

1 ***In vivo* nuclear RNA structurome reveals RNA-structure regulation of mRNA processing in**
2 **plants**

3 **Zhenshan Liu^{1,2}, Qi Liu^{1,2}, Xiaofei Yang¹, Yueying Zhang¹, Matthew Norris¹, Xiaoxi Chen¹,**
4 **Jitender Cheema¹, Yiliang Ding^{1*}**

5

6 ¹Department of Cell and Developmental Biology, John Innes Centre, Norwich Research Park,
7 Norwich NR4 7UH, United Kingdom

8 ²These authors contributed equally

9 ***Correspondence: yiliang.ding@jic.ac.uk**

10

11 **Abstract**

12 mRNA processing is critical for gene expression. A challenge in regulating mRNA processing is
13 how to recognize the actual mRNA processing sites, such as splice and polyadenylation sites,
14 when the sequence content is insufficient for this purpose. Previous studies suggested that RNA
15 structure affects mRNA processing. However, the regulatory role of RNA structure in mRNA
16 processing remains unclear. Here, we performed *in vivo* selective 2'-hydroxyl acylation analysed
17 by primer extension (SHAPE) chemical profiling on *Arabidopsis* nuclear RNAs and generated
18 the *in vivo* nuclear RNA structure landscape. We found that nuclear mRNAs fold differently from
19 cytosolic mRNAs. Notably, we discovered a two-nucleotide single-stranded RNA structure
20 feature upstream of 5' splice sites that is strongly associated with splicing and the selection of
21 alternative 5' splice sites. Moreover, we found the single-strandedness of branch point is also
22 associated with 3' splice site recognition. We also identified an RNA structure feature comprising
23 two close-by single-stranded regions that is specifically associated with both polyadenylation and
24 alternative polyadenylation events. Our work demonstrates an RNA structure regulatory
25 mechanism for mRNA processing.

26 **Introduction**

27 In eukaryotes, mRNAs undergo several processing steps including 5' capping, splicing, and
28 3' cleavage/polyadenylation to become functional mature mRNAs. Thus, mRNA processing plays
29 a critical role during gene expression^{1,2}. Over past decades, a key question is how mRNA
30 processing sites, such as polyadenylation and splice sites, are precisely recognized in the
31 transcriptome, particularly from surrounding sites with similar sequence content^{3,4}. For instance,
32 5' splice site recognition was found to be not always dependent on the sequence content of U1
33 snRNA binding motif. Some 5' splice sites were selected over those flanking sites with better
34 complementarity to U1 snRNA binding sequence⁴. In case-by-case studies, quite a number of
35 RNA binding proteins have been identified that contribute to the recognition of actual
36 polyadenylation and splice sites^{4,5}. However, a general regulatory mechanism that recognizes
37 actual sites during mRNA processing is lacking. As an intrinsic characteristic of RNA molecules,
38 RNA structure was suggested to be involved in mRNA processing⁶. Previous individual studies
39 suggested that RNA structure can affect polyadenylation and splicing⁷⁻¹³. Yet, how RNA structure
40 contributes to the recognition of polyadenylation and splice sites, in general, remains elusive.

41 With recent advances in RNA structure profiling¹⁴⁻¹⁶, more attention has been drawn toward
42 understanding how RNA structure influences mRNA processing. Previous *in vitro* enzymatic
43 RNA structure profiling (utilizing RNases that selectively cleave either single-stranded or double-
44 stranded nucleotides) in *Arabidopsis* nuclear RNAs, found that the 5' end of introns were more
45 double-stranded compared to upstream exons, and the 3' end of introns were more single-stranded
46 compared to upstream intron regions¹⁴. However, no significant structure signatures were

47 identified for either polyadenylation or alternative polyadenylation sites¹⁴. This may be due to
48 limitations imposed by using RNases, which are quite bulky and less sensitive in detecting
49 specific RNA structures, compared to the relatively small chemicals used for RNA structure
50 probing^{17,18}. Furthermore, several previous studies have shown that *in vitro* RNA structures were
51 not able to reflect the proper folding status of RNAs in living cells^{19,20}. A recent *in vivo* dimethyl
52 sulfate (DMS) RNA structure profiling study on human mature mRNAs identified RNA structure
53 features for polyadenylation (poly(A)) sites¹⁵. A more folded structure downstream of the
54 polyadenylation signal motif was identified that facilitated polyadenylation¹⁵. However,
55 mammalian RNAs were found to adopt different structure conformations in different cellular
56 compartments²¹. Thus, the structure of mature mRNAs in the cytosol is likely to be different from
57 the structure of pre-mRNA in the nucleus. If so, mature mRNA structures are unlikely to reveal
58 the role of RNA structure in polyadenylation. A notable limitation of this DMS method is the loss
59 of RNA structure information for the half transcriptome because DMS only detects structure
60 information of As and Cs, lacking the base-pairing status of Us and Gs.

61 Here, we studied the role of RNA structure in mRNA processing by performing *in vivo*
62 SHAPE (Selective 2' Hydroxyl Acylation analysed by Primer Extension) chemical probing on
63 *Arabidopsis thaliana* nuclear RNAs, to generate the first *in vivo* RNA structure landscape with all
64 four nucleotides in plants. We found that nuclear mRNA structures are globally different from
65 cytosolic mRNA structures in *Arabidopsis*. Our study further successfully dissected pre-mRNA
66 structure features before mRNA processing and determined the regulatory role of RNA structure
67 during mRNA maturation.

68 **Results**

69 **Nuclear SHAPE-Structure-Seq generates *in vivo* RNA structure landscape of *Arabidopsis*** 70 **nuclear RNAs with high coverage and accuracy**

71 To investigate the role of RNA structure in mRNA processing, we performed SHAPE chemical
72 probing²² on *Arabidopsis* nuclear RNAs and generated the first *in vivo* RNA structure profiles
73 with all four nucleotides in plants. Firstly, SHAPE reagent (2-methylnicotinic acid imidazolide,
74 NAI) treatment was applied on 5-day-old *Arabidopsis* seedlings²² (Fig. 1a). Intact nuclei were
75 isolated and nuclear RNAs were extracted. The intactness of isolated nuclei was confirmed by
76 microscopy imaging with DAPI staining²³ (Supplementary Fig. 1a). Enrichment of nuclear
77 histone H3 protein and absence of cytoplasmic protein PEPC (Phosphoenolpyruvate carboxylase)
78 in the isolated nucleus, further confirmed the high purity and quality of the isolated nuclei
79 (Supplementary Fig. 1b). We generated two independent biological replicates of (+)SHAPE
80 (samples with SHAPE treatment) and (-)SHAPE (control samples without SHAPE treatment)
81 Structure-Seq libraries for high-throughput sequencing^{24,25} (Fig. 1a and Supplementary Fig. 2, see
82 Methods). Given that interactions between RNA and RNA binding proteins can prevent the

83 SHAPE modification²², we also performed SHAPE treatment on nuclear RNAs after removing
84 proteins thus generating *deproteinized* nuclear SHAPE-Structure-Seq libraries in parallel (Fig. 1a)
85 to assess any effect on SHAPE modification signals caused by protein protection. The
86 *deproteinized* libraries were designed to preserve RNA secondary structure after cell lysis and
87 protein removal but not subjected to RNA denaturing under high temperature. Thus, the
88 *deproteinized* condition here is still *in vitro* condition. Over 616 million 100bp paired-end reads
89 per library were generated and further mapped onto *Arabidopsis* genome sequences (TAIR10)
90 with additional alternative spliced isoforms annotated from AtRTD2 database²⁶ (Supplementary
91 Table 1).

92 Nucleotide modification in both (+)SHAPE and (-)SHAPE libraries were highly concordant,
93 with slight enrichment in (+)SHAPE shown for As and Us over Cs and Gs, as expected, since As
94 and Us tend to be more single-stranded than Cs and Gs (Supplementary Fig. 3a). The high
95 correlation of mRNA abundance between the two biological replicates indicated the high
96 reproducibility of our nuclear SHAPE-Structure-Seq libraries (Supplementary Fig. 3b). To further
97 validate the reproducibility of our SHAPE structure probing, we compared SHAPE reactivity
98 profiles of U1 and U12 snRNA between the two biological replicates and noted a high correlation
99 between them (Pearson correlation coefficient=0.93-0.97) (Supplementary Fig. 3c). Thus, we
100 merged these two biological replicates for further RNA structure analysis.

101 We assessed both the sequencing reads coverage and reverse-transcription stop counts of our
102 nuclear SHAPE-Structure-Seq libraries. Notably, more than 20,752 genes had at least 10 reads
103 per nucleotide coverage (Supplementary Fig. 4a), among which more than 12,366 genes reached
104 the threshold of at least one reverse-transcription stop count per nucleotide for RNA structure
105 analysis (Supplementary Fig. 4b). To assess the accuracy of our RNA structure profiling, we
106 compared SHAPE reactivity profiles of U1 and U12 snRNAs with their phylogenetically derived
107 structures, which are evolutionarily conserved structures and are the closest models of *in vivo*
108 structure^{22,27}. Overall, the SHAPE reactivities were consistent with phylogenetically derived RNA
109 structures where high SHAPE reactivities were observed in single-stranded regions, while low
110 SHAPE reactivities were at double-stranded nucleotides (Fig. 1b,c, Supplementary Table 2). Both
111 U1 and U12 snRNAs interact with Sm proteins to form small nuclear ribonucleoparticle
112 structures^{22,27}. We also found that SHAPE reactivities at Sm protein binding sites of U1 and U12
113 snRNA were significantly higher in the *deproteinized* rather than *in vivo* condition (Fig. 1b,c). In
114 addition, global SHAPE reactivities were also found to be significantly higher in the *deproteinized*
115 condition compared to the *in vivo* condition suggesting that absence of protein protection in the
116 *deproteinized* condition allowed nucleotide modification by SHAPE (Supplementary Fig. 5).
117 Collectively, these results indicated that our nuclear SHAPE-Structure-Seq method can accurately
118 probe *in vivo* RNA structures of nuclear RNAs.

119 **Nuclear mRNA structures are globally different from cytosolic mRNA structures**

120 Cytosolic mRNAs are the processed products from nuclear mRNAs, thus they share the same
121 sequences. However, whether they share the same RNA structure features remains unclear. To
122 address this question, we generated *in vivo* SHAPE-Structure-Seq libraries of *Arabidopsis*
123 cytosolic mRNAs in parallel. We then compared these libraries with our *in vivo* nuclear SHAPE-
124 Structure-Seq libraries (**Supplementary Data 1 and 2**). Firstly, we compared average SHAPE
125 reactivities of exons between nuclear and cytosolic mRNAs. We found that exons in cytosolic
126 mRNAs had significantly higher average SHAPE reactivities than those in nuclear mRNAs,
127 suggesting exons in cytosolic mRNAs tended to be more single-stranded than those in nuclear
128 mRNAs (**Fig. 2a**). This result is also similar with that observed in mammals²¹. We further
129 compared average SHAPE reactivities in different genic regions of exons: the 5' untranslated
130 region (5'UTR), the coding region (CDS) and the 3' untranslated region (3'UTR), between
131 nuclear and cytosolic mRNAs. Notably, we found that average SHAPE reactivities in both 5'UTR
132 and 3'UTR were significantly higher in nuclear mRNAs than those in cytosolic mRNAs (**Fig. 2b**).
133 In contrast, significantly lower average SHAPE reactivities were observed in nuclear mRNA CDS
134 regions compared to those in cytosolic mRNAs (**Fig. 2b**). Previous studies on total mRNAs
135 dominated by cytosolic mRNAs observed unique structure features across translation start and
136 stop codons that were associated with translation^{24,28-30}. Consistent with these observations, we
137 also found higher SHAPE reactivities upstream of start codons, lower SHAPE reactivities
138 downstream of start codons, and higher SHAPE reactivities at stop codons compared to flanking
139 regions (**Fig. 2c,d**) in our cytosolic SHAPE-Structure-Seq libraries. We then compared average
140 SHAPE reactivity profiles between nuclear and cytosolic mRNAs across these two sites.
141 Significantly higher SHAPE reactivities downstream of start codons and significantly lower
142 SHAPE reactivities at stop codons in nuclear mRNAs were observed, compared to those in
143 cytosolic mRNAs (**Fig. 2c,d**). Taken together, nuclear mRNA structures are globally different
144 from cytosolic mRNA structures, which implies nuclear and cytosolic mRNAs might adopt
145 different structures to serve their respective biological functions, e.g. translation in the cytosol
146 and mRNA processing in the nucleus. Therefore, we further investigated how nuclear mRNA
147 structures are associated with mRNA processing.

148 **Distinctive pre-mRNA structure features are strongly associated with both splicing and** 149 **alternative splicing**

150 Splicing is a key mRNA processing step that was previously suggested to be influenced by RNA
151 structure⁸. Since only pre-mRNA structure before splicing (unspliced primary transcripts) can be
152 used for dissecting the mechanism underpinning splicing, we firstly assessed whether pre-mRNAs
153 were enriched in our nuclear SHAPE-Structure-Seq data. We found that the expression abundance
154 of constitutively spliced introns was much higher in our nuclear SHAPE-Structure-Seq libraries

155 compared to our cytosolic SHAPE-Structure-Seq libraries, indicating high enrichment of pre-
156 mRNAs in our nuclear SHAPE-Structure-Seq data (Supplementary Fig. 6). Since nuclear mRNAs
157 still contain spliced transcripts, we only used reads mapped across exon-intron junctions and in
158 intron regions for generating SHAPE reactivity profiles to obtain accurate RNA structure
159 information of pre-mRNAs before splicing (See details in Methods). Also, to eliminate any
160 ambiguous reads assignment at the conserved dinucleotide AG at 3' splice sites (3'ss), we only
161 calculated SHAPE reactivities across 5' splice site (5'ss) and the whole intron except for AG at
162 3'ss (See details in Methods).

163 In addition to generating RNA structure information of pre-mRNAs, we also calculated the
164 splicing efficiency for each intron to measure the outcome for splicing events (Supplementary Fig.
165 7, See details in Methods). Since most of the introns showed either very high ($\geq 90\%$) or very
166 low ($\leq 10\%$) splicing efficiencies, two groups of splicing events were classified: spliced events
167 (splicing efficiency $\geq 90\%$, 32,522 spliced events were identified, Supplementary Data 3) and
168 unspliced events (i.e. intron retention, splicing efficiency $\leq 10\%$, 4,056 unspliced events were
169 identified, Supplementary Data 3). We compared average SHAPE reactivity profiles between
170 these two groups of splicing events. Although the exon-intron regions of both spliced and
171 unspliced events shared similar nucleotide compositions (Supplementary Fig. 8), distinctive
172 SHAPE reactivity profiles were observed between these two groups (Fig. 3a,b). Specifically, we
173 found that *in vivo* SHAPE reactivities at the -1 position immediately upstream of 5'ss were notably
174 higher for spliced events compared to unspliced events (Fig. 3a). Similarly, SHAPE reactivities
175 at the -1 and -2 positions upstream of 5'ss were significantly higher in spliced events than those
176 in unspliced events for the *deproteinized* condition. These findings indicated that the -1 and -2
177 nucleotides upstream of 5'ss tended to be more single-stranded in spliced events compared to
178 unspliced events (Fig. 3b). We further assessed sequence content across 5'ss in both spliced and
179 unspliced events and found no sequence preference between these two groups (Supplementary
180 Fig. 8). Thus, our results suggested that this distinctive structure signature was associated with
181 splicing events, but not due to sequence preference.

182 We then assessed RNA structure features for branch sites and 3'ss regions, which are
183 important for 3'ss recognition during splicing¹. To assess RNA structure features at branch sites,
184 we predicted branch sites using SVM-BPfinder³¹. Higher SHAPE reactivities were observed at
185 branch points under both *in vivo* and *deproteinized* conditions for spliced events compared to
186 unspliced events, indicating single-strandedness at the branch point was associated with splicing
187 (Fig. 3a,b). SHAPE reactivities of regions immediately upstream of dinucleotide AG at 3'ss (from
188 -7 to -4 positions) were relatively lower than flanking regions (Fig. 3a,b). However, there was no
189 significant SHAPE reactivity difference between spliced and unspliced events at 3'ss regions,
190 indicating no direct association with splicing. Therefore, both RNA structure features upstream
191 of 5'ss and at the branch point in pre-mRNAs were associated with splicing.

192 We then explored whether these RNA structure features are also associated with splice site
193 selection in alternative splicing events. Firstly, we identified alternative 5'ss events from genome
194 annotation and selected those pre-mRNAs with two alternative 5'ss (5,116 alternative 5'ss events
195 were identified and used in the following analysis, [Supplementary Data 4](#)). We then classified
196 these two alternatives 5'ss as being either distal or proximal 5'ss, according to their relative
197 positions. Based on the expression levels of the corresponding isoforms, we then identified major
198 5'ss ($\geq 80\%$ of the total abundance of two isoforms) and minor 5'ss ($\leq 20\%$ of the total
199 abundance of two isoforms) (See details in Method). We found that SHAPE reactivities at the -1
200 and -2 positions upstream of 5'ss were significantly higher in the major 5'ss group than those in
201 the minor 5'ss group, regardless of distal or proximal positions ([Fig. 3c,d](#)). Therefore, the two-
202 nucleotide single-stranded RNA structure feature upstream of 5'ss was associated with the
203 selection of alternative 5'ss. We then performed the corresponding assessment for alternative 3'ss
204 events (9,237 alternative 3'ss events were identified and used in the following analysis,
205 [Supplementary Data 5](#)) and found SHAPE reactivity at the branch point was notably higher in the
206 major 3'ss group compared to the minor 3'ss group, regardless of distal or proximal positions ([Fig.](#)
207 [3e,f](#)). Thus, single-strandedness at the branch point was associated with the selection of alternative
208 3'ss. High SHAPE reactivity peaks were also observed at other positions around the branch point in
209 both major and minor 3'ss groups, suggesting these high SHAPE reactivities did not contribute to 3'ss
210 selection ([Fig. 3f](#)). Taken together, RNA structure features identified upstream of 5'ss and at the
211 branch point were also strongly associated with the recognition of alternative 5'ss and 3'ss in
212 alternative splicing events.

213 **The two-nucleotide single-stranded RNA structure feature upstream of 5'ss is sufficient to** 214 **regulate splicing**

215 A nucleotide with high GC content tends to be more double-stranded³². Thus, the distinctive
216 single-strandedness at the -1 nucleotide upstream of 5'ss, as a conserved G, is unexpected. In
217 addition, the -1 and -2 nucleotide positions lie within the nine nucleotide binding region of U1
218 snRNA (from -3 to +6 nt region of 5'ss) during splicing³³. If this splicing associated RNA structure
219 feature we observed, affected U1 snRNA binding, then a similar RNA structure feature should
220 have been observed across the whole binding site. However, high SHAPE reactivities were only
221 observed for two out of nine nucleotides rather than the whole binding site. Consequently, we
222 tested whether these two single-stranded nucleotides upstream of 5'ss were sufficient to regulate
223 splicing. We selected the first exon-intron-exon region of *AT5G56870* successfully spliced as a
224 representative example of the pre-mRNAs comprising this distinctive two-single-stranded RNA
225 structure feature upstream of 5'ss ([Fig. 4a](#)). We then made use of it for our functional validation.
226 To avoid disrupting base-pairing between 5'ss and U1 snRNA during splicing, we maintained the
227 U1 snRNA binding site sequence content and inserted a short sequence immediately upstream of

228 this U1 binding site to form a stable hairpin structure with the whole U1 binding site completely
229 base-paired (illustrated in Fig. 4b). Then, we introduced a series of mutations in the inserted
230 sequence that base-pair with the U1 binding site in order to disrupt the base-pairing status of
231 different nucleotides within this binding region (Fig. 4b). We assessed the splicing events on these
232 designed constructs through transient expression assays in *Nicotiana benthamiana* (Fig. 4c). First,
233 we confirmed that the native sequence construct was successfully spliced in tobacco leaves (Fig.
234 4c, lane 1). Splicing was completely inhibited when the whole U1 snRNA binding site was
235 completely base-paired with the inserted sequence upstream (Fig. 4c, lane 2). By introducing a
236 mutation “AA” to allow base-pairing disruption at -1 and -2 positions upstream of 5’ss, we found
237 splicing was rescued (Fig. 4c, lane 3). To avoid potential effects due to changing the sequence
238 content, we also mutated these two nucleotides to “GG” that also disrupted the base-pairing status
239 at -1 and -2 positions and found splicing was also rescued (Fig. 4c, lane 4). Furthermore, we
240 assessed the other mutations designed to disrupt other base-pairing sites across the whole U1
241 binding site (Fig. 4c, lane 5-12). Remarkably, structure disruptions of all other base-pairing sites,
242 even a three-nucleotide mutation, were not able to rescue splicing (Fig. 4c, lane 5-12). Hence, our
243 results indicated that only the two-nucleotide single-stranded RNA structure feature at -1 and -2
244 positions upstream of 5’ss was sufficient to regulate splicing.

245 **A unique RNA structure feature on pre-mRNAs is associated with polyadenylation and** 246 **alternative polyadenylation**

247 Another key step of mRNA processing is polyadenylation that starts with endonucleolytic
248 cleavage on pre-mRNAs followed by addition of a poly(A) tail at the cleavage site². Since only
249 the pre-mRNA structure before endonucleolytic cleavage can be used for elucidating the
250 mechanism underpinning polyadenylation, we assessed whether pre-mRNAs before
251 endonucleolytic cleavage were enriched in our nuclear SHAPE-Structure-Seq libraries. We
252 compared the sequencing reads coverage across cleavage sites (poly(A) sites) annotated in a
253 previous study³⁴ with both our nuclear SHAPE-Structure-Seq libraries and cytosolic SHAPE-
254 Structure-Seq libraries. The reads across poly(A) sites were highly enriched in our nuclear
255 SHAPE-Structure-Seq libraries compared to our cytosolic SHAPE-Structure-Seq data
256 (Supplementary Fig. 9). This indicated high enrichment of pre-mRNAs before polyadenylation in
257 our nuclear SHAPE-Structure-Seq libraries (Supplementary Fig. 9).

258 To accurately determine RNA structure features across poly(A) sites, only reads mapped
259 across poly(A) sites and in downstream flanking regions were used to generate SHAPE reactivity
260 profiles (3,077 and 551 poly(A) sites with ≥ 1 RT-stop per nucleotide under *in vivo* and
261 *deproteinized* conditions were used in the analysis, Supplementary Data 6 and 7). We found that
262 average SHAPE reactivities in two regions (from -28 nt to -17 nt upstream of the poly(A) site and
263 from -4 nt to +1 nt across the poly(A) site) were significantly higher compared to flanking regions

264 for both *in vivo* and *deproteinized* conditions (Fig. 5a,b), suggesting these two regions tended to
265 be more single-stranded than flanking regions. To eliminate the effect of nucleotide composition,
266 we identified control sites where nucleotide composition was similar to the sequence content
267 across poly(A) sites, but where polyadenylation did not occur (Supplementary Fig. 10a,b). We
268 found no significant RNA structure features across these control sites, indicating the two single-
269 stranded regions observed across the poly(A) sites above, were specifically associated with
270 polyadenylation (Fig. 5a,b). Furthermore, we assessed whether these two single-stranded regions
271 also appeared in alternative polyadenylation sites. Compared to constitutive poly(A) sites, we
272 found a similar but weaker structure feature across alternative polyadenylation sites
273 (Supplementary Fig. 11a,b, Supplementary Data 8 and 9). Notably, these structure features were
274 different to those identified from a previous RNA structurome study on mature mRNAs²⁴, further
275 indicating structure differences between pre-mRNAs and mature mRNAs. Therefore, this RNA
276 structure feature with two single-stranded regions may also be responsible for alternative
277 polyadenylation.

278 Further investigation of the sequence content in positions -28 nt to -17 nt upstream of poly(A)
279 sites showed that this region had an accumulation of the conventional polyadenylation signal
280 (PAS) motif “AAUAAA” (Supplementary Fig. 12, Supplementary Data 10). We then aligned
281 SHAPE reactivities across this conventional PAS motif “AAUAAA” upstream of poly(A) sites
282 and sorted pre-mRNAs by the distance between PAS and poly(A) sites (Fig. 5c). The
283 corresponding SHAPE reactivities across PAS and poly(A) sites for each pre-mRNA were then
284 plotted as a heatmap (Fig. 5c). We found that SHAPE reactivities were higher at both PAS sites
285 and across poly(A) sites compared to flanking regions (Fig. 5c). Thus, the conventional
286 polyadenylation signal (PAS) motif “AAUAAA” tended to be a single-stranded region.
287 Interestingly, this unique structure feature consistently appeared regardless of the distance
288 between PAS and poly(A) sites (Fig. 5c). Hence, our results suggested that the single-strandedness
289 of both PAS and poly(A) sites may serve as RNA structure signals for polyadenylation.

290 To understand what type of RNA structures could be formed with these two single-stranded
291 regions, we folded sequences across the poly(A) sites with the constraints of SHAPE reactivities
292 by using the *Vienna RNAfold* package³⁵. We then calculated the base-pairing probability (BPP) of
293 each nucleotide³⁵. Consistent with our SHAPE reactivity profiles, we found that the BPPs in these
294 two regions (from -28 nt to -17 nt upstream of the poly(A) site and from -4 nt to +1 nt across the
295 poly(A) site) were significantly lower compared to the flanking regions for both *in vivo* and
296 *deproteinized* conditions, confirming the single-strandedness of these two regions (Fig. 5d,e).
297 Furthermore, we found no obvious BPP features across the control sites, indicating this structure
298 feature was not due to preferential nucleotide composition (Fig. 5d,e). We also generated the
299 heatmap of BPPs across the conventional PAS motif “AAUAAA” and poly(A) sites. We found
300 that the BPPs were much lower at both PAS sites and poly(A) sites compared to flanking regions

301 (Fig. 5f), consistent with SHAPE reactivity profiles (Fig. 5c). In addition, we assessed the detailed
302 RNA structure elements across PAS and poly(A) sites using the *Forgi* utility³⁶. We found that
303 most RNA structures had both PAS and poly(A) sites located in single-stranded loop regions
304 including multiple loop, hairpin loop and internal loop (Fig. 5g). For instance, one type of RNA
305 structure comprised both PAS and poly(A) sites located in multiple loop regions and connected
306 by one hairpin structure (an example is illustrated in Fig. 5h-top). Another type of RNA structure
307 comprised the PAS site located in a multiple loop region with the poly(A) site located in a hairpin
308 loop region (an example is illustrated in Fig. 5h-bottom). Therefore, our results indicated that
309 diverse RNA structures were formed to maintain single-strandedness at both PAS and poly(A)
310 sites.

311 Discussion

312 For the first time, we generated the *in vivo* RNA structure landscape of *Arabidopsis* nuclear RNAs
313 with structure information for all four nucleotides by developing nuclear SHAPE-Structure-Seq.
314 Having achieved high coverage and high accuracy with our nuclear SHAPE-Structure-Seq, we
315 were able to investigate global RNA structure features of nuclear mRNAs and uncover the
316 regulatory role of RNA structure in mRNA processing.

317 Nuclear mRNA structures are globally different from cytosolic mRNA structures

318 Cytosolic mRNAs are the processed products from nuclear mRNAs, thus they share the same
319 sequences. An intriguing question is whether nuclear mRNA structures in these regions are the
320 same as cytosolic mRNA structures? A recent study in mammalian cells showed that nuclear
321 mRNAs were generally more folded than cytosolic mRNAs²¹. We found similar phenomena in
322 our study (Fig. 2a). However, by further dissecting different genic regions, we found that nuclear
323 mRNA structures in exons located in UTR regions tended to be more single-stranded than
324 cytosolic mRNA structures (Fig. 2b). Previous individual RNA structure studies showed that
325 strong RNA structures in 5'UTRs of mature mRNAs are required for recruitment of translation
326 initiation factors³⁷. Also, strong RNA structures in 3'UTRs are critical for mature mRNA
327 stability³⁸. Structure differences between nuclear and cytosolic mRNAs at 5' UTRs and 3' UTRs
328 might be associated with translation initiation and mature mRNA stability. Furthermore, we
329 observed that nuclear mRNAs tended to be more folded in CDS regions compared to cytosolic
330 mRNAs (Fig. 2b). Since ribosomes are known to remodel mature mRNAs by unwinding RNA
331 structures^{19,39}, more single-stranded features in cytosolic mRNA coding regions may be caused
332 by ribosome scanning. In other words, nuclear mRNA structures without interference from
333 ribosomes may remain more folded.

334 From our global to local assessment of RNA structure features, we found that RNA structure
335 features downstream of start codons and at stop codons were significantly different between

336 nuclear and cytosolic mRNAs. A previous *in vitro* study suggested that mature mRNAs might
337 require strong structures downstream of the start codon for increasing the 40S subunit “dwell
338 time”³⁷. Our observation (Fig. 2c) implied that stronger structures downstream of the start codon
339 in cytosolic mRNAs compared to nuclear mRNAs might relate to the ribosome pausing *in vivo*.
340 At stop codons, we found much higher SHAPE reactivities in cytosolic mRNAs (Fig. 2d). This
341 single-stranded structure feature was also observed in a previous RNA structurome study and was
342 suggested to facilitate translation termination⁴⁰. But in nuclear mRNAs, this structure feature was
343 much weaker (Fig. 2d), implying this single-stranded structure feature at stop codons in cytosolic
344 mRNAs might be specific for translation termination. Taken together, these structure feature
345 differences between nuclear and cytosolic mRNAs implied that mRNAs might undergo refolding
346 from the nucleus to the cytosol.

347 In addition to the effects on structure differences from translation, mRNA processing, e.g.
348 polyadenylation and splicing, might also impact the folding status of RNA structures in different
349 cellular compartments. Previous RNA structure profiling of mature mRNAs after polyadenylation
350 in human observed more folded structure features in the region downstream of PAS sites
351 compared to the region upstream of PAS, which were found to facilitate polyadenylation¹⁵.
352 However, we did not observe significant structure differences between these two regions in our
353 nuclear SHAPE-Structure-Seq, suggesting mRNAs might be refolded after polyadenylation (Fig.
354 5a,b). In addition, we found a distinctive single-stranded region across poly(A) sites (Fig. 5a,b),
355 demonstrating that our method had overcome the limitations of previous mature mRNA
356 structurome studies, which lacked structure information across poly(A) sites¹⁵. Furthermore, our
357 previous study on mature mRNAs in *Arabidopsis* revealed that significantly more folded structure
358 features formed upstream of alternative polyadenylation sites compared to flanking regions²⁴.
359 However, we found RNA structure features associated with alternative polyadenylation in the pre-
360 mRNAs before polyadenylation (Supplementary Fig. 10a,b) were different from those observed
361 in mature mRNAs²⁴. Additionally, our previous study on mature RNAs showed a stronger RNA
362 structure feature upstream of 5'ss in unspliced events²⁴. However, we did not observe similar
363 features in our nuclear SHAPE-Structure-Seq (Fig. 3a,b), indicating the RNA structure features
364 related to splicing are also different between pre-mRNAs and mature mRNAs²⁴. Thus, these
365 structure differences before and after mRNA processing implied that mRNAs may adopt different
366 structures for serving distinct biological processes. Many other factors, e.g. diverse protein
367 interactions, RNA modifications and distinct cellular conditions between the nucleus and cytosol,
368 may also contribute to these structure differences, which offers scope for future studies.

369 **Distinctive RNA structure features upstream of 5'ss and at the branch point are associated** 370 **with recognizing 5'ss and 3'ss respectively**

371 Distinct from mammalian splicing where exon skipping is the dominant type of alternative

372 splicing, intron retention is the most common alternative splicing event in plants and can result in
373 significant biological consequences⁴¹. Previous *in vitro* enzymatic RNA structure profiling in
374 *Arabidopsis* nuclear RNAs showed greater structure differences at the exon-intron junctions
375 where the 5' end of introns were much more double-stranded than upstream exons and 3' end of
376 introns were more single-stranded than flanking sequences¹⁴. However, we did not observe these
377 dramatic differences across exons and introns in our nuclear SHAPE-Structure-Seq data, further
378 confirming that *in vivo* RNA structures were different from *in vitro* RNA structures^{19,20}.

379 The recognition of both 5'ss and 3'ss are of great importance during splicing^{1,4}. The
380 consensus sequence motifs for both are so short that a large number of sites with matching
381 sequences are widely spread in the transcriptome⁴. How to distinguish actual splice sites from a
382 large number of false positives has been a primary challenge in elucidating the regulation of
383 splicing⁴. Previous individual studies in human suggested strong RNA structures at U1 and U2
384 snRNA binding sites can prevent the interactions with U1 and U2 snRNA, thus interfering with
385 the recruitment of U1 and U2 snRNPs during splicing⁴²⁻⁴⁴. In our transcriptome-wide study for
386 5'ss, we identified a two-nucleotide single-stranded RNA structure feature immediately upstream
387 of the 5'ss, which was associated with splicing events (Fig. 3a,b). Since the structure feature was
388 located within the U1 snRNA binding region (from -3 to +6 position across the 5'ss)³³, it is likely
389 that the single-strandedness of these two nucleotides promotes the binding of U1 snRNA in 5'ss
390 recognition. For 3'ss, we found the single-strandedness at the branch point was associated with
391 splicing events (Fig. 3a,b). Since U2 snRNA binds across the branch point through base-pairing¹,
392 the single-strandedness at the branch point might promote the binding of U2 snRNA in 3'ss
393 recognition. Alternatively, this single-strandedness might also be a consequence after binding
394 with U2 snRNA since the RNA-RNA base-pairing interaction leaves the branch point as an
395 internal bulge¹. Previous studies in yeast suggested that stem-loop structures between the branch
396 point and 3'ss could promote the recognition of 3'ss^{45,46}. We also found a 4nt low SHAPE
397 reactivity region upstream of AG dinucleotides at the 3'ss, which suggested the formation of a
398 stronger RNA structure between 3'ss and the branch site (Fig. 3a,b). However, this structure
399 feature was not associated with splicing events, and as such, might be linked with subsequent
400 steps after the recognition of 3'ss, such as docking the 3'ss into the reaction center to approach
401 5'ss⁴⁷. Notably, the two-nucleotide single-stranded RNA structure feature upstream of 5'ss and
402 the single-strandedness at the branch point were also strongly associated with the selection of
403 alternative 5'ss and 3'ss, respectively (Fig. 3c,d,e,f). These results further suggested that these
404 two *in vivo* RNA structure features might serve as general rules for determining actual 5'ss and
405 3'ss in splicing. Although we observed global SHAPE reactivity difference between *in vivo* and
406 *deproteinized* libraries (Supplementary Fig. 5), we found very similar structure features across
407 splicing sites under these two conditions (Fig. 3). A previous biophysics study suggested that
408 splicing occurred rapidly once splice sites were recognized⁴⁸. Therefore, our result suggested that

409 the *in vivo* RNA structure features we observed across splicing sites might represent the RNA
410 structures of pre-mRNAs before spliceosome assembly.

411 **The two-nucleotide single-stranded RNA structure feature upstream of 5'ss can regulate** 412 **splicing**

413 Previous studies of individual RNA structure suggested that strong RNA structures formed at 5'ss
414 can inhibit U1 snRNA binding, and subsequently repress splicing^{8,43,44}. However, the strong
415 structures in each case were so different that no general RNA structure features have been
416 identified for regulating splicing. From our nuclear SHAPE-Structure-Seq data, we were able to
417 sensitively determine that a very fine RNA structure feature showing single-strandedness at the -
418 1 and -2 positions upstream of 5'ss was associated with splicing at the transcriptome-wide scale
419 (Fig. 3a,b). Our functional assessment further confirmed that fine-tuning RNA structure by
420 switching the base-pairing status of only these -1 and -2 positions upstream of 5'ss was sufficient
421 to change the fate of splicing (Fig. 4).

422 One possible mechanism is the single-strandedness of the -1 and -2 positions upstream of
423 5'ss promoted splicing by facilitating the binding of U1 snRNA. U1 snRNA base-pairs with a
424 total of nine nucleotides (from -3 to +6 region of 5'ss) across 5'ss³³. Thus, any nucleotides within
425 this nine-nucleotide U1 binding site should have been able to affect splicing. However, we
426 observed that single-strandedness at all other nucleotide positions within the U1 binding site
427 (except for the -1 and -2 positions) were not able to rescue splicing events (Fig. 4b,c). Therefore,
428 our study revealed that the position of this two-nucleotide single-stranded RNA structure feature
429 was also important for regulating splicing. This phenomenon raised the possibility that the -1 and
430 -2 nucleotides upstream of 5'ss may be the first positions for the interaction with U1 snRNA.
431 Further biophysics studies might be able to assess this hypothesis. Furthermore, once the 5'ss is
432 recognized by base-pairing with U1 snRNA, the whole spliceosome is assembled onto the intron
433 region and the 5'ss-U1 interaction is replaced by interactions of 5'ss with U5 (from -3 to -1 region
434 of 5'ss) and U6 (from +4 to +8 region of 5'ss) snRNAs⁴⁹. It is possible that the single-strandedness
435 of the -1 and -2 positions may also promote interaction with U5 snRNA. Taken together, both our
436 transcriptome-wide RNA structure profiling and functional assessment indicated that the two-
437 nucleotide single-stranded structure feature at the -1 and -2 positions upstream of 5'ss can serve
438 as a general role in splicing regulation.

439 Since splicing is a fundamental biological process across eukaryotes, the regulatory motif for
440 splicing is likely to be conserved and highly selected during evolution. Previous identification of
441 the most conserved sequence motif required for 5'ss recognition is as short as only a dinucleotide
442 GU at 5'ss⁴. The sequence requirement of only two nucleotides might be minimized during
443 evolution selection. The short sequence length of the conserved nucleotides might provide the
444 plasticity for flanking nucleotides to contribute to other biological functions. Here, we postulate

445 that the very fine RNA structure feature we identified from the transcriptome is likely to have
446 evolved in a similar way as the sequence motif, in terms of the single-strandedness of only two
447 nucleotides being sufficient to regulate splicing. It will be of great interest to extend our study in
448 other species to investigate the generality of this regulatory mechanism.

449 **Two single-stranded regions upstream and across poly(A) sites are associated with both** 450 **polyadenylation and alternative polyadenylation**

451 Similar to the challenge of how to recognize splice sites, the recognition of poly(A) sites does not
452 always rely on sequence content. In particular, no unique sequence motif exists around poly(A)
453 sites in plants^{11,50}. Indeed, only ~10% of *Arabidopsis* genes contain the conventional PAS motif
454 “AAUAAA” upstream of poly(A) sites¹¹. Therefore, how to precisely determine actual poly(A)
455 sites has been a major question for improving our understanding of polyadenylation regulation.
456 A previous enzymatic probing study on *in vitro* nuclear RNAs in *Arabidopsis* had attempted to
457 investigate RNA structure features at poly(A) sites¹⁴. However, no structure features were
458 observed at either polyadenylation or alternative polyadenylation sites¹⁴, which may be due to the
459 low resolution of enzymatic probing or low comparability of single-stranded and double-stranded
460 RNase probing^{6,14}. Here, we identified two single-stranded regions (from -28 nt to -17 nt upstream
461 of the poly(A) sites and from -4 nt to +1 nt across the poly(A) sites) that were associated with
462 both polyadenylation and alternative polyadenylation (Fig. 5a,b,d,e and Supplementary Fig. 11).
463 These RNA structure features did not appear in the regions where the nucleotide composition was
464 similar but polyadenylation did not occur (Fig. 5a,b,d,e). Hence, these close-by two single-
465 stranded RNA structure features may serve as an additional signature for the recognition of poly(A)
466 sites.

467 Interestingly, most conventional PAS motifs “AAUAAA” are located within the region from
468 -28 nt to -17 nt upstream of the poly(A) sites (Supplementary Fig. 12). We did observe the
469 conventional PAS motif “AAUAAA” region was more single-stranded compared to flanking
470 regions (Fig. 5c,f), which suggested that the single-stranded region upstream of the poly(A) site
471 corresponded to the PAS motif site. Since sequence content is insufficient for predicting PAS
472 sites¹¹, the single-stranded region upstream of poly(A) sites could offer another signature for
473 recognizing the unconventional PAS motif. Moreover, the interactions of the PAS sites with
474 CPSF30 and WDR33 proteins are crucial during polyadenylation². Hence, PAS sites might adopt
475 this single-stranded structure feature to facilitate protein binding. Furthermore, the
476 endonucleolytic cleavage at poly(A) sites is catalyzed by CPSF73, which has been suggested to
477 prefer RNA single-strandedness⁵¹. Therefore, the single-stranded region across poly(A) sites
478 might facilitate the interaction between CPSF73 and poly(A) sites.

479 In summary, we generated the *in vivo* nuclear RNA structure landscape in *Arabidopsis*

480 achieving both high resolution and accuracy with our nuclear SHAPE-Structure-Seq method. We
481 revealed both global and local structure differences between nuclear and cytosolic mRNAs. We
482 successfully identified respective pre-mRNA structure features associated with splicing and
483 polyadenylation. Through functional validation we determined an RNA structure feature which
484 can regulate splicing. Our study unveiled a new RNA structure regulatory mechanism for mRNA
485 processing. Also, our work emphasized the importance of dissecting RNA populations from
486 different stages of the mRNA life cycle in order to investigate the relationship between RNA
487 structure and biological functions.

488

489 **Acknowledgements**

490 This research was supported in part by the NBIP Computing infrastructure for Science (CiS)
491 group through the provision of a High-Performance Computing Cluster. We thank Prof. Peter
492 Shaw and Prof. Chun Kit Kwok for their advice on the experimental design. We are also grateful
493 to Prof. Igor Vorechovsky for discussions. This work was supported by the Biotechnology and
494 Biological Sciences Research Council [BB/L025000/1], the Norwich Research Park Science
495 Links Seed Fund and a European Commission Horizon 2020 European Research Council (ERC)
496 Starting Grant [680324].

497 **Author contributions**

498 Y.D. conceived the research and designed the experiments; Q.L., X.Y., Y.Z., and X.C. performed
499 the experiments; Z.L. designed the data analysis and experimental validation; Z.L., M.N., and J.C.
500 performed the data analysis with assistance from Y.D.; Z.L. and Y.D. wrote the manuscript with
501 input from all authors. Z.L. and Q.L. contributed equally to this work.

502

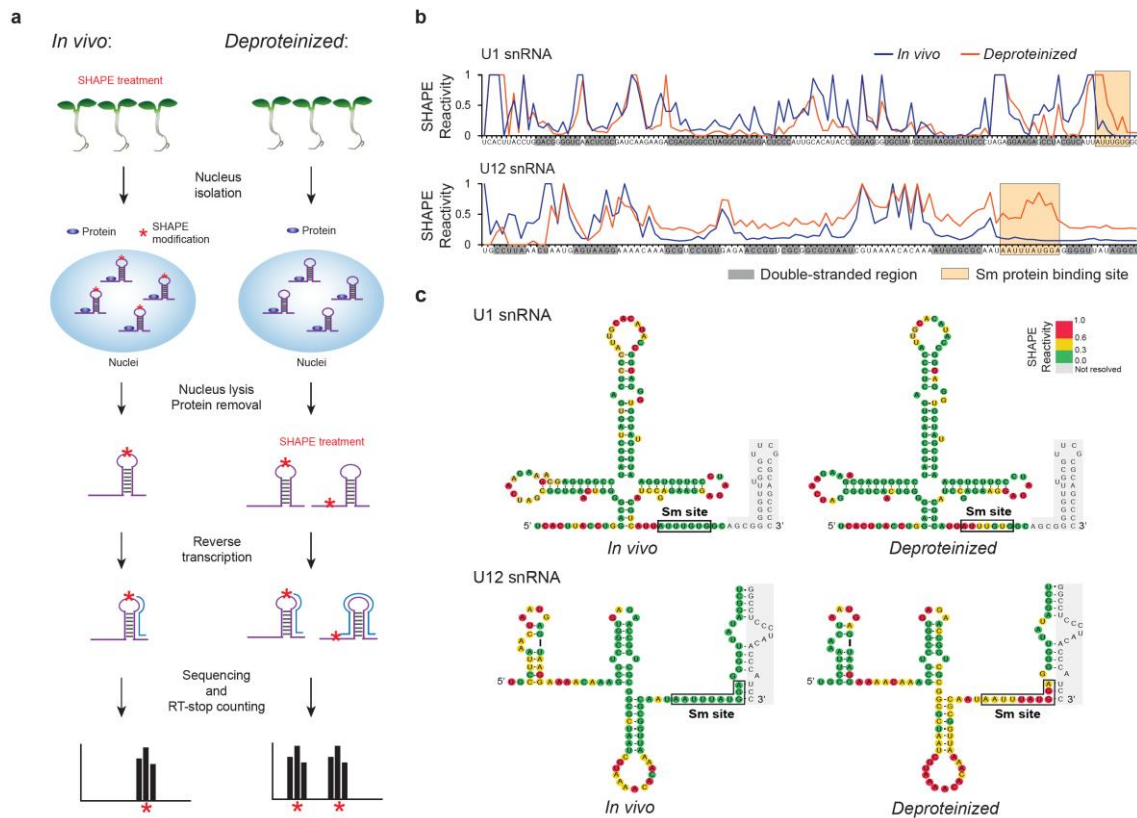
503 **References**

- 504 1 Hoskins, A. A. & Moore, M. J. The spliceosome: a flexible, reversible macromolecular machine.
505 *Trends Biochem Sci* **37**, 179-188, doi:10.1016/j.tibs.2012.02.009 (2012).
- 506 2 Neve, J., Patel, R., Wang, Z. Q., Louey, A. & Furger, A. M. Cleavage and polyadenylation: Ending
507 the message expands gene regulation. *RNA Biol.* **14**, 865-890,
508 doi:10.1080/15476286.2017.1306171 (2017).
- 509 3 Legendre, M. & Gautheret, D. Sequence determinants in human polyadenylation site selection.
510 *BMC Genomics* **4**, 7 (2003).
- 511 4 Roca, X., Krainer, A. R. & Eperon, I. C. Pick one, but be quick: 5' splice sites and the problems of
512 too many choices. *Genes Dev* **27**, 129-144, doi:10.1101/gad.209759.112 (2013).
- 513 5 Tian, B. & Manley, J. L. Alternative polyadenylation of mRNA precursors. *Nature reviews.*
514 *Molecular cell biology* **18**, 18-30, doi:10.1038/nrm.2016.116 (2017).
- 515 6 Bevilacqua, P. C., Ritchey, L. E., Su, Z. & Assmann, S. M. in *Genome-Wide Analysis of RNA*
516 *Secondary Structure* Vol. 50 *Annual Review of Genetics* 235-266 (Annual Reviews, 2016).

- 517 7 Rubtsov, P. M. Role of pre-mRNA secondary structures in the regulation of alternative splicing.
518 *Mol. Biol.* **50**, 823-830, doi:10.1134/s0026893316060170 (2016).
- 519 8 Warf, M. B. & Berglund, J. A. Role of RNA structure in regulating pre-mRNA splicing. *Trends*
520 *Biochem Sci* **35**, 169-178, doi:10.1016/j.tibs.2009.10.004 (2010).
- 521 9 Khaladkar, M., Smyda, M. & Hannehalli, S. Epigenomic and RNA structural correlates of
522 polyadenylation. *RNA Biol.* **8**, 529-537, doi:10.4161/rna.8.3.15194 (2014).
- 523 10 Darmon, S. K. & Lutz, C. S. Novel upstream and downstream sequence elements contribute to
524 polyadenylation efficiency. *RNA Biol.* **9**, 1255-1265, doi:10.4161/rna.21957 (2012).
- 525 11 Loke, J. C. *et al.* Compilation of mRNA polyadenylation signals in Arabidopsis revealed a new signal
526 element and potential secondary structures. *Plant physiology* **138**, 1457-1468,
527 doi:10.1104/pp.105.060541 (2005).
- 528 12 Phillips, C., Kyriakopoulou, C. B. & Virtanen, A. Identification of a stem-loop structure important
529 for polyadenylation at the murine IgM secretory poly(A) site. *Nucleic acids research* **27**, 429-438
530 (1999).
- 531 13 Klasens, B. I. F., Thiesen, M., Virtanen, A. & Berkhout, B. The ability of the HIV-1 AAUAAA signal
532 to bind polyadenylation factors is controlled by local RNA structure. *Nucleic acids research* **27**,
533 446-454, doi:10.1093/Nar/27.2.446 (1999).
- 534 14 Gosai, S. J. *et al.* Global analysis of the RNA-protein interaction and RNA secondary structure
535 landscapes of the Arabidopsis nucleus. *Molecular cell* **57**, 376-388,
536 doi:10.1016/j.molcel.2014.12.004 (2015).
- 537 15 Wu, X. & Bartel, D. P. Widespread Influence of 3'-End Structures on Mammalian mRNA Processing
538 and Stability. *Cell* **169**, 905-917 e911, doi:10.1016/j.cell.2017.04.036 (2017).
- 539 16 Yang, X., Yang, M., Deng, H. & Ding, Y. New Era of Studying RNA Secondary Structure and Its
540 Influence on Gene Regulation in Plants. *Frontiers in plant science* **9**, 671,
541 doi:10.3389/fpls.2018.00671 (2018).
- 542 17 Kwok, C. K. Dawn of the in vivo RNA structureome and interactome. *Biochemical Society*
543 *transactions* **44**, 1395-1410, doi:10.1042/BST20160075 (2016).
- 544 18 Bevilacqua, P. C. & Assmann, S. M. Technique Development for Probing RNA Structure In Vivo and
545 Genome-Wide. *Cold Spring Harbor perspectives in biology* **10**, doi:10.1101/cshperspect.a032250
546 (2018).
- 547 19 Mustoe, A. M. *et al.* Pervasive Regulatory Functions of mRNA Structure Revealed by High-
548 Resolution SHAPE Probing. *Cell* **173**, 181-195 e118, doi:10.1016/j.cell.2018.02.034 (2018).
- 549 20 Rouskin, S., Zubradt, M., Washietl, S., Kellis, M. & Weissman, J. S. Genome-wide probing of RNA
550 structure reveals active unfolding of mRNA structures in vivo. *Nature* **505**, 701-705,
551 doi:10.1038/nature12894 (2014).
- 552 21 Sun, L. *et al.* RNA structure maps across mammalian cellular compartments. *Nature structural &*
553 *molecular biology*, doi:10.1038/s41594-019-0200-7 (2019).
- 554 22 Kwok, C. K., Ding, Y., Tang, Y., Assmann, S. M. & Bevilacqua, P. C. Determination of in vivo RNA
555 structure in low-abundance transcripts. *Nat Commun* **4**, 2971, doi:10.1038/ncomms3971 (2013).
- 556 23 McKeown, P., Pendle, A. F. & Shaw, P. J. Preparation of Arabidopsis nuclei and nucleoli. *Methods*
557 *in molecular biology* **463**, 67-75, doi:10.1007/978-1-59745-406-3_5 (2008).
- 558 24 Ding, Y. *et al.* In vivo genome-wide profiling of RNA secondary structure reveals novel regulatory
559 features. *Nature* **505**, 696-700, doi:10.1038/nature12756 (2014).
- 560 25 Kwok, C. K., Marsico, G., Sahakyan, A. B., Chambers, V. S. & Balasubramanian, S. rG4-seq reveals
561 widespread formation of G-quadruplex structures in the human transcriptome. *Nature methods*

- 562 **13**, 841-844, doi:10.1038/nmeth.3965 (2016).
- 563 26 Zhang, R. *et al.* A high quality Arabidopsis transcriptome for accurate transcript-level analysis of
564 alternative splicing. *Nucleic acids research* **45**, 5061-5073, doi:10.1093/nar/gkx267 (2017).
- 565 27 Pomeranz Krummel, D. A., Oubridge, C., Leung, A. K., Li, J. & Nagai, K. Crystal structure of human
566 spliceosomal U1 snRNP at 5.5 Å resolution. *Nature* **458**, 475-480, doi:10.1038/nature07851
567 (2009).
- 568 28 Wan, Y. *et al.* Landscape and variation of RNA secondary structure across the human
569 transcriptome. *Nature* **505**, 706-709, doi:10.1038/nature12946 (2014).
- 570 29 Spitale, R. C. *et al.* Structural imprints in vivo decode RNA regulatory mechanisms. *Nature* **519**,
571 486-490, doi:10.1038/nature14263 (2015).
- 572 30 Kertesz, M. *et al.* Genome-wide measurement of RNA secondary structure in yeast. *Nature* **467**,
573 103-107, doi:10.1038/nature09322 (2010).
- 574 31 Corvelo, A., Hallegger, M., Smith, C. W. & Eyras, E. Genome-wide association between branch
575 point properties and alternative splicing. *Plos Comput Biol* **6**, e1001016,
576 doi:10.1371/journal.pcbi.1001016 (2010).
- 577 32 Deng, H. *et al.* Rice In Vivo RNA Structurome Reveals RNA Secondary Structure Conservation and
578 Divergence in Plants. *Molecular plant* **11**, 607-622, doi:10.1016/j.molp.2018.01.008 (2018).
- 579 33 Kondo, Y., Oubridge, C., van Roon, A. M. & Nagai, K. Crystal structure of human U1 snRNP, a small
580 nuclear ribonucleoprotein particle, reveals the mechanism of 5' splice site recognition. *Elife* **4**,
581 doi:10.7554/eLife.04986 (2015).
- 582 34 Sherstnev, A. *et al.* Direct sequencing of Arabidopsis thaliana RNA reveals patterns of cleavage
583 and polyadenylation. *Nature structural & molecular biology* **19**, 845-852, doi:10.1038/nsmb.2345
584 (2012).
- 585 35 Lorenz, R. *et al.* ViennaRNA Package 2.0. *Algorithms for molecular biology : AMB* **6**, 26,
586 doi:10.1186/1748-7188-6-26 (2011).
- 587 36 Kerpedjiev, P., Honer Zu Siederdisen, C. & Hofacker, I. L. Predicting RNA 3D structure using a
588 coarse-grain helix-centered model. *Rna* **21**, 1110-1121, doi:10.1261/rna.047522.114 (2015).
- 589 37 Hinnebusch, A. G., Ivanov, I. P. & Sonenberg, N. Translational control by 5'-untranslated regions
590 of eukaryotic mRNAs. *Science* **352**, 1413-1416, doi:10.1126/science.aad9868 (2016).
- 591 38 Moqtaderi, Z., Geisberg, J. V. & Struhl, K. Secondary structures involving the poly(A) tail and other
592 3' sequences are major determinants of mRNA isoform stability in yeast. *Microbial cell* **1**, 137-
593 139 (2014).
- 594 39 Xie, P. & Chen, H. Mechanism of ribosome translation through mRNA secondary structures.
595 *International journal of biological sciences* **13**, 712-722, doi:10.7150/ijbs.19508 (2017).
- 596 40 Shabalina, S. A., Ogurtsov, A. Y. & Spiridonov, N. A. A periodic pattern of mRNA secondary
597 structure created by the genetic code. *Nucleic acids research* **34**, 2428-2437,
598 doi:10.1093/nar/gkl287 (2006).
- 599 41 Marquez, Y., Brown, J. W., Simpson, C., Barta, A. & Kalyna, M. Transcriptome survey reveals
600 increased complexity of the alternative splicing landscape in Arabidopsis. *Genome Res* **22**, 1184-
601 1195, doi:10.1101/gr.134106.111 (2012).
- 602 42 Warf, M. B., Diegel, J. V., von Hippel, P. H. & Berglund, J. A. The protein factors MBNL1 and U2AF65
603 bind alternative RNA structures to regulate splicing. *Proc Natl Acad Sci U S A* **106**, 9203-9208,
604 doi:10.1073/pnas.0900342106 (2009).
- 605 43 Blanchette, M. & Chabot, B. A highly stable duplex structure sequesters the 5' splice site region
606 of hnRNP A1 alternative exon 7B. *Rna* **3**, 405-419 (1997).

- 607 44 Singh, N. N., Singh, R. N. & Androphy, E. J. Modulating role of RNA structure in alternative splicing
608 of a critical exon in the spinal muscular atrophy genes. *Nucleic acids research* **35**, 371-389,
609 doi:10.1093/nar/gkl1050 (2007).
- 610 45 Meyer, M., Plass, M., Perez-Valle, J., Eyras, E. & Vilardell, J. Deciphering 3'ss selection in the yeast
611 genome reveals an RNA thermosensor that mediates alternative splicing. *Molecular cell* **43**, 1033-
612 1039, doi:10.1016/j.molcel.2011.07.030 (2011).
- 613 46 Rogic, S. *et al.* Correlation between the secondary structure of pre-mRNA introns and the
614 efficiency of splicing in *Saccharomyces cerevisiae*. *BMC Genomics* **9**, 355, doi:10.1186/1471-
615 2164-9-355 (2008).
- 616 47 Wilkinson, M. E. *et al.* Postcatalytic spliceosome structure reveals mechanism of 3'-splice site
617 selection. *Science* **358**, 1283-1288, doi:10.1126/science.aar3729 (2017).
- 618 48 Oesterreich, F. C. *et al.* Splicing of Nascent RNA Coincides with Intron Exit from RNA Polymerase
619 II. *Cell* **165**, 372-381, doi:10.1016/j.cell.2016.02.045 (2016).
- 620 49 Yan, C., Wan, R., Bai, R., Huang, G. & Shi, Y. Structure of a yeast activated spliceosome at 3.5 Å
621 resolution. *Science* **353**, 904-911, doi:10.1126/science.aag0291 (2016).
- 622 50 Hunt, A. G. RNA regulatory elements and polyadenylation in plants. *Frontiers in plant science* **2**,
623 109, doi:10.3389/fpls.2011.00109 (2011).
- 624 51 Li, S. *et al.* Evidence that the DNA endonuclease ARTEMIS also has intrinsic 5'-exonuclease activity.
625 *The Journal of biological chemistry* **289**, 7825-7834, doi:10.1074/jbc.M113.544874 (2014).
626
627



628

629

630 **Fig. 1. SHAPE-Structure-Seq method can accurately probe the *in vivo* RNA structure of**

631 **nuclear RNAs.**

632 **a**, Schematic pipeline of nuclear SHAPE-Structure-Seq for both *in vivo* and *deproteinized*

633 conditions. Asterisks, SHAPE modification; blue oval, protein; RT, reverse transcription. For *in*

634 *in vivo* treatments (left), NAI was applied to *Arabidopsis thaliana* seedlings directly and single-

635 stranded nucleotides of RNA were modified. SHAPE treatment was also applied on the RNAs

636 after removing protein, which we termed the ‘*deproteinized condition*’ (right). Deep sequencing

637 was performed followed by the RT-stop counting. **b**, SHAPE reactivity profiles of U1 and U12

638 snRNAs. SHAPE reactivity profiles of both *in vivo* (blue) and *deproteinized* (orange) conditions

639 were shown. Double-stranded regions were shaded with grey. Sm protein binding sites were

640 highlighted with yellow boxes. At Sm protein binding sites, significantly higher SHAPE

641 reactivities were observed under the *deproteinized* condition rather than the *in vivo* condition for

642 both U1 and U12 snRNAs (Paired *t*-test, *P*-value= 6.8e-3 and 3.1e-6 for U1 and U12 snRNA

643 respectively). Higher *deproteinized* SHAPE reactivities were also observed at some double-

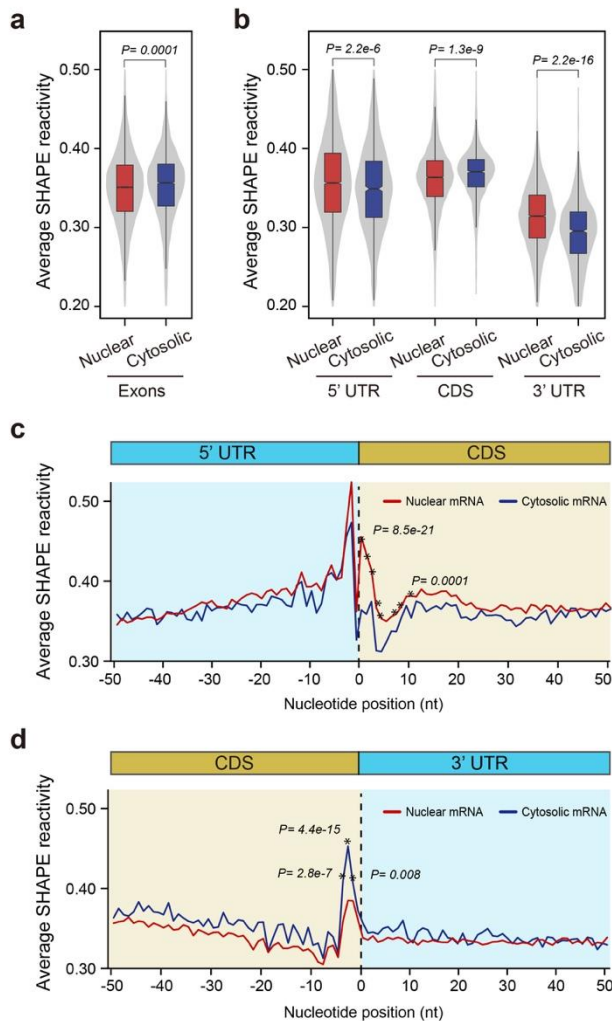
644 stranded regions of U12 snRNA, suggesting the structure of these regions might also be affected

645 by protein interaction. **c**, SHAPE reactivities are consistent with the phylogenetically derived U1

646 and U12 snRNA structures. Sm protein binding sites were highlighted with black boxes.

647 Nucleotides were colour-coded according to *in vivo* and *deproteinized* SHAPE reactivity values

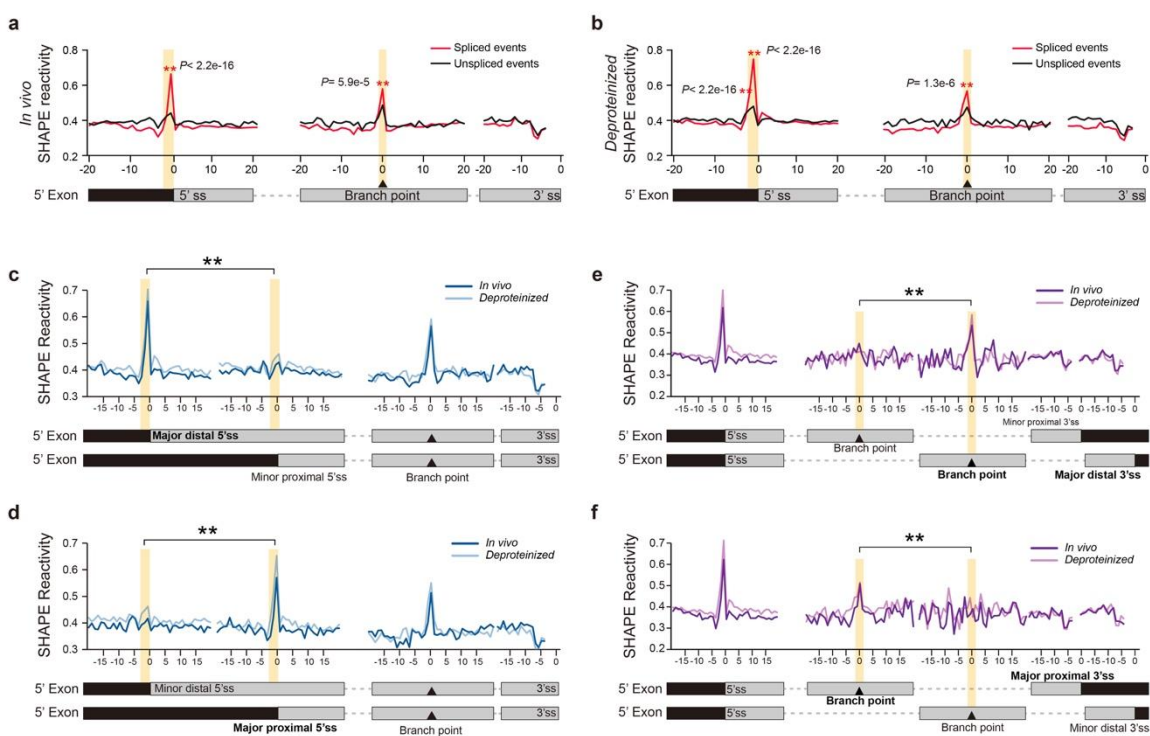
(SHAPE reactivity 0.6-1.0 marked in red, 0.3-0.6 marked in yellow, 0-0.3 marked in green).



648

649 **Fig. 2. *In vivo* nuclear mRNA structures are globally different from cytosolic mRNA**
 650 **structures.**

651 **a**, Comparison of SHAPE reactivities between the exon regions of nuclear and cytosolic mRNAs.
 652 The average SHAPE reactivity of exons in nuclear mRNAs is significantly lower than that in
 653 cytosolic mRNAs (Mann-Whitney test, the *P*-value is shown). **b**, Comparisons of average SHAPE
 654 reactivities between nuclear and cytosolic mRNAs for 5' UTR, CDS and 3'UTR. Average SHAPE
 655 reactivities in both 5'UTR and 3'UTR are significantly higher in nuclear mRNAs than those in
 656 cytosolic mRNAs, whereas average SHAPE reactivities in CDS are significantly lower in nuclear
 657 mRNA than those in cytosolic mRNAs (Mann-Whitney test, the *P*-values were shown). **c**,
 658 Comparison of average SHAPE reactivity profiles between nuclear and cytosolic mRNAs across
 659 the translation start codon. Average SHAPE reactivities downstream of the start codon are
 660 significantly higher in nuclear mRNAs compared to cytosolic mRNAs (Mann-Whitney test, the
 661 highest and lowest *P*-values for the first ten nucleotides of the CDS region are shown). **d**,
 662 Comparison of SHAPE reactivity profiles between nuclear and cytosolic mRNAs across the
 663 translation stop codon. Average SHAPE reactivities at the stop codon are significantly lower in
 664 nuclear mRNAs compared to cytosolic mRNAs (Mann-Whitney test, the *P*-values were shown).



665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

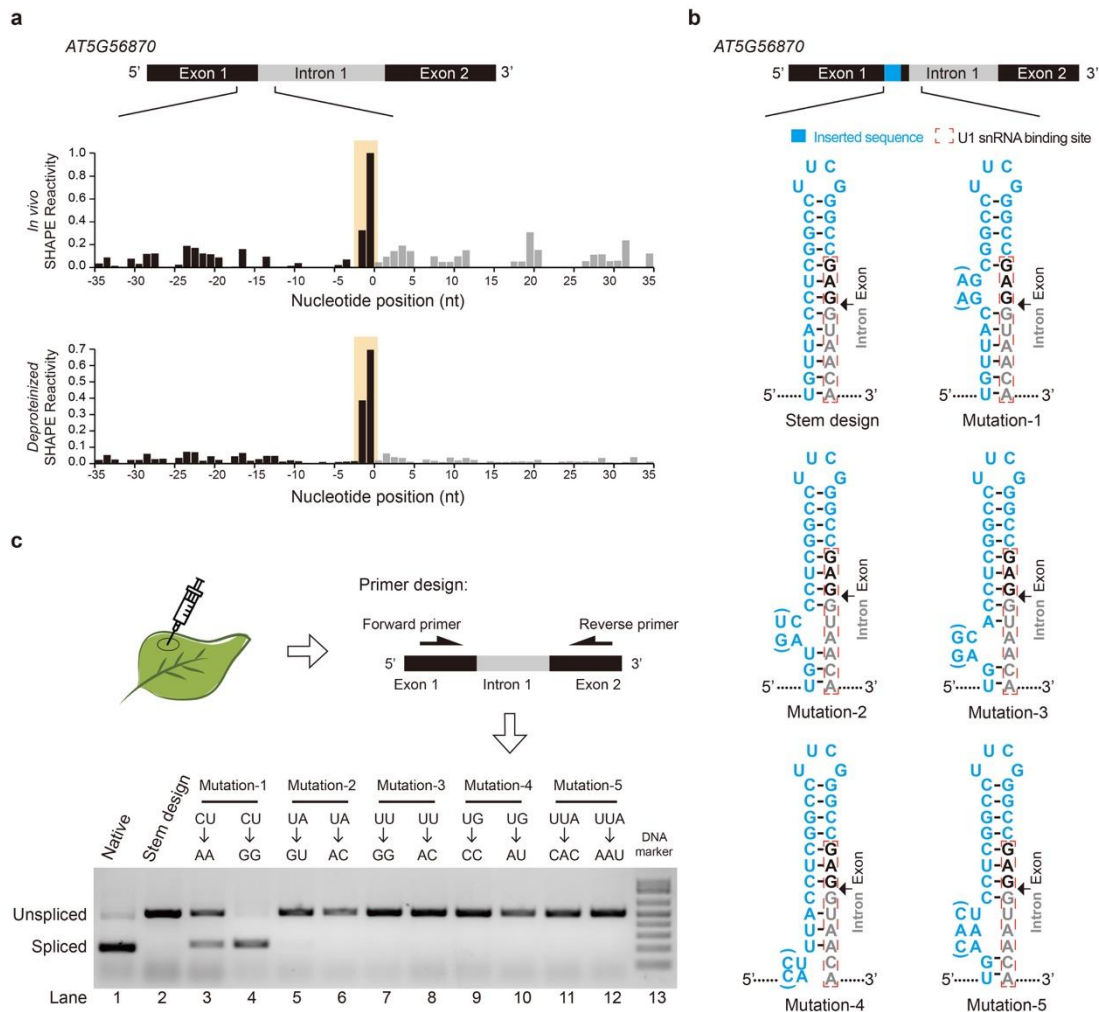
691

692

693

Fig. 3. pre-mRNA secondary structure features upstream of 5'ss and at the branch site are associated with splicing and alternative splice site selection

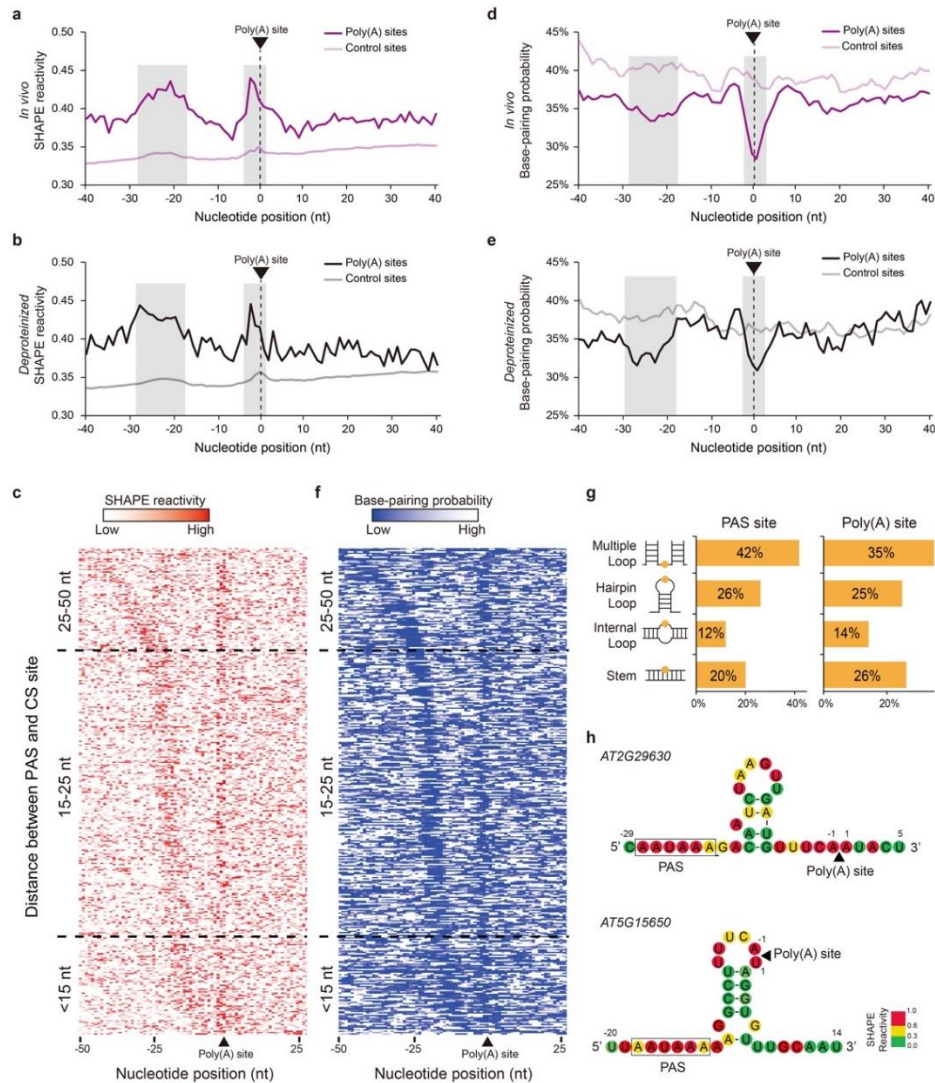
a, b, SHAPE reactivity profiles across 5'ss, branch point and 3'ss for *in vivo* (**a**) and *deproteinized* (**b**) conditions. Average SHAPE reactivity profiles for spliced (red) versus unspliced (black) events are shown. Significantly higher SHAPE reactivities are observed at the -1 and -2 nt positions of 5'ss and the branch point for spliced events rather than unspliced events (Marked with asterisks, Mann-Whitney test, P -values are shown). **c,** SHAPE reactivity profiles for alternative 5'ss events with the distal 5'ss as the major one. Average SHAPE reactivity profiles of both *in vivo* (dark blue) and *deproteinized* (light blue) conditions are shown. Gene models for the two alternative isoforms are shown at the bottom. Significantly higher SHAPE reactivities only appear at -1 and -2 positions upstream of the major distal 5'ss rather than the minor proximal 5'ss (Mann-Whitney test, P -value = 1.6×10^{-4} and $< 2.2 \times 10^{-16}$ at -1 and -2 positions under *in vivo* condition; P -value = 6.1×10^{-9} and $< 2.2 \times 10^{-16}$ at -1 and -2 positions under *deproteinized* condition). **d,** SHAPE reactivity profiles for alternative 5'ss events with the proximal 5'ss as the major one. The significantly higher SHAPE reactivities of -1 and -2 positions only appear upstream of the major proximal 5'ss rather than the minor distal 5'ss (Mann-Whitney test, P -value = 3.3×10^{-12} at -1 position under *in vivo* condition; no significant difference was detected at -2 position under *in vivo* condition; P -value = 3.1×10^{-5} and $< 2.2 \times 10^{-16}$ at -1 and -2 positions under *deproteinized* condition). **e,** SHAPE reactivity profiles for alternative 3'ss events with the distal 3'ss as the major one. Average SHAPE reactivity profiles of both *in vivo* (dark purple) and *deproteinized* (light purple) conditions across different 3'ss and the corresponding branch points are shown. Significantly higher SHAPE reactivity only appears at the branch point of the major distal 3'ss rather than the minor proximal 3'ss (Mann-Whitney test, P -value = 1.2×10^{-3} and 2.8×10^{-4} at branch point under *in vivo* and *deproteinized* conditions respectively). **f,** SHAPE reactivity profiles for alternative 3'ss events with the proximal 3'ss as the major one. The significantly higher SHAPE reactivity only appears at the branch point of the major proximal 3'ss rather than the minor distal 3'ss (Mann-Whitney test, P -value = 1.4×10^{-2} and 1.7×10^{-3} at the branch point under *in vivo* and *deproteinized* conditions respectively).



694

695 **Fig. 4. The two-nucleotide single-stranded RNA structure feature at -1 and -2 nt positions**
 696 **upstream of 5'ss can regulate splicing.**

697 **a**, SHAPE reactivity profiles across 5'ss of the first intron of AT5G56870. High SHAPE
 698 reactivities are observed at -1 and -2 nt positions (shaded in yellow) upstream of 5'ss under both
 699 *in vivo* (top) and *deproteinized* (bottom) conditions, which resemble the global SHAPE reactivity
 700 profiles for spliced events. **b**, Schematic of experimental design to validate the effect of single-
 701 strandedness at the -1 and -2 positions of 5'ss on splicing. A short sequence (blue) was inserted
 702 immediately upstream of the U1 snRNA binding site (red dashed box) to form a stable hairpin
 703 structure with the whole U1 binding site completely base-paired. The exon and intron sequences
 704 are colored in black and grey respectively. A series of mutations were introduced at different
 705 positions of the inserted sequence to disrupt the base-pairing status of different nucleotides within
 706 the U1 binding site. Two types of mutations (with/without bracket) were designed for each
 707 position to avoid potential effects due to changing the sequence content. **c**, Determination of
 708 splicing events by transient expression assay in *Nicotiana benthamiana*. The spliced and
 709 unspliced products were distinguished by semi-qPCR using the same pair of primers located
 710 upstream and downstream of the intron. Spliced and unspliced products are indicated by bands
 711 with different sizes. The construct with native sequence was successfully spliced (lane 1). The
 712 splicing was completely inhibited in the stem design (lane 2). The mutation "AA" or "GG"
 713 disrupted the base-pairing status at -1 and -2 positions upstream of 5'ss (Mutation-1) and rescued
 714 the splicing (lane 3 and 4). All other mutations (Mutation 2-5) designed to disrupt other base-
 715 pairing sites across the U1 binding site did not rescue the splicing (lanes 5-12). Lane 13, the DNA
 716 marker.



717

718 **Fig. 5. Two single-stranded regions on pre-mRNA are associated with polyadenylation.**

719 **a,b**, SHAPE reactivity profiles across poly(A) sites for *in vivo* (**a**) and *deproteinized* (**b**) conditions.

720 The X-axis represents the relative position to the poly(A) site. Average SHAPE reactivities in two

721 regions (from -28 nt to -17 nt upstream of the poly(A) site and from -4 nt to +1 nt position across

722 the poly(A) site) were significantly higher compared to flanking regions for both *in vivo* (purple)

723 and *deproteinized* (black) conditions (Fisher's exact test, P -value = $3.1e-14$ and $3.6e-6$ for *in vivo*;

724 P -values = $4.7e-12$ and $1.4e-3$ for *deproteinized*). The corresponding average SHAPE reactivity

725 profiles for the control sites are in light colours. **c**, Heatmap showing *in vivo* SHAPE reactivity

726 profiles across the PAS motif "AAUAAA" and poly(A) site. The pre-mRNAs are sorted by the

727 distance between PAS and poly(A) site. The gradient colour from light to dark red represents

728 SHAPE reactivity from low to high. The SHAPE reactivities are much higher at both the PAS and

729 poly(A) sites compared to flanking regions. **d, e**, Base-pairing probability (BPP) profiles across

730 poly(A) sites for *in vivo* (**d**) and *deproteinized* (**e**) conditions. Average BPPs in two regions (from

731 -28 nt to -17 nt upstream of the poly(A) site and from -4 nt to +1 nt position across the poly(A)

732 site) were significantly lower compared to flanking regions for both *in vivo* (purple) and

733 *deproteinized* (black) conditions (Fisher's exact test, P -value = $5.0e-12$ and $1.4e-7$ for *in vivo*;

734 P -values = $1.3e-8$ and $2.8e-6$ for *deproteinized*). The corresponding average BPPs for the control

735 sites are in light colours. **f**, Heatmap showing *in vivo* BPPs across the conventional PAS motif

736 "AAUAAA" and poly(A) site. **g**, Classification of RNA structure elements across the PAS and

737 poly(A) sites. The three different single-stranded types (multiple loop, hairpin loop and internal

738 loop) and the double-stranded stem type were assessed for all the PAS and poly(A) sites. The

739 percentage of each type is shown. Most of the PAS and poly(A) sites are located in the single-

740 stranded loop regions including multiple loop, hairpin loop and internal loop. **h**, Illustrations of

741 two individual pre-mRNA structures with both the PAS and poly(A) sites located in single-
742 stranded loop regions. Nucleotides were colour-coded according to the *in vivo* SHAPE reactivity
743 values (SHAPE reactivity 0.6-1.0 marked in red, 0.3-0.6 marked in yellow, 0-0.3 marked in green).