechnology* 3.

Chen, Feng-Chi, Sheng-Shun Wang, Chuang-Jong Chen, Wen-Hsiung Li, and Trees-Juen Chuang. 2006. "Alternatively and constitutively spliced exons are subject to different evolutionary forces." *Molecular biology and evolution* 23 (3):675-682.

Cheng, CHIA-YI, Vivek Krishnakumar, Agnes Chan, Seth Schobel, and Christopher D Town. 2016. "Araport11: a complete reannotation of the Arabidopsis thaliana reference genome." *bioRxiv*:047308.

Chuang, Trees-Juen, Min-Yu Yang, Chuang-Chieh Lin, Ping-Hung Hsieh, and Li-Yuan Hung. 2015. "Comparative genomics of grass EST libraries reveals previously uncharacterized splicing events in crop plants." *BMC plant biology* 15 (1):1.

Cotton, Joseph L, William P Wysocki, Lynn G Clark, Scot A Kelchner, J Chris Pires, Patrick P Edger, Dustin Mayfield-Jones, and Melvin R Duvall. 2015. "Resolving deep relationships of PACMAD grasses: a phylogenomic approach." *BMC plant biology* 15 (1):1.

D'Hont, Angélique, France Denoeud, Jean-Marc Aury, Franc-Christophe Baurens, Françoise Carreel, Olivier Garsmeur, Benjamin Noel, Stéphanie Bocs, Gaëtan Droc, and Mathieu Rouard. 2012. "The banana (*Musa acuminata*) genome and the evolution of monocotyledonous plants." *Nature* 488 (7410):213-217.

Darracq, Aude, and Keith L. Adams. 2013. "Features of evolutionarily conserved alternative splicing events between Brassica and Arabidopsis." *New Phytologist* 199 (1):252-263. doi: 10.1111/nph.12238.

- De Smet, Riet, Keith L Adams, Klaas Vandepoele, Marc CE Van Montagu, Steven Maere, and Yves Van de Peer. 2013. "Convergent gene loss following gene and genome duplications creates single-copy families in flowering plants." *Proceedings of the National Academy of Sciences* 110 (8):2898-2903.
- Droc, Gaetan, Delphine Lariviere, Valentin Guignon, Nabila Yahiaoui, Dominique This, Olivier Garsmeur, Alexis Dereeper, Chantal Hamelin, Xavier Argout, and Jean-François Dufayard. 2013. "The banana genome hub." *Database* 2013:bat035.
- Du, Hai, Zhe Liang, Sen Zhao, Ming-Ge Nan, LS Tran, Kun Lu, Yu-Bi Huang, and Jia-Na Li. 2014. "The Evolutionary History of R2R3-MYB Proteins Across 50 Eukaryotes: New Insights Into Subfamily Classification and Expansion." *Scientific reports* 5:11037-11037.
- Du, Zhou, Xin Zhou, Yi Ling, Zhenhai Zhang, and Zhen Su. 2010. "agriGO: a GO analysis toolkit for the agricultural community." *Nucleic acids research*:gkq310.
- Duvick, Jon, Ann Fu, Usha Muppирala, Mukul Sabharwal, Matthew D. Wilkerson, Carolyn J. Lawrence, Carol Lushbough, and Volker Brendel. 2008. "PlantGDB: a resource for comparative plant genomics." *Nucleic Acids Research* 36:D959-D965. doi: 10.1093/nar/gkm1041.
- Edgar, Robert C. 2004. "MUSCLE: multiple sequence alignment with high accuracy and high throughput." *Nucleic acids research* 32 (5):1792-1797.
- Emms, David M, and Steven Kelly. 2015. "OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy." *Genome biology* 16 (1):1-14.

- Feng, Jinlin, Jingjing Li, Zhaoxu Gao, Yaru Lu, Junya Yu, Qian Zheng, Shuning Yan, Wenjiao Zhang, Hang He, and Ligeng Ma. 2015. "SKIP confers osmotic tolerance during salt stress by controlling alternative gene splicing in Arabidopsis." *Molecular plant* 8 (7):1038-1052.
- Filichkin, S. A., H. D. Priest, S. A. Givan, R. K. Shen, D. W. Bryant, S. E. Fox, W. K. Wong, and T. C. Mockler. 2010. "Genome-wide mapping of alternative splicing in Arabidopsis thaliana." *Genome Research* 20 (1):45-58. doi: 10.1101/gr.093302.109.
- Fu, Yan, Oliver Bannach, Hao Chen, Jan-Hendrik Teune, Axel Schmitz, Gerhard Steger, Liming Xiong, and W Brad Barbazuk. 2009. "Alternative splicing of anciently exonized 5S rRNA regulates plant transcription factor TFIIIA." *Genome research* 19 (5):913-921.
- Gale, Michael D, and Katrien M Devos. 1998. "Comparative genetics in the grasses." *Proceedings of the National Academy of Sciences* 95 (5):1971-1974.
- Goff, Stephen A, Darrell Ricke, Tien-Hung Lan, Gernot Presting, Ronglin Wang, Molly Dunn, Jane Glazebrook, Allen Sessions, Paul Oeller, and Hemant Varma. 2002. "A draft sequence of the rice genome (*Oryza sativa* L. ssp. japonica)." *Science* 296 (5565):92-100.
- Goodstein, David M, Shengqiang Shu, Russell Howson, Rochak Neupane, Richard D Hayes, Joni Fazo, Therese Mitros, William Dirks, Uffe Hellsten, and Nicholas Putnam. 2012. "Phytozome: a comparative platform for green plant genomics." *Nucleic acids research* 40 (D1):D1178-D1186.
- Grabherr, M. G., B. J. Haas, M. Yassour, J. Z. Levin, D. A. Thompson, I. Amit, X. Adiconis, L. Fan, R. Raychowdhury, Q. D. Zeng, Z. H. Chen, E. Mauceli, N.

- Hacohen, A. Gnirke, N. Rhind, F. di Palma, B. W. Birren, C. Nusbaum, K. Lindblad-Toh, N. Friedman, and A. Regev. 2011. "Full-length transcriptome assembly from RNA-Seq data without a reference genome." *Nature Biotechnology* 29 (7):644-U130. doi: 10.1038/nbt.1883.
- Guan, Qingmei, Changlong Wen, Haitao Zeng, and Jianhua Zhu. 2013. "A KH domain-containing putative RNA-binding protein is critical for heat stress-responsive gene regulation and thermotolerance in Arabidopsis." *Molecular plant* 6 (2):386-395.
- Haas, Brian J, Arthur L Delcher, Stephen M Mount, Jennifer R Wortman, Roger K Smith Jr, Linda I Hannick, Rama Maiti, Catherine M Ronning, Douglas B Rusch, and Christopher D Town. 2003. "Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies." *Nucleic acids research* 31 (19):5654-5666.
- Hammond, Ming C, Andreas Wachter, and Ronald R Breaker. 2009. "A plant 5S ribosomal RNA mimic regulates alternative splicing of transcription factor IIIA pre-mRNAs." *Nature structural & molecular biology* 16 (5):541-549.
- Han, Fengming, Yong Peng, Lijia Xu, and Peigen Xiao. 2014. "Identification, characterization, and utilization of single copy genes in 29 angiosperm genomes." *BMC genomics* 15 (1):504.
- James, Allan B, Naeem Hasan Syed, Simon Bordage, Jacqueline Marshall, Gillian A Nimmo, Gareth I Jenkins, Pawel Herzyk, John WS Brown, and Hugh G Nimmo. 2012. "Alternative splicing mediates responses of the Arabidopsis circadian clock to temperature changes." *The Plant Cell* 24 (3):961-981.

- Jiang, Wen-kai, Yun-long Liu, En-hua Xia, and Li-zhi Gao. 2013. "Prevalent role of gene features in determining evolutionary fates of whole-genome duplication duplicated genes in flowering plants." *Plant physiology* 161 (4):1844-1861.
- Kalyna, Maria, Sergiy Lopato, Viktor Voronin, and Andrea Barta. 2006. "Evolutionary conservation and regulation of particular alternative splicing events in plant SR proteins." *Nucleic acids research* 34 (16):4395-4405.
- Laubinger, Sascha, and Ute Hoecker. 2003. "The SPA1-like proteins SPA3 and SPA4 repress photomorphogenesis in the light." *The Plant Journal* 35 (3):373-385.
- Lawrence, C. J., O. F. Dong, M. L. Polacco, T. E. Seigfried, and V. Brendel. 2004. "MaizeGDB, the community database for maize genetics and genomics." *Nucleic Acids Research* 32:D393-D397. doi: 10.1093/nar/gkh011.
- Le, Si Quang, Cuong Cao Dang, and Olivier Gascuel. 2012. "Modeling protein evolution with several amino acid replacement matrices depending on site rates." *Molecular biology and evolution*:mss112.
- Lee, Yeon, and Donald C Rio. 2015. "Mechanisms and regulation of alternative pre-mRNA splicing." *Annual review of biochemistry* 84:291.
- Li, Jigang, Xiaojuan Li, Lei Guo, Feng Lu, Xiaojie Feng, Kun He, Liping Wei, Zhangliang Chen, Li-Jia Qu, and Hongya Gu. 2006. "A subgroup of MYB transcription factor genes undergoes highly conserved alternative splicing in Arabidopsis and rice." *Journal of Experimental Botany* 57 (6):1263-1273.

- Li, Qin, Guanghui Xiao, and Yu-Xian Zhu. 2014. "Single-nucleotide resolution mapping of the *Gossypium raimondii* transcriptome reveals a new mechanism for alternative splicing of introns." *Molecular plant* 7 (5):829-840.
- Mandadi, Kranthi K, and Karen-Beth G Scholthof. 2015. "Genome-wide analysis of alternative splicing landscapes modulated during plant-virus interactions in *Brachypodium distachyon*." *The Plant Cell* 27 (1):71-85.
- Marquez, Y., J. W. S. Brown, C. Simpson, A. Barta, and M. Kalyna. 2012. "Transcriptome survey reveals increased complexity of the alternative splicing landscape in *Arabidopsis*." *Genome Research* 22 (6):1184-1195. doi: 10.1101/gr.134106.111.
- Marquez, Yamile, Markus Höpfler, Zahra Ayatollahi, Andrea Barta, and Maria Kalyna. 2015. "Unmasking alternative splicing inside protein-coding exons defines exons and their role in proteome plasticity." *Genome research*:gr. 186585.114.
- Mastrangelo, Anna M, Daniela Marone, Giovanni Laidò, Anna M De Leonardis, and Pasquale De Vita. 2012. "Alternative splicing: enhancing ability to cope with stress via transcriptome plasticity." *Plant Science* 185:40-49.
- Naika, Mahantesha, Khader Shameer, Oommen K Mathew, Ramanjini Gowda, and Ramanathan Sowdhamini. 2013. "STIFDB2: an updated version of plant stress-responsive transcription factor database with additional stress signals, stress-responsive transcription factor binding sites and stress-responsive genes in *Arabidopsis* and rice." *Plant and Cell Physiology* 54 (2):e8-e8.

- Nellore, Abhinav, Andrew E Jaffe, Jean-Philippe Fortin, José Alquicira-Hernández, Leonardo Collado-Torres, Siruo Wang, Robert A Phillips, Nishika Karbhari, Kasper D Hansen, Ben Langmead, and Jeffrey T Leek. 2016. "Human splicing diversity across the Sequence Read Archive." *bioRxiv*. doi: 10.1101/038224.
- Pan, Qun, Ofer Shai, Leo J Lee, Brendan J Frey, and Benjamin J Blencowe. 2008. "Deep surveying of alternative splicing complexity in the human transcriptome by high-throughput sequencing." *Nature genetics* 40 (12):1413-1415.
- Panahi, Bahman, Bahram Abbaszadeh, Mehdi Taghizadeghan, and Esmaeil Ebrahimie. 2014. "Genome-wide survey of alternative splicing in Sorghum bicolor." *Physiology and Molecular Biology of Plants* 20 (3):323-329.
- Paterson, Andrew H, John E Bowers, Remy Bruggmann, Inna Dubchak, Jane Grimwood, Heidrun Gundlach, Georg Haberer, Uffe Hellsten, Therese Mitros, and Alexander Poliakov. 2009. "The Sorghum bicolor genome and the diversification of grasses." *Nature* 457 (7229):551-556.
- Pertea, Mihaela, Geo M. Pertea, Corina M. Antonescu, Tsung-Cheng Chang, Joshua T. Mendell, and Steven L. Salzberg. 2015. "StringTie enables improved reconstruction of a transcriptome from RNA-seq reads." *Nature Biotechnology* 33 (3):290-+. doi: 10.1038/nbt.3122.
- Rambaldi, Davide, and Francesca D Ciccarelli. 2009. "FancyGene: dynamic visualization of gene structures and protein domain architectures on genomic loci." *Bioinformatics* 25 (17):2281-2282.
- Rauch, Hypaitia B, Tara L Patrick, Katarina M Klusman, Fabia U Battistuzzi, Wenbin Mei, Volker P Brendel, and Shailesh K Lal. 2014. "Discovery and

expression analysis of alternative splicing events conserved among plant SR proteins." *Molecular biology and evolution* 31 (3):605-613.

Reddy, Anireddy S. N., Mark F. Rogers, Dale N. Richardson, Michael Hamilton, and Asa Ben-Hur. 2012. "Deciphering the plant splicing code: experimental and computational approaches for predicting alternative splicing and splicing regulatory elements." *Frontiers in plant science* 3:18-18. doi: 10.3389/fpls.2012.00018.

Reddy, Anireddy SN, Yamile Marquez, Maria Kalyna, and Andrea Barta. 2013. "Complexity of the alternative splicing landscape in plants." *The Plant Cell* 25 (10):3657-3683.

Richardson, Dale N, Mark F Rogers, Adam Labadorf, Asa Ben-Hur, Hui Guo, Andrew H Paterson, and Anireddy SN Reddy. 2011. "Comparative analysis of serine/arginine-rich proteins across 27 eukaryotes: insights into sub-family classification and extent of alternative splicing." *PLoS One* 6 (9):e24542.

Schnable, James C, Michael Freeling, and Eric Lyons. 2012. "Genome-wide analysis of syntenic gene deletion in the grasses." *Genome biology and evolution* 4 (3):265-277.

Schnable, Patrick S, Doreen Ware, Robert S Fulton, Joshua C Stein, Fusheng Wei, Shiran Pasternak, Chengzhi Liang, Jianwei Zhang, Lucinda Fulton, and Tina A Graves. 2009. "The B73 maize genome: complexity, diversity, and dynamics." *science* 326 (5956):1112-1115.

Schöning, Jan C, Corinna Streitner, Irmtraud M Meyer, Yahong Gao, and Dorothee Staiger. 2008. "Reciprocal regulation of glycine-rich RNA-binding proteins via an interlocked feedback loop coupling alternative splicing to

nonsense-mediated decay in Arabidopsis." *Nucleic acids research* 36 (22):6977-6987.

Schöning, Jan C, Corinna Streitner, Damian R Page, Sven Hennig, Kenko Uchida, Eva Wolf, Masaki Furuya, and Dorothee Staiger. 2007. "Auto-regulation of the circadian slave oscillator component AtGRP7 and regulation of its targets is impaired by a single RNA recognition motif point mutation." *The Plant Journal* 52 (6):1119-1130.

Severing, Edouard I, Aalt DJ Dijk, Willem J Stiekema, and Roeland CHJ Ham. 2009. "Comparative analysis indicates that alternative splicing in plants has a limited role in functional expansion of the proteome." *BMC genomics* 10 (1):1.

Shannon, Paul, Andrew Markiel, Owen Ozier, Nitin S Baliga, Jonathan T Wang, Daniel Ramage, Nada Amin, Benno Schwikowski, and Trey Ideker. 2003. "Cytoscape: a software environment for integrated models of biomolecular interaction networks." *Genome research* 13 (11):2498-2504.

Shen, Yanting, Zhengkui Zhou, Zheng Wang, Weiyu Li, Chao Fang, Mian Wu, Yanming Ma, Tengfei Liu, Ling-An Kong, and De-Liang Peng. 2014. "Global dissection of alternative splicing in paleopolyploid soybean." *The Plant Cell* 26 (3):996-1008.

Shikata, Hiromasa, Kousuke Hanada, Tomokazu Ushijima, Moeko Nakashima, Yutaka Suzuki, and Tomonao Matsushita. 2014. "Phytochrome controls alternative splicing to mediate light responses in Arabidopsis." *Proceedings of the National Academy of Sciences* 111 (52):18781-18786.

- Singh, Rajinder, Meilina Ong-Abdullah, Eng-Ti Leslie Low, Mohamad Arif Abdul Manaf, Rozana Rosli, Rajanaidu Nookiah, Leslie Cheng-Li Ooi, Siew-Eng Ooi, Kuang-Lim Chan, and Mohd Amin Halim. 2013. "Oil palm genome sequence reveals divergence of interfertile species in Old and New worlds." *Nature* 500 (7462):335-339.
- Staiger, Dorothee, and John WS Brown. 2013. "Alternative splicing at the intersection of biological timing, development, and stress responses." *The Plant Cell* 25 (10):3640-3656.
- Stamatakis, Alexandros. 2014. "RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies." *Bioinformatics*:btu033.
- Streitner, Corinna, Craig G Simpson, Paul Shaw, Selahattin Danisman, John WS Brown, and Dorothee Staiger. 2013. "Small changes in ambient temperature affect alternative splicing in *Arabidopsis thaliana*." *Plant signaling & behavior* 8 (7):11240-55.
- Sturgill, David, John H Malone, Xia Sun, Harold E Smith, Leonard Rabinow, Marie-Laure Samson, and Brian Oliver. 2013. "Design of RNA splicing analysis null models for post hoc filtering of *Drosophila* head RNA-Seq data with the splicing analysis kit (Spanki)." *BMC bioinformatics* 14 (1):1.
- Su, Z. X., J. M. Wa, J. Yu, X. Q. Huang, and X. Gu. 2006. "Evolution of alternative splicing after gene duplication." *Genome Research* 16 (2):182-189. doi: 10.1101/gr.4197006.
- Su, Zhixi, and Xun Gu. 2012. "Revisit on the evolutionary relationship between alternative splicing and gene duplication." *Gene* 504 (1):102-106.
- Syed, Naeem H, Silvas J Prince, Raymond N Mutava, Gunvant Patil, Song Li, Wei Chen, Valliyodan Babu, Trupti Joshi, Saad Khan, and Henry T Nguyen.

2015. "Core clock, SUB1, and ABAR genes mediate flooding and drought responses via alternative splicing in soybean." *Journal of experimental botany* 66 (22):7129-7149.

Szklarczyk, Damian, Andrea Franceschini, Stefan Wyder, Kristoffer Forslund, Davide Heller, Jaime Huerta-Cepas, Milan Simonovic, Alexander Roth, Alberto Santos, and Kalliopi P Tsafou. 2014. "STRING v10: protein-protein interaction networks, integrated over the tree of life." *Nucleic acids research*:gku1003.

Tang, Ren-Jie, Fu-Geng Zhao, Veder J Garcia, Thomas J Kleist, Lei Yang, Hong-Xia Zhang, and Sheng Luan. 2015. "Tonoplast CBL-CIPK calcium signaling network regulates magnesium homeostasis in Arabidopsis." *Proceedings of the National Academy of Sciences* 112 (10):3134-3139.

Thatcher, Shawn R, Olga N Danilevskaya, Xin Meng, Mary Beatty, Gina Zastrow-Hayes, Charlotte Harris, Brandon Van Allen, Jeffrey Habben, and Bailin Li. 2016. "Genome-Wide Analysis of Alternative Splicing during Development and Drought Stress in Maize." *Plant Physiology* 170 (1):586-599.

Thatcher, Shawn R, Wengang Zhou, April Leonard, Bing-Bing Wang, Mary Beatty, Gina Zastrow-Hayes, Xiangyu Zhao, Andy Baumgarten, and Bailin Li. 2014. "Genome-wide analysis of alternative splicing in Zea mays: landscape and genetic regulation." *The Plant Cell* 26 (9):3472-3487.

Trapnell, Cole, Adam Roberts, Loyal Goff, Geo Pertea, Daehwan Kim, David R. Kelley, Harold Pimentel, Steven L. Salzberg, John L. Rinn, and Lior Pachter. 2012. "Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks." *Nature Protocols* 7 (3):562-578. doi: 10.1038/nprot.2012.016.

- Vogel, John P, David F Garvin, Todd C Mockler, Jeremy Schmutz, Dan Rokhsar, Michael W Bevan, Kerrie Barry, Susan Lucas, Miranda Harmon-Smith, and Kathleen Lail. 2010. "Genome sequencing and analysis of the model grass *Brachypodium distachyon*." *Nature* 463 (7282):763-768.
- Walley, Justin W, Ryan C Sartor, Zhouxin Shen, Robert J Schmitz, Kevin J Wu, Mark A Urich, Joseph R Nery, Laurie G Smith, James C Schnable, and Joseph R Ecker. 2016. "Integration of omic networks in a developmental atlas of maize." *Science* 353 (6301):814-818.
- Walters, Braden, Gengkon Lum, Gaurav Sablok, and Xiang Jia Min. 2013. "Genome-wide landscape of alternative splicing events in *Brachypodium distachyon*." *DNA research* 20 (2):163-171.
- Wang, Dapeng, Yubin Zhang, Zhang Zhang, Jiang Zhu, and Jun Yu. 2010. "KaKs_Calculator 2.0: a toolkit incorporating gamma-series methods and sliding window strategies." *Genomics, Proteomics & Bioinformatics* 8 (1):77-80.
- Wang, Zefeng, and Christopher B Burge. 2008. "Splicing regulation: from a parts list of regulatory elements to an integrated splicing code." *Rna* 14 (5):802-813.
- Wu, Hshin-Ping, YS Su, Hsiu-Chen Chen, Yu-Rong Chen, Chia-Chen Wu, Wen-Dar Lin, and Shih-Long Tu. 2014. "Genome-wide analysis of light-regulated alternative splicing mediated by photoreceptors in *Physcomitrella patens*." *Genome Biol* 15 (1):R10.
- Xing, Y., and C. Lee. 2005. "Evidence of functional selection pressure for alternative splicing events that accelerate evolution of protein subsequences." *Proceedings of the National Academy of Sciences of the*

United States of America 102 (38):13526-13531. doi:

10.1073/pnas.0501313102.

- Xing, Yi, and Christopher Lee. 2006. "Alternative splicing and RNA selection pressure—evolutionary consequences for eukaryotic genomes." *Nature Reviews Genetics* 7 (7):499-509.
- Xu, Peng, Yimeng Kong, Dongliang Song, Cheng Huang, Xuan Li, and Laigeng Li. 2014. "Conservation and functional influence of alternative splicing in wood formation of *Populus* and *Eucalyptus*." *BMC genomics* 15 (1):780.
- Zhang, Runxuan, Cristiane PG Calixto, Nikoleta A Tzioutziou, Allan B James, Craig G Simpson, Wenbin Guo, Yamile Marquez, Maria Kalyna, Rob Patro, and Eduardo Eyra. 2015. "AtRTD—a comprehensive reference transcript dataset resource for accurate quantification of transcript-specific expression in *Arabidopsis thaliana*." *New Phytologist* 208 (1):96-101.
- Zhao, Lei, Ning Zhang, Peng-Fei Ma, Qi Liu, De-Zhu Li, and Zhen-Hua Guo. 2013. "Phylogenomic analyses of nuclear genes reveal the evolutionary relationships within the BEP clade and the evidence of positive selection in Poaceae." *PloS one* 8 (5):e64642.

Table 1. Conserved alternative splicing across 9 species at the gene family level.

Percentage is based on total number of conserved alternative splicing clusters.

Number of species	AltA	AltD	ATE	ExonS	IntronR	Total/%
2	1,384	684	201	335	4,128	6,732/68.7
3	447	181	32	69	1,099	1828/18.6
4	218	112	6	26	351	713/7.3
5	117	36	2	21	109	285/2.9
6	58	29	0	10	52	149/1.5
7	22	12	0	3	32	69/0.7
8	7	3	0	4	7	21/0.2
9	2	2	0	1	2	7/0.1
Total/%	2,255/23.0	1,059/10.8	241/2.5	469/4.8	5,780/59.0	9,804

Table 2. Conserved alternative splicing clusters among different species and clades.

* Conserved at least one species in PCAMAD clade and one species in BEP clade.

	Total	IntronR	AltA	AltD	ExonS	ATE
Conserved between <i>Amborella</i> and other species	1,816	1,015	450	205	114	32
Conserved in <i>Amborella</i> , <i>Arabidopsis</i> , and at least one species of Monocots	129	64	43	11	10	1
Conserved between <i>Amborella</i> and <i>Arabidopsis</i> not in Monocots	80	34	31	11	4	0
Conserved across <i>Amborella</i> and seven Monocot not in <i>Arabidopsis</i>	9	2	3	2	2	0
Conserved in banana, African oil plum and at least one grass species	239	162	44	20	11	2
Conserved across five species in grass	224	60	106	37	21	0
Conserved across seven species in monocots	34	13	12	5	4	0
Conserved between PACMAD and BEP clades*	4333	2548	1052	480	177	76

Table 3. A comparison of intron length in the conserved intron retention events vs. non-conserved intron retention events across the species examined. (NS, Not significant; *, P < 0.05; ***, P < 0.001)

	<i>Amborella</i>	<i>Arabidopsis</i>	Banana	<i>Brachypodium</i>	Foxtail millet	Maize	Oil Palm	Rice	Sorghum
Conserved AS									
Mean	670	165	329	411	402	440	315	417	409
Median	390	96	143	233	237	207	131	269	256
No evidence for conserved AS									
Mean	651	186	327	474	410	471	287	486	437
Median	316	113	131	245	201	205	121	247	242
Wilcoxon test									
P value	***	***	NS	***	*	NS	*	NS	NS

Figure Legends

Figure 1. Pipeline to identify conserved alternative splicing events.

Up to 300bp on either side of a splice junction were extracted. These extracted sequences were binned together with all other similarly extracted and binned sequences flanking the same type of event for each species (eg. one bin with sequences flanking junctions involved in intron retention, one bin with sequences flanking junctions involved in exon skipping, etc.). Sequences within each bin were compared with tBLASTx to bins of the same event from other species to identify similar tags. Cases where the sequence tags were similar and the parent genes are from different species but represent orthologs as evidenced by being within the same orthogroup (see Methods) defined conserved alternative splicing events.

Figure 2. Genes with conserved alternative splicing events across nine species.

We calculate the number of genes with conserved alternative splicing events in each species. Tetraploidy and Hexaploidy are labeled on the phylogenetic trees. The length of the phylogenetic tree is not proportional to the phylogenetic distance.

Figure 3. Alternative splicing and conserved alternative splicing enriched across grass syntenic genes.

We calculated the percentage of AS in syntenic and non-syntenic genes (the percentage of AS genes are labeled in blue). In addition, we calculated the

percentage of conserved AS in syntenic AS genes and non-syntenic AS genes (the percentage of conserved AS genes are labeled in black).

Figure 4. Violin plot of Ka/Ks value from flanking upstream and downstream in different alternative splicing events across syntenic genes within grass species.

We calculated the pairwise Ka/Ks ratio within exonic regions upstream and downstream of conserved alternative splicing events in five grass species. Only those exonic regions aligned at $a \geq 120$ bp were considered in Ka/Ks calculation. The dotted lines within the violin plot represented the first, second and third quartile.

Figure 5. Phylogenetic tree and gene structure of SR, RS and RS2Z subfamilies in maize and sorghum.

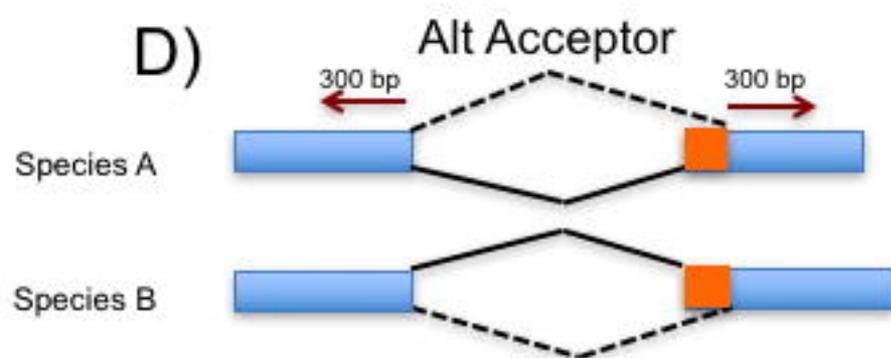
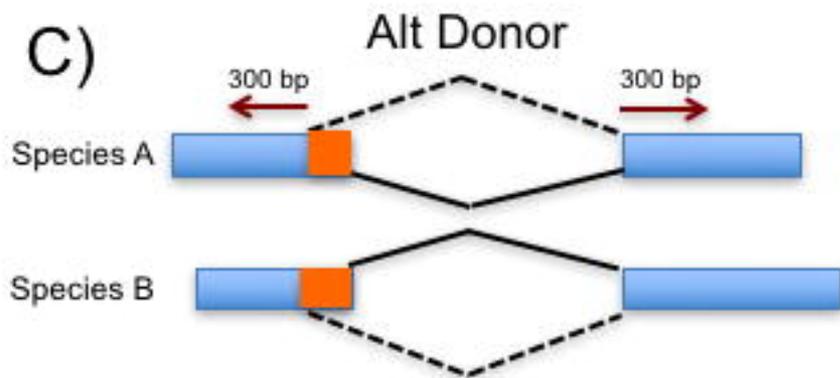
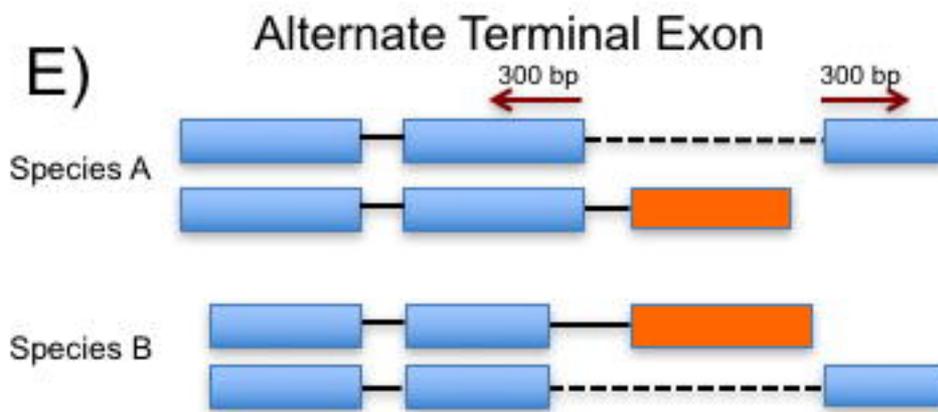
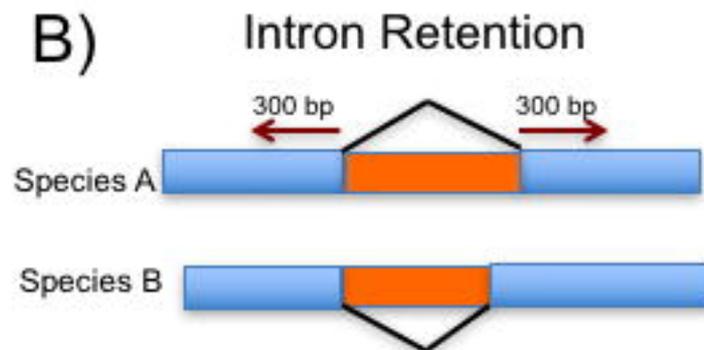
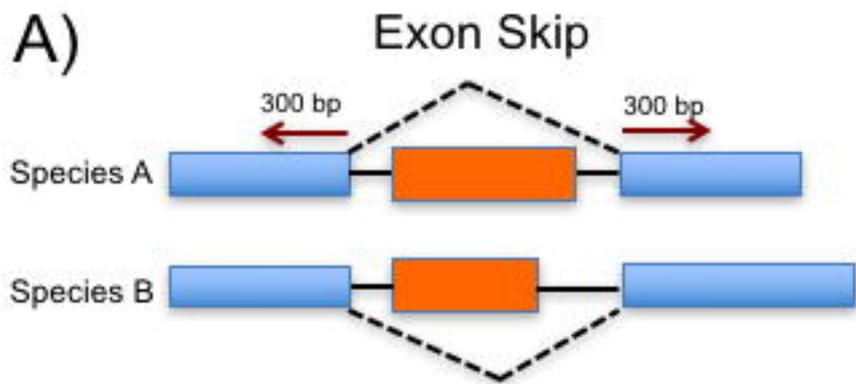
A) Phylogenetic tree representing the relationship of genes from SR and RS subfamilies in maize and sorghum; associated gene structures are illustrated by Fancygene to the right. Constitutive exons are shaded blue and alternative region are indicated with orange. Transcript structure is indicated by exons connected via dashed lines; orthologous exons are indicated with grey lines. B) Phylogenetic tree and gene structure of RS2Z subfamily in maize and sorghum are illustrated by Fancygene. Annotation and shading is described above; intron retention regions are indicated in clear boxes.

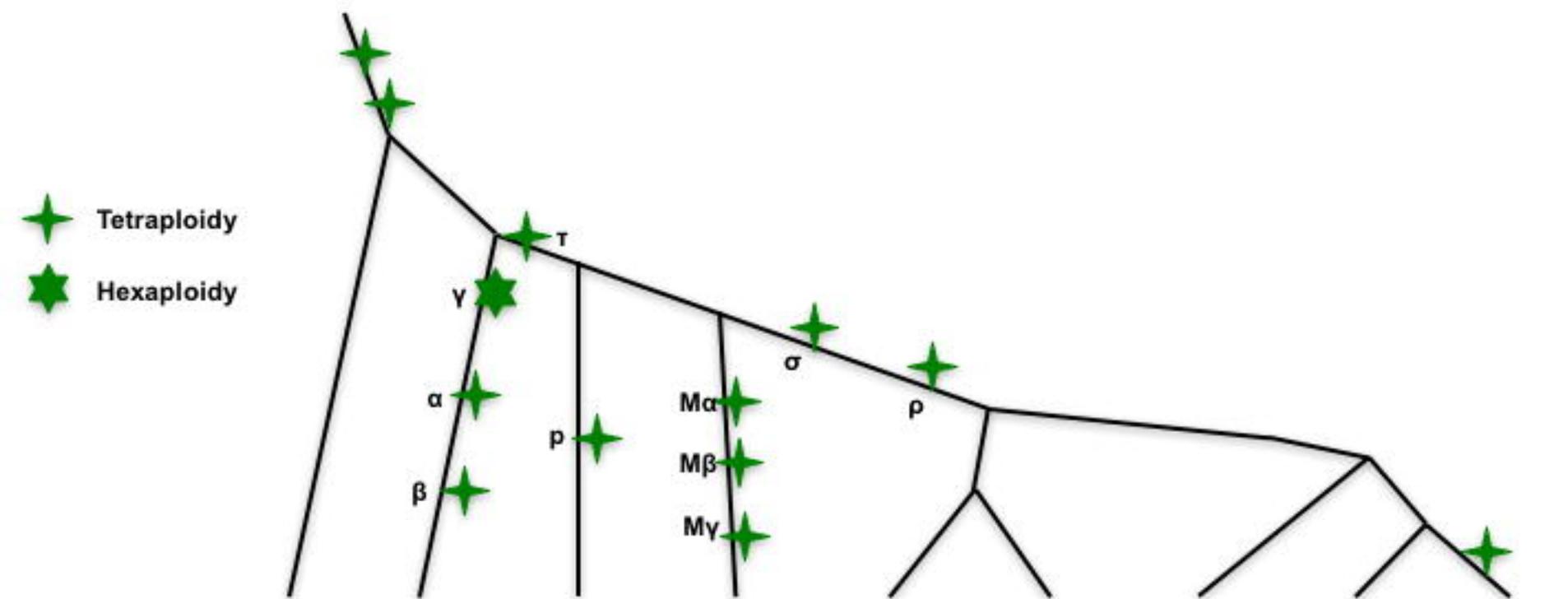
Figure 6. Protein-protein network of *Arabidopsis* genes with evidence of conserved alternative splicing with *Amborella* and at least one species of monocot.

The arrow in the network indicates the direction of the protein-protein network.

Figure 7. Sequence conservation within a gene encoding an RNA-binding KH domain-containing protein that produces an exon skipping even conserved across species.

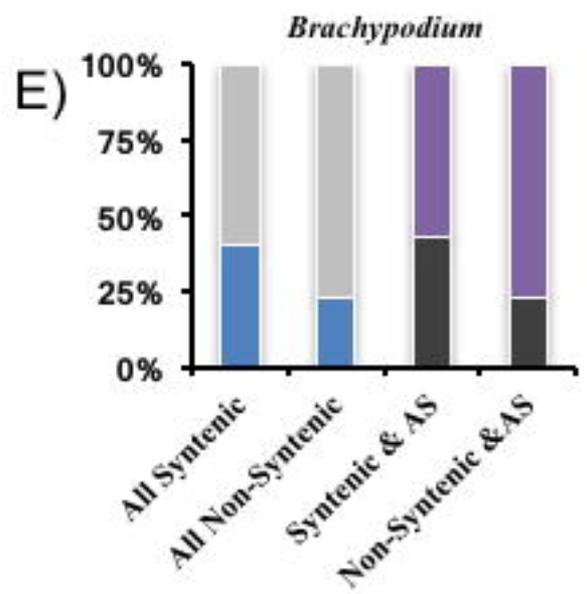
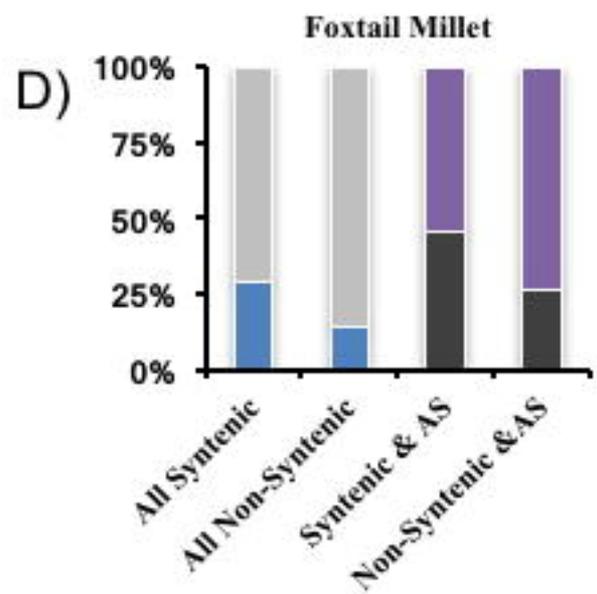
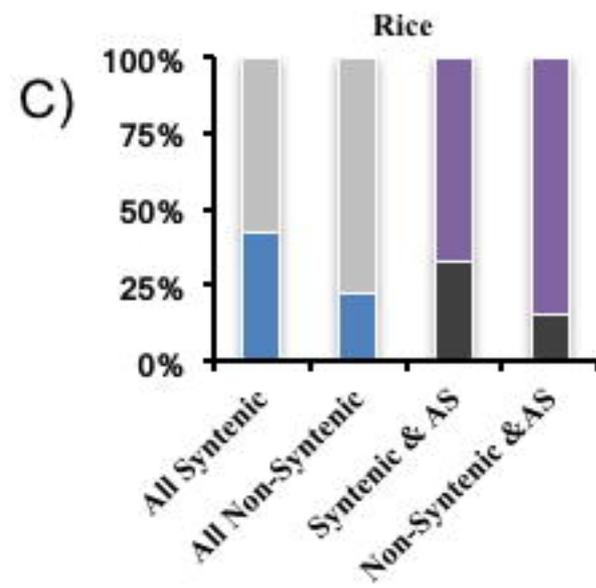
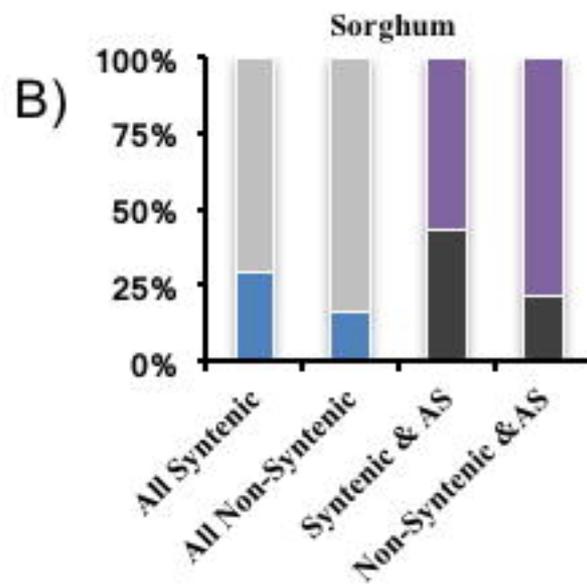
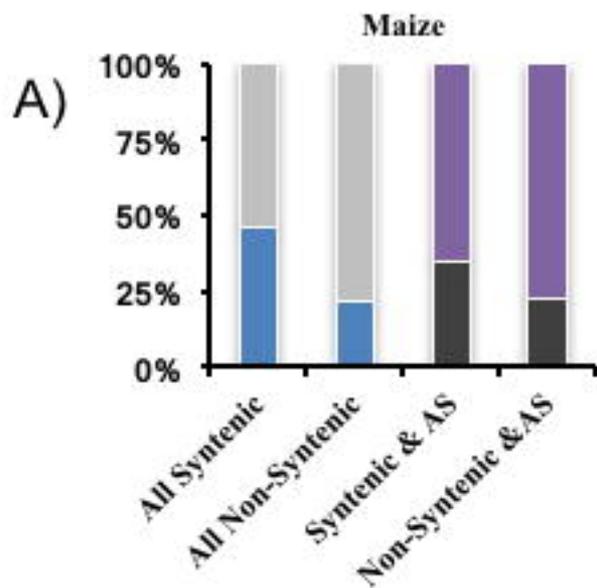
An example of sequence conservation of the alternative exons in a RNA-binding KH domain-containing protein. The exon skipping in *Arabidopsis* is not included in the Araport11, but detected in AtRTD. Vista plots of pair-wise alignments between the RNA-binding KH domain protein encoding gene from *Brachypodium* and its orthologue in *Amborella*, *Arabidopsis*, banana, African oil palm, rice, foxtail millet, sorghum and maize. The alternative exon (boxed in red) is highly conserved between *Brachypodium* and its orthologues in eight other plant species. Levels of sequence identity between *Brachypodium* and its orthologues in eight other plants species displayed are depicted as blue peaks. Pink bars signify segments that pass the alignment criteria of 70% identity over 100bp window.

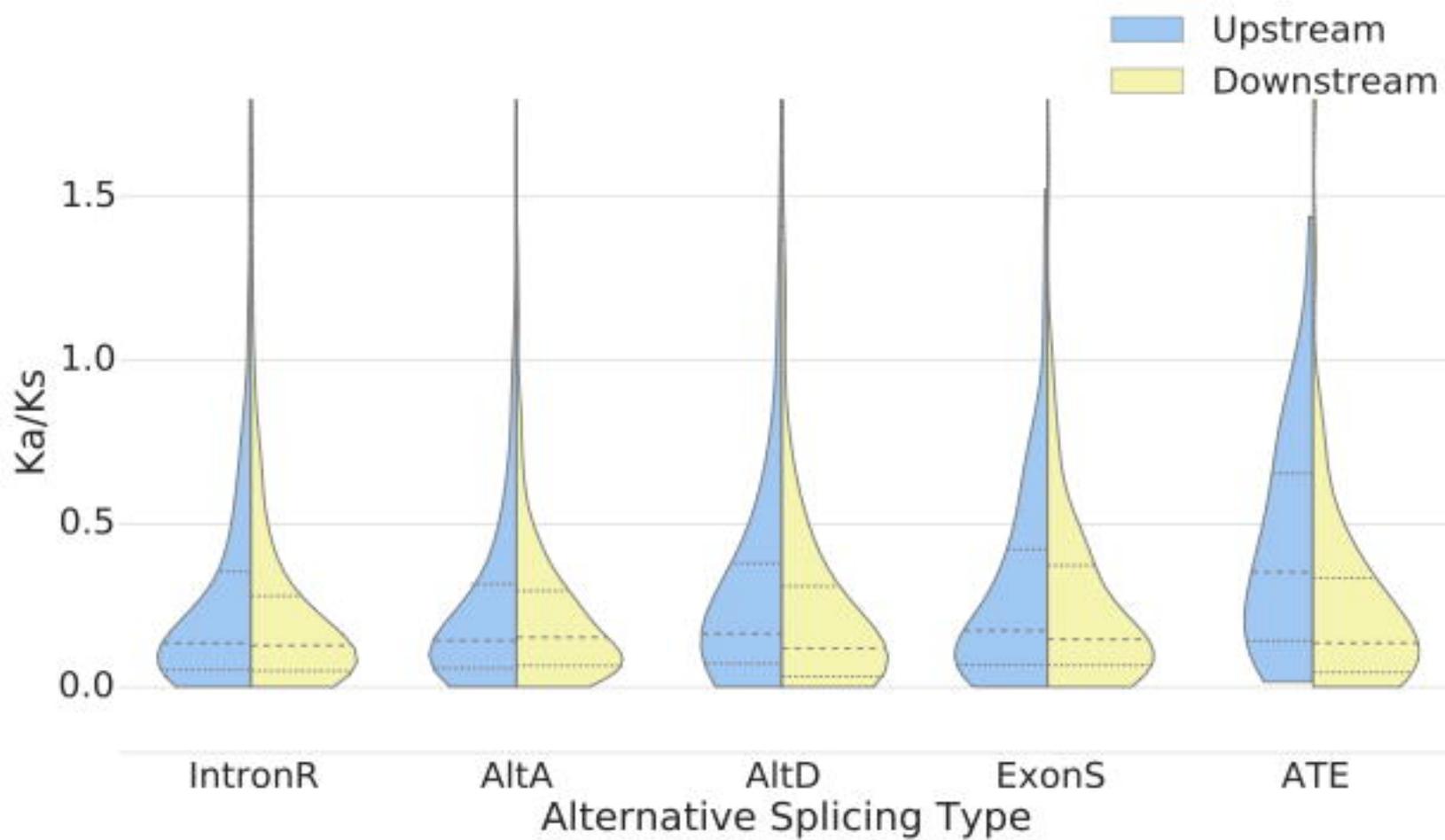




AS type	Total	Amborella gene counts	Arabidopsis gene counts	African oil palm gene counts	Banana gene counts	Rice gene counts	Brachypodium gene counts	Foxtail mille gene counts	Sorghum gene counts	Maize gene counts
AltA	6,152	444	222	197	409	867	931	1051	985	1046
AltD	2,925	207	105	72	187	427	414	473	491	549
AltTE	539	32	13	30	44	64	71	73	116	96
ExonS	1,249	116	43	68	113	180	173	163	201	192
IntronR	12,071	863	675	1019	974	1591	2238	1218	897	2596
Total (%)	19,235	1413 (18.3)	944 (12.0)	1269 (20.9)	1565 (17.4)	2579 (23.1)	3102 (33.2)	2440 (36.2)	2236 (32.7)	3687 (25.7)

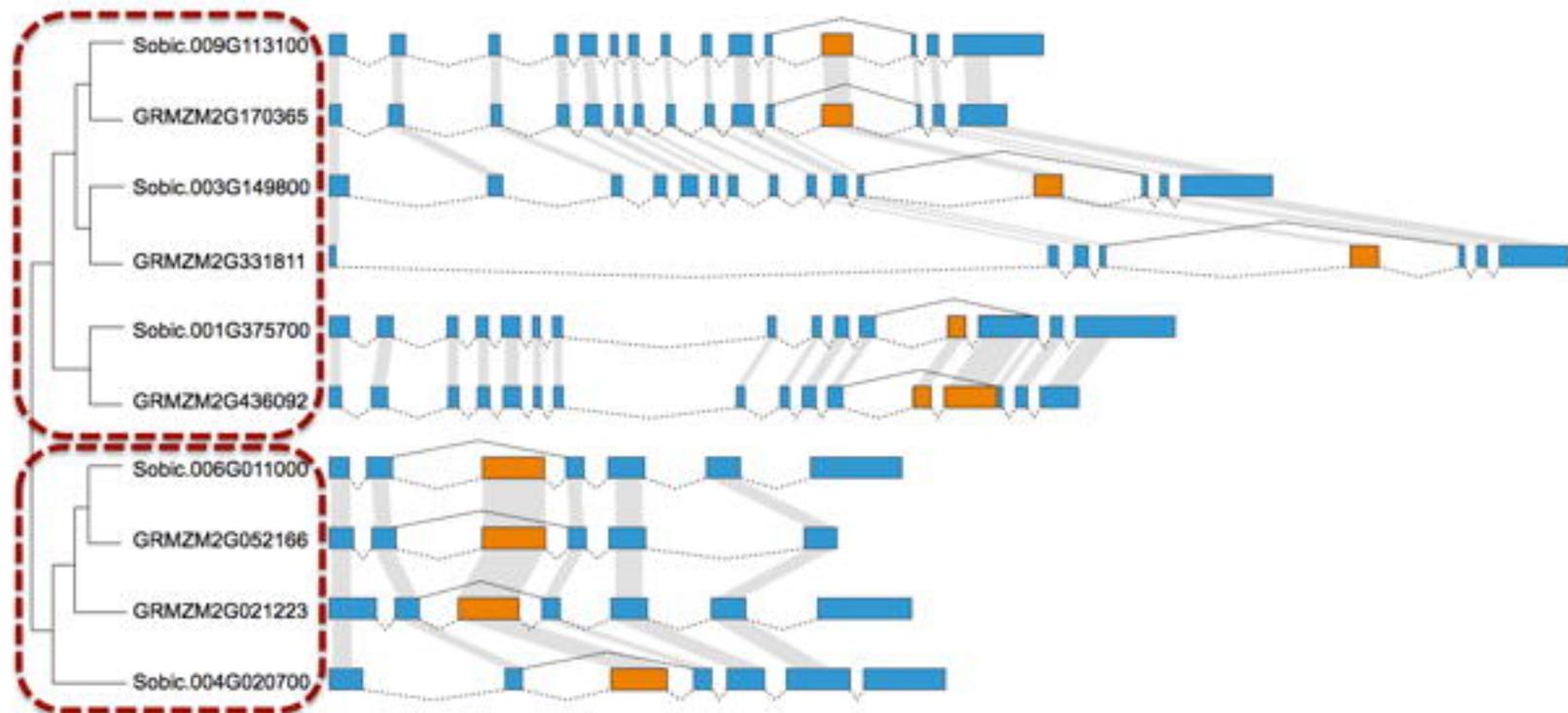
■ AS ■ Non-AS ■ Conserved AS ■ Non-Conserved AS





A)

SR



B)

RS2Z

