

A combination of mRNA features influence the efficiency of leaderless mRNA translation initiation

Mohammed-Husain M. Bharmal¹, Alisa Gega¹, and Jared M. Schrader^{1,*}

¹ Department of Biological Sciences, Wayne State University, Detroit, Michigan, 48202, USA

* To whom correspondence should be addressed. Tel: 1 (313) 577-0736; Fax: 1 (313) 577-6981; Email: Schrader@wayne.edu

ABSTRACT

Bacterial translation is thought to initiate by base-pairing of the 16S rRNA and the Shine-Dalgarno sequence in the mRNA's 5'UTR. However, transcriptomics has revealed that leaderless mRNAs, which completely lack any 5'UTR, are broadly distributed across bacteria and can initiate translation in the absence of the Shine-Dalgarno sequence. To investigate the mechanism of leaderless mRNA translation initiation, synthetic *in vivo* translation reporters were designed that systematically tested the effects of start codon accessibility, leader length, and start codon identity on leaderless mRNA translation initiation. Using this data, a simple computational model was built based on the combinatorial relationship of these mRNA features which can accurately classify leaderless mRNAs and predict the translation initiation efficiency of leaderless mRNAs. Thus, start codon accessibility, leader length, and start codon identity combine to define leaderless mRNA translation initiation in bacteria.

INTRODUCTION

Translation initiation is a critical step for fidelity of gene expression in which the ribosome initiation complex is formed on the start codon of the mRNA. Since the canonical start codon, AUG, compliments both initiator and elongator methionyl-tRNAs, the ribosome must distinguish the start AUG codon from elongator AUG codons. Incorrect initiation at an elongator AUG can lead to non-functional products that can be detrimental to cellular fitness (1-3). Canonical start codon selection is thought to occur by the base-pairing of the 16S rRNA with a Shine-Dalgarno (SD) sequence in the mRNA located 5nt upstream of the start codon (4-6). The base pairing between the 16S rRNA and mRNA was shown to be critical for initiation since mutation of the anti-SD (aSD) in the 16S rRNA is lethal (7), and translation of a gene lacking a canonical SD sequence could be restored when the 16S of the rRNA were mutated to a complimentary sequence (8). While the SD-aSD pairing clearly impacts translation initiation efficiency (TIE) in *E. coli*, other studies have found that the SD:aSD interaction is not essential for correct selection of the start codon (9,10). Indeed, "orthogonal" ribosomes with altered 16S rRNA aSD sequences were found to initiate at the normal start codons throughout the transcriptome (11). Interestingly, *E. coli* lacks SD sites within its genome in approximately 30% of its translation initiation regions (TIRs) with other species of bacteria containing SD sites in as few as 8% of their TIRs (12,13). Indeed, RNA-seq based transcription mapping experiments have found that many bacterial mRNAs are "leaderless" and begin directly at the AUG start codon (14-16), and that these mRNAs are abundant in pathogens such as *M. tuberculosis* and in the mammalian mitochondria (17).

39 To account for the lack of essentiality of the SD site, a “Unique accessibility model” was
40 proposed which posited that start codon selection occurs due to the TIR being accessible to initiating
41 ribosomes, while elongator AUGs are physically inaccessible due to RNA secondary structures (18).
42 This model was based upon a strong negative correlation observed between mRNA secondary
43 structure content in the TIR and TIE (19-21). This model is further supported by genomic analysis of
44 RNA secondary structure prediction of mRNA TIRs in which there’s a lower amount of secondary
45 structure in the TIR compared to elongator regions, which is conserved across all domains of life (22).
46 While the unique accessibility model is overly simplistic, more advanced computational approaches
47 have been able to combine TIR accessibility with SD strength, spacing, and standby sites to more
48 accurately predict TIE of leadered mRNAs (23). While TIR accessibility has been shown to be critical
49 in many leadered mRNAs, it has not yet been systematically tested for leaderless mRNAs.

50 Genome-wide RNA-seq transcript mapping experiments have revealed that leaderless
51 mRNAs are widespread across bacteria (14), yet little is known about their mechanism of translation
52 initiation. While very few leaderless mRNAs has been identified in *E. coli* (0.7% leaderless mRNAs
53 (24)), other bacteria and archaea contain a large majority of their transcripts as leaderless mRNAs (up
54 to 72% leaderless mRNAs (14,25)). Additionally, sizeable proportions of leaderless mRNAs have
55 been identified in bacteria of clinical significance, such as *Mycobacterium tuberculosis*, and of
56 industrial significance like *Corynebacterium glutamicum* (15,26). In the model bacterium *Caulobacter*
57 *crescentus* approximately 17% of mRNAs are leaderless (27), with the fastest doubling time known of
58 any bacterium with large numbers of leaderless mRNAs. In addition, *C. crescentus* has good genetic
59 tools, making it an ideal model to study translation initiation of leaderless mRNAs.

60 Importantly, the role of TIR accessibility has not been systematically tested for leaderless
61 mRNAs, however, some aspects of their initiation have been identified which are distinct from
62 leadered mRNAs. Mitochondrial leaderless mRNAs have been found to lack 5' secondary structure
63 (28), in support of a TIR accessibility model. Additionally, mutagenesis of the *Mycobacterium*
64 *smegmatis* *pafl* leaderless mRNA to perturb its secondary structure showed that secondary structure
65 content negatively correlated with this translation levels (29). However, the changes in codon usage
66 across the mutants make the relative impact of secondary structure and codon usage unknown for
67 this mRNA. In opposition to the canonical initiation mechanism, leaderless mRNAs can initiate with
68 70S ribosomes where IF2 is known to stimulate their translation, and IF3 can inhibit leaderless
69 translation (30,31). Additionally, AUG is the most efficient start codon in leaderless mRNAs in *E. coli*
70 or *Haloarchaea*, (32-36), while AUG or GUG are both efficient leaderless mRNA start codons in *M.*
71 *smegmatis* (16). In *E. coli*, suppressor tRNAs could restore initiation on non-AUG codons for leadered
72 RNAs, but not for leaderless RNAs (32), suggesting that for leaderless mRNAs an AUG start codon
73 has unique initiation properties independent of perfect codon-anticodon base-pairing. Indeed,
74 genomic prediction of leaderless mRNAs suggests a very high preference of AUG (79%) at the 5' end
75 of leaderless mRNAs; with a smaller percentage of GUG (10%), UUG (6%) and others (3%) (13). In
76 addition to the start codon identity, TIE of mRNAs with short leaders (<5nt) is significantly lower as

77 compared to their fully leaderless counterparts (34,35,37-39). Altogether, this suggests that leaderless
78 mRNAs strongly prefer AUG and are inhibited by having short leaders.

79 In order to understand the mRNA sequence features needed for leaderless translation
80 initiation, we systematically measured the effect of TIR accessibility, start codon identity, and leader
81 length on leaderless mRNA translation initiation in *C. crescentus*. Using synthetic *in vivo* translation
82 initiation reporters, we show that TIR accessibility, start codon identity, and leader length all
83 dramatically affect leaderless mRNA TIE. The dependencies of each mRNA feature on TIE were then
84 built into a simple computational model (TIE_{leaderless} model) that accurately predicts which RNAs in the
85 *C. crescentus* transcriptome would be initiated as leaderless RNAs with an area under the curve
86 (A.U.C.) of a Receiver Operator Characteristic (ROC) curve of 0.99. The TIE_{leaderless} model also
87 accurately predicts the translation initiation efficiency of *in vivo* leaderless mRNA reporters ($R^2=0.87$).
88 This therefore provides the first systematic analysis of mRNA features required for leaderless initiation
89 and the *C. crescentus* TIE_{leaderless} model will likely provide a foundation for our understanding of
90 leaderless mRNA translation initiation across bacteria.

91 MATERIAL AND METHODS

92 ***Computational predictions of start codon accessibility***

93 **Retrieving transcript sequences**

94 All the RNA sequences were retrieved from transcription start sites and translation start site
95 data available from RNA-seq and ribosome profiling respectively (27,40) using the *C. crescentus*
96 NA1000 genome sequence (41). For *M. smegmatis*, RNA-seq and ribosome profiling data were
97 downloaded from the European Nucleotide Archive
98 (<https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-2929/>) and for *M. tuberculosis*, RNA-seq
99 data was obtained from Gene Expression Omnibus (GEO) accession number GSE62152 and
100 analyzed using the CP000480.1 and NC_000962 genome sequences respectively (16). For *H.*
101 *volcanii*, RNA-seq and ribosome profiling were provided by The DiRuggerio lab, and analyzed with the
102 *H. volcanii* NCBI RefSeq genome (taxonomy identification [taxid] 2246; 1 chromosome, 4 plasmids)
103 (42). For *M. musculus* mitochondria, RNA-seq and ribosome profiling data were downloaded from (43)
104 and analyzed with the NC_005089 genome sequence. The TIR sequences were then extracted from
105 all open reading frames (ORFs) using 50 nt (25 nt upstream of start codon and 25 nt downstream
106 from start codon). If the 5' upstream untranslated region (UTR) was less than 25 nt, then 50 nt from
107 transcription start site was used for all TIR calculations. Classification of mRNA type (leaderless, non-
108 SD, or SD) were obtained from Schrader *et al.* PLOS Genetics 2014.

109

110 **Translation Efficiency Data**

111 Ribosome profiling and RNA-seq Translation efficiency data were obtained for *C. crescentus*
112 from (27). Ribosome profiling and RNA-seq Translation efficiency data were obtained for *H. volcanii*
113 from the group of Prof. Jocelyn DiRuggerio (44). Ribosome profiling and RNA-seq sequencing data

114 were obtained for *M. smegmatis* from the European Nucleotide Archive
115 (<https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-2929/>)(16). For both ribosome profiling
116 and RNA-seq data the sequencing reads were downloaded as fastq files and the adapter poly-A
117 sequences were trimmed using a custom python script. Trimmed reads were then depleted of rRNA
118 and tRNA reads by alignment with bowtie (45), and the remaining non-rRNA/tRNA reads were then
119 aligned to the *M. smegmatis* MC2 155 genome. RPKM values were then calculated based upon the
120 CP009494.1 annotation. To avoid confounding effects from initiating or terminating ribosomes, the
121 first 15nt and last 15nt of ORFs were omitted from the RPKM calculations. ORFs with less than 50
122 reads in a given sample were omitted from the RPKM calculation. The Translation Efficiency (TE)
123 was then calculated as the ratio of the $RPKM_{\text{Ribosome profiling}}/RPKM_{\text{RNA-seq}}$.

124 **Calculation of ΔG_{unfold}**

125 Start codon accessibility was computed similar to (46) by comparing the native TIR RNA
126 structure (ΔG_{mRNA}) to that of the same TIR bound by an initiating ribosome (ΔG_{init}). Since ribosome
127 binding requires a single-stranded region of the mRNA we approximated this by forcing the TIR to be
128 single stranded. The overall calculation was performed in three steps:

129 1. Calculation of ΔG_{mRNA} :

130 The minimum free energy (mfe) labelled as ΔG_{mRNA} was calculated using RNAfold web server
131 of the Vienna RNA websuite (47) at the growth temperature of each organism by inputting all
132 the TIR sequences in a text file using command line function 'RNAfold --temp "temp"
133 <input_sequences.txt >output.txt'. The output file was in the default RNAfold format with each
134 new sequence on one line followed by dot-bracket notation (Vienna format) in the next line.
135 RNAstructure (48) was used to generate ct files for each of the mfe structures predicted in
136 RNAfold which contained all the base pair indexes for each sequence.

137 2. Calculation of ΔG_{init} :

138 The base pairs in the TIR (from up to 12 nt upstream of the start codon to 13 nt downstream
139 of the start codon) were broken and forced to be single stranded including any pairs formed
140 from the TIR and outside. If the 5'UTR length was more than or equal to 25 nt, then the RBS
141 was selected from -12 to +13 nt (25 nt). If the 5'UTR length was less than 25, then the TIR
142 comprised of the entire 5'UTR to +13 nt. A new dot bracket file with these base-pairing
143 constraints was then used in the RNAfold program (47) with the same RNA sequence to
144 calculate the ΔG_{init} .

145 3. Calculation of ΔG_{unfold} :

146 Lastly, ΔG_{unfold} was calculated by subtracting ΔG_{mRNA} (mfe of mRNA in native state) from
147 ΔG_{init} (mfe of mRNA after ribosome binding) (eq 1. $\Delta G_{\text{unfold}} = \Delta G_{\text{init}} - \Delta G_{\text{mRNA}}$).

148 **Cell growth and media**

149 ***E. coli* culture**

150 For cloning, plasmids with the reporter gene were transformed in *E. coli* top10 competent
151 cells using heat shock method for 50-55 secs at 42°C. Luria-Bertani (LB) liquid media was used for

152 outgrowth and the colonies were plated on LB/kanamycin (50 µg/mL) agar plates. For miniprep, the *E.*
153 *coli* cultures were inoculated overnight(O/N) in liquid LB/kanamycin (30 µg/mL).

154 ***C. crescentus* culture**

155 For cloning, plasmids were transformed in NA1000 *C. crescentus* cells after sequence
156 verification using electroporation. The *C. crescentus* NA1000 cells were grown in Peptone Yeast
157 Extract (PYE) liquid medium. After transformation, for the outgrowth liquid PYE medium was used
158 (2mL) and then plated on PYE/kanamycin (25 µg/mL) agar plates. For imaging, the *C. crescentus*
159 culture were grown O/N at different dilutions in liquid PYE/kanamycin (5 µg/mL). Next day, the
160 cultures growing in log phase were diluted and induced in liquid PYE with kanamycin (5 µg/mL) and
161 Xylose (final concentration of 0.2%) such that the optical density (OD) was around 0.05 to 0.1.

162 ***Design and generation of translation reporters***

163 **Oligos and plasmid design**

164 For the design and generation of reporter assay, a plasmid with a reporter gene (yellow
165 fluorescent protein (YFP)), under the control of an inducible xylose promoter was used. The pBYFPC-
166 2 plasmid containing the kanamycin resistant gene was originally generated from (49). A list of oligos
167 used for generating plasmids with different 5' UTRs of YFP is attached as a supplementary table
168 (Table S6).

169 **Inverse PCR mutagenesis and Ligation**

170 The 5'UTR region and start codon of the YFP reporter protein was replaced with other TIR
171 sequences. This was done by inverse PCR, in which the leaderless TIR is attached to the reverse
172 primer as an overhang. Initial denaturation was done at 98°C for 5 mins. Followed by 30 cycles of
173 denaturation at 98°C for 10 secs, annealing at 60°C for 10 secs and extension at 72°C for 7 mins and
174 20 secs. After 30 cycles, final extension was done at 72°C for 5 mins. The polymerase used was
175 Phusion (Thermoscientific 2 U/µL). The PCR product was then DPNI treated to cut the template DNA
176 using DPNI enzyme (Thermoscientific 10 U/µL). The DPNI treated sample was then purified using
177 Thermo fisher GeneJET PCR Purification kit. The purified sample (50 ng) was then used for blunt end
178 ligation using T4 DNA Ligase (Thermoscientific 1 WeissU/µL).

179 **Transcription reporter design**

180 For the design of transcription reporter assay, a plasmid with a 28 nt mutant version of 5' UTR
181 (CCNA_03971) in front of reporter gene (yellow fluorescent protein (YFP)) was used. This reporter
182 gene had the nucleotide A at its +1 position and the reporter gene was under the control of an
183 inducible xylose promoter. The pBYFPC-2 plasmid containing the kanamycin resistant gene was
184 originally generated from (49). The +1 nucleotide was mutated to all other nucleotides (G, C or T) and
185 these 3 mutant plasmids were synthesized into DNA oligos and cloned by Genscript.

186
187 The insertion sequence from the +1 nt (underlined) to 28 nt including start codon (atg) to the RBS for
188 each construct is shown below:

189 A.) accgattaacgatggtggttgttctggc
190 C.) cccgattaacgatggtggttgttctggc
191 G.) gccgattaacgatggtggttgttctggc
192 T.) tccgattaacgatggtggttgttctggc

193

194 **Transformation in *E. coli* cells**

195 5 μ L of the ligation reaction was then added to 50 μ L of *E. coli* top10 competent cells. Then
196 the mixture was incubated in ice for 30 mins. Then heat shocked for 50-55 secs in the water bath at
197 42°C. Then immediately kept in ice for 5 mins, after which 750 μ L of LB liquid medium was added to
198 the cells for outgrowth and kept for incubation at 37°C for 1 hr at 200 rpm. After this, 200-250 μ L of
199 the culture was plated on LB/kanamycin (50 μ g/mL) agar plates.

200 **Colony screening and sequence verification**

201 The colonies grown on LB/kanamycin plates were screened by colony PCR to first screen for
202 the presence of the new TIR insert. The cloning results in the replacement of the larger 5'UTR region
203 of YFP with a smaller region containing a leaderless TIR, thus distinguished easily on an analytical
204 gel. The forward and reverse primer used for the screening results in approximately 180 base pairs,
205 whereas the original fragment amplified with the same oligos is 245 base pairs. The forward oligo
206 used was pxyl-for: cccacatgtagcgctaccaagtgc and reverse oligo is eGYC1: gtttacgtcgccgtccagctcgac.
207 Upon verification, a small aliquot (4 μ L) of the colony saved in Taq polymerase buffer was inoculated
208 in 5 mL of liquid LB/kanamycin (30 μ g/mL) and incubated overnight at 37°C at 200 rpm. Next day, the
209 culture was miniprep using Thermo fisher GeneJET Plasmid Miniprep kit. The concentration of
210 DNA in the miniprep samples were measured using Nanodrop 2000C from ThermoScientific. DNA
211 samples were sent to Genewiz for sanger sequencing to verify the correct insert DNA sequences
212 using the DNA primer eGYC1: gtttacgtcgccgtccagctcgac (49).

213 **Transformation in *C. crescentus* NA1000 cells**

214 After the sequences were verified, the plasmids were transformed in *C. crescentus* NA1000
215 cells. For transformation, the NA1000 cells were grown overnight at 28°C in PYE liquid medium at
216 200rpm. The next day, 5 mL of cells were harvested for each transformation, centrifuged and washed
217 three times with autoclaved milliQ water. Then, 1 μ L of sequence verified plasmid DNA was mixed
218 with the cells and electroporated using Bio-Rad Micropulser (program Ec1 set at voltage of 1.8 kV).
219 Then, the electroporated cells were immediately inoculated in 2 mL of PYE for 3 hours at 28°C at
220 200rpm. Then 10-20 μ L of culture was plated on PYE/ kanamycin agar plates. Kanamycin-resistant
221 colonies were grown in PYE/kanamycin media overnight and then stored as a freezer stock in the -
222 80°C freezer

223 **Cellular assay of translation reporters**

224 *C. crescentus* cells harboring reporter plasmids were serially diluted and grown overnight in
225 liquid PYE/kanamycin medium (5 μ g/mL). The next day, cells in the log phase were diluted with fresh
226 liquid PYE/kanamycin (5 μ g/mL) to have an optical density (OD) of 0.05-0.1. The inducer xylose was
227 then added in the medium such that the final concentration of xylose is 0.2%. The cells were grown
228 for 6 hours at 28°C at 200 rpm. After this, 2-5 μ L of the cultures were spotted on M2G 1.5% agarose
229 pads on a glass slide. After the spots soaked into the pad, a coverslip was placed on the pads and the
230 YFP level was measured using fluorescence microscopy using a Nikon eclipse NI-E with CoolSNAP
231 MYO-CCD camera and 100x Oil CFI Plan Fluor (Nikon) objective. Image was captured using Nikon
232 elements software with a YFP filter cube with exposure times of 30ms for phase-contrast images and

233 300 ms for YFP images respectively. The images were then analyzed using a plugin of software
234 ImageJ (50) called MicrobeJ (51).

235 **Three component model calculations and leader length/identity analysis**

236 For all RNA transcripts in the *C. crescentus* genome identified in (27,40), we computed their
237 capacity to initiate as a leaderless mRNA using equation 2: ($TIE_{\text{Leaderless mRNA}(k)} = \text{Max } TIE(1) - (1 -$
238 $TIE_{\Delta G_{\text{unfold}}}) - (1 - TIE_{\text{start codon identity}(j)}) - (1 - TIE_{\text{leader length}(i)})$) where k = a given RNA transcript, j =start codon
239 identity, and i =leader length(nt). To identify putative leaderless mRNA TIRs, we first asked if the 5'
240 end contained an AUG or near cognate start codon, and if not we scanned successively from the 5'
241 end for AUG trinucleotides within the first 8 nt. Near cognate start codons were omitted from positions
242 containing leader nucleotides since AUG codons yielded higher TIE values even in the presence of a
243 leader. We next asked if there is an AUG or near cognate start codon further downstream by
244 scanning 5' to 3' through the first 18 nt. If found, we calculated $TIE_{\text{leaderless mRNA}}$ with all different
245 possible cognate/near-cognate start codons along the TIR. Then of all the different possibilities, the
246 one having the highest $TIE_{\text{leaderless}}$ score was selected for further analysis (Fig 7A).

247 To utilize $TIE_{\text{leaderless mRNA}}$ for classification, each RNA was then categorized into two different
248 classes based on 5' end sequencing data and ribosome profiling based global assays ((27,40)): true
249 leaderless – RNAs that are known to initiate directly at a 5' start codon (judged by a complete lack of
250 a 5' UTR and a ribosome density >1/20 the downstream CDS (27)), and false leaderless – RNAs that
251 are not initiated at a 5' start codon. A small subset was classified as “unknown”, as they contain very
252 short leaders and lack SD sites, making their mode of translation initiation ambiguous. Using these
253 $TIE_{\text{leaderless mRNA}}$ values, a ROC curve was plotted using scikit-learn library in python (52) with the “true
254 leaderless” and “false leaderless” RNAs ($TIE_{\text{leaderless mRNA}}$ values for the *C. crescentus* transcriptome
255 can be found in Table S1).

256 To utilize $TIE_{\text{leaderless mRNA}}$ for prediction of translation initiation reporter levels, we first
257 converted all negative $TIE_{\text{leaderless mRNA}}$ scores to zero. Next, we compared the $TIE_{\text{leaderless mRNA}}$ scores
258 to the YFP levels of the translation initiation and performed a linear regression calculation using the
259 linest function in microsoft excel and libreoffice calc. For prediction of native leaderless mRNA
260 translation levels, TE measurements from ribosome profiling experiments (27) were compared to the
261 $TIE_{\text{leaderless mRNA}}$ scores.

262 **RESULTS**

263 **Computational prediction of *C. crescentus* start codon accessibility**

264 To assess the role of mRNA accessibility across mRNA types, ΔG_{unfold} calculations were
265 performed on all *C. crescentus* translation initiation regions (TIRs). ΔG_{unfold} represents the amount of
266 energy required by the ribosome to unfold the mRNA at the translation initiation region (TIR) and has
267 been identified as a metric that correlates with translation efficiency in *E. coli* (46). ΔG_{unfold} was
268 calculated for all TIRs by first predicting the minimum free energy of the 50 nt region of the mRNA
269 (ΔG_{mRNA}) around the start codon using RNAfold (47). ΔG_{init} was then calculated in which the TIR (25nt
270 surround the start codon), roughly equivalent to a ribosome footprint, was constrained to be single

271 stranded to approximate accessibility for the ribosome to initiation. ΔG_{unfold} was then calculated using
272 equation 1 (eq 1. $\Delta G_{\text{unfold}} = \Delta G_{\text{init}} - \Delta G_{\text{mRNA}}$) which represents the energy required to open the TIR to
273 facilitate translation initiation (Fig 1A). ΔG_{unfold} calculations were performed on all the CDSs in the
274 genome (Fig 1B) and classified into mRNA types based on transcriptome and ribosome profiling
275 maps of the *C. crescentus* genome (27). The transcripts were categorized into two major classes:
276 leaderless (no 5' UTR) and leadered (those containing a 5' UTR). Leadered mRNAs were further
277 categorized into subclasses based upon the presence of the Shine-Dalgarno (SD) sequence (27).
278 Shine-Dalgarno (SD) (containing a SD sequence in the 5' UTR) and nonSD (lacking an SD sequence
279 in the 5' UTR). Since it is also known that some polycistronic operons reinitiate translation between
280 CDSs without dissociation of the ribosomal subunits, we also examined the ΔG_{unfold} of TIRs occurring
281 downstream of the first CDS in polycistronic mRNAs (Operons). The average ΔG_{unfold} value of
282 leaderless mRNAs (5.6 kcal/mol) was significantly lower than SD (11.9 kcal/mol, $p= 1.5E-105$), nonSD
283 (10.3 kcal/mol, $p= 6.9E-71$) and internal operon TIRs (13.2 kcal/mol, $p= 1.1E-143$) as calculated by
284 pairwise 2-sided T-tests with unequal variance (Fig 1B). The lower ΔG_{unfold} values of nonSD TIRs may
285 be due to the loss of stabilization of TIRs from base pairing between the anti-SD site in the 16S rRNA
286 and the SD site in the mRNA. We also observed that average ΔG_{unfold} of nonSD TIRs was significantly
287 lower than SD TIRs ($p= 1.8E-14$) and operon TIRs ($p=1.4E-44$). The difference between the average
288 ΔG_{unfold} of SD and operon genes was also significant ($p=2.1E-09$). Because the ribosome is an
289 efficient RNA helicase, it is possible that the increased ΔG_{unfold} of operon TIRs may be tolerated by the
290 ribosome's ability to unwind such structures when terminating on the previous CDSs. We
291 hypothesized that the low ΔG_{unfold} observed for leaderless mRNAs was due to an intrinsic requirement
292 for their initiation, however, because the size of the leaderless mRNA footprint is significantly smaller
293 than a leadered mRNA footprint, the low ΔG_{unfold} observed for leaderless mRNAs could potentially be
294 explained by the smaller ribosome footprint size. To explore this possibility, we analyzed the ΔG_{unfold}
295 of leadered mRNA TIRs using the same footprint size and region as leaderless mRNAs (13nt) in the
296 ΔG_{unfold} calculation (Fig S1). We observed that the ΔG_{unfold} was still significantly lower for leaderless
297 mRNA TIRs, suggesting that the low ΔG_{unfold} for leaderless mRNA TIRs is not simply an artifact of the
298 smaller mRNA footprint size.

299 To explore whether low ΔG_{unfold} for leaderless mRNA TIRs a species-specific property of *C.*
300 *crescentus*, or a general property of leaderless mRNA TIRs, we calculated ΔG_{unfold} for TIRs in other
301 organisms identified to contain a significant number of leaderless mRNAs. We identified two
302 additional bacteria (*Mycobacterium smegmatis* and *Mycobacterium tuberculosis*) (16), one archaeal
303 species (*Haloferox volcanii*) (25,44), and one mitochondrial genome (*Mus musculus*) (43) which had
304 transcriptome information and ribosome profiling or mass spec data supporting a significant number
305 of leaderless mRNAs. Across bacteria and archaea, the leaderless mRNA ΔG_{unfold} remained quite low
306 as compared to their leadered mRNAs counterparts (leaderless average 5.6 to 7.6 kcal/mol, leadered
307 average 12.0-12.9 kcal/mol) (Fig 1C). In *M. musculus* mitochondria however, leaderless mRNAs and
308 leadered mRNAs were both observed to have low ΔG_{unfold} (Fig 1C), perhaps in part due to the
309 relatively low GC%. Since leaderless mRNAs showed a rather low ΔG_{unfold} , and lack complexities

310 associated with leadered mRNAs, such as SD or standby sites which are important for leadered
311 initiation (23), we further explored the functional role of ΔG_{unfold} in *C. crescentus* leaderless mRNAs.
312

313 **Systematic analysis of *C. crescentus* leaderless mRNA TIR determinants using *in vivo***
314 **translation reporters**

315 Leaderless mRNAs initiation is known to be strongly influenced by addition of nucleotides
316 prior to the start codon (leader nts) and by start codon identity (32-39); however, the role of TIR
317 accessibility has been poorly described in this class of mRNAs. To understand the role of these three
318 mRNA features we systematically tested each feature using *in vivo* leaderless mRNA translation
319 initiation reporters. Translation initiation reporters were designed in which the start codon of plasmid
320 pBXYFPC-2 was replaced with an AUG fused directly to the +1 nt of the xylose promoter(49). The
321 xylose promoter was chosen because it is one of the best characterized promoters in *Caulobacter*
322 and its TSS was mapped to the same nt by two independent methods(40). An additional 15-24 nt
323 after the 5' AUG was added to allow complete replacement of the 5' leader and start codon in
324 pBXYFPC-2 with a leaderless TIR. Since only the first 6-9 codons are altered across leaderless
325 mRNA mutants, and the vast majority of the YFP CDS is unaltered, this allows a sensitive system to
326 measure changes in translation initiation. As leaderless TIR mutants may also alter the amino acid
327 sequence, additional care was also taken to ensure that mutations would not alter the N-end rule
328 amino acid preferences of the resulting proteins (53). Using this *in vivo* translation initiation system,
329 we generated three different sets of leaderless TIR reporters to test the effect of ΔG_{unfold} , start codon
330 identity, and additional leader length on *C. crescentus* translation initiation.

331 As leaderless mRNAs were predicted to have TIRs with low ΔG_{unfold} values, we engineered
332 several RNA hairpins in the TIR to assess the role of ΔG_{unfold} on translation initiation (Fig 2A). Since
333 very few natural *C. crescentus* mRNAs contained RNA structure content in their TIRs (Fig 1B), six
334 synthetic hairpins were designed, varying in stem and loop sizes (Table S2). Into each construct, we
335 also introduced synonymous codon mutations designed to alter the secondary structure content,
336 yielding a range of ΔG_{unfold} values without altering the amino acid sequence within a given hairpin
337 (Table S2). Importantly, the entire range of ΔG_{unfold} values across the synthetic hairpins spans the
338 entire range calculated for natural leaderless mRNAs (Fig 1, Table S2). For all hairpins, we observed
339 that lowering ΔG_{unfold} and thereby increasing the accessibility of the start AUG led to an increase in
340 the level of YFP production (Fig 2B). Since, 6/7 of the hairpin mutant sets showed a relationship in
341 which hairpin codon usage frequency positively correlated with ΔG_{unfold} (Table S2), it is most likely that
342 the observed reduction in YFP reporter levels is a result of increased structure content and is not
343 likely to be caused by faster elongation of common codons in the TIR. Additionally, across all mutant
344 hairpins sets generated, we observed a strong negative correlation between the YFP reporter level
345 and the ΔG_{unfold} across a vast range of values with a linear correlation R^2 value of 0.84 (Fig 2B).
346 These data suggest that accessibility of the start codon is a critical feature for leaderless mRNA
347 translation initiation.

348 Next, we systematically tested the effect of the start codon identity on the *in vivo* translation
349 initiation reporters. In *C. crescentus*, natural leaderless mRNAs initiate with an AUG, GUG, or UUG
350 start codon (27,40). Since it is well established that start codon identity can affect leaderless mRNA
351 translation initiation (32-36) we generated variants with different start codon identities. Here, AUG was
352 mutated to other near cognate start codons GUG, CUG, UUG, AUC, AUU, AUA which are known to
353 be the start codons of other leadered mRNAs in *C. crescentus* (27). We also included a non-cognate
354 GGG codon as a negative control since no GGG start codons are known to occur in *C. crescentus*.
355 The results showed that replacing the original AUG codon with any of the other near cognate codons
356 drastically decreased the translation initiation reporter levels, while the GGG codon yielded the lowest
357 translation initiation reporter levels (Fig 3). To examine whether the mutation in +1 nt resulted from
358 lower transcription or from lower translation, we generated leadered mRNA reporters with all 4
359 possible +1 nts and tested their *in vivo* reporter levels as similarly has been performed in *M.*
360 *smegmatis* (16). A +1 G led to a mild reduction in reporter activity compared to a +1 A (Fig S5),
361 suggesting that most of the observed changes in the leaderless mRNA reporter levels likely come
362 from translation. These data show that the AUG triplet is by far the preferred start codon for *C.*
363 *crescentus* leaderless mRNAs.

364 Finally, we systematically tested the role of additional leader length on *C. crescentus*
365 leaderless mRNAs. In *E. coli*, even a single nucleotide before the AUG is known to inhibit initiation of
366 leaderless mRNAs (37). To test if *C. crescentus* leaderless mRNAs were negatively impacted by
367 leader nucleotides we generated a set of reporters with 0, 1, 2, 3, 5, 10, or 20 5' Adenosines before
368 the AUG start codon (Fig 4). An A-rich sequence was chosen as it lacks any possible SD sites and is
369 unlikely to form secondary structure, and ΔG_{unfold} values were not altered upon addition of these 5'
370 bases to the leaderless translation initiation reporter (Table S1). Across this set of mutants, additional
371 nucleotides showed a strong decrease in translation initiation reporter levels with increasing leader
372 length (Fig 4). The translation initiation reporter levels dropped by approximately 2-fold for each
373 additional A that was added to the 5' end ($\text{TIE}_{\text{leader length}} = 0.45 \times i^{-0.91}$, $R^2=0.92$, $i=\text{leader length (nt)}$). This
374 confirms that even a short leader can lead to a significant reduction in translation initiation of *C.*
375 *crescentus* leaderless mRNAs.

376 **Leaderless mRNA TIR determinants affect translation efficiency of natural leaderless mRNAs**

377 Because the *in vivo* translation initiation reporters were all synthetic constructs, we explored
378 the extent to which each mRNA feature (ΔG_{unfold} , start codon identity, and leader length) occur in
379 natural *C. crescentus* leaderless mRNAs. As noted previously, ΔG_{unfold} is significantly lower for
380 leaderless mRNAs than for other mRNA types (Fig 1B). To analyze the role of start codon selection,
381 we calculated the fraction of AUGs at the 5' end of all *C. crescentus* leaderless mRNAs and of the
382 random chance of finding each start codon based on the genomes' GC percentage. This analysis
383 revealed a strong enrichment of AUGs at the 5' end of *C. crescentus* leaderless mRNAs as compared
384 to random, and a slight enrichment of the GUG near cognate start codons (Fig 5A). While GUG TIR
385 reporters yielded similar TIR reporter levels to UUG and CUG, it is possible that the lack of U and C at
386 the +1 of *C. crescentus* leaderless mRNAs is due to their low abundance in TSSs (40). Of all the

387 leaderless mRNAs, only 4.4% (17/385) are initiating with non-AUG start codons as compared to the
388 leadered mRNAs of which 27.23% (989/3632) of genes initiate with non-AUG start codons (Table S3).
389 Since these near cognate start codons were translated much more poorly than AUG in our translation
390 initiation reporters, it's possible that for leaderless mRNAs there's a positive selection for the AUG
391 start codon and a negative selection for near-cognate start codons. Additionally, by exploring the
392 length of mRNAs, we noticed that there was a much greater occurrence of leaderless mRNAs than
393 mRNAs with short leaders <10nt (Fig 5B). Additional leader nucleotides were strongly inhibitory of
394 leaderless translation, and only 8 contain SD motifs, suggesting some of these short-leadered
395 mRNAs may be poorly initiated.

396 To estimate the effects of each mRNA feature (ΔG_{unfold} , start codon identity, and leader length)
397 on natural leaderless mRNA translation, we next analyzed ribosome profiling data of the *C.*
398 *crenscentus* mRNAs (27). Here, we utilized translation efficiency measurements which approximate the
399 relative number of ribosome footprints to mRNA fragments from the same cell samples (54). In total,
400 translation efficiency data for 191 leaderless mRNAs and 38 short leadered mRNAs (1-10 leader
401 length) were obtained for cells grown in PYE media (27). We separated leaderless mRNAs into three
402 groups based upon their ΔG_{unfold} values (0-5, 5-10, and >10 kcal/mol) and compared their translation
403 efficiency. The median translation efficiency was reduced as the ΔG_{unfold} increased (Fig 5C) (median=
404 1.2 for 0-5 kcal/mol, median= 0.89 for 5-10 kcal/mol, median= 0.54 for >10 kcal/mol), similar to the
405 dependence observed in the synthetic translation reporters (Fig 2B). For start codon identity, we
406 noticed that a majority of leaderless mRNAs with near-cognate start codons had translation
407 efficiencies that were not measurable because their genes contained an additional upstream TSS.
408 However, for the 7 GUG mRNAs whose translation efficiency was measured, the median (0.70) was
409 lower than that of the AUG initiated leaderless mRNAs median (0.97) (Fig 5D), in line with the findings
410 of the synthetic reporters (Fig 3). Finally, we compared the translation efficiency of leaderless mRNAs
411 with those with very short leaders (Fig 5E). Since 8 of these mRNAs with short leaders contain SD
412 sequences in the leader, we removed these RNAs from the analysis because we expect them to
413 initiate translation by the canonical mechanism. As leader length increases, we generally observed
414 that the TE tends to decrease (Fig 5E), again in line with the synthetic reporters (Fig 4). Overall these
415 data suggest that the effects of ΔG_{unfold} , start codon identity, and leader length observed in the
416 synthetic translation initiation reporters are also observed across natural *C. crescentus* leaderless
417 mRNAs.

418 Many RNAs present in the *C. crescentus* transcriptome are not initiated as leaderless mRNAs,
419 so we explored the relative fraction of 5' AUG trinucleotides in all classes of RNAs (Fig 6A). As noted
420 previously, leaderless mRNAs are highly enriched in AUG codons (Fig 5A). Surprisingly, leadered
421 mRNAs contain a similar fraction of 5' AUGs as would be predicted from the genome's GC%, which is
422 also observed in small non-coding RNAs (sRNAs), and anti-sense RNAs (asRNAs). Conversely,
423 tRNAs and rRNAs contain zero cases with a 5' AUG. To explore why these RNAs are not initiated as
424 leaderless mRNAs, we calculated the ΔG_{unfold} of each class of 5' AUG containing RNA (Fig 6B). If
425 these 5' AUGs found in non-leaderless RNAs were inaccessible to ribosomes, it would be permissible

426 for this sequence to be present at the 5' end without causing aberrant initiation. Indeed, for the RNAs
427 with 5' AUGs, we observe that leaderless mRNAs have a low ΔG_{unfold} (median= 5.0), while leadered
428 mRNAs (median= 9.5), sRNAs (median = 14), and asRNAs (median = 9.0) all contain a significantly
429 higher ΔG_{unfold} values. This suggests that RNAs with inaccessible 5' AUGs are blocked from
430 leaderless mRNA initiation.

431 Due to the strong involvement of ΔG_{unfold} , start codon identity, and leader length in *C.*
432 *crenscentus* leaderless mRNA TIRs, we also explored these features across organisms. As already
433 shown in figure 1C, the ΔG_{unfold} for leaderless mRNAs is markedly lower than leadered mRNAs in all
434 species analyzed with the exception of *M. musculus* mitochondria. The low ΔG_{unfold} in the
435 mitochondria may be due to alterations in the translation initiation mechanism of leadered mRNAs by
436 their highly proteinaceous ribosomes (55). Across these organisms, ribosome profiling and total-
437 RNA-seq performed in *M. smegmatis* (16) and *H. volcanii* (44), allowing the comparison of how
438 ΔG_{unfold} tracks with translation efficiency. As observed for *C. crescentus*, as ΔG_{unfold} increases, a drop
439 in the TE in leaderless mRNAs in both *M. smegmatis* and *H. volcanii* (Fig S4). We hypothesize that
440 the low overall ΔG_{unfold} observed for leaderless mRNAs across all species suggests that ribosome
441 accessibility is a key feature for leaderless mRNAs across organisms. In *M. smegmatis* and *H.*
442 *volcanii* the observed drop in TE from the 0-5 bin and 5-10 bins were smaller than observed for *C.*
443 *crenscentus*, which may be due to the higher growth temperatures of these organisms. Next, we
444 explored the distributions of leader lengths across species (Table S4). As observed in *C. crescentus*,
445 mRNAs with short leaders have a highly skewed distribution across species, with leaderless mRNAs
446 showing the largest peak, with a small fraction of mRNAs containing a 1nt leader, and a much smaller
447 population of mRNAs observed with a ≥ 2 nt leader (Table S4). The lower abundance of mRNAs with
448 very short leaders is likely to be due to their poorer translation levels observed across organisms
449 (FigS4) as this has even been observed with *E. coli* leaderless TIR reporters (37). To explore this
450 possibility, we compared TE across organisms. As observed in *C. crescentus*, *M. smegmatis* and *H.*
451 *volcanii* TE is also markedly decreased in mRNAs containing short leaders (Fig S4). While *C.*
452 *crenscentus* leaderless mRNA TE was not significantly lower than the 1-5nt bin, both *M. smegmatis*
453 and *H. volcanii* showed sharper drops in TE. While *C. crescentus* mRNAs with 6-10nt leaders
454 showed a significant drop in TE, neither *M. smegmatis* or *H. volcanii* showed a significant decrease.
455 This discrepancy may be explained by a low sample size in *M. smegmatis*, where only 3 mRNAs were
456 identified in the 5-10nt bin, while in *H. volcanii* the 5-10nt bin distribution contained a single outlier
457 whose TE was measured to be >30 . Finally, we examined start codon identities for leaderless mRNAs
458 across species (Table S5). Here only *H. volcanii* was similar to *C. crescentus* with a strong bias in the
459 AUG start codon for leaderless mRNAs (Table S5). *M. tuberculosis* and *M. smegmatis* both
460 contained GUG start codons in leaderless mRNAs with similar abundance to AUG (Table S5). In the
461 *M. musculus* mitochondria, AUG is the most common start codon across mRNAs, however, AUG is
462 less common in leaderless mRNAs, and the summed use of GUG, AUC, AUU, and AUA makes near-
463 cognate start codons more abundant than AUG initiated leaderless mRNAs (Table S5). Interestingly,
464 mitochondrial RNAP has been found to initiate transcription efficiently with NAD⁺ and NADH (56),

465 which has the potential to alter start codon selection by the translation machinery. As observed in *C.*
466 *crescentus*, *M. smegmatis* and *H. volcanii* both showed a lower TE for leaderless mRNAs starting with
467 GUG as compared to AUG (Fig S4). The magnitude of the reduced TE for GUG initiated leaderless
468 mRNAs is significantly smaller in *M. smegmatis* (1.1 AUG, 0.91 GUG) as compared to *H. volcanii* (2.5
469 AUG, 0.82 GUG), which is in line with previous data showing that GUG initiates with similar efficiency
470 to AUG in leaderless mRNAs in Mycobacteria (16). Overall, these data suggest that the effects of
471 ΔG_{unfold} , start codon identity, and leader length have similar effects on leaderless mRNA translation
472 across species. However, minor idiosyncratic differences in the frequency and magnitude of each
473 leaderless mRNA feature on TE were observed across species, likely arising from differences in the
474 translation machinery.

475 **Three component model describes leaderless mRNA start codon selection**

476 In order to understand how the mRNA determinants combine to dictate leaderless mRNA
477 translation, we built a computational model based upon the three features (ΔG_{unfold} , start codon
478 identity, and leader length) and explored its ability to describe leaderless mRNA start codon selection
479 and efficiency of leaderless mRNA translation initiation. From our synthetic *in vivo* translation initiation
480 reporters, we performed curve fitting to assess the relative effect of each feature on TIE. For each
481 feature (ΔG_{unfold} , start codon identity, and leader length) the highest reporter level measured in each
482 mutant set was normalized to 1 before curve fitting. ΔG_{unfold} data was fit to an exponential equation
483 ($\text{TIE}_{\Delta G_{\text{unfold}}} = e^{(-t^{+0.354})}$) where t is ΔG_{unfold} (kcal/mol), $R^2 = 0.78$), leader length data was fit to a power
484 equation ($\text{TIE}_{\text{leader length}} = 0.45 \times (i^{-0.92})$ where i is leader length >0 , $R^2 = 0.92$, and $\text{TIE}_{\text{leader length}}=1$ for $i=0$),
485 and $\text{TIE}_{\text{start codon}}$ was based directly on reporter levels for each near-cognate start codon (Fig 3) and all
486 other codons were given a value of 0 (Table S1). For each mRNA feature, we therefore generated a
487 function that could calculate the relative TIE of any RNA in *C. crescentus* based upon the mRNA
488 sequence. We then built a computational model in which the three features were assumed to be
489 independent from each other to calculate a summed TIE. In this model, we set the maximum TIE to 1,
490 and then subtracted the effects of the sequence feature as measured from the *in vivo* translation
491 reporters in equation 2 ($\text{TIE}_{\text{Leaderless mRNA}(k)} = \text{Max TIE} (1) - (1 - \text{TIE}_{\Delta G_{\text{unfold}}}) - (1 - \text{TIE}_{\text{start codon identity}(j)}) - (1 -$
492 $\text{TIE}_{\text{leader length}(i)})$ where k = a given RNA transcript, j =start codon identity, and i =leader length(nt). Using
493 equation 2 we predicted the TIE for each RNA in the *C. crescentus* transcriptome (Fig 7A). For all
494 RNAs, we successively scanned for the closest AUG or near cognate start codon to the 5' end and
495 used this for the TIE calculation. RNAs known to be initiated as leaderless mRNAs (27,40) yielded
496 higher TIE scores (median = 0.15, $\sigma = 0.35$), while TIE scores for all other RNAs were typically lower
497 (median = -0.95, $\sigma = 0.45$). To estimate the utility of this model at classifying leaderless mRNAs, we
498 used a ROC analysis (Fig 7B). The ROC area under the curve for the $\text{TIE}_{\text{leaderless}}$ model was equal to
499 0.99, which significantly outperforms identifying RNAs with 5' AUGs (ROC A.U.C. 0.68) suggesting
500 the $\text{TIE}_{\text{leaderless}}$ model can accurately classify those RNAs that are initiated as leaderless mRNAs with
501 high accuracy and precision. The success of this simple $\text{TIE}_{\text{leaderless}}$ model to classify leaderless
502 mRNAs based on the combinations of ΔG_{unfold} , start codon identity, and leader length suggests that
503 these mRNA features combinatorically control translation initiation on leaderless mRNAs.

504 In addition to the classification of RNAs as leaderless mRNAs we also explored how well the
505 $TIE_{\text{leaderless}}$ model predicted translation initiation efficiency. Here, the translation initiation reporters
506 generated were all scored with the $TIE_{\text{leaderless}}$ model and compared to their YFP fluorescence. Since
507 $TIE_{\text{leaderless}}$ scores below zero are not physically possible, those with negative $TIE_{\text{leaderless}}$ values were
508 set to zero to signify they are not predicted to be translated. As expected, the $TIE_{\text{leaderless}}$ score
509 correlates strongly to the YFP reporter levels ($R^2=0.87$) with a slope of 2050 A.U. We then compared
510 the $TIE_{\text{leaderless}}$ scores to the TE as measured by ribosome profiling of the natural leaderless mRNAs.
511 Since natural leaderless mRNAs encode many genes with diverse codon usages, a poorer correlation
512 was obtained with TE ($R^2=0.06$, slope=0.71 A.U.) than with the TIE reporters (Fig 7D). The correlation
513 of the $TIE_{\text{leaderless}}$ model at predicting ribosome profiling TE ($R=0.25$) is the same as observed for the
514 RBS calculator model of initiation and *E. coli* ribosome profiling data ($R=0.25$) (57). Since the TIE
515 reporters all code for YFP with near-identical codon usage, and the natural mRNAs have variable
516 codon usage frequencies, it is possible that translation elongation differences between natural ORFs
517 also impact translation efficiency. Indeed, translation elongation rates have been estimated to be rate
518 limiting *in vivo* in other bacteria (58,59). In addition, ribosome occupancy of stalled ribosomes can
519 complicate the analysis of ribosome profiling data, making the interpretation rather difficult. While it is
520 objectively harder to quantitatively predict translation levels, the $TIE_{\text{leaderless}}$ model performs rather well.

521 DISCUSSION

522 Here we provide the first systematic analysis of mRNA structure content, start codon identity,
523 and leader length on the initiation of leaderless mRNAs (Fig 7E). Importantly, this study was
524 performed using the bacterium *C. crescentus* which is adapted to efficient leaderless mRNA initiation
525 (27). As has been observed for leadered mRNAs (19,46), mRNA structure content at the leaderless
526 TIR hinders leaderless mRNA translation initiation, suggesting that ribosome accessibility is a key
527 feature for leaderless mRNAs. As previously observed in *E. coli*, changes in start codon identity from
528 the preferred “AUG” and presence of leader nucleotides leads to a significant reduction of TIE for *C.*
529 *crescentus* leaderless mRNAs. Using these quantitative data, we generated a combinatorial
530 $TIE_{\text{leaderless}}$ model that predicts the ability of an RNA to initiate as a leaderless mRNA from the
531 individual effects of these features which can be computed for any RNA in the transcriptome. This
532 $TIE_{\text{leaderless}}$ model both accurately and sensitively predicts the ability of all RNAs in the *C. crescentus*
533 transcriptome to initiate as leaderless mRNAs. While a 5' AUG is highly enriched in leaderless
534 mRNAs and only rarely observed in non-coding RNAs (Fig 6A), non-coding RNAs containing 5' AUGs
535 utilize a high ΔG_{unfold} to prevent aberrant translation initiation (Fig 6B). Additionally, very short leaders
536 which were found to inhibit leaderless mRNA initiation, are selected against in leaderless mRNAs and
537 are common in 5' regions of non-coding RNAs containing non-initiating AUGs. Finally, leaderless
538 mRNAs are much more selective for AUG start codons than are leadered mRNAs, suggesting that the
539 additional stabilization of the translation initiation complex provided by the SD-aSD base pairing helps
540 facilitate initiation on near-cognate start codons.

541 Leaderless mRNAs have been found to initiate translation in bacterial, archaeal, and both
542 cytoplasmic and mitochondrial eukaryotic ribosomes (17,28,60) suggesting that leaderless initiation is
543 an ancestral initiation mechanism. It is therefore possible that the TIE_{leaderless} model generated here in
544 *C. crescentus* may also perform well across organisms. Indeed, even a few nucleotides preceding the
545 AUG inhibit leaderless mRNA translation initiation in *C. crescentus*, *E. coli*, and mammalian
546 mitochondria (37,39). The strong inhibition of leaderless mRNA translation by TIR secondary structure
547 is likely why leaderless mRNAs in mitochondria have been found to lack 5' secondary structures (28).
548 *C. crescentus* shares a similar preference for 5' AUGs to *E. coli* for leaderless mRNA initiation (33).
549 Interestingly, in the *Mycobacteria*, GUG start codons are much more abundant in leaderless mRNAs
550 and tend to be initiated more similarly to AUG codons in this organism (16). *Mycobacterium* GUG
551 initiated leaderless mRNAs tend to code for short regulatory ORFs (16), as opposed to ORFs
552 encoding functional genes in *C. crescentus*. This suggests that there are likely to be some species-
553 specific differences in leaderless mRNA features arising from the differences in the translation
554 initiation machinery. Indeed, across prokaryotes, 79% of predicted leaderless genes contain AUG as
555 the start codon, whereas GUG, UUG and others are found with an average of 10%, 6% and 3%
556 respectively (13). Surprisingly, leaderless mRNAs across organisms appear to initiate with assembled
557 70S/80S ribosomes (31,61-63), further suggesting a conserved mechanism of initiation. Therefore, an
558 important goal moving forward will be to determine how broadly across organisms this TIE_{leaderless}
559 model might apply. Based upon the observations described here, it is likely that these features
560 (ΔG_{unfold} , start codon identity, and leader length) will combine similarly across species to define
561 leaderless mRNA TIRs. However, due to differences in the translation initiation machinery across
562 organisms, the specific hierarchy of mRNA features will need to be experimentally determined for a
563 given species in order to generate a TIE_{leaderless} model that accurately classifies leaderless mRNAs.

564 **AVAILABILITY**

565 Transcript architecture for the *Caulobacter crescentus* genome (Updated operon map (updated
566 4/24/2020)) was obtained from biochemicalphysics.com/resources.

567 **ACCESSION NUMBERS**

568 Not applicable.

569 **SUPPLEMENTARY DATA**

570 Supplementary Data are available at NAR online.

571 **ACKNOWLEDGEMENT**

572 We thank members of the Schrader lab for critical feedback.

573 **FUNDING**

574 This work was supported by the National Institutes of Health [R35GM124733 to J.M.S.]; and Start-up
575 funds from WSU to J.M.S. Funding for open access charge: National Institutes of Health.

576 CONFLICT OF INTEREST

577 No conflict of interests exist.

578 REFERENCES

- 579 1. Drummond, D.A. and Wilke, C.O. (2009) The evolutionary consequences of erroneous
580 protein synthesis. *Nat Rev Genet*, **10**, 715-724.
- 581 2. Kurland, C.G. and Ehrenberg, M. (1987) Growth-optimizing accuracy of gene expression.
582 *Annu Rev Biophys Biophys Chem*, **16**, 291-317.
- 583 3. Rodnina, M.V. and Wintermeyer, W. (2001) Fidelity of aminoacyl-tRNA selection on the
584 ribosome: kinetic and structural mechanisms. *Annu Rev Biochem*, **70**, 415-435.
- 585 4. Steitz, J.A. and Jakes, K. (1975) How Ribosomes Select Initiator Regions in Messenger-Rna -
586 Base Pair Formation between 3' Terminus of 16s Ribosomal-Rna and Messenger-Rna during
587 Initiation of Protein-Synthesis in Escherichia-Coli. *Proceedings of the National Academy of
588 Sciences of the United States of America*, **72**, 4734-4738.
- 589 5. Chen, H., Bjercknes, M., Kumar, R. and Jay, E. (1994) Determination of the optimal aligned
590 spacing between the Shine-Dalgarno sequence and the translation initiation codon of
591 Escherichia coli mRNAs. *Nucleic Acids Res*, **22**, 4953-4957.
- 592 6. Shine, J. and Dalgarno, L. (1974) The 3'-terminal sequence of Escherichia coli 16S ribosomal
593 RNA: complementarity to nonsense triplets and ribosome binding sites. *Proceedings of the
594 National Academy of Sciences of the United States of America*, **71**, 1342-1346.
- 595 7. Jacob, W.F., Santer, M. and Dahlberg, A.E. (1987) A single base change in the Shine-Dalgarno
596 region of 16S rRNA of Escherichia coli affects translation of many proteins. *Proceedings of
597 the National Academy of Sciences of the United States of America*, **84**, 4757-4761.
- 598 8. Hui, A. and de Boer, H.A. (1987) Specialized ribosome system: preferential translation of a
599 single mRNA species by a subpopulation of mutated ribosomes in Escherichia coli.
600 *Proceedings of the National Academy of Sciences of the United States of America*, **84**, 4762-
601 4766.
- 602 9. Calogero, R.A., Pon, C.L., Canonaco, M.A. and Gualerzi, C.O. (1988) Selection of the mRNA
603 translation initiation region by Escherichia coli ribosomes. *Proceedings of the National
604 Academy of Sciences of the United States of America*, **85**, 6427-6431.
- 605 10. Melancon, P., Leclerc, D., Destroismaisons, N. and Brakiergingras, L. (1990) The Anti-Shine-
606 Dalgarno Region in Escherichia-Coli 16s Ribosomal-Rna Is Not Essential for the Correct
607 Selection of Translational Starts. *Biochemistry*, **29**, 3402-3407.
- 608 11. Saito, K., Green, R. and Buskirk, A.R. (2020) Translational initiation in E. coli occurs at the
609 correct sites genome-wide in the absence of mRNA-rRNA base-pairing. *eLife*, **9**.
- 610 12. Chang, B., Halgamuge, S. and Tang, S.L. (2006) Analysis of SD sequences in completed
611 microbial genomes: non-SD-led genes are as common as SD-led genes. *Gene*, **373**, 90-99.
- 612 13. Srivastava, A., Gogoi, P., Deka, B., Goswami, S. and Kanaujia, S.P. (2016) In silico analysis of
613 5'-UTRs highlights the prevalence of Shine-Dalgarno and leaderless-dependent mechanisms
614 of translation initiation in bacteria and archaea, respectively. *J Theor Biol*, **402**, 54-61.
- 615 14. Beck, H.J. and Moll, I. (2018) Leaderless mRNAs in the Spotlight: Ancient but Not Outdated!
616 *Microbiol Spectr*, **6**.
- 617 15. Cortes, T., Schubert, O.T., Rose, G., Arnvig, K.B., Comas, I., Aebbersold, R. and Young, D.B.
618 (2013) Genome-wide mapping of transcriptional start sites defines an extensive leaderless
619 transcriptome in Mycobacterium tuberculosis. *Cell Rep*, **5**, 1121-1131.

- 620 16. Shell, S.S., Wang, J., Lapierre, P., Mir, M., Chase, M.R., Pyle, M.M., Gawande, R., Ahmad, R.,
621 Sarracino, D.A., Iroger, T.R. *et al.* (2015) Leaderless Transcripts and Small Proteins Are
622 Common Features of the Mycobacterial Translational Landscape. *PLoS Genet*, **11**, e1005641.
- 623 17. Montoya, J., Ojala, D. and Attardi, G. (1981) Distinctive features of the 5'-terminal sequences
624 of the human mitochondrial mRNAs. *Nature*, **290**, 465-470.
- 625 18. Nakamoto, T. (2006) A unified view of the initiation of protein synthesis. *Biochem Biophys*
626 *Res Commun*, **341**, 675-678.
- 627 19. de Smit, M.H. and van Duin, J. (1990) Secondary structure of the ribosome binding site
628 determines translational efficiency: a quantitative analysis. *Proceedings of the National*
629 *Academy of Sciences of the United States of America*, **87**, 7668-7672.
- 630 20. de Smit, M.H. and van Duin, J. (1994) Control of translation by mRNA secondary structure in
631 *Escherichia coli*. A quantitative analysis of literature data. *J Mol Biol*, **244**, 144-150.
- 632 21. Skripkin, E.A., Adhin, M.R., de Smit, M.H. and van Duin, J. (1990) Secondary structure of the
633 central region of bacteriophage MS2 RNA. Conservation and biological significance. *J Mol*
634 *Biol*, **211**, 447-463.
- 635 22. Gu, W., Zhou, T. and Wilke, C.O. (2010) A universal trend of reduced mRNA stability near the
636 translation-initiation site in prokaryotes and eukaryotes. *PLoS Comput Biol*, **6**, e1000664.
- 637 23. Salis, H.M. (2011) The ribosome binding site calculator. *Methods Enzymol*, **498**, 19-42.
- 638 24. Romero, D.A., Hasan, A.H., Lin, Y.F., Kime, L., Ruiz-Larrabeiti, O., Urem, M., Bucca, G.,
639 Mamanova, L., Laing, E.E., van Wezel, G.P. *et al.* (2014) A comparison of key aspects of gene
640 regulation in *Streptomyces coelicolor* and *Escherichia coli* using nucleotide-resolution
641 transcription maps produced in parallel by global and differential RNA sequencing. *Mol*
642 *Microbiol*.
- 643 25. Babski, J., Haas, K.A., Nather-Schindler, D., Pfeiffer, F., Forstner, K.U., Hammelmann, M.,
644 Hilker, R., Becker, A., Sharma, C.M., Marchfelder, A. *et al.* (2016) Genome-wide identification
645 of transcriptional start sites in the haloarchaeon *Haloferax volcanii* based on differential
646 RNA-Seq (dRNA-Seq). *BMC Genomics*, **17**, 629.
- 647 26. Pfeifer-Sancar, K., Mentz, A., Ruckert, C. and Kalinowski, J. (2013) Comprehensive analysis of
648 the *Corynebacterium glutamicum* transcriptome using an improved RNAseq technique. *BMC*
649 *Genomics*, **14**, 888.
- 650 27. Schrader, J.M., Zhou, B., Li, G.W., Lasker, K., Childers, W.S., Williams, B., Long, T., Crosson, S.,
651 McAdams, H.H., Weissman, J.S. *et al.* (2014) The coding and noncoding architecture of the
652 *Caulobacter crescentus* genome. *PLoS Genet*, **10**, e1004463.
- 653 28. Jones, C.N., Wilkinson, K.A., Hung, K.T., Weeks, K.M. and Spremulli, L.L. (2008) Lack of
654 secondary structure characterizes the 5' ends of mammalian mitochondrial mRNAs. *RNA*, **14**,
655 862-871.
- 656 29. Korman, M., Schlüssel, S., Vishkautzan, M. and Gur, E. (2019) Multiple layers of regulation
657 determine the cellular levels of the Pup ligase PafA in *Mycobacterium smegmatis*. *Mol*
658 *Microbiol*, **112**, 620-631.
- 659 30. Tedin, K., Moll, I., Grill, S., Resch, A., Graschopf, A., Gualerzi, C.O. and Blasi, U. (1999)
660 Translation initiation factor 3 antagonizes authentic start codon selection on leaderless
661 mRNAs. *Mol Microbiol*, **31**, 67-77.
- 662 31. O'Donnell, S.A. and Janssen, G.R. (2002) Leaderless mRNAs bind 70S ribosomes more
663 strongly than 30S ribosomal subunits in *Escherichia coli*. *Journal of Bacteriology*, **184**, 6730-
664 6733.
- 665 32. Van Etten, W.J. and Janssen, G.R. (1998) An AUG initiation codon, not codon-anticodon
666 complementarity, is required for the translation of unleadered mRNA in *Escherichia coli*.
667 *Molecular Microbiology*, **27**, 987-1001.

- 668 33. O'Donnell, S.M. and Janssen, G.R. (2001) The initiation codon affects ribosome binding and
669 translational efficiency in *Escherichia coli* of *cl* mRNA with or without the 5' untranslated
670 leader. *J Bacteriol*, **183**, 1277-1283.
- 671 34. Hering, O., Brenneis, M., Beer, J., Suess, B. and Soppa, J. (2009) A novel mechanism for
672 translation initiation operates in haloarchaea. *Mol Microbiol*, **71**, 1451-1463.
- 673 35. Chen, W.C., Yang, G.P., He, Y., Zhang, S.M., Chen, H.Y., Shen, P., Chen, X.D. and Huang, Y.P.
674 (2015) Nucleotides Flanking the Start Codon in *hsp70* mRNAs with Very Short 5'-UTRs
675 Greatly Affect Gene Expression in Haloarchaea. *Plos One*, **10**.
- 676 36. Brock, J.E., Pourshahian, S., Giliberti, J., Limbach, P.A. and Janssen, G.R. (2008) Ribosomes
677 bind leaderless mRNA in *Escherichia coli* through recognition of their 5'-terminal AUG. *Rna-*
678 *a Publication of the Rna Society*, **14**, 2159-2169.
- 679 37. Krishnan, K.M., Van Etten, W.J., 3rd and Janssen, G.R. (2010) Proximity of the start codon to
680 a leaderless mRNA's 5' terminus is a strong positive determinant of ribosome binding and
681 expression in *Escherichia coli*. *J Bacteriol*, **192**, 6482-6485.
- 682 38. Jones, R.L., 3rd, Jaskula, J.C. and Janssen, G.R. (1992) In vivo translational start site selection
683 on leaderless mRNA transcribed from the *Streptomyces fradiae* *aph* gene. *J Bacteriol*, **174**,
684 4753-4760.
- 685 39. Christian, B.E. and Spemulli, L.L. (2010) Preferential Selection of the 5'-Terminal Start
686 Codon on Leaderless mRNAs by Mammalian Mitochondrial Ribosomes. *Journal of Biological*
687 *Chemistry*, **285**, 28379-28386.
- 688 40. Zhou, B., Schrader, J.M., Kalogeraki, V.S., Abeliuk, E., Dinh, C.B., Pham, J.Q., Cui, Z.Z., Dill, D.L.,
689 McAdams, H.H. and Shapiro, L. (2015) The global regulatory architecture of transcription
690 during the *Caulobacter* cell cycle. *PLoS Genet*, **11**, e1004831.
- 691 41. Marks, M.E., Castro-Rojas, C.M., Teiling, C., Du, L., Kapatral, V., Walunas, T.L. and Crosson, S.
692 (2010) The genetic basis of laboratory adaptation in *Caulobacter crescentus*. *J Bacteriol*, **192**,
693 3678-3688.
- 694 42. Gelsinger, D.R. and DiRuggiero, J. (2018) Transcriptional Landscape and Regulatory Roles of
695 Small Noncoding RNAs in the Oxidative Stress Response of the Haloarchaeon *Haloferax*
696 *volcanii*. *J Bacteriol*, **200**.
- 697 43. Rudler, D.L., Hughes, L.A., Perks, K.L., Richman, T.R., Kuznetsova, I., Ermer, J.A., Abudulai,
698 L.N., Shearwood, A.M.J., Viola, H.M., Hool, L.C. *et al.* (2019) Fidelity of translation initiation is
699 required for coordinated respiratory complex assembly. *Sci Adv*, **5**.
- 700 44. Gelsinger, D.R., Dallon, E., Reddy, R., Mohammad, F., Buskirk, A.R. and DiRuggiero, J. (2020)
701 Ribosome profiling in archaea reveals leaderless translation, novel translational initiation
702 sites, and ribosome pausing at single codon resolution. *Nucleic Acids Res*, **48**, 5201-5216.
- 703 45. Langmead, B., Trapnell, C., Pop, M. and Salzberg, S.L. (2009) Ultrafast and memory-efficient
704 alignment of short DNA sequences to the human genome. *Genome Biol*, **10**, R25.
- 705 46. Mustoe, A.M., Corley, M., Laederach, A. and Weeks, K.M. (2018) Messenger RNA Structure
706 Regulates Translation Initiation: A Mechanism Exploited from Bacteria to Humans.
707 *Biochemistry*, **57**, 3537-3539.
- 708 47. Gruber, A.R., Lorenz, R., Bernhart, S.H., Neubock, R. and Hofacker, I.L. (2008) The Vienna
709 RNA websuite. *Nucleic Acids Res*, **36**, W70-74.
- 710 48. Reuter, J.S. and Mathews, D.H. (2010) RNAstructure: software for RNA secondary structure
711 prediction and analysis. *BMC Bioinformatics*, **11**, 129.
- 712 49. Thanbichler, M., Iniesta, A.A. and Shapiro, L. (2007) A comprehensive set of plasmids for
713 vanillate- and xylose-inducible gene expression in *Caulobacter crescentus*. *Nucleic Acids Res*,
714 **35**, e137.
- 715 50. Schindelin, J., Arganda-Carreras, I., Frise, E., Kaynig, V., Longair, M., Pietzsch, T., Preibisch, S.,
716 Rueden, C., Saalfeld, S., Schmid, B. *et al.* (2012) Fiji: an open-source platform for biological-
717 image analysis. *Nat Methods*, **9**, 676-682.

- 718 51. Ducret, A., Quardokus, E.M. and Brun, Y.V. (2016) MicrobeJ, a tool for high throughput
719 bacterial cell detection and quantitative analysis. *Nat Microbiol*, **1**, 16077.
- 720 52. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
721 Prettenhofer, P., Weiss, R., Dubourg, V. *et al.* (2011) Scikit-learn: Machine Learning in Python.
722 *Journal of Machine Learning Research*, **12**, 2825-2830.
- 723 53. Tobias, J.W., Shrader, T.E., Rocap, G. and Varshavsky, A. (1991) The N-end rule in bacteria.
724 *Science*, **254**, 1374-1377.
- 725 54. Ingolia, N.T., Ghaemmaghami, S., Newman, J.R.S. and Weissman, J.S. (2009) Genome-Wide
726 Analysis in Vivo of Translation with Nucleotide Resolution Using Ribosome Profiling. *Science*,
727 **324**, 218-223.
- 728 55. Amunts, A., Brown, A., Toots, J., Scheres, S.H.W. and Ramakrishnan, V. (2015) Ribosome. The
729 structure of the human mitochondrial ribosome. *Science*, **348**, 95-98.
- 730 56. Bird, J.G., Basu, U., Kuster, D., Ramachandran, A., Grudzien-Nogalska, E., Towheed, A.,
731 Wallace, D.C., Kiledjian, M., Temiakov, D., Patel, S.S. *et al.* (2018) Highly efficient 5' capping
732 of mitochondrial RNA with NAD(+) and NADH by yeast and human mitochondrial RNA
733 polymerase. *eLife*, **7**.
- 734 57. Li, G.W., Burkhardt, D., Gross, C. and Weissman, J.S. (2014) Quantifying absolute protein
735 synthesis rates reveals principles underlying allocation of cellular resources. *Cell*, **157**, 624-
736 635.
- 737 58. Racle, J., Picard, F., Girbal, L., Coccagn-Bousquet, M. and Hatzimanikatis, V. (2013) A genome-
738 scale integration and analysis of Lactococcus lactis translation data. *PLoS Comput Biol*, **9**,
739 e1003240.
- 740 59. Vieira, J.P., Racle, J. and Hatzimanikatis, V. (2016) Analysis of Translation Elongation
741 Dynamics in the Context of an Escherichia coli Cell. *Biophys J*, **110**, 2120-2131.
- 742 60. Baltz, R.H., Hegeman, G. and Skatrud, P.L. (1993) *Industrial microorganisms: basic and*
743 *applied molecular genetics*. American Society for Microbiology.
- 744 61. Moll, I., Hirokawa, G., Kiel, M.C., Kaji, A. and Blasi, U. (2004) Translation initiation with 70S
745 ribosomes: an alternative pathway for leaderless mRNAs. *Nucleic Acids Res*, **32**, 3354-3363.
- 746 62. Udagawa, T., Shimizu, Y. and Ueda, T. (2004) Evidence for the translation initiation of
747 leaderless mRNAs by the intact 70 S ribosome without its dissociation into subunits in
748 eubacteria. *J Biol Chem*, **279**, 8539-8546.
- 749 63. Andreev, D.E., Terenin, I.M., Dunaevsky, Y.E., Dmitriev, S.E. and Shatsky, I.N. (2006) A
750 leaderless mRNA can bind to mammalian 80S ribosomes and direct polypeptide synthesis in
751 the absence of translation initiation factors. *Mol Cell Biol*, **26**, 3164-3169.

752 TABLE AND FIGURES LEGENDS

753 **Figure 1. Leaderless mRNA translation initiation regions are more accessible than leadered**
754 **mRNAs.** A.) Predicted unfolding energy of mRNAs. The predicted mRNA minimum free energy
755 (ΔG_{mRNA}) is represented on the left. The orange translation initiation region indicates a ribosome
756 footprint surrounding the start codon (pink). The image on the right represents the mRNA upon
757 initiation (ΔG_{init}) where the orange initiation region is unfolded. The ΔG_{unfold} represents the amount of
758 energy required by the ribosome to unfold the translation initiation region of the mRNA. B.) Violin plots
759 of ΔG_{unfold} (right) calculated for all the mRNAs of each class (left) in the *Caulobacter crescentus*
760 genome based on the transcript architecture(27,40). P-values were calculated based on t-test (two
761 tailed, unequal variance),,

762 **Figure 2. ΔG_{unfold} strongly influences leaderless mRNA translation.** A.) A representative TIR
763 synthetic stem loop synonymous mutation set with varying ΔG_{unfold} values. The bases in the start
764 codon are colored pink, red bases highlight where mutations were introduced to disrupt base pairing.
765 B.) *In vivo* translation reporter levels the various leaderless RNA mutants. Each hairpin and its
766 synonymous codon mutant set are shown with the same color (Raw data can be found in Table S1).
767 Black points = leaderless set 1, grey points = leaderless set 2, dark blue points = leaderless set 3,
768 purple points = leaderless set 4, light blue points = leaderless set 5, red points = leaderless set 6, and
769 teal points = leaderless set 7. The natural log of the average YFP intensity per cell is shown and error
770 bars represent the standard deviation of three biological replicates. The dotted blue line represents a
771 linear curve fit with an R^2 value of 0.84 and a slope of -0.3.

772 **Figure 3. Leaderless mRNAs have a strong preference for AUG start codons.** Leaderless mRNA
773 *in vivo* translation reporters were generated with the start codons listed on the X-axis and their
774 average YFP intensity per cell were measured. On the right, is a zoomed in view of all non-AUG
775 codons tested. Error bars represent the standard deviation from three biological replicates.

776 **Figure 4. Leaderless mRNAs are inhibited by additional upstream nucleotides.** Leaderless
777 mRNA *in vivo* translation reporters were generated with variable number of leading nucleotides on the
778 X-axis and their average YFP intensity per cell were measured (Raw data can be found in Table S1).
779 Error bars represent three biological replicates.

780 **Figure 5. ΔG_{unfold} , start codon identity, and leader length correlate with translation efficiency
781 (TE) across native leaderless mRNAs.** A.) Bar graph showing the fraction of leaderless mRNAs
782 starting with AUG, GUG, UUG and CUG start codons. Also shown are the random chances of
783 trinucleotides being AUG, GUG, UUG and CUG calculated based on GC content (67%) of *C.*
784 *crenscentus* genome. P-values were calculated based on a two-tailed Z-test. B.) Bar graph showing
785 the fraction of leaderless mRNAs and mRNAs with 5' untranslated region (UTR) of length 1 to 10 (as
786 determined in (40)). mRNAs containing Shine-Dalagarno sites were excluded from this analysis. P-
787 values were calculated based on a two-tailed Z-test of each leader length compared to leader length 0.
788 C.) Violin plot of translation efficiency (TE) as measured by ribosome profiling(54) of natural
789 leaderless mRNAs binned in three groups depending on ΔG_{unfold} values (0-5, 5-10, and >10 kcal/mol).
790 P-values based on t-test (two tailed, unequal variances). D.) Violin plot of TE as measured by
791 ribosome profiling(54) of natural leaderless mRNAs starting with AUG and GUG. P-values were
792 calculated based on a t-test (2-tailed, unequal variance). E.) Violin plot showing the TE as measured
793 by ribosome profiling(54) on the Y-axis of leaderless mRNAs (green) and with leaders of varying
794 length (1-10) shown in grey. P-values were calculated based on t-test (2-tailed, unequal variance).

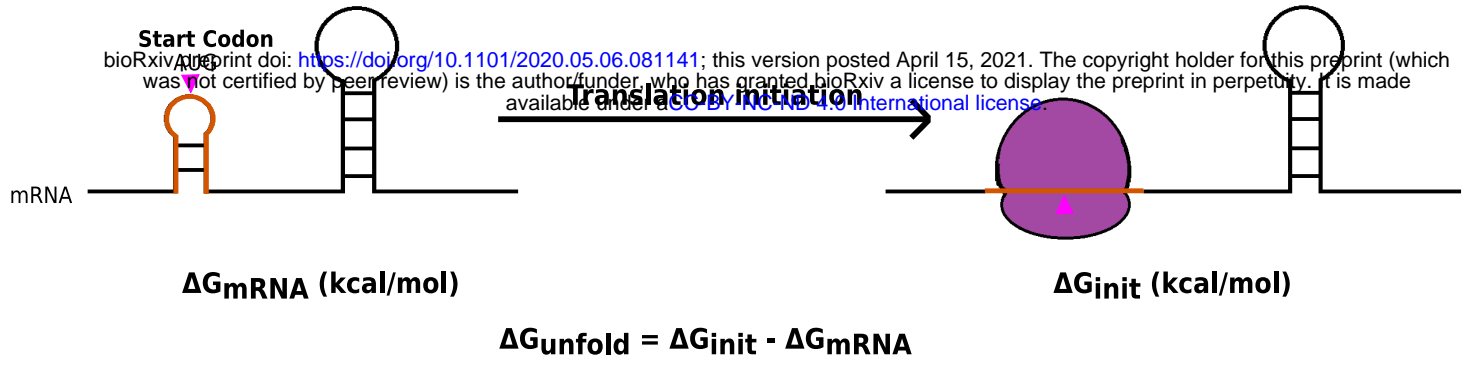
795 **Figure 6. Non-coding RNAs with 5' AUGs are rare and have higher ΔG_{unfold} .** A.) Bar graph
796 showing the fraction of natural leaderless mRNAs starting with trinucleotide AUG and other types of
797 RNAs starting with trinucleotide AUG, but not initiated at that AUG (leadered mRNAs, sRNAs, rRNAs,
798 tRNAs and asRNAs). Also shown is the random chance of trinucleotide being AUG out of 10000

799 nucleotides; calculated based on GC content of *C. crescentus* genome. P-values were calculated
800 using a two-tailed Z-test with each RNA class compared to the random probability of 5' AUG. B.)
801 Violin plot showing ΔG_{unfold} of natural leaderless mRNAs starting with AUG (green) and other types of
802 RNAs starting with AUG, but not initiated at that AUG (leadered mRNAs, RNAs and asRNAs) (shown
803 in grey). P-values were calculated based on a T-test (2-tailed, unequal variance).

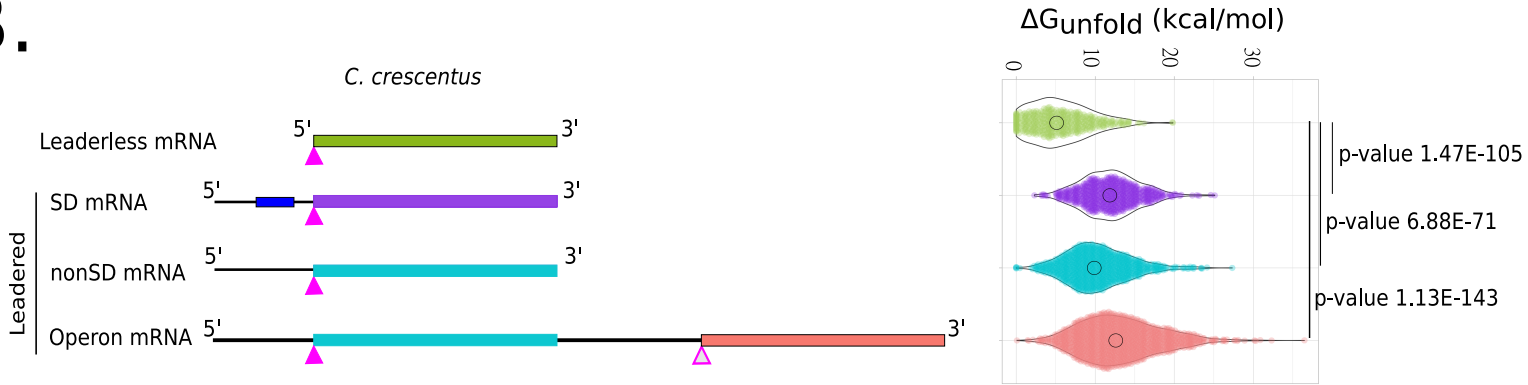
804 **Figure 7. A combinatorial model accurately predicts translation of leaderless mRNAs.** A.) Line
805 graph showing the predicted $\text{TIE}_{\text{leaderless}}$ scores on the X-axis and the number of RNAs on the Y axis.
806 The solid blue line represents natural leaderless mRNAs. The orange line represents the RNAs that
807 are not leaderless RNAs. The black dotted line represents all RNAs. RNAs with short leaders are
808 shown in Fig S2. B.) ROC curve (shown in solid blue, with “random” shown as a dotted line) with true
809 positive rate on Y-axis and false positive rate on X-axis. The area under curve (A.U.C.) was
810 calculated to be 0.99 for classification based on the $\text{TIE}_{\text{leaderless}}$ score and 0.68 for classification based
811 solely on presence of a 5' AUG (Fig S3). C.) TIE reporter levels compared to $\text{TIE}_{\text{leaderless}}$ scores. For
812 the leaderless TIE reporters tested (Table S1) the YFP reporter level (Y-axis) is plotted compared to
813 the $\text{TIE}_{\text{leaderless}}$ (X-axis). The trendline is the result of a least-squares fit yielding a slope of 2050 A.U.
814 with $R^2=0.87$. Error bars represent the standard deviation of at least three biological replicates. D.)
815 Translation efficiency (TE) of leaderless mRNAs (Y-axis) is plotted compared to $\text{TIE}_{\text{leaderless}}$ (X-axis).
816 The trendline is the result of a least-squares fit yielding a slope of 0.71 and $R^2=0.06$. E.) Model design
817 showing ribosome binding to the AUG trinucleotide (pink triangle) at the 5' end when it is highly
818 accessible as shown in the left. The ribosome binding is prevented when the region becomes more
819 structured and the accessibility decreases.

820

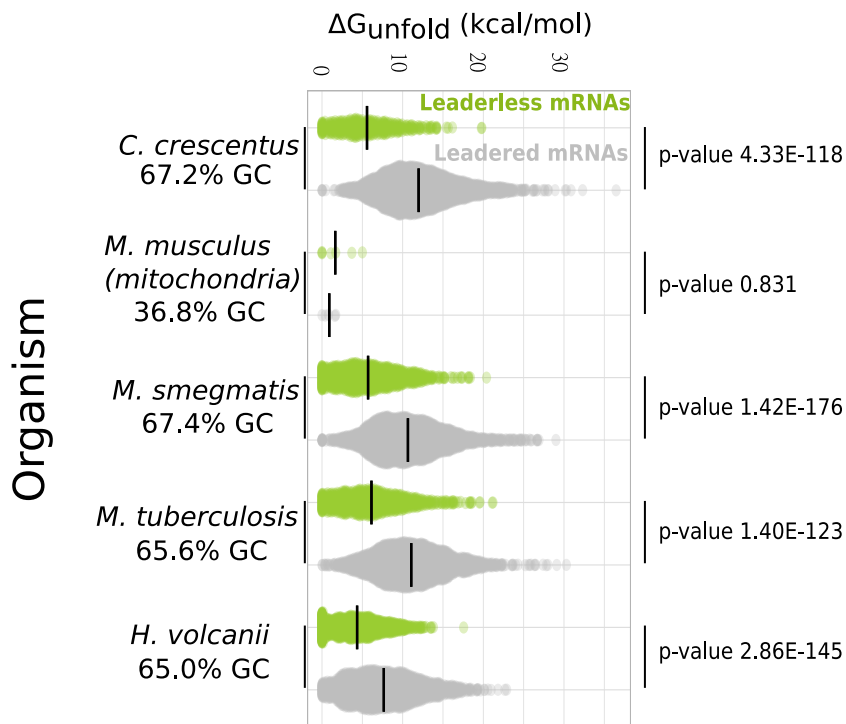
A.



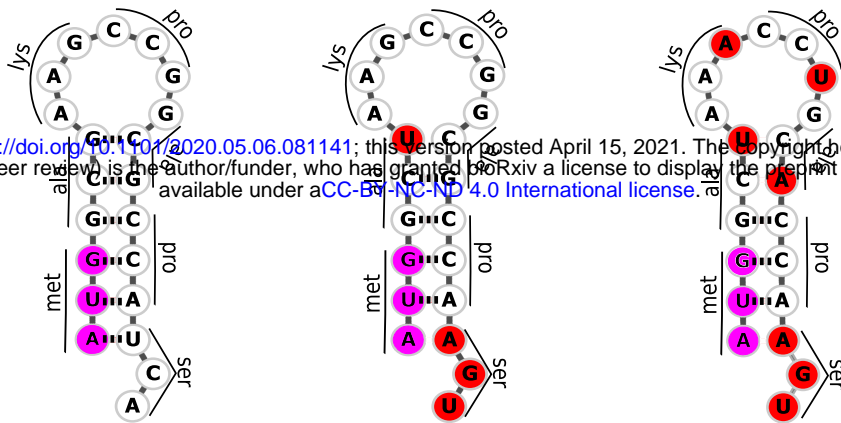
B.



C.



bioRxiv preprint doi: <https://doi.org/10.1101/2020.05.06.081141>; this version posted April 15, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



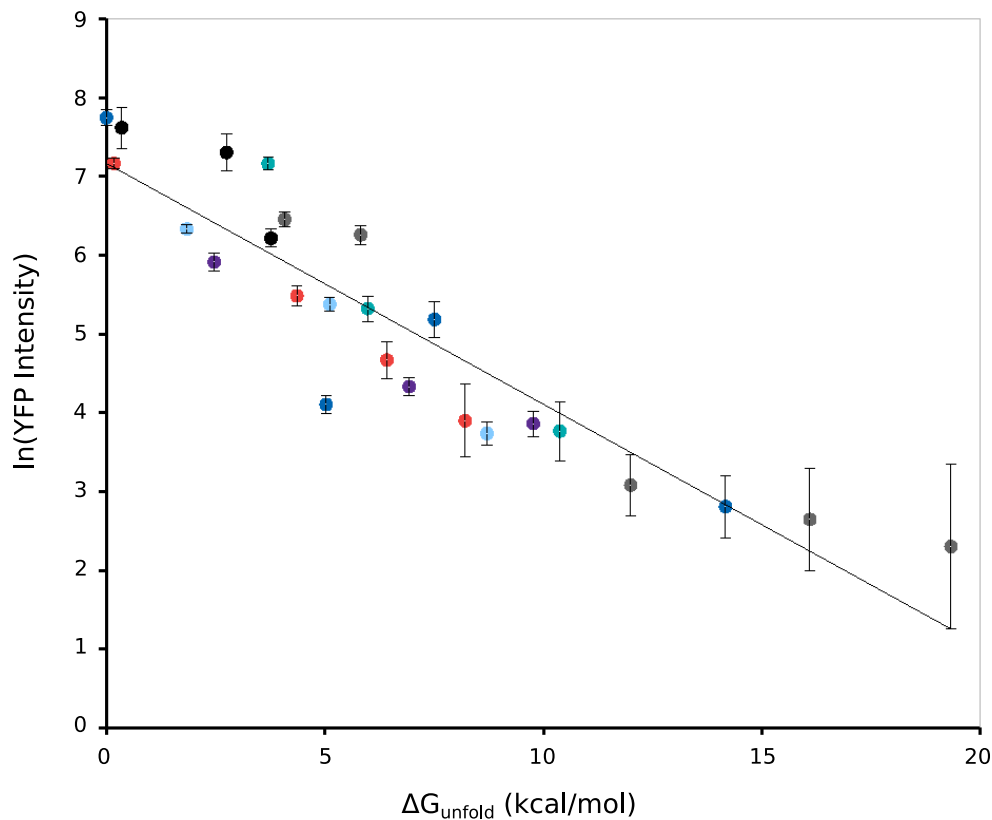
ΔG_{unfold} (kcal/mol) =

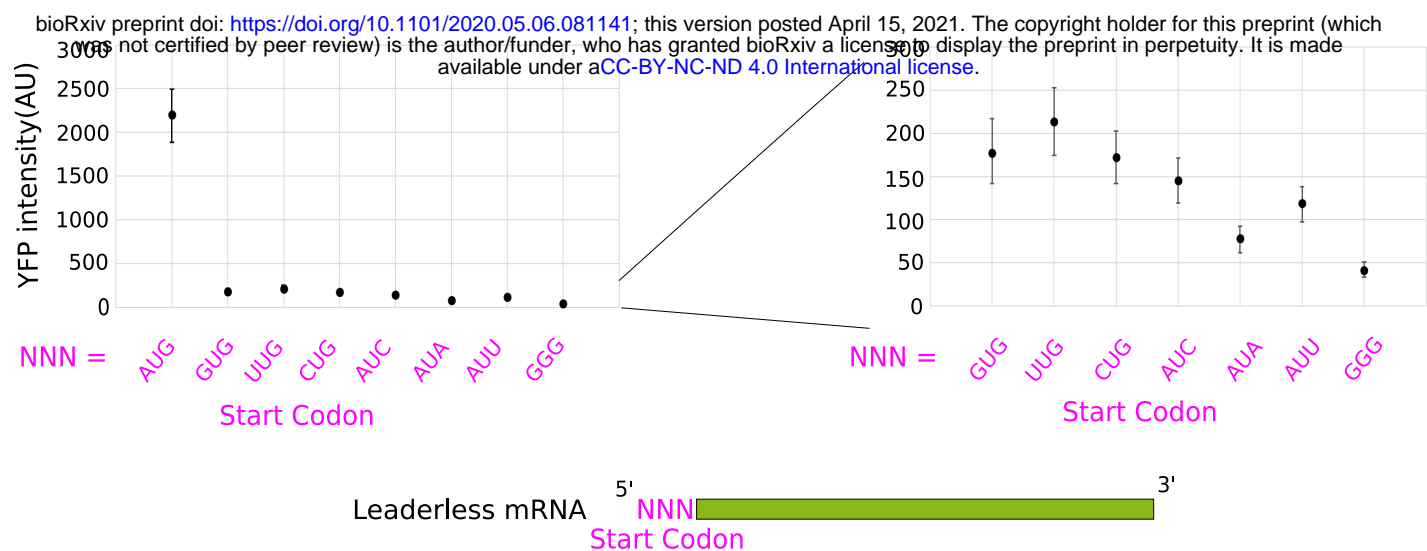
8.71

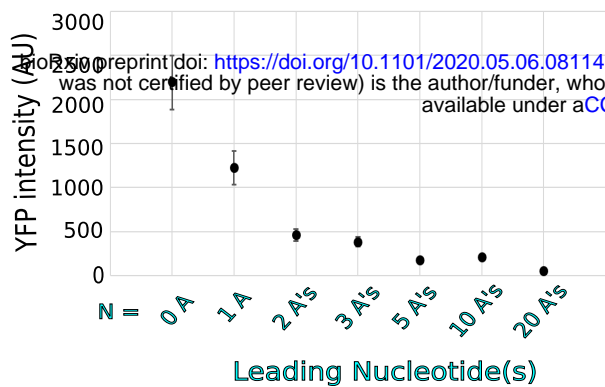
7.93

4.94

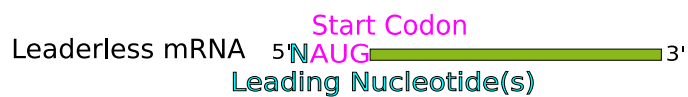
B.

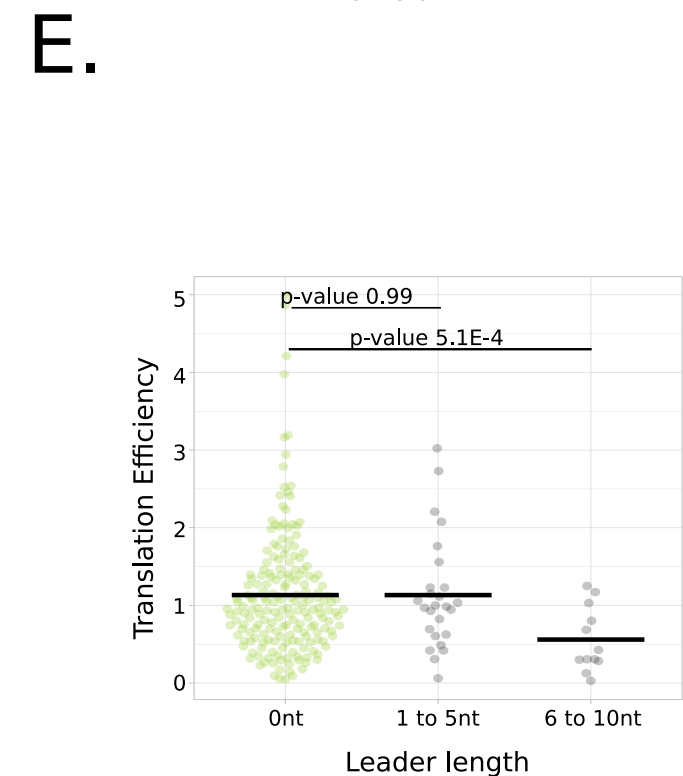
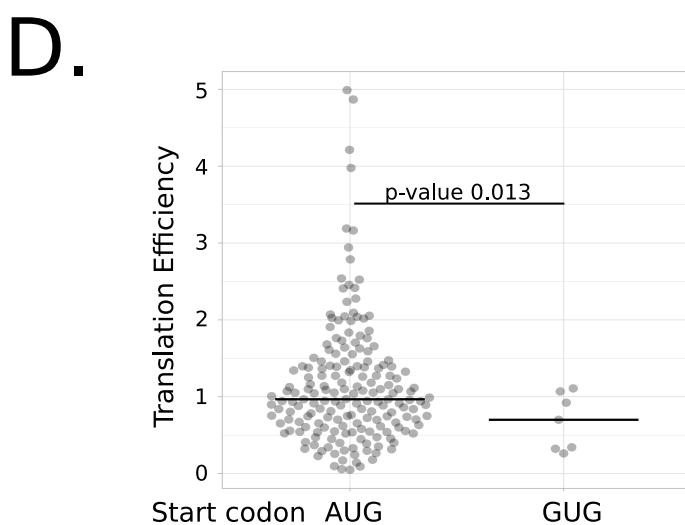
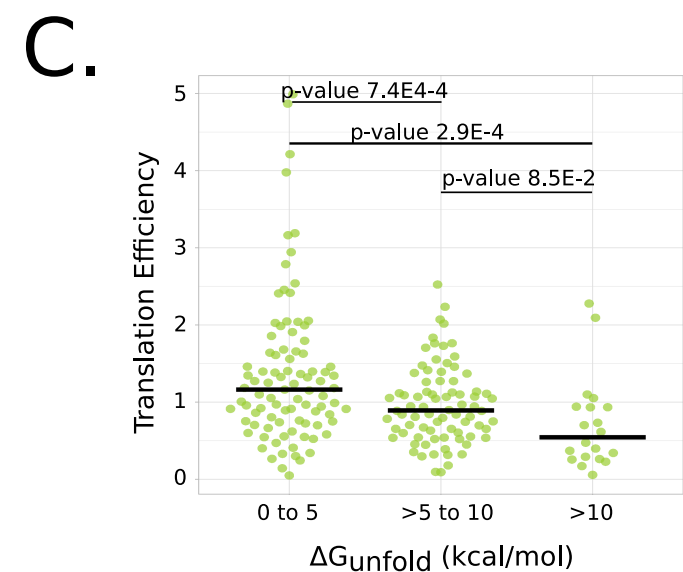
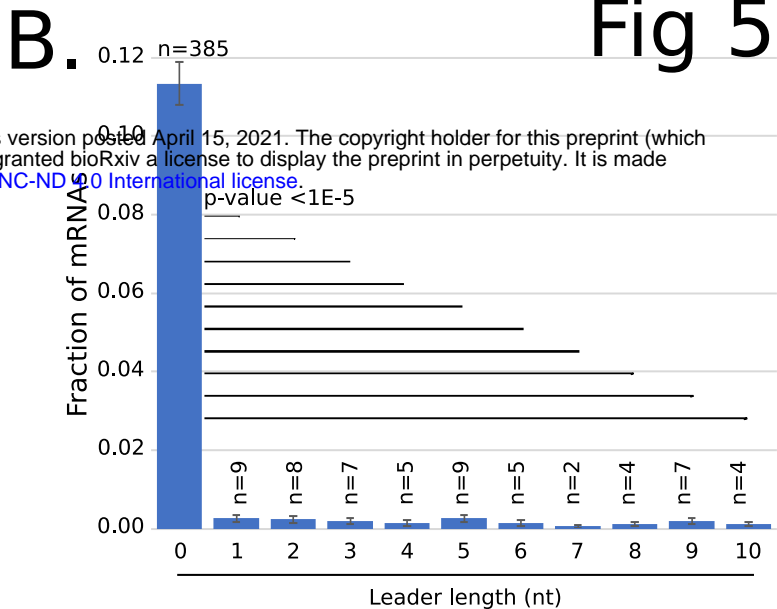
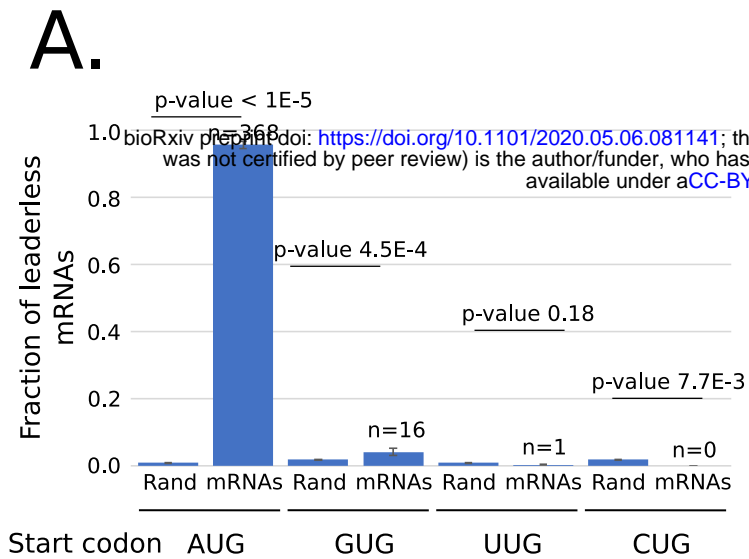






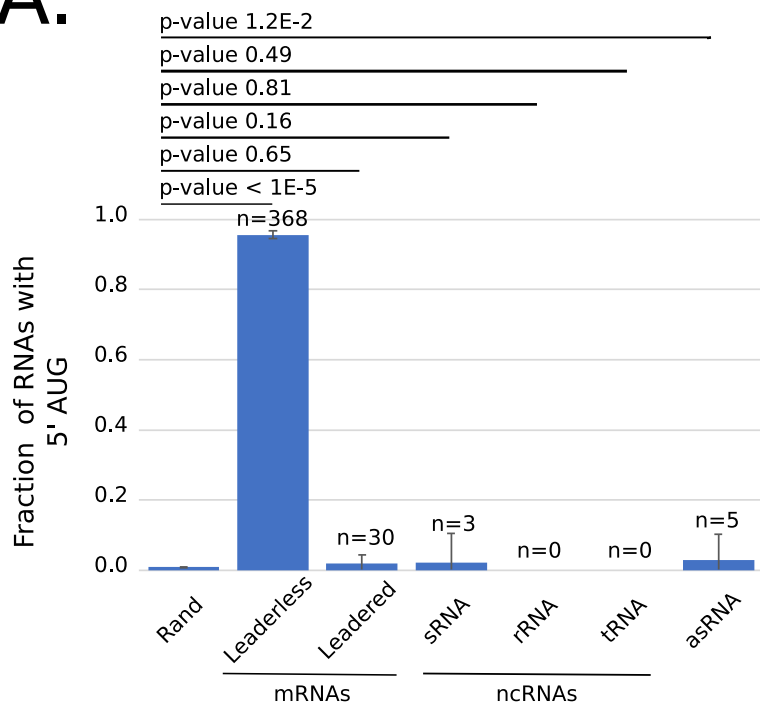
bioRxiv preprint doi: <https://doi.org/10.1101/2020.05.06.081141>; this version posted April 15, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).



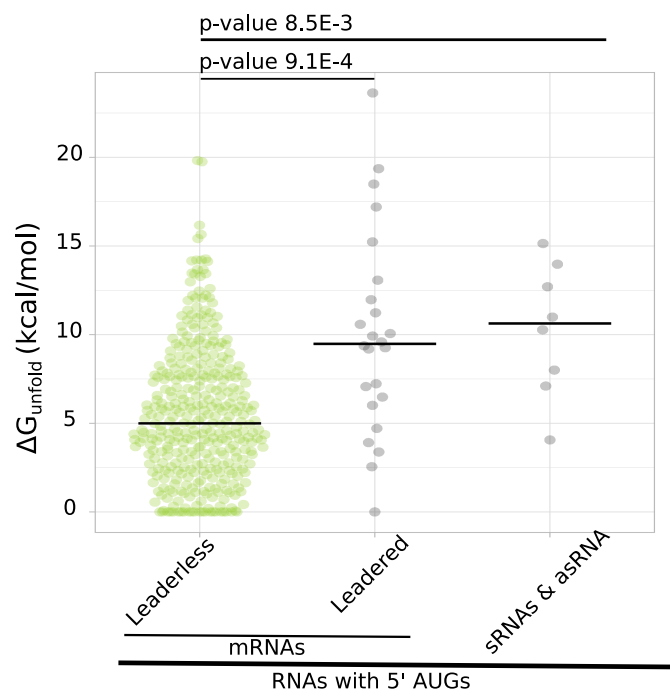


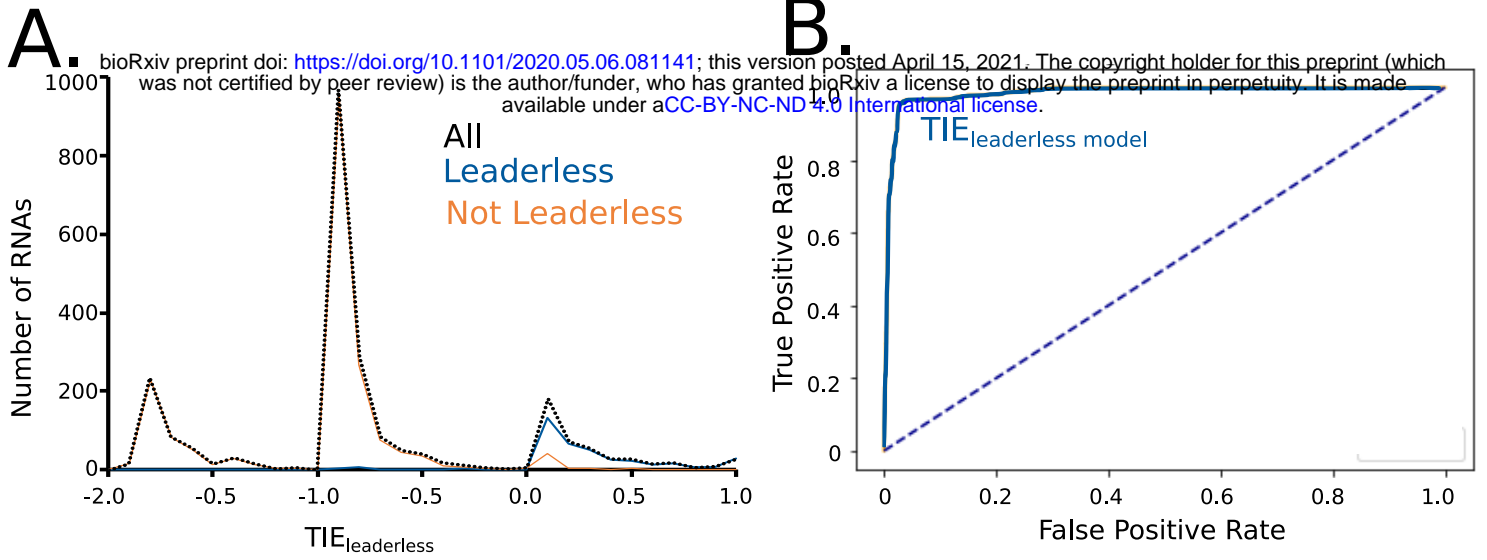
bioRxiv preprint doi: <https://doi.org/10.1101/2020.05.06.081141>; this version posted April 15, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

A.



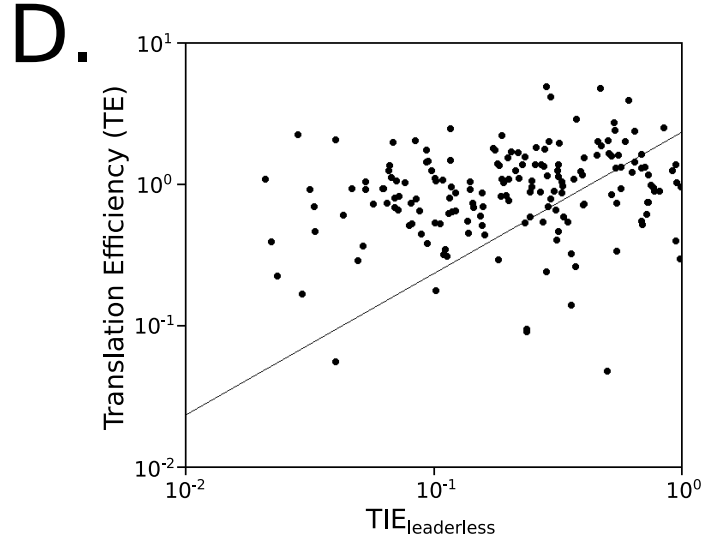
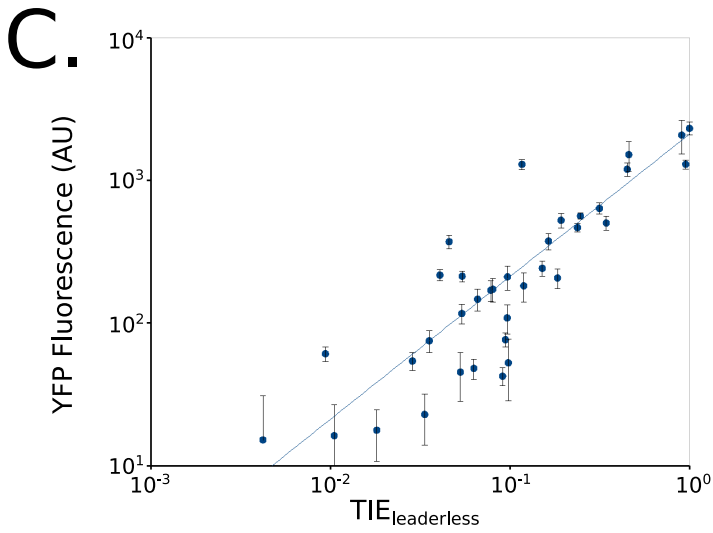
B.





Leaderless translation initiation reporter level

Translation level by ribosome profiling



E. Leaderless mRNAs

Other RNAs

