1  **A research paper for Genome Research**

2

3  **Title**: Sequence features around cleavage sites are highly conserved among

4  different species and a critical determinant for RNA cleavage position across

5  eukaryotes

6

7  **Running head**:

8  Conserved sequences related to RNA degradation

9

10  **Authors**: Daishin Ueno[1], Shotaro Yamasaki[1], Yuta Sadakiyo[1], Takumi

11  Teruyama[1], Taku Demura[1], and Ko Kato[1*]

12

13  **Author's Addresses**:

14  1. Graduate School of Science and Technology, Nara Institute of Science and

15  Technology, Ikoma, 630-0192 Japan

16  *Corresponding author: Ko Kato, E-mail, kou@bs.naist.jp; Tel: +81-743-72-

17  5461; Fax +81-743-72-5469.

18

19  **Keywords**: RNA degradation, degradome sequencing, machine learning

20

21

22

23

24

25

26

27

28  **ABSTRACT**

29  RNA degradation is critical for control of gene expression, and

30  endonucleolytic cleavage–dependent RNA degradation is conserved among

31  eukaryotes. Some cleavage sites are secondarily capped in the cytoplasm

32  and identified using the CAGE method. Although uncapped cleavage sites are

33  widespread in eukaryotes, comparatively little information has been obtained

34  about these sites using CAGE-based degradome analysis. Previously, we

35  developed the truncated RNA-end sequencing (TREseq) method in plant

36  species and used it to acquire comprehensive information about uncapped

37  cleavage sites; we observed G-rich sequences near cleavage sites. However,

38  it remains unclear whether this finding is general to other eukaryotes. In this

39  study, we conducted TREseq analyses in fruit flies (*Drosophila melanogaster*)

40  and budding yeast (*Saccharomyces cerevisiae*). The results revealed specific

41  sequence features related to RNA cleavage in *D. melanogaster* and *S.*

42  *cerevisiae* that were similar to sequence patterns in *Arabidopsis thaliana*.

43  Although previous studies suggest that ribosome movements are important

44  for determining cleavage position, feature selection using a random forest

45  classifier showed that sequences around cleavage sites were major

46  determinant for cleaved or uncleaved sites. Together, our results suggest that

47  sequence features around cleavage sites are critical for determining cleavage

48  position, and that sequence-specific endonucleolytic cleavage–dependent

49  RNA degradation is highly conserved across eukaryotes.

50

51

52 **INTRODUCTION**

53 Multiple mechanisms are involved in control of gene expression. RNA

54 degradation is a critical step in maintaining RNA homeostasis in eukaryotes

55 (Keene 2010). RNA degradation is initiated either by deadenylation followed

56 by removal of the cap and 5'-to-3' or 3'-to-5' exonucleolytic digestion, or by

57 endonucleolytic cleavage followed by exonucleolytic digestion (Parker 2012).

58 Deadenylation-dependent RNA degradation is well studied in yeast, plants,

59 and animals, and proteins and sequences related to deadenylation or

60 decapping have been characterized in detail (Wu and Brewer 2008; Chiba and

61 Green 2009; Parker 2012). Although microRNA- and siRNA-dependent forms

62 of RNA degradation are well studied in eukaryotes, endonucleolytic cleavage–

63 dependent RNA degradation has not been analyzed in detail.

64 　　　Previous studies developed several degradome sequencing methods to

65 detect cleavage sites in plants such as *A. thaliana*, mainly in the context of

66 microRNA- or siRNA-mediated cleavage (Addo-Quaye et al. 2008; German et

67 al. 2008; Gregory et al. 2008). Using these methods, RNA cleavage sites were

68 identified in plants (Zhou et al. 2010; Shamimuzzaman and Vodkin 2012;

69 Karlova et al. 2013; Hou et al. 2016; Anderson et al. 2018), yeast, and animal

70 models (Harigaya and Parker 2012; Bracken et al. 2011). However, because

71 these methods selected poly(A) RNA (Willmann et al. 2014; Addo-Quaye et

72 al. 2008), many of the 5' degradation intermediates they detected appeared

73 to be biased toward the 3'-ends of transcripts (Weinberg et al. 2016). In

74 addition, ligation of single-stranded 5' adapters using T4 RNA ligase induces

75 bias in the detected 5' RNA ends (Hou et al. 2014; Zhuang et al. 2012).

76 Specifically, ligation efficiencies are quite low if there are stable secondary

77 structures between adapters and 5'-proximal RNA sequences (e.g., high G/C

78 frequency). In addition, non-specific annealing occurs when using single-

3

79   stranded 5' adapters and subsequent PCR amplification, resulting in

80   sequence artifacts (Hou et al. 2014).

81       Alternatively, instead of previous degradome methods, Cap analysis of

82   gene expression (CAGE) can be used to detect cleavage sites (Mercer et al.

83   2010). To decrease the bias in 5' adapter ligation efficiency, the CAGE method

84   joins single-strand cDNAs to partially double-stranded adapter sequences in

85   the presence of polyethylene glycol (PEG) (Murata et al. 2014). The partially

86   double-stranded adapter sequence has phosphate modifications, allowing a

87   dramatic reduction in 5' adapter dimers while simultaneously increasing

88   adapter ligation efficiency (Takahashi et al. 2012). In addition, formation of

89   secondary structure between adapter and 5'-proximal RNA sequences is less

90   frequent for the partially double-stranded adapter sequences than for single-

91   stranded adapter sequences, and ligation efficiencies are significantly

92   increased by the addition of PEG to the ligation buffer (Song et al. 2014).

93   Furthermore, current CAGE methods do not involve PCR amplification or

94   restriction enzyme digestion (Murata et al. 2014), allowing CAGE to detect 5'

95   RNA ends more accurately than other methods (Adiconis et al. 2018).

96   However, because the CAGE method was developed for detecting

97   transcription start sites (TSSs), the majority of 5' RNA ends are TSSs, and

98   expected cleavage sites constituted less than 20% of all detected sites

99   (Mercer et al. 2010). Also, because some cleavage sites are secondarily

100  capped in the cytoplasm (Schoenberg and Maquat 2009) and the CAGE

101  method selects Cap RNA using Cap-trapping methods, the detected cleavage

102  sites are likely to be re-capped. Uncapped RNAs are prevalent in eukaryotes

103  (Ibrahim et al. 2018; Bracken et al. 2011) but the majority of uncapped

104  cleavage sites cannot be detected using CAGE methods.

105      In a previous study, we improved these features by developing truncated

4

106   RNA-end sequencing (TREseq) based on the no-amplification non-tagging

107   CAGE (nAnT-iCAGE) method in *A. thaliana* (Ueno et al. 2018). Initially, we

108   separated RNA into Cap and Cap-less RNA using Cap-trapping method and

109   constructed libraries for Cap-less RNA, in addition to Cap RNA. To increase

110   the concentration of Cap-less RNA, we removed the majority of expected Cap

111   RNA from the Cap-less RNA data during the mapping process. Because Cap

112   (7-methylguanosine) is not encoded in the genome and nucleotide of Cap

113   position in the genome is Y (C or T) in eukaryotes (Carninci et al. 2006;

114   Adiconis et al. 2018; Thieffry et al. 2020; Lu and Lin 2019), the majority of

115   Cap RNA in the Cap-less data could be removed by excluding Cap-less RNA

116   with mismatches at the 5' end of the read relative to the genome (Ueno et al.

117   2020). In addition, we decreased the 3' bias of the cleavage site using random

118   primers and conducting rRNA depletion (Ueno et al. 2018). By improving these

119   experimental procedures, we acquired comprehensive information about

120   uncapped cleavage sites and identified a relationship between cleavage

121   efficiencies and RNA stability. Although previous degradome methods cannot

122   identify such relationships, we observed that genes with high cleavage

123   efficiencies had short half-lives relative to those with lower cleavage

124   efficiencies (Ueno et al. 2018). When we focused on nucleotide frequencies,

125   we observed G-rich sequence motifs around the cleavage sites; moreover,

126   genes with G-rich sequence motifs accumulated lower RNA levels than genes

127   without motifs (Ueno et al. 2020). These tendencies were conserved not only

128   in *A. thaliana* but also among different plant species (Ueno et al. 2018, 2020).

129   In addition, cleavage sites were common around the start and stop codon,

130   and three-nucleotide periodicity was observed in coding region sequences

131   (CDSs) in plant species (Ueno et al. 2020). These tendencies are similar to

132   ribosome movements, which have been reported by ribosome profiling

5

133 methods, suggesting that the translation process is related to RNA cleavage

134 and that ribosome position is important for determining the position of

135 cleavage sites (5' truncated RNA ends) (Ibrahim et al. 2018; Yu et al. 2016;

136 Pelechano et al. 2015). Although specific sequence features around cleavage

137 sites are observed in plant species, it remains to be determined whether these

138 properties are shared with other eukaryotes. In addition, RNA cleavage

139 position seems to be determined by multiple determinants (e.g., sequence

140 features and translation), but previous studies have not performed an

141 integrated analysis of these determinants.

142 To address this issue, we conducted TREseq in fruit flies (*Drosophila*

143 *melanogaster*) and budding yeast (*Saccharomyces cerevisiae*), and obtained

144 TREseq data in *A. thaliana* from a previous study. We also obtained ribosome

145 profiling data and analyzed the effects of the translation process and

146 sequence features on RNA cleavage across species. In addition, we

147 conducted feature selection using a random forest classifier and evaluated

148 the significance of the effects of multiple determinants on RNA cleavage

149 positions in various species. Our research revealed that G-nucleotide

150 frequencies were higher around cleavage sites in *D. melanogaster* and *S.*

151 *cerevisiae*, similar to sequence features around cleavage sites in *A. thaliana,*

152 and that these sequence features contribute to control of RNA cleavage

153 across eukaryotes.

154

155

156

157

158

159

6

**RESULTS**

**Improvement of experiment procedures and validation of cleavage sites in TREseq**

We obtained TREseq data of *D. melanogaster* and *S. cerevisiae* in this study and *A. thaliana* (DRA005995) from previous study (Ueno et al. 2018) (Supplementary Figure S1). In a previous study, we developed the TREseq method to decrease the bias towards the 3' end of the transcript (Ueno et al. 2018). TREseq can detect cleavage sites both with [poly $(A)^+$] and without poly A tails [poly $(A)^-$] using random primers, decreasing the apparent bias. In this analysis, the cleavage sites were almost equally distributed across genes in *D. melanogaster* and *S. cerevisiae*, as in *A. thaliana* (Supplementary Figure S2). If this improvement of the bias towards the 3' end of the RNA was caused by selecting poly $(A)^-$ cleavage sites in addition to poly $(A)^+$ cleavage sites, we should observe similar tendencies when poly $(A)^-$ cleavage sites are isolated during library preparation. In the parallel analysis of RNA ends (PARE) method, poly $(A)^-$ cleavage sites are selected, and poly (A) tail are added to the 3' ends of cleavage sites before reverse transcription using oligo dT primer (Nagarajan et al. 2019) (Supplementary Figure 3). Using data obtained by the PARE method in *A. thaliana*, we compared the distribution patterns of cleavage sites in libraries containing RNA with and without poly(A) tails. The 3' bias of cleavage sites was decreased by selection of poly $(A)^-$ cleavage sites, even when oligo dT primer was used in reverse transcription. These results strongly suggest that the bias of detected cleavage sites toward the 3' end of transcripts was the result of selecting poly $(A)^+$ cleavage sites during library preparation.

In previous study, cleavage sites were detected using CAGE, but a low concentration of Cap-less RNA prevented high coverage (read depth) of

7

187    uncapped cleavage site data. Therefore, we sought to determine how much

188    TREseq could increase the concentration of Cap-less RNA. We added

189    synthetic Cap-less *R-luc* RNA to total RNA derived from *A. thaliana* and

190    separated RNA into Cap RNA or Cap-less RNA according to the method for

191    TREseq library preparation (Supplementary Figure S4A). Subsequently, we

192    quantified the abundance of Cap-less *R-luc* RNA in the Cap RNA and Cap-

193    less RNA libraries using qRT-PCR. Relative to the Cap RNA library, Cap-less

194    *Rluc* RNA was 10-fold more abundant in the Cap-less RNA library

195    (Supplementary Figure S4B). We also confirmed that RNA could be clearly

196    separated into Cap and Cap-less RNA by comparing a TREseq analysis in *A.*

197    *thaliana* that did not use the Cap-trapping method for library preparation

198    (Supplementary Figure S4C). The results revealed that Cap-less RNA was

199    highly concentrated using the Cap-trapping method, and that the degree of

200    Cap RNA contamination in Cap-less RNA data was low (Supplementary Figure

201    S4C). The same is true for Cap RNA data (Supplementary Figure S4C). These

202    results indicate that TREseq yields a higher concentration of Cap-less RNA

203    than the CAGE method.

204        To detect 5' RNA ends, several methods were developed in previous

205    studies, in particular for detection of transcription start sites (Adiconis et al.

206    2018). We can roughly classify methods into several groups: ligation of

207    partially double-stranded adapter (CAGE or TREseq), template switching

208    (RAMPAGE, STRT, and NanoCAGE XL), and ligation of single-stranded

209    adapter (previous degradome analysis) (Adiconis et al. 2018). Although

210    nanoCAGE (template switching), which was developed for detecting TSSs in

211    nanogram ranges of total RNA, included sequence artifacts, the experimental

212    procedures were improved in the present CAGE method, which can detect 5'

213    RNA ends more accurately than the other methods (Adiconis et al. 2018;

8

214 Murata et al. 2014). Because TREseq was developed based on CAGE (nAnT-
215 iCAGE) and uses partially double-strand adapter sequences (no PCR
216 amplification), sequence artifacts reported in previous degradome methods
217 (Hou et al. 2014) theoretically should not occur. Although expected cleavage
218 sites detected by CAGE are validated by 5' RACE (Carninci et al. 2006), and
219 cleavage positions in microRNA, siRNA, or long non-coding RNAs (lncRNAs)
220 detected by CAGE matched those in the database (Mercer et al. 2010; Nepal
221 et al. 2013; De Rie et al. 2017), we also confirmed that cleavage sites in
222 TREseq could be detected using different methods in *A. thaliana*. As a
223 different method, we used Nanopore sequencing (template-switching method).
224 We selected 3000 cleavage sites whose detected reads were highest in each
225 gene, and found that these sites were associated with high G-nucleotide
226 frequencies (Supplementary Figure S5A), similar to the sequence features
227 around sites with high cleavage efficiencies identified by TREseq in *A.
228 thaliana* (Ueno et al. 2018). A Nanopore sequencer can detect 5' RNA ends
229 similarly to a short-read sequencer (Parker et al. 2020), but accurate mapping
230 in continuous identical nucleotides (e.g., G-rich sequences) is more difficult
231 at one-nucleotide resolution. Although the peaks were slightly shifted, we
232 observed similar tendencies between TREseq and Nanopore sequencing
233 (Supplementary Figure S5B). These results suggest that the detected
234 cleavage sites are reliable, and that high G-nucleotide frequencies around
235 cleavage sites are not sequence artifacts of TREseq.

236      Based on these results, we concluded that TREseq can decrease the bias
237 of detected cleavage sites toward the 3' end of transcripts and improve the
238 concentration of Cap-less RNA. In addition, we validated these sites by a
239 different method (Nanopore sequencer), and suggesting that the cleavage
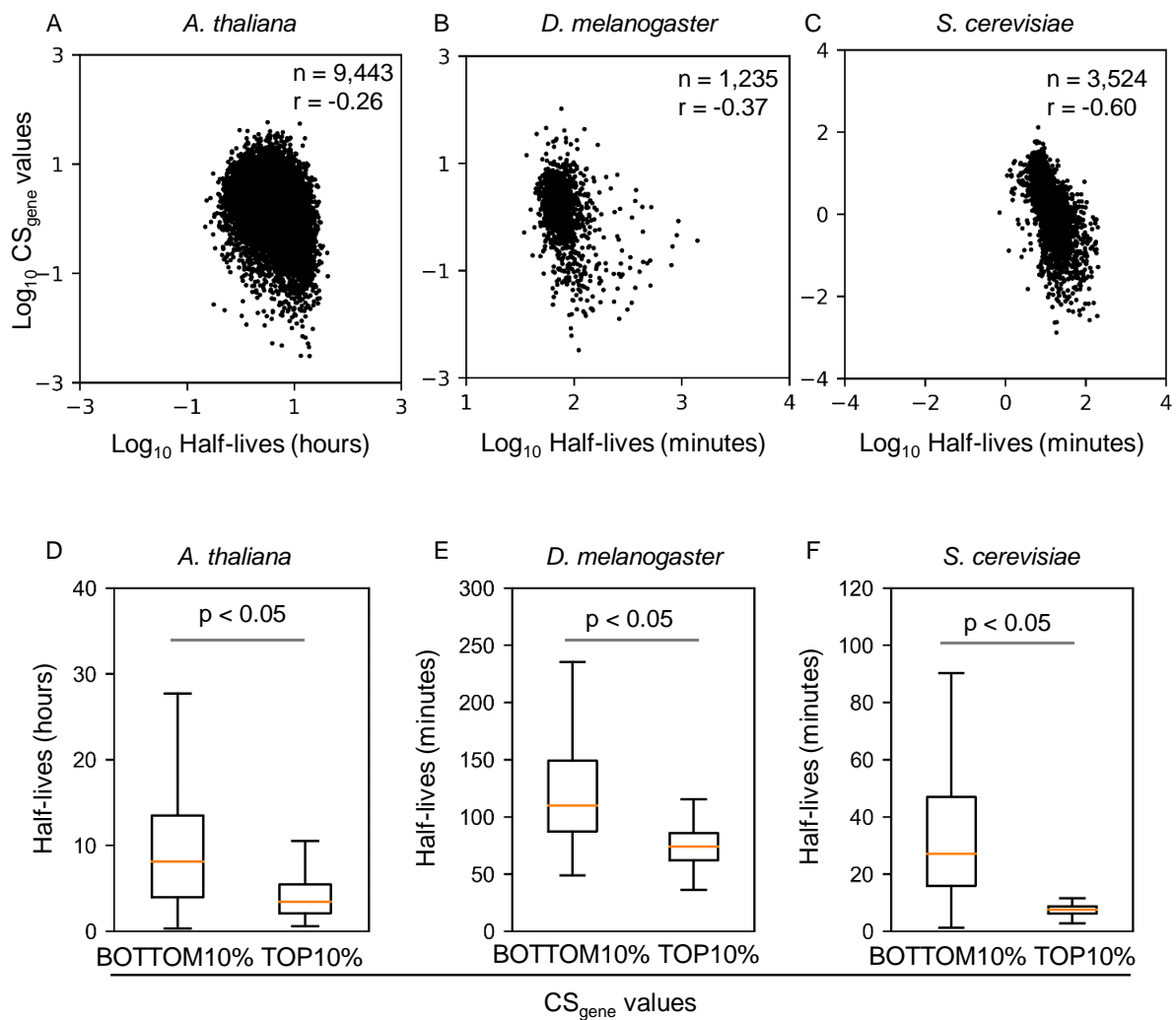240 sites detected by the TREseq method are reliable.

**241**    **Calculating the cleavage score (CS) in TREseq**

**242**    In degradome analysis, the term "cleavage score" was used to estimate RNA

**243**    cleavage efficiencies in previous studies (Anderson et al. 2018; Ueno et al.

**244**    2018). In this analysis, we defined cleavage score at the site level ($CS_{site}$) as

**245**    the number of 5' degradation intermediates normalized against RNA

**246**    abundance (Ueno et al. 2018, 2020). In addition, as with efficiencies of RNA

**247**    cleavage at the gene level, we defined the sum of $CS_{site}$ values in each gene

**248**    as the $CS_{gene}$ value. We calculated $CS_{site}$ and $CS_{gene}$ values in *S. cerevisiae*,

**249**    *D. melanogaster* and *A. thaliana* (Supplementary Figures S6-8). We confirmed

**250**    the reproducibility of $CS_{site}$ values in two biological replicates (Supplementary

**251**    Figure S9); the $CS_{site}$ values were almost normally distributed in both species,

**252**    as in *A. thaliana* (Supplementary Figure S9).

**253**

**254**    **Cleavage-dependent degradation mechanisms are important for**

**255**    **controlling RNA stability across species**

**256**    Initially, we compared the $CS_{gene}$ values with RNA half-lives obtained from

**257**    previous studies (Narsai et al. 2007; Neymotin et al. 2014; Burow et al. 2018).

**258**    $CS_{gene}$ values were related to half-life information in a previous study (Ueno

**259**    et al. 2018), and we confirmed that genes with high $CS_{gene}$ values had shorter

**260**    half-lives using different half-life measurement data from *A. thaliana* (Figure

**261**    1). $CS_{gene}$ values and half-lives were negatively correlated in *D. melanogaster*

**262**    and *S. cerevisiae*, as well as in *A. thaliana* (Figure 1). When we selected

**263**    genes with TOP and BOTTOM 10% $CS_{gene}$ values, similar tendencies were

**264**    observed in *A. thaliana* (median = 3.4 h in TOP 10%; median = 8.1 h in

**265**    BOTTOM 10%), *D. melanogaster* (median = 74 min in TOP 10%; median = 110

**266**    min in BOTTOM 10%) and *S. cerevisiae* (median = 7.5 min in TOP 10%;

**267**    median = 27 min in BOTTOM 10%). Gene ontology enrichment analysis

268    revealed that regulation of transcription, phosphorylation, signal transduction,

269    and cell cycle processes were enriched in genes with high $CS_{gene}$ values, and

270    that translation was enriched in genes with low $CS_{gene}$ values, across species

271    (Supplementary Tables S1, S2). These results are consistent with previous

272    RNA half-life analyses concluding that RNAs with short half-lives were quick

273    responses to stimuli, whereas RNAs with long half-lives tended to be involved

274    in processes related to constant maintenance of cellular physiology (Tani et

275    al. 2012; Narsai et al. 2007; Ueno et al. 2018). Although the previous method

276    could not identify the relationship between cleavage efficiencies and RNA

277    stability, we found that genes with high $CS_{gene}$ values have short half-lives.

278    These results were obtained due to our improvement of the experimental

279    procedures. Together, these results indicate that cleavage-dependent RNA

280    degradation is important for controlling RNA stability and contributes to

281    various biological processes across eukaryotes.

282

283

284

285

286

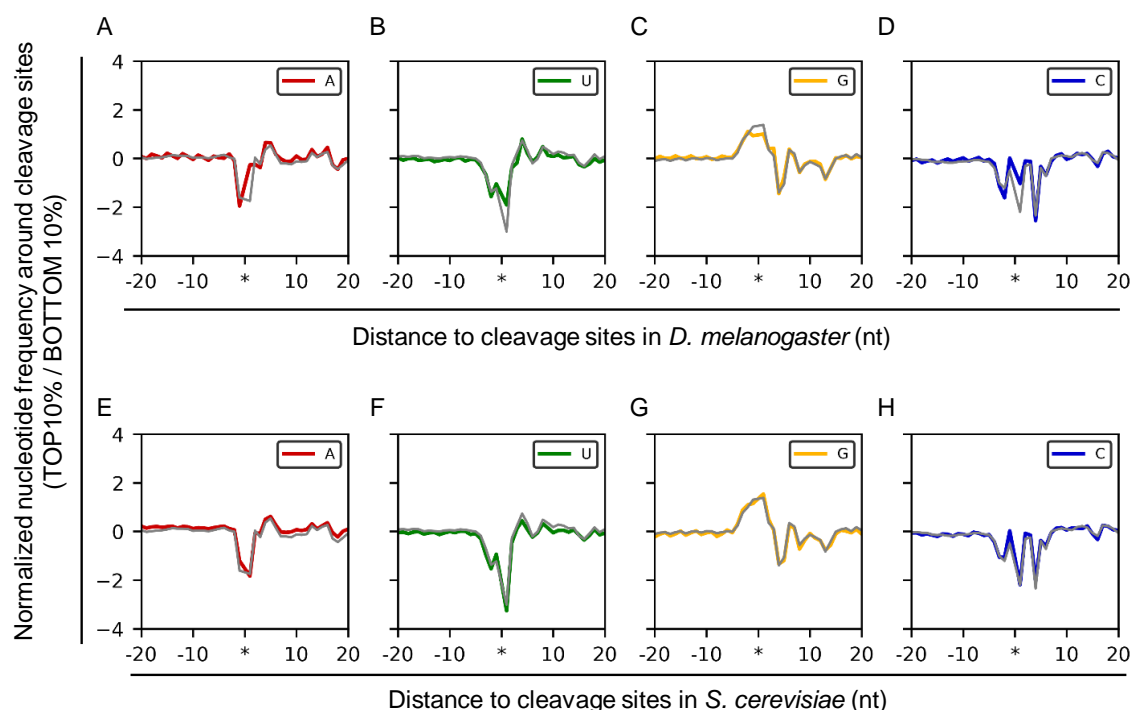287

288

289

290

291

292

293

294

295



296

**Figure 1. Relationship between $CS_{gene}$ values and RNA half-lives**

Scatter plots of half-life vs. $CS_{gene}$ values in *A. thaliana* (A), *D. melanogaster* (B), and *S. cerevisiae* (C). We selected the TOP and BOTTOM 10% of $CS_{gene}$ values and compared half-lives in each species (D–F). Welch's t-test (two-sided) was used to determine statistical significance. The number of genes analyzed in the TOP 10% and BOTTOM 10% were 944 genes in *A. thaliana*, 124 in *D. melanogaster*, and 352 in *S. cerevisiae*. Outliers were removed from box plots.

305

306

12

307 **Conservation of sequence features around cleavage sites**

308 Less information about untranslated regions (UTRs) is available in *D.*

309 *melanogaster* and *S. cerevisiae* than in *A. thaliana*. Therefore, in our

310 subsequent analysis, we compared the sequence features of cleavage sites

311 using CDS regions. When we focused on nucleotide frequencies around the

312 cleavage sites, we observed similar sequence features among species

313 (Figure 2). In addition, sequence features around cleavage sites could be

314 more clearly observed for sequences with high cleavage efficiencies in *D.*

315 *melanogaster* and *S. cerevisiae*, as in *A. thaliana* (Figure 2).

316

317

318

319

320

321

322

323

**Figure 2. Nucleotide frequency around the cleavage sites**

We selected the TOP and BOTTOM 10% cleavage sites based on $CS_{site}$ values and calculated normalized nucleotide frequency ($\log_2$ ratio of nucleotide frequencies between TOP and BOTTOM 10%) in *D. melanogaster* (60,096 sites analyzed in TOP and BOTTOM 10%) and *S. cerevisiae* (61,689 sites analyzed in TOP and BOTTOM 10%). We compared normalized nucleotide frequency around cleavage sites between *A. thaliana* and *D. melanogaster* (A–D) or *S. cerevisiae* (E–H). Gray line indicates nucleotide frequency in *A. thaliana*. Asterisks indicate positions between -1 and +1.

14

341      Subsequently, we discovered sequence motifs using sequences around

342    (from −10 to +10), upstream (from −20 to −5), or downstream (from +5 to +20)

343    of cleavage sites using DREME, and then used STAMP to carry out

344    phylogenetic analyses based on similarities among the discovered motifs

345    (Supplementary Figures S10−S12). We also obtained sequences around,

346    upstream, or downstream of cleavage sites in different plant species (*Lactuca*

347    *sativa* [lettuce], *Oryza sativa* [rice], and *Rosa hybrida* [rose]) reported in a

348    previous study (Ueno et al. 2020), and added sequence motifs for

349    phylogenetic analyses. Generally, similar G-rich sequence motifs were

350    observed around the cleavage sites (from −10 to +10). By contrast, A- or AU-

351    rich motifs were observed upstream or downstream of cleavage sites in plants

352    or yeast, and different sequence motifs were observed in *D. melanogaster*.

353    These results were consistent with sequence conservation in RNA 3'-

354    processing sites, indicating that similar sequences were observed around

355    processing sites, but sequences upstream or downstream of processing sites

356    were not shared among species (Tian and Graber 2012). Specifically,

357    sequences were quite similar between yeast and plants, but differed in

358    animals (Tian and Graber 2012). These results suggest that similar sequence

359    features around cleavage sites are conserved across eukaryotes, and that

360    different sequences in upstream or downstream regions, in addition to highly

361    conserved sequence around cleavage sites, also contribute to control of

362    cleavage efficiencies in different species.

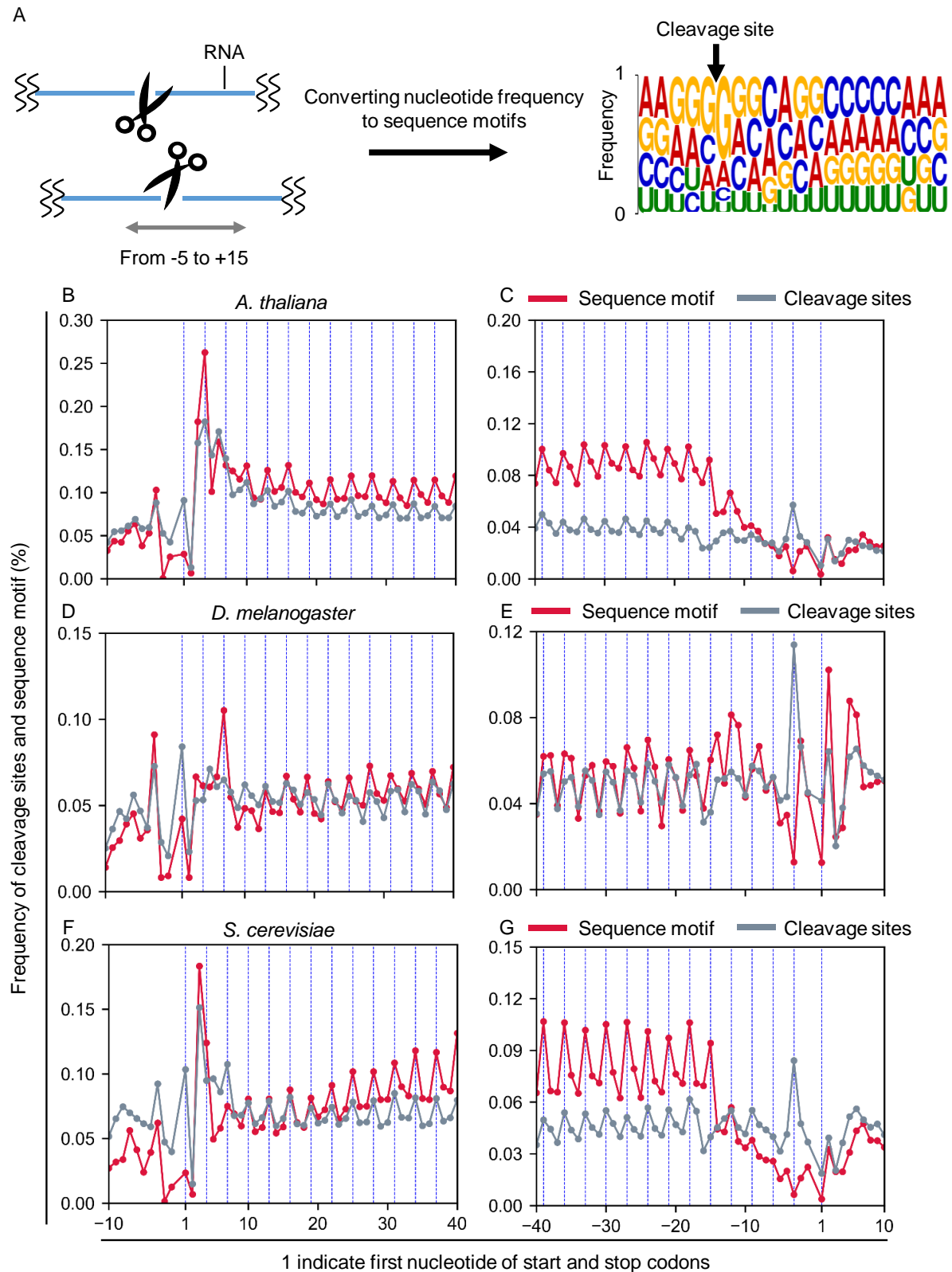363      Because G-rich sequences sometimes induce stable RNA structures,

364    which play important roles in biological process (Lipps and Rhodes 2009;

365    Subramanian et al. 2011; Millevoi et al. 2012), we subsequently focused on

366    the RNA structures around cleavage sites. We predicted paired or unpaired

367    site (dot-bracket notation) around cleavage sites using RNAfold and base-

15

368     paring frequencies at each position were calculated in each species. We

369     observed that base-pairing frequency in TOP 10% $CS_{site}$ values were higher

370     around cleavage sites compared to BOTTOM 10 % $CS_{site}$ values in all species

371     (Supplementary Figure S13A−C). In addition, pairing-frequency was

372     dependent on G-nucleotide frequencies around the cleavage sites

373     (Supplementary Figure S13D−E). These results suggest that RNA structure is

374     related to RNA cleavage, and that these structures are dependent on G-

375     nucleotide frequency around cleavage sites.

376        In our analysis, we observed specific sequence features around cleavage

377     sites and hypothesized that these sequences were important for determining

378     cleavage position. We extracted sequences from −5 to +15 nucleotides

379     around the cleavage sites, and then created sequence motifs (motif letter-

380     probability matrix lines) based on frequencies. Subsequently, we calculated

381     the distribution of sequence motifs in each gene using FIMO and compared it

382     with the distribution of cleavage sites (Figure 3A). If sequence features are

383     important for determining the cleavage position, we expected to observe the

384     same distribution patterns between sequence motifs and cleavage sites. We

385     found that sequence motif had similar three-nucleotide periodicities, and that

386     the phases were almost same (Figure 3B−G). Although previous studies

387     suggested that the three-nucleotide frequency is caused by ribosome

388     movement (Pelechano et al. 2015; Ibrahim et al. 2018; Yu et al. 2016), these

389     results suggest that sequence features are also related to the three-

390     nucleotide periodicity of cleavage sites and are important for determining RNA

391     cleavage position.

392        Together, these results indicate that sequence features around cleavage

393     sites are highly conserved across species and are important for determining

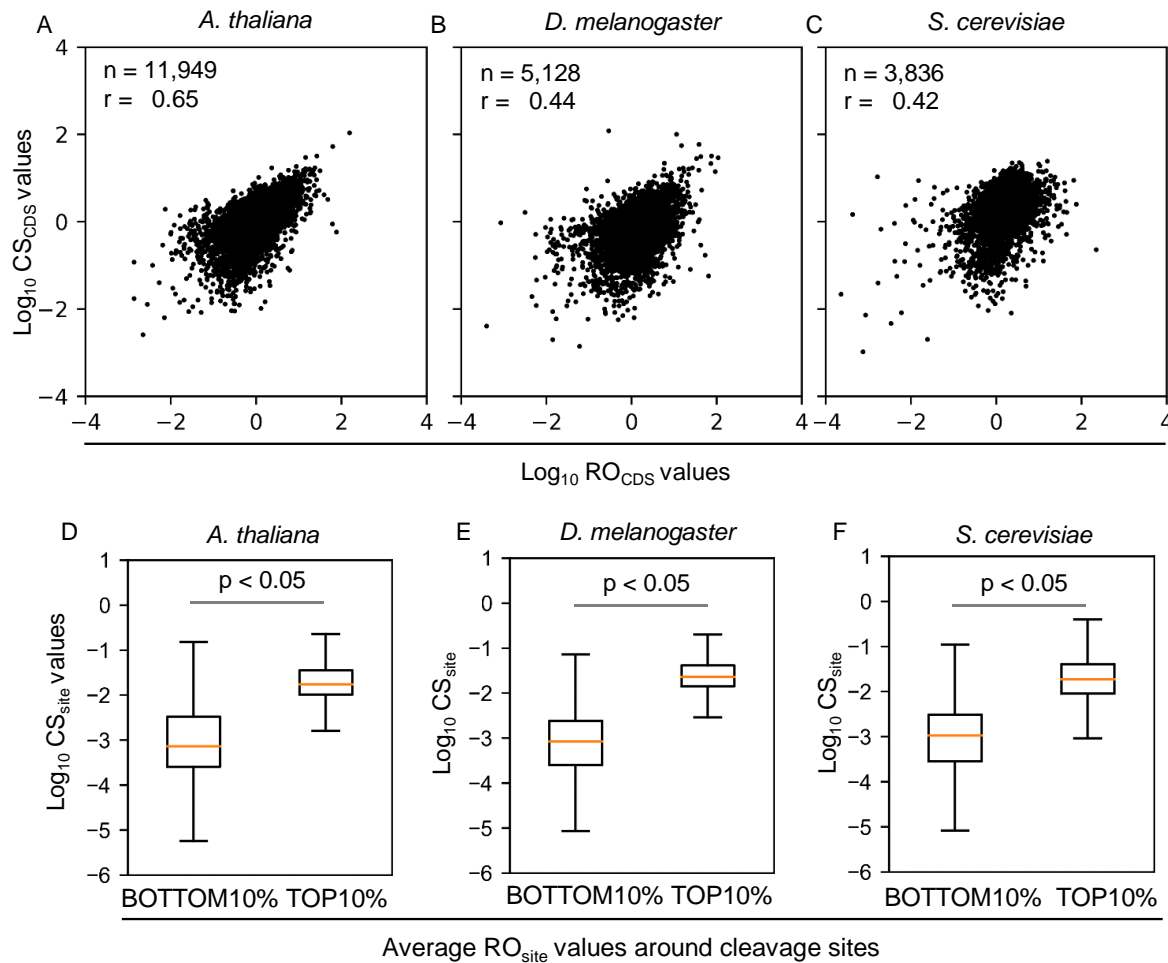394     both cleavage efficiencies and cleavage positions in eukaryotes.

16

399 **Figure 3. Distribution of cleavage sites and sequence motifs around**

400 **cleaved sites**

401 The nucleotide frequency information was converted into sequence motifs

402 (motif letter-probability matrix lines) in MEME (A), and distribution of

403 sequence motif was computed using FIMO. We compared the distribution

404 pattern of cleavage sites and sequence motifs around start (B, D, F) and stop

405 (C, E, G) codons in *A. thaliana* (B, C), *D. melanogaster* (D, E), and *S.*

406 *cerevisiae* (F, G). Because most UTR information in *S. cerevisiae* is not yet

407 defined, we expanded the annotation region of *S. cerevisiae* only in this

408 analysis.

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426 **Effects of the translation process on cleavage position and cleavage**

427 **efficiency**

428 Previous studies in yeast reported that changes in ribosome occupancy and

429 degradation efficiencies at the gene level were positively correlated between

430 oxidative stress and normal conditions (Pelechano et al. 2015). Hence, we

431 analyzed the relationship between cleavage efficiencies and ribosome

432 occupancy at the gene level among species. In addition, we analyzed the

433 effect of ribosome occupancy on cleavage efficiency at the site level. For *A.*

434 *thaliana*, we conducted ribosome profiling using the same conditions as for

435 our previous TREseq data. For *D. melanogaster* and *S. cerevisiae*, we

436 obtained ribosome profiling data from previous studies (Luo et al. 2018;

437 Gerashchenko and Gladyshev 2017). We defined the number of RPFs at each

438 site, normalized against RNA abundance, as the $RO_{site}$ values, and the sum

439 of $RO_{site}$ values at the CDS regions of each gene as the $RO_{CDS}$ values.

440 Similarly, we defined the sum of $CS_{site}$ values in the CDS region as the $CS_{CDS}$

441 values. At the gene level, we observed a positive correlation between $CS_{CDS}$

442 and $RO_{CDS}$ values among species (Figure 4). These results suggest that the

443 translation process has a positive effect on cleavage efficiencies at the gene

444 level.

445

446

447

448

449

450

451

452

19

**Figure 4. Relationship between ribosome occupancy and cleavage efficiencies**

Scatter plots of $CS_{CDS}$ and $RO_{CDS}$ values in *A. thaliana* (A), *D. melanogaster* (B), and *S. cerevisiae* (C). Both $CS_{CDS}$ and $RO_{CDS}$ values were normalized against RNA length. Average ribosome occupancy around cleavage sites in the CDS region was calculated and compared against $CS_{site}$ values in the TOP 10% (high ribosome occupancy) and BOTTOM 10% (low ribosome occupancy) (D–F). Welch's t-test (two-sided) was used to determine statistical significance. The number of genes analyzed in the TOP 10% and BOTTOM 10% were 177,248 in *A. thaliana*, 60,096 in *D. melanogaster*, and 61,689 in *S. cerevisiae*. Outliers were removed from box plots.

466     Next, we focused on the relationship between $CS_{site}$ values and

467     ribosome occupancies at the site level. We calculated the average $RO_{site}$

468     value around the cleavage sites (±50 nucleotides), and selected the TOP or

469     BOTTOM 10% based on occupancy around the cleavage site (Figure 4). Even

470     at the site level, cleavage sites with high ribosome occupancy (TOP 10%) had

471     higher cleavage efficiencies than the BOTTOM 10% across species. These

472     results suggested that ribosome occupancy has a positive effect on cleavage

473     efficiencies at the gene or site level. Although correlations between ribosome

474     occupancy and degradation efficiency were weaker under normal conditions

475     compared to correlations between changes in ribosome occupancy and

476     degradation efficiency between under normal condition and oxidative stress

477     in yeast (Pelechano et al. 2015), we could observe a positive correlation

478     among all species examined. Because the previous study conducted poly(A)

479     selection, resulting in a bias in detection level dependent on RNA length

480     (Weinberg et al. 2016), the correlation appeared to be lower than that

481     detected by TREseq analysis.

482     Because previous study suggested that ribosome movement is related to

483     cleavage position (5'-truncated RNA ends) (Pelechano et al. 2015; Yu et al.

484     2016; Ibrahim et al. 2018), we also compared the distributions of ribosomes

485     and cleavage positions (Figure 5). We observed a three-nucleotide periodicity

486     and found that both patterns were matched in some regions. However, the

487     phases were generally not matched between cleavage and ribosome positions

488     across different species. These results were consistent with a previous study

489     showing that ribosome and cleavage positions (5'-truncated RNA ends) have

490     three-nucleotide periodicities, but the phases are generally not matched

491     (Ibrahim et al. 2018; Pelechano et al. 2015). These results suggest that

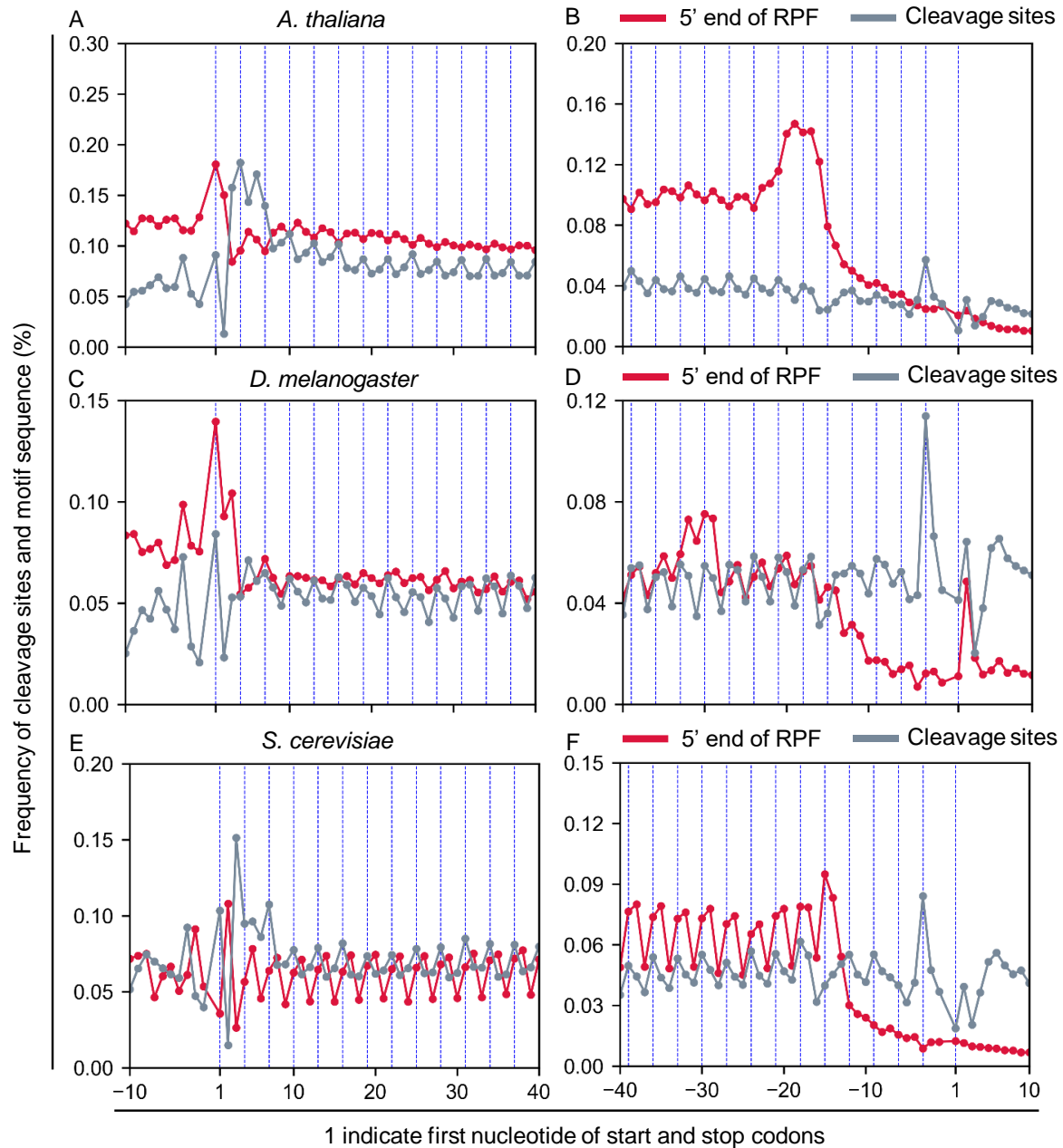492     ribosome position has effects on cleavage position (5'-truncated RNA ends)

493    at a certain distance (Pelechano et al. 2015) or is not a major determinant for

494    cleavage position. Together, these findings indicate that ribosome

495    occupancies are important for cleavage efficiencies, but the effect of

496    ribosome movement on cleavage position cannot be clarified solely by

497    comparing the distribution patterns of ribosome and cleavage position.

498

499

500

501

**Figure 5. Comparing distributions of ribosome and cleavage positions**

Distributions of 5' ends of RPFs (ribosome protected fragments) and cleavage sites around start (A, C, E) and stop (B, D, F) codons in *A. thaliana* (A, B), *D. melanogaster* (C, D), and *S. cerevisiae* (E, F). Because most UTR information in *S. cerevisiae* is not yet defined, we expanded the annotation region of *S. cerevisiae* only in this analysis.
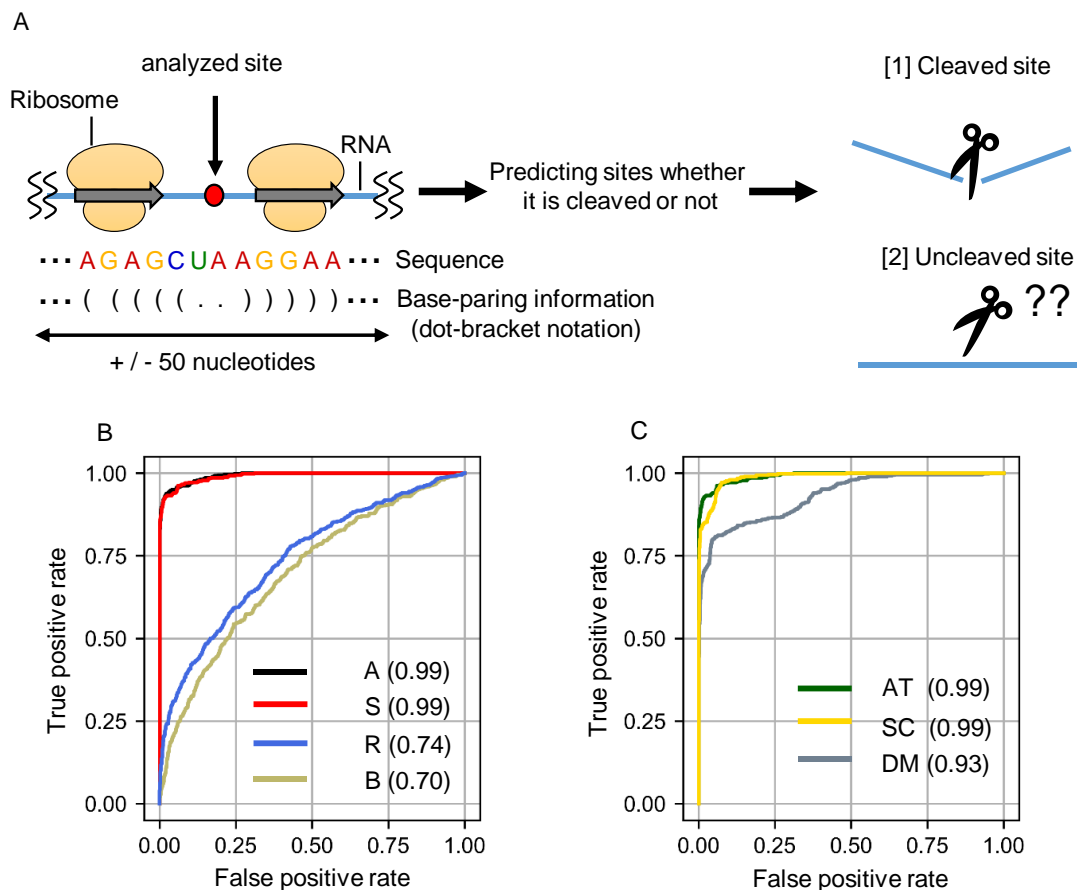
**Feature selection for determining cleavage position**

511

512 Previous studies suggest that ribosome position is important for

513 determination of 5'-truncated RNA ends (cleavage position) (Pelechano et al.

514 2015; Yu et al. 2016; Ibrahim et al. 2018). However, less information is

515 available about the relationship between ribosome and cleavage positions,

516 and we did not detect a strong effect of ribosome movement on cleavage

517 position in this analysis. To evaluate the influence of sequences or ribosome

518 movement on cleavage position, we conducted feature selection using a

519 random forest (RF) classifier, which allows us to evaluate several features

520 (determinants) according to the coefficient in the model. We selected sites

521 and predicted whether they were cleaved (Figure 6A). To simplify the

522 interpretation, we selected the site with the highest $CS_{site}$ value in each gene,

523 as well as sites with no $CS_{site}$ values. For our model, we used nucleotide

524 frequency, secondary structure, and ribosome position or ribosome

525 occupancy information (Supplementary Figure S14). Ribosome position

526 indicates only positional information, whereas ribosome occupancy includes

527 occupancy and positional information. To evaluate the effect of secondary

528 structure (base-paring information), we predicted paired or unpaired sites

529 (dot-bracket notation) around cleavage sites using RNAfold. We divided the

530 data into training and test sets. The model was constructed using the training

531 data, and model performance was evaluated using the test data

532 (Supplementary Figure S15). To estimate prediction accuracy in the test data,

533 we used ROC curve and AUC values (Singh et al. 2017; Lloréns-Rico et al.

534 2015). ROC curves indicate the true-positive and false-positive rates,

535 whereas AUC values indicate the area under the ROC curve. AUC values can

536 range from 0.5 (i.e., the method does not perform better than random) and 1

537 (the method classifies all samples perfectly with no misprediction) (Lloréns-

538     Rico et al. 2015). Initially, we focused on *A. thaliana* and made four models:

539     [1] all features, [2] only nucleotide frequency around cleavage sites, [3] only

540     secondary structure around cleavage sites, and [4] only ribosome position

541     and ribosome occupancy around the cleavage sites. AUC value (prediction

542     accuracy) was highest in classification model that used all features (Figure

543     6B). In addition, when we used only nucleotide frequencies, the prediction

544     accuracy was still high. By contrast, when we used only secondary structures

545     or ribosome position/occupancy information, prediction accuracy was lower.

546     These results were consistent with the Gini importance, which indicate the

547     importance of features for predicting whether site is cleaved or not, in RF

548     classifier using all features (Supplementary Figure S16). These tendencies

549     were also observed in *D. melanogaster* or *S. cerevisiae* (Supplementary

550     Figures S17–S19). Although ribosome movement also seemed to be related

551     to determination of cleavage position based on the Gini importance, these

552     results suggest that sequence features are the major determinant of cleavage

553     position across eukaryotes.

554       We hypothesized that prediction accuracy would still be high in *D.*

555     *melanogaster* and *S. cerevisiae* using a model that was constructed in *A.*

556     *thaliana* using only sequence information, so long as similar sequence

557     features were important for determining cleavage position across eukaryotes.

558     As shown in Figure 6C, AUC values were > 0.9 in *D. melanogaster* and *S.*

559     *cerevisiae*. Consistent with the sequence motif analysis, prediction accuracy

560     using models constructed in *A. thaliana* were higher in *S. cerevisiae* than in

561     *D. melanogaster*. These results suggest that common sequence features

562     involved in determining cleavage position were highly conserved among

563     eukaryotes, but that species-specific sequences also influence RNA cleavage.

564

**Figure 6. Feature selection using RF classifier**

Using information on ribosome occupancy, ribosome position, sequence, and secondary structure (base-pairing information) around cleavage sites, we predicted whether each site was cleaved (A). Data were divided into training and test data sets. We constructed models using the training data and evaluated the accuracy using the test data. Four models were constructed: [A] all features (sequence, ribosome position and ribosome occupancy, base-pairing), [S] only sequence, [R] only ribosome position or ribosome occupancy, and [B] only base-pairing information around the cleavage sites. Models were evaluated using ROC curves and AUC values (B). Parentheses indicate AUC values in each model. We predicted cleaved sites in *D. melanogaster* and *S. cerevisiae* using a model constructed exclusively from *A. thaliana* sequence (C). ROC curves and AUC values in *A. thaliana* are the same as in Figure 6B. AT, *A. thaliana*; SC, *S. cerevisiae*; DM, *D. melanogaster*.

26

579 **DISCUSSION**

580 **A High G nucleotide frequencies is a strong indicator for cleavage sites**

581 In TREseq, we observed high G-nucleotide frequency around cleavage sites.

582 Although these sequence features around cleavage sites were not clearly

583 observed in previous studies, weak trends were reported in previous

584 degradome analyses. In Akron-seq, higher G-nucleotide frequency (defined

585 as G-quadruplex) was observed downstream of the cleavage site in mammals,

586 and similar sequence features were also identified by PAREseq or 5PSeq

587 (Ibrahim et al. 2018). In particular, Akron-seq observed the expected G-

588 quadruplex structures in highly conserved regions in vertebrates.

589 Previous degradome sequencing methods involved ligation of single-

590 stranded adapters using T4 RNA ligase, making it difficult to detect 5'-

591 truncated RNAs with extensive base-pairing (Zhuang et al. 2012). This result

592 is consistent with a previous degradome analysis in which open structures

593 were observed around cleavage sites (Ibrahim et al. 2018). By contrast,

594 TREseq (CAGE) methods use ligation of different adapters, and open

595 structures were not observed around the cleavage sites (Supplementary

596 Figure S13). Therefore, G-rich sequences were observed around cleavage

597 sites, but not downstream of them. Also, because poly (A) RNA selection

598 caused the detected reads or sites to be enriched in 3' UTRs or around stop

599 codons (Weinberg et al. 2016), evaluation of cleavage efficiencies differed

600 from previous degradome analysis. Because of these improvements in the

601 experimental procedures, specific sequence features around cleavage sites

602 were clearly observed in TREseq. Although we cannot exclude the possibility

603 that G-quadruplex structures affect cleavage efficiencies, RF classifier

604 analyses indicated that the influence of RNA structure was weak (Figure 6).

605 Considering that cleavage sites with higher G-nucleotide frequencies were

27

606 detected by different methods in different species, we believe that these

607 sequence features are actually present around the cleavage sites, thus

608 indicting that sequence-specific RNA degradation mechanisms are highly

609 conserved in eukaryotes.

610 　　Although we focused on uncapped cleavage sites in this study, some

611 cleavage sites are secondarily capped in the cytosol (Schoenberg and Maquat

612 2009). Indeed, a recapping enzyme has been identified in animals, and some

613 cleaved RNAs that originate from no-go decay (NGD) are actually re-capped

614 (Schoenberg and Maquat 2009). However, because only the Cap RNA library

615 was used in CAGE method, the ratio of Cap to Cap-less RNA could not be

616 calculated. By contrast, we made Cap RNA and Cap-less RNA libraries in

617 TREseq, and can roughly estimate the proportion of re-capped RNA from

618 cleavage site data by comparing libraries. We expanded the annotation region

619 of the Cap RNA data following data processing in TREseq, and calculated the

620 expected re-capped RNA in cleavage site data using exon CDS regions from

621 different species (Supplementary Figure S20). Some cleavage sites in *A.*

622 *thaliana* appeared to be re-capped RNA, suggesting the existence of re-

623 capped cleavage sites in plants. However, the proportion of re-capped

624 cleavage sites was less than 10% in *A. thaliana, D. melanogaster*, and *S.*

625 *cerevisiae* (Supplementary Figure S20). These results are consistent with a

626 previous study reporting that thousands of cleavage sites in mammals are

627 Cap-less RNA (Malka et al. 2017), suggesting that the majority of cleavage

628 sites detected in genome-wide analysis are uncapped.

629

630 **Ribosome positions have little effect on RNA cleavage position**

631 A previous study reported that ribosome movement is important for RNA

632 stability (Presnyak et al. 2015). In particular, slow ribosome movement

633 (ribosome stalling or pausing) causes endonucleolytic cleavage (Doma and

634 Parker 2006; Simms et al. 2017). Although we do not know whether similar

635 trans-acting factors in NGD are related to cleavage sites detected in genome-

636 wide degradome analysis, we can refer to the relationships between ribosome

637 movement and cleavage sites. NGD can be induced by inserting several

638 codons into a reporter gene, and the resultant cleavage sites can be detected

639 in yeast (Simms et al. 2017; Doma and Parker 2006). These results suggest

640 that the translation process is related to cleavage probability.

641 Initially, these RNA cleavages were thought to occur near sites of

642 ribosome stalling or pausing detected in the NGD analysis (Doma and Parker

643 2006). However, RNA cleavage sites were widely identified in the gene,

644 including beside the ribosomes based on the analysis in recent study (Simms

645 et al. 2017). These results suggest that ribosome movement is important for

646 RNA cleavage, but not critical for determining cleavage position. In our

647 analyses, we cannot observe a strong relationship between ribosome and

648 cleavage position. Specifically, we observed that distribution patterns

649 between ribosome and cleavage sites have three-nucleotide periodicities, but

650 the phases were generally not matched (Figure 5). These results were

651 observed not only in TREseq analysis but also in another previous study

652 (Ibrahim et al. 2018; Pelechano et al. 2015). Similar results were observed in

653 feature selection using an RF classifier model in *A. thaliana, D. melanogaster,*

654 and *S. cerevisiae*: ribosome movement seemed to be related to cleavage

655 position, but was not a major determinant (Figure 6), suggesting that

656 ribosomes occupancy has a positive effect on cleavage efficiencies, but

657 ribosome position does not have a strong effect on cleavage position.

658 Because codon frequencies were generally the same in each translation frame

659 in the transcripts, we could observe the three-nucleotide periodicity of

29

660    nucleotide frequencies in CDS regions in *A. thaliana, D. melanogaster,* and

661    *S. cerevisiae* (Supplementary Figure S21 and S22) and these periodicities

662    appear to form the distribution patterns of cleavage sites.

663        Previous study reported that these periodicities in yeast and plant are

664    caused by exonucleolytic digestion because the amplitude of periodicity was

665    decreased in exonuclease mutant (Pelechano et al. 2015; Yu et al. 2016).

666    However, even in the exonuclease mutant, the periodicities were still

667    observed in animal (Ibrahim et al. 2018). In addition, disappearance of

668    periodicity appeared to be caused by low coverage of detected 5' degradation

669    intermediates and our previous study revealed that periodicity of 5'

670    degradation intermediates was almost identical in WT and exonuclease

671    mutants in *A. thaliana* (Ueno et al. 2020). Also, we obtained mapping data of

672    5Pseq in yeast (Pelechano et al. 2015) and reanalyzed whether the periodicity

673    disappeared around start or stop codon in exonuclease mutant

674    (Supplementary Figure S23). Because number of detected reads (read count)

675    at each 5' degradation intermediates were sometimes dependent on its RNA

676    accumulation levels, there was a possibility that the distribution pattern based

677    on the read count did not reflect the genome-wide tendency due to highly

678    expressed genes. Therefore, we calculated the number of 5' degradation

679    intermediates (detected sites in 5Pseq) at each RNA position by count and

680    frequency (relative to all detected 5' degradation intermediates). The

681    amplitude of periodicity was slightly decreased around stop codon in

682    exonuclease mutant (Supplementary Figure S23B and D), but the periodicity

683    around start codon was almost similar or more clearly observed in

684    exonuclease mutant compared to WT (Supplementary Figure S23A and C).

685    These results suggest that exonuclease is not major determinant for forming

686    three-nucleotide periodicity in yeast. Considering that periodicity was still

30

687     observed in exonuclease mutant in plant (Ueno et al. 2020), animal (Ibrahim

688     et al. 2018) and even in yeast (Supplementary Figure S23), periodicity appear

689     to be originated from endonucleolytic cleavage–dependent RNA degradation.

690         In summary, using TREseq, we observed that cleavage efficiencies play

691     important roles in RNA stability in various species. We observed that

692     sequence features and ribosome occupancy had a positive effect on cleavage

693     efficiencies, and that sequence features were more important than ribosome

694     movement for determining RNA cleavage positions across eukaryotes. These

695     results suggest that similar sequence-specific endonucleolytic cleavage

696     mechanisms are conserved across species, and that these mechanisms

697     contribute to the control of RNA stability among eukaryotes. On the other

698     hand, given that different sequences also seemed to be related to RNA

699     cleavage, based on the results of the sequence motifs and RF classifier

700     analyses, further studies are needed to elucidate endonucleolytic cleavage–

701     dependent RNA degradation in eukaryotes.

702

703 **MATERIAL AND METHODS**

704 **Growth conditions for *A. thaliana***

705 *A. thaliana* T87 cell suspension was cultured in modified Murashige–Skoog

706 medium as described previously (Matsui et al. 2011).

707

708 **Growth conditions for *D. melanogaster***

709 S2-R+ cells (Drosophila Genomics Resource Center #150) were grown in

710 Schneider's medium supplemented with 10% FBS and 100 U/ml penicillin–

711 streptomycin at 25°C.

712

713 **Growth conditions for *S. cerevisiae***

714 *S. cerevisiae* strain Σ1278b was cultured at 30°C to exponential growth phase

715 (OD$_{600}$ of 1.0) in YPD medium containing 2% (wt/vol) glucose, 1% (wt/vol)

716 Difco Bacto yeast extract (Thermo Fisher Scientific, Waltham, MA, USA), and

717 2% (wt/vol) Difco Bacto peptone (Thermo Fisher Scientific).

718

719 **Library construction for truncated RNA-end sequencing (TREseq)**

720 TREseq libraries were constructed using RNA from *D. melanogaster* and *S.*

721 *cerevisiae* as described previously (Ueno et al. 2018, 2020). In brief, RNA

722 was isolated and purified, and rRNA was depleted using the Ribo-Zero rRNA

723 Removal Kit (Human/Mouse/Rat) (Illumina, San Diego, CA, USA). After rRNA

724 depletion, total RNA was reverse-transcribed with a random primer, and RNA–

725 cDNA hybrids were separated into Cap RNA–cDNA hybrid (Cap RNA) and

726 Cap-less RNA–cDNA hybrid (Cap-less RNA) fractions using the Cap-trapping

727 method. Subsequently, we used different 5' adapter sequences in Cap and

728 Cap-less RNA (Table 1), followed by library preparation for nAnT-iCAGE

729 (Murata et al. 2014; Adiconis et al. 2018). We used Cap-less RNA to detect

32

730  5' degradation intermediates and Cap RNA to estimate transcription start site

731  (TSS) and RNA abundance using more than 50 reads for each gene. The

732  libraries were sequenced on a NextSeq 500 instrument (Illumina).

733      For the TREseq library without Cap-trapping method in *A. thaliana,* we

734  did not conduct Cap-trapping method in library preparation and used 5'

735  adapter sequence similar to Cap-less RNA library (Table 1). The libraries were

736  subjected to sequencing on a NextSeq 500 instrument (Illumina).

737

738

739

740

741  Table.1. 5' adapter sequence used in TREseq

| 5' adapter | Sequences |
|---|---|
| Cap RNA | 5' ACACGACGCTCTTCCGATCT(G/N)NNNNN-P 3'<br><br>3' P-TGTGCTGCGAGAAGGCTAGA-P 5' |
| Cap-less RNA | 5' ACACGACGCTCTTCCGATCTNNNNNN-P 3'<br><br>3' P-TGTGCTGCGAGAAGGCTAGA-P 5' |

742

743

744

745

746

747

748

749

750

751

33

**Data processing for TREseq**

We used a modified MOIRAI system to process TREseq data (Hasegawa et al. 2014; Ueno et al. 2018, 2020). For calculation of 5' degradation intermediates (Cap-less RNA), low-quality reads or reads originating from rRNA were removed. The remaining reads were mapped to the TAIR version 10 reference genome (www.arabidopsis.org), Flybase version 6.31 reference genome (https://flybase.org/), or Sigma1278b_MIT_2009_ACVY01000000 reference genome (https://www.yeastgenome.org/) using HISAT2. After converting BAM files to SAM files, we removed reads with mismatches within three nucleotides of the 5' end of a mapped read, except for TREseq library without Cap-trapping method. Because the 5' cap (7-methylguanosine) is not encoded by the genome, most Cap RNAs in the cleavage sites data were removed. Subsequently, the SAM files were converted to BED files, which were used to count only the first nucleotide (5' end) of each read. Because we used different mapping software than in previous TREseq analysis, data processing for *A. thaliana* is shown in the figures (Supplementary Figure S1). For calculation of TSS or RNA abundance (Cap RNA), the data were processed by CAGE (Murata et al. 2014; Adiconis et al. 2018). In the library preparation step, we could not completely exclude contamination (Cap RNA in Cap-less RNA library or Cap-less RNA in Cap RNA library). Hence, we removed Cap RNA from Cap-less RNA data and vice versa by comparing each library, as described in previous study (Ueno et al. 2020) (Supplementary Figures S6-S8). De-capped RNA was removed from Cap-less RNA data by comparing information for Cap and Cap-less RNA. As an indicator of cleavage efficiencies at each site, we defined the reads at each cleavage site normalized by RNA abundance as the cleavage score at the site level ($CS_{site}$), as described in previous studies (Ueno et al. 2020, 2018). For the gene level,

34

779    we defined the total $CS_{site}$ value at each gene or each CDS region as the

780    $CS_{gene}$ or $CS_{CDS}$ value, respectively.

781

782    **Library construction for ribosome profiling**

783    We selected ribosome-protected fragments (RPFs) as previously described

784    (Lei et al. 2015; Yamasaki et al. 2015). In brief, *A. thaliana* T87 cells were

785    harvested 3 days after inoculation and frozen in liquid nitrogen, followed by

786    homogenization with extraction buffer (200 mM Tris-HCl, pH8.5, 50 mM KCl,

787    25 mM $MgCl_2$, 2 mM EGTA, 100 µg/ml heparin, 100 µg/ml cycloheximide, 2%

788    polyoxyethylene 10-tridecyl ether, 1% sodium deoxycholate), and

789    centrifugation at 15,000 × g for 10 min at 4°C (Lei et al. 2015). We added 6

790    µl of RNase I (Thermo Fisher Scientific) for 30 min, and stopped the reaction

791    by adding 10 µl of RNase inhibitor (Thermo Fisher Scientific). Using a 26.25–

792    71.25% sucrose density gradient buffer (200 mM Tris-HCl, pH 8.5, 200 mM

793    KCl, 200 mM $MgCl_2$), monosomes were collected by sucrose density gradient

794    centrifugation at 55,000 r.p.m. for 50 min at 4°C in an SW55 rotor (Beckman

795    Coulter, Brea, CA, USA). After isolation of monosomes, RPFs were purified

796    using the TruSeq Ribo Profile kit (Illumina). The libraries were sequenced on

797    a NextSeq 500.

798

799    **Data processing for ribosome profiling**

800    Adapter sequences were trimmed (Luo et al. 2018; Gerashchenko and

801    Gladyshev 2017) and reads were mapped using the modified MOIRAI system

802    (Ueno et al. 2018, 2020). The remaining reads were mapped to the TAIR

803    version 10 reference genome (www.arabidopsis.org), Flybase version 6.31

804    reference genome (https://flybase.org/), or

805    Sigma1278b_MIT_2009_ACVY01000000 reference genome

806    (https://www.yeastgenome.org/) as appropriate, using BWA. After mapping,

807    we counted the first nucleotide (5' end) of each read using the BED files. As

808    an indicator of RPFs at each site, we defined the average of reads at each 5'

809    RPF, normalized against RNA abundance, as ribosome occupancy at the site

810    level ($RO_{site}$). For the gene level, we defined the total $RO_{site}$ value in the CDS

811    region at each gene as the $RO_{CDS}$ value.

812

813    **Library construction for Nanopore sequencing**

814    Total RNA was extracted using Trizol (Thermo Fisher Scientific), and mRNA

815    was isolated from total RNA using the Magnosphere UltraPure mRNA

816    Purification Kit (Takara Bio, Kusatsu, Shiga, Japan). Nanopore libraries were

817    prepared from poly(A)$^+$ RNA using the Nanopore PCR cDNA Sequencing Kit

818    (Oxford Nanopore Technologies, Oxford, UK). We conducted 12 cycles of PCR

819    amplification. RNA sequencing on the GridION platform (Oxford Nanopore

820    Technologies) was performed for 48 h on ONT R9.4 flow cells using a

821    minknow-core-gridion 3.1.8 software (Oxford Nanopore Technologies).

822

823    **Data processing for Nanopore profiling**

824    Reads were base-called with ont-guppy version 2.0.8-1 (Oxford Nanopore

825    Technologies). Fully sequenced reads were identified and oriented using

826    Pychopper2 (Oxford Nanopore Technologies). Reads were mapped to the

827    Arabidopsis TAIR10 genome using Minimap2. Subsequently, the SAM files

828    were converted to BED files, which were used to count only the first

829    nucleotide (5' end) of each read.

830

831    ***In vitro* synthesized RNA and quantitative RT-PCR (qRT-PCR) analysis**

832    *In vitro* synthesis of uncapped *Renilla* luciferase (*R-luc)* from plasmid pT3-

36

833   RL-pA was performed as described previously (Matsuura et al. 2008). We

834   added 1 ng of uncapped *R-luc* RNA to 5 µl of total RNA and conducted rRNA

835   depletion and reverse transcription as described for Cap-less RNA library

836   preparation in TREseq. We conducted qRT-PCR on a LightCycler 480 (Roche

837   Applied Science, Upper Bavaria, Germany). Gene-specific primers used for

838   qRT–PCR of *R-luc* RNA were 5'- GGATTCTTTTCCAATGCTATTGTT-3' and 5'-

839   AAGACCTTTTACTTTGACAAATTCAGT-3' (Matsuura et al. 2008).

840

841   **Data analysis**

842   For motif analysis, sequence motifs adjacent to the cleavage sites were

843   identified using MEME motif analysis tools (http://meme-suite.org). Motif

844   distribution was calculated using FIMO.

845       As for GO enrichment analysis, we used the Gene Ontology enRIchment

846   anaLysis   and   visuaLizAtion   tool   (GORILLA  )   (http://cbl-

847   gorilla.cs.technion.ac.il/).

848

849   **Construction of the random forest (RF) classifier**

850   For feature selection, we used the RandomForestClassifier from the Python

851   package scikit-learn. We defined the highest $CS_{site}$ value in each gene as the

852   cleavage site and selected a total of 1000 genes. For non-cleaved sites, we

853   selected sites from the selected 1000 genes with no $CS_{site}$ value at least

854   around +/- 30 nucleotide. We divided data into training (500 genes) and test

855   data (500 genes) (Table 2). In the models, we used three features: [1]

856   sequence, [2] secondary structure, and [3] ribosome position and ribosome

857   occupancy. We converted sequences or secondary structure information into

858   numeric numbers for the RF classifier. We optimized parameters of the RF

859   classifier based on the out-of-the-bag (OOB) score in the training data and

37

860  receiver operating characteristic (ROC) curves and area under the curve

861  (AUC) values in the test data, as described in previous studies (Lloréns-Rico

862  et al. 2015; Singh et al. 2017). The ROC curve specifies the true positive rate

863  and false-positive rate; a high AUC value indicates high prediction accuracy.

864  We constructed an RF classifier using the training data and predicted

865  accuracy using the test data. We also predicted cleavage sites in *D.*

866  *melanogaster* and *S. cerevisiae* using the RF classifier constructed in *A.*

867  *thaliana*.

868

869

870

871

872

873  Table 2. Sites analyzed in RF classifier model

| | *A. thaliana* | | *D. melanogaster* | | *S. cerevisiae* | |
|---|---|---|---|---|---|---|
| | Training | Test | Training | Test | Training | Test |
| Cleaved sites | 500 | 500 | 500 | 500 | 500 | 500 |
| Uncleaved sites | 1926 | 1732 | 3675 | 3327 | 2321 | 2314 |

874

875

876

38

877 **DATA ACCESS**

878 TREseq reads are available in the DDBJ Sequence Read Archive (DRA)

879 database with accession numbers DRA010805 (*Drosophila melanogaster*),

880 DRA010806 (*Saccharomyces cerevisiae*) and DRA010923 (*Arabidopsis*

881 *thaliana* without Cap-trapping method). Ribosome profiling and Nanopore

882 sequencing reads in *A. thaliana* are available with accession numbers

883 DRA010802 and DRA010803, respectively.

884

885 **COMPETING INTEREST STATEMENT**

886 The authors declare no competing interests.

887

888 **ACKNOWLEDGMENTS**

894

895

896

897

898

899 **REFERENCES**

900 Addo-Quaye C, Eshoo TW, Bartel DP, Axtell MJ. 2008. Endogenous siRNA

901 and miRNA Targets Identified by Sequencing of the Arabidopsis

902 Degradome. *Curr Biol* **18**: 758–762.

903 Adiconis X, Haber AL, Simmons SK, Levy Moonshine A, Ji Z, Busby MA, Shi

904 X, Jacques J, Lancaster MA, Pan JQ, et al. 2018. Comprehensive

905 comparative analysis of 5′-end RNA-sequencing methods. *Nat Methods*

906 **15**: 505–511.

907 Anderson SJ, Kramer MC, Gosai SJ, Yu X, Vandivier LE, Nelson ADL,

908 Anderson ZD, Beilstein MA, Fray RG, Lyons E, et al. 2018. N6-

909 Methyladenosine Inhibits Local Ribonucleolytic Cleavage to Stabilize

910 mRNAs in Arabidopsis. *Cell Rep* **25**: 1146–1157.

911 Bracken CP, Szubert JM, Mercer TR, Dinger ME, Thomson DW, Mattick JS,

912 Michael MZ, Goodall GJ. 2011. Global analysis of the mammalian RNA

913 degradome reveals widespread miRNA-dependent and miRNA-

914 independent endonucleolytic cleavage. *Nucleic Acids Res* **39**: 5658–

915 5668.

916 Burow DA, Martin S, Quail JF, Alhusaini N, Coller J, Cleary MD. 2018.

917 Attenuated Codon Optimality Contributes to Neural-Specific mRNA

918 Decay in Drosophila. *Cell Rep* **24**: 1704–1712.

919 Carninci P, Sandelin A, Lenhard B, Katayama S, Shimokawa K, Ponjavic J,

920 Semple CAM, Taylor MS, Engström PG, Frith MC, et al. 2006. Genome-

921 wide analysis of mammalian promoter architecture and evolution. *Nat*

922 *Genet* **38**: 626–635.

923 Chiba Y, Green PJ. 2009. mRNA degradation machinery in plants. *J Plant*

924 *Biol* **52**: 114–124.

925 De Rie D, Abugessaisa I, Alam T, Arner E, Arner P, Ashoor H, Åström G,

40

926      Babina M, Bertin N, Burroughs AM, et al. 2017. An integrated

927      expression atlas of miRNAs and their promoters in human and mouse.

928      *Nat Biotechnol* **35**: 872–878.

929   Doma MK, Parker R. 2006. Endonucleolytic cleavage of eukaryotic mRNAs

930      with stalls in translation elongation. *Nature* **440**: 561–564.

931   Gerashchenko M V., Gladyshev VN. 2017. Ribonuclease selection for

932      ribosome profiling. *Nucleic Acids Res* **45**: e6.

933   German MA, Pillay M, Jeong DH, Hetawal A, Luo S, Janardhanan P, Kannan

934      V, Rymarquis LA, Nobuta K, German R, et al. 2008. Global identification

935      of microRNA-target RNA pairs by parallel analysis of RNA ends. *Nat*

936      *Biotechnol* **26**: 941–946.

937   Gregory BD, O'Malley RC, Lister R, Urich MA, Tonti-Filippini J, Chen H,

938      Millar AH, Ecker JR. 2008. A Link between RNA Metabolism and

939      Silencing Affecting Arabidopsis Development. *Dev Cell* **14**: 854–866.

940   Harigaya Y, Parker R. 2012. Global analysis of mRNA decay intermediates

941      in Saccharomyces cerevisiae. *Proc Natl Acad Sci U S A* **109**: 11764–

942      11769.

943   Hasegawa A, Daub C, Carninci P, Hayashizaki Y, Lassmann T. 2014.

944      MOIRAI: A compact workflow system for CAGE analysis. *BMC*

945      *Bioinformatics* **15**: 144.

946   Hou CY, Lee WC, Chou HC, Chen AP, Chou SJ, Chen HM. 2016. Global

947      analysis of truncated RNA ends reveals new insights into Ribosome

948      Stalling in plants. *Plant Cell* **28**: 2398–2416.

949   Hou CY, Wu MT, Lu SH, Hsing YI, Chen HM. 2014. Beyond cleaved small

950      RNA targets: Unraveling the complexity of plant RNA degradome data.

951      *BMC Genomics* **15**: 15.

952   Ibrahim F, Maragkakis M, Alexiou P, Mourelatos Z. 2018. Ribothrypsis, a

953     novel process of canonical mRNA decay, mediates ribosome-phased

954     mRNA endonucleolysis. *Nat Struct Mol Biol* **25**: 302–310.

955   Karlova R, Van Haarst JC, Maliepaard C, Van De Geest H, Bovy AG,

956     Lammers M, Angenent GC, De Maagd RA. 2013. Identification of

957     microRNA targets in tomato fruit development using high-throughput

958     sequencing and degradome analysis. *J Exp Bot* **64**: 1863–1878.

959   Keene JD. 2010. Minireview: Global regulation and dynamics of ribonucleic

960     acid. *Endocrinology* **151**: 1391–1397.

961   Lei L, Shi J, Chen J, Zhang M, Sun S, Xie S, Li X, Zeng B, Peng L, Hauck A,

962     et al. 2015. Ribosome profiling reveals dynamic translational landscape

963     in maize seedlings under drought stress. *Plant J* **84**: 1206–1218.

964   Lipps HJ, Rhodes D. 2009. G-quadruplex structures: in vivo evidence and

965     function. *Trends Cell Biol* **19**: 414–422.

966   Lloréns-Rico V, Lluch-Senar M, Serrano L. 2015. Distinguishing between

967     productive and abortive promoters using a random forest classifier in

968     Mycoplasma pneumoniae. *Nucleic Acids Res* **43**: 3442–3453.

969   Lu Z, Lin Z. 2019. Pervasive and dynamic transcription initiation in

970     Saccharomyces cerevisiae. *Genome Res* **29**: 1198–1210.

971   Luo S, He F, Luo J, Dou S, Wang Y, Guo A, Lu J. 2018. Drosophila tsRNAs

972     preferentially suppress general translation machinery via antisense

973     pairing and participate in cellular starvation response. *Nucleic Acids*

974     *Res* **46**: 5250–5268.

975   Malka Y, Steiman-Shimony A, Rosenthal E, Argaman L, Cohen-Daniel L,

976     Arbib E, Margalit H, Kaplan T, Berger M. 2017. Post-transcriptional 3´-

977     UTR cleavage of mRNA transcripts generates thousands of stable

978     uncapped autonomous RNA fragments. *Nat Commun* **8**: 1–11.

979   Matsui T, Takita E, Sato T, Kinjo S, Aizawa M, Sugiura Y, Hamabata T,

980    Sawada K, Kato K. 2011. N-glycosylation at noncanonical Asn-X-Cys

981    sequences in plant cells. *Glycobiology* **21**: 994–999.

982  Matsuura H, Shinmyo A, Kato K. 2008. Preferential translation mediated by

983    Hsp81-3 5′-UTR during heat shock involves ribosome entry at the 5′-end

984    rather than an internal site in Arabidopsis suspension cells. *J Biosci*

985    *Bioeng* **105**: 39–47.

986  Mercer TR, Dinger ME, Bracken CP, Kolle G, Szubert JM, Korbie DJ,

987    Askarian-Amiri ME, Gardiner BB, Goodall GJ, Grimmond SM, et al.

988    2010. Regulated post-transcriptional RNA cleavage diversifies the

989    eukaryotic transcriptome. *Genome Res* **20**: 1639–1650.

990  Millevoi S, Moine H, Vagner S. 2012. G-quadruplexes in RNA biology. *Wiley*

991    *Interdiscip Rev RNA* **3**: 495–507.

992  Murata M, Nishiyori-Sueki H, Kojima-Ishiyama M, Carninci P, Hayashizaki Y,

993    Itoh M. 2014. Detecting expressed genes using CAGE. *Methods Mol Biol*

994    **1164**: 67–85.

995  Nagarajan VK, Kukulich PM, Von Hagel B, Green PJ. 2019. RNA

996    degradomes reveal substrates and importance for dark and nitrogen

997    stress responses of Arabidopsis XRN4. *Nucleic Acids Res* **47**: 9216–

998    9230.

999  Narsai R, Howell KA, Millar AH, O'Toole N, Small I, Whelan J. 2007.

1000    Genome-wide analysis of mRNA decay rates and their determinants in

1001    Arabidopsis thaliana. *Plant Cell* **19**: 3418–3436.

1002  Nepal C, Hadzhiev Y, Previti C, Haberle V, Li N, Takahashi H, Suzuki AMM,

1003    Sheng Y, Abdelhamid RF, Anand S, et al. 2013. Dynamic regulation of

1004    the transcription initiation landscape at single nucleotide resolution

1005    during vertebrate embryogenesis. *Genome Res* **23**: 1938–1950.

1006  Neymotin B, Athanasiadou R, Gresham D. 2014. Determination of in vivo

1007    RNA kinetics using RATE-seq. *RNA* **20**: 1645–1652.

1008    Parker MT, Knop K, Sherwood A V., Schurch NJ, Mackinnon K, Gould PD,

1009    Hall AJW, Barton GJ, Simpson GG. 2020. Nanopore direct RNA

1010    sequencing maps the complexity of arabidopsis mRNA processing and

1011    m6A modification. *Elife* **9**: e49658.

1012    Parker R. 2012. RNA degradation in Saccharomyces cerevisae. *Genetics*

1013    **191**: 671–702.

1014    Pelechano V, Wei W, Steinmetz LM. 2015. Widespread co-translational RNA

1015    decay reveals ribosome dynamics. *Cell* **161**: 1400–1412.

1016    Presnyak V, Alhusaini N, Chen YH, Martin S, Morris N, Kline N, Olson S,

1017    Weinberg D, Baker KE, Graveley BR, et al. 2015. Codon optimality is a

1018    major determinant of mRNA stability. *Cell* **160**: 1111–1124.

1019    Schoenberg DR, Maquat LE. 2009. Re-capping the message. *Trends*

1020    *Biochem Sci* **34**: 435–442.

1021    Shamimuzzaman M, Vodkin L. 2012. Identification of soybean seed

1022    developmental stage-specific and tissue-specific miRNA targets by

1023    degradome sequencing. *BMC Genomics* **13**: 310.

1024    Simms CL, Yan LL, Zaher HS. 2017. Ribosome Collision Is Critical for

1025    Quality Control during No-Go Decay. *Mol Cell* **68**: 361-373.

1026    Singh U, Khemka N, Rajkumar MS, Garg R, Jain M. 2017. PLncPRO for

1027    prediction of long non-coding RNAs (lncRNAs) in plants and its

1028    application for discovery of abiotic stress-responsive lncRNAs in rice

1029    and chickpea. *Nucleic Acids Res* **45**: e183.

1030    Song Y, Liu KJ, Wang TH. 2014. Elimination of ligation dependent artifacts

1031    in T4 RNA ligase to achieve high efficiency and low bias microRNA

1032    capture. *PLoS One* **9**: e94619.

1033    Subramanian M, Rage F, Tabet R, Flatter E, Mandel JL, Moine H. 2011. G-

44

1034    quadruplex RNA structure as a signal for neurite mRNA targeting. *EMBO*

1035    *Rep* **12**: 697–704.

1036  Takahashi H, Lassmann T, Murata M, Carninci P. 2012. 5' end-centered

1037    expression profiling using Cap-analysis gene expression (CAGE) and

1038    next-generation sequencing. *Nat Protoc* **7**: 542–561.

1039  Tani H, Mizutani R, Salam KA, Tano K, Ijiri K, Wakamatsu A, Isogai T,

1040    Suzuki Y, Akimitsu N. 2012. Genome-wide determination of RNA

1041    stability reveals hundreds of short-lived noncoding transcripts in

1042    mammals. *Genome Res* **22**: 947–956.

1043  Thieffry A, Vigh ML, Bornholdt J, Ivanov M, Brodersen P, Sandelin A. 2020.

1044    Characterization of arabidopsis thaliana promoter bidirectionality and

1045    antisense RNAs by inactivation of nuclear RNA decay pathways. *Plant*

1046    *Cell* **32**: 1845–1867.

1047  Tian B, Graber JH. 2012. Signals for pre-mRNA cleavage and

1048    polyadenylation. *Wiley Interdiscip Rev RNA* **3**: 385–396.

1049  Ueno D, Mukuta T, Yamasaki S, Mikami M, Demura T, Matsui T, Sawada K,

1050    Katsumoto Y, Okitsu N, Kato K. 2020. Different plant species have

1051    common sequence features related to mRNA degradation intermediates.

1052    *Plant Cell Physiol* **61**: 53–63.

1053  Ueno D, Yamasaki S, Demura T, Kato K. 2018. Comprehensive analysis of

1054    mRNA internal cleavage sites in Arabidopsis thaliana. *J Biosci Bioeng*

1055    **125**: 723–728.

1056  Weinberg DE, Shah P, Eichhorn SW, Hussmann JA, Plotkin JB, Bartel DP.

1057    2016. Improved Ribosome-Footprint and mRNA Measurements Provide

1058    Insights into Dynamics and Regulation of Yeast Translation. *Cell Rep*

1059    **14**: 1787–1799.

1060  Willmann MR, Berkowitz ND, Gregory BD. 2014. Improved genome-wide

1061       mapping of uncapped and cleaved transcripts in eukaryotes-GMUCT

1062       2.0. *Methods* **67**: 64–73.

1063    Wu X, Brewer G. 2008. The Regulation of mRNA Stability in Mammalian

1064       Cells: 2.0. *Gene* **23**: 10–21.

1065    Yamasaki S, Matsuura H, Demura T, Kato K. 2015. Changes in Polysome

1066       Association of mRNA Throughout Growth and Development in

1067       Arabidopsis thaliana. *Plant Cell Physiol* **56**: 2169–2180.

1068    Yu X, Willmann MR, Anderson SJ, Gregory BD. 2016. Genome-wide

1069       mapping of uncapped and cleaved transcripts reveals a role for the

1070       nuclear mrna cap-binding complex in cotranslational rna decay in

1071       arabidopsis. *Plant Cell* **28**: 2385–2397.

1072    Zhou M, Gu L, Li P, Song X, Wei L, Chen Z, Cao X. 2010. Degradome

1073       sequencing reveals endogenous small RNA targets in rice (Oryza sativa

1074       L. ssp. indica). *Front Biol (Beijing)* **5**: 67–90.

1075    Zhuang F, Fuchs RT, Sun Z, Zheng Y, Robb GB. 2012. Structural bias in T4

1076       RNA ligase-mediated 3′-adapter ligation. *Nucleic Acids Res* **40**: e54.

1077