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

2 Adam Giess[1], Elvis Ndah[2,3,4], Veronique Jonckheere[2,3], Petra Van Damme*[2,3], Eivind Valen*[1,5]

3

4 [1]Computational Biology Unit, Department of Informatics, University of Bergen, Bergen

5 5020, Norway

6 [2]Medical Biotechnology Center, VIB, B-9000 Ghent, Belgium

7 [3]Department of Biochemistry, Ghent University, B-9000 Ghent, Belgium

8 [4]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 [5]Sars International Centre for Marine Molecular Biology, University of Bergen, 5008

12 Bergen, Norway

13

14 **KEYWORDS**

15

16 ribosome profiling, bacterial translation initiation, machine learning, N-terminal

17 proteomics

18  **ABSTRACT**

19

20

21  While methods for annotation of genes are increasingly reliable the exact identification of

22  the translation initiation site remains a challenging problem. Since the N-termini of

23  proteins often contain regulatory and targeting information developing a robust method for

24  start site identification is crucial. Ribosome profiling reads show distinct patterns of read

25  length distributions around translation initiation sites. These patterns are typically lost in

26  standard ribosome profiling analysis pipelines, when reads from footprints are adjusted to

27  determine the specific codon being translated. Using these unique signatures we build a

28  model capable of predicting translation initiation sites and demonstrate its high accuracy

29  using N-terminal proteomics. Applying this to prokaryotic samples, we re-annotate

30  translation initiation sites and provide evidence of N-terminal truncations and elongations

31  of annotated coding sequences. These re-annotations are supported by the presence of

32  Shine-Dalgarno sequences, structural and sequence based features and N-terminal

33  peptides. Finally, our model identifies 61 novel genes previously undiscovered in the

34  genome.

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

49

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

59

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

69   from other, non-translating, footprints frequently observed in ribo-seq data. While extensive

70   progress has been made on finding translated regions, delineating their exact boundaries

71   has received less attention. Antibiotic treatment can be used to stall and capture footprints

72   from the initiating ribosome[14,24,25], but finding a suitable compound has been elusive in

73   prokaryotes with only one dataset to date[26].

74

75   Here we present a generally applicable method that does not depend on specialised

76   chemical treatment, but can be take advantage of such data (Figure 1a). Using N-terminal

77   proteomics we demonstrate its high accuracy and show that it is consistent with other

78   features linked to translation initiation. Applying the model we predict numerous novel

79   initiation sites in *Salmonella enterica* serovar Typhimurium.

80    **RESULTS**

81

82

83    **Translation Initiation Sites Carry a Unique Signature**

84

85    To investigate whether ribo-seq could aid in the accurate delineation of translated ORFs we

86    generated two ribo-seq libraries from monosome and polysome enriched fractions

87    originating from *S.* Typhimurium. The similarities in the profiles of the two libraries

88    (Supplementary Fig. S1e-f), taken with current literature reports of similarities in the

89    translational properties of polysome and monosome fractions[27], suggest that it is

90    reasonable to consider these libraries sufficiently similar to serve as replicates for the

91    purpose of initiation sites. The libraries were initially processed in a standard ribo-seq

92    work-flow, where trimmed footprints were aligned to a reference genome, and then

93    adjusted based on 5' read profiles to determine the specific codon under active translation

94    (Figure 1b, inset, Supplementary Fig. S1e-f, inset). When exploring the processed reads we

95    discovered that, consistent with previous reports[26,28], annotated start sites of ribosomes

96    treated with chloramphenicol carry a unique signature around the initiating codon (Figure

97    1b, inset). Examining the unprocessed reads we observed that the pattern is a consequence

98    of a specific distribution of fragment lengths (Figure 1b), information which is typically lost

99    in pipelines that pre-process the read signal by adjusting reads (Figure 1b, inset). More

100   specifically, heatmaps of 5' read profiles indicate that the pattern consists of an enrichment

101   of longer fragments (30-35 nucleotides(nt)) starting 14-19 nt upstream of the initiation

102   codon (a diagonal pattern), but ending at the same location, 15 nt downstream of the

103   initiation codon. A shorter set of fragments (23-24 nt) are enriched in the same region, but

104   have different end points, 7-9 nt downstream of the initiation codon. And finally, a strong

105   enrichment of 5' ends of reads of length 28-35 nt can be observed exactly over the start
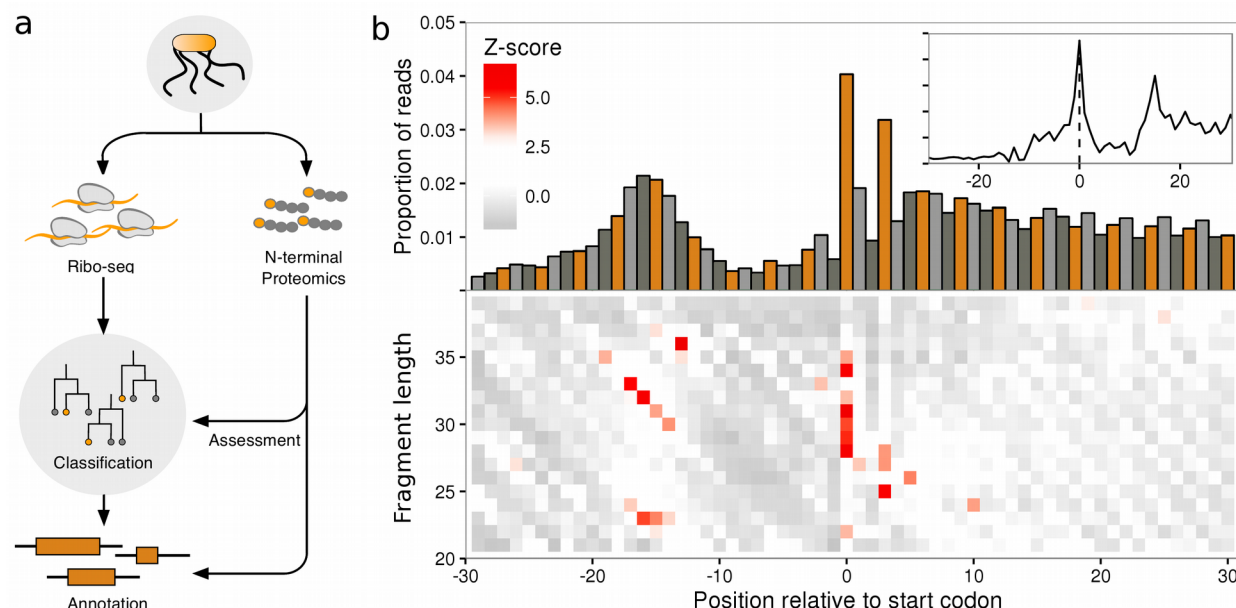
106   codon itself (Figure 1b).

107

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

109 (**a**) Schematic representation of the classification strategy. (**b**) Ribo-seq meta profiles in

110 windows around start codons for all annotated CDSs in the *S.* Typhimurium genome

111 (monosome sample, n=4187), contributions from each gene are scaled to a sum of one.

112 (**upper**) Proportion of 5' ribo-seq read counts per nucleotide position, coloured by codon

113 position. (**lower**) heatmaps of z-scores of 5' ribo-seq read counts per fragment length.

114 (**inset**) proportions of ribo-seq read counts per nucleotide position, after adjusting reads by

115 fragment length offsets (see methods).

116 **Ribosome profiling enables accurate annotation of translation initiation sites**

117

118 We trained a random forest model on TISs from the top 50% translated ORFs (see

119 methods), to recognise the patterns in 5' ribo-seq read lengths and sequence contexts in a

120 -20 to +10 nt window around start codons. In addition we encoded information about the

121 start codon position within the ORF and the read abundance upstream and downstream of

122 the start sites. The model was then used to predict TISs from all in-frame cognate and near

123 cognate (one edit distance from ATG) start codons around annotated genes in the *S.*

124 Typhimurium genome. Predictions on the two samples were highly accurate with area

125 under curve (AUC) values of 0.9958 and 0.9956 on independent validation sets for the

126 monosome and polysome sample, respectively (parameter importance for the models is

127 summarised in Supplementary Tables S1-2). In total 4610 (monosome) and 4601 (polysome)

128 TISs were predicted in the two sets. From these, we constructed a high confidence set from

129 predictions common to both replicates. In total this set contained 4272 predictions,

130 representing an 86.50% agreement between the replicates. The discrepancies

131 predominantly originate from genes with scarce translation. Of the high confidence TISs,

132 3853 matched annotated ORFs, 214 represented elongations and 205 truncations.

133 Examples of predicted elongated, truncated and matching ORFs are shown in figure 2.

134

135 As expected the predictions show the same codon usage distribution (Supplementary Fig.

136 S2), and carry the same read distribution signature as the annotated sites (Supplementary

137 Fig. S2). Consistent with annotated initiation sites an increase in ribosome protected

138 fragments can be seen downstream versus upstream of the predicted TIS (Figure 3a).

139 Furthermore, elongated ORFs exhibit a shift in ribo-seq density downstream relative to the

140 annotated TIS, consistent with the predicted elongation. Conversely, truncated ORFs

141 exhibit a shift in read density upstream relative to the annotated TIS and consistent with

142 the predicted truncation.

143

144 To further assess the predictions we compared the newly predicted TISs with the

145 previously, potentially erroneously, annotated TIS. A highly significant sequence feature of

146 translation initiation sites is the Shine-Dalgarno (SD) sequence which facilitate translation

147 initiation in prokaryotes[29]. The consensus sequence GGAGG is located approximately 10 nt

148 upstream of the start codon[30]. The predicted initiation sites show clear evidence of SD

149  sequences centred 9-10 nt upstream of the start codon (Figure 4a). Strikingly the

150  annotated TISs, in these same genes where our model has predicted novel sites, show an

151  absence of the SD sequence (Figure 4a). Since our model evaluates sequence context it is

152  unsurprising that the predictions carry this signature, but the absence of these motifs

153  around previously annotated start codons is notable.

154

155  Besides the presence of SD sequences, the guanine-cytosine (GC) content is commonly

156  used to identify CDSs in prokaryotes. The overall GC content of a genome or genomic

157  region is often highly optimised. In coding regions this optimisation can be achieved via

158  synonymous substitutions, predominantly at third codon positions[31], leading to a

159  pronounced bias in the GC content of third nucleotide positions in coding regions

160  compared to the rest of the genome. Interestingly at annotated sites, predicted elongations

161  exhibit an increase in GC content upstream of the annotated start codon consistent with

162  the location of the predicted site, conversely predicted truncations show a decrease

163  downstream of the annotated start codon. In contrast, at predicted sites both predicted

164  elongations and truncations fit closely to the expected distribution (Figure 4b upper).
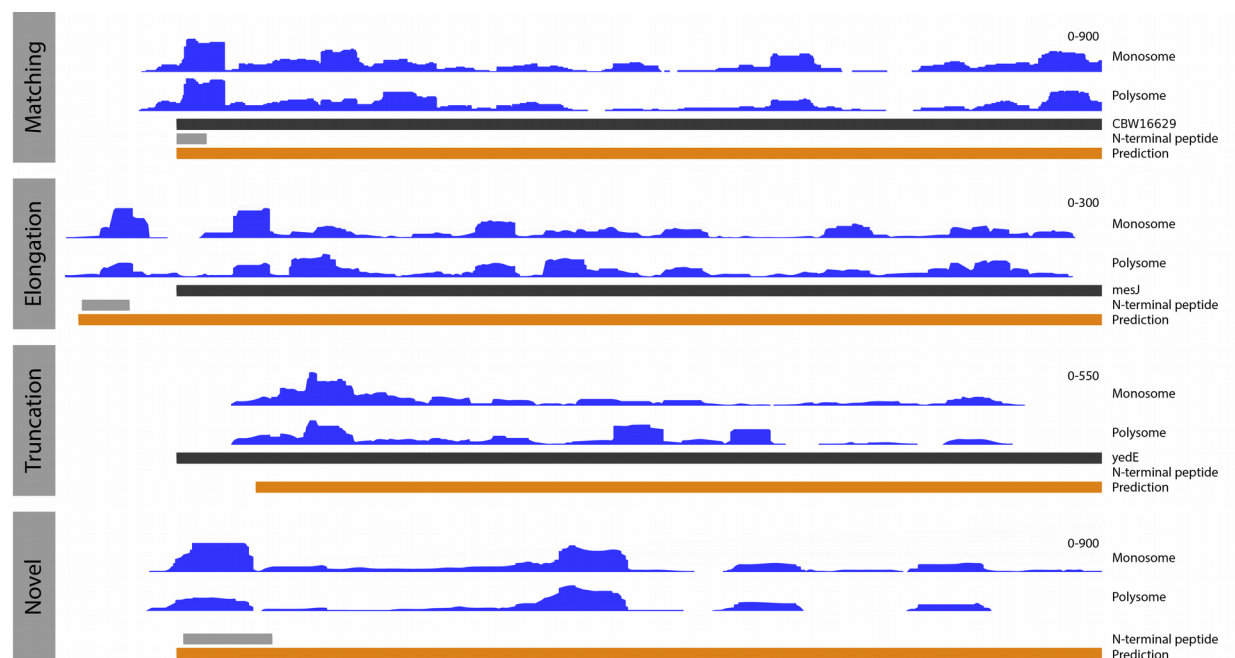
165

166  Another significant feature of prokaryotic translation initiation is the absence of intrinsic

167  structure in the region around the start codon enabling easier access for ribosomes to

168  bind[32]. We therefore calculated the average free energy over all predicted sites and

169  compared them to the previous annotation in the same genes. Consistent with GC content

170  patterns, the annotated sites display a lower propensity to form secondary structure

171  upstream of the start codon in elongated ORFs, and downstream of the start codon in

172  truncated ORFs. In the predicted sites these less-structured regions can clearly be

173  observed directly over the start codon, highly indicative of true initiation sites (Figure 4b

174  lower).

175

176  Ribosomes translocate along mRNAs three nucleotides at a time, corresponding to one

177  codon and amino acid (aa). Consequently, reads originating from bona fide translated

178  regions also exhibit a three nucleotide periodicity in adjusted read counts, with a bias

179  towards mapping to the first nucleotide in each codon[21]. At initiation sites read distribution

180  therefore switches from a random distribution upstream to a periodic, biased distribution

181  downstream. Comparing the density of reads falling into each of the three codon positions,

182  in elongated ORFs we observe increased read density at the first nucleotide position

183    upstream of annotated, but not predicted TISs. Similarly at truncated ORFs we see a

184    decrease in the density of reads at the first nucleotide position downstream of the

185    annotated TIS but not the predicted TIS (Figure 3c).

186

187    Taken together the patterns in read distribution, SD motifs, GC bias, unstructured regions

188    and triplicate periodicity, provide clear and consistent support that the TISs which we re-

189    annotate, show on average, a higher agreement with features indicative of canonically

190    translated prokaryotic ORFs, than their corresponding previously annotated counterparts.

191



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

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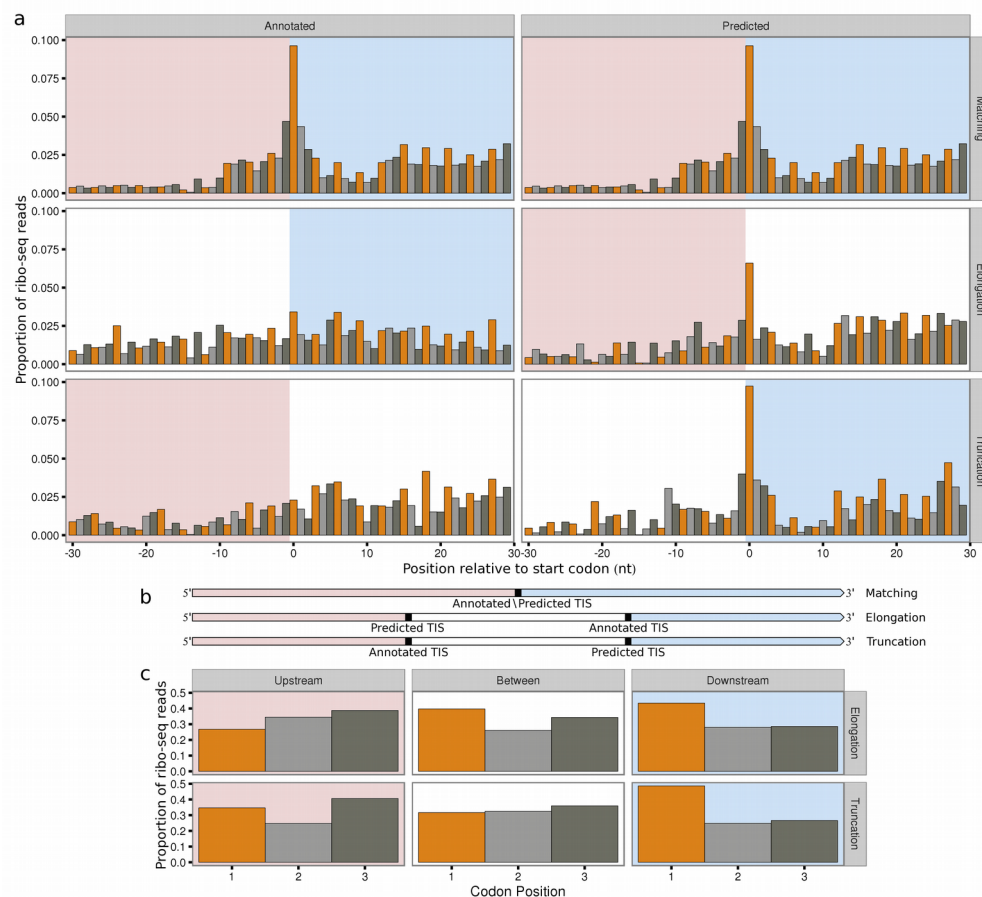201 supported by N-terminal evidence and ribo-seq coverage.

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

Bar colour indicates codon position. Downstream regions are highlighted in pink, upstream regions are highlighted in light blue. (**a**) Meta plots showing the proportion of scaled ribo-seq reads in relation to annotated or predicted translation initiation sites, for ORFs matching annotated genes (n=3853), predicted elongations (n=214) or predicted truncations (n=205). Contributions from each gene are scaled to a sum of one. Annotated TISs show increased ribo-seq density upstream (elongations), or downstream of start codons (truncations). (**b**) Transcript models. (**c**) Bar plots showing the sum of proportions of scaled ribo-seq read counts in each codon position. For truncations regions are 30 nt upstream of the annotated TIS, between the annotated and predicted TIS and 30 nt downstream of the predicted TIS. For elongations regions are 30 nt upstream of the predicted TIS, between the predicted and annotated TIS, and 30 nt downstream of annotated TIS. Three nt periodicity does not occur upstream of predicted TISs (truncations), but does occur upstream of annotated TISs (elongations).
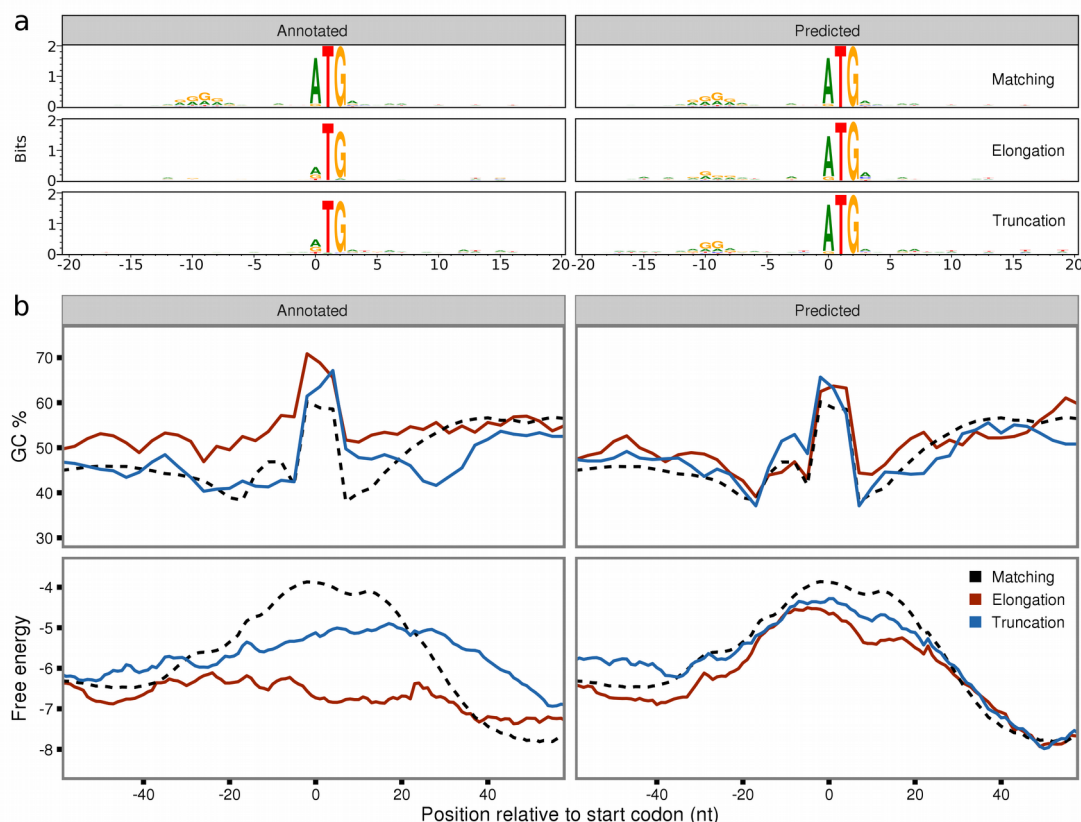
**Figure 4:** Sequence and structure features support re-annotation of translation initiation sites.

(**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.

262 **N-terminal proteomics confirms predicted sites**

263

264 In order to experimentally validate the accuracy of the predictions positional proteomics

265 analyses enriching for protein N-termini were performed. Blocked N-termini were

266 identified at 1040 *S.* Typhimurium ORFs, from which a high confidence subset of Nt-

267 formylated Met-starting N-termini was selected (see methods) and used for assessing the

268 accuracy of the model. In total 114 high confidence N-termini were identified supporting

269 102 annotated CDSs, three N-terminal CDS elongations, and nine N-terminal CDS

270 truncations. Because genomic positions with N-terminal peptide support were excluded

271 from the set used to train the random forest model, these high confidence TIS positions can

272 be used to determine the accuracy of the predictions. Of the 102 N-terminally supported

273 annotated genes, 96 were predicted by the model. Furthermore two of the elongations, and

274 four of the truncations were captured (Supplementary Table 3). Assuming that none of

275 these genes have multiple initiation sites the sensitivity of the model can be estimated to be

276 0.9444, the specificity to be 0.9991 and the positive predicted value to be 0.9444.

277

278 The remaining set of blocked N-termini supported a further 668 annotated CDSs, 50

279 suggested CDS elongations and 310 suggested CDS truncations. In addition, we found

280 peptides matching the predicted start positions from three distinct novel regions (defined

281 as ORFs at least 300 nt in length, in regions that were not overlapping with annotated

282 genes or regions at least 999 bp upstream of annotated genes). Comparing the predictions

283 to the wider blocked N-termini set we find support for 648 predictions that match

284 annotated TISs, 22 predicted elongations, 23 predicted truncated and three novel regions

285 (Supplementary Table S4).

286

287 **Translation initiation sites are predicted at novel genomic regions**

288

289 In order to discover potential novel translated ORFs we applied our prediction models to

290 look for TISs in genomic regions outside annotated ORFs. Novel ORFs that were similar in

291 size to known CDSs (> 100aa) and with ribo-seq coverage along a high proportion of the

292 ORF (>75 % coverage, see methods) were considered candidate translated novel ORFs. Of

293 the 219 (monosome) and 193 (polysome) ORFs under consideration, 104 and 115 novel

294 translated ORFs were predicted respectively. 61 of these novel translated ORFs were

295    common to both replicates (38.61% agreement) and used as a high confidence set of novel

296    predictions. Unlike the annotated genes, these novel ORFs are not previously confirmed as

297    translated regions and most had significantly lower read density (mean FPKM of 8) than

298    annotated genes (mean FPKM of 126). The higher discrepancy between the two replicates

299    is mainly a consequence of low-abundance start sites that did not pass the threshold in

300    either of the replicates.

301

302    Read density plots over the novel ORFs revealed features consistent with protein coding

303    regions, but with higher variance due to the low number of ORFs. Specifically, GC content

304    increases downstream of the initiation codon, the regions around the initiation codon have

305    less intrinsic structure potential and Shine-Dalgarno sequences are present upstream

306    (Supplementary Fig. S3). Additionally, three of the predicted novel translated ORFs were

307    supported by N-terminally enriched peptide evidence (a representative example is shown in

308    Figure 2). A further 22 showed high similarity to known protein sequences, four of which

309    contained functional protein domains (Supplementary Table S5).

310

311    **Tetracycline treated samples improve classifier accuracy**

312

313    While reads isolated from elongating ribosomes provide sufficient information to predict

314    the majority of translation start sites we set out to explore the full potential of our classifier

315    in combination with publicly available data from initiating ribosomes. A recent study on *E.*

316    *coli*[12] demonstrated the use of tetracycline as a translation inhibitor to enrich for footprints

317    from initiating ribosomes in prokaryotes. The tetracycline datasets show the pattern that

318    we expect to see from initiating ribosomes as a range of read lengths starting 28-14nt

319    upstream of the initiation codon (5' data), but ending at the same positions 14-15nt

320    downstream of the initiation codon (3' data). An additional pattern of shorter fragment

321    lengths can also be observed starting 26-18nt upstream, and ending 2nt downstream of the

322    initiation codon (Supplementary Fig. S1.a,b,g,h).

323

324    We trained separate classifiers on chloramphenicol (elongating) and tetracycline (initiating)

325    libraries from this dataset, using two replicates for each of the conditions (Supplementary

326    Table S6). Model performance was evaluated with receiver operating characteristic (ROC)

327    curves on the validation datasets for each replicate, the resulting AUC values of 0.9993 and

328    0.9994 in the tetracycline replicates were higher than those of chloramphenicol samples

329    (0.9992 and 0.9983). The parameter importance in each of the models is shown in

330    Supplementary Tables S7-10. The chloramphenicol models predicted a total of 3111 ORFs,

331    including 57 elongations and 53 truncations (Supplementary Table S11). In the tetracycline

332    dataset a total of 3711 ORFs were predicted, with 86 elongations and 79 truncations

333    (Supplementary Table S12).

334

335    *E. coli* predictions were assessed against the ecogene curated set of 923 experimentally

336    verified protein starts[33]. Genes within this dataset were excluded from the sets that were

337    used to train the random forest models, in order to provide a means of assessing the

338    accuracy of the ORF predictions. Five of the verified protein starts correspond to

339    pseudogenes without annotated CDSs, of the remaining 917 verified protein starts, 821

340    (89.53%) matched ORFs in the tetracycline predictions, with 24 (2.62%) predicted ORFs in

341    disagreement with the curated set (11 elongations, 13 truncations). In the chloramphenicol

342    predictions 760 (82.88%) were found to match ecogene start sites, and 27 (2.94%) were

343    found to be inconsistent (13 elongations, 14 truncations) with the verified protein starts

344    (assuming genes do not have multiples TIS) (Supplementary Tables S11-12). Based on the

345    experimentally verified starts the tetracycline-based classifier resulted in higher accuracy

346    (sensitivity 0.9194, specificity 0.9996, positive predictive value 0.9716) than the

347    chloramphenicol-based classifier (sensitivity 0.8539, specificity 0.9996, positive predictive

348    value 0.9657). Surprisingly, the difference was not major arguing that using initiating

349    ribosomes is not a requirement to obtain a good annotation of initiation sites.

350

351

352

353 **DISCUSSION**

354

355

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

369

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

375

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

385

386    In conclusion, this study demonstrates the utility of ribo-seq fragment length patterns for

387    TIS identification across multiple experimental conditions. These models provide a

388    significant step forward in experimental TIS discovery, facilitating the move towards

389    complete ORF annotation in both presumably well-annotated model organisms, as well as

390    the ever growing list of newly sequenced genomes.

391 **ONLINE METHODS**

392

393

394 **Preparation of ribo-seq libraries**

395

396 Overnight stationary cultures of wild type *S.* Typhimurium (*Salmonella enterica* serovar

397 Typhimurium - strain SL1344) grown in LB media at 37 °C with agitation (200 rpm) were

398 diluted at 1:200 in LB and grown until they reached and OD600 of 0.5 (i.e., logarithmic

399 (Log) phase grown cells). Bacterial cells were pre-treated for 5 min with chloramphenicol

400 (Sigma Aldrich) at a final concentration of 100 µg/ml before collection by centrifugation

401 (6000 × g, 5 min) at 4 °C. Collected cells were flash frozen in liquid nitrogen. The frozen

402 pellet of a 50 ml culture was re-suspended and thawed in 1 ml ice-cold lysis buffer for

403 polysome isolation (10 mM $MgCl_2$, 100 mM $NH_4Cl$, 20 mM Tris.HCl pH 8.0, 20 U/ml of

404 RNase-free DNase I (NEB 2 U/µl), 1mM chloramphenicol (or 300 µg/ml), 20 µl/ml lysozyme

405 (50mg/ml in water) and 100 u/ml SUPERase.In™ RNase Inhibitor (Thermo Fisher

406 Scientific, Bremen, Germany)), vortexed and left on ice for 2 min with periodical agitation.

407 Subsequently, the samples were subjected to mechanical disruption by two repetitive cycles

408 of freeze-thawing in liquid nitrogen, added 5 mM $CaCl_2$, 30µl 10% DOC and 1 × complete

409 and EDTA-free protease inhibitor cocktail (Roche, Basel, Switzerland) and left on ice for 5

410 min. Lysates were clarified by centrifugation at 16,000 x g for 10 min at 4 °C.

411

412 For the monosome sample, the supernatant was subjected to MNase (Roche diagnostics

413 Belgium) digestion using 600 U MNase (about~ 1000 U per mg of protein). Digestion of

414 polysomes proceeded for 1 h at 25 °C with gentle agitation at 400 rpm and the reaction

415 was stopped by the addition of 10 mM EGTA. Next, monosomes were recovered by

416 ultracentrifugation over a 1 M sucrose cushion in polysome isolation buffer without RNase-

417 free DNase I and lysozyme, and with 2 mM DTT added using a TLA-120.2 rotor for 4 hr at

418 75,000 rpm and 4 °C.

419

420 For the selective purification of monosomes from polysomes (polysome sample), the

421 supernatant was resolved on 10-55% (w/v) sucrose gradients by centrifugation using an

422 SW41 rotor at 35,000 rpm for 2.5 hr at 4 °C. The sedimentation profiles were recorded at

423 260 nm and the gradient fractionated using a BioComp Gradient Master (BioComp)

424    according to the manufacturer's instructions. Polysome-enriched fractions were pooled and

425    subjected to MNase digestion and monosome recovery as described above.

426

427    Ribosome-protected mRNA footprints with sizes ranging from 26-34 nt were selected and

428    processed as described previously[14] with some minor adjustments as previously described[35].

429    The resulting ribo-seq cDNA libraries of the monosome and polysome sample were

430    duplexed and sequenced on a NextSeq 500 instrument (Illumina) to yield 75 bp single-end

431    reads.

432

433    **Ribo-seq data processing**

434

435    Ribo-seq data were preprocessed with cutadapt[36] to remove sequencing adaptors,

436    discarding reads less that 20 nt in length after trimming. Trimmed reads were initially

437    aligned to the SILVA RNA database version 119[37], the remaining reads were then mapped

438    to either Salmonella enterica serovar Typhimurium - strain SL1344 (Assembly:

439    GCA_000210855.2) or Escherichia coli str. K-12 substr. MG1655 (Assembly:

440    GCA_000005845.2). Alignments were performed with bowtie2[38]. Reads were brought to

441    codon resolution by adjusting the 5' position of each read by a fixed distance offset, specific

442    to each fragment length, based on visual identification of periodicity meta plots of the read

443    counts per fragment length (Supplementary Figs. S5-6). In the *S.* Typhimurium dataset the

444    following fragments lengths were selected and adjusted by the values in brackets, in the

445    monosome sample 29 (13 nt), 30 (14 nt), 31 (15 nt), 32 (16 nt), 33 (17 nt), and in the

446    polysome sample 29 (13 nt), 30 (14 nt), 31 (15 nt), 32 (16 nt), 33 (17 nt) and 34 (18 nt).

447    Selected reads of the indicated lengths account for 39.98 and 48.69 % of total reads for the

448    monosome and polysome samples, respectively.

449

450    Recent publications reporting prokaryotic ribo-seq[26,28,39,40] suggest that read fragments from

451    libraries digested with micrococcal nuclease align more precisely to their 3' rather than 5'

452    ends. Consistent with this, we observe a modest increase in the periodicity of meta profiles

453    of the *S.* Typhimurium ribo-seq libraries when reads are brought to codon resolution from

454    the 3' end (Supplementary Fig. S1), however this does not hold true for the *E. coli* datasets,

455    where the use of  3' poly adenosine adaptors, results in a loss of resolution at the 3' end

456    after read trimming (Supplementary Fig. S1), making the use of 5' ends preferable.

457 Regardless, the protected read fragment patterns that we use in the input feature vectors

458 for the classifier takes both length and position into consideration. Consequently the

459 classifier is unaffected by this choice. However, to maintain consistency throughout this

460 study read counting for model predictors was performed from the 5' end for all libraries.

461

462 **Read distributions and heat maps**

463

464 Ribo-seq read distributions were summarised over all annotated start codons in the *S.*

465 Typhimurium and *E. coli* annotations respectively. 5' read counts were taken from regions

466 30nt upstream to 60nt downstream of the start codon, 3' read counts were taken from the

467 first nucleotide of the start codon up to 90 nucleotides downstream. All reads with a MAPQ

468 greater than 10, from the upper 90% of genes by total CDS expression were included. Total

469 counts were scaled to a sum of one per individual region, in order to not disproportionately

470 favour profiles from highly expressed genes. Meta plots were then produced to show the

471 proportion of read counts over the window across all genes. 3' and 5' heatmaps were

472 generated from the scaled regions, showing the number of standard deviations from the

473 row (fragment length) mean.

474

475 **Model implementation**

476

477 For each candidate TIS a feature vector was defined as each nucleotide in a -20 to +10nt

478 window around the position, the ribo-seq 5' FPKM (fragments per kilobase per million

479 mapped reads) between the current position and the next in-frame downstream stop codon,

480 the count of potential in-frame start sites (codons within one edit distance of ATG) from the

481 nearest in-frame upstream stop codon to the current position, the proportion of 5' ribo-seq

482 reads upstream in a 20 nt window, the proportion of 5' ribo-seq reads downstream in a 20

483 nt window, the ratio of 5' ribo-seq reads up and downstream and the proportion 5' ribo-seq

484 counts per fragment length for a fixed range of positions in relation to current site

485 (selected from visual inspection of 5' fragment length heatmaps (Supplementary Fig. S4).

486 In the *S.* Typhimurium samples fragment lengths 20-35 nt in positions -20 to -11 and 0 nt,

487 were used. In the *E. coli* datasets for the tetracycline samples fragment lengths 20-35 nt at

488 positions -25 to -16 nt, were selected and for chloramphenicol lengths 30 to 50 nt, at

489 positions -25 to -16 and -1 to +1nt were used.

490  Stop-to-stop windows were defined for each annotated gene as all in-frame positions

491  between the nearest in-frame upstream stop codon and the stop codon of the gene (with a

492  maximum length cut-off 999 nt upstream).

493

494  The H2o random forest implementation[41] was used and the models were trained with

495  positive examples of randomly selected annotated start codons from the upper 50% of

496  genes ranked by ribo-seq expression over the gene CDS. We additionally required that the

497  positive examples were not among the genes supported by N-terminal peptides in the *S.*

498  Typhimurium samples or included in the ecogene dataset for the *E. coli* samples, since

499  these were retained for model accuracy assessment. Negative examples were randomly

500  selected from in-frame codons in the stop-to-stop windows both upstream and downstream

501  of the annotated TIS. The *S.* Typhimurium models were trained on 1500 positive and 6000

502  negative positions, with an independent validation set of 200 positive and 800 negative

503  positions, the validation set was used for parameter training (number of trees monosome:

504  600, polysome: 600). The *E. coli* models were trained on 1100 positive and 4400 negative

505  positions, with an independent validation set of 200 positive and 800 negative positions for

506  parameter tuning (number of trees: CM1: 950, CM2: 950, TET2: 650 and TET3: 700).

507  Predictions were then run against all cognate and near cognate (defined as one edit

508  distance from ATG) in-frame positions, in the stop-to-stop regions. Novel predictions were

509  performed against all cognate and near cognate codons in stop-to-stop regions around

510  ORFs of at least 300nt in length, with a ribo-seq read coverage of 0.75 or more (ORF

511  coverage was defined as the proportion of nucleotides in each predicted ORF that at least

512  one ribo-seq read mapped to), that did not overlap with annotated exons. ORFs were

513  delineated by extending each candidate TIS to the closest in-frame stop codon. For a given

514  stop-to-stop region the model selected the TIS with the highest positive predicted score per

515  sample. Predictions from the replicates for each of the datasets were then compared,

516  discarding predictions that were unique to only one replicate.

517

518  **N-terminal proteomics**

519

520  Overnight stationary cultures of wild type *S.* Typhimurium (*Salmonella enterica* serovar

521  Typhimurium - strain SL1344) grown in LB media at 37 °C with agitation (200 rpm) were

522  diluted at 1:200 in LB and grown until they reached and OD600 of 0.5 (i.e., logarithmic

523  (Log) phase grown cells). Bacterial cells were collected by centrifugation (6000 × g, 5 min)

524    at 4 °C, flash frozen in liquid nitrogen and cryogenically pulverized using a liquid nitrogen

525    cooled pestle and mortar. The frozen pellet of a 50 ml culture was re-suspended and

526    thawed in 1 ml ice-cold lysis buffer (50 mm $NH_4HCO_3$ (pH 7.9) supplemented with a

527    complete protease inhibitor cocktail tablet (Roche Diagnostics GmbH, Mannheim,

528    Germany) and subjected to mechanical disruption by two repetitive freeze-thaw and

529    sonication cycles (i.e. 2 minutes of sonication on ice for 20-s bursts at output level 4 with a

530    40% duty cycle (Branson Sonifier 250; Ultrasonic Convertor)). The lysate was cleared by

531    centrifugation for 15 min at 16,000 × $g$ and the protein concentration measured using the

532    protein assay kit (Bio-Rad) according to the manufacturer's instructions. The lysate was

533    added Gu.HCl (4M f.c.) and subjected to N-terminal COFRADIC analysis as described

534    previously[42]. Free amines were blocked at the protein level making use of an N-

535    hydroxysuccinimide ester of (stable isotopic encoded) acetate (i.e. NHS esters of $^{13}C_2D_3$

536    acetate), which allows distinguishing in vivo and in vitro blocked N-terminal peptides[43]. The

537    modified protein sample was digested overnight with sequencing-grade modified trypsin

538    (1/100 (w/w trypsin /substrate)) at 37 °C and subsequent steps of the N-terminal

539    COFRADIC procedure were performed as previously described[42].

540

541    **LC-MS/MS analysis**

542

543    LC-MS/MS analysis was performed using an Ultimate 3000 RSLC nano HPLC (Dionex,

544    Amsterdam, the Netherlands) in-line connected to an LTQ Orbitrap Velos mass

545    spectrometer (Thermo Fisher Scientific, Bremen, Germany). The sample mixture was

546    loaded on a trapping column (made in-house, 100 μm I.D. × 20 mm, 5 μm beads C18

547    Reprosil-HD, Dr. Maisch). After back flushing from the trapping column, the sample was

548    loaded on a reverse-phase column (made in-house, 75 m I.D. × 150 mm, 5 μm beads C18

549    Reprosil-HD, Dr. Maisch). Peptides were loaded in solvent A' (0.1% trifluoroacetic acid, 2%

550    acetonitrile (ACN)) and separated with a linear gradient from 2% solvent A'' (0.1% formic

551    acid) to 50% solvent B' (0.1% formic acid and 80% ACN) at a flow rate of 300 nl/min

552    followed by a wash reaching 100% solvent B'. The mass spectrometer was operated in

553    data-dependent mode, automatically switching between MS and MS/MS acquisition for the

554    ten most abundant peaks in a given MS spectrum. Full scan MS spectra were acquired in

555    the Orbitrap at a target value of 1E6 with a resolution of 60,000. The 10 most intense ions

556    were then isolated for fragmentation in the linear ion trap, with a dynamic exclusion of 20

557    s. Peptides were fragmented after filling the ion trap at a target value of 1E4 ion counts.

558    Mascot Generic Files were created from the MS/MS data in each LC run using the Mascot

559    Distiller software (version 2.5.1.0, Matrix Science, www.matrixscience.com/Distiller.html).

560    To generate these MS/MS peak lists, grouping of spectra was allowed with a maximum

561    intermediate retention time of 30 s and a maximum intermediate scan count of 5. Grouping

562    was done with a 0.005 Da precursor tolerance. A peak list was only generated when the

563    MS/MS spectrum contained more than 10 peaks. There was no de-isotoping and the

564    relative signal-to-noise limit was set at 2.

565

566    The generated MS/MS peak lists were searched with Mascot using the Mascot Daemon

567    interface (version 2.5.1, Matrix Science). Searches were performed using a 6-FT database

568    of the *S.* Typhimurium genome combined with the Ensembl protein sequence database

569    (assembly AMS21085v2 version 86.1), which totalled 139,408 entries after removal of

570    redundant sequences. The 6-FT database was generated by traversing the entire genome

571    across the six reading frames and searching for all NTG (N=A,T,C,G) start codons and

572    extending each to the nearest in frame stop codon (TAA,TGA,TAG), discarding ORFs less

573    than 21nt in length. The Mascot search parameters were set as follows: Heavy acetylation

574    at lysine side-chains (Acetyl:2H(3)C13(2) (K)), carbamidomethylation of cysteine and

575    methionine oxidation to methionine-sulfoxide were set as fixed modifications. Variable

576    modifications were formylation, acetylation and heavy acetylation of N-termini

577    (Acetyl:2H(3)C13(2) (N-term)) and pyroglutamate formation of N-terminal glutamine (both

578    at peptide level). Endoproteinase semi-Arg-C/P (semi Arg-C specificity with Arg-Pro

579    cleavage allowed) was set as enzyme allowing for no missed cleavages. Mass tolerance was

580    set to 10 ppm on the precursor ion and to 0.5 Da on fragment ions. Peptide charge was set

581    to 1+, 2+, 3+ and instrument setting was put to ESI-TRAP. Only peptides that were ranked

582    one, have a minimum amino acid length of seven, scored above the threshold score, set at

583    95% confidence, and belonged to the category of *in vivo* or *in vitro* blocked N-terminal

584    peptides compliant with the rules of initiator methionine (iMet) processing[44] were withheld.

585    More specifically, iMet processing was considered in the case of iMet-starting N-termini

586    followed by any of the following amino acids; Ala, Cys, Gly, Pro, Ser, Thr, Met or Val and

587    only if the iMet was encoded by ATG or any of the following near-cognate start codons;

588    GTG and TTG (Supplementary Table S13). In contrast to eukaryotic nascent protein N-

589    termini, the typical lack of significant steady-state levels of N-terminal protein modification

590    (e.g. Nt-acetylation or Nt-formylation), warrant caution to unequivocally assign bacterial

591    protein N-termini as proxies of translation initiation. As such, a high confidence subset of

592    Nt-formylated Met-starting N-termini, was selected (Supplementary Table S14).

593 **Assessing model accuracy**

594

595 GC content was calculated at the third nucleotide positions for all annotated and predicted
596 ORFs and mean GC values were summarised for each subgroup of predicted ORFs
597 (matching annotations, truncations and elongations) in 9 nt sliding windows, over regions
598 57 nt upstream and 57 nt downstream of the annotated or predicted start sites.

599

600 -20 to + 20 nt nucleotide sequences were extracted around the predicted and annotated
601 TIS in the *S.* Typhimurium and *E. coli* genomes. Sequence logos were generated for each
602 subgroup of matching annotations, truncations, elongations and novel genes, using the
603 weblogo tool[45].

604

605 The minimum free energy of RNA secondary structure around predicted and annotated
606 ORFs was estimated with RNAfold version 2.1.9 from the ViennaRNA package[46]. Mean free
607 energy values were summarised for each ORF class in 39 nt sliding window across regions
608 57 nt up and downstream of the start codon.

609

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

618

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

625 ATG), in CDS regions of supported genes that were neither predicted nor supported were

626 considered true negatives.

627

628 Amino acid sequences of novel ORF were compared to known proteins in the nonredundant

629 protein database (Update date:2016/12/15) and protein domains (cdd.v.3.15) using

630 BLASTP[47] (version 2.5.1+). Hits with the greatest coverage of query sequence and lowest e-

631 value were selected. Hits were considered highly similar if they shared >95% identify to a

632 protein sequence, over 100% of the novel ORF sequence

633

634 **DATA AVAILABILITY**

635

636 The previously published *E. coli* ribo-seq dataset was downloaded from the NCBI SRA

637 (BioProject ID:PRJDB2960). *S.* Typhimurium ribo-seq sequencing data has been deposited

638 in NCBI's Gene Expression Omnibus[48] and is accessible through GEO Series accession

639 number GSE91066. *S.* Typhimurium mass spectrometry proteomics data have been

640 deposited to the ProteomeXchange Consortium via the PRIDE[49] partner repository with the

641 dataset identifier PXD005579 and 10.6019/PXD005579.

642

643 **ACKNOWLEDGMENTS**

644

650

651 **AUTHOR CONTRIBUTIONS**

652

653 A.G., E.V. and P.V.D. conceived the study and wrote the manuscript; A.G. performed the

654 computational analysis; P.V.D performed the proteomics experiment. E.N. and P.V.D.

655 performed proteomics analysis; P.V.D. and V.J. prepared the ribo-seq libraries.

656  **COMPETING FINANCIAL INTERESTS**

657

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

## REFERENCES

1. Hall, J., Hazlewood, G. P., Surani, M. A., Hirst, B. H. & Gilbert, H. J. Eukaryotic and prokaryotic signal peptides direct secretion of a bacterial endoglucanase by mammalian cells. J. Biol. Chem. 265, 19996–19999 (1990).

2. Kozak, M. Initiation of translation in prokaryotes and eukaryotes. Gene 234, 187–208 (1999).

3. Delcher, a L., Harmon, D., Kasif, S., White, O. & Salzberg, S. L. Improved microbial gene identification with GLIMMER. Nucleic Acids Res. 27, 4636–4641 (1999).

4. Brocchieri, L., Kledal, T. N., Karlin, S. & Mocarski, E. S. Predicting coding potential from genome sequence: application to betaherpesviruses infecting rats and mice. J. Virol. 79, 7570–96 (2005).

5. Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site identification. BMC Bioinformatics 11, 119 (2010).

6. Suzek, B. E., Ermolaeva, M. D., Schreiber, M. & Salzberg, S. L. A probabilistic method for identifying start codons in bacterial genomes. Bioinformatics 17, 1123–1130 (2001).

7. Zhu, H. Q., Hu, G. Q., Ouyang, Z. Q., Wang, J. & She, Z. S. Accuracy improvement for identifying translation initiation sites in microbial genomes. Bioinformatics 20, 3308–3317 (2004).

8. Ou, H. Y., Guo, F. B. & Zhang, C. T. GS-Finder: A program to find bacterial gene start sites with a self-training method. Int. J. Biochem. Cell Biol. 36, 535–544 (2004).

9. Tech, M., Morgenstern, B. & Meinicke, P. TICO: A tool for postprocessing the predictions of prokaryotic translation initiation sites. Nucleic Acids Res. 34, 588–590 (2006).

10. Hartmann, E. M. & Armengaud, J. N-terminomics and proteogenomics, getting off to a good start. Proteomics 14, 2637–2646 (2014).

11. Berry, I. J., Steele, J. R., Padula, M. P. & Djordjevic, S. P. The application of terminomics for the identification of protein start sites and proteoforms in bacteria. Proteomics 16, 257–272 (2016).

12. Nakahigashi, K. et al. Comprehensive identification of translation start sites by tetracycline-inhibited ribosome profiling. DNA Res. 23, 193–201 (2016).

13. Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S. & Weissman, J. S. Genome-wide analysis in vivo of translation with nucleotide resolution using ribosome profiling. Science 324, 218–23 (2009).

693    14. Ingolia, N. T., Lareau, L. F. & Weissman, J. S. Ribosome profiling of mouse embryonic
694    stem cells reveals the complexity and dynamics of mammalian proteomes. Cell 147, 789–
695    802 (2011).

696    15. Brar, G. A. et al. High-Resolution View of the Yeast Meiotic Program Revealed by
697    Ribosome Profiling. Science (80-. ). 335, 552–557 (2012).

698    16. Michel, A. M. et al. Observation of dually decoded regions of the human genome using
699    ribosome profiling data. Genome Res. 22, 2219–2229 (2012).

700    17. Chew, G. L. et al. Ribosome profiling reveals resemblance between long non-coding
701    RNAs and 5' leaders of coding RNAs. Development 140, 2828–34 (2013).

702    18. Crappé, J. et al. Combining in silico prediction and ribosome profiling in a genome-wide
703    search for novel putatively coding sORFs. BMC Genomics 14, 648 (2013).

704    19. Bazzini, A. A. et al. Identification of small ORFs in vertebrates using ribosome
705    footprinting and evolutionary conservation. EMBO J. 33, 981–993 (2014).

706    20. Pauli, A. et al. Toddler: an embryonic signal that promotes cell movement via Apelin
707    receptors. Science 343, 1248636 (2014).

708    21. Calviello, L. et al. Detecting actively translated open reading frames in ribosome
709    profiling data. Nat. Methods 13, 1–9 (2015).

710    22. Duncan, C. D. S. & Mata, J. The translational landscape of fission-yeast meiosis and
711    sporulation. Nat. Struct. Mol. Biol. 21, 641–7 (2014).

712    23. Ingolia, N. T. et al. Ribosome Profiling Reveals Pervasive Translation Outside of
713    Annotated Protein-Coding Genes. Cell Rep. 8, 1365–1379 (2014).

714    24. Fritsch, C. et al. Genome-wide search for novel human uORFs and N-terminal protein
715    extensions using ribosomal footprinting. Genome Res. 22, 2208–2218 (2012).

716    25. Lee, S. et al. Global mapping of translation initiation sites in mammalian cells at single-
717    nucleotide resolution. Proc. Natl. Acad. Sci. U. S. A. 109, E2424–E2432 (2012).

718    26. Nakahigashi, K. et al. Effect of codon adaptation on codon-level and gene-level
719    translation efficiency in vivo. BMC Genomics 15, 1115 (2014).

720    27. Heyer, E. E. & Moore, M. J. Redefining the Translational Status of 80S Monosomes. Cell
721    164, 757–769 (2016).

722    28. Woolstenhulme, C. J., Guydosh, N. R., Green, R. & Buskirk, A. R. High-Precision analysis
723    of translational pausing by ribosome profiling in bacteria lacking EFP. Cell Rep. 11, 13–21
724    (2015).

725    29. Shine, J. & Dalgarno, L. Determinant of cistron specificity in bacterial ribosomes.
726    Nature 254, 34–38 (1975).

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

762    46. Lorenz, R. et al. ViennaRNA Package 2.0. Algorithms Mol. Biol. 6, 26 (2011).

763    47. Altschul, S. F. et al. Gapped BLAST and PSI-BLAST:a new generation of protein

764    database search programs. Nucleic Acids Res 25, 3389–3402 (1997).

765    48. Edgar, R., Domrachev, M. & Lash, A. E. Gene Expression Omnibus: NCBI gene

766    expression and hybridization array data repository. Nucleic Acids Res 30, 207–210 (2002).

767    49. Vizcaino, J. A. et al. 2016 update of the PRIDE database and its related tools. Nucleic

768    Acids Res. 44, D447–D456 (2016).