

1 **A Unique Ribosome Signature Reveals Bacterial Translation Initiation Sites**

2 Adam Giess<sup>1</sup>, Elvis Ndah<sup>2,3,4</sup>, Veronique Jonckheere<sup>2,3</sup>, Petra Van Damme<sup>\*2,3</sup>, Eivind Valen<sup>\*1,5</sup>

3

4 <sup>1</sup>Computational Biology Unit, Department of Informatics, University of Bergen, Bergen  
5 5020, Norway

6 <sup>2</sup>Medical Biotechnology Center, VIB, B-9000 Ghent, Belgium

7 <sup>3</sup>Department of Biochemistry, Ghent University, B-9000 Ghent, Belgium

8 <sup>4</sup>Lab of Bioinformatics and Computational Genomics, Department of Mathematical  
9 Modelling, Statistics and Bioinformatics, Faculty of Bioscience Engineering, Ghent  
10 University, B-9000 Ghent, Belgium

11 <sup>5</sup>Sars International Centre for Marine Molecular Biology, University of Bergen, 5008  
12 Bergen, Norway

13

14 \*Address correspondence to E. Valen, Email: [eivind.valen@gmail.com](mailto:eivind.valen@gmail.com), or P. Van Damme,  
15 Email: [petra.vandamme@vib-ugent.be](mailto:petra.vandamme@vib-ugent.be)

16

17

18 **KEYWORDS**

19

20 ribosome profiling, bacterial translation initiation, machine learning, N-terminal  
21 proteomics

22 **ABSTRACT**

23

24

25 While methods for annotation of genes are increasingly reliable the exact identification of  
26 the translation initiation site remains a challenging problem. Since the N-termini of  
27 proteins often contain regulatory and targeting information developing a robust method for  
28 start site identification is crucial. Ribosome profiling reads show distinct patterns of read  
29 length distributions around translation initiation sites. These patterns are typically lost in  
30 standard ribosome profiling analysis pipelines, when reads from footprints are adjusted to  
31 determine the specific codon being translated. Using these unique signatures we build a  
32 model capable of predicting translation initiation sites and demonstrate its high accuracy  
33 using N-terminal proteomics. Applying this to prokaryotic samples, we re-annotate  
34 translation initiation sites and provide evidence of N-terminal truncations and elongations  
35 of annotated coding sequences. These re-annotations are supported by the presence of  
36 Shine-Dalgarno sequences, structural and sequence based features and N-terminal  
37 peptides. Finally, our model identifies 61 novel genes previously undiscovered in the  
38 genome.

39 Identification of translated open reading frames (ORFs) is a critical step towards  
40 annotation of genes and the understanding of a genome. In addition to providing functional  
41 information via the peptide sequence, regulatory and targeting information are often  
42 contained within protein N-termini<sup>1,2</sup>, this makes accurate identification of the beginning of  
43 ORFs essential. Whole genome ORF identification in prokaryotes is most commonly  
44 performed *in silico*, using a variety of sequence features, such as GC codon bias and motifs  
45 such as the ribosomal binding site or Shine-Dalgarno sequence<sup>3,4,5</sup> in order to differentiate  
46 those ORFs that are thought to be functional from those that occur in the genome by  
47 chance. While these techniques are able to identify genomic regions containing ORFs with  
48 a high accuracy<sup>5</sup>, predicting translation initiation sites (TISs), and thus the exact beginning  
49 of a protein coding sequence (CDS), is substantially more challenging. This has led to the  
50 development of a number of *in silico* based TIS identification methods relying on a variety  
51 of sequence features<sup>6-9</sup>, which are typically post processing tools applied after initial ORF  
52 annotation in order to re-annotate the often erroneously predicted TIS.

53

54 High throughput proteogenomics has the potential to enable identification of protein N-  
55 termini, and by extension TISs, from an entire proteome. In practice however variation in  
56 protein expression levels, physical properties, MS-incompatibility and the occurrence of  
57 protein modifications limit the number of detectable protein N-termini<sup>10,11</sup>. In prokaryotes  
58 N-terminal proteomics typically captures the corresponding peptides of hundreds to the  
59 low thousands of genes<sup>11</sup>. For example, a recent study identified N-terminal peptides of 910  
60 of the 4140 (22%) annotated genes in *Escherichia coli*<sup>12</sup>. Although falling short of providing  
61 full genome annotation, such datasets provide an effective means of experimental TIS  
62 validation.

63

64 Significantly higher coverage of TISs can be achieved by using sequencing based  
65 technologies. By specifically focusing on ribosome protected fragments, ribosome  
66 profiling<sup>13</sup> (ribo-seq) infers which parts of the transcriptome are actively undergoing  
67 translation. In this way, ribo-seq has been used to demonstrate translation of many RNAs  
68 and regions that were not thought to be associated with ribosomes<sup>14-20</sup>. Being able to  
69 identify translation on a transcriptome-wide scale has obvious application to ORF  
70 annotation and a number of methodologies have been developed for prediction of  
71 translated ORFs<sup>17,19,21-23</sup>. These methods rely on a number of features, like codon periodicity,  
72 read context and read lengths, in order to distinguish footprints indicative of translation

73 from other, non-translating, footprints frequently observed in ribo-seq data. While extensive  
74 progress has been made on finding translated regions, delineating their exact boundaries  
75 has received less attention. Antibiotic treatment can be used to stall and capture footprints  
76 from the initiating ribosome<sup>14,24,25</sup>, but finding a suitable compound has been elusive in  
77 prokaryotes with only one dataset to date<sup>26</sup>.

78

79 Here we present a generally applicable method that does not depend on specialised  
80 chemical treatment, but can be take advantage of such data (Figure 1a). Using N-terminal  
81 proteomics we demonstrate its high accuracy and show that it is consistent with other  
82 features linked to translation initiation. Applying the model we predict numerous novel  
83 initiation sites in *Salmonella enterica* serovar Typhimurium.

## 84 **RESULTS**

85

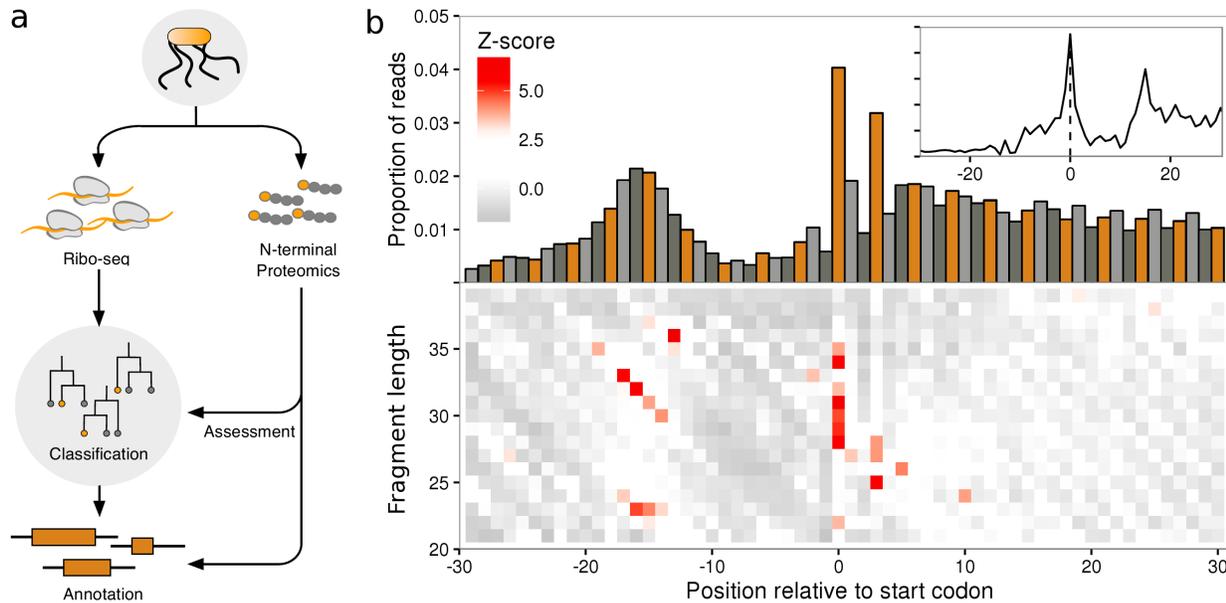
86

### 87 **Translation Initiation Sites Carry a Unique Signature**

88

89 To investigate whether ribo-seq could aid in the accurate delineation of translated ORFs we  
90 generated two ribo-seq libraries from monosome and polysome enriched fractions  
91 originating from *S. Typhimurium*. The similarities in the profiles of the two libraries  
92 (Supplementary Fig. S1e-f), taken with current literature reports of similarities in the  
93 translational properties of polysome and monosome fractions<sup>27</sup>, suggest that it is  
94 reasonable to consider these libraries sufficiently similar to serve as replicates for the  
95 purpose of initiation sites. The libraries were initially processed in a standard ribo-seq  
96 work-flow, where trimmed footprints were aligned to a reference genome, and then  
97 adjusted based on 5' read profiles to determine the specific codon under active translation  
98 (Figure 1b, inset, Supplementary Fig. S1e-f, inset). When exploring the processed reads we  
99 discovered that, consistent with previous reports<sup>26,28</sup>, annotated start sites of ribosomes  
100 treated with chloramphenicol carry a unique signature around the initiating codon (Figure  
101 1b, inset). Examining the unprocessed reads we observed that the pattern is a consequence  
102 of a specific distribution of fragment lengths (Figure 1b), information which is typically lost  
103 in pipelines that pre-process the read signal by adjusting reads (Figure 1b, inset). More  
104 specifically, heatmaps of 5' read profiles indicate that the pattern consists of an enrichment  
105 of longer fragments (30-35 nucleotides(nt)) starting 14-19 nt upstream of the initiation  
106 codon (a diagonal pattern), but ending at the same location, 15 nt downstream of the  
107 initiation codon. A shorter set of fragments (23-24 nt) are enriched in the same region, but  
108 have different end points, 7-9 nt downstream of the initiation codon. And finally, a strong  
109 enrichment of 5' ends of reads of length 28-35 nt can be observed exactly over the start  
110 codon itself (Figure 1b).

111



112 **Figure 1:** Translation initiation site classification with ribo-seq fragment length patterns.

113 (a) Schematic representation of the classification strategy. (b) Ribo-seq meta profiles in  
114 windows around start codons for all annotated CDSs in the *S. Typhimurium* genome  
115 (monosome sample, n=4187), contributions from each gene are scaled to a sum of one.

116 (upper) Proportion of 5' ribo-seq read counts per nucleotide position, coloured by codon  
117 position. (lower) heatmaps of z-scores of 5' ribo-seq read counts per fragment length.

118 (inset) proportions of ribo-seq read counts per nucleotide position, after adjusting reads by  
119 fragment length offsets (see methods).

## 120 **Ribosome profiling enables accurate annotation of translation initiation sites**

121

122 We trained a random forest model on TISs from the top 50% translated ORFs (see  
123 methods), to recognise the patterns in 5' ribo-seq read lengths and sequence contexts in a  
124 -20 to +10 nt window around start codons. In addition we encoded information about the  
125 start codon position within the ORF and the read abundance upstream and downstream of  
126 the start sites. The model was then used to predict TISs from all in-frame cognate and near  
127 cognate (one edit distance from ATG) start codons around annotated genes in the *S.*  
128 *Typhimurium* genome. Predictions on the two samples were highly accurate with area  
129 under curve (AUC) values of 0.9958 and 0.9956 on independent validation sets for the  
130 monosome and polysome sample, respectively (parameter importance for the models is  
131 summarised in Supplementary Tables S1-2). In total 4610 (monosome) and 4601 (polysome)  
132 TISs were predicted in the two sets. From these, we constructed a high confidence set from  
133 predictions common to both replicates. In total this set contained 4272 predictions,  
134 representing an 86.50% agreement between the replicates. The discrepancies  
135 predominantly originate from genes with scarce translation. Of the high confidence TISs,  
136 3853 matched annotated ORFs, 214 represented elongations and 205 truncations.  
137 Examples of predicted elongated, truncated and matching ORFs are shown in figure 2.

138

139 As expected the predictions show the same codon usage distribution (Supplementary Fig.  
140 S2), and carry the same read distribution signature as the annotated sites (Supplementary  
141 Fig. S2). Consistent with annotated initiation sites an increase in ribosome protected  
142 fragments can be seen downstream versus upstream of the predicted TIS (Figure 3a).  
143 Furthermore, elongated ORFs exhibit a shift in ribo-seq density downstream relative to the  
144 annotated TIS, consistent with the predicted elongation. Conversely, truncated ORFs  
145 exhibit a shift in read density upstream relative to the annotated TIS and consistent with  
146 the predicted truncation.

147

148 To further assess the predictions we compared the newly predicted TISs with the  
149 previously, potentially erroneously, annotated TIS. A highly significant sequence feature of  
150 translation initiation sites is the Shine-Dalgarno (SD) sequence which facilitate translation  
151 initiation in prokaryotes<sup>29</sup>. The consensus sequence GGAGG is located approximately 10 nt  
152 upstream of the start codon<sup>30</sup>. The predicted initiation sites show clear evidence of SD

153 sequences centred 9-10 nt upstream of the start codon (Figure 4a). Strikingly the  
154 annotated TISs, in these same genes where our model has predicted novel sites, show an  
155 absence of the SD sequence (Figure 4a). Since our model evaluates sequence context it is  
156 unsurprising that the predictions carry this signature, but the absence of these motifs  
157 around previously annotated start codons is notable.

158

159 Besides the presence of SD sequences, the guanine-cytosine (GC) content is commonly  
160 used to identify CDSs in prokaryotes. The overall GC content of a genome or genomic  
161 region is often highly optimised. In coding regions this optimisation can be achieved via  
162 synonymous substitutions, predominantly at third codon positions<sup>31</sup>, leading to a  
163 pronounced bias in the GC content of third nucleotide positions in coding regions  
164 compared to the rest of the genome. Interestingly at annotated sites, predicted elongations  
165 exhibit an increase in GC content upstream of the annotated start codon consistent with  
166 the location of the predicted site, conversely predicted truncations show a decrease  
167 downstream of the annotated start codon. In contrast, at predicted sites both predicted  
168 elongations and truncations fit closely to the expected distribution (Figure 4b upper).

169

170 Another significant feature of prokaryotic translation initiation is the absence of intrinsic  
171 structure in the region around the start codon enabling easier access for ribosomes to  
172 bind<sup>32</sup>. We therefore calculated the average free energy over all predicted sites and  
173 compared them to the previous annotation in the same genes. Consistent with GC content  
174 patterns, the annotated sites display a lower propensity to form secondary structure  
175 upstream of the start codon in elongated ORFs, and downstream of the start codon in  
176 truncated ORFs. In the predicted sites these less-structured regions can clearly be  
177 observed directly over the start codon, highly indicative of true initiation sites (Figure 4b  
178 lower).

179

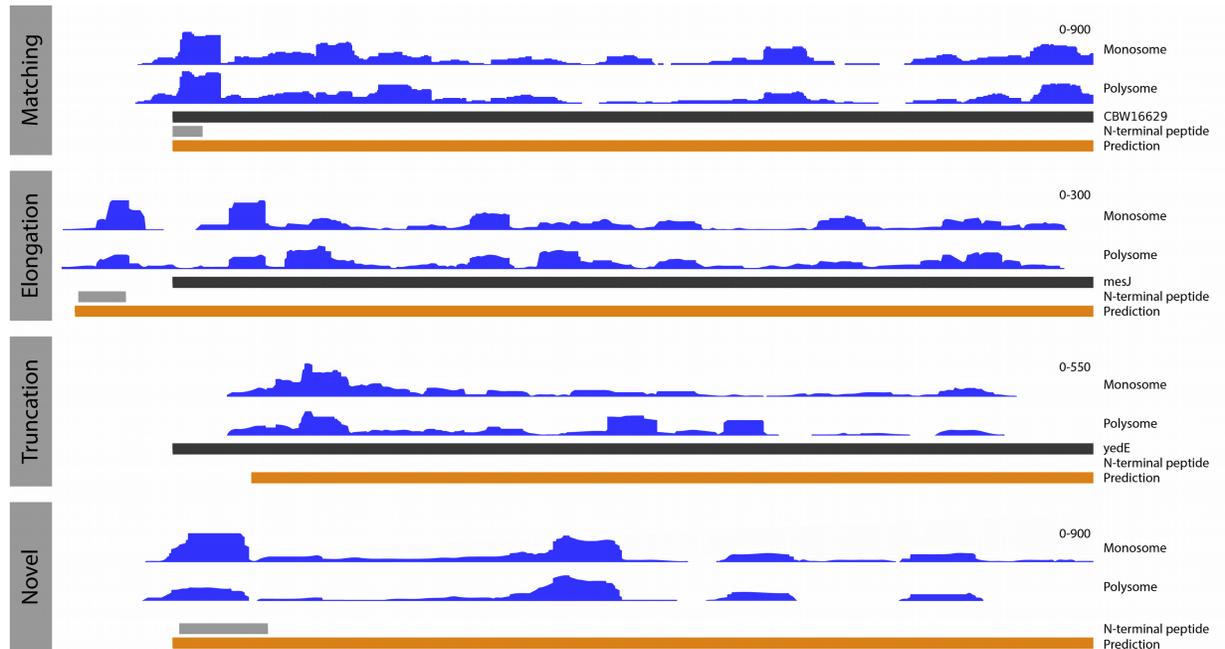
180 Ribosomes translocate along mRNAs three nucleotides at a time, corresponding to one  
181 codon and amino acid (aa). Consequently, reads originating from bona fide translated  
182 regions also exhibit a three nucleotide periodicity in adjusted read counts, with a bias  
183 towards mapping to the first nucleotide in each codon<sup>21</sup>. At initiation sites read distribution  
184 therefore switches from a random distribution upstream to a periodic, biased distribution  
185 downstream. Comparing the density of reads falling into each of the three codon positions,  
186 in elongated ORFs we observe increased read density at the first nucleotide position

187 upstream of annotated, but not predicted TISs. Similarly at truncated ORFs we see a  
188 decrease in the density of reads at the first nucleotide position downstream of the  
189 annotated TIS but not the predicted TIS (Figure 3c).

190

191 Taken together the patterns in read distribution, SD motifs, GC bias, unstructured regions  
192 and triplicate periodicity, provide clear and consistent support that the TISs which we re-  
193 annotate, show on average, a higher agreement with features indicative of canonically  
194 translated prokaryotic ORFs, than their corresponding previously annotated counterparts.

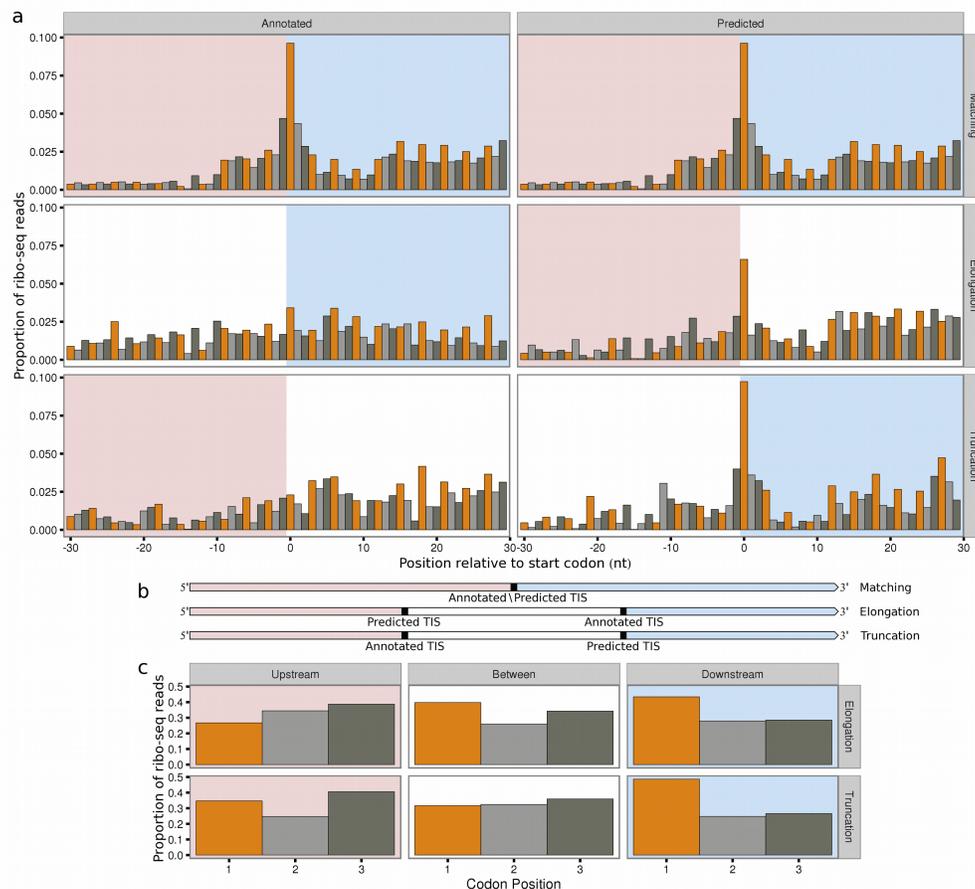
195



197 **Figure 2:** Examples of predicted translated ORFs.

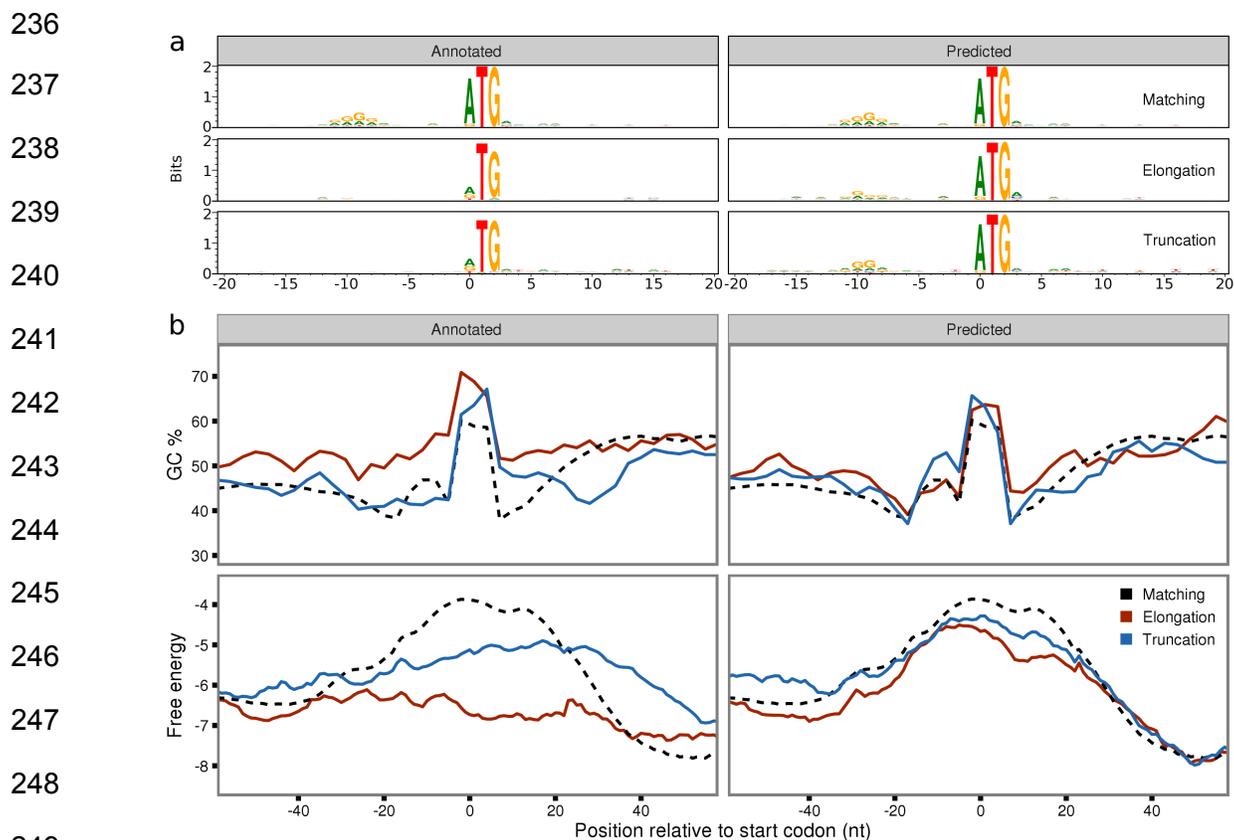
198 Showing genomic tracks of unadjusted ribo-seq read coverage in blue (y axis scale on the  
199 right hand side), annotated genes in black, predicted ORFs in orange and N-terminal  
200 peptides in grey. (**upper**) A predicted ORF in agreement with the annotated ORFs,  
201 supported by ribo-seq coverage and N-terminal evidence. (**middle upper**) A predicted  
202 elongation relative to the annotated ORF, with N-terminal evidence and ribo-seq coverage  
203 supporting the elongated prediction. (**middle lower**) a predicted truncation relative to the  
204 annotated ORF, with support from ribo-seq coverage. (**lower**) a novel predicted ORF,  
205 supported by N-terminal evidence and ribo-seq coverage.

206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220



221 **Figure 3:** Ribo-seq reads and periodicity are consistent with re-annotated translation  
222 initiation sites.

223 Bar colour indicates codon position. Downstream regions are highlighted in pink, upstream  
224 regions are highlighted in light blue. (a) Meta plots showing the proportion of scaled ribo-  
225 seq reads in relation to annotated or predicted translation initiation sites, for ORFs  
226 matching annotated genes (n=3853), predicted elongations (n=214) or predicted  
227 truncations (n=205). Contributions from each gene are scaled to a sum of one. Annotated  
228 TISs show increased ribo-seq density upstream (elongations), or downstream of start  
229 codons (truncations). (b) Transcript models. (c) Bar plots showing the sum of proportions  
230 of scaled ribo-seq read counts in each codon position. For truncations regions are 30 nt  
231 upstream of the annotated TIS, between the annotated and predicted TIS and 30 nt  
232 downstream of the predicted TIS. For elongations regions are 30 nt upstream of the  
233 predicted TIS, between the predicted and annotated TIS, and 30 nt downstream of  
234 annotated TIS. Three nt periodicity does not occur upstream of predicted TISs  
235 (truncations), but does occur upstream of annotated TISs (elongations).



**Figure 4:** Sequence and structure features support re-annotation of translation initiation sites.

(a) Sequence motifs relative to annotated or predicted translation initiation sites in the same genes. 'Matching' (n=3853) are identical, while predicted elongations (n=214) and truncations (n=205) have stronger SD sequences than their annotated counterparts. (b) Meta-profiles relative to annotated or predicted translation initiation sites, with lines representing ORFs matching annotated genes (dashed black), predicted elongations (red) and predicted truncations (blue). (upper) Mean GC content at third codon positions, averaged over 9nt sliding windows. Predicted TIS match the expected profile, showing an increase in GC content immediately after the start codon. Whereas predicted elongations and truncations show shifts down or upstream in annotated TIS. (lower) Meta-profiles of mean free energy averaged in 39 nt sliding windows. Peaks of low secondary structure potential, expected to occur over start codons, are centred over predicted TIS, but are clearly shifted down or upstream of annotated TIS, in predicted elongations and truncations.

## 266 **N-terminal proteomics confirms predicted sites**

267

268 In order to experimentally validate the accuracy of the predictions positional proteomics  
269 analyses enriching for protein N-termini were performed. Blocked N-termini were  
270 identified at 1040 *S. Typhimurium* ORFs, from which a high confidence subset of Nt-  
271 formylated Met-starting N-termini was selected (see methods) and used for assessing the  
272 accuracy of the model. In total 114 high confidence N-termini were identified supporting  
273 102 annotated CDSs, three N-terminal CDS elongations, and nine N-terminal CDS  
274 truncations. Because genomic positions with N-terminal peptide support were excluded  
275 from the set used to train the random forest model, these high confidence TIS positions can  
276 be used to determine the accuracy of the predictions. Of the 102 N-terminally supported  
277 annotated genes, 96 were predicted by the model. Furthermore two of the elongations, and  
278 four of the truncations were captured (Supplementary Table 3). Assuming that none of  
279 these genes have multiple initiation sites the sensitivity of the model can be estimated to be  
280 0.9444, the specificity to be 0.9991 and the positive predicted value to be 0.9444.

281

282 The remaining set of blocked N-termini supported a further 668 annotated CDSs, 50  
283 suggested CDS elongations and 310 suggested CDS truncations. In addition, we found  
284 peptides matching the predicted start positions from three distinct novel regions (defined  
285 as ORFs at least 300 nt in length, in regions that were not overlapping with annotated  
286 genes or regions at least 999 bp upstream of annotated genes). Comparing the predictions  
287 to the wider blocked N-termini set we find support for 648 predictions that match  
288 annotated TISs, 22 predicted elongations, 23 predicted truncated and three novel regions  
289 (Supplementary Table S4).

290

## 291 **Translation initiation sites are predicted at novel genomic regions**

292

293 In order to discover potential novel translated ORFs we applied our prediction models to  
294 look for TISs in genomic regions outside annotated ORFs. Novel ORFs that were similar in  
295 size to known CDSs (> 100aa) and with ribo-seq coverage along a high proportion of the  
296 ORF (>75 % coverage, see methods) were considered candidate translated novel ORFs. Of  
297 the 219 (monosome) and 193 (polysome) ORFs under consideration, 104 and 115 novel  
298 translated ORFs were predicted respectively. 61 of these novel translated ORFs were

299 common to both replicates (38.61% agreement) and used as a high confidence set of novel  
300 predictions. Unlike the annotated genes, these novel ORFs are not previously confirmed as  
301 translated regions and most had significantly lower read density (mean FPKM of 8) than  
302 annotated genes (mean FPKM of 126). The higher discrepancy between the two replicates  
303 is mainly a consequence of low-abundance start sites that did not pass the threshold in  
304 either of the replicates.

305

306 Read density plots over the novel ORFs revealed features consistent with protein coding  
307 regions, but with higher variance due to the low number of ORFs. Specifically, GC content  
308 increases downstream of the initiation codon, the regions around the initiation codon have  
309 less intrinsic structure potential and Shine-Dalgarno sequences are present upstream  
310 (Supplementary Fig. S3). Additionally, three of the predicted novel translated ORFs were  
311 supported by N-terminally enriched peptide evidence (a representative example is shown in  
312 Figure 2). A further 22 showed high similarity to known protein sequences, four of which  
313 contained functional protein domains (Supplementary Table S5).

314

### 315 **Tetracycline treated samples improve classifier accuracy**

316

317 While reads isolated from elongating ribosomes provide sufficient information to predict  
318 the majority of translation start sites we set out to explore the full potential of our classifier  
319 in combination with publicly available data from initiating ribosomes. A recent study on *E.*  
320 *coli*<sup>12</sup> demonstrated the use of tetracycline as a translation inhibitor to enrich for footprints  
321 from initiating ribosomes in prokaryotes. The tetracycline datasets show the pattern that  
322 we expect to see from initiating ribosomes as a range of read lengths starting 28-14nt  
323 upstream of the initiation codon (5' data), but ending at the same positions 14-15nt  
324 downstream of the initiation codon (3' data). An additional pattern of shorter fragment  
325 lengths can also be observed starting 26-18nt upstream, and ending 2nt downstream of the  
326 initiation codon (Supplementary Fig. S1.a,b,g,h).

327

328 We trained separate classifiers on chloramphenicol (elongating) and tetracycline (initiating)  
329 libraries from this dataset, using two replicates for each of the conditions (Supplementary  
330 Table S6). Model performance was evaluated with receiver operating characteristic (ROC)  
331 curves on the validation datasets for each replicate, the resulting AUC values of 0.9993 and

332 0.9994 in the tetracycline replicates were higher than those of chloramphenicol samples  
333 (0.9992 and 0.9983). The parameter importance in each of the models is shown in  
334 Supplementary Tables S7-10. The chloramphenicol models predicted a total of 3111 ORFs,  
335 including 57 elongations and 53 truncations (Supplementary Table S11). In the tetracycline  
336 dataset a total of 3711 ORFs were predicted, with 86 elongations and 79 truncations  
337 (Supplementary Table S12).

338

339 *E. coli* predictions were assessed against the ecogene curated set of 923 experimentally  
340 verified protein starts<sup>33</sup>. Genes within this dataset were excluded from the sets that were  
341 used to train the random forest models, in order to provide a means of assessing the  
342 accuracy of the ORF predictions. Five of the verified protein starts correspond to  
343 pseudogenes without annotated CDSs, of the remaining 917 verified protein starts, 821  
344 (89.53%) matched ORFs in the tetracycline predictions, with 24 (2.62%) predicted ORFs in  
345 disagreement with the curated set (11 elongations, 13 truncations). In the chloramphenicol  
346 predictions 760 (82.88%) were found to match ecogene start sites, and 27 (2.94%) were  
347 found to be inconsistent (13 elongations, 14 truncations) with the verified protein starts  
348 (assuming genes do not have multiples TIS) (Supplementary Tables S11-12). Based on the  
349 experimentally verified starts the tetracycline-based classifier resulted in higher accuracy  
350 (sensitivity 0.9194, specificity 0.9996, positive predictive value 0.9716) than the  
351 chloramphenicol-based classifier (sensitivity 0.8539, specificity 0.9996, positive predictive  
352 value 0.9657). Surprisingly, the difference was not major arguing that using initiating  
353 ribosomes is not a requirement to obtain a good annotation of initiation sites.

354

355

356

## 357 **DISCUSSION**

358

359

360 Our model shows that the distribution of ribo-seq footprint lengths can be used in  
361 conjunction with sequence features to accurately determine the translation initiation  
362 landscape of prokaryotes. These patterns are typically disrupted in standard ribo-seq  
363 analysis when reads of different fragment lengths are adjusted and merged to determine  
364 the specific codon under translation. The model is applicable across multiple organisms  
365 and experimental conditions and can be augmented with data from initiating ribosomes. It  
366 exhibits high accuracy as assessed by cross-validation, N-terminal proteomics and  
367 independent sequence-based metrics such as potential to form RNA structures.  
368 Interestingly, the predicted TISs exhibit known features of translation initiation which the  
369 previous annotations do not. In *S. Typhimurium*, our model provides evidence for 61 novel  
370 translated ORFs and the re-annotation of 419 genes. In particular, the current annotation  
371 includes 19 genes that lacked initiation codons, of which we were able to re-annotate 15  
372 (Supplementary Table S4).

373

374 As expected, models based on initiating reads perform better than models based on  
375 elongating ribo-seq reads, suggesting that an optimal strategy for TIS identification would  
376 favour the use of the more focused, initiating ribo-seq profiles. However the degree of  
377 improvement between the models was relatively small, confirming the suitability of both  
378 elongating and initiating ribo-seq libraries for the purposes of TIS and ORF detection.

379

380 While the mechanistic or experimental origin of the patterns that our model captures  
381 remain unexplored, it is interesting to note the importance the models place  
382 (Supplementary Tables S1-2, S9-10) on the shorter range of fragments of 23-25 nt (*S.*  
383 *Typhimurium*) or 21-26 nt (*E. coli* tetracycline) in length. These shorter reads are  
384 consistent with recent reports of ribosomal subunits in a variety of distinct configurations,  
385 observed from translation complex profiling in the eukaryote *Saccharomyces cerevisiae*<sup>34</sup>.  
386 Whether similar patterns of read length distributions can be observed in eukaryotic ribo-  
387 seq datasets remains to be determined, although the method that we describe in this  
388 article is, regardless, fully extendable to eukaryotic datasets.

389

390 In conclusion, this study demonstrates the utility of ribo-seq fragment length patterns for  
391 TIS identification across multiple experimental conditions. These models provide a  
392 significant step forward in experimental TIS discovery, facilitating the move towards  
393 complete ORF annotation in both presumably well-annotated model organisms, as well as  
394 the ever growing list of newly sequenced genomes.

395 **ONLINE METHODS**

396

397

398 **Preparation of ribo-seq libraries**

399

400 Overnight stationary cultures of wild type *S. Typhimurium* (*Salmonella enterica* serovar  
401 Typhimurium - strain SL1344) grown in LB media at 37 °C with agitation (200 rpm) were  
402 diluted at 1:200 in LB and grown until they reached and OD600 of 0.5 (i.e., logarithmic  
403 (Log) phase grown cells). Bacterial cells were pre-treated for 5 min with chloramphenicol  
404 (Sigma Aldrich) at a final concentration of 100 µg/ml before collection by centrifugation  
405 (6000 × g, 5 min) at 4 °C. Collected cells were flash frozen in liquid nitrogen. The frozen  
406 pellet of a 50 ml culture was re-suspended and thawed in 1 ml ice-cold lysis buffer for  
407 polysome isolation (10 mM MgCl<sub>2</sub>, 100 mM NH<sub>4</sub>Cl, 20 mM Tris.HCl pH 8.0, 20 U/ml of  
408 RNase-free DNase I (NEB 2 U/µl), 1mM chloramphenicol (or 300 µg/ml), 20 µl/ml lysozyme  
409 (50mg/ml in water) and 100 u/ml SUPERase.In™ RNase Inhibitor (Thermo Fisher  
410 Scientific, Bremen, Germany)), vortexed and left on ice for 2 min with periodical agitation.  
411 Subsequently, the samples were subjected to mechanical disruption by two repetitive cycles  
412 of freeze-thawing in liquid nitrogen, added 5 mM CaCl<sub>2</sub>, 30µl 10% DOC and 1 × complete  
413 and EDTA-free protease inhibitor cocktail (Roche, Basel, Switzerland) and left on ice for 5  
414 min. Lysates were clarified by centrifugation at 16,000 x g for 10 min at 4 °C.

415

416 For the monosome sample, the supernatant was subjected to MNase (Roche diagnostics  
417 Belgium) digestion using 600 U MNase (about~ 1000 U per mg of protein). Digestion of  
418 polysomes proceeded for 1 h at 25 °C with gentle agitation at 400 rpm and the reaction  
419 was stopped by the addition of 10 mM EGTA. Next, monosomes were recovered by  
420 ultracentrifugation over a 1 M sucrose cushion in polysome isolation buffer without RNase-  
421 free DNase I and lysozyme, and with 2 mM DTT added using a TLA-120.2 rotor for 4 hr at  
422 75,000 rpm and 4 °C.

423

424 For the selective purification of monosomes from polysomes (polysome sample), the  
425 supernatant was resolved on 10-55% (w/v) sucrose gradients by centrifugation using an  
426 SW41 rotor at 35,000 rpm for 2.5 hr at 4 °C. The sedimentation profiles were recorded at  
427 260 nm and the gradient fractionated using a BioComp Gradient Master (BioComp)

428 according to the manufacturer's instructions. Polysome-enriched fractions were pooled and  
429 subjected to MNase digestion and monosome recovery as described above.

430

431 Ribosome-protected mRNA footprints with sizes ranging from 26-34 nt were selected and  
432 processed as described previously<sup>14</sup> with some minor adjustments as previously described<sup>35</sup>.  
433 The resulting ribo-seq cDNA libraries of the monosome and polysome sample were  
434 duplexed and sequenced on a NextSeq 500 instrument (Illumina) to yield 75 bp single-end  
435 reads.

436

### 437 **Ribo-seq data processing**

438

439 Ribo-seq data were preprocessed with cutadapt<sup>36</sup> to remove sequencing adaptors,  
440 discarding reads less than 20 nt in length after trimming. Trimmed reads were initially  
441 aligned to the SILVA RNA database version 119<sup>37</sup>, the remaining reads were then mapped  
442 to either *Salmonella enterica* serovar Typhimurium - strain SL1344 (Assembly:  
443 GCA\_000210855.2) or *Escherichia coli* str. K-12 substr. MG1655 (Assembly:  
444 GCA\_000005845.2). Alignments were performed with bowtie2<sup>38</sup>. Reads were brought to  
445 codon resolution by adjusting the 5' position of each read by a fixed distance offset, specific  
446 to each fragment length, based on visual identification of periodicity meta plots of the read  
447 counts per fragment length (Supplementary Figs. S5-6). In the *S. Typhimurium* dataset the  
448 following fragment lengths were selected and adjusted by the values in brackets, in the  
449 monosome sample 29 (13 nt), 30 (14 nt), 31 (15 nt), 32 (16 nt), 33 (17 nt), and in the  
450 polysome sample 29 (13 nt), 30 (14 nt), 31 (15 nt), 32 (16 nt), 33 (17 nt) and 34 (18 nt).  
451 Selected reads of the indicated lengths account for 39.98 and 48.69 % of total reads for the  
452 monosome and polysome samples, respectively.

453

454 Recent publications reporting prokaryotic ribo-seq<sup>26,28,39,40</sup> suggest that read fragments from  
455 libraries digested with micrococcal nuclease align more precisely to their 3' rather than 5'  
456 ends. Consistent with this, we observe a modest increase in the periodicity of meta profiles  
457 of the *S. Typhimurium* ribo-seq libraries when reads are brought to codon resolution from  
458 the 3' end (Supplementary Fig. S1), however this does not hold true for the *E. coli* datasets,  
459 where the use of 3' poly adenosine adaptors, results in a loss of resolution at the 3' end  
460 after read trimming (Supplementary Fig. S1), making the use of 5' ends preferable.

461 Regardless, the protected read fragment patterns that we use in the input feature vectors  
462 for the classifier takes both length and position into consideration. Consequently the  
463 classifier is unaffected by this choice. However, to maintain consistency throughout this  
464 study read counting for model predictors was performed from the 5' end for all libraries.

465

## 466 **Read distributions and heat maps**

467

468 Ribo-seq read distributions were summarised over all annotated start codons in the *S.*  
469 *Typhimurium* and *E. coli* annotations respectively. 5' read counts were taken from regions  
470 30nt upstream to 60nt downstream of the start codon, 3' read counts were taken from the  
471 first nucleotide of the start codon up to 90 nucleotides downstream. All reads with a MAPQ  
472 greater than 10, from the upper 90% of genes by total CDS expression were included. Total  
473 counts were scaled to a sum of one per individual region, in order to not disproportionately  
474 favour profiles from highly expressed genes. Meta plots were then produced to show the  
475 proportion of read counts over the window across all genes. 3' and 5' heatmaps were  
476 generated from the scaled regions, showing the number of standard deviations from the  
477 row (fragment length) mean.

478

## 479 **Model implementation**

480

481 For each candidate TIS a feature vector was defined as each nucleotide in a -20 to +10nt  
482 window around the position, the ribo-seq 5' FPKM (fragments per kilobase per million  
483 mapped reads) between the current position and the next in-frame downstream stop codon,  
484 the count of potential in-frame start sites (codons within one edit distance of ATG) from the  
485 nearest in-frame upstream stop codon to the current position, the proportion of 5' ribo-seq  
486 reads upstream in a 20 nt window, the proportion of 5' ribo-seq reads downstream in a 20  
487 nt window, the ratio of 5' ribo-seq reads up and downstream and the proportion 5' ribo-seq  
488 counts per fragment length for a fixed range of positions in relation to current site  
489 (selected from visual inspection of 5' fragment length heatmaps (Supplementary Fig. S4)).  
490 In the *S. Typhimurium* samples fragment lengths 20-35 nt in positions -20 to -11 and 0 nt,  
491 were used. In the *E. coli* datasets for the tetracycline samples fragment lengths 20-35 nt at  
492 positions -25 to -16 nt, were selected and for chloramphenicol lengths 30 to 50 nt, at  
493 positions -25 to -16 and -1 to +1nt were used.

494 Stop-to-stop windows were defined for each annotated gene as all in-frame positions  
495 between the nearest in-frame upstream stop codon and the stop codon of the gene (with a  
496 maximum length cut-off 999 nt upstream).

497

498 The H2o random forest implementation<sup>41</sup> was used and the models were trained with  
499 positive examples of randomly selected annotated start codons from the upper 50% of  
500 genes ranked by ribo-seq expression over the gene CDS. We additionally required that the  
501 positive examples were not among the genes supported by N-terminal peptides in the *S.*  
502 *Typhimurium* samples or included in the ecogene dataset for the *E. coli* samples, since  
503 these were retained for model accuracy assessment. Negative examples were randomly  
504 selected from in-frame codons in the stop-to-stop windows both upstream and downstream  
505 of the annotated TIS. The *S. Typhimurium* models were trained on 1500 positive and 6000  
506 negative positions, with an independent validation set of 200 positive and 800 negative  
507 positions, the validation set was used for parameter training (number of trees monosome:  
508 600, polysome: 600). The *E. coli* models were trained on 1100 positive and 4400 negative  
509 positions, with an independent validation set of 200 positive and 800 negative positions for  
510 parameter tuning (number of trees: CM1: 950, CM2: 950, TET2: 650 and TET3: 700).  
511 Predictions were then run against all cognate and near cognate (defined as one edit  
512 distance from ATG) in-frame positions, in the stop-to-stop regions. Novel predictions were  
513 performed against all cognate and near cognate codons in stop-to-stop regions around  
514 ORFs of at least 300nt in length, with a ribo-seq read coverage of 0.75 or more (ORF  
515 coverage was defined as the proportion of nucleotides in each predicted ORF that at least  
516 one ribo-seq read mapped to), that did not overlap with annotated exons. ORFs were  
517 delineated by extending each candidate TIS to the closest in-frame stop codon. For a given  
518 stop-to-stop region the model selected the TIS with the highest positive predicted score per  
519 sample. Predictions from the replicates for each of the datasets were then compared,  
520 discarding predictions that were unique to only one replicate.

521

## 522 **N-terminal proteomics**

523

524 Overnight stationary cultures of wild type *S. Typhimurium* (*Salmonella enterica* serovar  
525 *Typhimurium* - strain SL1344) grown in LB media at 37 °C with agitation (200 rpm) were  
526 diluted at 1:200 in LB and grown until they reached and OD600 of 0.5 (i.e., logarithmic  
527 (Log) phase grown cells). Bacterial cells were collected by centrifugation (6000 × g, 5 min)

528 at 4 °C, flash frozen in liquid nitrogen and cryogenically pulverized using a liquid nitrogen  
529 cooled pestle and mortar. The frozen pellet of a 50 ml culture was re-suspended and  
530 thawed in 1 ml ice-cold lysis buffer (50 mM NH<sub>4</sub>HCO<sub>3</sub> (pH 7.9) supplemented with a  
531 complete protease inhibitor cocktail tablet (Roche Diagnostics GmbH, Mannheim,  
532 Germany) and subjected to mechanical disruption by two repetitive freeze-thaw and  
533 sonication cycles (i.e. 2 minutes of sonication on ice for 20-s bursts at output level 4 with a  
534 40% duty cycle (Branson Sonifier 250; Ultrasonic Converter)). The lysate was cleared by  
535 centrifugation for 15 min at 16,000 × *g* and the protein concentration measured using the  
536 protein assay kit (Bio-Rad) according to the manufacturer's instructions. The lysate was  
537 added Gu.HCl (4M f.c.) and subjected to N-terminal COFRADIC analysis as described  
538 previously<sup>42</sup>. Free amines were blocked at the protein level making use of an N-  
539 hydroxysuccinimide ester of (stable isotopic encoded) acetate (i.e. NHS esters of <sup>13</sup>C<sub>2</sub>D<sub>3</sub>  
540 acetate), which allows distinguishing in vivo and in vitro blocked N-terminal peptides<sup>43</sup>. The  
541 modified protein sample was digested overnight with sequencing-grade modified trypsin  
542 (1/100 (w/w trypsin /substrate)) at 37 °C and subsequent steps of the N-terminal  
543 COFRADIC procedure were performed as previously described<sup>42</sup>.

544

#### 545 **LC-MS/MS analysis**

546

547 LC-MS/MS analysis was performed using an Ultimate 3000 RSLC nano HPLC (Dionex,  
548 Amsterdam, the Netherlands) in-line connected to an LTQ Orbitrap Velos mass  
549 spectrometer (Thermo Fisher Scientific, Bremen, Germany). The sample mixture was  
550 loaded on a trapping column (made in-house, 100 µm I.D. × 20 mm, 5 µm beads C18  
551 Reprosil-HD, Dr. Maisch). After back flushing from the trapping column, the sample was  
552 loaded on a reverse-phase column (made in-house, 75 µm I.D. × 150 mm, 5 µm beads C18  
553 Reprosil-HD, Dr. Maisch). Peptides were loaded in solvent A' (0.1% trifluoroacetic acid, 2%  
554 acetonitrile (ACN)) and separated with a linear gradient from 2% solvent A' (0.1% formic  
555 acid) to 50% solvent B' (0.1% formic acid and 80% ACN) at a flow rate of 300 nl/min  
556 followed by a wash reaching 100% solvent B'. The mass spectrometer was operated in  
557 data-dependent mode, automatically switching between MS and MS/MS acquisition for the  
558 ten most abundant peaks in a given MS spectrum. Full scan MS spectra were acquired in  
559 the Orbitrap at a target value of 1E6 with a resolution of 60,000. The 10 most intense ions  
560 were then isolated for fragmentation in the linear ion trap, with a dynamic exclusion of 20  
561 s. Peptides were fragmented after filling the ion trap at a target value of 1E4 ion counts.

562 Mascot Generic Files were created from the MS/MS data in each LC run using the Mascot  
563 Distiller software (version 2.5.1.0, Matrix Science, [www.matrixscience.com/Distiller.html](http://www.matrixscience.com/Distiller.html)).  
564 To generate these MS/MS peak lists, grouping of spectra was allowed with a maximum  
565 intermediate retention time of 30 s and a maximum intermediate scan count of 5. Grouping  
566 was done with a 0.005 Da precursor tolerance. A peak list was only generated when the  
567 MS/MS spectrum contained more than 10 peaks. There was no de-isotoping and the  
568 relative signal-to-noise limit was set at 2.

569

570 The generated MS/MS peak lists were searched with Mascot using the Mascot Daemon  
571 interface (version 2.5.1, Matrix Science). Searches were performed using a 6-FT database  
572 of the *S. Typhimurium* genome combined with the Ensembl protein sequence database  
573 (assembly AMS21085v2 version 86.1), which totalled 139,408 entries after removal of  
574 redundant sequences. The 6-FT database was generated by traversing the entire genome  
575 across the six reading frames and searching for all NTG (N=A,T,C,G) start codons and  
576 extending each to the nearest in frame stop codon (TAA,TGA,TAG), discarding ORFs less  
577 than 21nt in length. The Mascot search parameters were set as follows: Heavy acetylation  
578 at lysine side-chains (Acetyl:2H(3)C13(2) (K)), carbamidomethylation of cysteine and  
579 methionine oxidation to methionine-sulfoxide were set as fixed modifications. Variable  
580 modifications were formylation, acetylation and heavy acetylation of N-termini  
581 (Acetyl:2H(3)C13(2) (N-term)) and pyroglutamate formation of N-terminal glutamine (both  
582 at peptide level). Endoproteinase semi-Arg-C/P (semi Arg-C specificity with Arg-Pro  
583 cleavage allowed) was set as enzyme allowing for no missed cleavages. Mass tolerance was  
584 set to 10 ppm on the precursor ion and to 0.5 Da on fragment ions. Peptide charge was set  
585 to 1+, 2+, 3+ and instrument setting was put to ESI-TRAP. Only peptides that were ranked  
586 one, have a minimum amino acid length of seven, scored above the threshold score, set at  
587 95% confidence, and belonged to the category of *in vivo* or *in vitro* blocked N-terminal  
588 peptides compliant with the rules of initiator methionine (iMet) processing<sup>44</sup> were withheld.  
589 More specifically, iMet processing was considered in the case of iMet-starting N-termini  
590 followed by any of the following amino acids; Ala, Cys, Gly, Pro, Ser, Thr, Met or Val and  
591 only if the iMet was encoded by ATG or any of the following near-cognate start codons;  
592 GTG and TTG (Supplementary Table S13). In contrast to eukaryotic nascent protein N-  
593 termini, the typical lack of significant steady-state levels of N-terminal protein modification  
594 (e.g. Nt-acetylation or Nt-formylation), warrant caution to unequivocally assign bacterial  
595 protein N-termini as proxies of translation initiation. As such, a high confidence subset of  
596 Nt-formylated Met-starting N-termini, was selected (Supplementary Table S14).

## 597 **Assessing model accuracy**

598

599 GC content was calculated at the third nucleotide positions for all annotated and predicted  
600 ORFs and mean GC values were summarised for each subgroup of predicted ORFs  
601 (matching annotations, truncations and elongations) in 9 nt sliding windows, over regions  
602 57 nt upstream and 57 nt downstream of the annotated or predicted start sites.

603

604 -20 to + 20 nt nucleotide sequences were extracted around the predicted and annotated  
605 TIS in the *S. Typhimurium* and *E. coli* genomes. Sequence logos were generated for each  
606 subgroup of matching annotations, truncations, elongations and novel genes, using the  
607 weblogo tool<sup>45</sup>.

608

609 The minimum free energy of RNA secondary structure around predicted and annotated  
610 ORFs was estimated with RNAfold version 2.1.9 from the ViennaRNA package<sup>46</sup>. Mean free  
611 energy values were summarised for each ORF class in 39 nt sliding window across regions  
612 57 nt up and downstream of the start codon.

613

614 Read distributions were created for each subgroup of predicted ORFs (matching  
615 annotations, truncations, elongations, and novel genes) and their corresponding annotated  
616 TIS. Distributions of ribo-seq reads adjusted to codon level resolution, were summarised in  
617 regions 30 nt upstream and downstream of the first nucleotide of the initiation codon, total  
618 counts of each individual region were scaled to a sum of one, in order to normalise profiles  
619 for differences in gene expression levels. Meta plots were then produced to show the  
620 proportion of reads over the window position from all predicted subgroups and their  
621 corresponding annotated start codons.

622

623 Sensitivity, specificity and positive predictive values were calculated from all genes that  
624 were supported by either high confidence n-terminal peptides (*S. Typhimurium*) or  
625 experimentally verified protein starts (*E. coli*). Supported predicted ORFs were considered  
626 true positives, predicted ORFs that disagreed with supported positions were classified as  
627 false positives. False negatives were assigned from supported genes where no ORF was  
628 predicted. All in-frame cognate and near cognate start codons (one edit distance from

629 ATG), in CDS regions of supported genes that were neither predicted nor supported were  
630 considered true negatives.

631

632 Amino acid sequences of novel ORF were compared to known proteins in the nonredundant  
633 protein database (Update date:2016/12/15) and protein domains (cdd.v.3.15) using  
634 BLASTP<sup>47</sup> (version 2.5.1+). Hits with the greatest coverage of query sequence and lowest e-  
635 value were selected. Hits were considered highly similar if they shared >95% identity to a  
636 protein sequence, over 100% of the novel ORF sequence

637

#### 638 **DATA AVAILABILITY**

639

640 The previously published *E. coli* ribo-seq dataset was downloaded from the NCBI SRA  
641 (BioProject ID:PRJDB2960). *S. Typhimurium* ribo-seq sequencing data has been deposited  
642 in NCBI's Gene Expression Omnibus<sup>48</sup> and is accessible through GEO Series accession  
643 number GSE91066. *S. Typhimurium* mass spectrometry proteomics data have been  
644 deposited to the ProteomeXchange Consortium via the PRIDE<sup>49</sup> partner repository with the  
645 dataset identifier PXD005579 and 10.6019/PXD005579.

646

#### 647 **ACKNOWLEDGMENTS**

648

649 We thank Gunnar Schulze and Kornel Labun for valuable discussions. Prof. Kris Gevaert for  
650 financial support of this research (Research Foundation - Flanders (FWO-Vlaanderen),  
651 project number G.0440.10). P.V.D. acknowledges support from the Research Foundation -  
652 Flanders (FWO-Vlaanderen), project number G.0269.13N. A.G and E.V. acknowledge  
653 support from the Bergen Research Foundation.

654

#### 655 **AUTHOR CONTRIBUTIONS**

656

657 A.G., E.V. and P.V.D. conceived the study and wrote the manuscript; A.G. performed the  
658 computational analysis; P.V.D performed the proteomics experiment. E.N. and P.V.D.  
659 performed proteomics analysis; P.V.D. and V.J. prepared the ribo-seq libraries.

660 **COMPETING FINANCIAL INTERESTS**

661

662 The authors declare that they have no competing financial interests.

663 **REFERENCES**

664

665

- 666 1. Hall, J., Hazlewood, G. P., Surani, M. A., Hirst, B. H. & Gilbert, H. J. Eukaryotic and  
667 prokaryotic signal peptides direct secretion of a bacterial endoglucanase by mammalian  
668 cells. *J. Biol. Chem.* 265, 19996–19999 (1990).
- 669 2. Kozak, M. Initiation of translation in prokaryotes and eukaryotes. *Gene* 234, 187–208  
670 (1999).
- 671 3. Delcher, a L., Harmon, D., Kasif, S., White, O. & Salzberg, S. L. Improved microbial gene  
672 identification with GLIMMER. *Nucleic Acids Res.* 27, 4636–4641 (1999).
- 673 4. Brocchieri, L., Kledal, T. N., Karlin, S. & Mocarski, E. S. Predicting coding potential from  
674 genome sequence: application to betaherpesviruses infecting rats and mice. *J. Virol.* 79,  
675 7570–96 (2005).
- 676 5. Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site  
677 identification. *BMC Bioinformatics* 11, 119 (2010).
- 678 6. Suzek, B. E., Ermolaeva, M. D., Schreiber, M. & Salzberg, S. L. A probabilistic method  
679 for identifying start codons in bacterial genomes. *Bioinformatics* 17, 1123–1130 (2001).
- 680 7. Zhu, H. Q., Hu, G. Q., Ouyang, Z. Q., Wang, J. & She, Z. S. Accuracy improvement for  
681 identifying translation initiation sites in microbial genomes. *Bioinformatics* 20, 3308–3317  
682 (2004).
- 683 8. Ou, H. Y., Guo, F. B. & Zhang, C. T. GS-Finder: A program to find bacterial gene start  
684 sites with a self-training method. *Int. J. Biochem. Cell Biol.* 36, 535–544 (2004).
- 685 9. Tech, M., Morgenstern, B. & Meinicke, P. TICO: A tool for postprocessing the predictions  
686 of prokaryotic translation initiation sites. *Nucleic Acids Res.* 34, 588–590 (2006).
- 687 10. Hartmann, E. M. & Armengaud, J. N-terminomics and proteogenomics, getting off to a  
688 good start. *Proteomics* 14, 2637–2646 (2014).
- 689 11. Berry, I. J., Steele, J. R., Padula, M. P. & Djordjevic, S. P. The application of terminomics  
690 for the identification of protein start sites and proteoforms in bacteria. *Proteomics* 16, 257–  
691 272 (2016).
- 692 12. Nakahigashi, K. et al. Comprehensive identification of translation start sites by  
693 tetracycline-inhibited ribosome profiling. *DNA Res.* 23, 193–201 (2016).
- 694 13. Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S. & Weissman, J. S. Genome-wide  
695 analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science*  
696 324, 218–23 (2009).

- 697 14. Ingolia, N. T., Lareau, L. F. & Weissman, J. S. Ribosome profiling of mouse embryonic  
698 stem cells reveals the complexity and dynamics of mammalian proteomes. *Cell* 147, 789–  
699 802 (2011).
- 700 15. Brar, G. A. et al. High-Resolution View of the Yeast Meiotic Program Revealed by  
701 Ribosome Profiling. *Science* (80-. ). 335, 552–557 (2012).
- 702 16. Michel, A. M. et al. Observation of dually decoded regions of the human genome using  
703 ribosome profiling data. *Genome Res.* 22, 2219–2229 (2012).
- 704 17. Chew, G. L. et al. Ribosome profiling reveals resemblance between long non-coding  
705 RNAs and 5' leaders of coding RNAs. *Development* 140, 2828–34 (2013).
- 706 18. Crappé, J. et al. Combining in silico prediction and ribosome profiling in a genome-wide  
707 search for novel putatively coding sORFs. *BMC Genomics* 14, 648 (2013).
- 708 19. Bazzini, A. A. et al. Identification of small ORFs in vertebrates using ribosome  
709 footprinting and evolutionary conservation. *EMBO J.* 33, 981–993 (2014).
- 710 20. Pauli, A. et al. Toddler: an embryonic signal that promotes cell movement via Apelin  
711 receptors. *Science* 343, 1248636 (2014).
- 712 21. Calviello, L. et al. Detecting actively translated open reading frames in ribosome  
713 profiling data. *Nat. Methods* 13, 1–9 (2015).
- 714 22. Duncan, C. D. S. & Mata, J. The translational landscape of fission-yeast meiosis and  
715 sporulation. *Nat. Struct. Mol. Biol.* 21, 641–7 (2014).
- 716 23. Ingolia, N. T. et al. Ribosome Profiling Reveals Pervasive Translation Outside of  
717 Annotated Protein-Coding Genes. *Cell Rep.* 8, 1365–1379 (2014).
- 718 24. Fritsch, C. et al. Genome-wide search for novel human uORFs and N-terminal protein  
719 extensions using ribosomal footprinting. *Genome Res.* 22, 2208–2218 (2012).
- 720 25. Lee, S. et al. Global mapping of translation initiation sites in mammalian cells at single-  
721 nucleotide resolution. *Proc. Natl. Acad. Sci. U. S. A.* 109, E2424–E2432 (2012).
- 722 26. Nakahigashi, K. et al. Effect of codon adaptation on codon-level and gene-level  
723 translation efficiency in vivo. *BMC Genomics* 15, 1115 (2014).
- 724 27. Heyer, E. E. & Moore, M. J. Redefining the Translational Status of 80S Monosomes. *Cell*  
725 164, 757–769 (2016).
- 726 28. Woolstenhulme, C. J., Guydosh, N. R., Green, R. & Buskirk, A. R. High-Precision analysis  
727 of translational pausing by ribosome profiling in bacteria lacking EFP. *Cell Rep.* 11, 13–21  
728 (2015).
- 729 29. Shine, J. & Dalgarno, L. Determinant of cistron specificity in bacterial ribosomes.  
730 *Nature* 254, 34–38 (1975).

- 731 30. Nakagawa, S., Niimura, Y., Miura, K. & Gojobori, T. Dynamic evolution of translation  
732 initiation mechanisms in prokaryotes. *Proc. Natl. Acad. Sci.* 107, 6382–6387 (2010).
- 733 31. Muto, A. & Osawa, S. The guanine and cytosine content of genomic DNA and bacterial  
734 evolution. *Proc. Natl. Acad. Sci. U. S. A.* 84, 166–9 (1987).
- 735 32. Del Campo, C., Bartholomäus, A., Fedyunin, I. & Ignatova, Z. Secondary Structure  
736 across the Bacterial Transcriptome Reveals Versatile Roles in mRNA Regulation and  
737 Function. *PLoS Genet.* 11, 1–23 (2015).
- 738 33. Zhou, J. & Rudd, K. E. EcoGene 3.0. *Nucleic Acids Res.* 41, 613–624 (2013).
- 739 34. Archer, S. K., Shirokikh, N. E., Beilharz, T. H. & Preiss, T. Dynamics of ribosome  
740 scanning and recycling revealed by translation complex profiling. *Nature* 535, 570–4  
741 (2016).
- 742 35. Gawron, D., Ndah, E., Gevaert, K. & Damme, P. Van. Positional proteomics reveals  
743 differences in N-terminal proteoform stability. *Mol. Syst. Biol.* 12, 858 (2016).
- 744 36. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing  
745 reads. *EMBnet.journal* 17, pp. 10–12 (2011).
- 746 37. Quast, C. et al. The SILVA ribosomal RNA gene database project: Improved data  
747 processing and web-based tools. *Nucleic Acids Res.* 41, 590–596 (2013).
- 748 38. Langmead, B. & Salzberg, S. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*  
749 9, 357–359 (2012).
- 750 39. Balakrishnan, R., Oman, K., Shoji, S., Bundschuh, R. & Fredrick, K. The conserved  
751 GTPase LepA contributes mainly to translation initiation in *Escherichia coli*. *Nucleic Acids*  
752 *Res.* 42, 13370–13383 (2014).
- 753 40. Mohammad, F., Woolstenhulme, C. J., Green, R. & Buskirk, A. R. Clarifying the  
754 Translational Pausing Landscape in Bacteria by Ribosome Profiling. *Cell Rep.* 14, 686–694  
755 (2016).
- 756 41. Aiello, S., Eckstrand, E., Fu, A., Landry, M. & Aboyoun, P. Machine Learning with R and  
757 H2O. (2016). at <<http://h2o.ai/resources>.>
- 758 42. Staes, A. et al. Selecting protein N-terminal peptides by combined fractional diagonal  
759 chromatography. *Nat. Protoc.* 6, 1130–1141 (2011).
- 760 43. Van Damme, P. et al. A review of COFRADIC techniques targeting protein N-terminal  
761 acetylation. *BMC Proc.* 3 Suppl 6, S6 (2009).
- 762 42. Frottin, F. et al. The Proteomics of N-terminal Methionine Cleavage. *Mol. Cell.*  
763 *Proteomics* 5, 2336–2349 (2006).
- 764 45. Crooks, G., Hon, G., Chandonia, J. & Brenner, S. WebLogo: a sequence logo generator.  
765 *Genome Res* 14, 1188–1190 (2004).

- 766 46. Lorenz, R. et al. ViennaRNA Package 2.0. *Algorithms Mol. Biol.* 6, 26 (2011).
- 767 47. Altschul, S. F. et al. Gapped BLAST and PSI-BLAST: a new generation of protein  
768 database search programs. *Nucleic Acids Res* 25, 3389–3402 (1997).
- 769 48. Edgar, R., Domrachev, M. & Lash, A. E. Gene Expression Omnibus: NCBI gene  
770 expression and hybridization array data repository. *Nucleic Acids Res* 30, 207–210 (2002).
- 771 49. Vizcaino, J. A. et al. 2016 update of the PRIDE database and its related tools. *Nucleic*  
772 *Acids Res.* 44, D447–D456 (2016).