1 **PlantRNA-FM: An Interpretable RNA Foundation Model for Exploration**

2 **Functional RNA Motifs in Plants**

3 Haopeng Yu[1,3†], Heng Yang[2,†], Wenqing Sun[1,3†], Zongyun Yan[1†], Xiaofei Yang[4,5], Huakun

4 Zhang[3*], Yiliang Ding[1,*] & Ke Li[2,*]

5 1 Department of Cell and Developmental Biology, John Innes Centre, Norwich Research Park,

6 Norwich NR4 7UH, UK

7 2 Department of Computer Science, University of Exeter, EX4 4QF, Exeter, UK

8 3 Key Laboratory of Molecular Epigenetics of the Ministry of Education, Northeast Normal

9 University, Changchun 130024, China

10 4 National Key Laboratory of Plant Molecular Genetics, CAS Center for Excellence in

11 Molecular Plant Sciences, Institute of Plant Physiology and Ecology, Chinese Academy of

12 Sciences, Shanghai 200032, China.

13 5 CAS-JIC Center of Excellence for Plant and Microbial Sciences, Institute of Plant Physiology

14 and Ecology, Chinese Academy of Sciences, Shanghai 200032, China.

15

16 * To whom correspondence should be addressed. Tel: +44 1392 724557; Email:

17 k.li@exeter.ac.uk.

18 Correspondence may also be addressed to yiliang.ding@jic.ac.uk or zhanghk045@nenu.edu.cn

19 † The authors wish it to be known that, in their opinion, the first four authors should be regarded

20 as Joint First Authors.

## ABSTRACT

The complex 'language' of plant RNA encodes a vast array of biological regulatory elements that orchestrate crucial aspects of plant growth, development, and adaptation to environmental stresses. Recent advancements in foundation models (FMs) have demonstrated their unprecedented potential to decipher complex 'language' in biology. In this study, we introduced PlantRNA-FM, a novel high-performance and interpretable RNA FM specifically designed based on RNA features including both sequence and structure. PlantRNA-FM was pre-trained on an extensive dataset, integrating RNA sequences and RNA structure information from 1,124 distinct plant species. PlantRNA-FM exhibits superior performance in plant-specific downstream tasks, such as plant RNA annotation prediction and RNA translation efficiency (TE) prediction. Compared to the second-best FMs, PlantRNA-FM achieved an $F1$ score improvement of up to 52.45% in RNA genic region annotation prediction and up to 15.30% in translation efficiency prediction, respectively. Our PlantRNA-FM is empowered by our interpretable framework that facilitates the identification of biologically functional RNA sequence and structure motifs, including both RNA secondary and tertiary structure motifs across transcriptomes. Through experimental validations, we revealed novel translation-associated RNA motifs in plants. Our PlantRNA-FM also highlighted the importance of the position information of these functional RNA motifs in genic regions. Taken together, our PlantRNA-FM facilitates the exploration of functional RNA motifs across the complexity of transcriptomes, empowering plant scientists with novel capabilities for programming RNA codes in plants.

## Introduction

The transcriptome contains a wide array of RNA motifs that impact diverse biological functions such as translation[1–5]. These RNA motifs encompass both RNA sequence and structure features. Previous individual studies have revealed the functional importance of RNA

46  sequence features such as the Kozak sequence motif[6]. Recently, our studies along with others

47  have suggested that both RNA secondary and tertiary structure motifs play important roles in

48  diverse biological processes[7–13]. Particularly in plants, the relatively low habitat temperatures

49  (~20 °C) favour the folding of RNA structure motifs, including RNA tertiary motifs such as

50  RNA G-quadruplex (rG4)[12]. However, systematically identifying functional RNA motifs

51  across transcriptomes remains a formidable challenge due to the high level of complexity

52  arising from astronomical combinations of the four nucleotide bases into tens of thousands of

53  transcripts[8,14]. For example, for a 50-nucleotide sequence, the number of artificially

54  synthesized sequences would be on the order of $4^{50}$ (approximately $1.27 \times 10^{30}$), which is

55  impossible to achieve experimentally. Additionally, the functional readouts using the reporter

56  gene assay for measuring biological functions such as translation may not be sensitive enough

57  to detect differences in individual single-nucleotide mutations[15].

58  The recent rapid advancements of foundation models (FMs) in artificial intelligence

59  (AI) are set to show exciting promise for supercharging scientific advances in life sciences[16].

60  FMs are distinguished by their massive scale, often encompassing millions to billions of

61  parameters. They are first pre-trained in a self-supervised manner on diverse forms of

62  unlabelled data. This makes them ideally suitable for bioscience, where acquiring abundant

63  labelled data is both prohibitively expensive and time-consuming. More importantly, FMs are

64  highly adaptable through fine-tuning and are poised to aid bioscientists in customising

65  generalist FMs in unravelling complex biological processes, paving the way for unprecedented

66  capabilities in modulating gene functions. For FMs on DNA sequences, DNABERT2 is one of

67  the FMs pre-trained on the genome sequences across 135 species, including mammals, fungi

68  and bacteria[17]. By pre-training on diverse human and non-human genomes, the Nucleotide

69  Transformers (NT) family learns transferable representations that enable accurate molecular

70  phenotype prediction with limited annotated data, while focusing on key genomic elements

71    without supervision[18]. FMs have also achieved success in protein sequences, also known as

72    protein language models. For example, ESM2 (Evolutionary Scale Modeling) has achieved

73    remarkable breakthroughs in atomic-level structure representations by pretraining on a vast

74    amount of protein sequences and structures[19].

75         For building RNA FM, several FM models were pre-trained using RNA sequence

76    information that has demonstrated great performance in RNA molecule design[20–22]. However,

77    RNA sequence information is not sufficient since RNA is capable of forming secondary or

78    tertiary structure motifs that are important for its functions[23,24]. Therefore, it is important to

79    generate an RNA FM including both RNA sequence and structure information to facilitate the

80    exploration of functional RNA motifs. Here, we developed PlantRNA-FM, a groundbreaking

81    RNA FM designed to globally identify functional RNA motifs including both RNA sequences

82    and structure motifs in plants (Fig. 1). By incorporating RNA sequences, annotations, and

83    structure information from 1,124 distinct plant species, PlantRNA-FM captures the extensive

84    diversity of plant transcriptomes (Fig. 1). We validate the superior performance of PlantRNA-

85    FM in downstream tasks compared to existing FMs. Furthermore, we also established an

86    interpretable framework based on our PlantRNA-FM to determine the critical regions across

87    the 5' untranslated regions (5' UTRs) that significantly impact translation. Remarkably,

88    PlantRNA-FM identifies RNA motifs at the transcriptome-wide scale that are functionally

89    important to translation including both RNA sequences, and secondary and tertiary structure

90    motifs. We further experimentally validated these identified RNA motifs in plants. The

91    development of our PlantRNA-FM represents a significant leap forward in our ability to

92    decipher hidden regulatory codes among the extensive complexity of nucleotides across the

93    transcriptome, opening new avenues for RNA-based gene regulation.

94    **1. Results**

## 1.1. Our PlantRNA-FM integrates both RNA sequence and structure information of the transcriptomes across 1,124 plant species.

The plant kingdom encompasses approximately 500,000 species, exhibiting remarkable diversity. The One Thousand Plant Transcriptomes Initiative (1KP) sequenced the transcriptomes of 1,124 species, capturing the extensive diversity of plant transcriptomes[14]. Here, we took advantage of this unique resource and generated the pre-training dataset for our PlantRNA-FM (Fig. 1). Different from existing FMs, our PlantRNA-FM was designed to capture and learn both RNA sequences and RNA structure motifs. We employed *RNAfold*[25] to predict RNA structures of individual RNA sequences across 1,124 transcriptomes and integrated them into the pre-training dataset. Our PlantRNA-FM has 35 million parameters, including 12 transformer network layers, 24 attention heads, and an embedding dimension of 480, optimised for RNA understanding rather than generation (**Methods**). Our tokenization approach surpasses the constraints of conventional $k$-mers and BPE methods, ensuring the preservation of RNA structure motifs as coherent units throughout the pre-training process (**Methods**). In addition, we incorporated RNA annotation information (CDS and UTRs) and employed advanced pre-training techniques, such as sequence truncation, filtering and masked nucleotide modeling (**Methods**).

To assess the effectiveness of our PlantRNA-FM in RNA structure prediction tasks, we evaluated its performance (Fig. S1, Table S1) using three benchmark datasets: bpRNA, ArchiveII, and RNAstralign[26–28]. The $F1$ score, which is the harmonic mean of precision and recall, was used to measure the model's predictive performance on these datasets. The $F1$ scores achieved by our PlantRNA-FM on these three datasets were 0.750, 0.924, and 0.981, respectively, while *RNAfold* alone only obtained $F1$ scores of 0.278, 0.759, and 0.748 (Fig. S1, Table S1). When compared to other state-of-the-art FMs, PlantRNA-FM outperformed the second-best model by 22.10%, 27.49%, and 17.38% on the respective datasets (Fig. S1,

120 Table S1). Therefore, the unique integration of RNA structure information equips our

121 PlantRNA-FM with the ability to predict RNA structure more accurately.

## 1.2. PlantRNA-FM demonstrates superior performance on plant-specific downstream tasks

124 To evaluate the performance of PlantRNA-FM, we curated a benchmark set consisting

125 of four other state-of-the-art FMs: DNABERT-2, Nucleotide Transformer, ESM2, and

126 cdsBERT. We assessed their performance in two plant-specific downstream tasks: genic region

127 annotation and translation efficiency (TE) prediction (Fig. 2a).

128 In the RNA genic region annotation prediction task, we aimed to identify and classify

129 different genic regions of given RNA sequences, such as the 5' UTR, coding sequence (CDS),

130 and 3' UTR. We used the transcriptomes of two model plant species, *Arabidopsis thaliana* (a

131 dicot model plant) and *Oryza sativa L.* ssp. *Japonica* (rice, a moncot model plant). Both of

132 them were not included in our pre-training dataset. For the RNA genic region annotation

133 prediction in these two species, our PlantRNA-FM outperformed other FM models, achieving

134 average $F1$ scores of 0.974 and 0.958 for *Arabidopsis* and rice, respectively, surpassing the

135 second-best model by 52.45% and 43.90% (Fig. 2b, Table 1).

136 For translation, one of the key RNA biological processes, previous research has

137 highlighted the critical role of the 5' UTR in regulating translation efficiencies[17–19,21,29–31]. To

138 evaluate the TE prediction performance of our PlantRNA-FM, we used the 5' UTR sequences

139 of both *Arabidopsis* and rice transcriptomes along with the corresponding TE values measured

140 by polysome profiling[8]. We first classified the TE datasets into high and low TE groups, using

141 the mean plus or minus the standard deviation as the threshold. In the TE prediction task,

142 PlantRNA-FM achieved $F1$ scores of 0.735 and 0.737 for *Arabidopsis* and rice, respectively,

143 outperforming the second-best model by 15.30% and 13.83% (Fig. 2c). Taken together, our

144 PlantRNA-FM is better suited for plant-specific downstream tasks compared to other FMs pre-

145 trained on non-plant datasets.

## 1.3. Interpretable PlantRNA-FM revealed RNA features important to translation

148       A general roadblock in applying AI models to biology is that, while these models

149 demonstrate strong predictive capabilities, the key to their successful application lies in

150 interpreting them to uncover the biological principles learned by the AI. In this paper, we

151 established an interpretable framework to derive an attention contrast matrix from our

152 PlantRNA-FM (**Methods**). In particular, we are interested in extracting the key RNA features

153 within the 5' UTR that significantly impact RNA translation, i.e., elucidating the RNA motifs

154 associated with translation (Fig. 3a). We developed two models in parallel: one is the true

155 model, denoted as PlantRNA-FM(+), trained using the real TE dataset, while the other one is

156 called the background model, PlantRNA-FM(-), altered using the same dataset but with

157 randomly assigned labels (Fig. 3a). The $F1$ score achieved by the background model is

158 approximately 50%, which is close to the random chance (mean $F1 = 0.522$), while the true

159 model attained a significantly higher mean $F1$ score of 0.737. This indicates that PlantRNA-

160 FM(+) has successfully learned the RNA features in the 5' UTR sequences associated with

161 translation.

162       By subtracting the attention matrices of the background model from those of the true

163 model, we obtained an attention contrast matrix that highlighted the significance of nucleotides

164 in the 5' UTR contributing to TE (Fig. 3a). Across the transcriptomes, we observed an increase

165 in attention contrast scores as the position approached the AUG start codon in both *Arabidopsis*

166 and rice (Fig. 3b). This result indicates that positions close to the start codon contribute the

167 most to the TE values. By underlining the RNA sequence contents with high contrast attention

168 score (identified by a z-score > 2.326), our PlantRNA-FM successfully identified the Kozak

169     sequence motifs in both *Arabidopsis* and rice transcriptomes that are associated with TE (Fig.

170     3c, 3d). This result demonstrates that our PlantRNA-FM successfully identifies evolutionarily

171     conserved RNA motifs that are important to translation (Fig. 3c, 3d).

## 1.4.    PlantRNA-FM globally identifies the translation-associated RNA secondary structure motifs

174     Since RNA structure is the unique RNA feature incorporated in our PlantRNA-FM, we

175     further identified the RNA secondary structure motifs important to translation through the

176     model's attention contrast matrix and an unsupervised hierarchical clustering strategy (Fig. 4a,

177     **Methods**). Overall, we identified 112 RNA secondary structure motifs that are important to

178     translation, including 63 low translation-associated and 49 high translation-associated RNA

179     secondary structure motifs (Table S2). Notably, we identified low translation-associated RNA

180     secondary structure motifs with high GC base pairs such as the RNA secondary structure motif

181     with four GC base pairs in the stem (Fig. 4b). Interestingly, we also identified high translation-

182     associated RNA structure motifs with a balanced ratio of GC and AU base pairs such as the

183     RNA structure motif with four base pairs formed by two repeats of ACGU (Fig. 4c).

184     To validate our identified RNA secondary structure motifs important to translation, we

185     conducted experimental validation using the dual luciferase reporter assay in plants[12]. For the

186     high translation-associated RNA secondary structure motif with four base pairs formed by two

187     repeats of ACGU, we changed the two AU base pairs to the two GC base pairs, resulting in a

188     significant decrease in TE with a reduction up to 5.3 -fold (Fig. 4d). In contrast, when we

189     exchanged the low translation-associated RNA secondary structure motif with four GC base

190     pair in the stem for the high translation-associated RNA secondary structure motif with a

191     balanced mix of GC and AU base pairs, we found a significant increase in TE (Fig. 4e).

192     Notably, when we completely disrupted this low translation-associated RNA structure motif,

193     resulting in complete single-strandedness, we observed an even greater enhancement of TE up

194     to 2.1-fold (Fig. 4f). Our results demonstrate that PlantRNA-FM is capable of determining

195     functional RNA secondary structure motifs in plants.

## 1.5. PlantRNA-FM globally identifies the translation-associated RNA tertiary structure motifs

198     RNA G-quadruplexes (rG4s) are one of the RNA tertiary structure motifs formed by

199     the stacking of two or more G-quartets, composed of four guanines held together by both

200     Watson-Crick and Hoogsteen hydrogen bonds[8,32,33]. Previous studies have demonstrated the

201     important role of individual rG4s in repressing translation[34]. However, it is impossible to

202     identify all the rG4 motifs important to translation from tens of thousands of rG4 motifs across

203     the transcriptome. Therefore, we took advantage of our PlantRNA-FM to identify the

204     translation-associated rG4s at the transcriptome-wide scale.

205     We first obtained all rG4 motifs in the 5' UTRs from our G4Atlas database[33].

206     Subsequently, we identified all rG4 motifs associated with translation using our model's

207     attention contrast matrix across the transcriptome (**Methods**). Notably, we only identified rG4

208     motifs associated with low TE, particularly with both GGA and GGU repeat (Table S3).

209     Therefore, our results indicate that rG4 serves as a translation repressor, which agrees with

210     previous studies on individual rG4s[35–37]. To validate our identified translation-associated rG4

211     motifs, we conducted the experimental validation using dual luciferase reporter assay in

212     plants[12]. We fused the 5'UTRs containing our identified rG4 motif and the corresponding

213     disrupted rG4 motif with the luciferase reporter genes[12]. We then measured the corresponding

214     TEs in plants and observed a significant increase of up to 5.8-fold in the disrupted rG4 motif

215     compared to the TE in the native rG4 motif (Fig. 4g). These results indicate that our PlantRNA-

216     FM is also capable of identifying functional RNA tertiary structure motifs such as translation-

217     associated rG4 motifs throughout the transcriptome.

## 2. Discussion

219  In this study, we developed PlantRNA-FM, a high-performance and interpretable plant-

220 specific RNA FM. PlantRNA-FM (Fig. 1) is designed for understanding RNA sequence and

221 structure information rather than generation. This state-of-the-art model was specifically

222 designed based on the extensive plant RNA information from 1,124 plant species, thereby

223 capturing the remarkable diversity of plant RNA features. From the perspective of the dataset,

224 we have incorporated RNA sequence information of all the RNAs from the transcriptomes

225 across 1,124 plant species. We also incorporated the corresponding RNA annotation

226 information. The integration of RNA structure information in our PlantRNA-FM achieves

227 superior performance in RNA structure prediction tasks compared to other FMs (Fig. S1).

228 Regarding the model architecture, we adopted a fine-grained tokenization method with single-

229 nucleotide resolution. This contrasts with commonly used tokenization methods, such as byte

230 pair encoding (BPE) and k-mers, which rely on frequency-based tokenization and may

231 inadvertently fragment RNA structure motifs into arbitrary pieces. This strategy ensures the

232 precise extraction and preservation of RNA structure motifs as coherent units throughout the

233 pre-training process, thereby maintaining the integrity of crucial structure information.

234 Additionally, PlantRNA-FM integrates rotational position embedding (RoPE), a technique that

235 has proven effective in enhancing the modeling capabilities for long tokens in large FMs[38]. The

236 implementation of RoPE leads to a approximately 30% reduction in the number of parameters

237 in the embedding layer, consequently improving the efficiency of RNA tokenisation and

238 modeling.

239  The superior performance of PlantRNA-FM can be further demonstrated in the plant-

240 specific downstream tasks (Fig. 2a). Our PlantRNA-FM achieved the best *F1* scores of 0.974

241 and 0.958 for the genic region annotation in *Arabidopsis* and rice, while our PlantRNA-FM

242 also achieved much better performance in predicting TE compared to other FMs (Fig. 2b, 2c).

243 The outperformance of our PlantRNA-FM is likely due to the combination of both RNA

244  sequence and structure information in our pre-training dataset, highlighting the importance of

245  RNA structure, a key RNA feature, in regulating RNA biological processes.

246      Notably, we developed an interpretable framework for our PlantRNA-FM to explore

247  the RNA features within the 5' UTR that influence translation (Fig. 3a). Using the attention

248  contrast matrices, we found that the nucleotides in the regions close to the start codon affect

249  the translation the most, emphasizing the importance of positional information of functional

250  RNA motifs (Fig. 3b). In contrast to conventional meta-gene analysis, our PlantRNA-FM is

251  capable of providing positional information of RNA motifs across transcriptomes, which is

252  critical for biological regulatory functions. Furthermore, the Kozak sequence, an evolutionary

253  conserved translation-associated sequence motif across translation initiation sites was

254  successfully identified in both *Arabidopsis* and rice using our PlantRNA-FM (Fig. 3c, 3d). This

255  result successfully demonstrates the capability of our PlantRNA-FM in identifying the RNA

256  sequence motifs important to translation across the transcriptomes. By using an unsupervised

257  hierarchical clustering strategy to explore our attention contrast matrix, we further

258  systematically identified RNA secondary and tertiary structure motifs that are functionally

259  important to translation (Fig. 4a). Notably, we identified both high translation-associated and

260  low translation-associated RNA secondary structure motifs where their differences are mainly

261  in the strengths of the base pairs (Fig. 4b, 4c). This suggests that RNAs may adopt different

262  RNA structure motifs with diverse folding strengths in regulating biological processes such as

263  translation. In contrast to conventional meta-gene analysis, our PlantRNA-FM is capable of

264  delivering a comprehensive understanding of functional RNA motifs such as the type of RNA

265  motifs, the genic position of the RNA motifs, the positive or negative effects of the RNA motifs

266  on their functions, and the exact contributions of the RNA motifs to their functions. For

267  instance, high GC content in the 5' UTR has been shown to be anti-correlated with translation

268  efficiency[39–41]. However, these correlations are not able to facilitate understanding of which

269 type of regulatory motifs with high GC content repress translation. Here, our PlantRNA-FM

270 revealed diverse RNA structure motifs such as the RNA secondary structure motif with four

271 GC base pairs in the stem and rG4s, serving as low translation-associated RNA motifs. This

272 suggests the diversity of RNA regulatory motifs across the transcriptomes (Fig. 4b).

273 In summary, we have built the first interpretable RNA FM with both RNA sequence

274 and structure information. Our PlantRNA-FM was pre-trained using 1,124 plant transcriptomes.

275 We have demonstrated that our PlantRNA-FM is capable of identifying functional RNA motifs

276 such as translation-associated sequence and structure motifs across the transcriptomes.

277 Through our experimental validations, we have elucidated novel translation-associated RNA

278 motifs in plants. Our FM model can be extended to explore functional RNA motifs in other

279 kingdoms and investigate RNA motifs important for other biological functions such as RNA

280 decay and maturation. Our PlantRNA-FM is poised to transform the way we determine RNA

281 motifs for regulating gene expression, opening new horizons for programming RNA codes to

282 facilitate crop improvements and RNA-based applications.

283 ## 3. Methods

284 ### 3.1. Pre-training datasets curation

285 The plant transcriptome data used for pre-training PlantRNA-FM was obtained from

286 the one thousand plant transcriptomes project (1KP)[14]. Note that modeling genomic sequences

287 differs significantly from natural language modeling. For instance, while RNA sequences are

288 one-dimensional, they strictly follow biological genomic patterns and depend heavily on

289 certain structural characteristics. In contrast, natural language models are more resilient and

290 can tolerate linguistic errors such as typos and grammar mistakes. Thus, effective RNA

291 sequence curation is crucial to minimize the impact of noisy data and enhance modeling

292 performance. Specifically, our data curation protocol is as follows.

293 • **Sequence truncation and filtering**: We truncated RNA sequences exceeding 512

294 nucleotides to comply with the model's maximum length capacity and filtered out sequences

295 shorter than 20 nucleotides to eliminate noise, such as RNA fragment sequences.

296 • **RNA secondary structure annotation:** Given the significant impact of RNA

297 secondary structures on sequence function, we annotated the local RNA structures of all RNA

298 sequences using ViennaRNA (with parameters maxBPspan = 30)[25].

299 • **Annotation of CDS and UTR sequences:** After obtaining the assembled transcripts

300 and translated RNA regions from the dataset, we retrieve the CDS (translated RNA), 5' UTR,

301 and 3' UTR sequences (upstream and downstream of the translated RNA).

## 3.2. Model architecture

303 In this study, we developed PlantRNA-FM, a specialised language model based on the

304 transformer architecture (Fig. 1). PlantRNA-FM has 35 million parameters, including 12

305 transformer network layers, 24 attention heads, and an embedding dimension of 480. We

306 applied layer normalisation and residual connections both before and after the encoder block.

307 As our focus is on RNA understanding rather than generation, we only utilised the encoder

308 component of the transformer architecture. PlantRNA-FM is capable of processing sequences

309 up to 512 nucleotides in length, making it compatible with consumer-grade GPUs, such as the

310 Nvidia RTX 4090, with a batch size of 16. The model was trained on four A100 GPUs over a

311 period of three weeks, completing 3 epochs.

## 3.3. Pretraining strategies of PlantRNA-FM

313 To develop an RNA FM for exploiting all potential patterns within RNA sequences, we

314 investigated the biological domain knowledge of RNA sequences and propose three self-

315 supervised pre-training objectives to enhance the foundational model.

### 3.3.1. Pretraining with Masked nucleotides modeling

317 Inspired by the concept of masked language modelling (MLM) in NLP, we introduced

318 masked nucleotide modelling (MNM) for RNA sequences. This approach involves randomly

319 masking a portion of nucleotides and leveraging the model itself to reconstruct these masked

320 nucleotides. Note that the ability to accurately reconstruct masked nucleotides indicates that

321 the model is empowered with the capability of understanding RNA sequence. MNM

322 dynamically selects 20% of nucleotides for masking in each input sequence, as opposed to the

323 fixed 15% masking used in the classic MLM objective designed for shorter natural language

324 sentences. This increased masking ratio is chosen to enhance MNM's modeling capability,

325 considering that RNA sequences typically contain around one thousand bases. Specifically,

326 10% are replaced with a '<mask>' token, 5% with random nucleotides, and the remaining 5%

327 are left as is. This approach, which aims for token classification, employs cross-entropy as the

328 loss function to enhance the model's predictive accuracy for masked or replaced nucleotides.

329 The loss function $L_{MLM}(\theta)$ for MLM is defined as follows:

330
$$L_{MLM}(\theta) = -\frac{1}{|m|} \sum_{i \in m} \log p_\theta(x_i \mid x_{\setminus i}),$$

331 where $\theta$ and $m$ are the parameter set inside the FM and the number of masked nucleotides.

332 $p_\theta(x_i \mid x_{\setminus i})$ indicates the probability of predicting the masked nucleotide $x_i$ based on its

333 context $(x_{\setminus i})$.

334 **3.3.2. Pretraining with RNA Structure Prediction**

335 We hypothesise that effectively aligning RNA sequences with their corresponding

336 secondary structures is important during the pre-training phase. In practice, we annotated the

337 secondary structures within the 1KP dataset, which comprises 50 billion nucleotides. This

338 establishes a robust foundation for our model to recognise the critical role of secondary

339 structures. Based on these annotated data, we utilized cross-entropy as the loss function to

340 predict the RNA secondary structure:

341
$$L_{SSP}(\theta) = -\sum_{i=1}^{N} \sum_{c=1}^{C} y_{i,c} \log p(y_{i,c} \mid x; \theta),$$

342    where $N$ is the length of the RNA sequence, i.e., the total number of nucleotides in the

343    sequence; $C$ denotes the number of prediction for each nucleotide (e.g., '(', ')', '.'); $y_{i,c}$ is the

344    prediction of the $i$-th nucleotide $c$, and $p(y_{i,c} \mid x; \theta)$ is the probability predicted by the model

345    parameterised by $\theta$. $L_{SSP}(\theta)$ is the loss function that quantifies the discrepancy between the

346    model's predicted probabilities for each nucleotide's secondary structure and the actual

347    structure, with the aim of minimising this loss to improve the model's accuracy in secondary

348    structure prediction.

### 3.3.3. Pretraining with RNA annotation prediction

350    RNA sequences exhibit significant variation across different regions, each serving

351    distinct functions within an organism. Beyond the two aforementioned training objectives, the

352    third one focuses on classifying regions within RNA sequences. The loss function is as follows:

$$L_{CLS}(\theta) = -\sum_{i=1}^{N} \sum_{r=1}^{R} y_{i,r} \log p(y_{i,r} \mid x; \theta),$$

354    where $N$ is the length of the RNA sequence, i.e., the total number of nucleotides or segments

355    considered for classification. $R$ represents the number of region categories we are classifying,

356    including CDS, 3' UTR, and 5' UTR. $y_{i,r}$ is the prediction of the $i$-th nucleotide $r$.

357    $p(y_{i,r} \mid x; \theta)$ is the probability predicted by the model, with parameters $\theta$, for the $i$-th

358    nucleotide given the RNA sequence $x$. $L_{CLS}(\theta)$ is the cross-entropy loss function aimed at

359    training the model to identify different regions.

### 3.4.    Fine-tuning of downstream tasks

361    After the pre-training phase, our FM can be fine-tuned to adapt to various downstream

362    tasks. The fine-tuning phase consists of three steps. First, we gathered an annotated dataset

363    specific to each downstream task, which consists of sequences and their corresponding labels.

364    Note that we pre-sliced any sequences that exceed the model's maximum length, to ensure

365    compatibility. Next, using the pre-trained FM as a starting point, we adapted the output layer

366    to accommodate the requirements of RNA modelling tasks, which may include outputting

367   sequences, labels, or scalar values. Finally, the training and inference processes are tailored to

368   the demands of each downstream task by selecting task-specific optimisers, loss functions, and

369   tuning hyperparameters to achieve optimal performance. The source code for our training and

370   inference can be found in our repository.

## 3.5.   Polysome profiling mapping and data processing

372   Raw polysome profiling sequencing data for *A. thaliana* were obtained from published

373   research[12]. For rice, we performed polysome-seq using the same protocol as *Arabidopsis*[8]. The

374   genomes and annotation files of *O. sativa* and *A. thaliana* were obtained from Phytozome v13

375   with version of *Oryza sativa* v7.0 and TAIR10[44]. After extracting the transcriptome sequence

376   through the reference genome and annotation files, clean polysome profiling and RNA-Seq

377   reads were mapped to the reference transcriptome using HISAT2 and followed by library

378   normalisation and quantification using DESeq2[45,46]. Next, genes with an RPKM of less than

379   1 were removed, and the TE of each gene was calculated by dividing the polysome-associated

380   RNA levels (polysome profiling RNA-seq) by the corresponding RNA levels (RNA-Seq)[12].

381   Subsequently, the dataset was classified as high or low TE, using the mean plus or minus the

382   standard deviation as a threshold, and were respectively assigned the labels 1 and 0 for high

383   and low TE.

## 3.6.   RNA structure motif identification approach

### 3.6.1.   Extraction of the attention contrast matrix

386   To facilitate better model interpretation, we created two additional models. One is the

387   true model, denoted as PlantRNA-FM(+), trained using the real TE labels, while the other one

388   is the background model, PlantRNA-FM(-), altered using the same dataset but with randomly

389   assigned labels. Specifically, we fine-tuned the pre-trained PlantRNA-FM (+) and (-) on each

390   dataset for 100 epochs, using regular hyperparameter settings. To avoid overfitting, we

391   employed an early stopping strategy to terminate the fine-tuning process when the best *F1* score

392   remained unchanged for 30 epochs. Once the fine-tuning was completed, we used the fine-

393    tuned models to predict each dataset and derive the raw attention score matrices corresponding

394    to each RNA sequence. Since the raw attention score matrices are five-dimensional, we

395    reshaped them through average-based downsampling to generate attention contrast matrices.

396    Finally, we subtracted the attention contrast matrices of PlantRNA-FM (+) from those of

397    PlantRNA-FM (-). Furthermore, we padded any negative values in the attention contrast

398    matrices with zeros for better visualisation.

### 399    3.6.2.  Generation of the RNA structure motif seed library

400    To identify RNA structure motifs, we first generate a library of that contains RNA

401    structure motif seeds derived from RNA sequences across the transcriptomes. In this work, we

402    apply the Zuker algorithm from the Vienna RNA package to obtain all suboptimal RNA

403    structure foldings for each RNA in our dataset[35,36]. We restrict the length of the RNA structure

404    motifs to a maximum of thirty[48]. The folded RNA structures are then annotated using

405    "bpRNA". Subsequently, all RNA structure motifs are extracted to generate a seed library of

406    RNA structure motifs for the plant transcriptomes[26]. In order to obtain reliable RNA structure

407    motifs, we set the range of RNA structure stems from 4 to 7, and the loop length from 4 to 9.

### 408    3.6.3.  Identification of translation-associated RNA secondary structure motifs

409    From the previous step, we obtained all potential foldings of the RNA structure motif

410    in the 5' UTR and aligned them with the attention contrast matrix. For each RNA structure

411    motif, we evaluated it using a paired $t$-test to obtain a $p$-value. Then, we corrected the obtained

412    $p$-value using the Benjamini-Hochberg (BH) method. RNA structure motifs with $p$-values less

413    than 0.01 were considered significant and extracted as the high-attention RNA structure

414    motifs. Then we extracted their corresponding RNA sequence and converted them into

415    numerical matrices using the one-hot encoding method. Subsequently, we applied an

416    unsupervised hierarchical clustering strategy to classify the nucleotides corresponding to the

417    positions of the RNA structure pairs into 2 to 100 clusters[49]. For each cluster containing a

418    minimum of 30 high-attention RNA structure motifs, the significance was assessed using

419   Fisher's exact test. RNA motifs with an odds ratio over 1 and a $p$-value below 0.05 were

420   identified as high translation-associated motifs. On the contrary, those with an odds ratio less

421   than 1 and a $p$-value below 0.05 were associated with low TE. Additionally, we calculated the

422   mean information content of all bases, defined as the "Average Positional Information Content"

423   (APIC). RNA motifs with an APIC below 1.5 were excluded from further analysis.

### 424   3.6.4. Identification of translation-associated rG4s

425   We obtained all potential rG4 in rice from our G4Atlas database[33]. Next, we aligned

426   the rG4 sequences with the corresponding attention contrast matrix and employed the paired $t$-

427   test to assess the statistical significance. For each length of rG4, we adjusted its $p$-value using

428   the Benjamini-Hochberg (BH) correction method and selected rG4s with a $p$-value less than

429   0.01 as the high attention rG4s.

## 430   4. Data availability

431   The polysome-seq sequence data of *A. thaliana* was obtained from the Sequence Read Archive

432   (SRA) (https://www.ncbi.nlm.nih.gov/sra) under BioProject ID number PRJNA762705[8]. The

433   raw sequence data of *O. sativa* has been deposited in the Sequence Read Archive (SRA)

434   (https://www.ncbi.nlm.nih.gov/sra) under BioProject ID number PRJNA1112739.

## 435   5. Code availability

436   The source code of this study is freely available at Huggingface

437   (https://huggingface.co/yangheng/PlantRNA-FM).

438

## 439   6. References

440   1. Piao, M., Sun, L. & Zhang, Q. C. RNA Regulations and Functions Decoded by

441   Transcriptome-wide RNA Structure Probing. *Genomics, Proteomics & Bioinformatics* **15**,

442   267–278 (2017).

443   2. Komatsu, K. R. *et al.* RNA structure-wide discovery of functional interactions with

444   multiplexed RNA motif library. *Nature Communications* **11**, 6275 (2020).

445   3. Espah Borujeni, A., Channarasappa, A. S. & Salis, H. M. Translation rate is controlled by

446   coupled trade-offs between site accessibility, selective RNA unfolding and sliding at

447   upstream standby sites. *Nucleic Acids Research* **42**, 2646–2659 (2014).

448   4. Gorochowski, T. E., Ignatova, Z., Bovenberg, R. A. L. & Roubos, J. A. Trade-offs between

449   tRNA abundance and mRNA secondary structure support smoothing of translation

450   elongation rate. *Nucleic Acids Research* **43**, 3022–3032 (2015).

451   5. Mortimer, S. A., Kidwell, M. A. & Doudna, J. A. Insights into RNA structure and function

452   from genome-wide studies. *Nat Rev Genet* **15**, 469–479 (2014).

453   6. Kozak, M. An analysis of vertebrate mRNA sequences: intimations of translational control.

454   *Journal of Cell Biology* **115**, 887–903 (1991).

455   7. Ding, Y. *et al.* In vivo genome-wide profiling of RNA secondary structure reveals novel

456   regulatory features. *Nature* **505**, 696–700 (2014).

457   8. Yang, X. *et al.* RNA G-quadruplex structure contributes to cold adaptation in plants. *Nat*

458   *Commun* **13**, 6224 (2022).

459   9. Xu, B. *et al.* Recent advances in RNA structurome. *Sci. China Life Sci.* **65**, 1285–1324

460   (2022).

461   10.   Yang, M. *et al.* Intact RNA structurome reveals mRNA structure-mediated regulation

462   of miRNA cleavage in vivo. *bioRxiv* 2019.12.21.885699 (2020) doi:10/ghccqf.

463   11.   Yang, M. *et al.* In vivo single-molecule analysis reveals COOLAIR RNA structural

464   diversity. *Nature* 1–6 (2022) doi:10.1038/s41586-022-05135-9.

465   12.   Xiaofei Yang & Haopeng Yu. Wheat in vivo RNA structure landscape reveals a

466   prevalent role of RNA structure in modulating translational subgenome expression

467   asymmetry. 26 (2021).

468    13.    Deng, H. *et al.* Rice In Vivo RNA Structurome Reveals RNA Secondary Structure

469    Conservation and Divergence in Plants. *Molecular Plant* **11**, 607–622 (2018).

470    14.    One Thousand Plant Transcriptomes Initiative. One thousand plant transcriptomes and

471    the phylogenomics of green plants. *Nature* **574**, 679–685 (2019).

472    15.    Cao, J. *et al.* High-throughput 5′ UTR engineering for enhanced protein production in

473    non-viral gene therapies. *Nat Commun* **12**, 4138 (2021).

474    16.    Consens, M. E. *et al.* To Transformers and Beyond: Large Language Models for the

475    Genome. Preprint at https://doi.org/10.48550/arXiv.2311.07621 (2023).

476    17.    Zhou, Z. *et al.* DNABERT-2: Efficient Foundation Model and Benchmark For Multi-

477    Species Genome. Preprint at https://doi.org/10.48550/arXiv.2306.15006 (2023).

478    18.    Dalla-Torre, H. *et al.* The Nucleotide Transformer: Building and Evaluating Robust

479    Foundation Models for Human Genomics. 2023.01.11.523679 Preprint at

480    https://doi.org/10.1101/2023.01.11.523679 (2023).

481    19.    Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a

482    language model. *Science* **379**, 1123–1130 (2023).

483    20.    Chu, Y. *et al.* A 5' UTR Language Model for Decoding Untranslated Regions of

484    mRNA and Function Predictions. 2023.10.11.561938 Preprint at

485    https://doi.org/10.1101/2023.10.11.561938 (2023).

486    21.    Hallee, L., Rafailidis, N. & Gleghorn, J. P. cdsBERT - Extending Protein Language

487    Models with Codon Awareness. *bioRxiv* 2023.09.15.558027 (2023)

488    doi:10.1101/2023.09.15.558027.

489    22.    Chen, K. *et al.* Self-supervised learning on millions of primary RNA sequences from

490    72 vertebrates improves sequence-based RNA splicing prediction. *Briefings in*

491    *Bioinformatics* **25**, bbae163 (2024).

492  23.  Yang, X., Yang, M., Deng, H. & Ding, Y. New Era of Studying RNA Secondary

493      Structure and Its Influence on Gene Regulation in Plants. *Front. Plant Sci.* **9**, (2018).

494  24.  Zhang, H., Chung, B. Y.-W., Wang, Z. & Ding, Y. Editorial: Plant RNA structure.

495      *Front. Plant Sci.* **14**, (2023).

496  25.  Lorenz, R. *et al.* ViennaRNA Package 2.0. *Algorithms Mol Biol* **6**, 26 (2011).

497  26.  Danaee, P. *et al.* bpRNA: large-scale automated annotation and analysis of RNA

498      secondary structure. *Nucleic Acids Research* **46**, 5381–5394 (2018).

499  27.  Sloma, M. F. & Mathews, D. H. Exact calculation of loop formation probability

500      identifies folding motifs in RNA secondary structures. *RNA* **22**, 1808–1818 (2016).

501  28.  Tan, Z., Fu, Y., Sharma, G. & Mathews, D. H. TurboFold II: RNA structural

502      alignment and secondary structure prediction informed by multiple homologs. *Nucleic*

503      *Acids Research* **45**, 11570–11581 (2017).

504  29.  Hardy, E. C. & Balcerowicz, M. Untranslated yet indispensable—UTRs act as key

505      regulators in the environmental control of gene expression. *Journal of Experimental*

506      *Botany* erae073 (2024) doi:10.1093/jxb/erae073.

507  30.  Dever, T. E., Ivanov, I. P. & Hinnebusch, A. G. Translational regulation by uORFs and

508      start codon selection stringency. *Genes Dev.* **37**, 474–489 (2023).

509  31.  Evfratov, S. A. *et al.* Application of sorting and next generation sequencing to study

510      5′-UTR influence on translation efficiency in Escherichia coli. *Nucleic Acids Research* **45**,

511      3487–3502 (2017).

512  32.  Lyu, K., Chow, E. Y.-C., Mou, X., Chan, T.-F. & Kwok, C. K. RNA G-quadruplexes

513      (rG4s): genomics and biological functions. *Nucleic Acids Research* **49**, 5426–5450 (2021).

514  33.  Yu, H., Qi, Y., Yang, B., Yang, X. & Ding, Y. G4Atlas: a comprehensive

515      transcriptome-wide G-quadruplex database. *Nucleic Acids Research* **51**, D126–D134

516      (2023).

34. Song, J., Perreault, J.-P., Topisirovic, I. & Richard, S. RNA G-quadruplexes and their potential regulatory roles in translation. *Translation* **4**, e1244031 (2016).

35. Kumari, S., Bugaut, A., Huppert, J. L. & Balasubramanian, S. An RNA G-quadruplex in the 5′ UTR of the NRAS proto-oncogene modulates translation. *Nat Chem Biol* **3**, 218–221 (2007).

36. Beaudoin, J.-D. & Perreault, J.-P. 5′-UTR G-quadruplex structures acting as translational repressors. *Nucleic Acids Research* **38**, 7022–7036 (2010).

37. Jia, L. *et al.* Decoding mRNA translatability and stability from the 5′ UTR. *Nat Struct Mol Biol* **27**, 814–821 (2020).

38. Su, J. *et al.* RoFormer: Enhanced transformer with Rotary Position Embedding. *Neurocomputing* **568**, 127063 (2024).

39. Araujo, P. R. *et al.* Before It Gets Started: Regulating Translation at the 5′ UTR. *International Journal of Genomics* **2012**, e475731 (2012).

40. Leppek, K., Das, R. & Barna, M. Functional 5′ UTR mRNA structures in eukaryotic translation regulation and how to find them. *Nat Rev Mol Cell Biol* **19**, 158–174 (2018).

41. van der Velden, A. W. & Thomas, A. A. M. The role of the 5′ untranslated region of an mRNA in translation regulation during development. *The International Journal of Biochemistry & Cell Biology* **31**, 87–106 (1999).

42. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).

43. Verkuil, R. *et al.* Language models generalize beyond natural proteins. 2022.12.21.521521 Preprint at https://doi.org/10.1101/2022.12.21.521521 (2022).

44. Goodstein, D. M. *et al.* Phytozome: a comparative platform for green plant genomics. *Nucleic Acids Research* **40**, D1178–D1186 (2012).

541    45.    Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low

542         memory requirements. *Nat Methods* **12**, 357–360 (2015).

543    46.    Love, M. I., Huber, W. & Anders, S. *Moderated Estimation of Fold Change and*

544         *Dispersion for RNA-Seq Data with DESeq2*. http://biorxiv.org/lookup/doi/10.1101/002832

545         (2014) doi:10.1101/002832.

546    47.    Zuker, M. On Finding All Suboptimal Foldings of an RNA Molecule. *Science* **244**,

547         48–52 (1989).

548    48.    Fish, L. *et al.* A prometastatic splicing program regulated by SNRPA1 interactions

549         with structured RNA elements. *Science* **372**, eabc7531 (2021).

550    49.    Steinbach, M., Karypis, G. & Kumar, V. A Comparison of Document Clustering

551         Techniques. (2000).

552
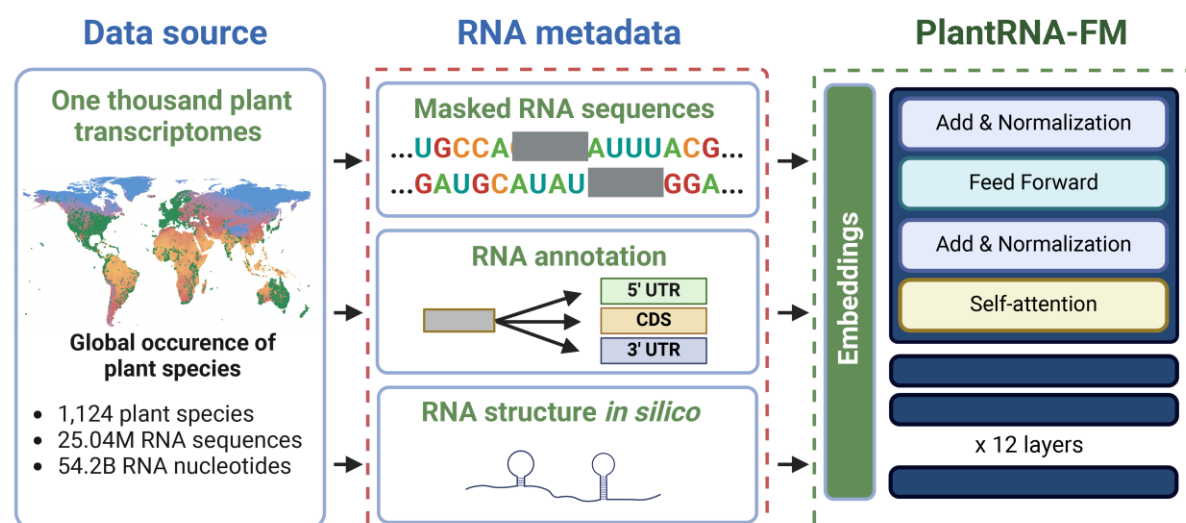553

## 7. Funding

566
567

# 8. TABLE AND FIGURES



**Fig. 1. Schematic overview of the Pre-training Phase of PlantRNA-FM.** The pre-training dataset comprises transcriptomic sequences from 1,124 plant species, consisting of approximately 25.0M RNA sequences and 54.2B RNA bases. The green dots on the global mean temperature map represent the geographical distribution of these plant species across the world.
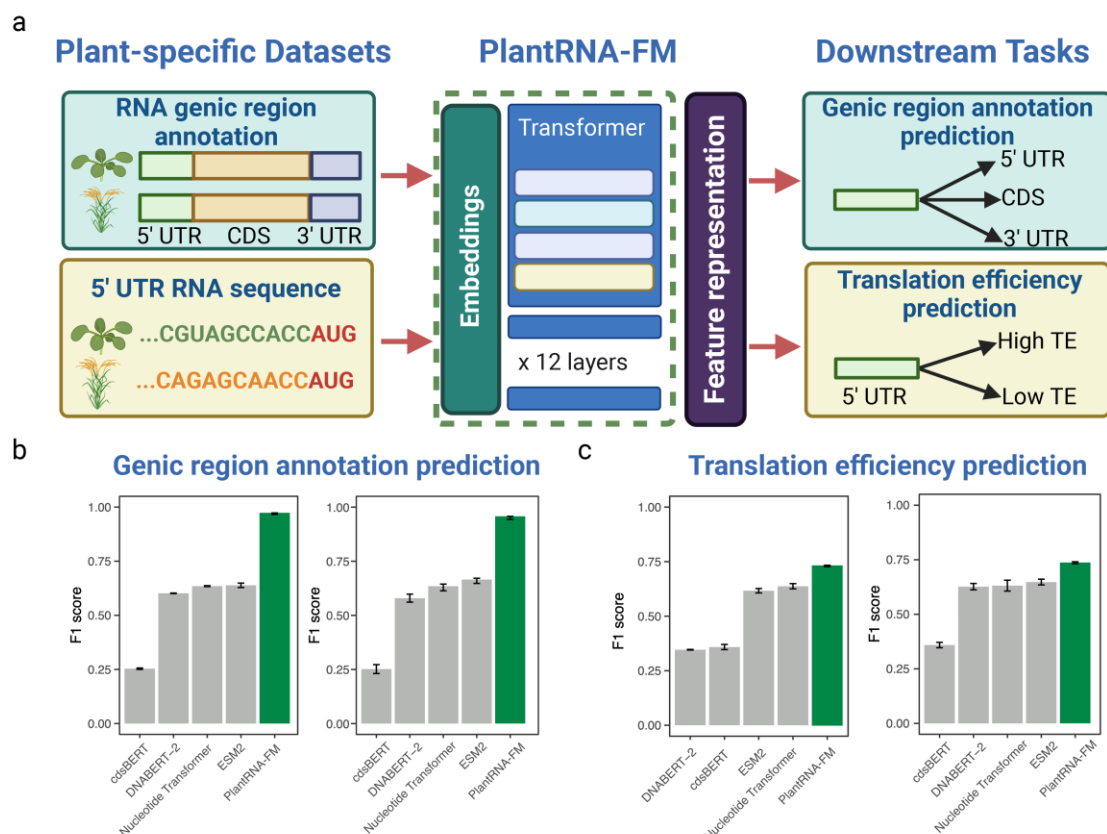
**Fig. 2. Fine-tuning PlantRNA-FM on plant-specific datasets.** a, Overview of fine-tuning PlantRNA-FM for RNA genic region annotation prediction and RNA translation efficiency (TE) prediction tasks. *A. thaliana* and *O. sativa* were selected as representative plant species. For the RNA genic region annotation prediction task, RNA sequences from these two species were included, along with three labels: 5' UTR, CDS, and 3' UTR. For the RNA TE prediction task, 5'UTR sequences from these two species were included, along with TE labels (high TE and low TE). b, c, Comparison of the model performance of different pre-trained models on RNA genic region annotation prediction and RNA TE prediction tasks. The error bars represent the standard deviation of the *F1* scores obtained from three fine-tuning replicates.
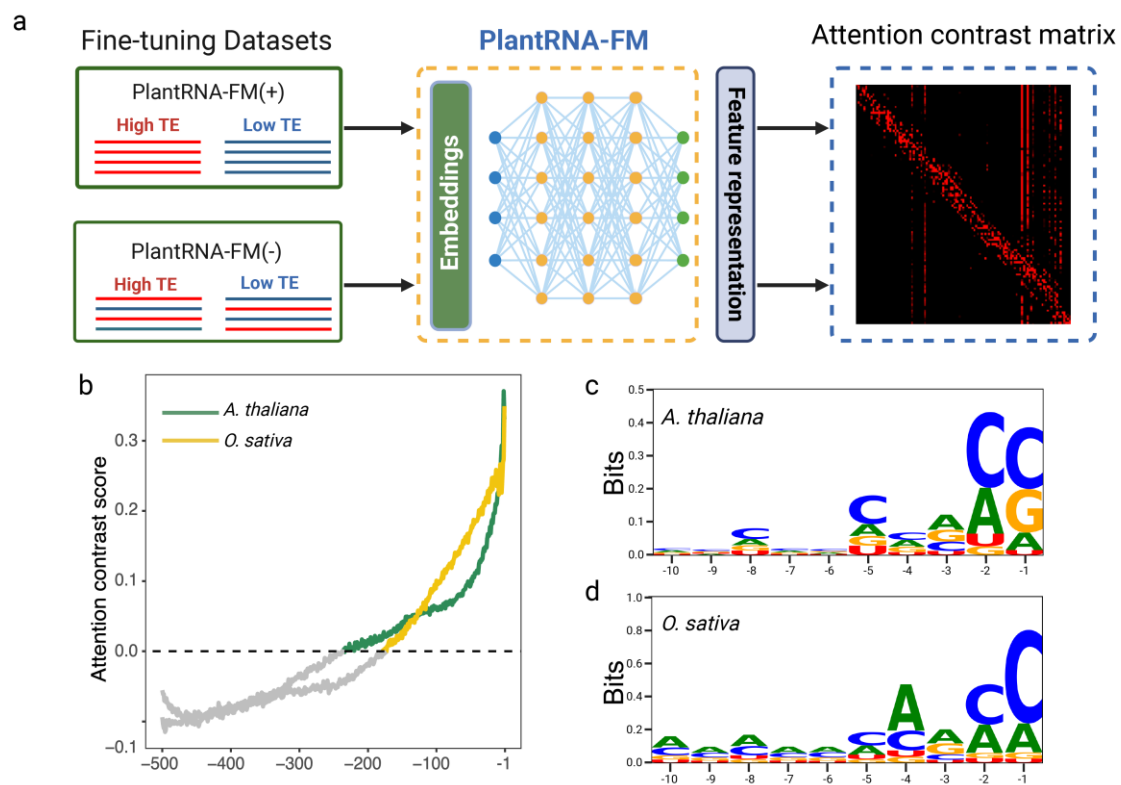
586

**Fig. 3. Our model interpretable framework reveals translation-associated RNA features.**

a, Schematic of the model interpretability approach. b, Transcriptome-wide attention contrast scores. The -1 position represents the first site upstream of the AUG. Different species are distinguished by colours. c, d, The information content of the 10 high-attention bases closest to the AUG start codon.
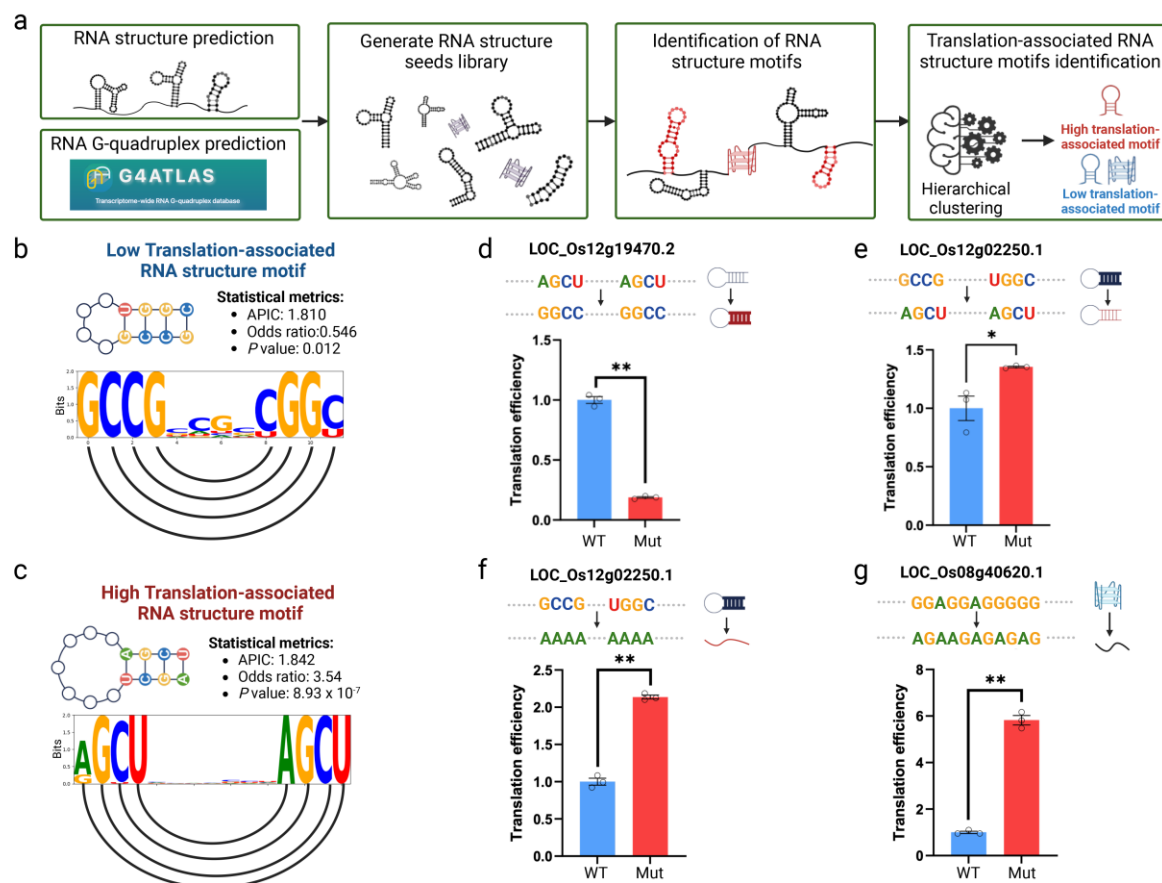
**Fig 4. RNA structure motif identification approach reveals translation-associated RNA structure motifs.** a, Overview of the RNA structure motif identification approach. RNA structures are predicted using RNAfold with a maximum length of 30 nucleotides to obtain RNA structure seeds. Predicted RNA G-quadruplexes were obtained from the G4Atlas database. b, c, Schematic diagram of high translation-associated RNA structure motifs and low translation-associated RNA structure motifs. Sequence logos show the information content of each nucleotide, with semicircles connecting paired bases. APIC stands for average positional information content. The *p-value* is derived from Fisher's exact test. d,e,f,g, Experimental validation of high and low translation-associated RNA structure motifs and low translation-associated RNA G-quadruplex. The bar plot represents the translational efficiency of the original (WT) and RNA structure-mutated (Mut) constructs from the dual luciferase reporter assay in plants. It represents the change from high translation-associated RNA structure motifs

606    to low translation-associated RNA structure motifs (d), the change from low translation-

607    associated RNA structure motifs to high translation-associated RNA structure motifs (e), the

608    complete disruption of low translation-associated RNA structure motifs (f), and the complete

609    disruption of low translation-associated rG4 (g). * indicates $P < 0.05$, ** indicates $P < 0.01$, by

610    Student's t-test, n = 3, error bars indicate se.

611 Table 1. Comparison of *F1* scores achieved by different pre-trained models on benchmark

612 datasets.

613

| Tasks | Species | PlantRNA-FM | cdsBERT | DNABERT-2 | Nucleotide Transformer | ESM2 |
|---|---|---|---|---|---|---|
| RNA genic region annotation prediction | *A. thaliana* | 0.974±0.003 | 0.254±0.003 | 0.602±0.001 | 0.635±0.002 | 0.639±0.008 |
| | *O. sativa* | 0.958±0.006 | 0.252±0.017 | 0.580±0.015 | 0.635±0.013 | 0.665±0.010 |
| RNA translation efficiency prediction | *A. thaliana* | 0.735±0.003 | 0.359±0.010 | 0.346±0.001 | 0.637±0.010 | 0.617±0.008 |
| | *O. sativa* | 0.737±0.004 | 0.359±0.010 | 0.627±0.012 | 0.631±0.020 | 0.649±0.011 |

614

615

616