

1 **Intron-driven gene expression in the absence of a core**
2 **promoter in *Arabidopsis thaliana***

3
4

5 Jenna E Gallegos¹ & Alan B Rose^{1,*}

6
7

8 Author affiliation:

9 1. Department of Molecular and Cellular Biology, University of California,
10 Davis, 1 Shields Avenue, Davis, CA, 95616

11 * Corresponding Author

12

13 Keywords: intron, gene expression, transcription initiation, intron-mediated
14 enhancement, promoter

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36 Abstract

37
38 In diverse eukaryotes, certain introns increase mRNA accumulation through the poorly
39 understood mechanism of intron-mediated enhancement (IME). A distinguishing feature
40 of IME is that these introns have no effect from upstream or more than 1 Kb
41 downstream of the transcription start site (TSS). To more precisely define the intron
42 position requirements for IME in Arabidopsis, we tested the effect of the *UBQ10* intron
43 on gene expression from 6 different positions surrounding the TSS of a *TRP1:GUS*
44 fusion. The intron strongly increased expression from all transcribed positions, but had
45 no effect when 204 nt or more upstream of the 5'-most TSS. When the intron was
46 located in the 5' UTR, the TSS unexpectedly changed, resulting in longer transcripts.
47 Remarkably, deleting 303 nt of the core promoter, including all known TSS's and all but
48 18 nt of the 5' UTR, had virtually no effect on the level of gene expression as long as a
49 stimulating intron was included in the gene. When the core promoter was deleted,
50 transcription initiated in normally untranscribed sequences the same distance upstream
51 of the intron as when the promoter was intact. Together, these results suggest that
52 certain introns play unexpectedly large roles in directing transcription initiation and
53 represent a previously unrecognized type of downstream regulatory elements for genes
54 transcribed by RNA polymerase II. This study also demonstrates considerable flexibility
55 in the sequences surrounding the TSS, indicating that the TSS is not determined by
56 promoter sequences alone. These findings are relevant in practical applications where
57 introns are used to increase gene expression and contribute to our general
58 understanding of gene structure and regulation in eukaryotes.
59

60 Introduction

61
62 The transcription start site (TSS) of genes transcribed by RNA Polymerase II in
63 eukaryotes is thought to be primarily determined by the assembly of the pre-initiation
64 complex on recognizable sequences in the core promoter (reviewed in (1, 2, 3)). This
65 process starts with binding of the TATA-binding protein subunit of the general
66 transcription factor TFIID to the TATA box sequences 30-40 nt upstream of the TSS.
67 Core promoters often contain additional recognizable sequences including the initiator,
68 which encompasses the TSS, the TFIIB recognition element, and the downstream
69 promoter element, which, along with the initiator, binds to TFIID (4, 1). However, there
70 are no universally conserved core promoter elements, and even the TATA box is found
71 in a minority of genes in plants (5), yeast (6), and humans (7). The common
72 heterogeneity in start sites within a single gene further suggests that there is
73 considerable flexibility in the sequences that can support transcription initiation (5). In
74 addition to the core promoter, many genes rely on proximal promoter elements (usually
75 less than 1 Kb from the TSS) and distal enhancer elements to regulate gene expression
76 (3).

77
78 There are many examples where a fully intact promoter with all its transcription factor
79 binding sites is inactive unless one or more of the gene's endogenous introns are
80 included (8, 9, 10, 11, 12). The expression of other genes that are weakly active without

81 an intron can be increased tenfold or more by the addition of an intron (8, 13,
82 14, 15, 16). Furthermore, introns can change the spatial expression patterns of tissue-
83 specific genes (12, 17, 18). This demonstrates that some introns significantly boost
84 expression, even in the absence of prior promoter activity.
85
86 In some cases, the mechanism through which an intron increases gene expression is
87 well understood. For example, introns can contain enhancers or alternative promoters
88 (19, 20, 21). In addition, interactions between the splicing machinery and other factors
89 involved in mRNA synthesis and maturation, as well as the exon-junction complex
90 proteins deposited on the mRNA during splicing, assist in mRNA production, stability,
91 export, and translation (22, 23, 24, 25).
92
93 Certain introns that stimulate gene expression exhibit properties that suggest that they
94 must operate by a different and poorly understood mechanism. First, the varying
95 abilities of efficiently spliced introns to increase mRNA accumulation indicate that only
96 certain introns boost mRNA levels beyond the general effects caused by the splicing
97 machinery or exon junction complexes (14). Second, these introns are unlike enhancers
98 because they must be located within 1kb of the TSS to have an effect (26). Third, the
99 sequences responsible for increasing mRNA accumulation are redundant and dispersed
100 throughout these stimulating introns (27), unlike the discrete sequences to which
101 transcription factors bind in promoters and enhancers. The increase in mRNA
102 accumulation caused by specific introns only when in transcribed sequences near the
103 TSS will be referred to here as intron-mediated enhancement (IME) (28).
104
105 An algorithm known as the IMETER assigns a score based on the oligomer composition
106 of a given intron compared to promoter-proximal introns genome-wide (23). The utility of
107 the IMETER in predicting the stimulating ability of introns has been confirmed in
108 *Arabidopsis* (27), rice (29), soybeans (30), and other angiosperms (31). The sequence
109 TTNGATYTG is over-represented in introns with high IMETER scores (27) and is
110 sufficient to convert the previously non-stimulating *COR15a* intron into one that strongly
111 boosts mRNA accumulation. Introns containing either 6 or 11 copies of this motif,
112 named *COR15a6L* and *COR15a11L*, increase mRNA accumulation 15-fold and 24-fold
113 respectively (32).
114
115 While the ability of an intron to stimulate expression is known to vary with its location
116 within a gene, the exact positional requirements for IME have not been fully determined.
117 Early studies demonstrated that most introns that boost expression from the 5' end of a
118 gene have no effect when inserted into the 3' UTR (8, 17, 33). Experiments varying the
119 location of an intron in *TRP1:GUS* coding sequences revealed that their ability to
120 stimulate mRNA accumulation declines when moved from roughly 250 nt to 550 nt
121 downstream of the major TSS, and is completely gone when 1100 nt or more from the
122 major TSS (26). In most previous studies in which the intron was placed upstream of the
123 promoter (8, 34, 14), the distance between the intron and the TSS was more than 550
124 nt, so the introns may have been too far away to affect expression. While these
125 experiments clearly demonstrated that introns are unlike enhancer elements whose
126 influence extends several kilobases, they did not establish if introns must be transcribed

127 to have an effect, or if they must simply be near the TSS to boost expression whether
128 transcribed or not.

129
130 Furthermore, it is unknown if the stimulating ability of an intron continues to increase as
131 it gets closer to the TSS, or if introns have their maximum effect approximately 200 nt
132 downstream of the major TSS where genome-wide average intron IMEter scores peak
133 (35). Pinpointing the ideal location for a stimulating intron will help to maximize gene
134 expression in industrial and research applications. Additionally, the differing effects of a
135 stimulating intron at various locations could yield insight into the mechanism of IME.

136
137 In this study, we completed the first gene-scale mapping of the effect of intron location
138 on IME by comparing the expression of *TRP1:GUS* fusions containing the *UBQ10* intron
139 at different locations around the transcription start site. In the process, we discovered an
140 unexpected role for introns in determining the site of transcription initiation.

141

142 Results

143

144 **The 5'-limit of intron position for IME.** To determine the 5'-limit from which an
145 intron can stimulate expression by IME, the first intron of the *UBQ10* gene was tested at
146 six locations upstream of the normal second exon of a *TRP1:GUS* reporter gene
147 (designated positions 1-6 in Figure 1a and Supplementary Figure 1). Previously
148 generated (26) transgenic *TRP1:GUS* plants with the first intron of *UBQ10* at three
149 additional locations, designated position 0, -1, and -2, were included in figure 1 to show
150 the 3' limits of intron position for IME. Position 0, -1, and -2 had been previously
151 designated "259", "551", and "1136" respectively according to their distance from the
152 most frequently used transcription start site. A different numbering system was used
153 here because a description of intron location relative to the TSS is complicated by the
154 presence of multiple TSSs in this gene. The main transcription start site is at -41 relative
155 to the ATG, as determined by primer extension and RNAse protection (36), RNAseq,
156 and the 5' ends of ESTs and cDNAs in the TAIR database
157 (<https://www.arabidopsis.org/>). There are also a few minor transcription start sites, of
158 which the one at -117 is the furthest from the start codon and is the annotated TSS for
159 this gene.

160

161 Introns inserted at position 3 should be in the 5'-UTR of all *TRP1:GUS* transcripts, while
162 position 4 is upstream of the major TSS but downstream of the 5'-most TSS. The intron
163 at position 5 is upstream of all known TSSs, and roughly midway between
164 the *TRP1* start codon and the stop codon of the upstream gene. This gene (*At5g17980*)
165 is a 3.15 Kb intronless gene of which only 1.8 Kb from the 3' end is present in
166 the *TRP1* promoter fragment in the *TRP1:GUS* fusions. The intron at position 6 is in
167 coding sequences of this gene 198 nt from the stop codon.

168

169 To insert the intron, a *PstI* site was added to the 5' end of the intron, and the last six
170 nucleotides of the intron were converted into a *PstI* site. The intron was then cloned as
171 a *PstI* fragment into a *PstI* site created by site directed mutagenesis at the desired
172 location (Supplemental Figure 1) as described (37). Introns inserted as *PstI* sites leave

173 no extraneous nucleotides in the mRNA (26), and are efficiently spliced, presumably
174 because the 5' splice site is unchanged, and the 3' *Pst*I site (CTGCAG) conforms to the
175 major 3' splice site consensus (NNNYAG) in dicots.

176
177 For intron positions in coding sequences (positions 1 and 2) there were no places where
178 *Pst*I sites could be made without changing an amino acid, so sites were chosen such
179 that changes were unlikely to disrupt function. The *Pst*I sites introduced at positions 1
180 and 2 convert an isoleucine and a glycine respectively into alanines, all of which are in
181 the family of nonpolar amino acids. Further, the first two exons of the *TRP1* gene
182 encode the chloroplast transit peptide (36). Transit peptides have non-stringent
183 sequence requirements and can usually tolerate moderate changes without losing
184 function (38). They are also cleaved off as the protein is imported into the chloroplast,
185 so conservative changes here are unlikely to affect the enzymatic activity of the *GUS*
186 reporter (37).

187
188 To qualitatively assess the positional effect of the *UBQ10* intron on expression of the
189 *TRP1:GUS* reporter, pooled T₂ transgenic seedlings of unknown transgene copy
190 number were histochemically stained for GUS activity with X-gluc (Figure 1b). The
191 *UBQ10* intron clearly stimulated expression from all four transcribed locations but not
192 from either position upstream of the TSS where they may reduce expression. This is
193 consistent with previous observations that introns acting by IME fail to increase
194 expression from upstream of the promoter (8), and suggests that introns may need to
195 be transcribed to affect expression, even if they are within a few hundred nucleotides of
196 the TSS.

197
198 To quantitatively compare expression of reporter genes containing the intron at all
199 transcribed positions, GUS enzyme activity and mRNA accumulation were measured in
200 single-copy transgenic *A. thaliana* lines (Figure 1c, and Supplementary Tables 1 and 2).
201 Expression levels between independent single-copy lines of each construct were similar
202 indicating that, for this transgene, the site of integration had little effect. This has been
203 demonstrated for other *TRP1:GUS* fusions (37, 26), and likely stems from the fact that
204 the *TRP1* promoter is in the middle of the T-DNA insertion and is isolated from flanking
205 plant sequences by at least 2 Kb of sequence on each side.

206
207 *GUS* enzyme assays and RNA gel blots both indicated that *TRP1:GUS* mRNA levels
208 were similar when the *UBQ10* intron was located at position 0, 1, 2, 3, or 4 (Figure 1,
209 and Supplementary Tables 1 and 2). The activity of the intron at position 4 was
210 surprising because this location is 43 nt upstream of the major TSS, meaning that most
211 of the transcripts were not predicted to contain the intron. Inserting 304 nt of sequence
212 between the core promoter and the major TSS was expected to reduce expression.

213
214 **The effect of intron position on TSS location.** To determine if any transcripts
215 initiated within the intron, cDNA from pooled seedlings was amplified by 5'-RACE,
216 cloned, and sequenced. None of the sequenced products contained any intron
217 sequences, but they did reveal an effect on the site of initiation. When the intron was
218 located at positions 1 or 2 (within the first exon), the 5' ends of most 5'-RACE products,

219 and presumably the TSS, mapped to within 10 nt of the major TSS (Figure 2). However,
220 when the intron was located at position 3 (within the 5'-UTR), transcription began almost
221 exclusively within 5 nt of the furthest upstream TSS (Figure 2). When the intron was
222 located at position 4 (upstream of the major TSS), transcription began at the 5'-most
223 TSS or at locations further upstream where initiation does not normally occur, but not at
224 any site downstream of the intron including the major TSS (Figure 2). This suggests that
225 introns may affect the site of transcription initiation and explains how the gene could still
226 be highly expressed when the intron was located upstream of the major TSS.

227
228 **Expression in the absence of the core promoter.** The observation that
229 inserting the intron into the 5'-UTR changed the location at which transcription
230 predominantly began suggests that transcription does not always start a fixed distance
231 downstream of conserved sequences in the core promoter. To determine the extent of
232 the flexibility in sequences that can support initiation, 303 nt of the *TRP1* core promoter
233 were deleted from the *TRP1:GUS* fusion (from position 3 to 5 in Figure 1a). The deletion
234 encompassed all previously recognized transcription start sites and all but 18 nt of the
235 5'-UTR (Supplementary Figure 1). The deletion was created in reporter gene fusions
236 containing no intron, the non-stimulating *COR15a* intron, or one of three stimulating
237 introns: *UBQ10*, *COR15a6L*, or *COR15a11L*. Remarkably, deleting the core promoter
238 had no obvious effect on gene expression, as determined by histochemical staining for
239 *GUS* activity in seedlings of unknown transgene copy number (Figure 3a). To quantify
240 potential subtle differences in expression, mRNA levels were compared in single-copy
241 transgenic *A. thaliana* lines. Deleting the promoter did not diminish mRNA levels
242 from *TRP1:GUS* fusions containing any of the three stimulating introns (Figure 3b).
243 However, deleting the promoter reduced but did not eliminate expression of the
244 constructs containing the non-stimulating *COR15a* intron or no intron (Figure 3b). The
245 size of the mature mRNA was not changed by the promoter deletion regardless of which
246 intron was present.

247
248 **TSS location in the absence of the core promoter.** To determine where
249 transcription was initiating in the deletion-containing constructs, and to verify that
250 transcription of the *TRP1:GUS* fusion was similar to that of the endogenous *TRP1* gene,
251 cDNA from pooled seedlings was amplified by 5'-RACE, cloned, and sequenced. For
252 *TRP1:GUS* constructs containing the *COR15a11L* intron in which the promoter was
253 intact, the apparent transcription start sites mapped within the range of -114 to +1
254 relative to the ATG (Figure 4), consistent with the start sites of the endogenous *TRP1*
255 gene.

256
257 When the promoter was deleted, most of the apparent start sites from constructs
258 containing either the *COR15a11L* or the *UBQ10* intron mapped upstream of the deletion
259 in normally untranscribed sequences (Figure 4). Even though these initiation sites are
260 far upstream of the normal TSS, they are a similar distance upstream of the intron (223-
261 311 nt) as when the promoter was intact. This can be seen in the average length of the
262 5'-UTRs in the sequenced 5'-RACE clones, which did not significantly differ between the
263 promoter-deletion constructs containing either the *Cor15a11L* or *UBQ10* intron, the
264 intact promoter with the *Cor15a11L* intron, or the previously sequenced cDNAs and

265 ESTs from the endogenous *TRP1* gene.

266
267 To verify that deleting the core promoter caused transcription to start in a new location,
268 RNA gel blots were probed with a 709 nt fragment spanning between the location of
269 *PstI* sites at positions 4 and 6 (Figure 3c). The probe hybridized much more strongly
270 with mRNA from lines in which the promoter was deleted and a stimulating intron was
271 present than with the promoter intact regardless of whether or not an intron was
272 present. This confirms that transcripts derived from promoter deletion constructs with
273 stimulating introns contain upstream sequences that are not transcribed when the
274 promoter is intact. The smaller band present in all lanes is of unknown origin.

275
276 **Expression in the absence of the entire intergenic region.** To further test
277 the limits of sequences that can support transcript initiation, a larger deletion (spanning
278 from position 3 to position 6 in Figure 1a) was created in *TRP1:GUS* fusions containing
279 no intron or the *UBQ10* intron. This deletion removed the entire intergenic region from
280 18 nt upstream of the *TRP1* ATG into coding sequences near the 3'-end of the
281 upstream gene, and encompassed the whole region of DNase sensitivity associated
282 with the endogenous *TRP1* promoter (Supplementary Figure 2). Transgenic plants
283 containing *TRP1:GUS* constructs with this complete promoter deletion had undetectably
284 low levels of *GUS* activity even when a stimulating intron was included (Supplementary
285 Figure 3). This suggests that even though the promoter sequences that can support
286 initiation in response to a stimulating intron are surprisingly flexible, there are some
287 features of the *TRP1* promoter, either sequences or chromatin structure, that are
288 absolutely required for expression.

289 290 Discussion

291
292 **Intron position and gene expression.** With the results presented here, the
293 effect on expression of the *UBQ10* intron has been measured from a total of 14
294 locations within the *TRP1:GUS* reporter gene. The intron increased mRNA
295 accumulation from only six positions near the start of the gene. The 5'- and 3'-most
296 positions from which the intron boosted expression are 594 nt apart. Of the six active
297 positions, the effect of the intron on mRNA accumulation was greatest at -18 from the
298 start codon and least from +510, but the effect at positions 0-4 differed from the average
299 by less than 25%. The intron position for maximum IME was closer to the start of
300 transcription than predicted from average IMEter scores, which peak several hundred
301 nucleotides downstream of the TSS. However, this could be due in part to TSS
302 annotation errors in the genomic data used to analyze IMEter score distributions. In
303 practical applications where introns are used to maximize gene expression, for both
304 efficacy and ease it might be best to insert the intron into the 5'-UTR near the start
305 codon.

306
307 It was previously unclear whether introns could affect gene expression if they were
308 upstream of but near the TSS. Here we showed that the *UBQ10* intron, which strongly
309 affects expression from 627 nt downstream of the 5'-most TSS (position 1), clearly did
310 not stimulate expression when located 204 nt upstream of the 5'-most TSS (position 5).

311 While this finding is consistent with the idea that introns must be transcribed to have an
312 effect, this conclusion is complicated by the possibility that stimulating introns might
313 cause initiation upstream of themselves, as discussed below.

314
315 **Intron position and the TSS.** The lack of universal promoter elements (2)
316 illustrates a high degree of flexibility in the sequences that can support transcription
317 initiation. Nonetheless, transcription tends to initiate within the same region for most
318 genes. In genes with many TSSs, it remains unclear whether each TSS lies a fixed
319 distance from multiple functional promoters, or if transcription starts at varying distances
320 from a single promoter because polymerase scans a flexible distance before initiating.
321 Our finding that moving a stimulating intron into the 5'-UTR, or deleting the core
322 promoter, causes transcription to begin in sequences that do not normally support
323 initiation suggests that start sites are not determined by promoter sequence alone.
324 There must be additional conditions that can be influenced by an intron to allow initiation
325 at competent sites that are normally inactive.

326
327 The transcription stimulated by the *UBQ10* intron did not initiate a fixed distance
328 upstream of the intron, as moving the intron through coding sequences towards the start
329 of the gene did not alter the TSS until the intron reached the 5'-UTR. Relocating the
330 intron over a range of nearly 500 nt (at positions +29, +118, +218, and +510 relative to
331 the ATG) in coding sequences of constructs with intact promoters did not appreciably
332 change the TSS, as determined by the size of transcripts on RNA gel blots and
333 sequencing 5'-RACE products (Figures 1, 2, and (26)). When at position 4, the intron
334 separated the promoter from the main TSS, possibly leading to the preferential use of
335 TSSs upstream of the intron at -117 and other locations. However, the predominant
336 TSS also shifted when the intron was at position 3, downstream of the main TSS. In
337 yeast, splicing occurs during transcription very soon after the nascent RNA emerges
338 from RNA polymerase II (41), so moving the intron into the 5'-UTR might push the
339 predominant site of initiation further upstream if there is steric interference between the
340 bulky spliceosome and the transcription machinery as it assembles on the promoter for
341 the next round of transcription.

342
343 **Promoter deletions.** Deletions in the *TRP1* core promoter resulted in comparable
344 levels of transcript initiation 220-310 nt upstream of the intron, as when the promoter
345 was intact. The sequences between positions 5 and 3 are clearly dispensable for
346 expression, even though this region contains all known TSSs for this gene. The larger
347 deletion extending to position 6 in the upstream gene eliminated all expression
348 regardless of the presence of introns. This indicates that the sequences between
349 positions 5 and 6 are capable of supporting initiation, while those upstream of position 6
350 are not.

351
352 The region between positions 5 and 6 may have distinguishing features beyond specific
353 sequence composition, such as a propensity to form open chromatin. The sequences
354 between positions 3 and 6 are noticeably more AT-rich than sequences upstream of
355 position 6, and histone associations are favored by GC-richness (42). Further, the
356 DNase sensitive region that includes the *TRP1* promoter does not extend as far

357 upstream as position 6.

358

359 These results may help to explain why some promoters are completely dependent on
360 introns for expression. A subset of intron-dependent genes may not have core
361 promoters in the traditional sense but rather rely on introns to initiate transcription
362 upstream of themselves. The ability of stimulating introns to override the tissue
363 specificity of promoters further suggests that the mechanism of IME does not depend on
364 prior promoter activity, and inherently leads to constitutive expression throughout the
365 plant.

366

367 If introns can drive expression in the absence of a core promoter, it is puzzling that
368 promoters have not been found more generally dispensable in the large number of
369 publications reporting promoter deletion analyses. One possible explanation is that
370 introns are rarely included in the genes used to assess promoter activity. One exception
371 is a study of the *unc-54* gene of *Caenorhabditis elegans* in which the expression of
372 different versions of the gene was measured by their ability to rescue an *unc-54* null
373 mutation (43). While *unc-54* genomic DNA rescued efficiently, *unc-54* cDNA under the
374 control of the *unc-54* promoter did not. An intron-containing but promoterless gene
375 rescued, and transcripts derived from this construct initiated in the plasmid sequences
376 newly fused upstream of the start codon. Either some bacterial sequences can
377 fortuitously act as a promoter in *C. elegans*, or the first intron of *unc-54* causes initiation
378 upstream of itself.

379

380 There is suggestive evidence from the ENCODE project that some introns might cause
381 initiation upstream of themselves in humans as well. Using unique EST ends as a
382 genome-wide indicator of transcription start sites, initiation was observed not only at the
383 annotated transcription start site but also 250-300 nt upstream of the first intron of
384 genes with long first exons (44). Additionally, genes with a naturally occurring intron
385 very near the transcription start site tend to be more actively transcribed, as
386 demonstrated by association with RNA Polymerase II and TFIID (44). Thus, the
387 apparent ability of introns to stimulate transcript initiation upstream of themselves and
388 for promoter proximal introns to increase gene expression may be conserved across
389 kingdoms.

390

391 **Model.** While the mechanism through which introns increase gene expression remains
392 unclear, the results presented here indicate that introns can have an unexpectedly large
393 influence on determining the location at which transcription initiates. Any model of IME
394 must account for the locations from which introns affect expression, and the ability of
395 some but not all upstream sequences to substitute for the normal transcription start
396 sites of the *TRP1* gene. The following speculation incorporates the results presented
397 here with previous data regarding the effect of introns on gene expression in plants and
398 other eukaryotes. This discussion focuses largely on the *UBQ10* intron because it is the
399 Arabidopsis intron for which the most data are available. From the number of introns
400 with high IMEter scores, we predict the expression of as many as 15% of genes in
401 Arabidopsis may be influenced by a similar mechanism (28). There are likely other ways

402 in which different introns increase expression.

403
404 The main requirement for sequences to act as the promoter of intron-regulated genes
405 such as *TRP1* might be a region of open chromatin that allows access of the
406 transcription machinery to the DNA. Within that region of open chromatin there may be
407 DNA sequences that favor initiation, but others can readily substitute when the preferred
408 sites are deleted. Introns may also help to create the local chromatin state necessary for
409 initiation through the interaction between DNA sequences within the intron and histone
410 modifying or nucleosome positioning factors. Transcription from many or most
411 promoters may be inherently bidirectional (45). In yeast, introns have been shown to
412 affect the proportion of transcripts heading towards coding sequences by recruiting
413 termination factors that decrease transcription in the opposite direction (46). Introns also
414 favor re-initiation by facilitating DNA looping that brings the 3' end of the gene, and thus
415 the transcription machinery, in proximity to the promoter (47). In short, stimulating
416 introns may boost mRNA production by creating local chromatin conditions that favor
417 initiation, decreasing the proportion of PolII molecules that transcribe away from the
418 gene, and positive feedback that amplifies the accumulation of mRNA through re-
419 initiation.

420
421 The limitations to the positions from which the *UBQ10* intron is capable of stimulating
422 mRNA accumulation could be caused by a combination of the relatively short distances
423 (1 Kb or less) over which the postulated intron-driven changes in chromatin structure
424 extend, and the location of potential translational start codons (as described (28)).
425 Ribosomes usually scan from the cap structure at the 5' end of an mRNA and initiate
426 translation at the first ATG, and mRNAs that contain a termination codon too far from
427 their 3' end are rapidly degraded by nonsense-mediated mRNA decay (48, 49, 50).
428 Transcripts that initiate either upstream or downstream of the 220 nt window in which
429 the *TRP1* start codon is the first ATG are therefore likely to be unstable (Supplementary
430 Figure 4). The intron at all locations might increase transcription initiation upstream of
431 itself, but mRNA accumulates only when the intron is near the 5' end of the gene
432 because only then does the *TRP1:GUS* open reading frame occupy almost the entire
433 mRNA. This might also explain why there was actually a slight decrease in expression
434 when the intron was located at position 5 or 6 if transcription was initiating preferentially
435 upstream of the intron. Both the deletion between sites 5 and 3 that permits expression,
436 and the deletion between sites 6 and 3 that eliminates it, leave windows (112 nt and 187
437 nt respectively) in which transcription could initiate to generate functional *TRP1:GUS*
438 protein. The apparent lack of change in TSS when the *UBQ10* intron is moved through
439 *TRP1:GUS* coding sequences is consistent with the idea that only transcripts that start
440 within the 220 nt window upstream of the ATG accumulate, and that within that range
441 the preferred site for initiation is at -41.

442
443 **Conclusion.** The ability of some introns to increase mRNA accumulation from more
444 than 500 nt downstream of the TSS, but not from more than 1100 nt away, and to
445 stimulate expression of a gene lacking its core promoter, indicates that introns represent
446 a previously unrecognized type of regulatory element for genes transcribed by RNA
447 polymerase II. The ability of an intron to effect gene expression over a range of several

448 hundred nucleotides suggests that introns are unlike the downstream transcription
449 factor binding sites found in genes transcribed by RNA Polymerase III. The flexibility in
450 sequences that can support transcript initiation in response to an intron may contribute
451 to the difficulty of identifying conserved promoter elements. Further characterization of
452 the sequences that support initiation in response to an intron, and the role of chromatin
453 structure, should lead to better understanding of the mechanism through which introns
454 increase gene expression.

455

456 **Materials and Methods**

457

458 **Cloning of reporter gene fusions.** The starting template for all constructs was
459 an intronless *TRP1:GUS* fusion containing 1.8 Kb of sequence stretching from the 3'
460 end of the gene upstream of *TRP1* (*At5g17980*) through the first 8 amino acids of the
461 3rd exon of *TRP1* fused to the *E. coli uidA* (*GUS*) gene (39). *PstI* sites were introduced
462 at 6 locations using PCR mutagenesis, and confirmed by sequencing (26). For intron
463 position experiments, the first intron of the *UBQ10* gene was inserted as a *PstI* fragment
464 at each of the 6 locations (26). To delete the core promoter, a promoter fragment whose
465 3' end was the *PstI* site at position 5 (in the intergenic region) was ligated to an exon
466 fragment whose 5' end was the *PstI* site at position 3 in the 5'-UTR, thereby deleting the
467 sequences between positions 5 and 3. The same procedure was used to generate the
468 full promoter deletion, except that the distal fragment had the *PstI* site at position 6 (in
469 the upstream gene). All fusions were then cloned into the binary vector pEND4K,
470 transformed into *Agrobacterium tumefaciens* by electroporation, and then introduced
471 into *Arabidopsis thaliana* ecotype Columbia (Col) by floral dip as described (37).

472

473 **Qualitative *GUS* expression assays.** For qualitative comparisons of expression
474 (as in Figure 1b and Figure 3a), T₂ seedlings for multiple lines for each construct were
475 histochemically stained for *GUS* activity. The plate of seedlings in buffer containing the
476 substrate 5-bromo-4-chloro-3-indolyl β-D-glucuronic acid (Calbiochem, La Jolla, CA,
477 USA) was incubated at 37° for 30 minutes to 2 hours, the plants were rinsed in water,
478 and chlorophyll was removed with ethanol.

479

480 **Identification of single-copy transgenic lines.** Single-copy transgenic lines
481 were identified for quantitative measurements of *GUS* enzyme activity and mRNA
482 levels. For each construct, 18-72 T₂ lines were screened for segregation ratios by
483 kanamycin resistance (26). Genomic DNA was extracted from lines that exhibited a 3:1
484 segregation ratio (kanamycin resistant:sensitive) indicative of a single locus of
485 transgene insertion. The DNA was digested with *BamHI* or *PstI* and probed with
486 the *GUS* gene. All lines for which both enzymes generated a single band were
487 propagated to the T₃ or T₄ generation, and homozygous plants were used for
488 quantitative comparisons. Expression levels were compared to the analogous intronless
489 control pAR281 (37).

490

491 **Mapping transcription start sites with 5'-RACE.** The 5' ends
492 of *TRP1:GUS* mRNAs were mapped by 5'-RACE (as described (40)). Briefly, RNA

493 extracted from transgenic plants was reverse transcribed (Reverse Transcription
494 System; Promega, Madison, WI, USA) using the *GUS*-specific primer OAR37 (5'-
495 TAACGCGCTTTCCACCAACG-3'). The resulting cDNA was polyadenylated with
496 terminal deoxynucleotidyltransferase and used as the template for PCR with the primers
497 OAR37, QT (5'-
498 CCAGTGAGCAGAGTGACGAGGACTCGAGCTCAAGCTTTTTTTTTTTTTTTTTTTT-3') and
499 QO (5'-CCAGTGAGCAGAGTGACG-3'). A dilution of the first round PCR product was
500 used as the template for a second round of PCR with the nested primers QI (5'-
501 GAGGACTCGAGCTCAAGC-3') and OAR29 (5'-GGTTGGGGTTTCTACAGGACG-3').
502 The 2nd round PCR product was cleaned up with a PCR Purification Kit (Qiagen,
503 Hilden, Germany), digested with restriction enzymes *XhoI* and *BamHI*, and ligated into
504 the vector pBluescript KS+. The DNA from transformants was screened by digestion
505 with *SacI*. Clones with inserts long enough to contain the *SacI* site that begins 24 nt
506 downstream of the *TRP1* start codon
507

508 Acknowledgements

509 This research was supported by funding from the UC Davis PI Bridge Program. Jenna E
510 Gallegos was supported by an NSF Graduate Research Fellows Program Grant
511 1148897.

512 References

- 513
- 514 1. Thomas MC, Chiang CM (2006) The general transcription machinery and general
515 cofactors. *Crit Rev Biochem Mol Biol* **41**(3):105–78.
 - 516 2. Danino YM, Even D, Ideses D, Juven-Gershon T (2015) The core promoter: At the
517 heart of gene expression. *Biochim Biophys Acta* **1849**(8):1116–31.
 - 518 3. Vernimmen D, Bickmore WA (2015) The hierarchy of transcriptional activation: from
519 enhancer to promoter. *Trends Genet* **31**(12):696–708.
 - 520 4. Butler JE, Kadonaga JT (2002) The RNA polymerase II core promoter: a key
521 component in the regulation of gene expression. *Genes Dev* **16**(20):2583–92.
 - 522 5. Morton T, et al. (2014) Paired-end analysis of transcription start sites in Arabidopsis
523 reveals plant-specific promoter signatures. *Plant Cell* **26**(7):2746–60.
 - 524 6. Basehoar AD, Zanton SJ, Pugh BF (2004) Identification and distinct regulation of
525 yeast TATA box-containing genes. *Cell* **116**(5):699–709.
 - 526 7. Yang C, Bolotin E, Jiang T, Sladek FM, Martinez E (2007) Prevalence of the initiator
527 over the TATA box in human and yeast genes and identification of DNA motifs enriched
528 in human TATA-less core promoters. *Gene* **389**(1):52–65.
 - 529 8. Callis J, Fromm M, Walbot V (1987) Introns increase gene expression in cultured
530 maize cells. *Genes Dev* **1**(10):1183–200.

- 531 9. Palmiter RD, Sandgren EP, Avarbock MR, Allen DD, Brinster RL (1991)
532 Heterologous introns can enhance expression of transgenes in mice. *Proc Natl Acad*
533 *Sci USA* **88**(2):478–82.
- 534 10. Duncker BP, Davies PL, Walker VK (1997) Introns boost transgene expression in
535 *Drosophila melanogaster*. *Mol Gen Genet* **254**(3):291–6.
- 536 11. Juneau K, Miranda M, Hillenmeyer ME, Nislow C, Davis RW (2006) Introns regulate
537 RNA and protein abundance in yeast. *Genetics* **174**(1):511–8.
- 538 12. Emami S, Arumainayagam D, Korf I, Rose AB (2013) The effects of a stimulating
539 intron on the expression of heterologous genes in *Arabidopsis thaliana*. *Plant Biotechnol*
540 *J* **11**(5):555–63.
- 541 13. Clancy M, Vasil V, Hannah LC, Vasil IK (1994) Maize *Shrunken-1* intron and exon
542 regions increase gene expression in maize protoplasts. *Plant Science* **98**(2):151–161.
- 543 14. Rose AB (2002) Requirements for intron-mediated enhancement of gene
544 expression in *Arabidopsis*. *RNA* **8**(11):1444–53.
- 545 15. Morello L, Gianì S, Troina F, Breviario D (2011) Testing the IMEter on rice introns
546 and other aspects of intron-mediated enhancement of gene expression. *J Exp Bot*
547 **62**(2):533–44.
- 548 16. Kempe K, Rubtsova M, Riewe D, Gils M (2013) The production of male-sterile
549 wheat plants through split barnase expression is promoted by the insertion of introns
550 and flexible peptide linkers. *Transgenic Res* **22**(6):1089–105.
- 551 17. Jeong YM, Mun JH, Lee I, Woo JC, Hong CB, K SG (2006) Distinct roles of the first
552 introns on the expression of *Arabidopsis* profilin gene family members. *Plant Physiol*
553 **140**(1):196–209.
- 554 18. Gianì S, Altana A, Campanoni P, Morello L, Breviario D (2009) In transgenic rice,
555 alpha- and beta-tubulin regulatory sequences control *GUS* amount and distribution
556 through intron mediated enhancement and intron dependent spatial expression.
557 *Transgenic Res* **18**(2):151–62.
- 558 19. Kim MJ, et al. (2010) Seed-expressed casein kinase I acts as a positive regulator of
559 the *SeFAD2* promoter via phosphorylation of the *SebHLH* transcription factor. *Plant Mol*
560 *Biol* **73**(4-5):425–37.
- 561 20. Deyholos MK, Sieburth LE (2000) Separable whorl-specific expression and negative
562 regulation by enhancer elements within the *AGAMOUS* second intron. *Plant Cell*
563 **12**(10):1799–1810.
- 564 21. Morello L, Bardini M, Sala F, Breviario D (2002) A long leader intron of the *Ostub16*
565 rice beta-tubulin gene is required for high-level gene expression and can autonomously
566 promote transcription both in vivo and in vitro. *Plant J* **29**(1):33–44.
- 567 22. Le HH, Gatfield D, Izaurralde E, Moore MJ (2001) The exon-exon junction complex
568 provides a binding platform for factors involved in mRNA export and nonsense-
569 mediated mRNA decay. *EMBO J* **20**(17):4987–97.

- 570 23. Maniatis T, Reed R (2002) An extensive network of coupling among gene
571 expression machines. *Nature* **416**(6880):499–506.
- 572 24. Dahan O, Gingold H, Pilpel Y (2011) Regulatory mechanisms and networks couple
573 the different phases of gene expression. *Trends Genet* **27**(8):316–22.
- 574 25. Wiegand HL, Lu S, Cullen BR (2003) Exon junction complexes mediate the
575 enhancing effect of splicing on mRNA expression. *Proc Natl Acad Sci USA*
576 **100**(20):11327–32.
- 577 26. Rose AB (2004) The effect of intron location on intron-mediated enhancement of
578 gene expression in *Arabidopsis*. *Plant J* **40**(5):744–51.
- 579 27. Rose AB, Elfersi T, Parra G, Korf I (2008) Promoter-proximal introns in *Arabidopsis*
580 *thaliana* are enriched in dispersed signals that elevate gene expression. *Plant Cell*
581 **20**(3):543–51.
- 582 28. Gallegos JE, Rose AB (2015) The enduring mystery of intron-mediated
583 enhancement. *Plant Sci* **237**:8–15.
- 584 29. Morello L, Giani S, Troina F, Breviario D (2010) Testing the IMEter on rice introns
585 and other aspects of intron-mediated enhancement of gene expression. *Journal of*
586 *Experimental Botany* **62**(2):533–544.
- 587 30. Zhang N, McHale LK, Finer JJ (2015) Isolation and characterization of GmScream
588 promoters that regulate highly expressing soybean (*Glycine max* Merr.) genes. *Plant Sci*
589 **241**:189–98.
- 590 31. Aguilar-Hernández V, Guzmán P (2013) Spliceosomal introns in the 5' untranslated
591 region of plant BTL RING-H2 ubiquitin ligases are evolutionary conserved and required
592 for gene expression. *BMC Plant Biol* **13**(1):179.
- 593 32. Rose AB, Carter A, Korf I, Kojima N (2016) Intron sequences that stimulate
594 expression in *Arabidopsis*. *Plant Mol Bio* **92**(3):337–346
- 595 33. Snowden KC, Buchholz WG, Hall TC (1996) Intron position affects expression from
596 the *tpi* promoter in rice. *Plant Mol Biol* **31**(3):689–92.
- 597 34. Jeon JS, et al. (2000) Tissue-preferential expression of a rice alpha-tubulin gene,
598 *OsTubA1*, mediated by the first intron. *Plant Physiol* **123**(3):1005–14.
- 599 35. Parra G, Bradnam K, Rose AB, Korf I (2011) Comparative and functional analysis of
600 intron-mediated enhancement signals reveals conserved features among plants.
601 *Nucleic Acids Res* **39**(13):5328–37.
- 602 36. Rose AB, Casselman AL, Last RL (1992) A phosphoribosylanthranilate transferase
603 gene is defective in blue fluorescent *Arabidopsis thaliana* tryptophan mutants. *Plant*
604 *Physiol* **100**(2):582–92.
- 605 37. Rose AB, Beliakoff JA (2000) Intron-mediated enhancement of gene expression
606 independent of unique intron sequences and splicing. *Plant Physiol* **122**(2):535–42.
- 607 38. Inoue K, Glaser E (2014) Processing and Degradation of Chloroplast Extension
608 Peptides. *Plastid Biology* (Springer Science, New York), pp 305–323.

- 609 39. Rose AB, Last RL (1997) Introns act post-transcriptionally to increase expression of
610 the *Arabidopsis thaliana* tryptophan pathway gene *PAT1*. *The Plant Journal* **11**(3):455–
611 464.
- 612 40. Scottolavino E, Du G, Frohman MA (2007) 5' end cDNA amplification using classic
613 RACE. *Nat Protoc* **1**(6):2555–2562.
- 614 41. Carrillo OF, et al. (2016) Splicing of nascent RNA coincides with intron exit from
615 RNA polymerase II. *Cell* **165**(2):372–81.
- 616 42. Segal E, Widom J (2009) What controls nucleosome positions? *Trends Genet*
617 **25**(8):335–43.
- 618 43. Okkema PG, Harrison SW, Plunger V, Aryana A, Fire A (1993) Sequence
619 requirements for myosin gene expression and regulation in *Caenorhabditis elegans*.
620 *Genetics* **135**(2):385–404.
- 621 44. Bieberstein NI, Carrillo OF, Straube K, Neugebauer KM (2012) First exon length
622 controls active chromatin signatures and transcription. *Cell Rep* **2**(1):62–8.
- 623 45. Bagchi DN, Iyer VR (2016) The determinants of directionality in transcriptional
624 initiation. *Trends Genet* **32**(6):322–33.
- 625 46. Agarwal N, Ansari A (2016) Enhancement of transcription by a splicing-competent
626 intron is dependent on promoter directionality. *PLoS Genet* **12**:e1006047.
- 627 47. Moabbi AM, Agarwal N, El KB, Ansari A (2012) Role for gene looping in intron-
628 mediated enhancement of transcription. *Proc Natl Acad Sci USA* **109**(22):8505–10.
- 629 48. Brogna S, McLeod T, Petric M (2016) The meaning of NMD: translate or perish.
630 *Trends Genet* **32**(7):395–407.
- 631 49. Maquat L (2004) Nonsense-mediated mRNA decay: a comparative analysis of
632 different species. *CG* **5**(3):175–190.
- 633 50. Kertész S, et al. (2006) Both introns and long 3'-UTRs operate as cis-acting
634 elements to trigger nonsense-mediated decay in plants. *Nucleic Acids Res*
635 **34**(21):6147–57.

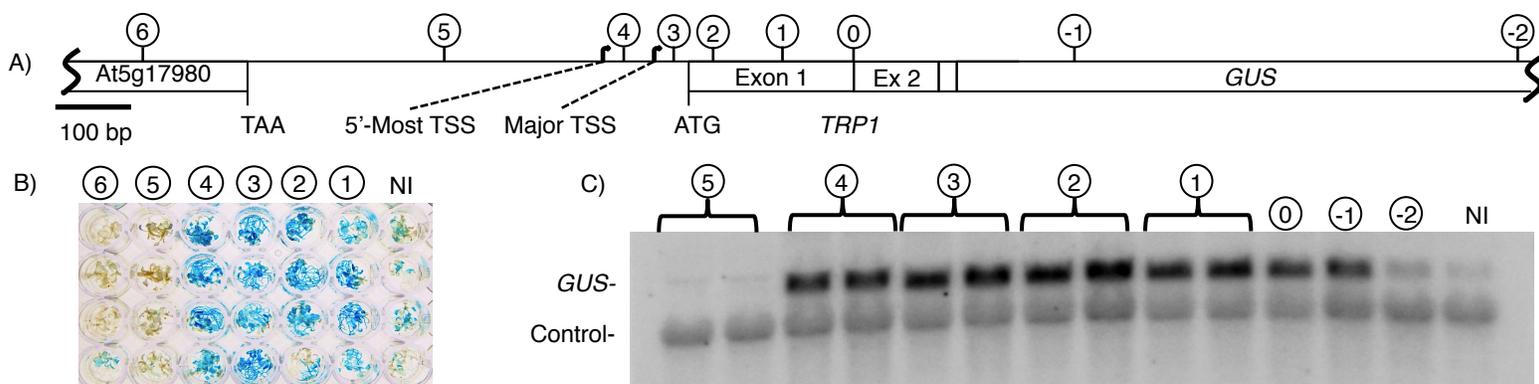


Figure 1. Comparing the stimulating ability of the *UBQ10* intron at locations near the TSS. **A)** The first intron of *UBQ10* was inserted at one of six locations (numbered 1-6) around the TSS. Previously generated constructs with the same intron at position 0, -1, and -2 were also included to show the full limits of intron position for IME. Arrows mark the most commonly used TSS (-41 relative to the start codon), and the 5'-most TSS (-117 relative to the start codon). **B)** Histochemical staining for *GUS* activity in transgenic plants that contain the *TRP1*:*GUS* fusion with the *UBQ10* intron at the indicated position or no intron (NI). Each well in a vertical column contains five T_2 seedlings from an independent line of unknown copy number. **C)** RNA gel blot probed with *GUS* and a loading control (the endogenous *TRP1* gene). Each adjacent lane with the same label represents an independent single-copy homozygous line.

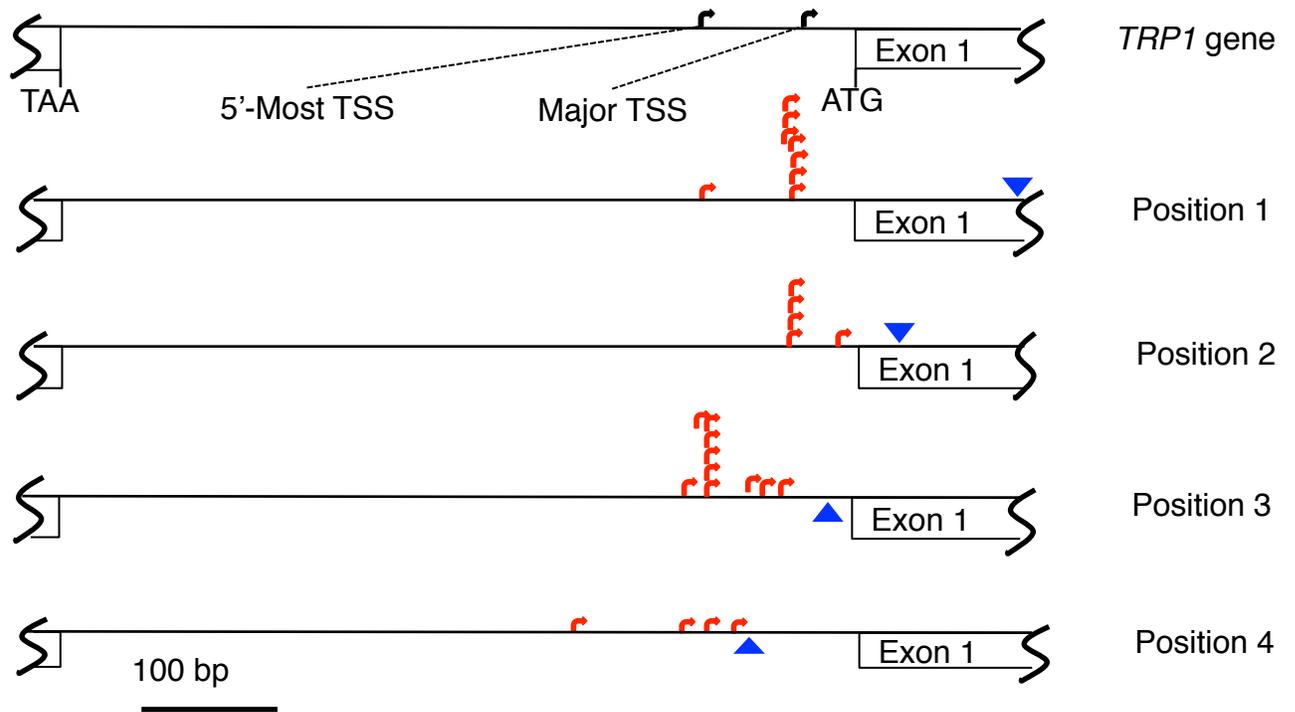


Figure 2. Locations of the 5'-end of *TRP1:GUS* 5'-RACE products when the *UBQ10* intron was at position 1, 2, 3, or 4 (intron position marked with blue triangles). Each red arrow represents an individual sequenced cDNA.

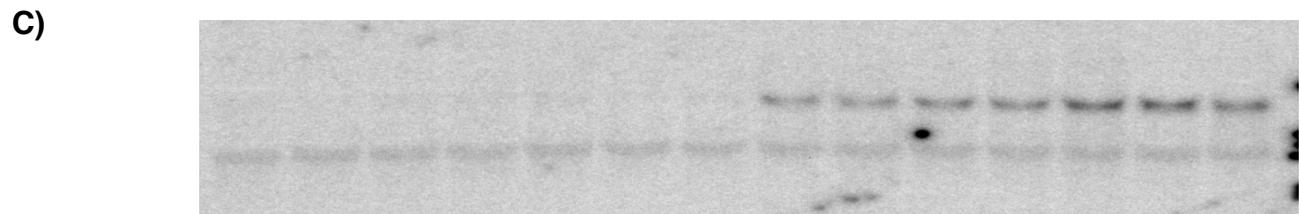
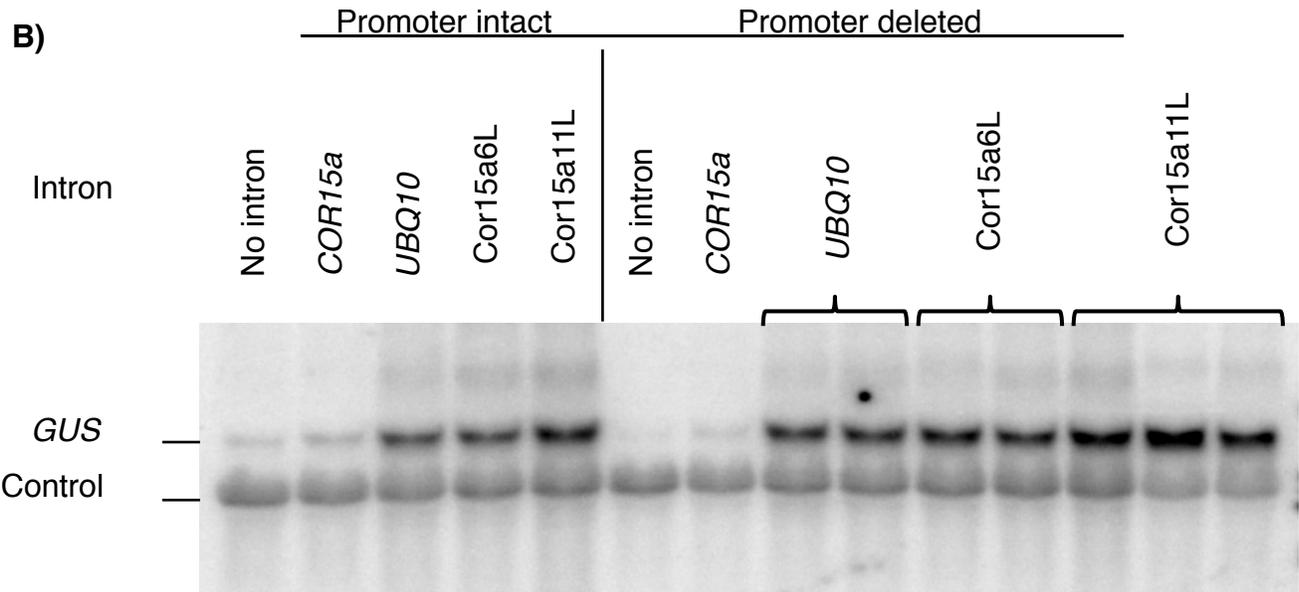
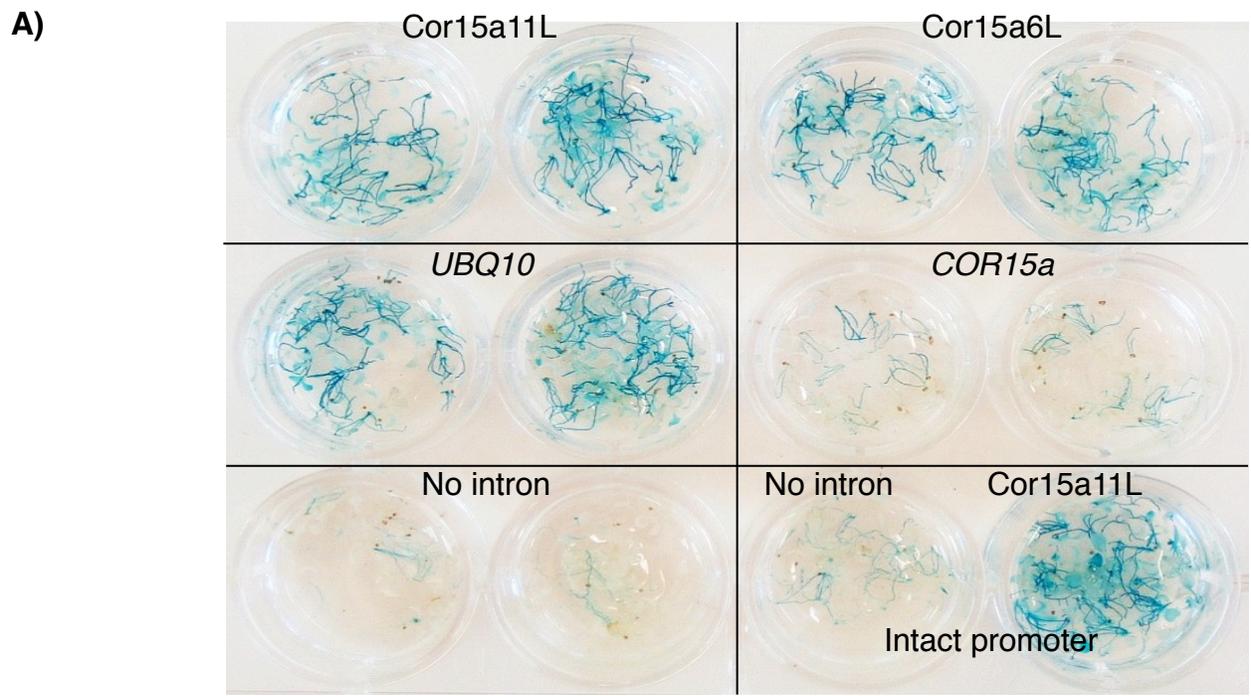


Figure 3. Introns stimulate expression in the absence of the core promoter. **A)** Histochemical staining of transgenic plants containing *TRP1:GUS* fusions with the indicated introns. The promoters in all constructs except the two on the bottom right have had all sequences between positions 5 and 3 deleted. Each circular well contains 20 unrelated T_1 seedlings from an independent transformation. **B)** RNA gel blot probed with *GUS* and a loading control (the endogenous *TRP1* gene). Each lane with the same label contains RNA from an independent single-copy homozygous T_3 line. **C)** Same blot as in **B)** probed with a *TRP1* promoter fragment extending from positions 4 to 6.

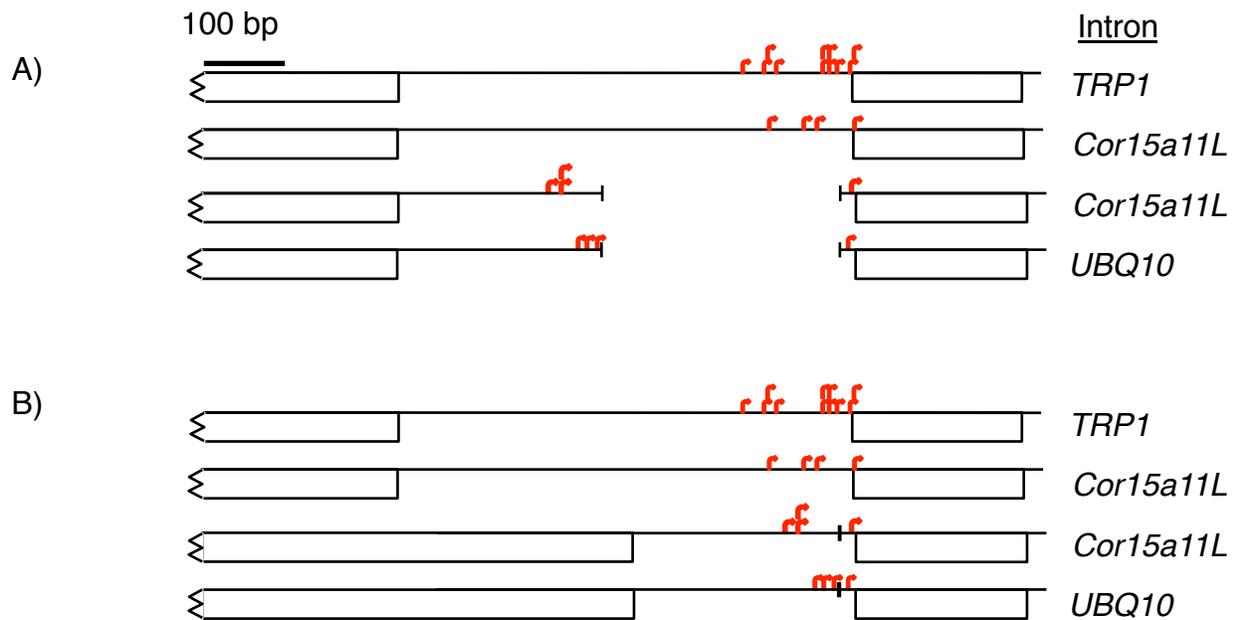


Figure 4. Deleting the core promoter changes the TSS. **A)** Locations of TSSs of the *TRP1* gene (top line) and *TRP1:GUS* fusions containing the indicated introns, with the promoter intact (second line) or deleted (bottom two lines). The sequences are aligned to the genome, so the promoter deletion appears as a gap. The *TRP1* TSSs were determined by primer extension and RNase protection (36) and from cDNAs and ESTs in the TAIR database (<https://www.arabidopsis.org/>). The *TRP1:GUS* TSSs were determined by 5'-RACE, with each red arrow representing an individual sequenced cDNA. **B)** As in **A)**, but with the gap closed to show the distances of the TSS's from the start of the intron.