1 **Alternative splicing preferentially increases transcript diversity associated with stress responses in the**

2 **extremophyte *Schrenkiella parvula***

3

4 **Running title: Splicing increases isoform diversity under stress**

5

6

7 Chathura Wijesinghege, Kieu-Nga Tran, Maheshi Dassanayake*

8

9 Department of Biological Sciences, Louisiana State University, Baton Rouge, LA 70803, USA

10 *Address correspondence to: maheshid@lsu.edu

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

**Abstract**

Alternative splicing extends the coding potential of genomes by creating multiple isoforms from one gene. Isoforms can render transcript specificity and diversity to initiate multiple responses required during transcriptome adjustments in stressed environments. Although the prevalence of alternative splicing is widely recognized, how diverse isoforms facilitate stress adaptation in plants that thrive in extreme environments are unexplored. Here we examine how an extremophyte model, *Schrenkiella parvula*, coordinates alternative splicing in response to high salinity compared to a salt-stress sensitive model, *Arabidopsis thaliana*. We use Iso-Seq to generate full length reference transcripts and RNA-seq to quantify differential isoform usage in response to salinity changes. We find that single-copy orthologs where *S. parvula* has a higher number of isoforms than A. thaliana as well as S. parvula genes observed and predicted using machine learning to have multiple isoforms are enriched in stress associated functions. Genes that showed differential isoform usage were largely mutually exclusive from genes that were differentially expressed in response to salt. *S. parvula* transcriptomes maintained specificity in isoform usage assessed via a measure of expression disorderdness during transcriptome reprogramming under salt. Our study adds a novel resource and insight to study plant stress tolerance evolved in extreme environments.

Keywords: Extremophyte, Salt stress, Alternative splicing, Disorderdness of transcripts, Isoform usage

**Introduction**

Alternative splicing produces different mature RNAs from a single gene. Its impact on increasing transcript diversity has continued to broaden our understanding of gene regulatory mechanisms since it was first observed in 1977 [1,2]. The potential to create novel transcript diversity via alternative splicing is immense. The Drosophila *DSCAM* gene, which functions as an axon guidance receptor, is an extreme example of alternative splicing. It contains 115 exons and is estimated to give rise to more than 38,000 isoforms that are spatiotemporally regulated to achieve specific regulation in Drosophila neural development [3,4]. Alternative splicing events can be observed in more than 95% of human genes [5]. High throughput proteomics and ribosome bound mRNA sequencing (Ribo-Seq) studies show that a significant fraction of alternative splice variants are translated into protein isoforms [6,7]. Additionally, an ever-increasing array of transcriptome sequencing has revealed the existence of novel non-coding RNAs generated through alternative splicing suggesting their importance in gene regulatory circuits in all eukaryotic clades [8,9]. Differential splicing in closely related species have shown to reflect their divergent adaptive strategies not readily detectable at the primary gene expression level [10]. Tissue- and species-specific splicing is more divergent than promoter level divergence among closely related species facilitating independent evolutionary trajectories fitting to each species as highlighted by Calarco et al.[11] Therefore, genome wide discovery of new transcripts produced via alternative splicing becomes a critical initiative to understand the diversity of gene products and systematically assess their role in evolutionary innovations.

72    Alternative splicing increases proteome diversity as well as regulatory complexity in plants [12–14]. In

73    the model plant *Arabidopsis thaliana,* majority of the genes (>60%) undergo alternative splicing. There are

74    more than 70,000 non-redundant transcript isoforms reported for *A. thaliana* [15,16]. Similar reports on maize [17],

75    sorghum [18], and cotton [19] demonstrate that alternative splicing is prevalent in plants. Differential splicing has

76    also been targeted in crop breeding as shown with sunflowers [20].

77    Large scale changes in alternative splicing have been reported to allow transcriptional adjustments in

78    response to abiotic stresses including salt [21], cold [22], hypoxia [23], and heat stress [24]. Targeted functional studies

79    have also highlighted the significance of alternative splicing in responses to abiotic stresses. For example, the

80    heat shock protein gene, *hsf2* in *A. thaliana* produces an alternatively spliced transcript resulting in a truncated

81    protein that in turns binds to the *hsf2* promoter to enhance transcription of *hsf2* during heat stress [25]. While the

82    majority of published studies converge on alternative splicing being a key mechanism for environmental stress

83    adaptation in plants, such studies are limited to abiotic stress sensitive crop models or to *A. thaliana.*

84    Compared to crop plants, extremophytes that are naturally found in extreme environments are

85    equipped with evolutionary innovations that give them the ability to cope with multiple and extreme levels of

86    environmental stresses [26]. Therefore, extremophytes could show how the expanded transcriptome diversity via

87    alternative splicing may render additional paths for abiotic stress adaptations absent in stress sensitive models.

88    In this study, we have used the extremophyte model, *Schrenkiella parvula* (formerly *Thellungiella parvula* and

89    *Eutrema parvulum*) [27,28]. to examine its transcriptome diversity augmented by alternative splice variants. *S.*

90    *parvula* shares a highly co-linear genome with *A. thaliana* [29,30]. Yet, *S. parvula* is uniquely adapted to multiple

91    abiotic stresses reflecting its natural habitats often associated with hypersaline lakes in the Irano-Turanian

92    region [31,32].

93    Previous studies have shown that the *S. parvula* genome is enriched with duplicated genes associated

94    with abiotic stress responses and stress responsive genes show constitutive high expression as a stress

95    preadaptation compared to *A. thaliana* [29,32,33]. Alternative splicing plays a complementary role to gene

96    duplications and provides an additional path to increase transcript diversity [34]. Therefore, we aimed to test the

97    overall hypothesis that alternative splicing leads to increased diversity of stress responsive transcripts in *S.*

98    *parvula.*

99    In this study we investigated the complexity of the alternative splicing landscape in roots and shoots in

100    response to salt stress and how alternative splicing may provide transcript diversity associated with adaptations

101    to environmental stress in the model extremophyte *S. parvula*. We used PacBio Iso-Seq sequencing to identify

102    and annotate alternative splice variants and Illumina short reads to quantify isoform abundance. We find that

103    the *S. parvula* transcriptome is enriched in stress-associated isoforms. It shows specific isoform usage in a less

104    disordered state compared to the stress-sensitive model *A. thaliana* in response to high salinity.

105    **Materials and Methods**

106    **Plant material**

107       *Schrenkiella parvula* (ecotype Lake Tuz, Turkey; Arabidopsis Biological Resource 575 Center/ABRC

108   germplasm CS22663) seeds were grown hydroponically as previously described [33]. Briefly, plants were grown

109   at a light/dark cycle of 12/12 hr, 100 - 120 mM·m$^{-2}$s$^{-1}$ photon intensity, 20-22 °C, and 1/5$^{th}$ Hoagland's

110   solution for four weeks. These were treated with a with a combination of 250 mM NaCl, 250 mM KCl, 30 mM

111   LiCl, and 15 mM $H_3BO_3$ for three days to generate tissue samples used to create a reference transcriptome with

112   PacBio Iso-seq sequencing. Shoots and roots were harvested separately. RNA was extracted using QIAGEN

113   RNeasy Plant Mini Kit (QIAGEN, Hilden, Germany) with column digestion to remove DNA contamination.

114   About 4 µg of total RNA per tissue type at a quality of RNA integrity number ≥ 8 based on a Agilent 2100

115   Bioanalyzer (Agilent Technologies, CA, USA) were used to generate RNA-Seq libraries.

116       Shoot and root (1 µg) extracted as described above were used for cDNA synthesis using the

117   SuperScript cDNA Synthesis Kit (Invitrogen, Massachusetts, USA) following manufacturer's instructions to

118   test the presence of multiple isoforms independent from Iso-Seq for a randomly selected gene set expected to

119   express multiple isoforms. Isoform specific PCR primers (Supplementary Table 1) that span the alternative

120   splice sites were designed to use with an amplification protocol (initial denaturation at 95 °C for 3 min; 30

121   cycles of 95 °C for 30 s, 50-56 °C for 30 s, 72 °C for 2.30 to 3 min ; 72 °C for 10 min) run on a Bio-Rad T100

122   Thermal Cycler (Hercules, CA, USA) with a PCR Master mix Solution i-MAX II (iNtRON Biotechnology, S

123   Korea). PCR products were separated on a 1% agarose gel.

124       To quantify isoform abundance in response to high salinity compared to control conditions, RNA was

125   extracted from hydroponically grown *S. parvula* and *A. thaliana* (Col-0) as described in Tran et al. (2021).

126   These plants were treated with 150 mM NaCl for 24 hours and harvested together with samples hydroponically

127   grown without added NaCl as a control condition. The hydroponic growth conditions except for the specific

128   salt treatment was kept equivalent to growth conditions given to plants used for reference transcript generation

129   with Iso-Seq. At least 5 plants were used per biological replicate and three biological replicates were used for

130   each root and shoot sample for *S. parvula* and *A. thaliana* to yield a minimum of 1 µg of total RNA per sample

131   used for standard RNA-seq library preparation.

132

**Transcriptome sequencing**

134       For Iso-Seq based long read sequencing, cDNA synthesis, sequencing library preparation, and PacBio

135   sequencing were conducted at the Arizona Genomics Institute, University of Arizona, USA. Two Iso-seq

136   sequencing SMRT libraries were constructed following size selection from ≤ 4 kb and ≥ 3.5 kb per each tissue

137   and ran on two Pacific Biosciences Sequel cells with v2.1 Chemistry. For RNA-seq based short read

138   sequencing, mRNA enriched cDNA synthesis, library preparation, and sequencing were conducted at the Roy

139   K. Carver Biotechnology Center, University of Illinois Urbana-Champaign, USA. Briefly, True-Seq strand

140   specific libraries (Illumina, San Diego, CA, USA) were multiplexed and sequenced on an Illumina HiSeq4000

141   platform to generate >15 million 50-nucleotide single-end reads per sample.

142

143 **Identification and annotation of full-length transcript models for *S. parvula***

144       Raw Sequel data were processed using isoseq_sa5.1 pipeline

145 (https://github.com/PacificBiosciences/IsoSeq_SA3nUP. Circular consensus sequences (CCS) were generated

146 from subread BAM files with following parameters:  minLength□=□50, -noPolish --minLength=50, --

147 maxLength=15000, --minPasses=1, --minPredictedAccuracy=0.8, --minZScore=-999 --maxDropFraction=0.8.

148 CCS reads were selected as full length reads if it contained the 5′ and 3′ primers and a poly(A) signal

149 preceding the 3′ primer without additional copies of adapters. The full length consensus transcripts were

150 further clustered using ICE (Iterative Clustering for Error Correction) to obtain high-quality isoforms with

151 post-correction accuracy above 99% using Quiver. Error corrected full length reads were mapped to the

152 *Schrenkiella parvula* reference genome v2.2 (Phytozome genome ID: 574) to annotate isoforms assigned to

153 gene models and further select a set of high confidence transcript models. An isoform is annotated as a full

154 length transcript mapped to a genomic locus that has a single gene model assigned. If more than one isoform is

155 mapped to a gene model, the second and subsequent isoforms are considered products of alternative splicing.

156 First, TAPIS [17] was used to map isoforms and further error correct the isoforms. To map reads to the genome

157 GMAP [36] was used with parameters,  --no-chimeras, --cross-species --expand-offsets 1, -K 3000. Then,

158 SQANTI [37] was used with default parameters to identify the isoform that matched the primary gene model in

159 the genome and to assign additional isoforms that may be derived from that gene model as alternatively spliced

160 isoforms if both types of full-length isoforms were present in our processed full length data. Canonical splice

161 sites were defined as AG at the acceptor site and GT at the donor site. All the other splice sites were

162 categorized as non-canonical splice sites. Custom python script was used to count canonical and non-canonical

163 splice sites. Finally, we selected non-redundant structurally distinct isoform models that also contained a

164 complete and uninterrupted open reading frame as a selected set of putative protein coding transcript models.

165 Only isoforms that are likely to code for proteins were used for downstream analyses in the current study due

166 to the high uncertainty of functional significance and limited annotation resources available for newly

167 identified non-coding isoforms.

168       Functional annotations were assigned using PANTHER [38] and *A. thaliana* Gene Ontology (GO)

169 annotations (version release date 2020-07-16; DOI:10.5281/zenodo.3954044). Test for enriched functions were

170 performed using BiNGO [39]. Further clustering of enriched functions were performed using GOMCL [40] (with

171 parameters: -gosize 1500 -Ct 0.7 -I 1.5 -hm -nw -d -hg 0 -hgt –ssd) to get a non-redundant set of representative

172 functional annotations at p-values ≤0.05 adjusted for false discovery rate.

173

174 **Transcript and gene expression quantification**

175       Following quality checks using FastQC (http://www.bioinformatics.babraham.ac. uk/projects/fastqc/).

176 RNA-seq reads were mapped to gene models for *A. thaliana* (TAIR10) or *S. parvula* (Reference v2.2) as well

177 as transcript models obtained from AtRTD2 [41] or *S. parvula* Iso-seq supplemented transcript models using

178 Salmon [42] with parameters, "--type quasi -k 31" for indexing and "--gcBias -l A" for quantification. Ortholog

179 pairs between *S. parvula* and *A. thaliana* were assigned based on Oh & Dassanayake 2019 [43]. RNA-seq reads

180 mapped to gene models were used to identify differently expressed genes. A custom python script was used to

181 count uniquely mapped reads to each gene model. Differentially expressed genes between control and salt

182 treatments within each species were identified using DESeq2 [44] RNA-seq reads mapped to transcript models

183 were used for generating expression values for isoforms as well as quantify alternative splicing event

184 frequency. Expression counts for isoforms were converted to TPM (Transcript Per Million) and in

185 comparisons where an isoform was counted as expressed had $\geq$ 0.5 TPM normalized expression per isoform

186 independent from the expression quantified at the gene level. Isoform ratio per ortholog pairs was calculated

187 based on the number of isoforms per *S. parvula* ortholog divided the number of isoforms detected in the *A.*

188 *thaliana* ortholog.

189      Differential splicing was assessed using SUPPA2 [45]. Briefly, alternative splice events were identified

190 using generateEvents program and differential isoform expression was calculated based on the total expressed

191 number of isoforms per gene using psiPerIsoform included in SUPPA2 together with diffSplice to compare

192 differences in isoform expression between two conditions.

193

194 **Shannon entropy calculation for isoform specific transcriptome responses**

195 Isoform expression shifts between conditions or species were quantified using PSI values (proportion of

196 spliced isoforms) assigned for each alternatively spliced isoform per gene as given in the equation below. We

197 used the PSI values to calculate Shannon entropy per gene as described by Ritchie et al. (2008)[46] and used

198 normalized values between 0 and 1 for between species comparisons as described in Kumar et al. (1986) [47]

$$PSI_{isoform\ i\ of\ GeneA} = \frac{TPM_{isoform\ i\ of\ GeneA}}{\sum TPM_{all\ isoforms\ of\ GeneA}}$$

$$Normalized\ Shannon\ entropy_{GeneA} = -\frac{1}{\log N}\sum_{i=1}^{N=\#\ of\ isoforms\ of\ GeneA} PSI_{isoform\ i}\left(\log PSI_{isoform\ i}\right)$$

199      We calculated Shannon entropy values for genes expressed in control and salt treated samples for *S.*

200 *parvula* and *A. thaliana*. Genes with PSI values less than 0.01 or higher than 0.99 (expected when an isoform

201 is rarely expressed or dominates approximating zero alternative splicing for that gene) were removed from our

202 analysis to test for isoform expression shifts. Further, gene models which were not represented by at least two

203 isoforms were removed from the analysis.

204

205 **Splice site prediction for the *S. parvula* genome**

206      We used a deep-neural network, SpliceAi [48] to predict genome wide splice sites for *S. parvula* from

207 primary gene model sequences. The network model was trained first with *A. thaliana* gene models from

208 chromosome 1 to 4 and validated with chromosome 5 gene models described in Araport11 [49]. We provided

209 200 nucleotides upstream and downstream of a given base scanning all bases per gene in all genes models to

210 predict whether that site is a splice site donor, acceptor or not a splice site. Model prediction was assigned a

211 probability score between 0 and 1 for a given site with values closer to 1 representing the probability of that

212 site being a splice site. We used a probability score of ≥0.6 for the selection of potential splice sites. We used

213 this trained network to predict splice sites for the *S. parvula* genome v2. We compared the predicted splice

214 sites to observed splice sites and identified new splice sites. If new splice sites were predicted for a gene model

215 we had identified more than one isoform, the prediction of a novel splice site or sites for that gene was

216 considered as one additional predicted putative isoform.

217

218 **Results**

219 **Improvement of isoform annotation in *S. parvula***

220        Prior to this study, the *S. parvula* reference gene models (v2.2) were predicted based on *ab inito*

221 methods as well as RNA-seq evidence based prediction derived from non-stressed conditions [29]. To maximize

222 the identification of transcripts that may be conditionally expressed under stress, we used 4-week-old *S.*

223 *parvula* plants treated with multiple salts (NaCl, KCl, LiCl, and $H_3BO_3$) that are found at high levels in its

224 native soils [50] for PacBio Iso-Seq sequencing. We obtained 500,265 error corrected circular consensus

225 sequences (CCS) as our primary source of sequence reads to create an isoform specific reference transcriptome

226 and to supplement the genome-based transcript annotation for *S. parvula*. We identified putative full-length

227 transcripts based on 338,812 high quality CCS reads that contained 5' and 3' primers and polyA tails (Figure

228 1A). Following iterative clustering, error correction, and mapping to the *S. parvula* reference genome, we

229 annotated 16,828 (corresponding to 11,348 genomic loci) structurally distinct putative protein coding

230 transcripts expressed in *S. parvula* tissues exposed to multiple salts (see Methods for details). This added 7,732

231 new protein coding transcript models to the *S. parvula* reference genome to provide a total of 34,582 reference

232 protein coding transcripts (Table 1).

233        We were able to improve the *S. parvula* reference genome to include full length transcripts inclusive

234 of 5' and 3' UTR regions with Iso-Seq reads. The average length of new Iso-Seq supported reference transcript

235 models was greater than the corresponding length of transcript models in the *S. parvula* v2.2 reference genome

236 annotation (Fig. 1B). The increase in transcript lengths was largely due to the identification of 5' and 3' UTR

237 sequences that were previously missed in transcript model predictions in the reference genome. This

238 refinement of reference transcript models generated UTR length distribution comparable to that of *A. thaliana*

239 reference genome (Araport11) (Figure S1) and significantly increased the percentage of standard RNA-seq

240 reads mapped to the reference transcriptome (Fig. 1C). This is expected to improve estimates of gene

241 expression counts when using short-read RNA-seq data.

242        New genes previously not reported for *S. parvula* was added with Iso-Seq supported transcripts. The

243 current reference *S. parvula* v2.2 genome includes 26,847 total protein coding primary gene models. The Iso-

244 Seq supported transcripts mapped to 11,348 (42%) of those genes (Table 1). We additionally identified 301

7

245 novel gene models that were missed (i.e. sequence present in the genome but annotation absent) in the *S.*

246 *parvula* reference genome. For example, the putative ortholog of the Arabidopsis *Magnesium/proton*

247 *exchanger* (*MHX,* similar to *At2G47600)* in the *S. parvula* genome was annotated on chromosome 4 between

248 *Sp4g29520* and *Sp4g29540*, using an Iso-seq based transcript model detected in this study (Figure S2). The

249 novel transcript models further improved the reference genome annotation by adding multiple isoforms

250 assigned to gene models, alternative transcription start and end sites for existing models, and UTR sequences

251 (Table 1). The improved gene models, isoform specific expression, and Iso-Seq reads are available at

252 Bioproject ID PRJNA63667.

253 We identified 5,911 alternatively spliced events, resulting in structurally different protein coding

254 regions from the primary transcript models in the *S. parvula* genome from this study. These splice variants

255 were categorized into intron retention, alternative 3' acceptor, alternative 5' donor, exon skipping, use of

256 alternative first exon, use of alternative last exon, and use of mutually exclusive exons based on their

257 frequency (Table 2). Intron retention was the most prevalent (55.2%) alternatively spliced event in *S. parvula*.

258 We observed that two or more distinctly spliced isoforms could be co-expressed in either shoots or roots

259 (Figure S3) when multiple isoforms were checked for their expression using RT-PCR for a select set of genes.

260

261 **Salt stress associated genes show a higher isoform diversity in *S. parvula* compared to *A. thaliana***

262 Alternative splicing can increase the repertoire of transcripts that are available to respond to abiotic

263 stresses more efficiently and dynamically, independent of gene copy number variation [51]. Therefore, we

264 hypothesized that *S. parvula* would have a higher diversity of alternatively spliced isoforms for genes related

265 to abiotic stress tolerance, specifically salinity tolerance, than in the less-tolerant species *A. thaliana*. To test

266 this, we calculated the isoform ratio per ortholog pair in *S. parvula* and *A. thaliana* using the *S. parvula*

267 reference isoforms identified in this study and *A. thaliana* reference isoforms obtained from AtRTD2 database

268 [16]. To avoid missing data or lack of expression of a certain gene in mature shoots or roots in one species being

269 inferred as lack of isoform diversity in that species*,* we limited our comparison to genes expressed in our study

270 that were represented by at least one transcript model in both species. We identified 10,859 *A. thaliana - S.*

271 *parvula* ortholog pairs that had one or more isoforms per ortholog in each species (Figure 2A; Supplementary

272 Table 2). Among them there were 6,874 ortholog pairs showing more isoforms in *A. thaliana* while only 1,201

273 pairs had a higher isoform number in *S. parvula* (Fig. 2A). Ortholog pairs annotated as "Response to stress"

274 (GO:0006950) and "Transport" (GO:0006810) had a higher isoform diversity in *S. parvula*, while ortholog

275 pairs annotated under "Nitrogen metabolism" (GO:0034641) had a higher isoform diversity in *A. thaliana* (Fig.

276 2B). As a control, we examined the distribution of isoforms in all ortholog pairs and found that these

277 distributions were not significantly different between the two species (Fig. 2B).

278 The genes that had a higher isoform diversity in *S. parvula* included some of the most highly

279 conserved and key stress responsive genes in plants including the $Na^+/H^+$ antiporter, *SOS1* known for its role

280 in excluding $Na^+$ from roots during salt stress [52] and *P5CS1* that codes for delta1-pyrroline-5-carboxylate

8

281 synthase, the rate-limiting enzyme in proline biosynthesis known for its role in oxidative and osmotic stress

282 responses [53]. Notably, both *SOS1* and *P5CS1* are represented by single copy orthologs in *S. parvula* and *A.*

283 *thaliana*. Five *SOS1* (out of 8 detected) and 8 *P5CS1* (out of 22 detected) isoforms for *S. parvula* were

284 expressed at $\geq$ 0.5 TPM in both shoots and roots in control as well as salt treated conditions (Figure S4). The

285 AtRTD2 database reported three *SOS1* and six *P5CS1* isoforms for *A. thaliana* [16].

286

287 **Isoform usage is less disordered in *S. parvula* compared to *A. thaliana* during salt stress**

288 Diversity and conditional expression (i.e. specificity) of isoforms can be assessed using the Shannon

289 entropy based information theory applied to transcriptomes [54]. Stressed compared to growth optimal conditions

290 are known to have higher transcriptome entropy and disorderdness with an increased number of alternative

291 splice events when assessed using Shannon Entropy [46]. We hypothesized that *S. parvula* transcriptomes will

292 show a smaller entropy increase in its isoform usage when transitioning from control to salt stressed treatments

293 compared to the salt-sensitive model *A. thaliana.* To test if isoform usage from control to stressed conditions

294 went through a measurable entropy transition distinctive of the species, we used RNA-seq data from root and

295 shoot samples to quantify the isoform abundance in *S. parvula* and *A. thaliana* and calculated the Shannon

296 entropy (see Methods). We used *A. thaliana-S. parvula* ortholog pairs that were represented by at least two

297 expressed isoforms with a normalized expression $\geq$ 0.5 TPM per ortholog within a species to avoid incomplete

298 comparisons due to rare isoforms difficult to quantify in one species. This resulted in a total of 1,678 and 1,592

299 ortholog pairs expressed in roots and shoots. Roots had 3,832 and 5,239 isoforms for *S. parvula* and *A.*

300 *thaliana* while shoots had 3658 and 4431 isoforms respectively. We found that both *S. parvula* and *A. thaliana*

301 root transcript distributions increased mean entropy in response to salt stress (Figure 3A). This is aligned with

302 the expectation that stress conditions create higher transcript diversity, lower specificity, and more

303 disorderdness in transcript expression compared to a stress-neutral control condition [55]. *A. thaliana* shoots

304 showed a significant increase in entropy when transitioning from control to salt stressed conditions (Fig. 3A).

305 The change in entropy for *S. parvula* was less in both roots and shoots suggesting a less disordered state of

306 isoform usage compared to the relatively stress-sensitive *A. thaliana* when responding to stress conditions. We

307 observed that the isoform usage in response to salt was highly species specific. The number of ortholog pairs

308 that showed increased or decreased isoform usage as a shared response to salt stress in both species roots (156

309 expressed orthologs) and shoots (142 expressed orthologs) were much fewer than those orthologs (977 in roots

310 and 937 shoots) that had a specific usage change in one species (Fig. 3B). Orthologs that showed high isoform

311 usage specificity (i.e. maintained or lowered entropy) in response to salt stress in *S. parvula* roots compared to

312 *A. thaliana* were enriched in functions largely associated with salt stress (Fig. 3C). In shoots, genes that

313 maintained isoform usage specificity under salt stress in both species were enriched in salt stress associated

314 functions (Fig. 3C).

315

316 **Distinct regulation between different isoform usage and differential expression in response to salt stress**

9

317    We next examined if the differently expressed genes in response to salt stress were also subjected to

318    changes in their isoform usage under high salinity. Supplementary Table 3 lists all genes identified as

319    differently spliced or differently expressed. Genes that were differently expressed as well as differently spliced

320    in response to salt stress were rare in *S. parvula* and *A. thaliana* (≤ 3%) (Figure 4A). Moreover, the shared

321    orthologs that were differently spliced in response to salt stress between species either in roots or shoots were

322    also low (~3%) (Fig. 4B). Multiple genes differently expressed under salt stress in *A. thaliana* are found to be

323    only differently spliced in response to salt in *S. parvula* (Supplementary Table 3). Figure S5 further highlights

324    the high degree of species-specific regulation in differential isoform usage in response to stress. However,

325    there is high convergence in the enriched functions represented by differently spliced isoforms in response to

326    salt stress in both species (Fig. 4C).

327

328    **Non-canonical splice sites are enriched in stress associated genes**

329    Majority of splice sites in plants are marked by GU at the 5' and AG at the 3' sites in introns [56].

330    Although less common, plant genes are spliced at alternative sites termed as non-canonical splice sites and

331    alternative splicing at non-canonical sites are associated with abiotic stress responses [56–58]. We investigated

332    whether the expression of transcripts with non-canonical splice sites (Supplementary Table 4) increased under

333    salt stress in *S. parvula* compared to *A. thaliana*. We found that *S. parvula* did not show any significant

334    difference in mean expression strength between non-canonical and canonical transcripts in both roots and

335    shoots while *A. thaliana* shoots showed an increased expression in transcripts that had non-canonical splice

336    sites when treated with salt (Figure 5A).

337    Previous studies have reported increases in non-canonical splicing in plants under abiotic stresses [59,60].

338    Therefore, we examined if usage of transcripts with non-canonical splice sites significantly increased under

339    salt stress compared to control conditions in *S. parvula* differently from *A. thaliana*. Similar to previous

340    reports, non-canonically spliced transcripts are less frequent than canonically spliced transcripts regardless of

341    the condition tested (≤ 10%; Fig. 5B). However, our analysis does not find a significant increase in non-

342    canonically spliced isoforms from control to salt treated conditions in either species (Fig. 5B).

343    Next, we tested if genes with non-canonical splice sites were enriched for stress associated functions

344    in *S. parvula*. We found 424 genes out of 25,145 multi-exon coding genes to be enriched in non-canonical

345    splice sites in the *S. parvula* genome (Fig. 5C). Some of these are notable genes associated with stress

346    regulatory pathways (for example, *PAL1, PAL2, P5CS1* and *HSC70-1*) (Fig. 5C). *S. parvula* genes enriched for

347    non-canonical splice sites were indeed primarily enriched in stress response pathways (Fig. 5D). Further sub-

348    clustering of the functional group annotated under "stress responses" (cluster C1 of Fig. 5D) showed that genes

349    in salt/metal ion and osmotic stress were specifically contributing to this cluster.

350

351    **Predicted isoforms for the *S. parvula* genome is enriched for stress responsive genes**

352    It is likely that we may have missed to detect stress responsive isoforms expressed in *S. parvula* in this

353 study because exhaustive searches for conditionally expressed isoforms are impractical for emerging model

354 organisms. Therefore, we sought to employ a machine learning approach to predict alternative splicing sites in

355 the *S. parvula* genome as an alternative. We applied the deep neural network, SpliceAI which is expected to

356 yield high confidence predictions among recent tools developed to predict splice events using genomic

357 sequences [48,61,62]. We used the known splice site information from *A. thaliana* chromosomes 1-4 to train the

358 SpliceAI network and received an average precision of 0.92 when tested with *A. thaliana* chromosome 5

359 (Figure 6A). We then predicted splice sites from 26,847 *S. parvula* pre-mRNA sequences and obtained

360 214,901 splice site predictions including 114,284 novel splice sites (Fig. 6B). Twenty-six percent of splice

361 sites previously observed were also predicted using SpliceAI and we found 7,302 genes with at least one

362 newly, predicted isoform. Prediction probability scores were highest for splice sites within the gene compared

363 to those in the first and the last introns (Figure S6).

364 With the current analysis, we have identified 16,061 potential protein coding isoforms (observed or

365 predicted) for 9,033 genes in the *S. parvula* genome (Supplementary Table 5). Interestingly, stress and

366 transport associated functions are enriched among those genes that are observed or predicted to have more than

367 one isoform (Fig. 6C). Stress and transport related functions deduced from GO annotations account for 35% of

368 genes that are alternatively spliced in the *S. parvula* genome.

369

370 **Discussion**

371 **Alternative isoforms of stress related genes from an extremophyte model as a resource in environmental**

372 **stress adaptations**

373 Alternative splicing allows genes to acquire new functions independent from gene duplications and

374 promoter evolution. Previous studies have shown that duplicated genes are enriched in stress associated

375 functions in *S. parvula* and other extremophytes facilitating their stress adapted lifestyles more than in stress-

376 sensitive sister species [29,63,64]. However, extremophyte gene diversity represented by alternatively spliced

377 isoforms is underexplored [65]. Certain genes are regulated only at the alternative splicing level with no change

378 at the gene expression level that have led to the increasing recognition of the importance of isoform specific

379 reference transcript datasets in gene expression studies[66]. In this study we examined the possibility of

380 diversifying gene functions through alternative splicing and specially focused on isoforms differently used

381 during salt stress in one of the leading model extremophytes [26].

382 *Schrenkiella parvula* and *A. thaliana* genomes have similar gene numbers (~27,000) and similar

383 genome sizes (~120 MB) [29]. A recent study that explored alternative splicing in *A. thaliana* using full length

384 transcript sequencing based on Iso-Seq reports the discovery of isoforms in similar proportions to our study

385 with intron retention being the most common alternative splicing event [67] (Table 2). This suggests that *S.*

386 *parvula* is not an exception in highly increased or decreased transcript diversity through alternative splicing

387 although the recorded number of isoforms for the model plant through aggregate studies using multiple tissues,

388 developmental stages, and treatments are much higher (Zhang et al., 2017). Given the genomic similarities

389  between *S. parvula* and *A. thaliana*, their transcriptome adjustments with differential splicing in response to

390  salt stress were remarkably distinct from one another when an identical salt treatment was given to mature

391  plants (same age and tissues tested in both species) (Figs. 4 and S5).

392       In support of our hypothesis that extremophytes would diversify their response to stress via alternative

393  splicing in selected gene groups, we observed that *S. parvula* orthologs had a higher number of isoforms

394  compared to *A. thaliana* in genes associated with stress and transport functions (Fig. 2). Stress and transport

395  functions were also enriched among duplicated genes in *S. parvula* compared to *A. thaliana* [68]. We found that

396  differently expressed genes and genes that showed differential isoform usage were largely mutually exclusive

397  within species as well as in one-to-one ortholog pairs between *S. parvula* and *A. thaliana* (Figs. 4 and S5).

398  Further, when we combine both observed and predicted splice sites in the *S. parvula* genome, the potential

399  protein coding isoform pool is enriched in functions associated with stress tolerance (Fig. 6). These

400  observations together indicate that genes expressed in response to stress are highly diversified and non-

401  overlapping in their mode of function, but converge on common functions associated with stress tolerance in *S.*

402  *parvula*. Therefore, our study provides a novel resource for assessing functional significance of stress tolerance

403  genes in the extremophyte model. It allows selection of target genes that could be tested at the isoform level

404  when expression modulation via promoter modifications or single gene-knockouts of essential genes do not

405  offer optimal methods to test novel gene functions contributing to stress tolerance.

406

407  **Isoform usage and the specificity of their expression in response to salt**

408       Our current study in agreement with a previous study on *A. thaliana* have shown that most differently

409  spliced genes were not differently expressed in response to salt stress representing an independent layer of

410  gene regulation in response to stress [21]. Compared to animals, plants tend to use alternative splicing biased to

411  environmental stress responses more than for tissue-specific responses [69]. Multiple studies have reported

412  specific associations of alternative splicing and environmental stress in plants [22,57,59,70,71]. However, fewer

413  studies have examined the presence of non-specific alternative splicing leading to increased number of

414  differently spliced isoforms under abiotic stress [12,72]. Additionally, components of the spliceosome are

415  differently expressed leading to differential splicing of target genes in *A. thaliana* during stress conditions [60].

416  In animals, stressed conditions are reported to have increased amount of alternatively spliced isoforms with

417  high non-specific expression, thus creating a higher level of disorderdness in isoform expression [46] which can

418  be quantified using Shannon entropy [73,74]. We predicted that plants will show a similar trend in increased

419  disorderdness in isoform expression at the transcriptome level during stress conditions. Furthermore, we

420  expected to see a smaller change in entropy in the extremophyte when transitioning to a salt treated condition

421  compared to the stress-sensitive species. Indeed, this prediction was supported by the shoot transcriptomic

422  response we observed for *S. parvula* and *A. thaliana* (Fig. 3). Notably, the genes that shifted to lower entropy

423  values representing shifts to specific isoform in their expression specificity under stress were enriched for

424  stress associated functions in both roots and shoots in *S. parvula* (Fig. 3). Our study cannot test if the tendency

425  to increase transcriptome disorderdness via less specifically expressed isoforms per gene is indicative of

426  aberrant splicing under stress. Yet, the comparison between *S. parvula* and *A. thaliana* suggests that the

427  extremophyte is more prepared to respond to salt stress by specific isoforms mostly expressed for stress

428  associated genes.

429      In conclusion, this study provides a novel resource for a leading extremophyte model and expands our

430  knowledge on the ability to respond to stress via differential isoform usage independently from differential

431  gene expression. Stress associated functions were enriched among genes observed or predicted to have

432  multiple isoforms in *S. parvula*; one-to-one orthologs where *S. parvula* has a higher number of isoforms than

433  *A. thaliana*; genes that showed differential isoform usage in response to stress in *S. parvula*; *S. parvula* genes

434  that were enriched in non-canonical splice sites; and *S. parvula* genes that maintained or lowered their

435  disorderdness by expression of specific isoforms under stress. These findings contribute to how we understand

436  stress tolerance evolved in an extremophyte. Differential isoform usage offers a complementary path to

437  increase the coding potential of the *S. parvula* genome that cannot be fully explained by gene duplication or

438  promoter evolution alone. Future studies on other extremophytes exploring isoform diversity will facilitate the

439  identification of convergent traits in isoform usage evolved in stress-adapted plants. Such a resource will be

440  influential in deducing diverse stress responsive networks and identifying transferable stress responsive genes

441  into crops.

442

443  **Acknowledgements**

450

451  **Author Contributions**

452  CW and KT conducted wet lab experiments. CW performed bioinformatics analyses. MD developed the

453  experimental design and supervised the overall project. CW, KT, and MD interpreted results and wrote the

454  article.

455

456  **References**

457  1.    Berget, S. M., Moore, C., and Sharp, P. A. 1977, Spliced segments at the 5' terminus of adenovirus 2 late

458        mRNA. *Proc. Natl. Acad. Sci. U. S. A.*, **74**, 3171–5.

459  2.    Ule, J., and Blencowe, B. J. 2019, Alternative Splicing Regulatory Networks: Functions, Mechanisms, and

460        Evolution. *Mol. Cell*, **76**, 329–45.

461   3.   Schmucker, D., Clemens, J. C., Shu, H., et al. 2000, Drosophila Dscam Is an Axon Guidance Receptor
462        Exhibiting Extraordinary Molecular Diversity modified by neuronal activity (reviewed by Albright et al.,
463        2000). Extraordinary progress has been made in identifying. *Cell*, **101**, 671–84.

464   4.   Celotto, A. M., and Graveley, B. R. 2001, Alternative splicing of the Drosophila Dscam pre-mRNA is both
465        temporally and spatially regulated. *Genetics*, **159**, 599–608.

466   5.   Pan, Q., Shai, O., Lee, L. J., Frey, B. J., and Blencowe, B. J. 2008, Deep surveying of alternative splicing
467        complexity in the human transcriptome by high-throughput sequencing. *Nat. Genet.*, **40**, 1413–5.

468   6.   Furlanis, E., Traunmüller, L., Fucile, G., and Scheiffele, P. 2019, Landscape of ribosome-engaged transcript
469        isoforms reveals extensive neuronal-cell-class-specific alternative splicing programs. *Nat. Neurosci.*, **22**,
470        1709–17.

471   7.   Reixachs-Solé, M., Ruiz-Orera, J., Albà, M. M., and Eyras, E. 2020, Ribosome profiling at isoform level
472        reveals evolutionary conserved impacts of differential splicing on the proteome. *Nat. Commun.*, **11**, 1768.

473   8.   Morgan, J. T., Fink, G. R., and Bartel, D. P. 2019, Excised linear introns regulate growth in yeast. *Nature*,
474        **565**, 606–11.

475   9.   Gil, N., and Ulitsky, I. 2020, Regulation of gene expression by cis-acting long non-coding RNAs. *Nat. Rev.*
476        *Genet.*, pp. 102–17.

477   10.  Barbosa, C., Peixeiro, I., and Romão, L. 2013, Gene Expression Regulation by Upstream Open Reading
478        Frames and Human Disease. *PLoS Genet.*, **9**, 1–12.

479   11.  Calarco, J. A., Xing, Y., Cáceres, M., et al. 2007, Global analysis of alternative splicing differences between
480        humans and chimpanzees. *Genes Dev.*, **21**, 2963–75.

481   12.  Kalyna, M., Simpson, C. G., Syed, N. H., et al. 2012, Alternative splicing and nonsense-mediated decay
482        modulate expression of important regulatory genes in Arabidopsis. *Nucleic Acids Res.*, **40**, 2454–69.

483   13.  Yang, X., Zhang, H., and Li, L. 2012, Alternative mRNA processing increases the complexity of
484        microRNA-based gene regulation in Arabidopsis. *Plant J.*, **70**, 421–31.

485   14.  Chaudhary, S., Jabre, I., Reddy, A. S. N., Staiger, D., and Syed, N. H. 2019, Perspective on Alternative
486        Splicing and Proteome Complexity in Plants. *Trends Plant Sci.*, **24**, 496–506.

487   15.  Filichkin, S. A., Priest, H. D., Givan, S. A., et al. 2010, Genome-wide mapping of alternative splicing in
488        Arabidopsis thaliana. *Genome Res.*, **20**, 45–58.

489   16.  Zhang, R., Calixto, C. P. G., Marquez, Y., et al. 2017, A high quality Arabidopsis transcriptome for accurate
490        transcript-level analysis of alternative splicing. *Nucleic Acids Res.*, **45**, 5061–73.

491   17.  Wang, B., Tseng, E., Regulski, M., et al. 2016, Unveiling the complexity of the maize transcriptome by
492        single-molecule long-read sequencing. *Nat. Commun.*, **7**.

493   18.  Wang, B., Regulski, M., Tseng, E., et al. 2018, A comparative transcriptional landscape of maize and
494        sorghum obtained by single-molecule sequencing. *Genome Res.*, **28**, 921–32.

495   19.  Wang, M., Wang, P., Liang, F., et al. 2018, A global survey of alternative splicing in allopolyploid cotton:
496        landscape, complexity and regulation. *New Phytol.*, **217**, 163–78.

497   20.  Smith, C. C. R., Tittes, S., Paul Mendieta, J., et al. 2018, Genetics of alternative splicing evolution during

498          sunflower domestication. *Proc. Natl. Acad. Sci. U. S. A.*, **115**, 6768–73.

499    21.    Ding, F., Cui, P., Wang, Z., Zhang, S., Ali, S., and Xiong, L. 2014, Genome-wide analysis of alternative

500          splicing of pre-mRNA under salt stress in Arabidopsis. *BMC Genomics*, **15**, 1–14.

501    22.    Calixto, C. P. G., Guo, W., James, A. B., et al. 2018, Rapid and dynamic alternative splicing impacts the

502          arabidopsis cold response transcriptome[CC-BY]. *Plant Cell*, **30**, 1424–44.

503    23.    Chen, M. X., Zhu, F. Y., Wang, F. Z., et al. 2019, Alternative splicing and translation play important roles in

504          hypoxic germination in rice. *J. Exp. Bot.*, **70**, 885–95.

505    24.    Kannan, S., Halter, G., Renner, T., and Waters, E. R. 2018, Patterns of alternative splicing vary between

506          species during heat stress. *AoB Plants*, **10**, 1–11.

507    25.    Liu, J., Sun, N., Liu, M., et al. 2013, An autoregulatory loop controlling Arabidopsis HsfA2 expression:

508          Role of heat shock-induced alternative splicing. *Plant Physiol.*, **162**, 512–21.

509    26.    Kazachkova, Y., Eshel, G., Pantha, P., Cheeseman, J. M., Dassanayake, M., and Barak, S. 2018,

510          Halophytism: What have we learnt from arabidopsis thaliana relative model systems? *Plant Physiol.*, **178**,

511          972–88.

512    27.    Oh, D.-H., Dassanayake, M., Bohnert, H. J., and Cheeseman, J. M. 2012, Life at the extreme: lessons from

513          the genome. *Genome Biol. 2012 133*, **13**, 127–30.

514    28.    Zhu, J.-K., Jessica Whited, Andrei Seluanov, et al. 2015, The Next Top Models. *Cell*, **163**, 18–20.

515    29.    Dassanayake, M., Oh, D.-H., Haas, J. S., et al. 2011, The genome of the extremophile crucifer Thellungiella

516          parvula. *Nat. Genet.*, **43**, 913–8.

517    30.    Oh, D. H., and Dassanayake, M. 2019, Landscape of gene transposition-duplication within the Brassicaceae

518          family. *DNA Res.*, **26**, 21–36.

519    31.    Gul Nilhan, T., Ahmet Emre, Y., and Osman, K. 2008, Soil Determinants for Distribution of Halocnemum

520          strobilaceum Bieb. (Chenopodiaceae) Around Lake Tuz, Turkey. *Pakistan J. Biol. Sci.*, **11**, 565–70.

521    32.    Oh, D. H., Hong, H., Lee, S. Y., Yun, D. J., Bohnert, H. J., and Dassanayake, M. 2014, Genome structures

522          and transcriptomes signify niche adaptation for the multiple-ion-tolerant extremophyte Schrenkiella parvula.

523          *Plant Physiol.*, **164**, 2123–38.

524    33.    Wang, G., DiTusa, S. F., Oh, D., et al. 2021, Cross species multi□omics reveals cell wall sequestration and

525          elevated global transcript abundance as mechanisms of boron tolerance in plants. *New Phytol.*, **230**, 1985–

526          2000.

527    34.    Iñiguez, L. P., and Hernández, G. 2017, The evolutionary relationship between alternative splicing and gene

528          duplication. *Front. Genet.*, **8**, 1–7.

529    35.    Tran, K.-N., Wang, G., Oh, D.-H., Larkin, J. C., Smith, A. P., and Dassanayake, M. 2021, Multiple paths

530          lead to salt tolerance - pre-adaptation vs dynamic responses from two closely related extremophytes.

531          *bioRxiv*, 2021.10.23.465591.

532    36.    Wu, T. D., and Watanabe, C. K. 2005, GMAP: A genomic mapping and alignment program for mRNA and

533          EST sequences. *Bioinformatics*, **21**, 1859–75.

534    37.    Tardaguila, M., De La Fuente, L., Marti, C., et al. 2018, SQANTI: Extensive characterization of long-read

15

535    transcript sequences for quality control in full-length transcriptome identification and quantification.

536    *Genome Res.*, **28**, 396–411.

537    38.    Mi, H., Muruganujan, A., and Thomas, P. D. 2013, PANTHER in 2013: Modeling the evolution of gene

538        function, and other gene attributes, in the context of phylogenetic trees. *Nucleic Acids Res.*, **41**, 377–86.

539    39.    Maere, S., Heymans, K., and Kuiper, M. 2005, BiNGO: A Cytoscape plugin to assess overrepresentation of

540        Gene Ontology categories in Biological Networks. *Bioinformatics*, **21**, 3448–9.

541    40.    Wang, G., Oh, D. H., and Dassanayake, M. 2020, GOMCL: A toolkit to cluster, evaluate, and extract non-

542        redundant associations of Gene Ontology-based functions. *BMC Bioinformatics*, **21**, 1–9.

543    41.    Zhang, R., Calixto, C. P. G., Marquez, Y., et al. 2017, A high quality Arabidopsis transcriptome for accurate

544        transcript-level analysis of alternative splicing. *Nucleic Acids Res.*, **45**, 5061–73.

545    42.    Patro, R., Duggal, G., Love, M. I., Irizarry, R. A., Kingsford, C., and Biology, C. 2017, Salmon, **14**, 417–9.

546    43.    Oh, D. H., and Dassanayake, M. 2019, Landscape of gene transposition-duplication within the Brassicaceae

547        family. *DNA Res.*, **26**, 21–36.

548    44.    Love, M. I., Huber, W., and Anders, S. 2014, Moderated estimation of fold change and dispersion for RNA-

549        seq data with DESeq2. *Genome Biol.*, **15**, 1–21.

550    45.    Trincado, J. L., Entizne, J. C., Hysenaj, G., et al. 2018, SUPPA2: Fast, accurate, and uncertainty-aware

551        differential splicing analysis across multiple conditions. *Genome Biol.*, **19**, 1–11.

552    46.    Ritchie, W., Granjeaud, S., Puthier, D., and Gautheret, D. 2008, Entropy measures quantify global splicing

553        disorders in cancer. *PLoS Comput. Biol.*, **4**, 1–9.

554    47.    Kumar, U., Kumar, V., and Kapur, J. N. 1986, Normalized measures of entropy. *Int. J. Gen. Syst.*, **12**, 55–

555        69.

556    48.    Jaganathan, K., Kyriazopoulou Panagiotopoulou, S., McRae, J. F., et al. 2019, Predicting Splicing from

557        Primary Sequence with Deep Learning. *Cell*, **176**, 535-548.e24.

558    49.    Cheng, C. Y., Krishnakumar, V., Chan, A. P., Thibaud-Nissen, F., Schobel, S., and Town, C. D. 2017,

559        Araport11: a complete reannotation of the Arabidopsis thaliana reference genome. *Plant J.*, **89**, 789–804.

560    50.    Helvaci, C., Mordogan, H., Çolak, M., and Gündogan, I. 2004, Presence and distribution of lithium in borate

561        deposits and some recent lake waters of west-central turkey. *Int. Geol. Rev.*, **46**, 177–90.

562    51.    Pleiss, J. A., Whitworth, G. B., Bergkessel, M., and Guthrie, C. 2007, Rapid, Transcript-Specific Changes in

563        Splicing in Response to Environmental Stress. *Mol. Cell*, **27**, 928–37.

564    52.    Shi, H., Ishitani, M., Kim, C., and Zhu, J. K. 2000, The Arabidopsis thaliana salt tolerance gene SOS1

565        encodes a putative Na+/H+ antiporter. *Proc. Natl. Acad. Sci. U. S. A.*, **97**, 6896–901.

566    53.    Székely, G., Ábrahám, E., Cséplő, Á., et al. 2008, Duplicated P5CS genes of Arabidopsis play distinct roles

567        in stress regulation and developmental control of proline biosynthesis. *Plant J.*, **53**, 11–28.

568    54.    Martínez, O., and Reyes-Valdés, M. H. 2008, Defining diversity, specialization, and gene specificity in

569        transcriptomes through information theory. *Proc. Natl. Acad. Sci. U. S. A.*, **105**, 9709–14.

570    55.    Bou Sleiman, M., Frochaux, M. V., Andreani, T., Osman, D., Guigo, R., and Deplancke, B. 2020, Enteric

571        infection induces Lark-mediated intron retention at the 5′ end of Drosophila genes. *Genome Biol.*, **21**, 1–19.

572   56.   Pucker, B., and Brockington, S. F. 2018, Genome-wide analyses supported by RNA-Seq reveal non-
573         canonical splice sites in plant genomes. *BMC Genomics*, **19**, 1–13.

574   57.   Zhu, F. Y., Chen, M. X., Ye, N. H., et al. 2017, Proteogenomic analysis reveals alternative splicing and
575         translation as part of the abscisic acid response in Arabidopsis seedlings. *Plant J.*, **91**, 518–33.

576   58.   Li, J., Li, X., Guo, L., et al. 2006, A subgroup of MYB transcription factor genes undergoes highly
577         conserved alternative splicing in Arabidopsis and rice. *J. Exp. Bot.*, **57**, 1263–73.

578   59.   Laloum, T., Martín, G., and Duque, P. 2018, Alternative Splicing Control of Abiotic Stress Responses.
579         *Trends Plant Sci.*, **23**, 140–50.

580   60.   Feng, J., Li, J., Gao, Z., et al. 2015, SKIP Confers Osmotic Tolerance during Salt Stress by Controlling
581         Alternative Gene Splicing in Arabidopsis. *Mol. Plant*, **8**, 1038–52.

582   61.   Liu, H., Dai, J., Li, K., et al. 2022, Performance evaluation of computational methods for splice-disrupting
583         variants and improving the performance using the machine learning-based framework. *Brief. Bioinform.*, **23**,
584         bbac334.

585   62.   Eraslan, G., Avsec, Ž., Gagneur, J., and Theis, F. J. 2019, Deep learning: new computational modelling
586         techniques for genomics. *Nat. Rev. Genet.*, **20**, 389–403.

587   63.   Wu, H. J., Zhang, Z., Wang, J. Y., et al. 2012, Insights into salt tolerance from the genome of Thellungiella
588         salsuginea. *Proc. Natl. Acad. Sci.*, **109**, 12219–24.

589   64.   Liu, C., Wu, Y., Liu, Y., et al. 2021, Genome-wide analysis of tandem duplicated genes and their
590         contribution to stress resistance in pigeonpea (Cajanus cajan). *Genomics*, **113**, 728–35.

591   65.   Chaudhary, S., Khokhar, W., Jabre, I., et al. 2019, Alternative splicing and protein diversity: Plants versus
592         animals. *Front. Plant Sci.*, **10**, 1–14.

593   66.   Brown, J. W. S., Calixto, C. P. G., and Zhang, R. 2017, High quality reference transcript datasets hold the
594         key to transcript specific RNA sequencing analysis in plants. *New Phytol.*, **213**, 525–30.

595   67.   Huang, C. K., Lin, W. D., and Wu, S. H. 2022, An improved repertoire of splicing variants and their
596         potential roles in Arabidopsis photomorphogenic development. *Genome Biol.*, **23**, 1–28.

597   68.   Oh, D. H., Hong, H., Lee, S. Y., Yun, D. J., Bohnert, H. J., and Dassanayake, M. 2014, Genome structures
598         and transcriptomes signify niche adaptation for the multiple-ion-tolerant extremophyte Schrenkiella parvula.
599         *Plant Physiol.*, **164**, 2123–38.

600   69.   Martín, G., Márquez, Y., Mantica, F., Duque, P., and Irimia, M. 2021, Alternative splicing landscapes in
601         Arabidopsis thaliana across tissues and stress conditions highlight major functional differences with animals.
602         *Genome Biol.*, **22**, 1–26.

603   70.   Filichkin, S., Priest, H. D., Megraw, M., and Mockler, T. C. 2015, Alternative splicing in plants: Directing
604         traffic at the crossroads of adaptation and environmental stress. *Curr. Opin. Plant Biol.*, **24**, 125–35.

605   71.   Mastrangelo, A. M., Marone, D., Laidò, G., De Leonardis, A. M., and De Vita, P. 2012, Alternative
606         splicing: Enhancing ability to cope with stress via transcriptome plasticity. *Plant Sci.*, **185**–**186**, 40–9.

607   72.   Drechsel, G., Kahles, A., Kesarwani, A. K., et al. 2013, Nonsense-mediated decay of alternative precursor
608         mRNA splicing variants is a major determinant of the Arabidopsis steady state transcriptome. *Plant Cell*, **25**,

17

609    3726–42.

610    73.    Sterne-Weiler, T., Weatheritt, R. J., Best, A. J., Ha, K. C. H., and Blencowe, B. J. 2018, Efficient and

611          Accurate Quantitative Profiling of Alternative Splicing Patterns of Any Complexity on a Laptop. *Mol. Cell*,

612          **72**, 187-200.e6.

613    74.    Dankó, B., Szikora, P., Pór, T., Szeifert, A., and Sebestyén, E. 2022, SplicingFactory - Splicing diversity

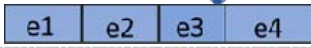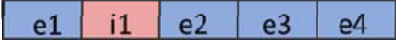614          analysis for transcriptome data. *Bioinformatics*, **38**, 384–90.

615

616
617
618
619    **Tables**
620    **Table 1.** Summary of *S. parvula* V2 genome updated with transcript models supported by Iso-Seq

| Category | Number of genes or transcripts |
| --- | --- |
| Genes supported by an Iso-Seq based transcript model | 11,348 |
| Total transcript models identified from Iso-Seq | 16,828 |
| Gene models identified with at least one new isoform | 7,028 |
| New genes identified from Iso-Seq | 301 |
| Gene models supplemented with UTR information using Iso-Seq | 11,348 |
| Genes with alternative splicing identified using Iso-Seq | 3,470 |
| Genes with alternative starts identified using Iso-Seq | 4,760 |
| Genes with alternative ends identified using Iso-Seq | 4,756 |
| Total number of protein coding transcripts annotated in the genome | 34,582 |

621

622    **Table 2. Major alternative splicing events identified using Iso-Seq reads in *S. parvula*.** Blue boxes
623    represent exons while red boxes or dash lines represent introns.

18

| Event type | # of events (%) |
|---|---|
| Reference gene model | 5,911 (100%) |
| Intron retention | 3,266 (55.2%) |
| Alternative 3' acceptor | 1,451 (24.5%) |
| Alternative 5' donor | 693 (11.7%) |
| Exon skipping | 283 (4.7%) |
| Alternative first exon | 189 (3.19%) |
| Alternative last exon | 18 (0.30%) |
| Mutually exclusive exons | 11 (0.18%) |

624
625
626    **Figure legends**

627    **Figure 1. Improved *S. parvula* gene models using full length transcripts. [A]** Majority of error corrected

628    circular consensus reads (CCS) contain full length reads with polyA sequences and an average length of 1.7

629    kb. **[B]** Length distribution of Iso-Seq based transcript models and *S. parvula* Reference genome v2.2 gene

630    models. **[C]** Percentage of mapped reads to the current genome and genome updated with Iso-Seq transcript

631    models. Data = mean ± SD. Dots represent 4 independent RNA-Seq datasets used in this study as biological

632    replicates (n = 12). Asterisk indicates significant difference ($p \leq 0.05$), determined by one-sided t-test.

633

634    **Figure 2. Genes with higher isoform diversity are associated with stress responses in *S. parvula***

635    **compared to** *A. thaliana*. [A] Number of isoforms of *S. parvula* and *A. thaliana* per single-copy ortholog

636    pairs given as a ratio. Blue and pink shaded areas indicate ortholog pairs where one species has more isoforms

637    than the other. [B] Enriched functions associated with ortholog pairs that show at least 2-fold difference in

638    isoform ratio between *S. parvula* and *A. thaliana*. Center line in the boxplots indicates median; box indicates

19

639      interquartile range (IQR); whiskers show $1.5 \times$ IQR. Asterisks indicate significant difference between isoform

640      distributions of the two species, measured by Wilcoxon rank sum test at $p$-value cutoff $\leq 0.05$.

641

642      **Figure 3. Isoform usage specificity between control and salt treated conditions. [A]** Shannon entropy

643      distribution of 1,678 ortholog pairs with at least two isoforms expressed per ortholog per species. Center line

644      in the boxplots indicates median; box indicates interquartile range (IQR); whiskers show $1.5 \times$ IQR. Each

645      treatment was compared to control according to Student's t test with $p$-values indicated above the relevant

646      pairs. **[B]** Shannon entropy change between salt and control conditions in *S. parvula* and *A. thaliana* ortholog

647      pairs in roots. Each dot represents an ortholog pair. Black lines indicate 0.5 entropy differences. Frequency

648      distribution of data are shown on the marginal plot. **[C]** Functionally enriched processes represented by

649      ortholog pairs in distinct categories of entropy shifts. A node in each cluster represents a gene ontology (GO)

650      term; size of a node represents the number of genes included in that GO term; the clusters that represent similar

651      functions share the same color and are given a representative cluster name and ID; and the edges between

652      nodes show shared genes between functions. All clusters included in the network have adj $p$-values $\leq 0.05$ with

653      false discovery rate correction applied.

654

655      **Figure 4. Genes differently spliced and differently expressed in response to salt stress. [A]** Number of

656      genes in *S. parvula* and *A. thaliana* that are differently regulated under salt stress. **[B]** Number of orthologs

657      that show differential splicing in *S. parvula* and *A. thaliana* root and shoot in response to salt. **[C]** Functionally

658      enriched processes represented by differently spliced genes in *S. parvula* and *A. thaliana*.

659

660      **Figure 5. Use of non-canonical splice sites in transcripts expressed under stress. [A]** Expression

661      distribution of transcripts that contain only canonical splice sites and transcripts with at least one non-canonical

662      splice site in roots and shoots of *S. parvula* (left panel) and *A. thaliana* (right panel). Asterisks indicate

663      significant difference ($p \leq 0.05$) of expression distributions between control and salt treated condition

664      measured by two-sided t-test. **[B]** Number of expressed non-canonically spliced transcripts as a % out of total

665      transcripts expressed in *S. parvula* and *A. thaliana* in response to salt. Significant differences between control

666      and stress conditions were tested using Fisher's exact test. **[C]** *S. parvula* genes that were enriched in non-

667      canonical splice sites. The y axis shows the $-\log_{10}$p-value for a test of excess of non-canonical splice sites

668      computed using a binomial test, where the probability of enrichment is calculated as the total non-canonical

669      splice sites divided by the total number of splice sites per gene, ordered in the chromosomal order (x-axis) for

670      the *S. parvula* genome. Genes with a high enrichment for non-canonical splicing are labeled. Red line indicates

671      the $-\log_{10}$p corresponding to adjusted $p$-value of 0.05. **[D]** Functional processes enriched in genes detected to

672      be non-canonically spliced in [C]. A node in each cluster represents a gene ontology (GO) term; size of a node

673      represents the number of genes included in that GO term; the clusters that represent similar functions share the

20

674    same color and are given a representative cluster name and ID; and the edges between nodes show the

675    connectivity of genes between functions. All clusters included in the network have adj $p$-values ≤0.05 with

676    false discovery rate correction applied. More significant values are represented by darker node colors. The

677    right panel shows the sub-clustered functions represented by the largest cluster C1 in the left panel.

678

679    **Figure 6. Genome wide prediction of splice sites for *S. parvula* using a deep neural network. [A]** Training

680    and testing with *A. thaliana* and application of the SpliceAi model to *S. parvula.* **[B]** Overlap between the

681    observed and predicted splice sites for *S. parvula* protein coding gene models. A probability score ≥ 0.6 was

682    used for the predicted splice sites. **[C]** Functional processes enriched in genes observed and predicted to have

683    more than one isoform in the *S. parvula* genome. A node in each cluster represents a gene ontology (GO) term;

684    size of a node represents the number of genes included in that GO term; the clusters that represent similar

685    functions share the same color and are given a representative cluster name; and the edges between nodes show

686    the connectivity of genes between functions. All clusters included in the network have adj $p$-values ≤0.05 with

687    false discovery rate correction applied.

688

689    **Figure S1. UTR length distribution of transcript models in *S. parvula* and *A. thaliana*. [A]** 5' UTR and

690    **[B]** 3' UTR length distributions. The distributions were obtained from 16,828 *S. parvula* and 41,064 *A.*

691    *thaliana* transcript models

692

693    **Figure S2. Novel gene *SP4G29525* (*SpMHX1*) annotated using Iso-Seq transcript models.** Two transcript

694    models were detected with Iso-Seq full length reads for both *MHX* and *PHR2* gene models.

695

696    **Figure S3. Independent detection of selected isoforms first identified with Iso-Seq in *S. parvula***

697    **transcriptome. [A]** Transcript models of selected genes with isoform IDs. Coding region and UTR regions are

698    indicated in dark and light shades. Locations of primer binding sites are shown by arrows. Exons are

699    numbered. Introns are given as connecting lines between exons. Identical exon structures past the 2$^{nd}$ exon

700    between isoforms are represented by dashed lines. **[B]** Predicted protein coding regions and functional domains

701    of the corresponding transcripts. Functional domains for each transcript are marked as colored blocks. **[C]** Gel

702    electrophoresis images of amplified transcripts obtained from RT-PCR using primers indicated in A. Arrow

703    heads indicate the expected size of the amplified product.

704

705    **Figure S4. *SOS1* and *P5CS1* isoform diversity in *S. parvula*. [A]** *SpSOS1* isoforms expressed above 0.5

706    TPM in all conditions. *SpSOS1 (v2)* serves as the primary gene model annotated in the current genome

707    annotation. **[B]** *SpP5CS1* isoforms expressed above 0.5 TPM in all conditions. *SpP5CS1 (v2)* serves as the

708    primary gene model annotated in the current genome annotation.  Data = mean ± SD (n = 3).

709

710     **Figure S5. Differential regulation of orthologs in *S. parvula* and *A. thaliana* in response to salt stress. [A]**

711     Root and **[B]** shoot. UpSet plot numbers represent number of orthologs. DS - Differently spliced; DE –

712     Differently expressed; Sp - *S. parvula*; At - *A. thaliana*.

713

714     **Figure S6. Probability of splice sites identified using SpliceAi from 5' to 3' for *S. parvula* gene models**.

715     Dashed line indicates the probability thresholds used to predict a splice site.

716
717
718     **Supplementary Table 1.** Primers used for RT-PCR.

719
720     **Supplementary Table 2.** Ortholog pairs between *S. parvula* and *A. thaliana* isoforms, annotations, and
721     isoform ratio.

722
723     **Supplementary Table 3.** Differently spliced and differently expressed genes in *S. parvula* and *A. thaliana*.

724
725     **Supplementary Table 4.** Enrichment for non-canonical over canonical splice in *S. parvula*.

726
727     **Supplementary Table 5.** Predicted and observed splice sites in the *S. parvula* genome. Link:
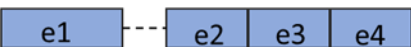728     https://github.com/wchathura/Iso-Seq_Dataset/blob/main/Supplemntry_table_5.txt
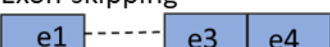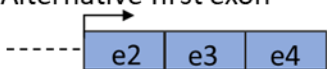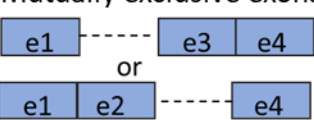729

**Figure 1. Improved *S. parvula* gene models using full length transcripts. [A]** Majority of error corrected circular consensus reads (CCS) contain full length reads with polyA sequences and an average length of 1.7 kb. **[B]** Length distribution of Iso-Seq based transcript models and *S. parvula* Reference genome v2.2 gene models. **[C]** Percentage of mapped reads to the current genome and genome updated with Iso-Seq transcript models. Data = mean ± SD. Dots represent 4 independent RNA-Seq datasets used in this study as biological replicates (n = 12). Asterisk indicates significant difference ($p \leq 0.05$), determined by one-sided t-test.
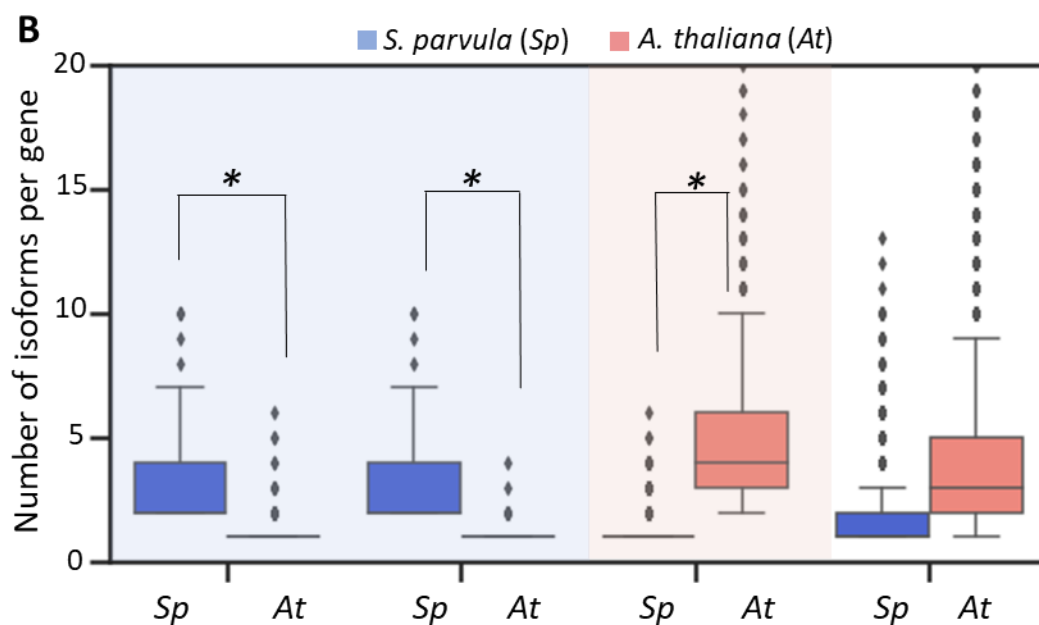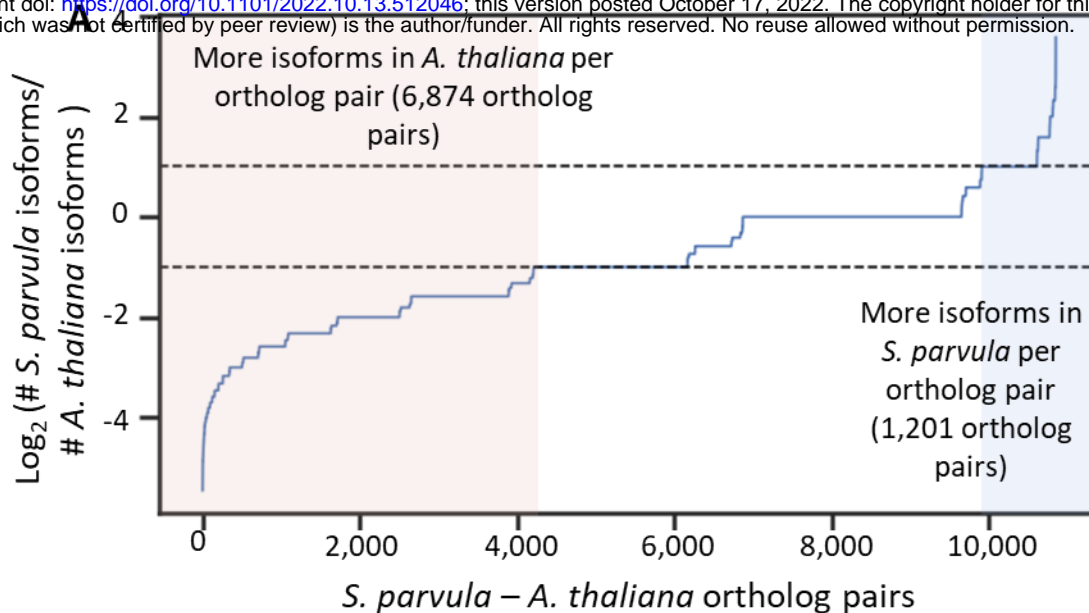
**Table 1.** Summary of *S. parvula* V2 genome updated with transcript models supported by Iso-Seq.

| Category | Number of genes or transcripts |
|---|---|
| Genes supported by an Iso-Seq based transcript model | 11,348 |
| Total transcript models identified from Iso-Seq | 16,828 |
| Gene models identified with at least one new isoform | 7,028 |
| New genes identified from Iso-Seq | 301 |
| Gene models supplemented with UTR information using Iso-Seq | 11,348 |
| Genes with alternative splicing identified using Iso-Seq | 3,470 |
| Genes with alternative starts identified using Iso-Seq | 4,760 |
| Genes with alternative ends identified using Iso-Seq | 4,756 |
| Total number of protein coding transcripts annotated in the genome | 34,582 |

**Table 2. Major alternative splicing events identified using Iso-Seq reads in *S. parvula*. Blue boxes represent exons while red boxes or dash lines represent introns.**
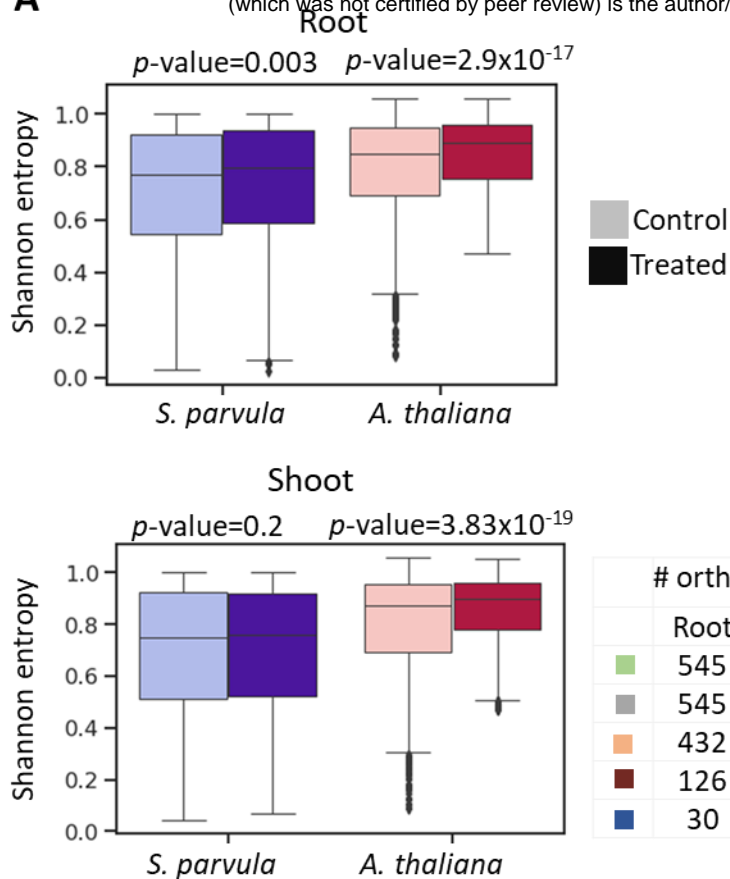
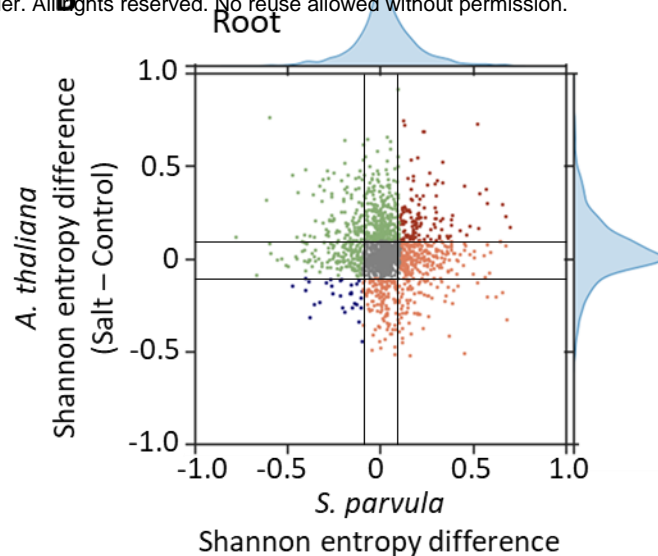| Event type | # of events (%) |
|---|---|
| Reference gene model | 5,911 (100%) |
| Intron retention | 3,266 (55.2%) |
| Alternative 3' acceptor | 1,451 (24.5%) |
| Alternative 5' donor | 693 (11.7%) |
| Exon skipping | 283 (4.7%) |
| Alternative first exon | 189 (3.19%) |
| Alternative last exon | 18 (0.30%) |
| Mutually exclusive exons | 11 (0.18%) |

**Figure 2. Genes with higher isoform diversity are associated with stress responses in *S. parvula* compared to *A. thaliana*.** [A] Number of isoforms of *S. parvula* and *A. thaliana* per single-copy ortholog pairs given as a ratio. Blue and pink shaded areas indicate ortholog pairs where one species has more isoforms than the other. [B] Enriched functions associated with ortholog pairs that show at least 2-fold difference in isoform ratio between *S. parvula* and *A. thaliana*. Center line in the boxplots indicates median; box indicates interquartile range (IQR); whiskers show 1.5 × IQR. Asterisks indicate significant difference between isoform distributions of the two species, measured by Wilcoxon rank sum test at *p*-value cutoff ≤ 0.05.
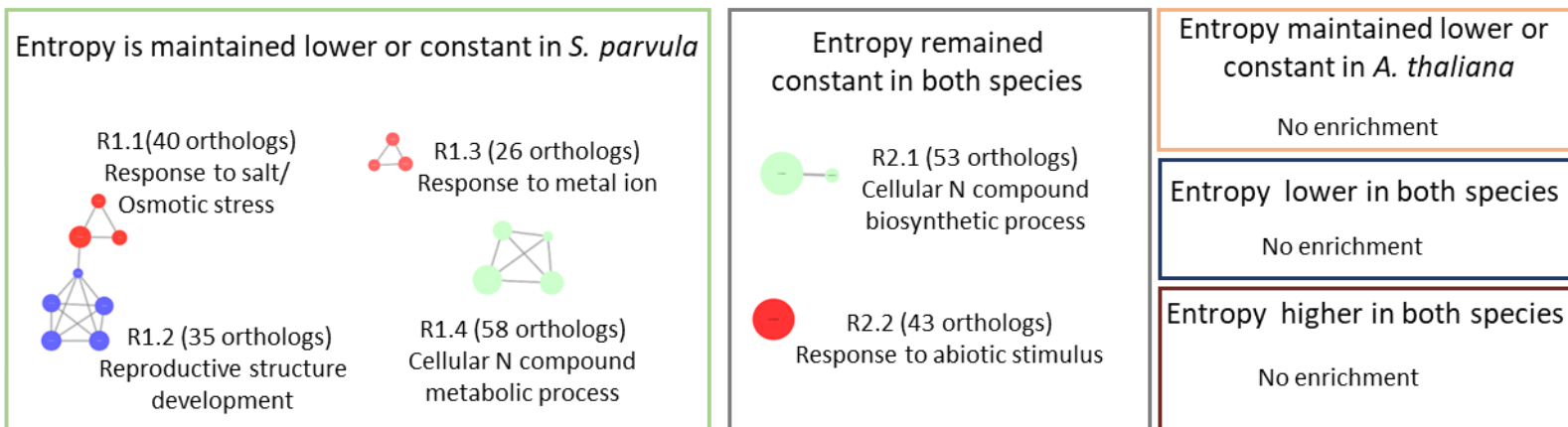
**A** Root

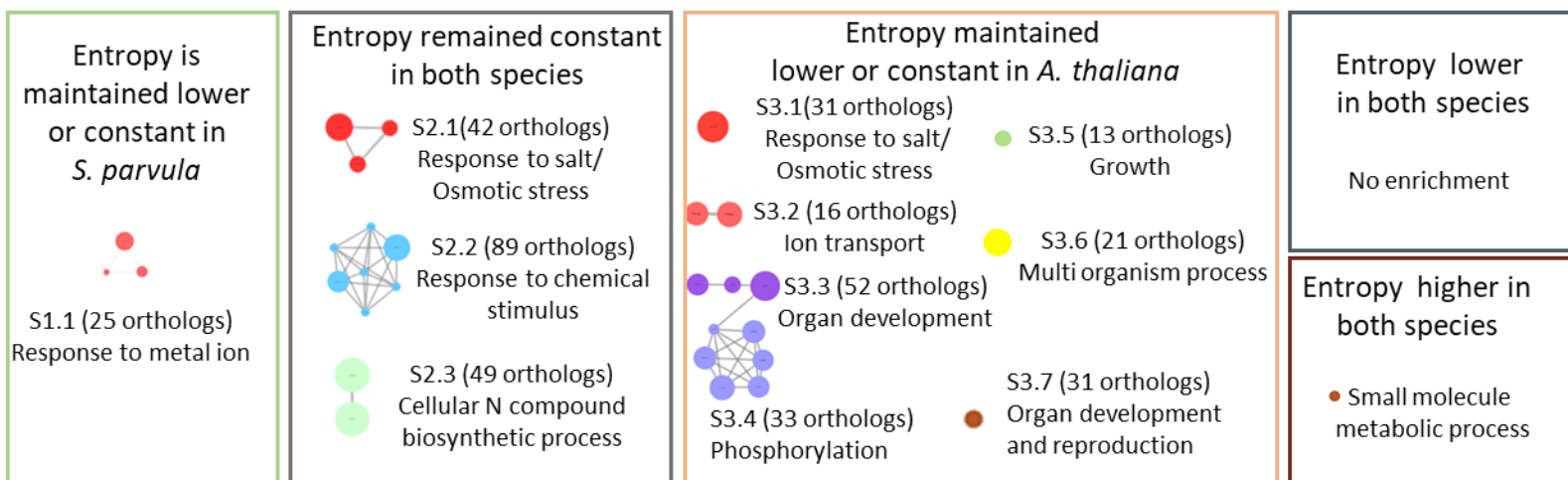p-value=0.003    p-value=2.9x10$^{-17}$

Control
Treated

Shoot

p-value=0.2    p-value=3.83x10$^{-19}$

**B** Root

S. parvula
Shannon entropy difference
(Salt − Control)

| | # ortholog pairs | | |
|---|---|---|---|
| | Root | Shoot | |
| (green) | 545 | 362 | Entropy is maintained lower or constant in S. parvula |
| (grey) | 545 | 513 | Entropy remained constant in both species |
| (orange) | 432 | 575 | Entropy maintained lower or constant in A. thaliana |
| (dark red) | 126 | 93 | Entropy higher in both species |
| (blue) | 30 | 49 | Entropy lower in both species |

**C** Root

Entropy is maintained lower or constant in S. parvula

R1.1(40 orthologs)
Response to salt/
Osmotic stress

R1.3 (26 orthologs)
Response to metal ion

R1.2 (35 orthologs)
Reproductive structure
development

R1.4 (58 orthologs)
Cellular N compound
metabolic process

Entropy remained
constant in both species

R2.1 (53 orthologs)
Cellular N compound
biosynthetic process

R2.2 (43 orthologs)
Response to abiotic stimulus

Entropy maintained lower or
constant in A. thaliana

No enrichment

Entropy lower in both species

No enrichment

Entropy higher in both species

No enrichment

Shoot

Entropy is
maintained lower
or constant in
S. parvula

S1.1 (25 orthologs)
Response to metal ion

Entropy remained constant
in both species

S2.1(42 orthologs)
Response to salt/
Osmotic stress

S2.2 (89 orthologs)
Response to chemical
stimulus

S2.3 (49 orthologs)
Cellular N compound
biosynthetic process

Entropy maintained
lower or constant in A. thaliana

S3.1(31 orthologs)
Response to salt/
Osmotic stress

S3.5 (13 orthologs)
Growth

S3.2 (16 orthologs)
Ion transport

S3.6 (21 orthologs)
Multi organism process

S3.3 (52 orthologs)
Organ development

S3.4 (33 orthologs)
Phosphorylation

S3.7 (31 orthologs)
Organ development
and reproduction

Entropy lower
in both species

No enrichment
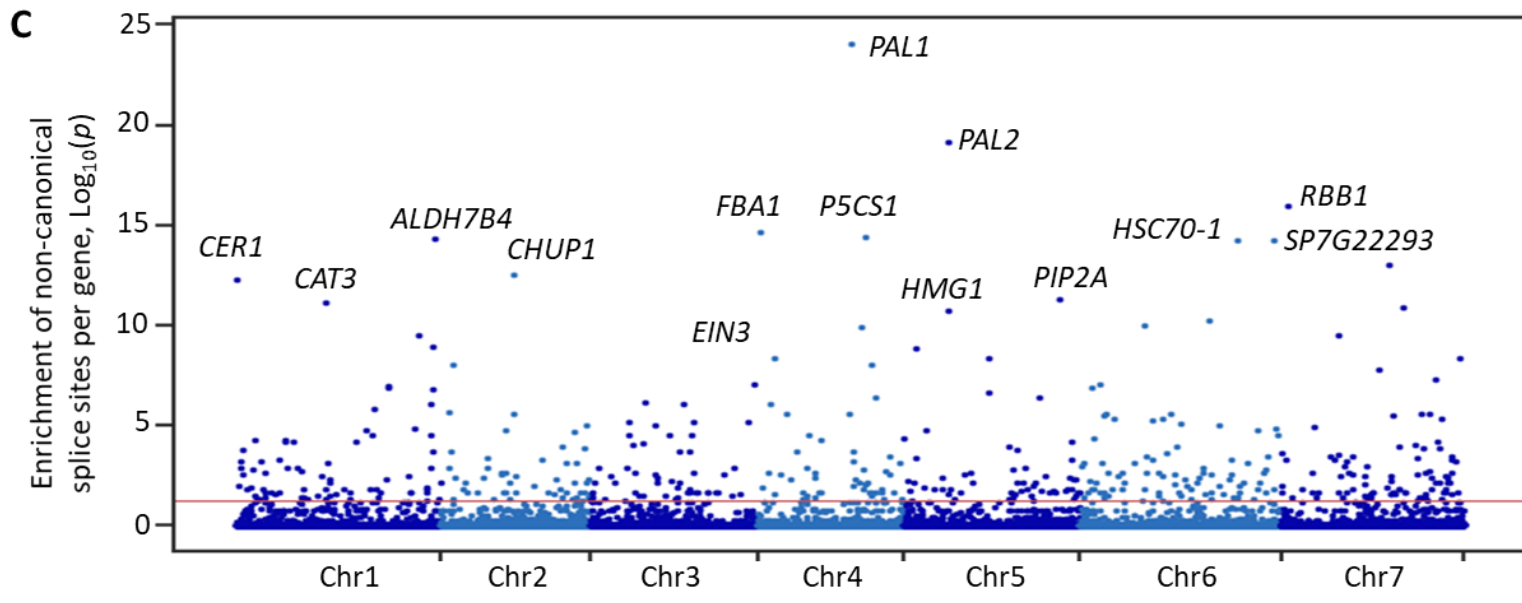
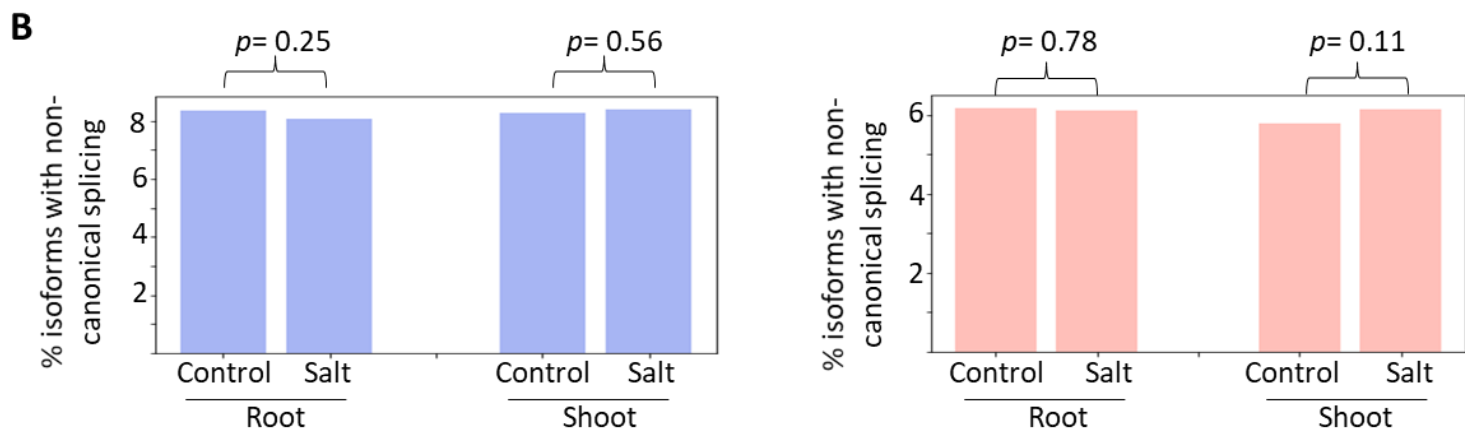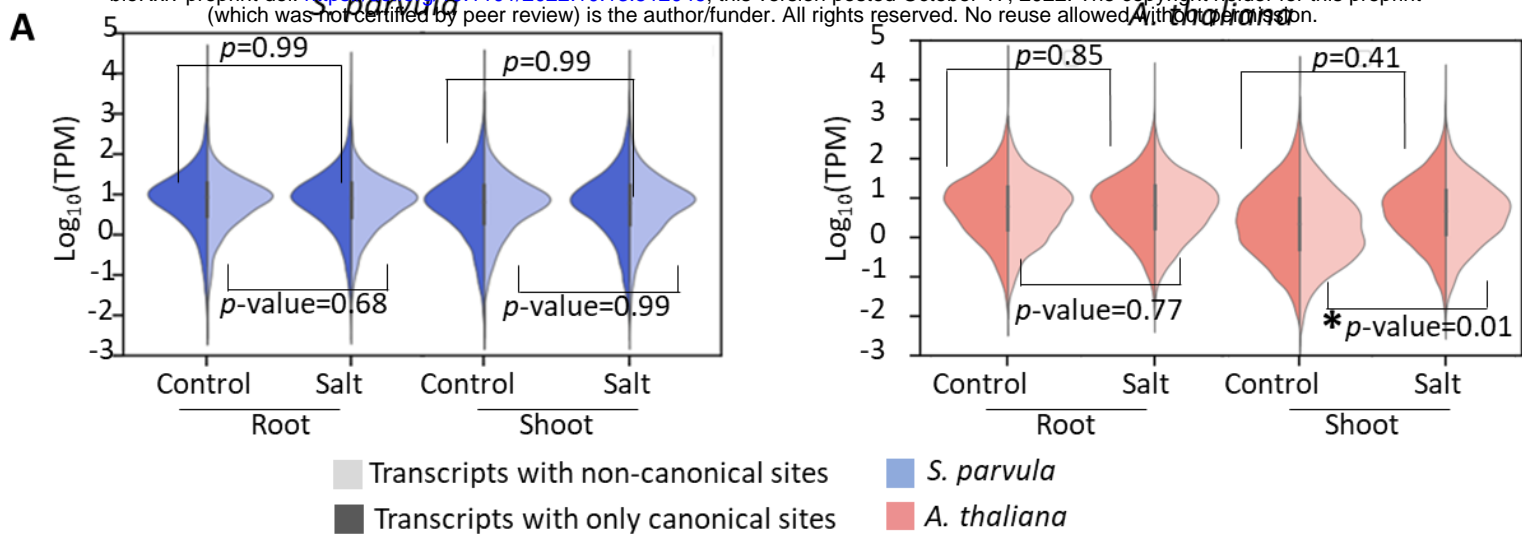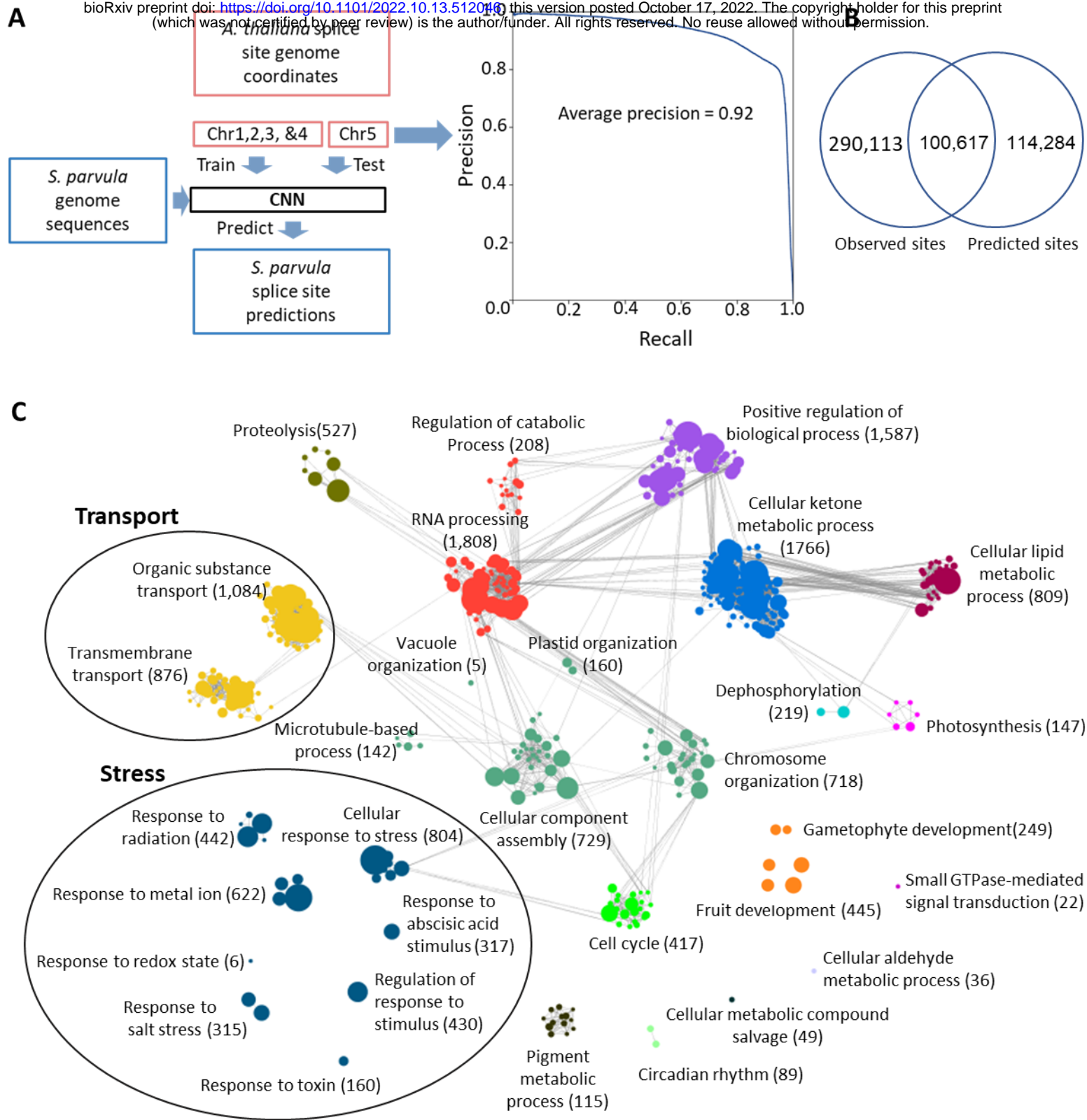Entropy higher in
both species

• Small molecule
metabolic process

**Figure 3. Isoform usage specificity between control and salt treated conditions. [A]** Shannon entropy distribution of

1,678 ortholog pairs with at least two isoforms expressed per ortholog per species. Center line in the boxplots indicates median; box indicates interquartile range (IQR); whiskers show 1.5 × IQR. Each treatment was compared to control according to Student's t test with p-values indicated above the relevant pairs. **[B]** Shannon entropy change between salt and control conditions in *S. parvula* and *A. thaliana* ortholog pairs in roots. Each dot represents an ortholog pair. Black lines indicate 0.5 entropy differences. Frequency distribution of data are shown on the marginal plot. **[C]** Functionally enriched processes represented by ortholog pairs in distinct categories of entropy shifts. A node in each cluster represents a gene ontology (GO) term; size of a node represents the number of genes included in that GO term; the clusters that represent similar functions share the same color and are given a representative cluster name and ID; and the edges between nodes show shared genes between functions. All clusters included in the network have adj p-values ≤0.05 with false discovery rate correction applied.

**Figure 4. Genes differently spliced and differently expressed in response to salt stress. [A]** Number of genes in *S. parvula* and *A. thaliana* that are differently regulated under salt stress. **[B]** Number of orthologs that show differential splicing in *S. parvula* and *A. thaliana* root and shoot in response to salt. **[C]** Functionally enriched processes represented by differently spliced genes in *S. parvula* and *A. thaliana*.

**Figure 5. Use of non-canonical splice sites in transcripts expressed under stress. [A]** Expression distribution of transcripts that contain only canonical splice sites and transcripts with at least one non-canonical splice site in roots and shoots of *S. parvula* (left panel) and *A. thaliana* (right panel). Asterisks indicate significant difference of expression distributions between control and salt treated condition measured by two-sided t-test at p-value ≤ 0.05. **[B]** Number of expressed non-canonically spliced transcripts as a % out of total transcripts expressed in *S. parvula* and *A. thaliana* in response to salt. Significant differences between control and stress conditions were tested using Fisher's exact test. **[C]** *S. parvula* genes that were enriched in non-canonical splice sites. The y axis shows the $-\log_{10}$p-value for a test of excess of non-canonical splice sites computed using a binomial test, where the probability of enrichment is calculated as the total non-canonical splice sites divided by the total number of splice sites per gene, ordered in the chromosomal order (x-axis) for the *S. parvula* genome. Genes with a high enrichment for non-canonical splicing are labeled. Red line indicates the $-\log_{10}$p corresponding to adjusted p-value of 0.05. **[D]** Functional processes enriched in genes detected to be non-canonically spliced in [C]. A node in each cluster represents a gene ontology (GO) term; size of a node represents the number of genes included in that GO term; the clusters that represent similar functions share the same color and are given a representative cluster name and ID; and the edges between nodes show the connectivity of genes between functions. All clusters included in the network have adj p-values ≤0.05 with false discovery rate correction applied. More significant values are represented by darker node colors. The right panel shows the sub-clustered functions represented by the largest cluster C1 in the left panel.

**Figure 6. Genome wide prediction of splice sites for *S. parvula* using a deep neural network. [A]** Training and testing with *A. thaliana* and application of the SpliceAi model to *S. parvula.* **[B]** Overlap between the observed and predicted splice sites for *S. parvula* protein coding gene models. *A* probability score of ≥0.6 was used for the predicted splice sites. **[C]** Functional processes enriched in genes observed and predicted to have more than one isoform in the *S. parvula* genome. A node in each cluster represents a gene ontology (GO) term; size of a node represents the number of genes included in that GO term; the clusters that represent similar functions share the same color and are given a representative cluster name; and the edges between nodes show the connectivity of genes between functions. All clusters included in the network have adj p-values ≤0.05 with false discovery rate correction applied.