

# 1 **Ultradeep characterisation of translational sequence determinants** 2 **refutes rare-codon hypothesis and unveils quadruplet base pairing** 3 **of initiator tRNA and transcript**

4 Simon Höllerer<sup>1</sup> and Markus Jeschek<sup>1,2\*</sup>

5

6 <sup>1</sup>Department of Biosystems Science and Engineering, Swiss Federal Institute of Technology - ETH  
7 Zurich, Basel, CH-4058, Switzerland

8 <sup>2</sup>Institute of Microbiology, Synthetic Microbiology Lab, University of Regensburg, D-93053  
9 Regensburg, Germany

10

11 \*To whom correspondence should be addressed: E-mail: [markus.jeschek@bsse.ethz.ch](mailto:markus.jeschek@bsse.ethz.ch)

12

## 13 **ABSTRACT**

14 Translation is a key determinant of gene expression and an important biotechnological engineering  
15 target. In bacteria, 5'-untranslated region (5'-UTR) and coding sequence (CDS) are well-known mRNA  
16 parts controlling translation and thus cellular protein levels. However, the complex interaction of 5'-UTR  
17 and CDS has so far only been studied for few sequences leading to non-generalisable and partly  
18 contradictory conclusions. Herein, we systematically assess the dynamic translation from over 1.2  
19 million 5'-UTR-CDS pairs in *Escherichia coli* to investigate their collective effect using a new method for  
20 ultradeep sequence-function mapping. This allows us to disentangle and precisely quantify effects of  
21 known and hypothetical sequence determinants of translation. We find that 5'-UTR and CDS individually  
22 account for 53% and 20% of variance in translation, respectively, and show conclusively that, contrary  
23 to a common hypothesis, tRNA abundance does not explain expression changes between CDSs with  
24 different synonymous codons. Moreover, the obtained large-scale data clearly point to a base-pairing  
25 interaction between initiator tRNA and mRNA beyond the anticodon-codon interaction, an effect that is  
26 often masked for individual sequences and therefore inaccessible to low-throughput approaches. Our  
27 study highlights the indispensability of ultradeep sequence-function mapping to accurately determine  
28 the contribution of parts and phenomena involved in gene regulation.

## 29 INTRODUCTION

30 Translation is a key step of gene expression and an important engineering target in synthetic biology.  
31 To this end, genetic parts that influence translation are modified to alter absolute and relative expression  
32 levels to engineer biosystems through control of individual genes, pathways and even entire metabolic  
33 networks (1-3). In prokaryotes, initiation of translation is the rate-limiting step in the translational process,  
34 during which ribosomes assemble on the mRNA to start the templated elongation of the nascent  
35 polypeptide (4-7). At the onset of this step, the 30S ribosomal subunit attaches to the ribosome binding  
36 site (RBS) in the 5'-untranslated region (5'-UTR) upstream of the coding sequence (CDS). The 3'-end  
37 of the 16S rRNA hybridises with the Shine-Dalgarno (SD) motif, a conserved five to eight nucleotide (nt)  
38 sequence located upstream of the start codon, which facilitates translation (8-12). However, since Shine  
39 and Dalgarno's discovery in 1973 (10), various additional influencing factors and sequence  
40 determinants affecting translation initiation were identified. For example, the distance between SD motif  
41 and start codon, the type of start codon, and interactions between distant 5'-UTR parts and the ribosome  
42 play important roles (13-21). Remarkably, in some cases SD-like motifs are not required for translation,  
43 an observation hinting at the existence of other mechanisms besides "canonical" translation initiation  
44 (22-26). Further, the influence of mRNA secondary structures was studied under the hypothesis that  
45 the required unfolding of such structures during translation initiation might decrease expression  
46 (14,20,27-40). For example, stable secondary structures around the start codon were found to hinder  
47 translation, while structures further up- or downstream had less pronounced effects (35).  
48 Moreover, codon usage was found to influence translation. Genome-wide analyses of *E. coli* and other  
49 organisms revealed an overrepresentation of rare codons in the first five to ten triplets of the CDS in  
50 native genes, and their occurrence in this region was found to coincide with high expression (29,41-44).  
51 These observations led to two different hypotheses that differ fundamentally in terms of the underlying  
52 causality. The first hypothesis is related to the fact that cellular tRNA concentrations correlate with the  
53 occurrence frequency of their cognate codons (45-47). It was postulated that rare codons (with low-  
54 abundant cognate tRNAs) may have been evolutionary selected for within the N-terminal CDS to slow  
55 down early translation elongation and reduce premature termination due to clashing ribosomes (38,48-  
56 52). These "translational ramps" were postulated to be causally responsible for elevated expression of  
57 genes rich in rare codons at the CDS's 5'-end. As an alternative explanation independent of tRNA  
58 abundance, a second hypothesis has been proposed based on the fact that many rare codons are (or  
59 happen to be) AT-rich (29,52). Their occurrence is therefore associated with a lower tendency to form  
60 stable mRNA secondary structures (29,33,43), which are known to hinder translation initiation.  
61 In the context of these two hypotheses, several studies have been conducted to investigate the impact  
62 of codon usage on expression focussing either on the N-terminal codons alone (29,33,37,43,44) or the  
63 entire CDS (29) while applying different metrics of codon usage such as the codon adaptation index  
64 (CAI) (41), the frequency of "optimal" codons pairing with the most abundant tRNAs (45), and the tRNA  
65 adaptation index (tAI) (53), as discussed in detail elsewhere (44,54-56). Remarkably, while there is  
66 clear evidence for a high degree of interactivity between 5'-UTR and CDS, these two mRNA parts were  
67 handled separately in these studies: commonly only one of the two parts (either 5'-UTR or CDS) was  
68 diversified at a time, and systematic testing of larger numbers of 5'-UTR-CDS combinations to assess

69 their interaction was not performed (15,20,33). Thus, due to the strong interdependence the measured  
70 effects could not be clearly assigned to individual sequence parameters, and their contribution to overall  
71 expression could not be accurately quantified. Moreover, many early studies relied on experimental  
72 testing of only a few “hand-picked” sequences (usually less than 100 variants) due to limitations in  
73 experimental throughput or library generation (note that the CDS cannot be freely mutated, since of  
74 amino acid substitution may result in change or loss of reporter protein activity). Although valuable  
75 contributions, such empiric efforts have proven insufficient to establish generalisable rules and  
76 quantitative measurements for the potential effects of sequence parameters, which in some cases even  
77 led to contradictory conclusions. For example, the question of whether tRNA abundance has a  
78 significant impact on translation or whether the observed effect is caused by mRNA secondary  
79 structures alone remains inconclusively answered (29, 54). Enabled by advances in DNA synthesis and  
80 sequencing, some recent works assessed larger numbers of 5'-UTRs or CDSs, again only diversifying  
81 one of the two sequence parts at a time (20,33,36,38,57,58). In a recent study, Arkin and co-workers  
82 combined full-factorial *in silico* design with DNA synthesis on arrays to evaluate the principles of  
83 sequence design for translation in a systematic manner (37). They tested synthetic sequences  
84 combining a single bicistronic 5'-UTR (15,59) with 244,000 CDSs using fluorescence-activated cell  
85 sorting combined with next-generation sequencing (NGS). Several relevant sequence parameters such  
86 as AT-content, codon usage, and mRNA folding were varied and combined in a statistically full-factorial  
87 manner. This was achieved using a sophisticated modular design approach based on *a priori*  
88 hypotheses, which, however, bears the risk of introducing “user-borne” bias.

89 Herein, we describe our efforts to overcome the prevailing lack of knowledge about the impact of  
90 different mRNA parts and sequence parameters on translation with the goal to assess and accurately  
91 quantify their effect. We combine randomly generated 5'-UTRs and CDSs following different assembly  
92 strategies to obtain libraries of random, combinatorial and full-factorial 5'-UTR-CDS combinations.  
93 Using a recently developed method for ultradeep sequence-function mapping (58), we dynamically  
94 assess translation of more than 1.2 million 5'-UTR-CDS pairs in more than 8.8 million sequence-function  
95 data points and different genetic backgrounds. The extremely high throughput and the modular  
96 assembly strategy applied herein allow us to systematically disentangle and assess individual and  
97 combined effects of 5'-UTR and CDS, and to quantify the contribution of various sequence parameters  
98 including individual bases and positions, mRNA secondary structures, 16S-rRNA hybridisation, and  
99 codon usage.

100

## 101 MATERIAL AND METHODS

### 102 Reagents

103 All chemicals were obtained from Sigma Aldrich (Buchs, Switzerland). Restriction enzymes were  
104 obtained from New England Biolabs (Ipswich, USA). PCR was performed using Q5 DNA polymerase  
105 from New England Biolabs (Ipswich, USA). Oligonucleotides (**Suppl. Tab. 1**) were obtained from  
106 Microsynth AG (Balgach, Switzerland). All primers containing degenerate bases were ordered PAGE-  
107 purified. Custom duplex DNA adapters and gene fragments were obtained from Integrated DNA  
108 Technologies (Leuven, Belgium). Plasmid DNA for cloning was extracted with the ZR Plasmid Miniprep  
109 kit from Zymo research (Irvine, USA). Plasmid DNA from cultures used for subsequent sample  
110 preparation for NGS was extracted with the QIAprep Spin Miniprep kit from Qiagen (Hilden, Germany).  
111 Gel extraction of DNA was performed using Zymoclean Gel DNA Recovery Kits from Zymo research  
112 (Irvine, USA).

113

### 114 Strains, cultivation conditions and growth analysis

115 *Escherichia coli* (*E. coli*) TOP10  $\Delta rhaA$  (L-rhamnose isomerase) was used throughout the study. The  
116 generation of these rhamnose utilisation-deficient strain is described elsewhere (58). For experiments  
117 with plasmid-borne variants of tRNA<sup>fMet</sup>, the strain *E. coli* TOP10  $\Delta rhaA \Delta metZ WV$  was generated by  
118 additional replacement of the chromosomal *metZ WV* locus with a spectinomycin resistance cassette  
119 using the method described by Datsenko and Wanner (60). The spectinomycin resistance cassette was  
120 PCR-amplified from a commercial gene fragment (**Suppl. Note 1**) using primers p1 and p2 (**Suppl.**  
121 **Tab. 1**) to generate the linear fragment for transformation complementary to 41 bp both up- and  
122 downstream of the chromosomal *metZ WV* locus. Transformants were verified for successful integration  
123 by colony PCR using primers p3 and p4 and subsequent Sanger sequencing. The exact genotypes of  
124 both *E. coli* strains are provided in **Supplementary Table 2**. *E. coli* cells were generally cultivated in  
125 lysogeny broth (LB) supplemented with 50 mg L<sup>-1</sup> kanamycin, 50 mg L<sup>-1</sup> streptomycin and 10 g L<sup>-1</sup> D-  
126 glucose for repression of the rhamnose-inducible promoter. 15 g L<sup>-1</sup> agar were added for plate cultures.  
127 Cells were grown at 37 °C in an incubator (plates) or shaking incubator at 200 rpm (shake flasks  
128 cultivations). Doubling times of strains with different tRNA<sup>fMet</sup> variants were determined in biological  
129 triplicate cultures as follows. *E. coli* TOP10  $\Delta rhaA$  ("WT") and *E. coli* TOP10  $\Delta rhaA \Delta metZ WV$   
130 (" $\Delta metZ WV$ ") were transformed with pSEVA361 (empty vector), ptRNA<sup>fMet-A37</sup>, ptRNA<sup>fMet-A37G</sup> or  
131 ptRNA<sup>fMet-A37U</sup>, respectively. Sequence-verified transformants of each strain were used to inoculate an  
132 overnight pre-culture in LB (34 mg L<sup>-1</sup> chloramphenicol; 12.5 mg L<sup>-1</sup> spectinomycin for  $\Delta metZ WV$ ). After,  
133 120 mL main cultures in baffled shake flasks (1 L) were inoculated to a starting OD<sub>600</sub> of 0.01 and  
134 incubated shaking (37 °C, 200 rpm). The OD<sub>600</sub> was measured in intervals of 15-30 minutes and  
135 doubling times were determined by dividing ln(2) by the specific growth rate during exponential growth.

136

### 137 Plasmid and library construction

138 A list of plasmids used in this study is provided in **Supplementary Table 3**. Plasmids were constructed  
139 by conventional restriction-ligation cloning. To enable facile library cloning, plasmid pASPIre4  
140 (**Suppl. Fig. 1**) was generated as a derivative of the previously published pASPIre3 (58). pASPIre4

141 additionally contains a *SpeI* restriction site within the CDs of *bxb1* to enable diversification of the 5'-  
142 UTR and codons 2-16 of *bxb1*.

143 Library inserts were generated by PCR with degenerate primers to diversify the respective regions and  
144 inserted into the pASPIre4 backbone thereafter. The fully randomised 5'-UTR-CDS library was  
145 generated via PCR using pASPIre4 as template and primers p5 and p6. After, the PCR product and  
146 pASPIre4 were digested with *SpeI* and *PstI* (37 °C, 3 h), gel purified and ligated (16 °C, T4 ligase,  
147 overnight). The ligation mixture was purified and used to electroporate freshly prepared *E. coli* TOP10  
148  $\Delta rhaA$  cells (61). After 60 min recovery at 37 °C in LB with 10 g L<sup>-1</sup> D-glucose, transformants were plated  
149 in different dilutions for colony counting on LB agar plates (50 mg L<sup>-1</sup> kanamycin, 50 mg L<sup>-1</sup> streptomycin  
150 and 10 g L<sup>-1</sup> D-glucose). After overnight incubation (37 °C), 10 mL LB were added to the plates and  
151 approximately 400,000 colonies were scraped off with a spatula. Glycerol was added to the cell  
152 suspension to a final concentration of 150 g L<sup>-1</sup> and the optical density at 600 nm (OD<sub>600</sub>) of the glycerol  
153 stock was adjusted to 5.0 before freezing of aliquots in liquid nitrogen and storage at -80 °C. This pool  
154 of clones was designated Lib<sub>random</sub> and the corresponding plasmid architecture was termed pASPIre4<sub>lib</sub>  
155 (**Suppl. Fig. 2**). For the uASPIre with mutated tRNA<sup>fMet</sup> variants, a glycerol stock of Lib<sub>random</sub> was plated  
156 on LB agar and plasmid DNA of approximately 50,000 clones was extracted and subsequently used to  
157 transform *E. coli* bearing the respective plasmids for the expression of tRNA<sup>fMet</sup> (see below).

158 Combinatorial and full factorial libraries combining different 5'-UTRs and CDSs were generated in a  
159 stepwise procedure as illustrated in **Supplementary Figure 3**. First, 5'-UTR and CDS half-libraries  
160 (**Suppl. Figs. 4, 5**) were cloned separately as described above. The 5'-UTR half-library was generated  
161 by PCR with primers p5 and p7 on pASPIre4 as template and subsequently inserted into the pASPIre4  
162 backbone using *PstI* and *NotI*. Primer p7 introduces degeneracy in the 5'-UTR and a *BbsI* site between  
163 the randomised 5'-UTR and the *NotI* site (**Suppl. Fig. 3**). The CDS half-library was generated by PCR  
164 with primers p8 and p9 on pASPIre4 as template and inserted into the pASPIre4 backbone using *PstI*  
165 and *NotI*. Primer p8 introduces degeneracy in the CDS and a *BbsI* site between the CDS and the *PstI*  
166 site (**Suppl. Fig. 3**). Transformants of both half-libraries were plated separately in various dilutions.  
167 Depending on the libraries to be created afterwards, a desired number of colonies was scraped off with  
168 a spatula and plasmid DNA was extracted: for Lib<sub>comb1</sub>, approximately 1,000 colonies of the 5'-UTR half-  
169 library and approximately 1,000 colonies of the CDS half-library; for Lib<sub>comb2</sub>, approximately 100 colonies  
170 of the 5'-UTR half-library and approximately 10,000 colonies of the CDS half-library. For Lib<sub>fact</sub>, ten  
171 plates of approximately 100 colonies each of the 5'-UTR half-library and ten plates of approximately  
172 100 colonies each of the CDS half-library were scraped off. In a second step, 5'-UTR and CDS half-  
173 libraries were combined to generate libraries Lib<sub>comb1</sub>, Lib<sub>comb2</sub> and Lib<sub>fact</sub>. To achieve this, plasmid DNA  
174 from the different 5'-UTR half-libraries was PCR-amplified with primers p9 and p10 and the PCR product  
175 was digested with *BbsI* and *PvuI*. Subsequently, these half-libraries were ligated into plasmid  
176 backbones isolated from the individual CDS half-libraries via digestion with *PvuI* and *BbsI*. Note that  
177 the *BbsI* type IIS restriction site enables scarless joining of 5'-UTR and CDS half-libraries using ATGC  
178 (start codon ATG + first downstream base) as sticky ends for ligation. Lib<sub>comb1</sub> (approx. 1,000 5'-UTRs  
179 combined with approx. 1,000 CDSs) and Lib<sub>comb2</sub> (approx. 100 5'-UTRs combined with approx. 10,000  
180 CDSs) were used to transform *E. coli* TOP10  $\Delta rhaA$  yielding approximately 1.5 million and 2.3 million

181 colonies, respectively. Lib<sub>fact</sub> was transformed in ten separate batches (ten times 100 5'-UTRs combined  
182 with 100 CDSs) yielding ten full-factorial sub-libraries. Each of these should contain a maximum of  
183 approximately 10,000 different 5'-UTR-CDS combinations, amongst which theoretically all 5'-UTRs are  
184 combined with all CDSs and *vice versa*. Colonies of these ten sub-libraries were scraped off plates and  
185 pooled to equivalent cell densities according to their OD<sub>600</sub>.

186 All plasmids for overexpression of tRNA<sup>fMet</sup> variants are derivatives of pSEVA361 (62). We selected the  
187 chromosomal *metY* locus including promoters and terminators of *E. coli* TOP10 as a scaffold since it is  
188 monocistronic and therefore simpler to mutate compared to the *metZ WV* locus. In this scaffold we  
189 introduced an A-to-G point mutation at position 47 of the tRNA<sup>fMet</sup> to match the sequence of *metZ WV*  
190 (note that the *metY*-derived tRNA differs by this one base from *metZ WV* tRNAs, which are three  
191 identical tRNA<sup>fMet</sup> copies). The resulting monocistronic design was obtained as commercial gene  
192 fragment in four versions containing the wild-type base (A) as well as three mutants (C, G and T) at  
193 position 37 of tRNA<sup>fMet</sup>, respectively. The gene fragments were cloned into pSEVA361 (p15A replicon,  
194 chloramphenicol resistance) via KpnI and SpeI sites using standard procedures and sequence verified.  
195 The resulting plasmids were designated ptRNA<sup>fMet-A37</sup>, ptRNA<sup>fMet-A37C</sup>, ptRNA<sup>fMet-A37G</sup> and ptRNA<sup>fMet-A37U</sup>  
196 (**Suppl. Fig. 6, Suppl. Tab. 3**) and used to transform *E. coli* TOP10  $\Delta rhaA$  and *E. coli* TOP10  $\Delta rhaA$   
197  $\Delta metZ WV$ . Note that transformants of ptRNA<sup>fMet-A37C</sup> failed to grow and could thus not be included in  
198 further experiments. To assess the effect of tRNA<sup>fMet</sup> mutations, *E. coli* TOP10  $\Delta rhaA$  and *E. coli* TOP10  
199  $\Delta rhaA \Delta metZ WV$  bearing the plasmids for tRNA overexpression were each co-transformed with the  
200 pool of 50,000 variants of Lib<sub>random</sub> (see above).

201

## 202 **Library cultivation, sample preparation and NGS**

203 The different libraries were separately grown in independent shake flask cultivations. Lib<sub>fact</sub> was  
204 cultivated in two biological replicates. Cultivations were conducted in 600 mL LB with 50 mg L<sup>-1</sup>  
205 kanamycin and, in case of tRNA<sup>fMet</sup> overexpression, 34 mg L<sup>-1</sup> chloramphenicol in 5 L baffled shake  
206 flasks. Pre-warmed (37 °C) LB was inoculated from glycerol stocks of the respective libraries to an initial  
207 OD<sub>600</sub> of 0.05. Cultures were grown at 37 °C in a shaking incubator at 200 rpm. At an OD<sub>600</sub> of  
208 approximately 0.5, expression of *bxh1* was induced by addition of 2 g L<sup>-1</sup> L-rhamnose. Samples were  
209 drawn at 0, 95, 225, 290, 360 and 480 minutes after induction and immediately diluted in an excess of  
210 ice-cold PBS. Cell suspensions were centrifuged (4,000 g, 10 min, 4 °C) and pellets were snap frozen  
211 on dry ice. Afterwards, plasmid DNA was extracted and digested with SpeI and NcoI (4 h, 37 °C). Target  
212 fragments containing the 5'-UTR-CDS region and the Bxb1 recombination substrate were purified via  
213 gel electrophoresis (2.5% agarose). Afterwards, duplex DNA adapters for Illumina NGS with sample-  
214 specific indices (**Suppl. Tab. 4**) were ligated to the target fragments and full-length ligation products  
215 were purified via gel electrophoresis (2% MetaPhor agarose, Lonza, Basel, Switzerland). Purity and  
216 concentration of extracted fragments were determined using capillary electrophoresis (Fragment  
217 Analyser, Agilent) and samples were pooled in equimolar ratios. The pool was spiked with 15% PhiX  
218 DNA to increase sample diversity and afterwards sequenced on an Illumina NovaSeq6000 platform (SP  
219 flowcell, paired-end reading with at least 30 cycles forward and 100 cycles reverse read). Primary

220 sequencing data were processed with Illumina RTA version V3.4.4 and bcl2fastq to obtain \*.fastq files  
221 for further processing (see below).

222

### 223 **NGS data processing**

224 NGS raw data analysis was performed using a combination of *bash* and *R* scripts (R version 4.1.2)  
225 running on a Red Hat Enterprise Linux Server (release 7.9). Annotated scripts for raw data processing  
226 will be made available upon final publication.

227 In brief, forward and reverse reads from \*.fastq files were paired. From the forward reads, the identity  
228 of the sample-specific index (six options) and the state of the Bxb1 substrate (either unflipped or flipped),  
229 were extracted through alignment against all possible twelve combinations allowing a maximum of three  
230 mismatches between read and reference to avoid data loss due to sequencing errors. Afterwards, a  
231 similar procedure was applied to the reverse reads to identify the second sample-specific index (six  
232 options). Next, the sample-specific combination of forward and reverse indices was used to split the  
233 data and assigning reads to the different libraries and sampling time points (**Suppl. Tab. 5**).

234 Next, NGS reads with a frameshift within the CDS (e.g. due to sequencing errors or undesired mutations)  
235 were removed by filtering for the correct positioning of the constant first five nucleotides (ATGCG) of  
236 the *bx1* CDS. Then, all 40 randomised nucleotides of 5'-UTR (25 nt) and CDS (each third nucleotide  
237 in codons 2-16; in total 15 nt) were extracted for each read, serving as unique identifier for each variant  
238 (i.e. 5'-UTR-CDS combination). To rescue reads with sequencing errors in the variable regions (less  
239 than 5% of total reads), a clustering procedure was applied to Lib<sub>comb1</sub>, Lib<sub>comb2</sub> and Lib<sub>fact</sub> to map them  
240 to actual (i.e. physically present) variants. This clustering can be applied since the extremely large  
241 theoretical sequence space of these variable regions (40 nt randomised;  $>10^{23}$  possible permutations)  
242 renders the occurrence of highly similar sequences virtually impossible. First, variants were sorted  
243 based on their total read number across all time points. Then, starting with the most frequent variant,  
244 all other variants with a Hamming distance of 1 (i.e. maximum of one substitution) were mapped back  
245 to this variant. This procedure was continued with the next most abundant variant until all remaining  
246 variants were further than one substitution apart from all others. 5'-UTRs and CDSs were treated  
247 separately to keep the computational complexity manageable. For Lib<sub>random</sub>, clustering was omitted  
248 since all 5'-UTRs and CDSs in this library are unique rendering the mapping process computationally  
249 infeasible. Afterwards, the number of reads with unflipped and flipped Bxb1 substrates was counted for  
250 the remaining variants and for each time sample to obtain time-resolved flipping profiles.

251 Lastly, an additional filtering step was performed to ensure high data quality, which excludes variants  
252 with less than 10 reads in at least one of the six time points. Moreover, variants containing an  
253 unintended non-synonymous codon mutation in the CDS were removed (227 variants).

254 This data processing procedure resulted in 1,214,438 high-quality variants split across the 4 libraries  
255 with an average of 464.3 reads per variant or 77.4 reads per variant and time point. For the uASPIre of  
256 tRNA<sup>fMet</sup> mutants, this procedure resulted in 44,289 high-quality variants. In total, this amounts to  
257 8,881,032 sequence-function pairs obtained from three NGS runs. The relative trapezoidal area under  
258 the flipping curve (termed "integral of the flipping profile", IFP) was calculated for each variant. For Lib<sub>fact</sub>,

259

260 the average IFP of the two biological replicates was used. Processed data and annotated scripts for  
261 data processing will be made available upon final publication.

262

### 263 **Correlation of Bxb1 recombination with cellular Bxb1-sfGFP levels**

264 To convert Bxb1-catalysed flipping into relative cellular Bxb1 concentrations, we used the same  
265 approach as described previously, which relies on translational fusion of Bxb1 to the superfolder green  
266 fluorescent protein (sfGFP) and the use of internal standard RBSs (58). In brief, we first recorded the  
267 sfGFP fluorescence of 31 manually constructed RBSs controlling translation of the Bxb1-sfGFP fusion.  
268 These RBSs span a wide range of RBS strengths (from low to high) as previously shown in triplicate  
269 shake flask cultivations (58). A pool of these 31 standard RBSs was cultivated in a separate shake flask  
270 in parallel to the cultivations of Lib<sub>random</sub>, Lib<sub>comb1</sub> and Lib<sub>comb2</sub> and processed alongside the different  
271 libraries as described above. From the resulting NGS data, we obtained the IFP for the standard RBSs  
272 and constructed a calibration curve between IFP and the aforementioned sfGFP fluorescence  
273 measurements (58). A LOESS fit (locally estimated scatterplot smoothing) was used to correlate the  
274 IFP with the slope of the cell-specific sfGFP signal between 0 and 290 minutes after induction (slope  
275 GFP<sub>0-290min</sub>) using the function *loess* from the *R* package *stats*. Relying on the LOESS function, the IFP  
276 values of all library members were converted into the corresponding slope GFP<sub>0-290min</sub>. The resulting  
277 values were normalised to the maximum slope GFP<sub>0-290min</sub> in the entire data and the normalised slope  
278 GFP<sub>0-290min</sub> was designated relative translation rate (rTR) and used for all further analyses. Code and  
279 parameters of the LOESS fit will be made available upon final publication.

280

### 281 **Splitting of full-factorial sub-libraries**

282 Since Lib<sub>fact</sub> consists of ten full-factorial sub-libraries that were sequenced in bulk, the resulting data  
283 had to be computationally split into the sub-libraries for further analysis. Therefore, we sequenced at  
284 least three clones (reference variants) from each sub-library by Sanger sequencing covering both the  
285 randomised 5'-UTR and CDS regions. From the resulting reference sequences, we reconstructed and  
286 split the ten individual sub-libraries as follows: all variants that shared either the 5'-UTR or CDS with  
287 one of the reference sequences were assigned to the corresponding sub-library. To obtain full-factorial  
288 sub-libraries (i.e. libraries in which the majority of 5'-UTRs is combined with each CDS and vice versa),  
289 we further removed all variants with a 5'-UTR that occurred in combination with less than 50 CDSs as  
290 well as all variants with a CDS that occurred in combination with less than 50 5'-UTRs.

291

### 292 **Data analyses**

293 Data analysis was conducted in *R* (version 4.1.2) and figures were produced using the package *ggplot2*.  
294 Scripts will be made available upon final publication.

295 For ANOVA of positional effects, variants from Lib<sub>random</sub> were split according to their respective base in  
296 each of the 40 randomised positions within 5'-UTR and CDS (i.e. 40 splits for 40 position). After, type II

297

298

299



300 ANOVA was performed using the R function *Anova* (package *car*) treating each positional group as  
301 covariate to determine the contribution of each covariate/position to the variance of the rTR in the entire  
302 library assuming additive behaviour. For the assessment of effects of single bases, we calculated the  
303 average rTR of all variants in  $Lib_{random}$  with a given base at a given position and divided the resulting  
304 value by the average rTR of all variants with any other base at this position. For example, the effect of  
305 U at 5'-UTR position -1 was calculated by dividing the average rTR of all variants with U at 5'-UTR  
306 position -1 (0.185) by the average rTR of all other variants (0.150). The resulting value (example: 1.233)  
307 represents the average relative in- or decrease in rTR for a given base and position. In the example  
308 above this means that the rTR of variants with U at 5'-UTR position -1 is on average 23.3% increased  
309 over the rest of the library. To assess the enrichment of bases amongst strong variants, variants in  
310  $Lib_{random}$  were first split into two groups with  $rTR \geq 0.5$  (strong) and  $rTR < 0.5$  (weak). After, the relative  
311 occurrence of each base at each position was calculated within each group. The ratio between the  
312 occurrences in the two groups represents the relative enrichment/depletion of a given base in a given  
313 position amongst strong variants over weak variants.

314 For calculations related to mRNA folding, bash scripts were used. Minimum free energy (mfe),  
315 ensemble free energy (efe) and mRNA accessibility (acc) were each calculated using two models for  
316 base pairing, the turner energy model (T) and the CONTRAfold model (C) (65,66), resulting in six  
317 different metrics (mfeT, mfeC, efeT, efeC, accT, accC). For mfeT and efeT, *RNAfold* (*ViennaRNA*  
318 package, version 2.4.18) and default parameters were used (67). For mfeC and efeC, and default  
319 parameters were applied. For accT and accC, the *Raccess* program was used (68). Next, Spearman's  
320 correlation was calculated between each metric and the rTR. Note that Spearman's correlation was  
321 used since rTR values do not follow a normal distribution (p-value of  $1.11 \times 10^{-79}$  according to Shapiro-  
322 Wilk normality test). Squared Spearman's coefficient ( $\rho^2$ ) is reported as a measure of correlation  
323 between the respective folding metric and the ranked the rTR. Accordingly, the higher  $\rho^2$  of a metric,  
324 the more it explains the observed variance in the rTR. To identify the optimal mRNA sequence window  
325 that leads to the highest correlation between folding and rTR, mfeT and efeT were calculated for all  
326 possible sequence windows of lengths between 10 and 200 nucleotides within the first 200 positions of  
327 the mRNA. For computational reasons, this analysis was performed only on the 10'000 variants of  
328  $Lib_{random}$  with the highest number of NGS reads. The best correlation between folding energy and rTR  
329 was achieved using the first 80 positions of the mRNA (i.e. between positions -27 and +53) (**Suppl. Fig.**  
330 **7**). This "optimal" sequence window was then used to calculate mfeT, mfeC, efeT, efeC, accT and accC  
331 for all variants in all libraries. For accT and accC, the access length was set to 80 nucleotides in *Raccess*.  
332 For accessibility scanning, the correlation between the accessibility of each position and the rTR was  
333 determined applying an access length of 10 nucleotides in *Raccess* (accT<sub>10nt</sub> and accC<sub>10nt</sub>).

334 To calculate 16S rRNA hybridisation energies, *RNA duplex* from the *ViennaRNA* package (67) was used,  
335 which only allows intermolecular base pairing. Allowing intramolecular base pairing would favour 5'-  
336 UTR-internal folds and thus disregard interactions with the 16S rRNA. Specifically, hybridisation energy  
337 was calculated between 5'-UTR (positional window: -18 to -4) and the 16S rRNA 3'-end  
338 (5'-ACCUCCUUA-3'). As an alternative, we also calculated a positional hybridisation energy between  
339 16S rRNA 3'-end and a 9-nt sliding window along the entire mRNA.

340 The minimum edit distance was determined using the *stringdist* function of the R package *stringdist* and  
341 corresponds to the Levenshtein distance between the 7-bp long canonical SD motif AGGAGGU and a  
342 sliding 7-nt window within 5'-UTR positions -18 and -4. Levenshtein distance is the minimum number  
343 of operations (substitutions, deletions, and insertions) to transform one string into another.

344 The random forest model was built using *h2o.randomForest* from the R package *h2o*  
345 (<https://github.com/h2oai/h2o-3>). Variants of Lib<sub>random</sub> were split into a randomly selected training set  
346 (90%) and a test set (10%), which was strictly held out during training. Sequences were encoded using  
347 one-hot encoding, a position-wise accessibility score accC<sub>1nt</sub> (compare above), GC-content, minimum  
348 edit distance to the SD motif AGGAGGU, 16S rRNA hybridisation energy, the position of 16S rRNA  
349 hybridisation on the mRNA, as well as the folding metrics mfeT, mfeC, efeT, efeC, accT and accC (see  
350 above). Using tenfold cross-validation, the model was then trained with default parameters using 50  
351 trees, and its performance was validated on the strictly held-out test set.

352 To quantify the contributions of UTR and CDS, we first grouped variants from Lib<sub>comb1</sub>, Lib<sub>comb2</sub> and  
353 Lib<sub>fact</sub> by their 5'-UTR and then calculated the average rTR of all CDSs in each group (i.e. rTR<sub>UTR</sub>).  
354 Similarly, we also grouped variants by their CDS and calculated the average rTR of all 5'-UTRs in each  
355 group (i.e. rTR<sub>CDS</sub>).

356 Codon adaptation index (CAI) and tRNA adaptation index (tAI) were calculated using the *cai* function  
357 from the R package *seqinr*. Codon weights and frequencies (**Suppl. Tab. 6**) were used as presented in  
358 Sharp *et al.* (41) and dos Reis *et al.*, respectively (53).

359 All sequence variants and their calculated parameters were combined into a single dataset and further  
360 analysed. This data set will be made available upon final publication.

361

### 362 **Data availability**

363 Time series data including IFP and cellular Bxb1-sfGFP values (rTR, see above) for each variant  
364 including annotated scripts for data processing, statistical analyses and plotting will be made  
365 available upon final publication.

366

367

368

## 369 RESULTS

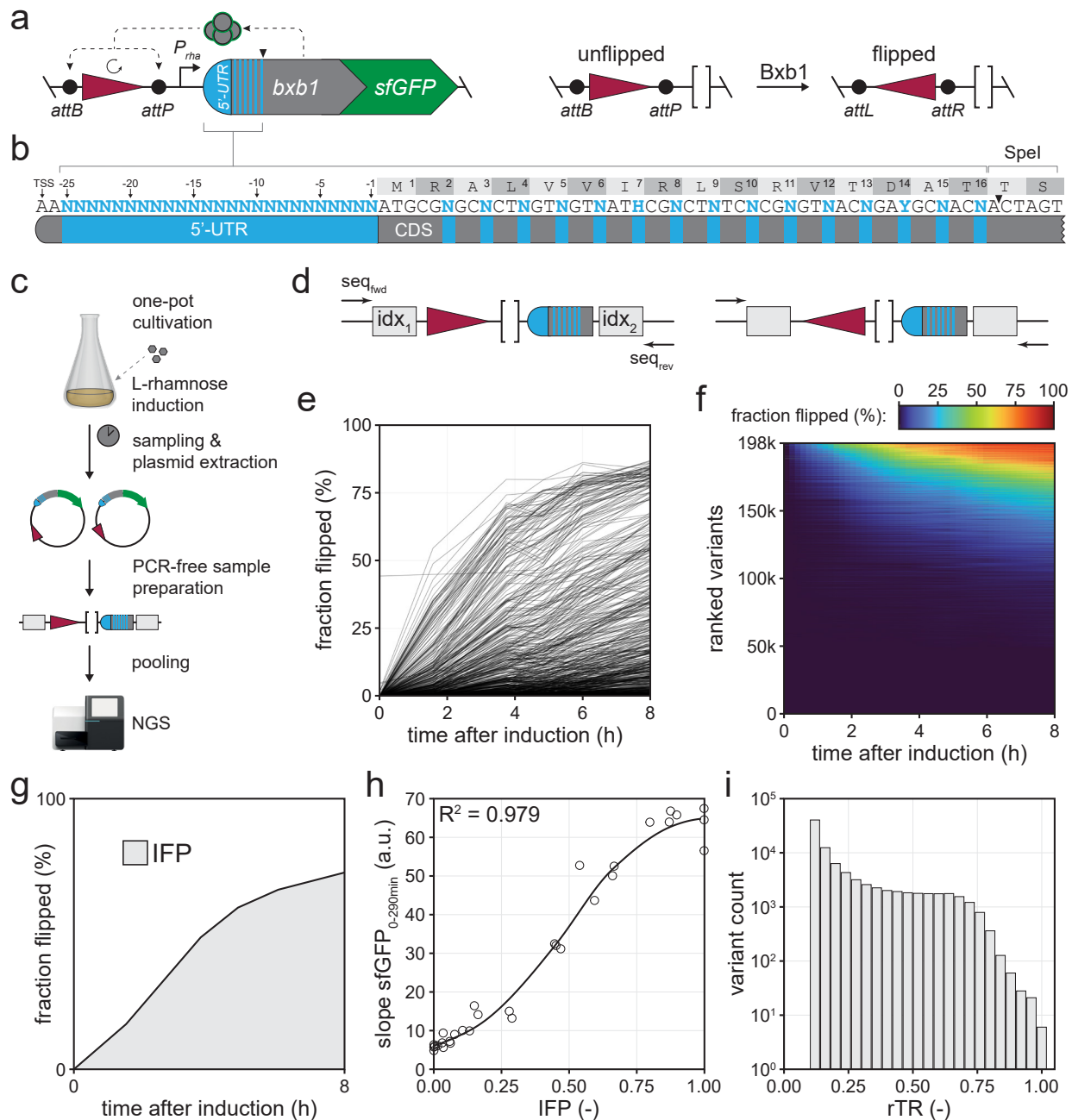
### 370 High-throughput characterisation of 5'-UTR-CDS combinations

371 It is challenging to investigate the impact of different mRNA parts on translation due to the vast  
372 sequence space of possible variants. For instance, even for a comparably short 5'-UTR of twelve  
373 nucleotides, more than 16 million ( $4^{12}$ ) sequences are possible. The sequence space becomes even  
374 larger if different parts are diversified simultaneously, which is required to analyse interactions and  
375 combined effects. Such combinatorial complexity cannot be addressed appropriately by measuring the  
376 expression of a few handpicked sequences. Instead, it requires high-throughput methodology capable  
377 of linking sequences to corresponding expression levels at large scale. To achieve this for 5'-UTR-CDS  
378 combinations, we capitalise herein on a recently developed technology for ultradeep Acquisition of  
379 Sequence-Phenotype Interrelations (uASPIre) (58). Briefly, uASPIre uses the phage recombinase Bxb1  
380 to record functional information in DNA. This DNA-recorder enables, for instance, to determine both  
381 sequence and corresponding gene expression of gene regulatory elements via NGS at extremely high  
382 throughputs, which we have recently demonstrated in a proof-of-concept study (58).

383 To make uASPIre amenable for the characterisation of 5'-UTR-CDS combinations, we created the  
384 plasmid architecture shown in **Figure 1a**, which contains a gene encoding a Bxb1-sfGFP fusion (58)  
385 controlled by an L-rhamnose-inducible promoter ( $P_{rha}$ ) and a 150-bp stretch of silent DNA flanked by  
386 Bxb1's cognate attachment sites *attB* and *attP* in opposite orientation (62). Furthermore, a *SpeI* site is  
387 introduced in codons 17 and 18 of the *bx1* CDS via silent mutation (**Fig. 1a/b**), which enables facile  
388 exchange of the 5'-UTR and the first 16 codons of the *bx1* CDS as well as NGS sample preparation  
389 (**Methods**). Once expressed, Bxb1-sfGFP converts its *attB*-/*P*-flanked DNA substrate from its initial  
390 ("unflipped" hereafter) to an inverted ("flipped" hereafter) state (**Fig. 1a**). Thus, Bxb1-sfGFP expression  
391 can be read out by determining the state of the substrate DNA by sequencing. Importantly, the flipping  
392 rate directly correlates with the cellular Bxb1-sfGFP concentration, and sequencing of many copies of  
393 this architecture via NGS can be used to determine the fraction of flipped DNA substrates ("fraction  
394 flipped" hereafter) amongst all copies of a given variant. This "oversampling" facilitates a precise,  
395 quantitative readout for Bxb1-sfGFP expression, whose resolution solely depends on the sequencing  
396 depth (i.e. number of reads obtained per variant) as we have previously shown (58).

397 Next, we generated a first library through simultaneous diversification of the 5'-UTR and CDS of *bx1*-  
398 *sfGFP* with the goal to characterise the impact on bacterial translation in a highly parallelised fashion  
399 relying on uASPIre (**Fig. 1b, Methods**). We mutated the 25 nucleotides directly upstream of the start  
400 codon applying full randomisation (i.e.  $N_{25}$ -mer, N: equimolar mixture of A, C, G and T). This  
401 corresponds to the entire 5'-UTR in our setup except for two consecutive A's at the 5'-end of the mRNA,  
402 which were fixed to match the native transcriptional start of  $P_{rha}$  and thus avoid changes in transcription  
403 rates (69). Further, we mutated the third positions of codons 2-16 downstream of the start codon (ATG  
404 itself was kept constant) to additionally diversify the CDS. We selected this region since the first 30-50  
405 nucleotides of CDSs reportedly affect translation whereas sequence changes further downstream show  
406 only negligible effects on expression (29,33). Importantly, in this region we only allowed synonymous  
407 ("silent") codon replacements to maintain the same Bxb1 amino acid sequence and hence specific

408 recombination activity for all library members, which is crucial to study only translational effects. This  
409 library is designated Lib<sub>random</sub> hereafter pointing to the full randomisation of 5'-UTR and N-terminal CDS.  
410 Lib<sub>random</sub> was used to transform *E. coli* yielding approximately 400,000 individual transformants.  
411 Specifically, we used the rhamnose-utilisation deficient strain TOP10  $\Delta rhaA$  to ensure temporally stable  
412 induction due to the lack of inducer consumption (58). Afterwards, transformants were pooled and  
413 cultivated in a single shake flask (**Fig. 1c**). In parallel, we cultivated 31 5'-UTR variants ("standard RBSs"  
414 hereafter) controlling the same *bxb1-sfGFP* fusion, which were constructed and characterised in a  
415 previous study (58). These standard RBSs span a wide range of expression levels and serve as internal  
416 standard sequences to compare different experiments. Further, they are used to convert the fraction  
417 flipped time series into practically more relevant metrics for protein expression relying on calibration  
418 curves generated from individual sfGFP fluorescence measurements (see below, **Methods**) (58). After  
419 induction by addition of L-rhamnose, six samples each were drawn over the course of eight hours from  
420 both cultures (Lib<sub>random</sub> and standard RBSs), and plasmid DNA was extracted followed by NGS sample  
421 preparation (**Methods**). Note that sample preparation was carried without PCR amplification, which  
422 avoids non-linear PCR bias (58). The final target DNA fragments are flanked by NGS adapters with  
423 sample-specific indices and contain the DNA substrate modifiable by Bxb1 and the randomised 5'-UTR-  
424 CDS region. NGS adapters, substrate and 5'-UTR-CDS region were sequenced in an Illumina platform  
425 yielding approximately  $10^8$  paired-end reads for Lib<sub>random</sub>. (**Fig. 1d**)  
426 Next, we processed the NGS data to obtain time series of Bxb1-mediated flipping ("flipping profiles")  
427 using a previously developed computational pipeline adapted to the new plasmid architecture  
428 (**Methods**) (58). This procedure yielded flipping profiles for 198,174 5'-UTR-CDS pairs above an applied  
429 minimal threshold of ten reads per time point and variant (i.e. high-quality data, average of 433.7 reads  
430 per variant). The base composition in Lib<sub>random</sub> was homogeneously distributed across all diversified  
431 positions (**Suppl. Fig. 8**). Library members showed a diverse range of translational activities from low  
432 to high and a skew towards weaker variants as to be expected for full randomisation of the 5'-UTR  
433 (**Fig. 1e, f**) (70). Notably, the behaviour of the standard RBSs correlated strongly with results from our  
434 previous study even though the experiments were carried out approximately two years apart from each  
435 other (**Suppl. Fig. 9**) (58). This confirms the validity of the recorded data and indicates a high  
436 reproducibility and robustness of the uASPIre method in general. Next, we calculated the trapezoid  
437 integral of the flipping profiles (IFP, **Fig. 1g**), which constitutes a robust metric correlating well with rates  
438 of cellular Bxb1-sfGFP accumulation as previously shown (58). Indeed, the IFP of the 31 standard RBSs  
439 as determined in this study correlated well with the linear slope of the cell-specific Bxb1-sfGFP  
440 fluorescence between 0 and 290 minutes after induction (slope sfGFP<sub>0-290min</sub>, **Fig. 1h, Methods**).  
441 Therefore, IFP values can be converted into the slope sfGFP<sub>0-290min</sub> relying on a fit applied between the  
442 two metrics for the standard RBSs. Specifically, we performed locally estimated scatterplot smoothing  
443 (LOESS) (**Fig. 1h**), and used the resulting fit function to convert the IFPs of Lib<sub>random</sub> members into the  
444 corresponding slope sfGFP<sub>0-290min</sub> normalised to the strongest variant found in this study (**Fig. 1i**,  
445 **Methods**). This normalised parameter was designated relative translation rate (rTR) and used for all  
446 further analyses, because it represents a practically more relevant metric for translational activity  
447 directly corresponding to cell-specific protein accumulation.



448

449 **Figure 1: Ultradeep characterisation of 5'-UTR-CDS combinations.** **a)** Plasmid architecture for the

450 uASPIre of 5'-UTR-CDS pairs. A *bxb1-sfGFP* gene (translational fusion) controlled by *P<sub>rha</sub>* is placed on

451 the same DNA molecule as the substrate modifiable by Bxb1-sfGFP, which is flanked by Bxb1

452 attachment sites (*attB/I/P*). A *Spel* site in codons 17 and 18 of *bxb1-sfGFP* allows for seamless exchange

453 of 5'-UTR and N-terminal CDS. Once expressed, Bxb1-sfGFP inverts its substrate from an unflipped

454 into a flipped state creating recombined attachment sites (*attL/R*).

455 **b)** Design of *Lib<sub>random</sub>*. The 25 nucleotides preceding the start codon are fully randomized. Additionally, the third positions of codons

456 2-16 are mutated allowing only synonymous codon replacements. Sequences follow the IUPAC

457 nucleotide code (N: A/C/G/T, H: A/C/T, Y: C/T). TSS: transcriptional start site of *P<sub>rha</sub>*.

458 **c)** Experimental workflow for the uASPIre of 5'-UTR-CDS pairs. Pooled transformants of *Lib<sub>random</sub>* are grown in LB and

459 *bxb1-sfGFP* expression is induced by L-rhamnose addition. After, samples are taken at different time

460 points followed by plasmid extraction and preparation of NGS fragments followed by pooling of samples

461 and NGS (**Methods**). NGS fragments are flanked by duplex adapters with sample-specific index  
462 combinations (grey boxes). **d**) Close-up view of target fragments for paired-end NGS using forward  
463 ( $seq_{fwd}$ ) and reverse ( $seq_{rev}$ ) sequencing primers. Forward reads are used to identify the first index ( $idx_1$ )  
464 and the state of the recombinase substrate. Reverse reads are used to obtain the second index ( $idx_2$ )  
465 and the sequence of 5'-UTR and CDS. **e**) Representative flipping profiles of 5'-UTR-CDS variants from  
466  $Lib_{random}$ . For clarity, only the 1,000 most abundant variants are displayed. **f**) Flipping profiles of all  
467 198,174  $Lib_{random}$  members above high-quality read-count threshold (**Methods**). Horizontal lines are  
468 time series of individual variants coloured according to the fraction flipped and ranked by the average  
469 fraction flipped across all time points from high (top) to low (bottom). **g**) Illustration of the IFP (grey area),  
470 i.e. the normalised trapezoidal integral of the flipping profile. **h**) Correlation between IFP and slope  
471  $sfGFP_{0-290min}$  as shown for 31 standard RBSs (**Methods**). A LOESS function (black line) can be used to  
472 interconvert IFP and slope  $sfGFP_{0-290min}$  with high confidence. **i**) Histogram of the rTR of all variants  
473 from  $Lib_{random}$ .

474

### 475 **Analysis of positional and base-specific effects on translation**

476 Relying on the data generated for  $Lib_{random}$ , we investigated the impact of different positions, nucleotides,  
477 and sequence motifs on expression. To assess positional effects, we performed analysis of variance  
478 (ANOVA) treating each variable position in the 5'-UTR (-25 to -1) and CDS (third positions of codons 2-  
479 16) as a covariate and calculated the contribution to the observed variance in rTR (**Fig. 2a, Methods**).  
480 Individual positions in the 5'-UTR explain between 0.3 and 1.5% of the variance. The most pronounced  
481 effect was observable for positions -13 to -8, which corresponds to an anticipated SD region, and, more  
482 unexpectedly, position -1. Within the CDS, the impact of codons decreases with increasing distance  
483 from the start codon with codon 2 showing the highest contribution (2.1%). Codons 2 to 8 show a  
484 marked effect, which strongly decreases to a negligible degree thereafter. Notably, the cumulative  
485 contribution of all 40 randomised positions only amounts to about 25% of which about 17.5% and 7.4%  
486 are attributed to 5'-UTR and CDS, respectively (**Suppl. Fig. 10**). The remaining high fraction of  
487 unexplained variance (about 75%) points towards a strong interaction between positions leading to non-  
488 additive behaviour. Next, we calculated the effect of specific bases at the variable positions by dividing  
489 the average rTR of variants with a given base at a position by the average rTR of all other variants  
490 (**Fig. 2b**). Generally, C and G tend to have a negative, and A and U a positive effect on translation,  
491 which is stronger in the 5'-UTR and weaker in the CDS decreasing with increasing distance to the start  
492 codon. A striking exception to that end are positions -14 to -7 (SD region), for which the effect of G is  
493 highly positive. The strongest negative effect is observable for CGG as the 2<sup>nd</sup> codon (Arg) with  
494 corresponding variants being on average 26.3% weaker than those with CGA, CGC or CGU in this  
495 codon. The strongest positive impact is associated with U at 5'-UTR position -1 amounting to a mean  
496 rTR increase of 23.3%. Finally, to identify characteristic sequence determinants in particular of strong  
497 variants, we split the data from  $Lib_{random}$  into two sets of strong variants (i.e.  $rTR \geq 0.5$ ; 11,212 sequences)  
498 and weaker variants (i.e.  $rTR < 0.5$ ; 186,962 sequences) and calculated the relative enrichment or  
499 depletion of each base at each position in the strong over the weaker subset (**Fig. 2c, Methods**). This  
500 analysis confirmed that both 5'-UTR and CDS of strong variants are generally enriched for A and U,



## 521 **Quantification of sequence parameters and their effect on translation**

522 Since less than 30% of variance in translation could be explained by global analysis of individual  
523 positions, we sought to examine the impact of different sequence parameters on the level of individual  
524 variants. Specifically, we computed several parameters known or hypothesised to influence rTR for all  
525 members of Lib<sub>random</sub> and calculated their correlation with rTR. This analysis included parameters related  
526 to GC-content, hybridisation between mRNA and 16S rRNA, mRNA folding and other features. Since  
527 rTR values follow a non-normal distribution ( $p$ -value =  $1.11 \times 10^{-79}$ , Shapiro-Wilk normality test) and  
528 some sequence parameters are likely to non-linearly correlate with rTR, we also report Spearman's  
529 correlation (coefficient  $\rho$ ) as a metric of rank correlation between parameters and rTR.

530 Overall GC-content shows significant correlation with the rTR ( $\rho^2 = 18.6\%$ ,  $R^2 = 11.3\%$ ) and its impact  
531 is higher in the 5'-UTR than the CDS (**Fig. 3a, Suppl. Fig. 11**). In particular high GC-content is strongly  
532 associated with low rTRs (**Suppl. Fig. 11**), likely due to a tendency of GC-rich sequences to form stable  
533 secondary structures, which are known to counteract translation (27). Further, we determined the  
534 minimum free energy (mfe), ensemble free energy (efe) and mRNA accessibility (acc) using two models  
535 for base pairing, the Turner energy model (T) and the CONTRAfold (C) model (65,66), resulting in six  
536 metrics related to mRNA folding: mfeT, mfeC, efeT, efeC, accT and accC (**Fig. 3a, Methods**). In brief,  
537 mfe and efe are energies required for the unfolding of the most likely and the ensemble of possible  
538 mRNA secondary structure(s), respectively, whereas acc is an accessibility score for a sliding window  
539 along the mRNA corresponding to the probability of this window being embedded within a secondary  
540 structure (18). Folding of mRNA showed a clear impact on rTR across all tested metrics (**Fig. 3a**). The  
541 latter show a positive correlation with the rTR, which is stronger than for GC-content and highest for  
542 efeC ( $\rho^2 = 30.8\%$ ,  $R^2 = 12.6\%$ ) and accC ( $\rho^2 = 30.4\%$ ,  $R^2 = 12.2\%$ ) (**Fig. 3a, b**). In particular very strong  
543 folding (e.g. efeC <  $-15 \text{ kcal} \times \text{mol}^{-1}$ ) completely abolishes efficient translation (**Fig. 3b**). We investigated  
544 further the impact of the positioning of secondary structures by calculating mRNA accessibility within a  
545 sliding window of ten nucleotides. Correlation of the resulting scores (accT/C<sub>10nt</sub>) with rTR is highest  
546 around the first few codons followed by the SD region, and sharply decreases further downstream in  
547 the CDS (**Fig. 3c**).

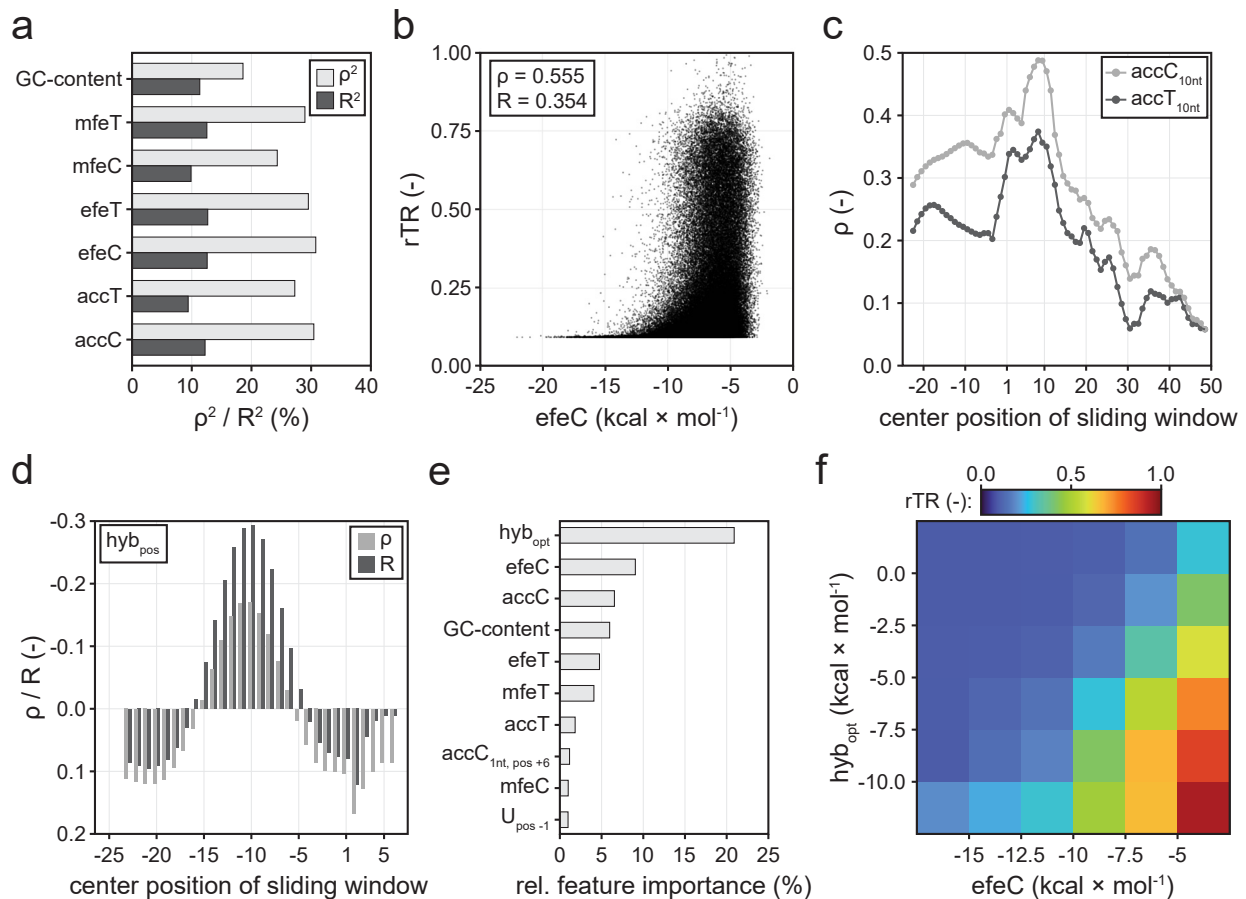
548 Next, we investigated the impact of interactions between mRNA and 16S rRNA. As expected, the  
549 hybridisation energy  $\text{hyb}_{\text{SD}}$  between *E. coli*'s 16S rRNA (sequence: 5'-ACCUCCUUA-3') and the  
550 approximate SD region in the 5'-UTR (window between positions -18 and -4) shows a clear correlation  
551 with the rTR (**Suppl. Fig. 12, Methods**) (67). This observation is further corroborated by the fact that  
552 similarity with the canonical SD motif AGGAGGU in this window is strongly associated with high rTRs  
553 (**Suppl. Fig. 13**). Since the position of hybridisation is known to be critical for efficient translation, we  
554 further calculated positional hybridisation energies  $\text{hyb}_{\text{pos}}$  sliding the 9-nt 16S rRNA sequence along the  
555 mRNA (**Fig. 3d, Methods**). We found that  $\text{hyb}_{\text{pos}}$  is negatively correlated with rTR between 5'-UTR  
556 positions -15 and -6 indicating that stronger hybridisation (i.e. lower  $\text{hyb}_{\text{pos}}$ ) has a translation-favouring  
557 effect in this region. Outside of this window, a negative effect on rTR is observable. The 9-nt  
558 hybridisation window with the strongest correlation to rTR is centred around position -10 corresponding  
559 to a binding of the 16S rRNA 3'-end to the 5'-UTR between positions -14 and -6. A more systematic  
560 analysis of hybridisation windows and positions (**Suppl. Tab. 7**) revealed the mean of hybridisation



561 energies at positions -11 and -10 ( $hyb_{opt}$ ) as the parameter with the highest correlation with rTR  
562 ( $\rho^2 = 2.9\%$ ,  $R^2 = 8.9\%$ ).

563 Based on those findings, we sought to quantify the utility of different sequence parameters for predictive  
564 modelling. To this end, we used the data from  $Lib_{random}$  to train a random forest regressor with the goal  
565 to predict the rTR from different features including primary sequence information as well as the above-  
566 mentioned secondary parameters (**Methods**). The model was trained using tenfold cross-validation  
567 (**Suppl. Fig. 14**) and its performance was evaluated on a test set strictly held out during training  
568 (randomly selected 10% of data). The resulting model predicts rTR values with good confidence ( $R^2 =$   
569  $58\%$ , **Suppl. Fig. 15**). More importantly, we extracted the relative importance of features of the random  
570 forest (**Fig. 3e**). Remarkably, while the 16S rRNA hybridisation parameter  $hyb_{opt}$  had shown only  
571 moderate correlation coefficients  $\rho$  and  $R$ , it was by far the most important model feature (20.9%)  
572 followed by the folding parameters  $efeC$  (9.1%) and  $accC$  (6.5%). The over-proportional importance of  
573  $hyb_{opt}$  could imply that successful hybridisation with the 16S rRNA must be fulfilled to obtain strong  
574 translation initiation rendering  $hyb_{opt}$  a critical, early decision criterion for the model. Furthermore, U at  
575 5'-UTR position -1 ranked 10<sup>th</sup> (1.0%) amongst the total of 248 encodings constituting the most  
576 important single-nucleotide feature. The majority of features (227) exhibited a relative importance below  
577 0.5% pointing towards the multifactorial, interactive nature of the translation (initiation) process and  
578 likely to a high degree of redundancy between the tested encodings.

579 Lastly, we binned the variants from  $Lib_{random}$  according to the two most important features of the random  
580 forest,  $hyb_{opt}$  and  $efeC$ , and calculated the average rTR of each bin (**Fig. 3f**). Interestingly, we found  
581 that the appearance of very high rTRs (i.e.  $> 0.5$ ) is co-dependent on strong 16S rRNA hybridisation  
582 and weak mRNA folding. Variants with strong secondary structures ( $efeC < -15 \text{ kcal} \times \text{mol}^{-1}$ ) only exhibit  
583 significant translation initiation if they hybridise well with the 16S rRNA. By contrast, variants with low  
584 folding energy can exhibit intermediate-to-strong translation even in the absence of SD motifs.



585  
 586 **Figure 3: Effect of different sequence parameters on translation in  $Lib_{\text{random}}$ .** **a)** Correlation of GC-  
 587 content and different mRNA folding metrics with rTR. Spearman's  $\rho^2$  and Pearson's  $R^2$  are displayed.  
 588 **b)** Scatterplot between rTR and the best-correlating mRNA folding parameter efeC. **c)** Correlation of  
 589 rTR with local mRNA accessibility. Parameters  $accT_{10nt}$  and  $accC_{10nt}$  correspond to the mRNA  
 590 accessibility of a 10-nt window centered around the mRNA position specified on the horizontal axis.  
 591 Endings C and T denote base pairing calculated by two different energy models (**Methods**). **d)**  
 592 Correlation of hybridisation energy between 16S rRNA and different mRNA positions with rTR.  
 593 Positional hybridisation energy ( $hyb_{\text{pos}}$ ) is displayed for 9-bp windows centered around the indicated  
 594 mRNA position (horizontal axis). **e)** Relative feature importance of a random forest model trained on  
 595  $Lib_{\text{random}}$ . The ten most important of 248 features are displayed.  $hyb_{\text{opt}}$ : best-correlating hybridisation  
 596 parameter (see main text).  $accC_{1nt, \text{pos}+6}$ : AccC score for position +6 of the mRNA.  $U_{\text{pos}-1}$ : one-hot  
 597 encoded U at position -1 of the mRNA. **f)** Mean rTR of variants in  $Lib_{\text{random}}$  as grouped by the two most  
 598 predictive features of the random forest,  $hyb_{\text{opt}}$  and efeC. Tick labels mark the boundaries of the  
 599 respective bins (boxes).

600

601

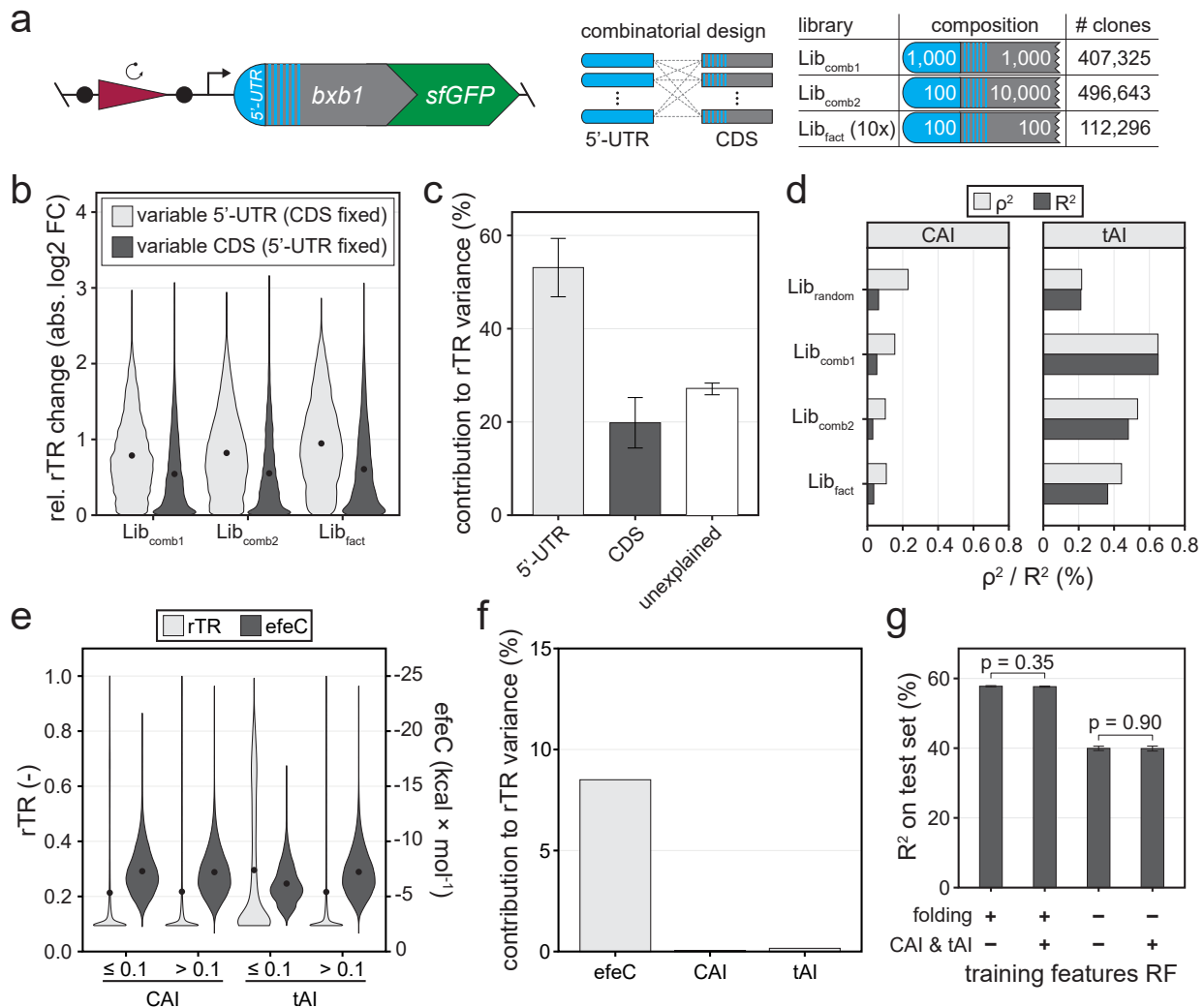
### 602 Codon usage and interaction between 5'-UTR and CDS

603 A long-standing question is how strong the impact of the CDS on translation is, both in absolute terms  
 604 and relative to the 5'-UTR. Changes in the CDS affect critical determinants of translation initiation such  
 605 as codon usage and mRNA folding. Importantly, testing many different CDSs in combination with a  
 606 single 5'-UTR (as amply done in previous studies) is insufficient to unambiguously assign observed

607 effects to different sequence parameters and to quantify their contribution in a precise fashion, since  
608 some parameters also depend on and change with the 5'-UTR in place. Thus, it remains unclear if and  
609 how strong any observed effect is causally related to a sequence parameter change in a generalisable  
610 fashion, or whether it is merely a context-specific artefact only occurring for the selected 5'-UTR.  
611 Similarly, full randomization (as in Lib<sub>random</sub> in this work) only delivers unique pairs of 5'-UTRs and CDSs,  
612 which again prohibits unambiguous attribution of effects to either of the two mRNA parts (5'-UTR or  
613 CDS). This problem can only be circumvented by testing large numbers of 5'-UTR-CDS combinations  
614 in a combinatorial manner with sufficient overlap allowing to average out case-specific artefacts.  
615 Therefore, to investigate the individual impact of 5'-UTR and CDS independently, we generated three  
616 additional libraries of combinatorial (Lib<sub>comb1</sub>, Lib<sub>comb2</sub>) and full-factorial (Lib<sub>fact</sub>) 5'-UTR-CDS pairs, which  
617 were constructed through combination of defined half-libraries (**Fig. 4a, Methods**): Lib<sub>comb1</sub> combines  
618 about 1,000 5'-UTRs with about 1,000 CDSs, Lib<sub>comb2</sub> is a combination of approximately 100 5'-UTRs  
619 with approximately 10,000 CDSs, and Lib<sub>fact</sub> features ten independently cloned batches of about 100  
620 5'-UTRs combined with about 100 CDSs each. Note that Lib<sub>fact</sub> was designed such that in each batch  
621 every 5'-UTR is combined with every CDS and *vice versa* (i.e. full-factorial design). Next, we recorded  
622 the activity of variants from the three libraries applying the same uASPIre workflow as described for  
623 Lib<sub>random</sub> above. Processing of NGS data yielded time series for 407,325, 496,643, 112,296 unique  
624 variants above high-quality read count threshold for Lib<sub>comb1</sub>, Lib<sub>comb2</sub> and Lib<sub>fact</sub>, respectively. For Lib<sub>fact</sub>,  
625 two independent biological replicates were tested. We then grouped variants according to the 5'-UTR  
626 (or CDS) in place and analysed the diversity of the rTR amongst all CDSs (or 5'-UTRs) appearing with  
627 the respective fixed 5'-UTR (or CDS). Exchanging either 5'-UTR or CDS (while maintaining the other)  
628 can lead to strong up- and downshifts in expression (**Fig. 4b**). Shifts are on average much stronger for  
629 an exchange of the 5'-UTR than of the CDS, and in many cases cover a large fraction of the rTR range  
630 (**Fig. 4b, Suppl. Fig. 16**). We further quantified the individual impact of 5'-UTR and CDS performing an  
631 ANOVA with the mean rTRs of all 5'-UTRs and CDSs (**Fig. 4c, Methods**). This analysis was performed  
632 exclusively on Lib<sub>fact</sub>, since full-factorial design is required to exclude case-specific artefacts and achieve  
633 a precise quantification of each part's individual contribution (see above). We observed that the 5'-UTR  
634 explains on average  $53.12 \pm 6.3\%$  and the CDS  $19.8 \pm 5.4\%$  of rTR variance confirming the higher  
635 impact of the 5'-UTR compared to the CDS. The  $27.0 \pm 1.3\%$  of variance remain unexplained in the  
636 additive model and must therefore be caused by non-linear interactions between 5'-UTR and CDS  
637 confirming a high degree of interdependence between both parts.

638 A controversially discussed sequence feature of the CDS is codon usage, which is well known to  
639 influence translation (initiation). To this end, the appearance of rare codons within the first few triplets  
640 of the CDS was found to coincide with high expression (29,41-44). Thus, we first analysed the impact  
641 of two commonly used metrics for codon usage, CAI and tAI (**Suppl. Tab. 6**) (41,53), on rTR, which  
642 indicated a weak ( $R^2$  and  $p^2$  consistently below 0.7%) yet significant correlation in all libraries (**Fig. 4d**).  
643 However, it remains unclear whether this is caused by differential abundance of the corresponding  
644 tRNAs in the cell or by changes in mRNA folding. Since folding is also co-dependant on the 5'-UTR in  
645 place, combinatorial testing of 5'-UTR-CDS pairs is also essential in this case to unambiguously test if  
646 and to which extent the two aforementioned hypotheses are correct. Accordingly, we first compared the

647 rTR of Lib<sub>fact</sub> variants rich in rare codons (i.e. CAI/tAI  $\leq$  10%) with the other variants (i.e. CAI/tAI  $>$  10%).  
648 Variants with low CAI exhibit a mean rTR of 0.213, which is virtually indifferent from high-CAI variants  
649 (0.217) (**Fig. 4e**). This is further corroborated by the fact that the mean rTRs of CDSs and CAI do not  
650 correlate significantly (p-value = 0.256, one sample t-test) in Lib<sub>fact</sub> (**Suppl. Fig. 17**). Low-tAI variants,  
651 by contrast, exhibits on average a higher rTR than the control group (**Fig. 4e, Suppl. Fig. 17**). At the  
652 same time, however, mRNA folding is significantly weaker (p-value  $<$   $10^{-300}$ , one-sided Welch two  
653 sample t-test) in low- versus high-tAI variants, which is not the case for the corresponding CAI groups  
654 (p-value = 1.0, **Fig. 4e**). Moreover, the codon frequency of *E. coli* showed only very small and  
655 inconsistent effects on the rTR for the randomised codons (**Suppl. Fig. 18**). Therefore, we further  
656 analysed to which extent the dependence of the rTR on codon usage can be explained by mRNA folding.  
657 An ANOVA with only efeC, CAI and tAI as covariates indicated that the overwhelming majority of  
658 variance in rTR explainable by these parameters is attributed to efeC (8.5%), whereas the contribution  
659 CAI and tAI was about 155- and 53-fold lower, respectively (**Fig. 4f**). Furthermore, we re-trained the  
660 former random forest model (see above) with different sets of sequence parameters including CAI and  
661 tAI (**Fig. 4g**). Remarkably, while removal of mRNA folding parameters led to a substantial decrease in  
662 model performance, addition of CAI and tAI did neither increase accuracy of the initial random forest  
663 nor was it able to compensate for the performance loss in the absence of folding parameters.  
664 Accordingly, the relative feature importance of CAI and tAI was very low (**Suppl. Fig. 19**). Collectively,  
665 these findings strongly suggest that any influence of codon usage on rTR can be virtually completely  
666 explained by mRNA folding. On the contrary, a causal connection to cellular tRNA abundance or the  
667 previously postulated translational ramps could not be established and is either insignificant or  
668 negligible amongst the over 1.2 million sequences tested in this study.



669

670 **Figure 4: Overall impact of 5'-UTR, CDS and codon usage on translation.** **a)** Three additional  
671 libraries of combinatorial (Lib<sub>comb1</sub>, Lib<sub>comb2</sub>) and full-factorial (Lib<sub>fact</sub>) design were assessed via uASPIre.  
672 Lib<sub>comb1</sub>: combinatorial combination of about 1,000 5'-UTRs and 1,000 CDSs. Lib<sub>comb2</sub>: combinatorial  
673 combination of about 100 5'-UTRs and 10,000 CDSs. Lib<sub>fact</sub>: ten independent batches, each a full  
674 factorial combination of approx. 100 5'-UTRs and 100 CDSs. Lib<sub>fact</sub> was tested in two independent  
675 biological replicates. The number of analysed clones is indicated for each library. **b)** Impact of the  
676 exchange of 5'-UTRs or CDSs on translation. The rTR change (absolute value) of a given 5'-UTR upon  
677 exchanging its CDS versus the mean rTR of all variants with that same 5'-UTR is displayed (and *vice*  
678 *versa*). Black circles within violins are mean relative rTR changes. **c)** ANOVA with the mean rTRs of all  
679 5'-UTRs and CDSs in Lib<sub>fact</sub>. Error bars: standard deviation between ten independent batches of Lib<sub>fact</sub>.  
680 **d)** Correlation of codon usage indices CAI and tAI with rTR. **e)** Comparison of rTRs and folding energies  
681 (efeC) of variants with low ( $\leq 10\%$ ) and high ( $> 10\%$ ) CAI/tAI in all libraries. Black circles within violins  
682 are mean rTR/efeC values. **f)** Contribution of efeC, CAI and tAI to the rTR variance in all libraries  
683 according to an ANOVA with only the three parameters as covariates. **g)** Impact of folding and codon  
684 usage metrics on the performance of random forest (RF) models trained on Lib<sub>random</sub>. Sequence  
685 parameters for mRNA folding (mfeT, mfeC, efeT, efeC, accT and acc) and codon usage (CAI and tAI)  
686 were added or omitted during training. Error bars: Standard deviation of five training repeats with 10-  
687 fold cross-validation each. p-value were calculated with Welch two sample t-tests.

## 688 **Assessment of translational anomalies of arginine codon 2 and 5'-UTR position -1**

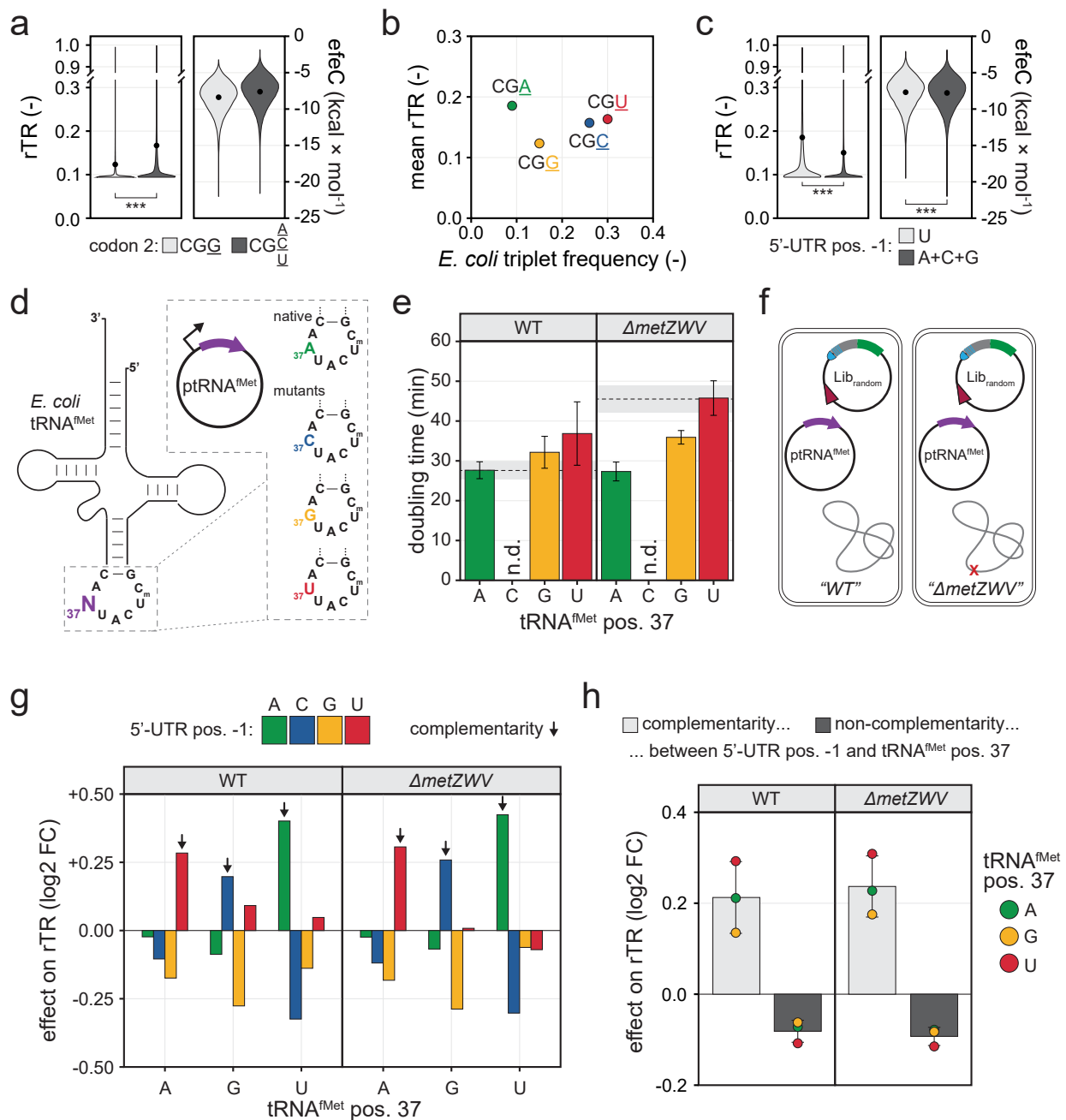
689 Lastly, we sought to decipher the reasons for the unexpected behaviour of two mRNA positions  
690 observable in our data (see above). To this end, the presence of G in the third position of arginine codon  
691 2 and U in position -1 of the 5'-UTR exhibit a profound impact on the rTR, which is negative in the former  
692 and positive in the latter case (compare **Fig. 2**). Variants with CGG as the second codon show an  
693 average decrease in rTR of 26.3% compared to variants carrying A, C or U in the third position (**Fig. 5a**).  
694 This different behaviour is likely not caused by codon frequencies or tRNA availability, since both  
695 arginine codons with higher (CGC, CGU) and lower (CGA) frequency show significantly higher mean  
696 rTRs (**Fig. 5b**). By contrast, the average folding energy of CGG-bearing variants is significantly lower  
697 ( $\Delta_{\text{efeC}} = -0.76 \text{ kcal} \times \text{mol}^{-1}$ ) than for the other codons (**Fig. 5a**), pointing again to mRNA folding (and  
698 not tRNA availability) as the mechanistic reason for the differential expression of synonymous codons.  
699 For variants with U at position -1 in the 5'-UTR, the mean rTR is 23.3% higher than for those with any  
700 other base in this position (**Fig. 5c**). In this case, however, the average folding energy is even slightly  
701 increased for U ( $\Delta_{\text{efeC}} = +0.10 \text{ kcal} \times \text{mol}^{-1}$ ) excluding mRNA folding as the reason (**Fig. 5c**). As an  
702 alternative explanation, we suspected that an interaction of this U with the initiator tRNA (tRNA<sup>fMet</sup>) could  
703 be responsible for the observed effect. In *E. coli*, initiator tRNAs are encoded by one monocistronic  
704 (*metY*) and one tricistronic (*metZWV*) transcriptional unit, and their sequences are identical except for  
705 position 46 (G in *metY*, A in *metZWV*). Importantly, methionine elongator tRNAs (*metT*, *metU*) do not  
706 initiate translation (71), and all tRNA<sup>fMet</sup> copies carry an A in position 37 directly 3' to the CAU anticodon,  
707 which could preferentially hybridise with mRNAs carrying a U directly 5' to the start codon.

708 Several previous studies have postulated or shown that the presence of U in this position favours  
709 formation of the prokaryotic ribosomal initiation complex and/or translation of the corresponding genes  
710 *in vitro* and *in vivo* (36,58,72-80). These effects were attributed to a proposed interaction of A37 in  
711 tRNA<sup>fMet</sup> and U in 5'-UTR position -1, for which further evidence was later provided in algal chloroplasts  
712 through compensatory mutation of tRNA<sup>fMet</sup> position 37 (81). Furthermore, structural analyses have  
713 shown that A37 is released from internal base pairing upon reaching the ribosomal P-site (82), which  
714 would render this position available for Watson-Crick base pairing with nucleotide(s) upstream of the  
715 start codon. Collectively, these prior works highlight the importance of bases directly upstream of the  
716 start codon and point to a potential interaction of mRNA and tRNA<sup>fMet</sup> beyond the codon-anticodon  
717 hybridisation. A causal link between any observed impact on translation and an interaction with the  
718 5'-UTR position -1 was, however, so far not conclusively established. A potential reason for this could  
719 be that only few mRNA sequence variants were tested prohibiting generalisable statements due to the  
720 high context dependence of translation initiation and statistical error.

721 We therefore investigated whether the proposed interaction between mRNA and tRNA<sup>fMet</sup> could be  
722 substantiated relying on systematic high-throughput sequence-function mapping. We first constructed  
723 plasmids for the overexpression of tRNA<sup>fMet</sup> with the native A37 as well as the mutants A37C, A37G  
724 and A37U (**Fig. 5d, Methods**). To reduce the background from the chromosomal tRNA<sup>fMet</sup> copies, we  
725 further deleted the *metZWV* locus of *E. coli* TOP10  $\Delta rhaA$  ("WT") yielding strain TOP10  $\Delta rhaA$   
726  $\Delta metZWV$  (" $\Delta metZWV$ "), and transformed both strains with the tRNA plasmids. Note that simultaneous  
727 knockout of *metZWV* and *metY* failed in our hands despite complementation via plasmid-borne tRNA<sup>fMet</sup>.

728 Remarkably, transformants of  $\text{tRNA}^{\text{fMet-A37C}}$  showed severe growth inhibition (colonies visible only few  
729 days after transformation), whereas the native  $\text{tRNA}^{\text{fMet-A37}}$  and the other two mutants ( $\text{tRNA}^{\text{fMet-A37G}}$ ,  
730  $\text{tRNA}^{\text{fMet-A37U}}$ ) were tolerated with minor effects on growth in both strains (**Fig. 5e**). While in the case of  
731 the WT strain a small increase of doubling times was observable,  $\Delta\text{metZ WV}$  showed an improvement  
732 of growth upon overexpression of all  $\text{tRNA}^{\text{fMet}}$  variants, likely due to compensation of the reduced level  
733 of chromosomally-derived  $\text{tRNA}^{\text{fMet}}$  copies in this strain. The apparent toxicity of  $\text{tRNA}^{\text{fMet-A37C}}$  could  
734 stem from global dysregulation of translational, and due to its prohibitively slow growth we excluded this  
735 variant from further experiments. Next, we tested approximately 50,000 variants from  $\text{Lib}_{\text{random}}$  in both  
736 strains (WT,  $\Delta\text{metZ WV}$ ) in presence of the remaining  $\text{tRNA}^{\text{fMet}}$  plasmids via uASPIre (**Fig. 5f**, **Suppl.**  
737 **Fig. 20**, **Methods**). We analysed the resulting NGS data comparing 44,289 common 5'-UTR-CDS  
738 variants above high-quality read count threshold that appeared in all six conditions (i.e. two strains with  
739 three plasmids). Specifically, we determined for each condition the effects of 5'-UTR position -1 by  
740 dividing the mean rTR of variants with a given base at this position by the mean rTR of all other variants  
741 (**Fig. 5g**). This analysis confirmed the strong, base-specific impact of this position, and, beyond that,  
742 revealed a significant dependence of the effect on the base present in position 37 of  $\text{tRNA}^{\text{fMet}}$ . To this  
743 end, we observed a strong increase in the rTR for variants whose base upstream of the start codon is  
744 complementary to position 37 of the overexpressed  $\text{tRNA}^{\text{fMet}}$  in both the WT and  $\Delta\text{metZ WV}$  strain. Non-  
745 complementarity consistently leads to a lower expression compared to the complementarity case  
746 across both strains and all  $\text{tRNA}^{\text{fMet}}$  variants (**Fig. 5h**). Similarly, a small yet significant positive impact  
747 on rTR is observable for the major wobble base pair G-U/U-G, which appears consistently for both  
748 directions of interaction (G in position 37 of  $\text{tRNA}^{\text{fMet}}$  with U in 5'-UTR position -1 and *vice versa*) and  
749 both strains (**Suppl Fig. 21**). Interestingly, a U at 5'-UTR position -1 leads to a small rTR-boosting effect  
750 also in presence of the non-complementary initiators  $\text{tRNA}^{\text{fMet-A37G}}$  and  $\text{tRNA}^{\text{fMet-A37U}}$  only in the WT  
751 strain (**Fig. 5g**). This can be explained by the presence of chromosomally encoded, endogenous  
752  $\text{tRNA}^{\text{fMet-A37}}$  copies, since this positive effect is neutralised or slightly inverted in the  $\Delta\text{metZ WV}$  strain.  
753 The effects at all other randomised positions in the mRNA were similar to the ones obtained for  $\text{Lib}_{\text{random}}$   
754 without overexpression of  $\text{tRNA}^{\text{fMet}}$  variants (**Suppl. Fig. 22** compare **Fig. 2b**).

755 These findings strongly suggest a direct base-pairing interaction of 5'-UTR position -1 with the  
756 nucleotide following the anticodon in  $\text{tRNA}^{\text{fMet}}$  (position 37), which leads to a significant positive effect  
757 on translation initiation upon successful hybridisation. Thus, our analyses confirm previous hypotheses  
758 to that end in a statistically solid manner based on more than 132,000 mRNA- $\text{tRNA}^{\text{fMet}}$  combinations,  
759 which were kinetically assessed in two different genetic backgrounds.



760

761

762

763

764

765

766

767

768

769

770

771

772

**Figure 5: Assessment of translational anomalies of arginine codon 2 and 5'-UTR position -1 in Lib<sub>random</sub>.** **a**) Effect of different synonymous codons in the second triplet of the CDS on rTR and mRNA folding energy (efeC). Black circles within violins are mean rTR/efeC values. \*\*\* denote p-values < 10<sup>-16</sup> in a Welch two sample t-test. **b**) Relationship between relative triplet frequency in *E. coli* and rTR for the four synonymous triplets in arginine codon 2. **c**) Effect of different bases in 5'-UTR position -1 on rTR and mRNA folding energy (efeC). Black circles within violins are mean rTR/efeC values. \*\*\* denote p-values < 10<sup>-16</sup> in a Welch two sample t-test. **d**) Plasmids for the overexpression of native initiator tRNA<sup>fMet</sup> and mutants thereof (**Suppl. Fig. 6, Methods**). Position 37 (3'-adjacent to the CAU anticodon) of tRNA<sup>fMet</sup> is mutated from A to C, G, or T/U. **e**) Growth of *E. coli* strains carrying plasmids for tRNA<sup>fMet</sup> overexpression in shake flask cultivations (LB, 37 °C). Bars are mean doubling times of independent biological triplicate cultivations with standard deviation as error bars. Dashed lines are the mean doubling time of the respective strain without tRNA overexpression (i.e. empty vector control) with



773 standard deviation as grey shaded areas. For tRNA<sup>fMet-A37C</sup>, doubling times were not determined (n.d.)  
774 due to severe growth inhibition (see main text). **f**) Approximately 50,000 variants of Lib<sub>random</sub> were tested  
775 in the presence of overexpressed tRNA<sup>fMet</sup> variants in *E. coli* strains containing (WT) and lacking  
776 ( $\Delta metZ WV$ ) the chromosomal *metZ WV* locus. **g**) Impact of tRNA<sup>fMet</sup> mutations on the rTR of variants  
777 from Lib<sub>random</sub>. Displayed effects are log<sub>2</sub>-transformed fold-changes (log<sub>2</sub> FC) of the average rTR of  
778 variants with a given base at 5'-UTR position -1 over the average rTR of variants with any other base  
779 at this position. Black arrows indicate complementarity between 5'-UTR position -1 and position 37 of  
780 the tRNA<sup>fMet</sup> variant. **h**) Impact of complementarity between 5'-UTR position -1 and tRNA<sup>fMet</sup> position 37.  
781 Circles are log<sub>2</sub>-transformed fold-changes (log<sub>2</sub> FC) of the average rTR of variants with  
782 complementarity or non-complementarity between mRNA and tRNA over the mean rTR of all variants  
783 in the same group (i.e. same tRNA<sup>fMet</sup> variant and strain). Bars are the mean log<sub>2</sub> FCs of the three  
784 tRNA<sup>fMet</sup> variants for each case and strain with standard deviation as error bars.

785

786

## 787 DISCUSSION

788 In this study, we systematically investigated the impact of 5'-UTR and N-terminal CDS on translation  
789 through mapping of more than 1.2 million mRNA sequence variants to their corresponding expression  
790 levels in *E. coli*. In combination with random and combinatorial library design, the ultrahigh throughput  
791 of our approach allowed us to critically assess sequence parameters known or supposed to influence  
792 translation efficiency. Furthermore, the generated large data basis enabled a precise quantification of  
793 effect sizes and correction for sequence-specific artefacts via statistically solid analyses.

794 To this end, we assessed mean effects of individual bases and positions in 5'-UTR and CDS along with  
795 various higher-order sequence parameters of the mRNA. We found that 25% of variance in our data  
796 could be explained by individual nucleotides and that GC-content, hybridisation with the 16S-rRNA and  
797 mRNA folding are the most significant determinants of translation confirming findings from previous  
798 studies (e.g. (8-12,20,27-30,32-37,39,40,83)). Using a simplistic machine learning approach, we  
799 compared the predictive potential of 248 parameters, which ranked 16S-rRNA hybridisation highest  
800 (20.9%) followed by various mRNA folding features (between 1.0% and 9.1%) and GC-content (5.97%)  
801 and pointed to a high degree of interaction and redundancy amongst parameters (**Fig. 5e**).

802 Furthermore, we found an unexpectedly large, base-specific contribution of two individual nucleotides,  
803 the negative impact of G in the third position of arginine codon 2 and the positive effect of U in position  
804 -1 of the 5'-UTR (**Fig. 2**). Follow-up analyses revealed that the former is not causally related to tRNA  
805 availability in the cell but can likely be attributed to a stronger tendency of variants with CGG as second  
806 codon to form mRNA secondary structures (**Fig. a, b**). Notably, mRNA accessibility at this position  
807 ranked amongst the most important features ( $accC_{1nt, pos. +6}$ ) of a predictive random forest model  
808 (**Fig. 3e**), which confirms the relation of the observed effect to mRNA folding. The positive effect of U  
809 directly upstream of the start codon, by contrast, was not linked to folding or any other mRNA parameter  
810 (**Fig. 5c**), which prompted further experiments to that end. Specifically, we assessed whether a base-  
811 pairing interaction of 5'-UTR position -1 with the base in 3' to the anticodon in initiator tRNA<sup>fMet</sup> (position  
812 A37) could be responsible for the effect. This hypothesis could be confirmed through compensatory

813 mutation of tRNA<sup>fMet</sup> position 37, which led to a translation-favouring effect in all cases of  
814 complementarity between tRNA and 5'-UTR (**Fig. 5g, h**). Several previous studies had shown that a U  
815 upstream of the start codon favours ribosome assembly and/or translation *in vitro* and *in vivo* (36,58,72-  
816 80). A link of these effects to the aforementioned base-pairing interaction, however, was only postulated  
817 and not experimentally confirmed in these studies. Esposito *et al.* (81) attempted to confirm the  
818 interaction in algal chloroplasts by substitution of position A37 in tRNA<sup>fMet</sup>. Notably, the variant  
819 tRNA<sup>fMet-A37C</sup> was not generated in their study, which showed severe growth inhibition and was thus  
820 excluded also in our work. For the tested reporter gene (*petA*), substitution of A37 indeed led to a  
821 translation-favouring effect only in cases of complementarity between tRNA<sup>fMet</sup> and the base upstream  
822 of the start codon. However, this observation was made on the basis of only three 5'-UTR position -1  
823 variants of *petA* carrying a non-native weak UAA start codon and could not be confirmed for several  
824 other analysed genes. Whether the impact for *petA* is specific to this gene (context) or the weak start  
825 codon, or indeed related to an interaction between tRNA<sup>fMet</sup> and the base upstream of the start codon  
826 therefore remains unclear. In this study, we assessed 45,258 mRNA sequences tested with three  
827 tRNA<sup>fMet</sup> variants and in two strains of normal and reduced endogenous expression of native tRNA<sup>fMet-A37</sup>.  
828 This did not only confirm the proposed quadruplet interaction in a statistically firm fashion but allowed  
829 to even quantify comparably subtle phenomena such as wobble base pairing (**Suppl. Fig. 21**), which  
830 can be masked for individual sequences and thus are inaccessible to low-throughput approaches.  
831 Lastly, we constructed and assessed more than a million combinatorial and full-factorial 5'-UTR-CDS  
832 combinations, which, in view of the high degree of interactivity, is indispensable to correctly assign  
833 observed effects to different mRNA parts and sequence parameters, and to precisely measure their  
834 contribution. This allowed us to quantify the mean individual contribution of the 5'-UTR and CDS to  
835 translational variance in a manner that would not be possible otherwise (e.g. using fully random  
836 libraries), which amount to 53% and 20%, respectively. Moreover, we capitalised on the combinatorial  
837 libraries and the large data basis to revise different hypotheses on the causal relationship between  
838 translation efficiency and codon usage. Similar to previous studies (e.g. (33,37)), our data confirmed a  
839 strong dependence of the rTR on the N-terminal CDS and a decreasing impact of codons with  
840 increasing distance to the start codon (**Figs. 2, 4b, c**). While this dependence unquestionably exists,  
841 the underlying mechanistic reasons remain less clear and were linked to both differences in mRNA  
842 folding and cellular tRNA abundance in the past. In our data, we found a small ( $R^2/p^2 < 0.7\%$ ) yet  
843 significant correlation of the rTR with codon usage metrics (**Fig. 4d**). However, the majority of the  
844 corresponding variance of the rTR can be explained by mRNA folding to an overwhelming degree while  
845 the contribution of codon usage metrics is extremely low (**Fig. 4f**). This low impact is further  
846 corroborated by the fact that none of the codon usage metrics was capable to increase the prediction  
847 accuracy of a random forest model, whereas mRNA folding had a very large impact (**Fig. 4g**). In  
848 summary, amongst the 1.2 million unique 5'-UTR-CDS combinations tested in this study the influence  
849 of different codons is virtually fully explainable by mRNA folding, whereas a causal connection to cellular  
850 tRNA abundance was either insignificant or negligibly small. The small apparent correlation between  
851 codon usage indices and rTR thus likely stems from differences in GC-content between rare and  
852 frequent codons, which leads to different tendencies to form secondary mRNA structures.

853 Consequently, our study highlights the importance of ultradeep sequence-function mapping for the  
854 accurate determination of the contribution of parts and phenomena involved in gene regulation. It  
855 should be mentioned that several other factors are known to influence translation (initiation), which  
856 have not been addressed in this study. These include the use of different start codons, long-range  
857 interactions between ribosome and 5'-UTR, and limitations of translation elongation (e.g. related to  
858 protein folding). Nonetheless, the presented methodology can be applied to scrutinise these additional  
859 factors, which, together with the results from this study, could serve as a basis to improve on  
860 inaccuracies of currently available models for the prediction and forward design of prokaryotic protein  
861 expression.

862

863

864

865

#### 866 **ACKNOWLEDGEMENT**

867 We thank Dr. Christian Beisel (ETH Zurich) for his support with NGS, and Dr. Nico Claassens and Dr.  
868 Thijs Nieuwkoop for critical reading of the manuscript.

869

#### 870 **AUTHOR CONTRIBUTIONS**

871 M.J. and S.H. conceived the study and planned experiments. S.H. performed experiments and  
872 computational works. S.H. and M.J. analysed data. M.J. coordinated the study. S.H. and M.J. wrote the  
873 manuscript.

874

#### 875 **FUNDING**

876 This work was supported by the European Commission [grant number 766975], and the Swiss National  
877 Science Foundation under the NCCR "Molecular Systems Engineering".

878

#### 879 **CONFLICT OF INTEREST**

880 The authors declare no conflict of interest.

881

882

883

884

885

886

887

888

889

890

891

892

## 893 REFERENCES

- 894 1. Wang, H.H., Isaacs, F.J., Carr, P.A., Sun, Z.Z., Xu, G., Forest, C.R. and Church, G.M. (2009)  
895 Programming cells by multiplex genome engineering and accelerated evolution. *Nature*, **460**,  
896 894-898.
- 897 2. Pullmann, P., Ulpinnis, C., Marillonnet, S., Gruetzner, R., Neumann, S. and Weissenborn, M.J.  
898 (2019) Golden Mutagenesis: An efficient multi-site-saturation mutagenesis approach by Golden  
899 Gate cloning with automated primer design. *Sci Rep*, **9**, 10932.
- 900 3. Xu, W., Klumbys, E., Ang, E.L. and Zhao, H. (2020) Emerging molecular biology tools and  
901 strategies for engineering natural product biosynthesis. *Metab Eng Commun*, **10**, e00108.
- 902 4. Lovmar, M. and Ehrenberg, M. (2006) Rate, accuracy and cost of ribosomes in bacterial cells.  
903 *Biochimie*, **88**, 951-961.
- 904 5. Vellanoweth, R.L. and Rabinowitz, J.C. (1992) The influence of ribosome-binding-site elements  
905 on translational efficiency in *Bacillus subtilis* and *Escherichia coli* in vivo. *Mol Microbiol*, **6**, 1105-  
906 1114.
- 907 6. Hersch, S.J., Elgamal, S., Katz, A., Ibba, M. and Navarre, W.W. (2014) Translation initiation  
908 rate determines the impact of ribosome stalling on bacterial protein synthesis. *J Biol Chem*, **289**,  
909 28160-28171.
- 910 7. Tietze, L. and Lale, R. (2021) Importance of the 5' regulatory region to bacterial synthetic  
911 biology applications. *Microb Biotechnol*, **14**, 2291-2315.
- 912 8. Jacob, W.F., Santer, M. and Dahlberg, A.E. (1987) A single base change in the Shine-Dalgarno  
913 region of 16S rRNA of *Escherichia coli* affects translation of many proteins. *Proc Natl Acad Sci U S A*, **84**, 4757-4761.
- 914 9. Dalboge, H., Carlsen, S., Jensen, E.B., Christensen, T. and Dahl, H.H. (1988) Expression of  
915 recombinant growth hormone in *Escherichia coli*: effect of the region between the Shine-  
916 Dalgarno sequence and the ATG initiation codon. *DNA*, **7**, 399-405.
- 917 10. Shine, J. and Dalgarno, L. (1974) The 3'-terminal sequence of *Escherichia coli* 16S ribosomal  
918 RNA: complementarity to nonsense triplets and ribosome binding sites. *Proc Natl Acad Sci U*  
919 *S A*, **71**, 1342-1346.
- 920 11. Steitz, J.A. and Jakes, K. (1975) How ribosomes select initiator regions in mRNA: base pair  
921 formation between the 3' terminus of 16S rRNA and the mRNA during initiation of protein  
922 synthesis in *Escherichia coli*. *Proc Natl Acad Sci U S A*, **72**, 4734-4738.
- 923 12. Li, G.W., Oh, E. and Weissman, J.S. (2012) The anti-Shine-Dalgarno sequence drives  
924 translational pausing and codon choice in bacteria. *Nature*, **484**, 538-541.
- 925 13. Chen, H., Bjercknes, M., Kumar, R. and Jay, E. (1994) Determination of the optimal aligned  
926 spacing between the Shine-Dalgarno sequence and the translation initiation codon of  
927 *Escherichia coli* mRNAs. *Nucleic Acids Res*, **22**, 4953-4957.
- 928 14. Salis, H.M., Mirsky, E.A. and Voigt, C.A. (2009) Automated design of synthetic ribosome  
929 binding sites to control protein expression. *Nat Biotechnol*, **27**, 946-950.
- 930 15. Mutalik, V.K., Guimaraes, J.C., Cambray, G., Lam, C., Christoffersen, M.J., Mai, Q.A., Tran,  
931 A.B., Paull, M., Keasling, J.D., Arkin, A.P. *et al.* (2013) Precise and reliable gene expression  
932 via standard transcription and translation initiation elements. *Nat Methods*, **10**, 354-360.
- 933 16. Osterman, I.A., Evfratov, S.A., Sergiev, P.V. and Dontsova, O.A. (2013) Comparison of mRNA  
934 features affecting translation initiation and reinitiation. *Nucleic Acids Res*, **41**, 474-486.
- 935 17. Espah Borujeni, A., Channarasappa, A.S. and Salis, H.M. (2014) Translation rate is controlled  
936 by coupled trade-offs between site accessibility, selective RNA unfolding and sliding at  
937 upstream standby sites. *Nucleic Acids Res*, **42**, 2646-2659.
- 938 18. Bonde, M.T., Pedersen, M., Klausen, M.S., Jensen, S.I., Wulff, T., Harrison, S., Nielsen, A.T.,  
939 Herrgard, M.J. and Sommer, M.O. (2016) Predictable tuning of protein expression in bacteria.  
940 *Nat Methods*, **13**, 233-236.
- 941 19. Hecht, A., Glasgow, J., Jaschke, P.R., Bawazer, L.A., Munson, M.S., Cochran, J.R., Endy, D.  
942 and Salit, M. (2017) Measurements of translation initiation from all 64 codons in *E. coli*. *Nucleic*  
943 *Acids Res*, **45**, 3615-3626.
- 944 20. Kuo, S.T., Jahn, R.L., Cheng, Y.J., Chen, Y.L., Lee, Y.J., Hollfelder, F., Wen, J.D. and Chou,  
945 H.D. (2020) Global fitness landscapes of the Shine-Dalgarno sequence. *Genome Res*, **30**, 711-  
946 723.
- 947 21. Komarova, E.S., Chervontseva, Z.S., Osterman, I.A., Evfratov, S.A., Rubtsova, M.P., Zatsepin,  
948 T.S., Semashko, T.A., Kostryukova, E.S., Bogdanov, A.A., Gelfand, M.S. *et al.* (2020) Influence  
949 of the spacer region between the Shine-Dalgarno box and the start codon for fine-tuning of the  
950 translation efficiency in *Escherichia coli*. *Microb Biotechnol*, **13**, 1254-1261.
- 951

- 952 22. Fargo, D.C., Zhang, M., Gillham, N.W. and Boynton, J.E. (1998) Shine-Dalgarno-like  
953 sequences are not required for translation of chloroplast mRNAs in *Chlamydomonas reinhardtii*  
954 chloroplasts or in *Escherichia coli*. *Mol Gen Genet*, **257**, 271-282.
- 955 23. Zheng, X., Hu, G.Q., She, Z.S. and Zhu, H. (2011) Leaderless genes in bacteria: clue to the  
956 evolution of translation initiation mechanisms in prokaryotes. *BMC Genomics*, **12**, 361.
- 957 24. Beck, H.J., Fleming, I.M. and Janssen, G.R. (2016) 5'-Terminal AUGs in *Escherichia coli*  
958 mRNAs with Shine-Dalgarno Sequences: Identification and Analysis of Their Roles in Non-  
959 Canonical Translation Initiation. *PLoS One*, **11**, e0160144.
- 960 25. Nakagawa, S., Niimura, Y. and Gojobori, T. (2017) Comparative genomic analysis of translation  
961 initiation mechanisms for genes lacking the Shine-Dalgarno sequence in prokaryotes. *Nucleic*  
962 *Acids Res*, **45**, 3922-3931.
- 963 26. Saito, K., Green, R. and Buskirk, A.R. (2020) Translational initiation in *E. coli* occurs at the  
964 correct sites genome-wide in the absence of mRNA-rRNA base-pairing. *Elife*, **9**.
- 965 27. de Smit, M.H. and van Duin, J. (1994) Control of translation by mRNA secondary structure in  
966 *Escherichia coli*. A quantitative analysis of literature data. *J Mol Biol*, **244**, 144-150.
- 967 28. Voges, D., Watzel, M., Nemetz, C., Wizemann, S. and Buchberger, B. (2004) Analyzing and  
968 enhancing mRNA translational efficiency in an *Escherichia coli* in vitro expression system.  
969 *Biochem Biophys Res Commun*, **318**, 601-614.
- 970 29. Kudla, G., Murray, A.W., Tollervey, D. and Plotkin, J.B. (2009) Coding-sequence determinants  
971 of gene expression in *Escherichia coli*. *Science*, **324**, 255-258.
- 972 30. Simonetti, A., Marzi, S., Jenner, L., Myasnikov, A., Romby, P., Yusupova, G., Klaholz, B.P. and  
973 Yusupov, M. (2009) A structural view of translation initiation in bacteria. *Cell Mol Life Sci*, **66**,  
974 423-436.
- 975 31. Na, D., Lee, S. and Lee, D. (2010) Mathematical modeling of translation initiation for the  
976 estimation of its efficiency to computationally design mRNA sequences with desired expression  
977 levels in prokaryotes. *BMC Syst Biol*, **4**, 71.
- 978 32. Seo, S.W., Yang, J.S., Kim, I., Yang, J., Min, B.E., Kim, S. and Jung, G.Y. (2013) Predictive  
979 design of mRNA translation initiation region to control prokaryotic translation efficiency. *Metab*  
980 *Eng*, **15**, 67-74.
- 981 33. Goodman, D.B., Church, G.M. and Kosuri, S. (2013) Causes and effects of N-terminal codon  
982 bias in bacterial genes. *Science*, **342**, 475-479.
- 983 34. Reeve, B., Hargest, T., Gilbert, C. and Ellis, T. (2014) Predicting translation initiation rates for  
984 designing synthetic biology. *Front Bioeng Biotechnol*, **2**, 1.
- 985 35. Espah Borujeni, A., Cetnar, D., Farasat, I., Smith, A., Lundgren, N. and Salis, H.M. (2017)  
986 Precise quantification of translation inhibition by mRNA structures that overlap with the  
987 ribosomal footprint in N-terminal coding sequences. *Nucleic Acids Res*, **45**, 5437-5448.
- 988 36. Yus, E., Yang, J.S., Sogues, A. and Serrano, L. (2017) A reporter system coupled with high-  
989 throughput sequencing unveils key bacterial transcription and translation determinants. *Nat*  
990 *Commun*, **8**, 368.
- 991 37. Cambray, G., Guimaraes, J.C. and Arkin, A.P. (2018) Evaluation of 244,000 synthetic  
992 sequences reveals design principles to optimize translation in *Escherichia coli*. *Nat Biotechnol*,  
993 **36**, 1005-1015.
- 994 38. Verma, M., Choi, J., Cottrell, K.A., Lavagnino, Z., Thomas, E.N., Pavlovic-Djuranovic, S.,  
995 Szczesny, P., Piston, D.W., Zaher, H.S., Puglisi, J.D. *et al.* (2019) A short translational ramp  
996 determines the efficiency of protein synthesis. *Nat Commun*, **10**, 5774.
- 997 39. Terai, G. and Asai, K. (2020) Improving the prediction accuracy of protein abundance in  
998 *Escherichia coli* using mRNA accessibility. *Nucleic Acids Res*, **48**, e81.
- 999 40. Cetnar, D.P. and Salis, H.M. (2021) Systematic Quantification of Sequence and Structural  
1000 Determinants Controlling mRNA stability in Bacterial Operons. *ACS Synth Biol*, **10**, 318-332.
- 1001 41. Sharp, P.M. and Li, W.H. (1987) The codon Adaptation Index--a measure of directional  
1002 synonymous codon usage bias, and its potential applications. *Nucleic Acids Res*, **15**, 1281-  
1003 1295.
- 1004 42. Eyre-Walker, A. and Bulmer, M. (1993) Reduced synonymous substitution rate at the start of  
1005 enterobacterial genes. *Nucleic Acids Res*, **21**, 4599-4603.
- 1006 43. Bentele, K., Saffert, P., Rauscher, R., Ignatova, Z. and Bluthgen, N. (2013) Efficient translation  
1007 initiation dictates codon usage at gene start. *Mol Syst Biol*, **9**, 675.
- 1008 44. Hanson, G. and Collier, J. (2018) Codon optimality, bias and usage in translation and mRNA  
1009 decay. *Nat Rev Mol Cell Biol*, **19**, 20-30.

- 1010 45. Ikemura, T. (1981) Correlation between the abundance of Escherichia coli transfer RNAs and  
1011 the occurrence of the respective codons in its protein genes. *Journal of Molecular Biology*, **146**,  
1012 1-21.
- 1013 46. Ikemura, T. (1985) Codon usage and tRNA content in unicellular and multicellular organisms.  
1014 *Mol Biol Evol*, **2**, 13-34.
- 1015 47. Dong, H., Nilsson, L. and Kurland, C.G. (1996) Co-variation of tRNA abundance and codon  
1016 usage in Escherichia coli at different growth rates. *J Mol Biol*, **260**, 649-663.
- 1017 48. Mitarai, N., Sneppen, K. and Pedersen, S. (2008) Ribosome collisions and translation efficiency:  
1018 optimization by codon usage and mRNA destabilization. *J Mol Biol*, **382**, 236-245.
- 1019 49. Zhang, G. and Ignatova, Z. (2009) Generic algorithm to predict the speed of translational  
1020 elongation: implications for protein biogenesis. *PLoS One*, **4**, e5036.
- 1021 50. Dobrzynski, M. and Bruggeman, F.J. (2009) Elongation dynamics shape bursty transcription  
1022 and translation. *Proc Natl Acad Sci U S A*, **106**, 2583-2588.
- 1023 51. Tuller, T., Carmi, A., Vestsigian, K., Navon, S., Dorfan, Y., Zaborske, J., Pan, T., Dahan, O.,  
1024 Furman, I. and Pilpel, Y. (2010) An evolutionarily conserved mechanism for controlling the  
1025 efficiency of protein translation. *Cell*, **141**, 344-354.
- 1026 52. Tuller, T., Waldman, Y.Y., Kupiec, M. and Ruppin, E. (2010) Translation efficiency is  
1027 determined by both codon bias and folding energy. *Proc Natl Acad Sci U S A*, **107**, 3645-3650.
- 1028 53. dos Reis, M., Wernisch, L. and Savva, R. (2003) Unexpected correlations between gene  
1029 expression and codon usage bias from microarray data for the whole Escherichia coli K-12  
1030 genome. *Nucleic Acids Res*, **31**, 6976-6985.
- 1031 54. Plotkin, J.B. and Kudla, G. (2011) Synonymous but not the same: the causes and  
1032 consequences of codon bias. *Nat Rev Genet*, **12**, 32-42.
- 1033 55. Quax, T.E., Claassens, N.J., Soll, D. and van der Oost, J. (2015) Codon Bias as a Means to  
1034 Fine-Tune Gene Expression. *Mol Cell*, **59**, 149-161.
- 1035 56. Nieuwkoop, T., Finger-Bou, M., van der Oost, J. and Claassens, N.J. (2020) The Ongoing  
1036 Quest to Crack the Genetic Code for Protein Production. *Mol Cell*, **80**, 193-209.
- 1037 57. Kosuri, S., Goodman, D.B., Cambray, G., Mutalik, V.K., Gao, Y., Arkin, A.P., Endy, D. and  
1038 Church, G.M. (2013) Composability of regulatory sequences controlling transcription and  
1039 translation in Escherichia coli. *Proc Natl Acad Sci U S A*, **110**, 14024-14029.
- 1040 58. Hollerer, S., Papaxanthos, L., Gumpinger, A.C., Fischer, K., Beisel, C., Borgwardt, K.,  
1041 Benenson, Y. and Jeschek, M. (2020) Large-scale DNA-based phenotypic recording and deep  
1042 learning enable highly accurate sequence-function mapping. *Nat Commun*, **11**, 3551.
- 1043 59. Claassens, N.J., Finger-Bou, M., Scholten, B., Muis, F., de Groot, J.J., de Gier, J.W., de Vos,  
1044 W.M. and van der Oost, J. (2019) Bicistronic Design-Based Continuous and High-Level  
1045 Membrane Protein Production in Escherichia coli. *ACS Synth Biol*, **8**, 1685-1690.
- 1046 60. Datsenko, K.A. and Wanner, B.L. (2000) One-step inactivation of chromosomal genes in  
1047 Escherichia coli K-12 using PCR products. *Proc Natl Acad Sci U S A*, **97**, 6640-6645.
- 1048 61. Sambrook, J.F. and Russel, D.W. (2001) *Molecular Cloning: A Laboratory Manual*. Cold Spring  
1049 Harbor Laboratory Press.
- 1050 62. Silva-Rocha, R., Martinez-Garcia, E., Calles, B., Chavarria, M., Arce-Rodriguez, A., de Las  
1051 Heras, A., Paez-Espino, A.D., Durante-Rodriguez, G., Kim, J., Nickel, P.I. et al. (2013) The  
1052 Standard European Vector Architecture (SEVA): a coherent platform for the analysis and  
1053 deployment of complex prokaryotic phenotypes. *Nucleic Acids Res*, **41**, D666-675.
- 1054 63. R Core Team. (2017). 4.0.3 ed. R Foundation for Statistical Computing, Vienna, Austria.
- 1055 64. Wickham, H. (2009), *Elegant Graphics for Data Analysis*. Springer-Verlag New York, pp. VIII,  
1056 213.
- 1057 65. Do, C.B., Woods, D.A. and Batzoglou, S. (2006) CONTRAfold: RNA secondary structure  
1058 prediction without physics-based models. *Bioinformatics*, **22**, e90-98.
- 1059 66. Mathews, D.H., Sabina, J., Zuker, M. and Turner, D.H. (1999) Expanded sequence  
1060 dependence of thermodynamic parameters improves prediction of RNA secondary structure. *J*  
1061 *Mol Biol*, **288**, 911-940.
- 1062 67. Lorenz, R., Bernhart, S.H., Honer Zu Siederdisen, C., Tafer, H., Flamm, C., Stadler, P.F. and  
1063 Hofacker, I.L. (2011) ViennaRNA Package 2.0. *Algorithms Mol Biol*, **6**, 26.
- 1064 68. Kiryu, H., Terai, G., Imamura, O., Yoneyama, H., Suzuki, K. and Asai, K. (2011) A detailed  
1065 investigation of accessibilities around target sites of siRNAs and miRNAs. *Bioinformatics*, **27**,  
1066 1788-1797.
- 1067 69. Egan, S.M. and Schleif, R.F. (1993) A regulatory cascade in the induction of rhaBAD. *J Mol*  
1068 *Biol*, **234**, 87-98.

- 1069 70. Jeschek, M., Gerngross, D. and Panke, S. (2016) Rationally reduced libraries for combinatorial  
1070 pathway optimization minimizing experimental effort. *Nat Commun*, **7**, 11163.
- 1071 71. Seong, B.L. and RajBhandary, U.L. (1987) Mutants of Escherichia coli formylmethionine tRNA:  
1072 a single base change enables initiator tRNA to act as an elongator in vitro. *Proc Natl Acad Sci*  
1073 *U S A*, **84**, 8859-8863.
- 1074 72. Ganoza, M.C., Fraser, A.R. and Neilson, T. (1978) Nucleotides contiguous to AUG affect  
1075 translational initiation. *Biochemistry*, **17**, 2769-2775.
- 1076 73. Eckhardt, H. and Luhrmann, R. (1981) Recognition by initiator transfer ribonucleic acid of a  
1077 uridine 5' adjacent to the AUG codon: different conformational states of formylatable  
1078 methionine-accepting transfer ribonucleic acid at the ribosomal peptidyl site. *Biochemistry*, **20**,  
1079 2075-2080.
- 1080 74. Ganoza, M.C., Sullivan, P., Cunningham, C., Hader, P., Kofoid, E.C. and Neilson, T. (1982)  
1081 Effect of bases contiguous to AUG on translation initiation. *J Biol Chem*, **257**, 8228-8232.
- 1082 75. Hui, A., Hayflick, J., Dinkelspiel, K. and de Boer, H.A. (1984) Mutagenesis of the three bases  
1083 preceding the start codon of the beta-galactosidase mRNA and its effect on translation in  
1084 Escherichia coli. *EMBO J*, **3**, 623-629.
- 1085 76. Ganoza, M.C., Marliere, P., Kofoid, E.C. and Louis, B.G. (1985) Initiator tRNA may recognize  
1086 more than the initiation codon in mRNA: a model for translational initiation. *Proc Natl Acad Sci*  
1087 *U S A*, **82**, 4587-4591.
- 1088 77. Gross, G., Mielke, C., Hollatz, I., Blocker, H. and Frank, R. (1990) RNA primary sequence or  
1089 secondary structure in the translational initiation region controls expression of two variant  
1090 interferon-beta genes in Escherichia coli. *J Biol Chem*, **265**, 17627-17636.
- 1091 78. Esposito, D., Hicks, A.J. and Stern, D.B. (2001) A role for initiation codon context in chloroplast  
1092 translation. *Plant Cell*, **13**, 2373-2384.
- 1093 79. Krishnan, K.M., Van Etten, W.J., 3rd and Janssen, G.R. (2010) Proximity of the start codon to  
1094 a leaderless mRNA's 5' terminus is a strong positive determinant of ribosome binding and  
1095 expression in Escherichia coli. *J Bacteriol*, **192**, 6482-6485.
- 1096 80. Krishnan, K.M. (2010).
- 1097 81. Esposito, D., Fey, J.P., Eberhard, S., Hicks, A.J. and Stern, D.B. (2003) In vivo evidence for  
1098 the prokaryotic model of extended codon-anticodon interaction in translation initiation. *EMBO*  
1099 *J*, **22**, 651-656.
- 1100 82. Barraud, P., Schmitt, E., Mechulam, Y., Dardel, F. and Tisne, C. (2008) A unique conformation  
1101 of the anticodon stem-loop is associated with the capacity of tRNA<sup>fMet</sup> to initiate protein  
1102 synthesis. *Nucleic Acids Res*, **36**, 4894-4901.
- 1103 83. Na, D. and Lee, D. (2010) RBSDesigner: software for designing synthetic ribosome binding  
1104 sites that yields a desired level of protein expression. *Bioinformatics*, **26**, 2633-2634.
- 1105