1 **Title: Identification of mRNAs that undergo stop codon readthrough in**

2 ***Arabidopsis thaliana***

3 Running title: stop codon readthrough in *Arabidopsis thaliana*

4

5 Sarthak Sahoo[1,2,#], Divyoj Singh[1,#], Anumeha Singh[2,#] and Sandeep M Eswarappa[2,*].

6

7 [1]Undergraduate Program, [2] Department of Biochemistry, Indian Institute of Science,

8 Bengaluru, India.

9

10 [#], equal contribution

11

12 *Correspondence:

13 Sandeep M Eswarappa

14 ORCID: 0000-0002-7903-5198

15 Tel: +91-80-22932881

16 Email: sandeep@iisc.ac.in

17

18

19

20

21

22

23

24

**ABSTRACT**

A stop codon ensures termination of translation at a specific position on an mRNA. Sometimes, termination fails as translation machinery recognizes a stop codon as a sense codon. This leads to stop codon readthrough (SCR) resulting in the continuation of translation beyond the stop codon, generating protein isoforms with C-terminal extension. SCR has been observed in viruses, fungi, and multicellular organisms including mammals. However, SCR is largely unexplored in plants. In this study, we have analyzed ribosome profiling datasets to identify mRNAs that undergo SCR in *Arabidopsis thaliana*. Analyses of the ribosome density, ribosome coverage and three-nucleotide periodicity of the ribosome profiling reads, in the mRNA region downstream of the stop codon, provided strong evidence for SCR in mRNAs of 144 genes. This process generates putative peroxisomal targeting signal, nuclear localization signal, prenylation signal, transmembrane helix and intrinsically disordered regions in the C-terminal extension of several of these proteins. Gene ontology (GO) functional enrichment analysis revealed that these 144 genes belong to three major functional groups - translation, photosynthesis and abiotic stress tolerance. Finally, using a luminescence-based assay, we experimentally demonstrate SCR in representative mRNAs belonging to these functional classes. Based on these observations, we propose that SCR plays an important role in plant physiology by regulating the protein localization and function.

**AUTHOR SUMMARY**

Protein synthesis executed by macromolecular complexes, termed ribosomes, starts and stops at specific locations on a messenger RNA (mRNA). This fidelity is critical for the normal functioning of cells. However, sometimes ribosomes don't stop translation at the stop signal (termed stop codon) on an mRNA resulting in longer proteins with properties different from those of the canonical shorter protein. This process called stop codon readthrough (SCR) has been observed in viruses, fungi, and multicellular organisms including mammals. However, it remains largely unexplored in plants. In this study, we report evidence of SCR in 144 genes of *Arabidopsis thaliana*, a small flowering weed widely used as a model system to study plant biology. These genes are involved in protein synthesis, photosynthesis and stress tolerance in plants. We have also experimentally demonstrated SCR in a few genes that represent these functional classes. Our analysis shows that SCR can change the localization and functional properties of these proteins. We propose that SCR plays an important role in plant physiology.

**INTRODUCTION**

A stop codon (UGA/UAA/UAG) on an mRNA signals the translating ribosomes to terminate the process of translation. However, in certain mRNAs, ribosomes fail to terminate at the canonical stop codon and continue translation till the next in-frame stop codon. This is caused by recoding of stop codons by a near-cognate tRNA or a suppressor tRNA. This process of stop codon readthrough (SCR) generates protein isoforms with extended C-terminus, thus contributing to proteome expansion [1]. Because of the extended C-terminus, the protein isoform generated by SCR can be different from the canonical isoform in terms of its localization, function, or stability [2-5]. Since this process occurs at the translational level, SCR enables cells to swiftly respond to environmental cues.

SCR has been observed in bacteria, yeast, insects, mammals, and viruses. It is well-studied in plant viruses [6,7]. For example, tobacco necrosis virus-D (TNV-D) expresses its polymerase, and potato leafroll virus generates a minor capsid protein by SCR [8,9]. It enables viruses to maximize the coding potential of their compact genome. Since plant viruses utilize the translation machinery of the host, it is likely that some plant mRNAs also undergo SCR. However, so far, there is only one report of a plant mRNA undergoing SCR. *Arabidopsis* eRF1-1 mRNA undergoes SCR, which regulates its expression by protecting the mRNA from non-sense mediated decay [10]. A wide range of other translation regulation mechanisms have been observed during plant development, light-dark cycle, viral infections, and environmental stresses [11]. A genome-wide analysis of SCR, which is also a translation regulation mechanism, is lacking to understand its role in plant physiology.

94    Ribosome profiling technique, based on the deep sequencing of ribosome

95    protected RNA fragments, has revolutionized our understanding of the process of

96    translation and its regulation. Because it reveals ribosome-occupied regions on an

97    mRNA, ribosome profiling has the potential to identify novel translation events such as

98    SCR [12]. In this study, we analyzed ribosome profiling datasets and identified 144

99    genes of *Arabidopsis thaliana* as targets of SCR. Further, we experimentally confirmed

100   this process in 4 candidate genes.

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

## RESULTS AND DISCUSSION

### Selection and curation of ribosome profiling datasets

The presence of translating ribosomes after the canonical stop codon of an mRNA strongly indicates SCR [12]. We analyzed ribosome profiling data generated using *A. thaliana,* which are available at Sequence Read Archive (SRA) of National Center for Biotechnology Information (NCBI). We retrieved 14 *A. thaliana* ribosome profiling datasets from SRA and processed them as described in Methods.

Ribosomal footprints obtained from translating ribosomes exhibit frame bias. i.e., they show a fixed distribution of reads across the three frames of the coding sequence. This three-nucleotide periodicity (or phasing) is a sign of translation on the corresponding region of an mRNA (in our case, the proximal 3′UTR). This kind of spatial resolution along mRNAs is required in ribosome profiling datasets to claim unusual translation events such as SCR. Therefore, we first analyzed the three-nucleotide periodicity profile of the ribosome profiling datasets. We chose 9 ribosome profiling datasets based on a clear three-nucleotide periodicity of ribosome profiling (ribo-seq) reads corresponding to the coding sequences of all genes. A representative profile is shown in Fig 1A.

These 9 datasets were derived from seedlings, root, shoot, flower, and from a cell line of *A. thaliana* (Table S1) [13-20]. We then analyzed the length distribution of the reads in each dataset. Length distribution was consistent with footprints of 80S ribosomes, which is a signature of translation. A representative length distribution is shown in Fig 1B. Based on this distribution, we chose the reads of 3 most abundant lengths for further analyses.

140    After removing the reads that map onto non-coding RNAs, we aligned the rest of

141    the reads with *A. thaliana* protein-coding mRNAs. Only those reads that align 100%

142    (i.e., without any mismatch) with an mRNA were considered for the analysis. As our aim

143    was to identify SCR, we focused on the ribosome footprints in the proximal part of the

144    3′UTR - from the canonical stop codon to the downstream in-frame stop codon (Fig S1).

145    This region was termed inter-stop codon region (ISR). mRNAs without downstream in-

146    frame stop codon were not included in the analysis.

147

148    **215 mRNAs from 144 genes of *A. thaliana* show evidence of SCR**

149    We subjected the mRNAs of *A. thaliana* to a stringent four-level screening to identify the

150    targets of SCR (Fig 1C and Fig S2). It is important to distinguish ribosome profiling

151    (ribo-seq) reads due to SCR from reads resulting from non-translating events [21]. This

152    was achieved by comparing the ribosome densities in different regions of an mRNA.

153    The average ribosome density in the 3′UTR (untranslated region) of mRNAs indicates

154    ribosome occupancy due to events not related to translation. mRNAs that showed at

155    least 4-fold higher ribosome density in the ISR compared to the rest of the 3′UTR were

156    considered for further analysis. 1144 mRNAs were identified in this first level of

157    screening (Fig S2 and Table S1).

158    It is possible that the increased ribosome density can be due to a specific

159    segment of the ISR with a strong RNA structure or a strong interaction with a protein (or

160    any *trans*-acting molecule). To exclude such events, we subjected the mRNAs to a

161    coverage-based filtering. We eliminated mRNAs with < 50% coverage in the ISR and >

162    25% coverage in the rest of the 3′UTR. 550 mRNAs satisfied this criterion and all of

163 them had at least one ribo-seq read spanning the canonical stop codon (Fig S2 and

164 Table S1).

165       The three-nucleotide periodicity profile of the ribo-seq reads assigned to ISR,

166 similar to that of the reads assigned to the coding sequence (CDS), provides a strong

167 evidence for SCR. 236 mRNAs satisfied this third level of screening. It is possible that

168 the remaining mRNAs may include ribosomal frameshifting (also known as translational

169 frameshifting) candidates. We did not pursue frameshifting events as the focus of our

170 study was SCR, where the frame of the ribosomes translating the ISR is same as that in

171 the CDS. The ribosome density profile and the three-nucleotide periodicity profile of the

172 ISRs of four representative genes - *RPS15AD, CURT1B, CAM1* and *MUB6* - are shown

173 in Fig 2.

174       The mRNA sequences of *A. thaliana* were retrieved from Ensembl Plants. The

175 annotation of the coding sequence on mRNAs can vary across the databases, providing

176 false-positive evidence for SCR. To rule this out, we performed BLAST analysis for the

177 peptides encoded by the ISRs of 236 mRNAs that passed the screening described

178 above, against NCBI's protein database for *A. thaliana*. We found 21 matches, which

179 were eliminated in this fourth level screening.

180       Thus, 215 mRNAs encoded by 144 genes of *A. thaliana* passed our stringent

181 four-level screening, and they were designated as SCR-positive genes (Fig S2 and

182 Table S2). The average ribosome density in the ISR increased with each screening step

183 (Fig 3A). Also, the average ribosome density in the ISR of SCR-positive mRNAs was

184 10-fold higher compared to the same in all mRNAs. This difference was not observed in

185 the ribosome density of the 3′UTR excluding the ISR (Fig 3B).

186    These results show that our screening methods were able to identify translational

187    event immediately after the stop codon, which constitutes SCR. Though translational

188    frameshifting could also result in increased ribosomal density in the 3′UTR, the three-

189    nucleotide periodicity-based 3rd screen will remove frameshifting events as described

190    above. We have not allowed a single mismatch while assigning reads to different

191    regions of an mRNA. Also, all SCR-positive candidates have at least one ribosome

192    profiling read mapping on to the region spanning the canonical stop codon (the junction

193    of the coding sequence and the ISR). These two conditions rule out RNA editing and

194    polymorphism at the stop codon as reasons for ribosome footprints after the stop codon.

195    In our analyses, we have excluded genes which have ISRs and/or 3′UTRs < 45

196    nucleotides and genes with < 30 reads in their ISR. Also, our method does not reveal

197    candidates that undergo SCR under specific physiological or pathological conditions not

198    included in the ribosome profiling studies. Hence, 144 SCR-positive genes is an

199    underestimate; it is likely that more mRNAs undergo SCR in *A. thaliana*. For example,

200    we did not observe any ribo-seq footprints in the ISR of *eRF1-1*, which has been

201    demonstrated to undergo SCR in *A. thaliana* [10].

202

203    **The stop codon TGA is enriched in SCR-positive genes**

204    Since the identity of the stop codons can influence the efficiency of translation

205    termination [22], we examined the distribution of the three stop codons among the SCR-

206    positive mRNAs at their canonical termination position. We observed a 25% higher

207    occurrence of TGA stop codon in SCR-positive mRNAs compared to the expected

208  frequency (Fig 4A). Interestingly, TGA is the leakiest among the three stop codons,

209  which facilitates the process of SCR.

210       The context of stop codons, especially the nucleotides immediately before (-1)

211  and after (+1) the stop codon, can also influence the efficiency of translation termination

212  [23]. Hence, we examined if there are any conserved sequences around the stop codon

213  in SCR-positive mRNAs. We used WebLogo, a sequence logo generator, to visualize

214  the extent of conservation around the stop codon [24]. Here, the height of the stack

215  indicates the extent of conservation at that particular position. Interestingly, nucleotides

216  just before (-1) and after (+1) and the stop codon showed higher conservation

217  compared to other positions. A and U were more frequently observed in these positions

218  than the other two nucleotides (Fig 4B). These conserved residues are possibly

219  important to provide SCR-permissive context in SCR-positive mRNAs.

220

221  **Gene ontology analysis: Genes involved in translation, photosynthesis, and**

222  **stress response are enriched in SCR-positive genes**

223  To gain some insight into the functional significance of SCR in *A. thaliana*, we

224  performed gene ontology (GO) functional enrichment analysis on 144 SCR-positive

225  genes using PANTHER web server [25]. Among biological processes, we observed that

226  genes involved in translation, photosynthesis, and abiotic stress response were

227  enriched in the list of SCR-positive genes. For instance, 21 genes involved in translation

228  were part of this list. In consistence with this, genes encoding proteins localized in

229  ribosomes and nucleolus (site of ribosome assembly) were enriched. With respect to

230  molecular functions, there was an enrichment of genes encoding components of

231  ribosomes and mRNA-binding proteins. Interestingly, 34 SCR-positive genes encode

232  RNA-binding proteins. Proteins encoded by 20 of them localize in chloroplast and 18 of

233  them are ribosomal proteins. Together, these observations indicate that SCR could play

234  an important role in regulating the process of translation, photosynthesis, and abiotic

235  stress response in *A. thaliana*.

236

237  **SCR can change the localization of the proteins**

238  SCR has been shown to change the localization of the protein product in some cases.

239  For example, the SCR product of mammalian *MTCH2* is localized to the cytoplasm

240  while the canonical isoform is a mitochondrial membrane protein [3]. SCR products of

241  mammalian malate dehydrogenase and lactate dehydrogenase have a peroxisomal

242  targeting sequence (PTS) at the C-terminus, which directs them to peroxisomes.

243  However, the canonical isoforms are found in the cytoplasm or mitochondria  [5,26,27].

244  Since PTS is usually found at the C-terminus of a protein, SCR provides a mechanism

245  to generate peroxisomal protein isoforms. We analyzed the extended C-termini of 144

246  *A. thaliana* proteins potentially generated after SCR for a possible PTS using

247  PredPlantPTS1 tool [28]. We found four of them with a PTS at their extended C-

248  terminus - AT1G09310 (unknown function), *CHS* (encodes chalcone synthase), *AGP15*

249  (encodes an arabinogalactan protein) and *GGH2* (encodes a gamma-glutamyl

250  hydrolase). None of the canonical isoforms of these four gene-products is known to be

251  located in the peroxisomes (Table 1). Thus, SCR can drive the protein isoforms

252  generated by these four genes to peroxisomes and regulate their functions.

253         Nuclear localization signal (NLS) at the C-terminus of the SCR products can

254     drive them to the nucleus. We searched for the presence of NLS in the ISR of the 144

255     SCR products using SeqNLS and NLStradamus [29,30]. Three of them showed strong

256     NLS with a score > 0.7 - *CURT1B, KCS12* and *AT5G56200*. Among these, the

257     canonical protein isoforms of *CURT1B* (The P subunit of Photosystem I) and *KCS12* (3-

258     ketoacyl-CoA synthase) are not localized in the nucleus. Our analysis predicts that their

259     SCR isoforms are localized in the nucleus, possibly with a moonlighting function (Table

260     2).

261         We then searched for a transmembrane helix at the C-terminus of 144 SCR

262     isoforms using TMHMM Server v. 2.0 [31], as this can be a mechanism to change the

263     localization of the protein to a membrane. 23 of the SCR products showed a

264     transmembrane helix at their extended C-terminus (Table S3). The canonical isoforms

265     of 17 of these 23 genes are not known to be membrane proteins. Our analysis suggests

266     that SCR can potentially regulate their function by driving them to the cell membrane.

267     We also searched for endoplasmic reticulum retention signal, KDEL, in the peptides

268     encoded by the ISR of SCR-positive genes. None of them possess this signal.

269         Isoprenylation is a post-translational modification that occurs at a cysteine

270     residue in the C-terminus. This modification can lead to anchoring of the protein to the

271     cell membrane. We analyzed the peptides encoded by the ISR of 144 SCR-positive

272     genes for potential prenylation signal using the PrePS tool [32]. The peptide encoded by

273     the ISR of metallothionein 1C (MT1C; AT1G07610) showed a potential prenylation

274     signal at its C-terminus – RNYQHGLKPRKMGKK<u>CVLC</u>. CVLC is the predicted

275     prenylation signal. Metallothioneins are cysteine-rich proteins required for tolerance to

276     heavy metals, an abiotic stress. Prenylation of the extended isoform of MT1C can

277     potentially alter its localization from the cytosol to membranes regulating its function.

278         In case of mammalian *VEGFA* and *AGO1*, SCR results in a C-terminus with

279     intrinsically disordered region (IDR). This changes the functional properties of their SCR

280     isoforms [2,4]. We analyzed the products of 144 SCR-positive genes for possible IDRs

281     at the C-terminus using the IUPred2A tool [33]. We observed IDR in the C-terminus of

282     the products of 6 SCR-positive genes. They are involved in the organization of

283     cytoskeleton, redox homeostasis, fatty acid synthesis, auxin and hypersensitive

284     response (Table 3). IDR in their C-terminus can potentially alter the functional properties

285     of these six SCR products.

286

287     **Experimental validation of SCR**

288     We performed *in vitro* translation experiments using wheat germ extract to validate SCR

289     in four genes: *RPS15AD* encodes a ribosomal protein; *CURT1B* encodes the P subunit

290     of Photosystem I; *CAM1* encodes calmodulin, which is involved in abiotic stress

291     response [34]. These three genes represent three functional classes that are enriched

292     in SCR-positive genes - translation, photosynthesis, and abiotic stress response. We

293     also selected one more gene, *MUB6* (encodes membrane-anchored ubiquitin-fold

294     protein 6), which does not belong to any of these three classes, but is one of the 144

295     SCR-positive genes. As described above, ribosome footprints were observed after the

296     stop codon in the ISR of these mRNAs. Also, the three-nucleotide periodicity of the

297     ribosomal footprints on the ISR was comparable to that of the coding sequence, but not

298     to that of the 3′UTR (Fig 3).

13

299    Luminescence-based SCR assays were performed as described previously [4].

300    We cloned the cDNAs of these genes upstream of and in-frame with the cDNA of firefly

301    luciferase without its start codon. Luminescence will be observed only if the translation

302    continues across the canonical stop codon of the test cDNA (Schematic in Fig 6). Thus,

303    luminescence in these assays indicates SCR. *In vitro* transcription followed by *in vitro*

304    translation using wheat germ extract revealed significant luminescence activity in

305    mRNAs of all four genes, much above the background level. Constructs without the

306    corresponding ISRs were used to know the background level of luciferase activity. A

307    construct without a stop codon between the test cDNA and the firefly luciferase cDNA

308    was used to measure the efficiency of SCR. This analysis revealed 8%, 50%, 25%, and

309    6.5% SCR in *RPS15AD*, *CURT1B, CAM1,* and *MUB6*, respectively.

310    Overall, our analysis of ribosome profiling datasets provides strong evidence for

311    SCR in mRNAs of 144 genes of *A. thaliana*. Using a similar analysis of ribosome

312    profiling data, mRNAs of 350 *Drosophila* genes and 42 human genes have been

313    predicted to undergo SCR [12].  The advent of ribosome profiling technique has

314    revealed previously unknown (or lesser-known) mechanisms of translational regulation,

315    including SCR. Since this technique is based on experimentally generated ribosome

316    footprints on mRNAs, it is superior to evolutionary conservation-based computational

317    screening methods to detect SCR, which will miss SCR events that have emerged

318    relatively recently during evolution. Furthermore, the nucleotide resolution of ribosome

319    profiling enables us to decipher the frame of translation at the ISRs. The distribution of

320    length of ribosome profiling reads will have a signature of 80S ribosome occupancy.

321    These features are important to distinguish SCR from ribosomal frameshifting and non-

322 translational events (e.g., protein binding and RNA structures) [35]. Thus, ribosome

323 profiling is a powerful tool to identify SCR events at the transcriptome level.

324 It would be remarkable if SCR does change the properties of the proteins in

325 multiple ways as predicted by our analyses – by introducing peroxisomal targeting

326 signal, nuclear localization signal, prenylation signal, transmembrane helices and

327 intrinsically disordered region in the ISR-encoded C-terminal extension. Other

328 mechanisms such as post-translational modification and degradation, which cannot be

329 predicted with high confidence, might also occur at ISR-encoded extensions. This is not

330 very surprising as random peptide sequences have been shown to have functional

331 motifs. For example, 1/5th of randomly generated peptide sequences carry export signal

332 in *Saccharomyces cerevisiae* [36]. Also, in another study involving *S. cerevisiae*, 8 out

333 of 28 randomly generated peptides showed multiple organellar localization signals [37].

334 Therefore, for a gene the chances of acquiring novel functions by SCR are high.

335 As shown in mammalian and viral SCR processes, it is likely that the nucleotide

336 sequence of ISR is responsible for driving the SCR via a *cis*-acting RNA motif or *trans*-

337 acting molecule [2,4,38,39]. Thus, ISR likely possesses a dual function – driving the

338 SCR and altering the properties of the SCR product. It will be interesting to study how

339 natural selection will shape such genomic regions with constraints at both nucleotide

340 (ability to induce SCR) and amino acid level (novel function).

341 The GO analysis suggests that SCR in *A. thaliana* influences three major

342 physiological processes in plants – protein synthesis, photosynthesis and stress

343 tolerance. Our *in vitro* translation experiments performed using a plant-based system

344 show that the efficiency of SCR is much above the basal error rate, suggesting that

345    these are programmed events with physiological consequences. This is consistent with

346    various functional motifs identified in the C-terminal extensions. We anticipate that more

347    studies will follow to characterize individual SCR events in order to understand the

348    mechanism of SCR as well as its physiological significance in plants.

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

**Materials and Methods**

**Curation of *A. thaliana* transcriptome**

*A. thaliana* has 55,398 mRNAs derived from 27,655 protein-coding and 6,563 non-coding genes. (http://plants.ensembl.org/Arabidopsis_thaliana/Info/Annotation/#assembly). Sequences of the mRNAs were downloaded from Ensembl Plants. From this, we created a file containing sequences of rRNAs, tRNAs, snRNAs, snoRNAs, and miRNAs, which were later used to remove ribosomal footprints that aligned to these sequences. Using cDNA sequences (downloaded from the same source), we noted the positions of the start codon, the canonical stop codon and the first in-frame stop codon (if any). mRNAs with inter-stop codon region (ISR) < 45 nucleotides or rest of the 3′UTR < 45 nucleotides were removed. This is because sequences with shorter length will not give enough statistical power to draw any conclusions from the ribosomal density differences between them (i.e., ISR and 3′UTR).  mRNAs whose ISR sequence was matching with > 24 nucleotide sequence of any other coding sequence were removed. This was done because reads cannot be mapped onto an ISR if its sequence matches with a coding sequence. After these filtrations, we were left with 14,732 protein coding mRNAs for our analysis.

**Preprocessing and sequence alignment of ribosome profiling datasets**

Sequence Read Archive (SRA)-formatted ribosome profiling datasets of 9 studies on *A. thaliana* were downloaded from SRA (Table S1). They were converted to FASTQ format files using the prefetch and fastq-dump command of SRAToolkit (https://github.com/ncbi/sra-tools). The adapter sequences were removed (if not

391  removed already) from the datasets using fastp. Additionally, 3 nucleotides from the 5'

392  end of all reads were also trimmed using fastp as these nucleotides were generally

393  found to be of a low-quality score. Reads that aligned to non-coding RNA sequences

394  (rRNA, tRNA, snRNA, snoRNA, and miRNA) were removed using Bowtie2 (version:

395  2.3.4.1). The FASTQ files were then aligned with a list of protein-coding mRNAs to

396  create BAM (Binary alignment map) files.

397

398  **Mapping of ribosome profiling (ribo-seq) reads onto the coding sequence (CDS),**

399  **the ISR and the 3′UTR of mRNAs**

400  We first analyzed the length distribution of the ribo-seq reads in each dataset. Based on

401  this distribution, we chose the reads of 3 most abundant lengths for further analyses.

402  Ribo-seq reads were assigned to different regions of an mRNA (i.e., CDS, ISR and

403  3′UTR) based on the alignment of the beginning of the read to any of these regions (Fig

404  S1). To avoid ambiguity during the assignment of the ribo-seq reads to different regions

405  of an mRNA, we followed these criteria:

406  (i) For coding sequence (CDS): reads that align to the region from 12 nucleotides

407  upstream of the start codon till $22^{nd}$ nucleotide upstream of the canonical stop codon.

408  (ii) For ISR: reads that align to the region from 12 nucleotides upstream of the canonical

409  stop codon till $22^{nd}$ nucleotide upstream of the downstream in-frame stop codon.

410  (iii)  For 3′UTR: reads that align to the region from 12 nucleotides upstream of the

411  downstream in-frame stop codon till the end of the mRNA.

412  Only those reads that showed 100% alignment to an mRNA region were considered

413  (Even a single mismatch was not allowed).

**Identification of potential SCR candidates**

*Selection based on ribo-seq read density*:

The mRNAs with higher density of reads in their ISR than that in the coding sequence were removed from the analysis as this feature is not consistent with SCR. On the contrary, mRNAs with a 4-fold higher density of reads in their ISR than that in the rest of the 3′UTR were included as this is a signature of translational readthrough.

*Selection based on ribo-seq read coverage*:

To increase the stringency of this screening process, we applied three more selection criteria based on the coverage:

(i) at least 30 ribo-seq reads should map onto the ISR of a given mRNA

(ii) > 50% of ISR and < 25% 3′UTR should be covered by ribo-seq reads

(iii) there should be at least one ribo-seq read spanning the canonical stop codon.

*Selection based on three-nucleotide periodicity*:

We looked at a 62-nucleotide window length around the start and the stop codons to ensure that the ribo-seq datasets showed three-nucleotide periodicity. To quantify these frame biases, all genes with at least 200 reads in the coding sequence were considered. Reads were assigned to three frames based on which frame the first nucleotide aligns with on an mRNA sequence. We computed the mean and the standard deviation of the fraction of reads that fell in each frame across all the codons of the CDS region. This was then used as the reference distribution against which each of the SCR candidates (filtered based on density and coverage) was compared. The candidates which satisfy the following criterion were selected: the fraction of reads that

436    fell in each frame in the CDS and the ISR regions (test distributions) are within two

437    standard deviations of the reference distribution, in at least one frame.

438

439    All codes used in this study are available at:

440    https://github.com/Divyoj-Singh/Stop_codon_readthrough_pipeline

441

442    **Experimental validation**

443    *Plasmid constructs*: Luciferase constructs for luminescence-based SCR assay were

444    generated in pcDNA 3.1 backbone. The coding sequence of the test gene along with

445    the canonical stop codon and the ISR was cloned upstream of and in-frame with the

446    coding sequence of the firefly luciferase (FLuc) between *Hind*III and *BamH*I sites (*MUB6*

447    and *RPS15AD*) or *Kpn*I and *BamH*I sites (*CAM1* and *CURT1B*). A linker sequence

448    (GGCGGCTCCGGCGGCTCCCTCGTGCTCGGG) was included upstream of the FLuc

449    coding sequence.

450    *In vitro transcription and translation*: The plasmid DNA was linearized using *Not*I

451    enzyme, and 2 μg of the linearized DNA was transcribed *in vitro* using T7 RNA

452    polymerase (Thermo Fisher Scientific). The resultant RNA was purified using GeneJET

453    RNA purification kit (Thermo Fisher Scientific). The concentration and quality of the

454    RNA were measured using BioPhotometer (Eppendorf). 2-3 μg of the purified RNA was

455    *in vitro* translated using wheat germ extract (Promega) at 25 °C for 2 h as per the

456    manufacturer's instructions. Luciferase activity was then measured using the Luciferase

457    Assay System (Promega Corporation) in the GloMax Explorer System (Promega

458    Corporation).

**Acknowledgements**

**Author contributions**

**SS:** Conceptualization, Data Curation, Formal Analysis, Investigation, Methodology, Visualization, Software. **DS:** Data Curation, Investigation, Methodology, Visualization, Software. **AS:** Investigation, Methodology. **SME**: Conceptualization, Funding Acquisition, Formal analysis, Visualization, Project Administration, Resources, Supervision, Writing – Original Draft Preparation.

482 **Figure Legends**

483 **Figure 1. Selection and analysis of ribosome profiling datasets**

484 (A) Heat map showing the three-nucleotide periodicity profile of the dataset

485 SRP074840. Ribosomal footprints on all coding sequences were analyzed to get this

486 profile for read lengths 24, 25, and 26. Reads were assigned to three frames based on

487 which frame the first nucleotide aligns with on an mRNA sequence. The start codon

488 ATG is indicated by the position 0 on the x-axis.

489 (B) Distribution of ribo-seq read lengths. The graph shown is from the dataset

490 SRP074840

491 (C) Flow chart showing the four-level screening method to identify mRNAs that show

492 SCR in *A. thaliana*. Another flow chart with more details is shown in Fig S2.

493

494 **Figure 2. Ribosomal density and three-nucleotide periodicity at the ISR of 4 SCR-**

495 **positive mRNAs – *RPS15AD, CURT1B, CAM1, MUB6*.**

496 (A) Graphs showing ribo-seq reads in the ISR of four genes. Red arrows indicate the

497 position of the two stop codons. Some parts of the coding sequence and the 3′UTR are

498 also shown for comparison.

499 (B) Three-nucleotide periodicity. Graphs show fraction of ribo-seq reads in three

500 translation frames. The three-nucleotide periodicity profile of coding sequences of all

501 protein-coding genes is shown for comparison (All CDS).

502 The data shown in (A) and (B) are from the dataset SRP074840. CDS, coding

503 sequence; ISR, inter-stop codon region; UTR, untranslated region.

504

505 **Figure 3. Ribosomal density in the ISR of SCR-positive mRNAs**

506 (A) The graph shows the increase in ribosome density in the ISRs after each round of

507 screening. *, P < 0.001 Mann-Whitney Rank Sum Test (compared to 'All mRNAs'). (B)

508 Graph shows the comparison of ribosomal density in the ISR vs that in the 3′UTR. This

509 comparison is shown for SCR-positive mRNAs and for all mRNAs. P values were

510 calculated using Mann-Whitney Rank Sum Test.

511 Numbers at the bottom of the graph indicate the mean value. The box represents 25%

512 and 75% values, and the horizontal line within the box shows the median value. The

513 analysis shown is for the dataset SRP074840.

514

515 **Figure 4. The canonical stop codon and its context in SCR-positive mRNAs**

516 (A) Distribution of the three stop codons in SCR-positive mRNAs. Expected values were

517 obtained based on their occurrence in all mRNAs of *A. thaliana*.

518 (B) Sequence logo of the stop codon context of SCR-positive mRNAs. The analysis was

519 performed using WebLogo.

520

521 **Figure 5. Gene ontology analysis of SCR-positive genes of *A. thaliana***

522 Results of gene ontology (GO) functional enrichment analysis on SCR-positive genes

523 using PANTHER web server. The X-axis shows false discovery rate and the Y-axis

524 shows multiple functional classes enriched in SCR-positive genes. Color of the circle

525 indicates fold enrichment. The number of SCR-positive genes showing enrichment in a

526 functional group is shown next to the circle. Size of the circle is proportional to this

527 number. SCR-positive genes showing more than 4-fold enrichment are shown here.

528 **Figure 6. Experimental validation of SCR in four *A. thaliana* mRNAs – *RPS15AD*,**

529 ***CURT1B*, *CAM1* and *MUB6*.**

530 Luminescence-based SCR assay. cDNA of a test gene along with the ISR was cloned

531 upstream of and in-frame with the cDNA of firefly luciferase (FLuc) such that FLuc is

532 expressed only if there is SCR across the stop codon of the test cDNA (see the

533 schematic). Constructs were subjected to *in vitro* transcription followed by *in vitro*

534 translation as described in Methods. Expression of Fluc was measured by its

535 luminescence activity, which is shown in the graphs. Constructs without ISR were used

536 to measure background signal (first bar), and constructs without any stop codon

537 between the test cDNA and the FLuc were used to measure the maximum

538 luminescence activity (third bar). Statistical significance (two-sided *P*-value) was

539 obtained using Student's t-test. Input RNA obtained by *in vitro* transcription is shown

540 below the graphs.

541

542

543

544

545

546

547

548

549

550

24

## REFERENCES

1. Schueren F, Thoms S (2016) Functional Translational Readthrough: A Systems Biology Perspective. PLoS Genet 12: e1006196.

2. Eswarappa SM, Potdar AA, Koch WJ, Fan Y, Vasu K, et al. (2014) Programmed translational readthrough generates antiangiogenic VEGF-Ax. Cell 157: 1605-1618.

3. Manjunath LE, Singh A, Sahoo S, Mishra A, Padmarajan J, et al. (2020) Stop codon read-through of mammalian MTCH2 leading to an unstable isoform regulates mitochondrial membrane potential. J Biol Chem 295: 17009-17026.

4. Singh A, Manjunath LE, Kundu P, Sahoo S, Das A, et al. (2019) Let-7a-regulated translational readthrough of mammalian AGO1 generates a microRNA pathway inhibitor. EMBO J 38: e100727.

5. Schueren F, Lingner T, George R, Hofhuis J, Dickel C, et al. (2014) Peroxisomal lactate dehydrogenase is generated by translational readthrough in mammals. Elife 3: e03640.

6. Dreher TW, Miller WA (2006) Translational control in positive strand RNA plant viruses. Virology 344: 185-197.

7. Miras M, Miller WA, Truniger V, Aranda MA (2017) Non-canonical Translation in Plant RNA Viruses. Front Plant Sci 8: 494.

8. Xu Y, Ju HJ, DeBlasio S, Carino EJ, Johnson R, et al. (2018) A Stem-Loop Structure in Potato Leafroll Virus Open Reading Frame 5 (ORF5) Is Essential for Readthrough Translation of the Coat Protein ORF Stop Codon 700 Bases Upstream. J Virol 92.

9. Newburn LR, Nicholson BL, Yosefi M, Cimino PA, White KA (2014) Translational readthrough in Tobacco necrosis virus-D. Virology 450-451: 258-265.

10. Nyiko T, Auber A, Szabadkai L, Benkovics A, Auth M, et al. (2017) Expression of the eRF1 translation termination factor is controlled by an autoregulatory circuit involving readthrough and nonsense-mediated decay in plants. Nucleic Acids Res 45: 4174-4188.

11. Urquidi Camacho RA, Lokdarshi A, von Arnim AG (2020) Translational gene regulation in plants: A green new deal. Wiley Interdiscip Rev RNA 11: e1597.

12. Dunn JG, Foo CK, Belletier NG, Gavis ER, Weissman JS (2013) Ribosome profiling reveals pervasive and regulated stop codon readthrough in Drosophila melanogaster. Elife 2: e01179.

13. Merchante C, Brumos J, Yun J, Hu Q, Spencer KR, et al. (2015) Gene-specific translation regulation mediated by the hormone-signaling molecule EIN2. Cell 163: 684-697.

14. Hsu PY, Calviello L, Wu HL, Li FW, Rothfels CJ, et al. (2016) Super-resolution ribosome profiling reveals unannotated translation events in Arabidopsis. Proc Natl Acad Sci U S A 113: E7126-E7135.

15. Willems P, Ndah E, Jonckheere V, Stael S, Sticker A, et al. (2017) N-terminal Proteomics Assisted Profiling of the Unexplored Translation Initiation Landscape in Arabidopsis thaliana. Mol Cell Proteomics 16: 1064-1080.

16. Bazin J, Baerenfaller K, Gosai SJ, Gregory BD, Crespi M, et al. (2017) Global analysis of ribosome-associated noncoding RNAs unveils new modes of translational regulation. Proc Natl Acad Sci U S A 114: E10018-E10027.

17. Kurihara Y, Makita Y, Kawashima M, Fujita T, Iwasaki S, et al. (2018) Transcripts from downstream alternative transcription start sites evade uORF-mediated inhibition of gene expression in Arabidopsis. Proc Natl Acad Sci U S A 115: 7831-7836.

18. Liu MJ, Wu SH, Wu JF, Lin WD, Wu YC, et al. (2013) Translational landscape of photomorphogenic Arabidopsis. Plant Cell 25: 3699-3710.

19. Chotewutmontri P, Barkan A (2018) Multilevel effects of light on ribosome dynamics in chloroplasts program genome-wide and psbA-specific changes in translation. PLoS Genet 14: e1007555.

20. Waltz F, Nguyen TT, Arrive M, Bochler A, Chicher J, et al. (2019) Small is big in Arabidopsis mitochondrial ribosome. Nat Plants 5: 106-117.

21. Guydosh NR, Green R (2014) Dom34 rescues ribosomes in 3' untranslated regions. Cell 156: 950-962.

22. Palma M, Lejeune F (2021) Deciphering the molecular mechanism of stop codon readthrough. Biol Rev Camb Philos Soc 96: 310-329.

612    23. Floquet C, Hatin I, Rousset JP, Bidou L (2012) Statistical analysis of readthrough
613         levels for nonsense mutations in mammalian cells reveals a major determinant of
614         response to gentamicin. PLoS Genet 8: e1002608.

615    24. Crooks GE, Hon G, Chandonia JM, Brenner SE (2004) WebLogo: a sequence logo
616         generator. Genome Res 14: 1188-1190.

617    25. Mi H, Muruganujan A, Ebert D, Huang X, Thomas PD (2019) PANTHER version 14:
618         more genomes, a new PANTHER GO-slim and improvements in enrichment
619         analysis tools. Nucleic Acids Res 47: D419-D426.

620    26. Stiebler AC, Freitag J, Schink KO, Stehlik T, Tillmann BA, et al. (2014) Ribosomal
621         readthrough at a short UGA stop codon context triggers dual localization of
622         metabolic enzymes in Fungi and animals. PLoS Genet 10: e1004685.

623    27. Hofhuis J, Schueren F, Notzel C, Lingner T, Gartner J, et al. (2016) The functional
624         readthrough extension of malate dehydrogenase reveals a modification of the
625         genetic code. Open Biol 6.

626    28. Reumann S, Buchwald D, Lingner T (2012) PredPlantPTS1: A Web Server for the
627         Prediction of Plant Peroxisomal Proteins. Front Plant Sci 3: 194.

628    29. Lin JR, Hu J (2013) SeqNLS: nuclear localization signal prediction based on
629         frequent pattern mining and linear motif scoring. PLoS One 8: e76864.

630    30. Nguyen Ba AN, Pogoutse A, Provart N, Moses AM (2009) NLStradamus: a simple
631         Hidden Markov Model for nuclear localization signal prediction. BMC
632         Bioinformatics 10: 202.

633    31. Krogh A, Larsson B, von Heijne G, Sonnhammer EL (2001) Predicting
634         transmembrane protein topology with a hidden Markov model: application to
635         complete genomes. J Mol Biol 305: 567-580.

636    32. Maurer-Stroh S, Eisenhaber F (2005) Refinement and prediction of protein
637         prenylation motifs. Genome Biol 6: R55.

638    33. Meszaros B, Erdos G, Dosztanyi Z (2018) IUPred2A: context-dependent prediction
639         of protein disorder as a function of redox state and protein binding. Nucleic Acids
640         Res 46: W329-W337.

641    34. Virdi AS, Singh S, Singh P (2015) Abiotic stress responses in plants: roles of
642         calmodulin-regulated proteins. Front Plant Sci 6: 809.

643   35. Ingolia NT, Ghaemmaghami S, Newman JR, Weissman JS (2009) Genome-wide
644       analysis in vivo of translation with nucleotide resolution using ribosome profiling.
645       Science 324: 218-223.
646   36. Kaiser CA, Preuss D, Grisafi P, Botstein D (1987) Many random sequences
647       functionally replace the secretion signal sequence of yeast invertase. Science
648       235: 312-317.
649   37. Kaiser CA, Botstein D (1990) Efficiency and diversity of protein localization by
650       random signal sequences. Mol Cell Biol 10: 3163-3173.
651   38. Houck-Loomis B, Durney MA, Salguero C, Shankar N, Nagle JM, et al. (2011) An
652       equilibrium-dependent retroviral mRNA switch regulates translational recoding.
653       Nature 480: 561-564.
654   39. Firth AE, Wills NM, Gesteland RF, Atkins JF (2011) Stimulation of stop codon
655       readthrough: frequent presence of an extended 3' RNA structural element.
656       Nucleic Acids Res 39: 6679-6691.

657

658

659

660

661

662

663

664

665

666

667

**Table 1. List of SCR-positive genes whose products exhibit peroxisomal targeting sequence after SCR$^\$$**

| Gene | Function* | Location of the canonical protein* | Peptide encoded by the ISR# |
|---|---|---|---|
| AT1G09310 | Unknown | apoplast, cytosol, extracellular region, nucleus | SAQLQQIKETRIFKCT<span style="color:red">SRI</span> |
| AT5G13930 (*CHS*) | chalcone synthase involved in the biosynthesis of flavonoids | cytoplasm, endoplasmic reticulum, nucleus, plant-type vacuole membrane | ERLPSICLPTY<span style="color:red">AKL</span> |
| AT5G11740 (*AGP15*) | arabinogalactan protein | plasma membrane | VTVMVISYRDCFCGIGHSS LFVVSCVFR<span style="color:red">SSL</span> |
| AT1G78680 (*GGH2*) | a gamma-glutamyl hydrolase acting specifically on monoglutamates. | Vacuole | NGGFCRIGYDEVYIFTQQR<span style="color:red">SLL</span> |

$ The analysis was done using the PredPlantPTS1 tool

*Function and location information were obtained from The Arabidopsis Information

Resource (TAIR)

# Experimentally verified plant PTS (peroxisome targeting sequence) tripeptides are

highlighted in red

29

**Table 2. List of SCR-positive genes whose products exhibit Nuclear localization**

**signal (NLS) in the extended C-terminus[$]**

| Gene | Function[*] | Peptide encoded by ISR[#] |
|---|---|---|
| AT2G46820 (*CURT1B*) | The P subunit of Photosystem I | <span style="color:red">IKGGRRRRR</span>AFLRPFMNWNE GYQKNLTQRPRPSFNLSFL (0.727) |
| AT2G28630 (*KCS12*) | 3-ketoacyl-CoA synthase (involved in the biosynthesis of very long chain fatty acids) | <span style="color:red">NVYAQKRKRKRK</span>NNT RIELVKTCLAIGKPNKCV (0.895) |
| AT5G56200 | Encodes a transcription factor expressed in the female gametophyte. | ETYICKQVIF<span style="color:red">LTLKKKKTKK</span> (0.862) |

[$]The analysis was done using SeqNLS and NLStradamus

[*]Function information was obtained from The Arabidopsis Information Resource (TAIR)
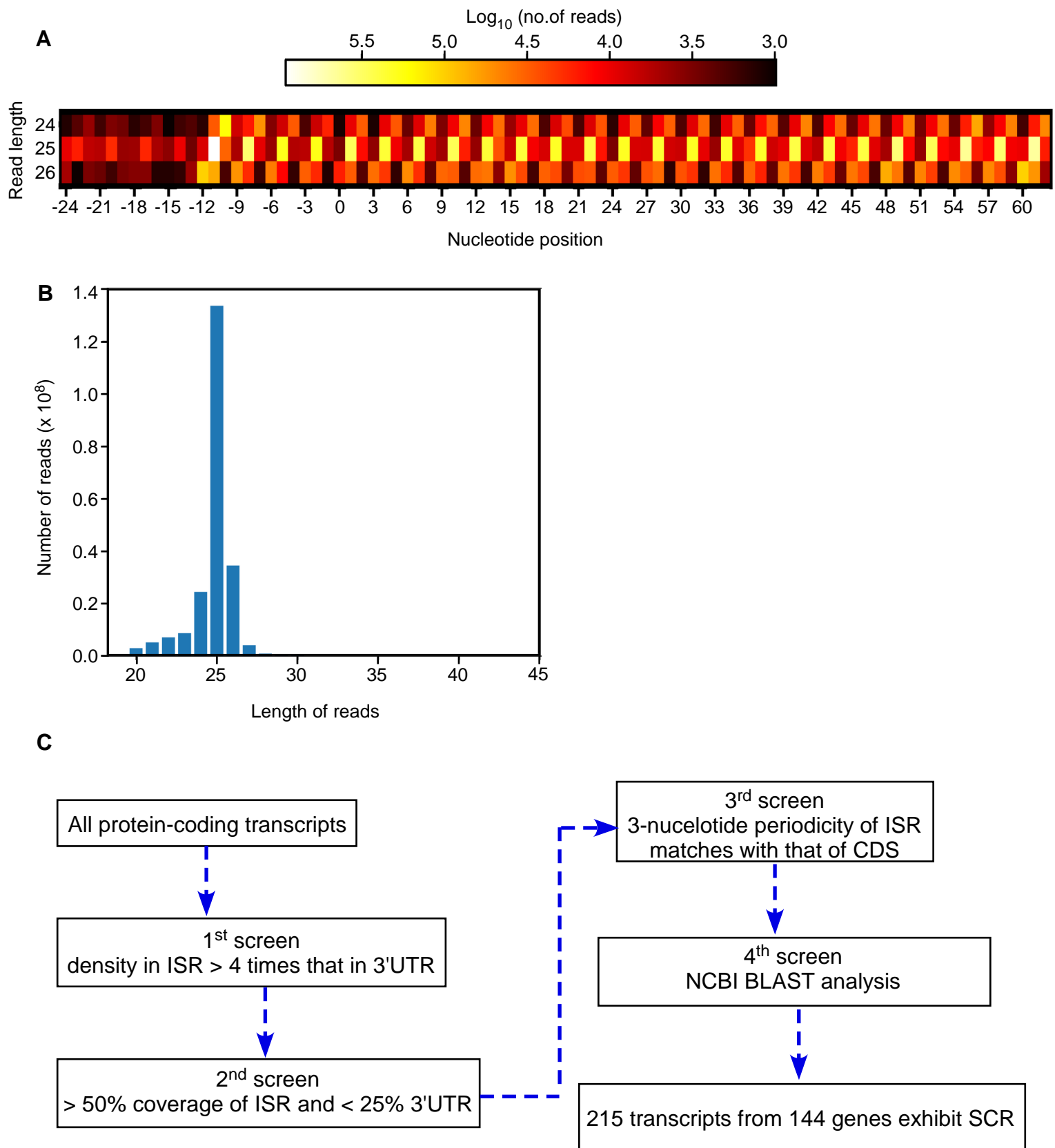
[#]NLSs are highlighted in red and the numbers indicate the SeqNLS score

30

**Table 3. List of SCR-positive genes whose products exhibit intrinsically disordered regions in the extended C-terminus after SCR$^\$$**

| Gene | Function* | Peptide encoded by ISR# |
|---|---|---|
| AT2G01910 (*MAP65-6*) | Binds microtubules. Induces a crisscross mesh of microtubules. | LDSLFHRICGVMLMVKK EGSEEE<span style="color:red">GRRLVNTEGD</span> |
| AT3G23030 (*IAA2*) | Auxin inducible gene expressed in the nucleus | <span style="color:red">SREAENLLSKKEMMTMIDE</span> |
| AT3G14415 (*GOX2*) | Glycolate oxidase | <span style="color:red">RRKKKQRTETTRHQNVFIF</span> |
| AT5G43470 (*RPP8*) | Hypersensitive response to turnip crinkle virus | <span style="color:red">QERPRSEPNSLILGDID</span>AA STESSADQQVFPKNIWYCL |
| AT2G28630 (*KCS12*) | 3-ketoacyl-CoA synthase (involved in the biosynthesis of very long chain fatty acids) | <span style="color:red">NVYAQKRKRKRK</span>NNT RIELVKTCLAIGKPNKCV |
| AT5G14740 (*BETA CA2*) | Beta carbonic anhydrase | <span style="color:red">TNTSPSPSLPPPSQTSSSSSSS</span> |

$^\$$The analysis was done using IUPred2A

*Function information was obtained from The Arabidopsis Information Resource (TAIR)

#Intrinsically disordered regions are highlighted in red

**Figure 1**

**A**
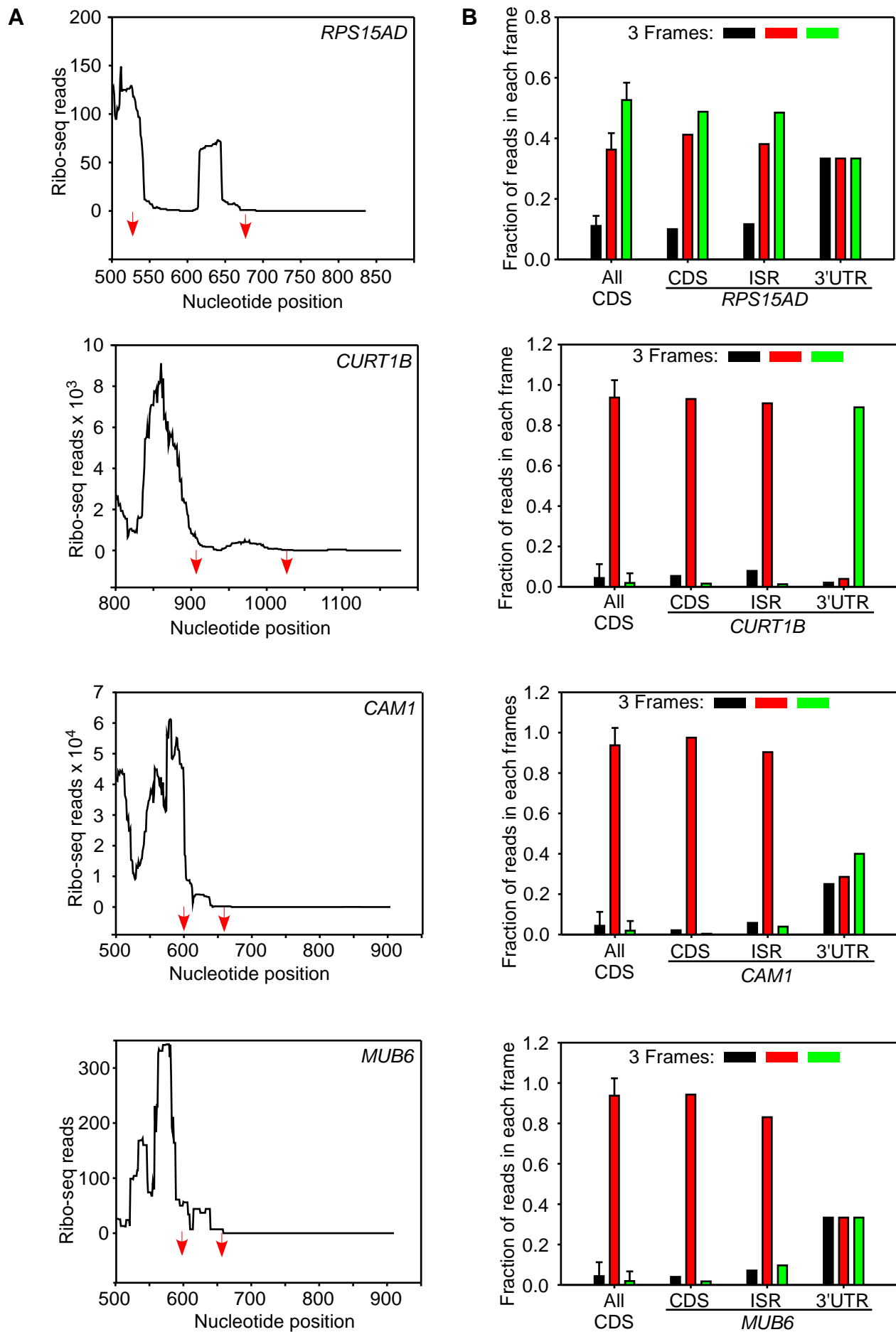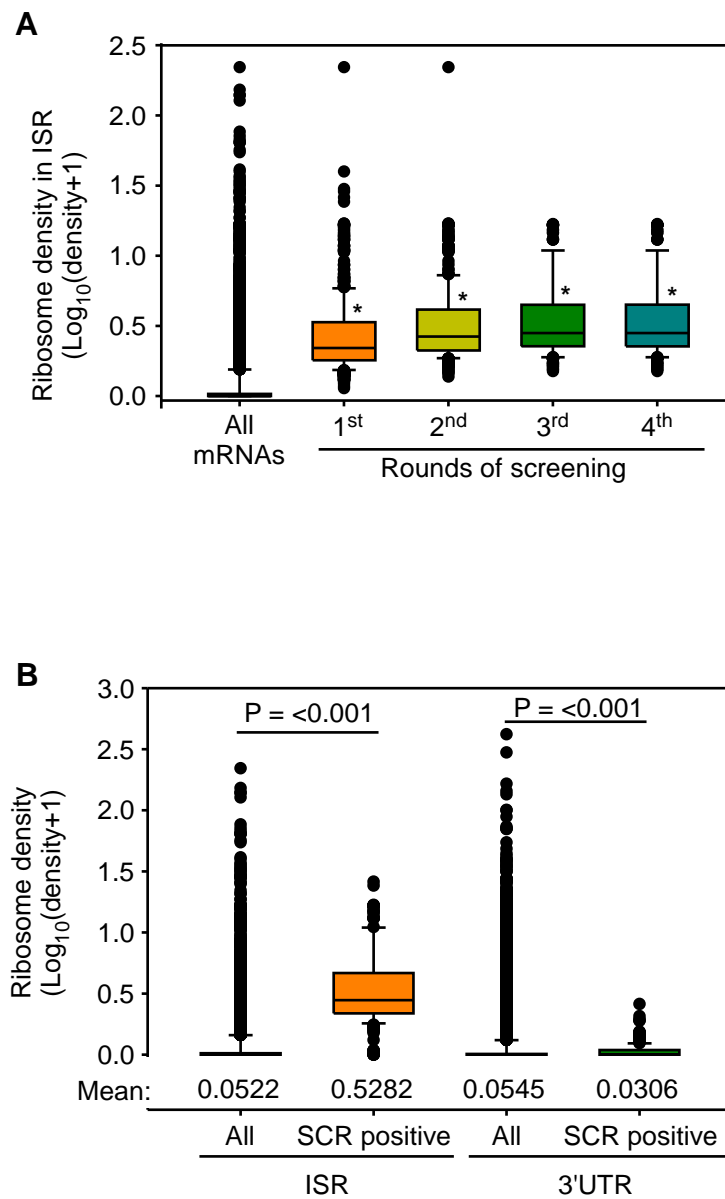


**B**



**C**

## Figure 2

# Figure 3

# Figure 4

**A**
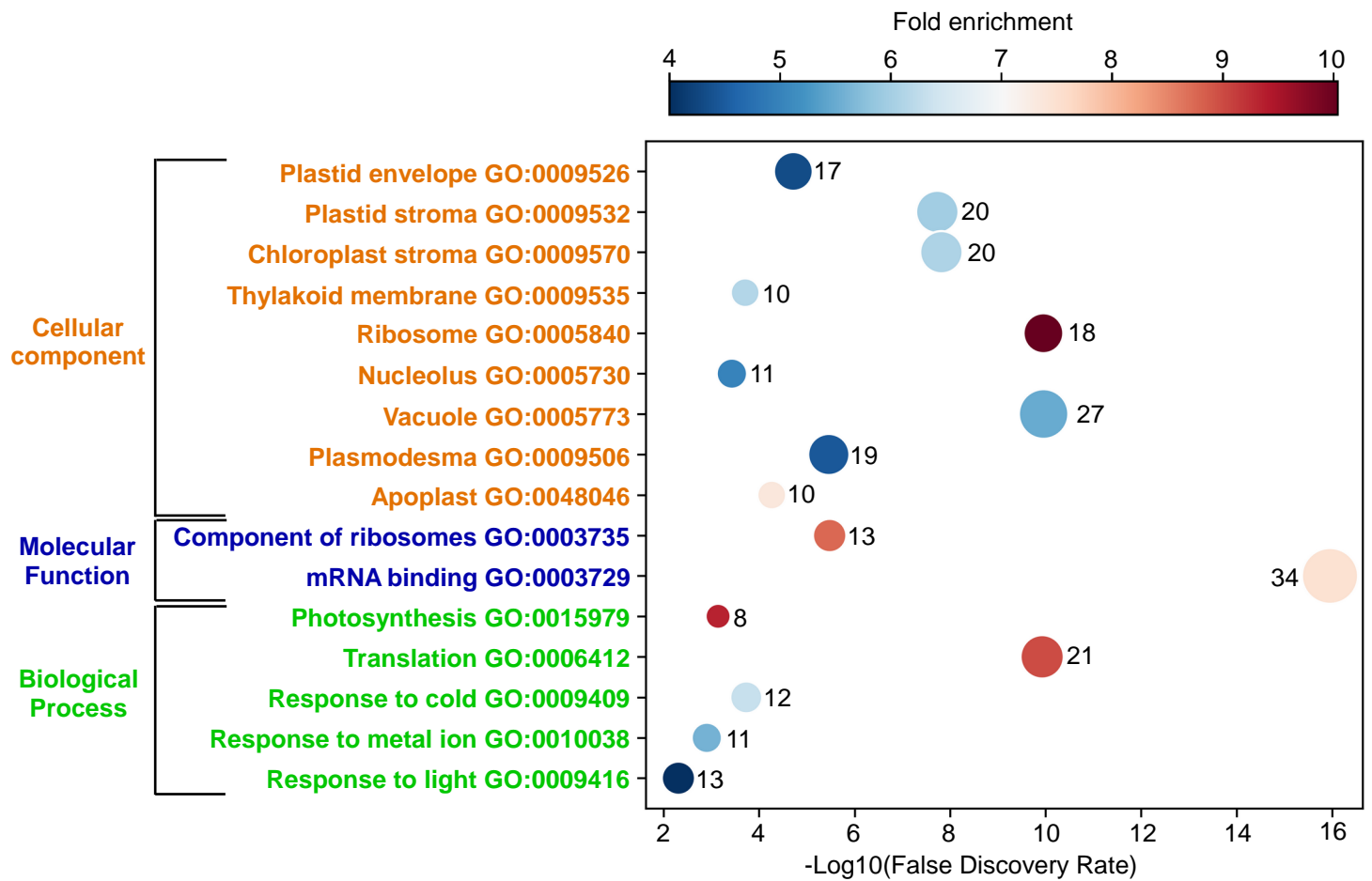


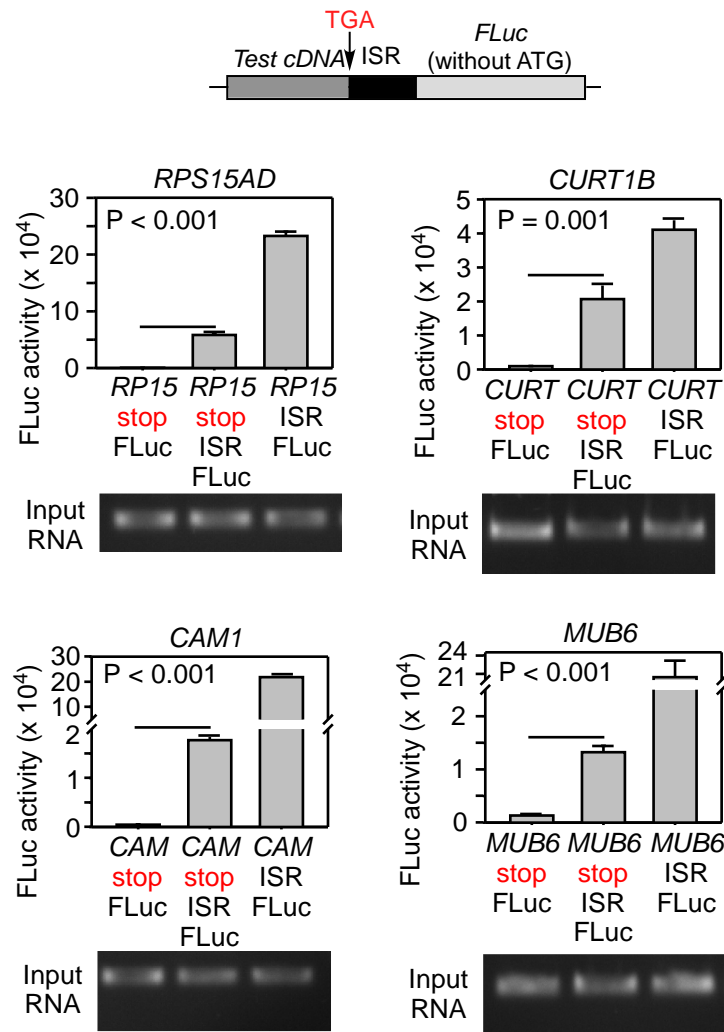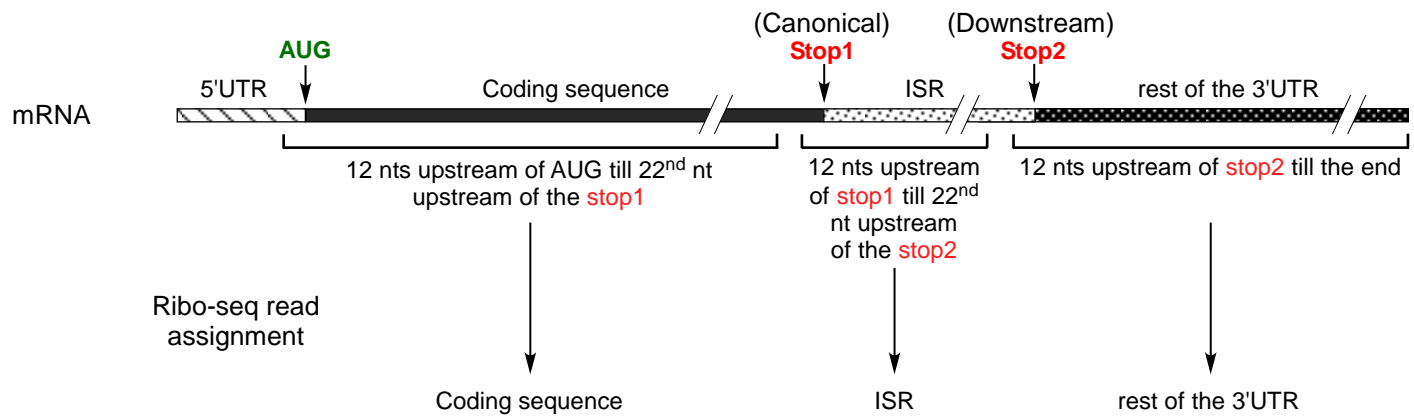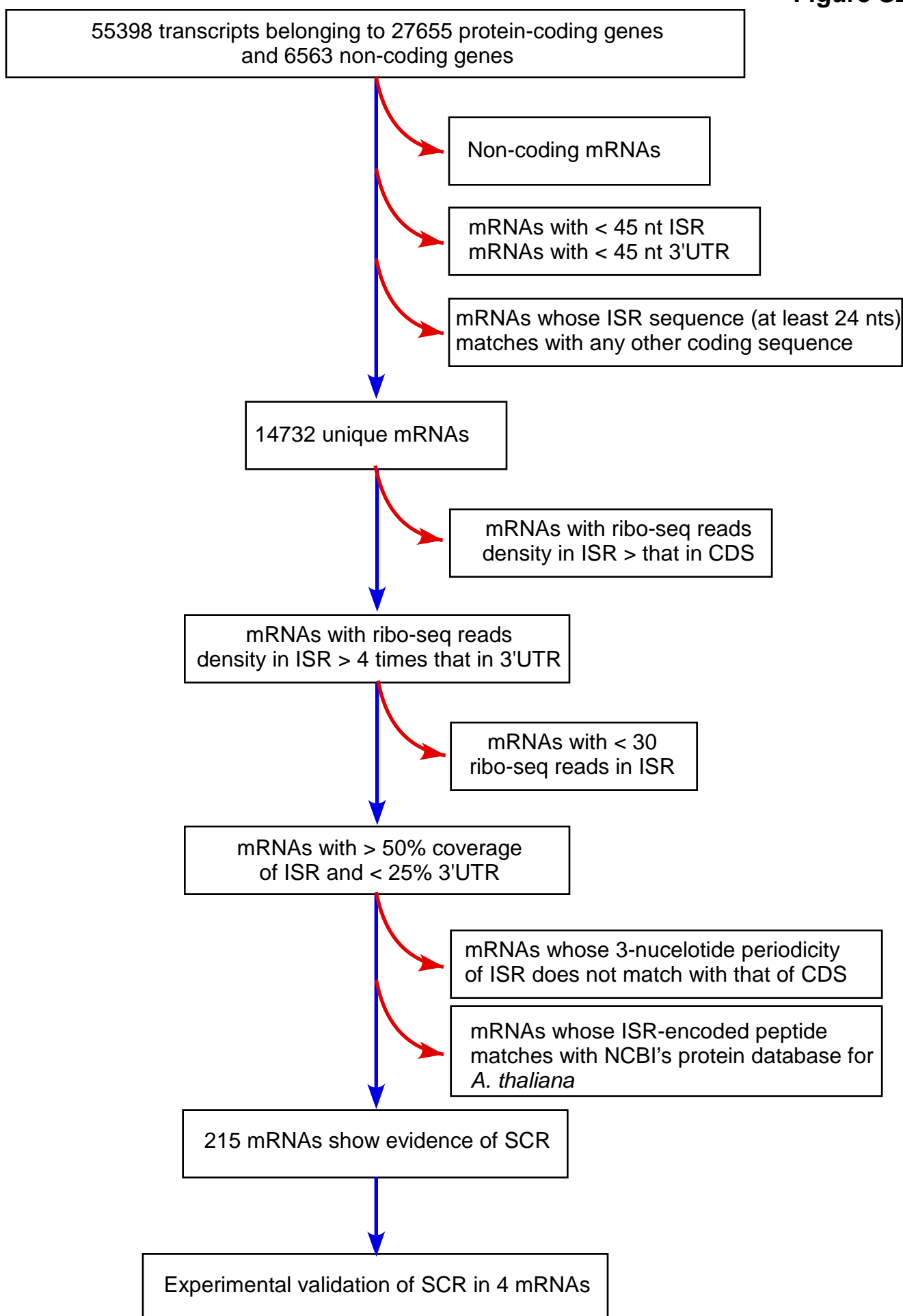**B**

# Figure 5

**Figure 6**

# Figure S1

**Figure S2**

**Legends to supplementary figures**

**Figure S1**. Schematic to explain the assignment of ribo-seq reads to various regions of an mRNA.

**Figure S2**.  Flow chart showing the four-level screening method to identify mRNAs that show SCR in *A. thaliana.* Blue arrows indicate inclusion and red arrows indicate exclusion.