

Quantifying Selection on Codon Usage in Signal Peptides: Gene Expression and Amino Acid Usage Explain Apparent Selection for Inefficient Codons

Alexander L. Cope¹, Robert L. Hettich^{1,2}, and Michael A. Gilchrist^{1,3,4}

¹Genome Science and Technology, University of Tennessee, Knoxville

²Chemical Sciences Division, Oak Ridge National Laboratory, Oak Ridge,

TN

³Department of Ecology and Evolutionary Biology, University of

Tennessee, Knoxville

⁴National Institute for Mathematical and Biological Synthesis, Knoxville,

TN

Last compiled on Thursday 14th June, 2018 at 18:44.

1 **Abstract**

2 Secreted proteins play central roles across all taxa. Although secretion mechanism can vary
3 across taxa, all taxa share the Sec secretion pathway. A critical and distinct feature shared
4 by Sec secreted proteins is the signal peptide. Researchers claim signal peptides contain a
5 bias for translation inefficient codons in signal peptides, leading researchers to suggest selec-
6 tion favors translation inefficiency in this region. We investigate codon usage in the signal
7 peptides of *E. coli* using the Codon Adaptation Index (CAI) and tRNA Adaptation Index
8 (tAI), and the ribosomal overhead cost formulation of the stochastic evolutionary model
9 of protein production rates (ROC-SEMPPR). Initial comparisons between signal peptides
10 and 5'-ends of non-signal peptide genes using CAI and tAI are consistent with translation-
11 ally inefficient codons being preferred in signal peptides. However, simulations reveal these
12 differences are due to amino acid usage and gene expression – we find evidence for novel
13 selection disappears when accounting for both of these factors. In contrast, ROC-SEMPPR,
14 a mechanistic population genetics model capable of separating the effects of selection and
15 mutation bias, shows codon usage bias (CUB) of the signal peptides is indistinguishable from
16 the 5'-coding regions of cytoplasmic proteins. Additionally, we find CUB in the 5'-coding
17 regions is weaker than later segments of the gene. Results illustrate the value in using mod-
18 els grounded in population genetics to interpret genetic data. In summary, we show failure
19 to account for mutation bias and the effects of gene expression on the efficacy of selection
20 against translation inefficiency can lead to a misinterpretation of codon usage patterns.

21 **Introduction**

22 A secreted protein can broadly be defined as any protein entering a secretory pathway
23 for transport through a cellular membrane. These proteins serve important cellular func-
24 tions, including metabolism and antibiotic resistance (Green and Mecsas, 2016; Saier, 2006).
25 Secreted proteins also play essential roles in the virulence of pathogenic bacteria (Green

26 and Mecsas, 2016). Numerous secretion systems exist and vary between and within taxa
27 (Bendtsen *et al.*, 2005; Green and Mecsas, 2016; Saier, 2006). Despite the diversity of secre-
28 tion pathways, the general secretion pathway, also commonly referred to as the Sec pathway,
29 is found across all domains of life (Green and Mecsas, 2016; Natale *et al.*, 2008). In brief,
30 proteins are transported to the SecYEG translocon located in the membrane via the SecA/B
31 chaperone-dependent (SecA/B and SRP) or chaperone-independent manner (Natale *et al.*,
32 2008; Tsirigotaki *et al.*, 2017). All SecA/B-dependent proteins and chaperone-independent,
33 as well as some SRP-dependent proteins, contain a short peptide chain located at the N-
34 terminus of the protein (Green and Mecsas, 2016; Natale *et al.*, 2008; Tsirigotaki *et al.*, 2017)
35 This short peptide chain is called the signal peptide which is essential for the protein to enter
36 into the SecYEG translocon. Although signal peptides do vary in their amino acid sequences,
37 signal peptides have distinct physicochemical properties which constrain their amino acid
38 usage. A signal peptide generally consists of 3 regions: a positively charged N-terminus,
39 a hydrophobic core, and a polar C-terminus, where the signal peptide is cleaved from the
40 rest of the protein (a.k.a. the mature peptide) (Natale *et al.*, 2008; Tsirigotaki *et al.*, 2017;
41 Zalucki *et al.*, 2009).

42 The ability to accurately predict signal peptides is useful for identifying secreted proteins
43 in non-model organisms; this has led to the development of machine learning approaches
44 to predict signal peptides which take advantage of the distinct physicochemical properties
45 of signal peptides, such as SignalP (Petersen *et al.*, 2011). Although the physicochemical
46 properties of signal peptides are consistent, previous work found altering the N-terminus
47 has a range of effects on protein secretion: from a decrease in secretion to no effect (Inouye
48 *et al.*, 1982; Nesmeyanova *et al.*, 1997; Puziss *et al.*, 1989; Vlasuk *et al.*, 1983). These varying
49 effects led some researchers to suspect other mechanisms also contribute to the efficacy of
50 protein secretion (Zalucki *et al.*, 2009, 2011a).

51 Numerous studies suggest codon usage bias (CUB) – the non-uniform usage of synony-
52 mous codons – contributes to effective protein secretion in *E. coli* (Burns and Beacham,

53 1985; Power *et al.*, 2004; Zalucki and Jennings, 2007; Zalucki *et al.*, 2008, 2010, 2011b).
54 Power *et al.* (2004) found *E. coli* K12 MG1655 signal peptides are biased for translation
55 inefficient codons, which are predicted to be translated slower than their synonymous coun-
56 terparts. This is in stark contrast to the rest of the *E. coli* proteome, where *E. coli* is biased
57 towards the most efficient codons (Ikemura, 1981; Power *et al.*, 2004). Li *et al.* (2009); Liu
58 *et al.* (2017); Mahlab and Linial (2014) examined the usage of inefficient codons in signal
59 peptides of *S. coelicolor*, *S. cerevisiae*, and various multicellular eukaryotes and came to
60 similar conclusions when applying codon usage indices such as the Codon Adaptation Index
61 (CAI, Sharp and Li, 1987) and tRNA Adaptation Index (tAI, dos Reis *et al.*, 2004). Consis-
62 tent across this work is the interpretation that selection is driving the apparent increase in
63 inefficient codon usage in signal peptides. Similar studies concluded an overabundance of the
64 lysine codon AAA at the second position in the signal peptide promoted efficient translation
65 initiation (Zalucki *et al.*, 2007).

66 Researchers proposed an adaptive role for inefficient codons in the protein secretion pro-
67 cess in which the combination of efficient translation initiation and inefficient translation
68 resulted in reduced distance between sequential ribosomes translating the mRNA of a pro-
69 tein containing a signal peptide (Zalucki *et al.*, 2009, 2011a). They argued this would lead
70 to more efficient recycling of the chaperones involved in the secretion process. Other expla-
71 nations for the observed increase in inefficient codons include the inability of *E. coli* SRP to
72 induce a translational pause following signal peptide recognition (Powers and Walter, 1997;
73 Zalucki *et al.*, 2009) and slowing down the co-translational folding of the protein, as a folded
74 protein cannot be translocated through the SecYEG translocon (Power *et al.*, 2004; Zalucki
75 and Jennings, 2007; Zalucki *et al.*, 2008, 2011a). If signal peptides have a different CUB
76 relative to the rest of the genome, then codon-level information could be incorporated into
77 signal peptide prediction tools.

78 In contrast, a recent analysis of ribosome profiling data found no difference in the ri-
79 bosome densities of the signal peptides and the 5'-ends of nonsecretory genes in various

80 eukaryotes (Liu *et al.*, 2017). If selection were acting on codon usage in signal peptides to
81 slow down translation, we would expect to see higher ribosome densities in these regions.
82 Additionally, while both Mahlab and Linial (2014) and Liu *et al.* (2017) examined codon
83 usage in relation to secretion in *H. sapiens* using a metric based on tAI, only Mahlab and
84 Linial (2014) found results consistent with increased frequencies of inefficient codons in sig-
85 nal peptides. From a population genetics perspective, it is surprising statistically significant
86 results were obtained in a mammal, which usually have little adaptive CUB due to their
87 lower effective population sizes (Charlesworth, 2009; Lynch *et al.*, 2016). More recently,
88 Samant *et al.* (2014) found codon optimization of a signal peptide improved localization of
89 the protein to the periplasm of *E. coli*, seemingly contradicting a general role for inefficient
90 codon usage in signal peptides. A potential reason for these contradictions is the previous
91 analyses of signal peptide codon usage by Li *et al.* (2009); Liu *et al.* (2017); Mahlab and
92 Linial (2014); Power *et al.* (2004) did not adequately account for the evolutionary forces
93 shaping codon usage (Bulmer, 1990; Gilchrist *et al.*, 2015; Shah and Gilchrist, 2011; Wallace
94 *et al.*, 2013).

95 We re-examined CUB in signal peptides of *E. coli* using CAI, tAI, and ROC-SEMPPR - a
96 population genetics model which accounts for selection, mutation bias, and gene expression
97 - to determine if selection on codon usage in signal peptides differs from the 5'-ends of
98 genes. Although we find significant differences in codon usage using CAI and tAI, we present
99 evidence these differences are due to signal peptide-specific amino acid biases and differences
100 in the gene expression distributions of signal peptide and non-signal peptide genes. When
101 comparing signal peptides and the 5'-ends of non-signal peptide genes with ROC-SEMPPR,
102 we find signal peptide codon usage is consistent with the 5'-ends of genes not containing a
103 signal peptide. We find selection on codon usage favors the efficient codons, but the strength
104 of selection is weaker at the 5'-ends, corroborating previous analyses (Eyre-Walker, 1996;
105 Gilchrist and Wagner, 2006; Gilchrist, 2007; Power *et al.*, 2004; Qin *et al.*, 2004).

106 Our work demonstrates the value of analyzing CUB from a formal population genetics

107 framework, as well as highlights potential limitations with using more common metrics such
108 as CAI for analyzing codon usage on relatively small regions of the genome. Failure to ac-
109 count for variation in the strength of selection due to variation in gene expression can lead
110 to conflating mutation bias with selection, resulting in a misinterpretation of observed codon
111 usage patterns. Our work also illustrates the importance of considering non-adaptive forces
112 in shaping biological phenomenon before invoking adaptive explanations (Gould and Lewon-
113 tin, 1979). We believe this is particularly important in the modern genomic-age when the
114 combination of large datasets, misinterpretation of p-values, and and inherent bias towards
115 adaptationist interpretations could mislead researchers.

116 **Materials and Methods**

117 **Signal Peptide Prediction**

118 Signal peptides were predicted using SignalP 4.1 (Petersen *et al.*, 2011) using both the
119 default cutoff D-score of 0.51 and a more conservative D-score of 0.75. In brief, SignalP
120 consists of two neural networks, one for determining the amino acid sequence similarity to
121 signal peptides and the other for identifying the most likely cleavage site. The results of
122 both neural networks are combined into one value, called the D-score. The D-score ranges
123 between 0 and 1. Setting the cutoff D-score closer to 1 results in a lower false positive rate. A
124 set of confirmed signal peptides for *E. coli* K12 MG1655 was taken from The Signal Peptide
125 Website. All analyses in the main text will focus on the set of signal peptides with $D \geq 0.51$
126 as this set provides us with the most data; analyses of the $D > 0.75$ and set of confirmed
127 signal peptides give similar results (see Supplementary Material).

128 **ROC-SEMPPR**

129 Given a set of protein-coding genes, ROC-SEMPPR employs a Markov Chain Monte Carlo
130 (MCMC) to estimate codon specific parameters for mutation bias ΔM and pausing times $\Delta \eta$

131 for each codon within a synonymous codon family (Table 1). In previous work, $\Delta\eta$ was scaled
132 relative to the most efficient codon, which had $\Delta\eta$ and ΔM values fixed at 0. To avoid the
133 choice of reference codon affecting our comparisons of CUB between regions, all $\Delta\eta$ values
134 in this paper are re-scaled such that these values are centered around 0 for each amino acid.
135 The $\Delta\eta$ values reflect the strength and direction of selection against translation inefficiency
136 in a set of protein-coding regions (e.g. the signal peptides). A region with stronger selection
137 against translation inefficiency will have higher $\Delta\eta$ values on average than a region with
138 weaker selection. Similarly, a region which favors translation inefficiency would be expected
139 to have $\Delta\eta$ values which negatively correlate with a region which favors translation efficiency.

140 ROC-SEMPPR also estimates an average protein production rate ϕ for each gene (Table
141 1). We find ROC-SEMPPR estimated ϕ values correlate well with empirical measurements of
142 protein production rates for *E. coli* (see Supplementary Methods: Assessing ROC-SEMPPR
143 Model Adequacy and Figures S1 - S2). If changes in synonymous codon usage alter the
144 efficiency at which a protein is translated, then such a change will have the largest impact
145 on the energetic costs of proteins with high production rates, making ϕ a more appropriate
146 gene expression metric than say, mRNA abundance or protein abundance. Thus, we use
147 protein production rates ϕ as our metric of gene expression. For more details on ROC-
148 SEMPPR, see Gilchrist *et al.* (2015). Analysis of CUB with ROC-SEMPPR was performed
149 using AnaCoDa (Landerer *et al.*, 2018).

Parameters	Description
$\Delta\eta_i$	Cost of translating codon i relative to reference codon
ΔM_i	Mutation bias towards codon i relative to the reference codon
ϕ_k	Average Protein Production Rate of gene k

Table 1: Description of ROC-SEMPPR parameters used in this paper.

150 CAI and tAI

151 Analysis of CUB was also performed using CAI (Sharp and Li, 1987) and tAI (dos Reis
152 *et al.*, 2004). Both CAI and tAI quantify CUB by assigning weights to the 61 sense codons.
153 For CAI, each codon is assigned a weight based on its relative frequency to its synonymous
154 counterparts in a reference set of highly expressed genes, such as ribosomal protein coding
155 genes. The key assumption of CAI is the most frequent codons in the reference set are
156 the most efficient codons (Sharp and Li, 1987). In contrast, tAI assigns weights based on
157 tRNA abundances corresponding to a codon, as well as accounting for codon-anticodon
158 interactions. The key assumption of tAI is the most efficient codons are usually those with
159 the most abundant tRNA (dos Reis *et al.*, 2004).

160 CAI and tAI both range between 0 and 1. A CAI score closer to 1 represents a sequence
161 which more closely resembles the codon usage of the reference set of genes, while a tAI
162 closer to 1 indicates a sequence is more closely adapted to the genomic tRNA pool (dos Reis
163 *et al.*, 2004; Sharp and Li, 1987). Calculations for CAI were performed using the AnaCoDa
164 (Landerer *et al.*, 2018), while tAI was calculated using the R package tAI (dos Reis, 2016).

165 Generating Datasets

166 Previous analysis of the *E. coli* genome found a set of genes with CAI values that had a
167 negative correlation with their gene expression estimates (dos Reis *et al.*, 2003). It is expected
168 many of these genes were the result of horizontal gene transfer and had not yet reached
169 evolutionary equilibrium with respect to their CUB. We repeated the analysis described in
170 dos Reis *et al.* (2003) on the current *E. coli* K12 MG1655 genome (version 3 ,NC_000913.3).
171 Briefly, correspondence analysis was performed using CodonW (Peden, 1999), followed by
172 clustering based on the principle axis scores using the CLARA algorithm (Maechler *et al.*,
173 2018) in R. Our analysis was consistent with the findings of (dos Reis *et al.*, 2003), revealing
174 782 genes with a CUB deviating significantly from the majority of the *E. coli* genome. We
175 will refer to this set of 782 genes as the “exogenous” component of the genome and the rest of

176 the *E. coli* genome as the “endogenous” for simplicity. All analyses presented will consider
177 only “endogenous” genes because the “exogenous” genes may violate the assumptions of
178 ROC-SEMPPR, CAI, and tAI.

179 Proteins with a signal peptide were split into the signal peptide and the mature peptide
180 – the segment of the peptide chain after the signal peptide. On average, the signal peptides
181 were 23 codons long. For comparisons to the 5'-ends of nonsecretory genes – defined here
182 as those lacking a signal peptide – the first 23 codons of the nonsecretory genes were used.
183 We note the nonsecretory genes have an average protein production rate ϕ lower than that
184 of the signal peptide genes ($\bar{\phi} = 0.992$ and $\bar{\phi} = 1.08$, respectively, Figure S3).

185 As the strength of selection on CUB scales with protein production rate ϕ , we created a
186 control group that eliminates differences in the distribution of ϕ for the nonsecretory genes
187 and signal peptide genes. Specifically, the nonsecretory genes were selected using acceptance-
188 rejection sampling to create the “pseudo-secreted proteins”. In brief, acceptance-rejection
189 sampling is a procedure for sampling from a population such that its distribution of a metric
190 for one population mirrors the distribution of the same metric for another population. In
191 this case, the pseudo-secreted proteins were sampled such that the mean and variance of the
192 $\log(\phi)$ values reflected those of the genes with a signal peptide. The CUB signature of a
193 gene varies with protein production rate ϕ ; thus we can be more confident any differences
194 seen between genes with a signal peptide and pseudo-signal peptide genes are not due to
195 differences in their respective ϕ distributions. All pseudo-secreted proteins were split into two
196 regions we will refer to as the “pseudo-signal peptides” and the “pseudo-mature peptides”
197 (the first 23 codons and the remainder of the gene, respectively).

198 To assess the performance of CAI and tAI when comparing regions with differences in the
199 distributions of protein production rates ϕ and amino acid biases, simulated sequences were
200 used. Sequences based on the 5'-ends of nonsecretory genes, pseudo-signal peptides, and
201 signal peptides were simulated using the AnaCoDa package (Landerer *et al.*, 2018). To nor-
202 malize for amino acid usage, sequences 23 amino acids in length were randomly generated to

203 match the amino acid frequencies of the signal peptides. The codon usage of these sequences
204 was also simulated in AnaCoDa, assuming either the ϕ distribution of the nonsecretory genes
205 or the pseudo-secreted proteins. All sequences were simulated using the pausing times $\Delta\eta$
206 and mutation bias ΔM parameters estimated from the 5'-end of endogenous nonsecretory
207 genes.

208 CUB analyses

209 We estimated protein production rates ϕ by fitting ROC-SEMPPR to the complete protein-
210 coding sequences in the *E. coli* K12 MG1655 genome. Analysis of intragenic (eg. signal
211 vs. mature peptides) and intergenic (eg. pseudo-signal peptides vs. real signal peptides)
212 CUB was carried out using the mixture distribution functionality available in the AnaCoDa
213 implementation of ROC-SEMPPR (Landerer *et al.*, 2018). Each group of regions (eg. signal
214 peptides, mature peptides, etc.) was assumed to have an independent CUB, allowing pausing
215 time $\Delta\eta$ estimates to vary between them. We assumed mutation bias was consistent for the
216 entire genome; thus, we forced mutation bias ΔM parameters to be equal across the groups
217 of regions. ϕ was fixed for each region at the value estimated from the region's corresponding
218 complete protein-coding sequence. This is done for two reasons: (a) shorter regions, such as
219 the signal peptide, likely have insufficient information to accurately estimate ϕ and (b) this
220 guarantees our gene expression metric has the same impact on the estimates of $\Delta\eta$ and ΔM
221 for intragenic regions, such as a signal peptide and its corresponding mature peptide.

222 A Model-II regression was used to compare pausing times $\Delta\eta$ between regions. Unlike
223 ordinary least squares, Model-II regression, or errors-in-variables regression, accounts for
224 errors in both the x and y variables (Sokal and Rohlf, 1995). When both variables are
225 subject to error, which is the case for the $\Delta\eta$ estimates, the use ordinary least squares leads
226 to downwardly biased parameter estimates. A Model-II regression slope $\beta = 1$ (or the $y = x$
227 line) will serve as the null hypothesis, as this indicates both the strength and direction of
228 selection between two regions are the same. The intercept parameter was fixed at $\alpha = 0$

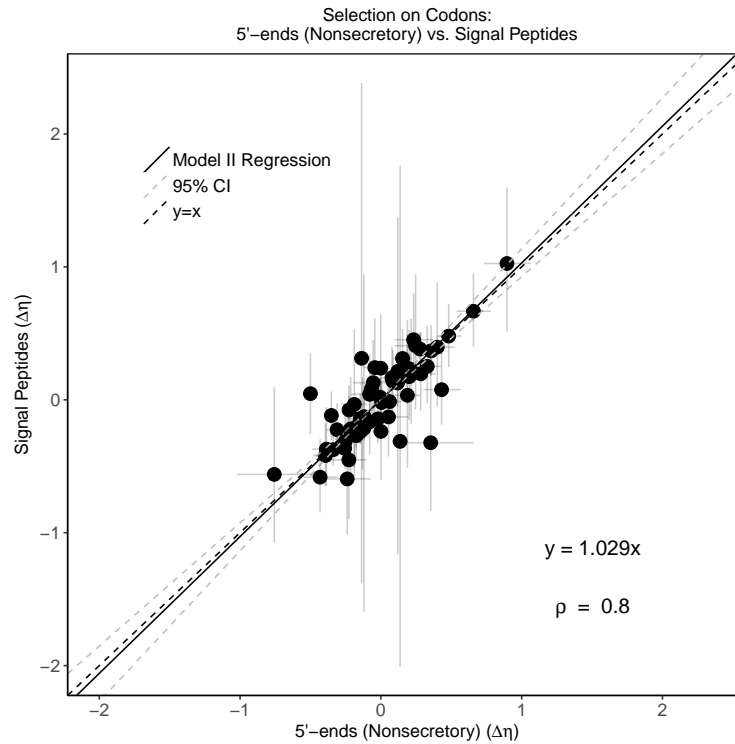
229 because the $\Delta\eta$ estimates are scaled such that the mean value of $\Delta\eta$ is 0. We note that when
230 we allowed the α parameter to vary, it was as expected, approximately 0. For more details
231 on our use of Model-II regression, see Supplementary Methods.

232 CAI and tAI were used to compare codon usage between signal peptides, 5'-ends, and
233 pseudo-signal peptides (dos Reis *et al.*, 2003, 2004; Sharp and Li, 1987). As recommended
234 by Sharp and Li (1987), methionine and tryptophan were not included when normalizing for
235 the length of the gene in our calculations of CAI. Statistical significance was assessed using
236 a one-tailed Welch's t-test in R (R Core Team, 2018). R and Python scripts used for this
237 paper can be found at https://github.com/acope3/Signal_Peptide_Scripts.

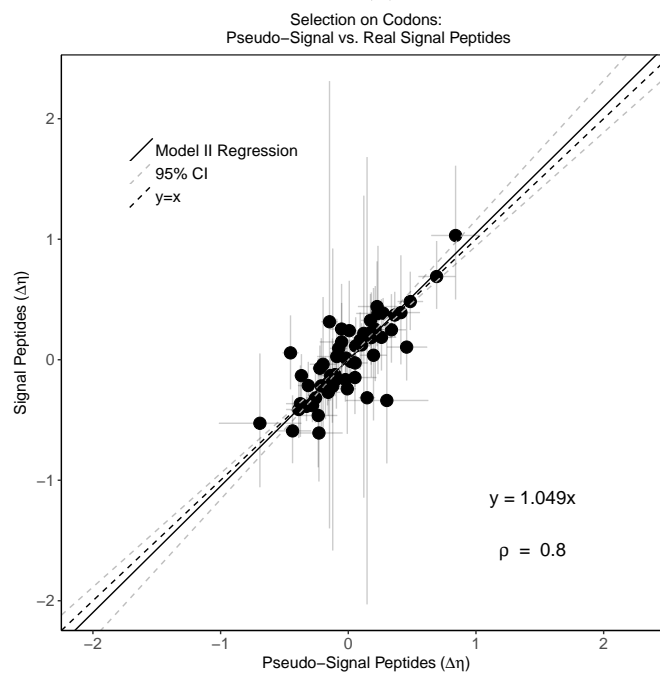
238 Results

239 Our analysis of CUB in signal peptides and the 5'-ends of nonsecretory genes using ROC-
240 SEMPPR revealed these regions to be highly similar. Qualitatively, the expected codon
241 frequencies for the 5'-ends of nonsecretory genes and the signal-peptides based on the pausing
242 time $\Delta\eta$ and mutation bias ΔM values estimated from these regions are similar (Figure S4).
243 Notable exceptions appear to be cysteine, aspartic acid, lysine, glutamine, and tyrosine;
244 however, the 95% posterior probability intervals of cysteine and glutamine are the only
245 ones which fail to overlap with $y = x$ line. When comparing the pausing times $\Delta\eta$ of
246 signal peptides to the 5'-ends of nonsecretory genes using a Model-II regression, we find no
247 significant difference from the $y = x$ line (slope β 95% confidence interval: 0.923 – 1.128,
248 Figure 1a). To determine if differences were not detected due to underlying differences in
249 the distributions of ϕ , we compared $\Delta\eta$ estimated from signal peptides and pseudo-signal
250 peptides. Again, no statistically significant difference from the $y = x$ line was found and the
251 expected codon frequencies are similar (β 95% confidence interval: 0.939 – 1.149, Figure 1b
252 and S5). Similar results are obtained using the signal peptides with a D-score greater than
253 0.75 or the confirmed signal peptides (Figures S6 - S7). We also see no significant result when

254 using empirically estimated ϕ values ($\beta = 0.908$, 95% confidence interval: 0.67 – 1.16, Figure
255 S8), although these results show much more variability. The increased variability in the $\Delta\eta$
256 values and corresponding regression line is unsurprising given the empirically estimated ϕ
257 values are subject to significant noise (Figure S2), but are, in this case, treated as fixed
258 values representative of the true average protein production rate for a gene.



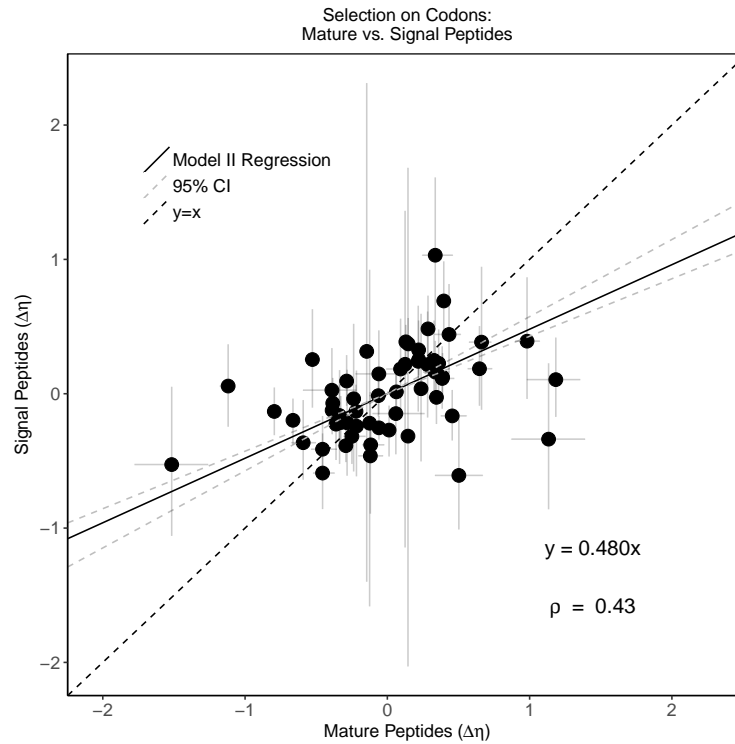
(a)



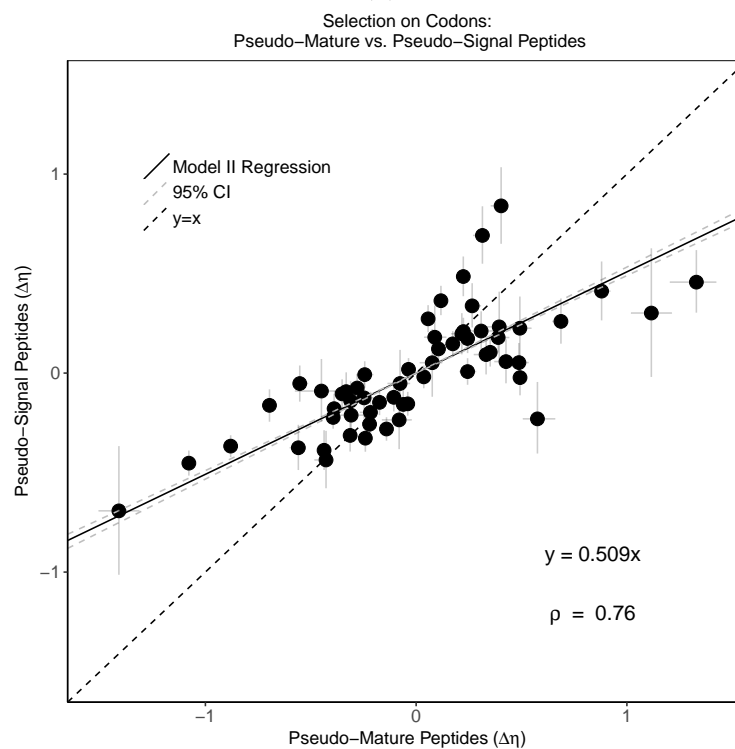
(b)

Figure 1: Comparing the pausing time estimates $\Delta\eta$ between (a) the 5'-ends of nonsecretory genes or (b) pseudo-signal peptides to signal peptides. Grey dashed lines represent the 95% confidence intervals of the regression line. Results clearly show a strong positive linear relationship ($\rho = 0.80$) between the regions and a regression line not significantly different from $y = x$.

259 The Model-II regression lines estimated from the mature vs. signal peptide comparison
260 and the pseudo-mature vs. pseudo-signal peptide comparison are similar, which serves as
261 further evidence the selection on codon usage in signal peptides and the 5'-ends of nonse-
262 cretory genes is the same (Figure 2). The mature vs. signal peptide comparison produces
263 a regression line with slope $\beta = 0.480$ (95% confidence interval: 0.428 - 0.574), while the
264 pseudo-mature vs. pseudo-signal peptide comparison produces a regression line with slope
265 $\beta = 0.496$ (95% confidence interval: 0.490 - 0.533). If selection on codon usage differs in
266 signal peptides from pseudo-signal peptides, we would not expect to see similar regression
267 lines.



(a)



(b)

Figure 2: (a) Comparing the codon pausing time estimates $\Delta\eta$ between mature peptides and signal peptide regions. Grey dashed lines represent the 95% confidence intervals of the regression line. Results show a positive linear relationship ($\rho = 0.43$) between the $\Delta\eta$ estimates for the two regions. This indicates codons favored in one region tend to be favored in the other. (b) Same comparison for pseudo-signal peptide genes. Regression estimates are similar to those estimated for the mature and signal peptide comparison.

268 Noting CAI and tAI do not account for the effects of gene expression, mutation bias, drift,
269 or amino acid biases, we found signal peptides have lower CAI and tAI values compared to
270 the first 23 codons of nonsecretory genes (one-tailed Welch's t-test, $p < 10^{-5}$). This was also
271 the case when looking at the pseudo-signal peptides, which normalizes for protein production
272 rates ϕ . These results with CAI and tAI can potentially be explained by either the preferred
273 use of inefficient codons in signal peptides *or* as artifacts of amino acid biases. Signal peptides
274 have a different amino acid composition from the 5'-end due to the required physicochemical
275 properties of this region (Figure S9). We examined the robustness of tAI and CAI as a
276 means of quantifying differences in selection on codon usage when underlying differences
277 between amino acid composition and ϕ exists using data simulated under the same mutation
278 bias ΔM and pausing time $\Delta\eta$ parameters. When comparing simulated signal peptides to
279 simulated 5'-end of nonsecretory genes and simulated pseudo-signal peptides using CAI, the
280 simulated signal peptides are found to have a significantly lower mean CAI (Welch's t-test,
281 $p < 0.05$) 100% of the time (Figure 3a-b), despite the fact the $\Delta\eta$ and ΔM parameters used
282 to simulate these regions were the same. This suggests the amino acid usage is biasing the
283 signal peptides towards a lower CAI.

284 When using simulated 5'-ends of nonsecretory genes which have amino acid composition
285 consistent with the signal peptides, the p-values were heavily skewed towards 1. (Figure
286 3c). This odd behavior is due to the differences in the ϕ distribution differences of the signal
287 peptide and nonsecretory genes. As the former has a higher mean ϕ , the signal peptides on
288 average will have a stronger CUB after normalizing for the amino acid biases. A one-tailed
289 Welch's t-test with the alternative hypothesis being signal peptides have a lower mean CAI,
290 when in reality they likely have a larger mean CAI, would skew the p-value distribution
291 towards 1. Importantly, ROC-SEMPPR did not detect significant differences between signal
292 peptides and the 5'-ends of non-secretory genes, despite differences in the ϕ distributions
293 (Figure 1a). When normalizing for both amino acid usage and ϕ , significant differences in
294 CAI are found approximately 4% of the time, which is close to the expected number of false

295 positives at the 0.05 significance level (Figure 3d). Similar results are seen when using tAI
296 (Figure S10). Our results indicate CAI and tAI are prone to inflating differences in CUB
297 between two regions when differences in ϕ and amino acid usage are not accounted for.

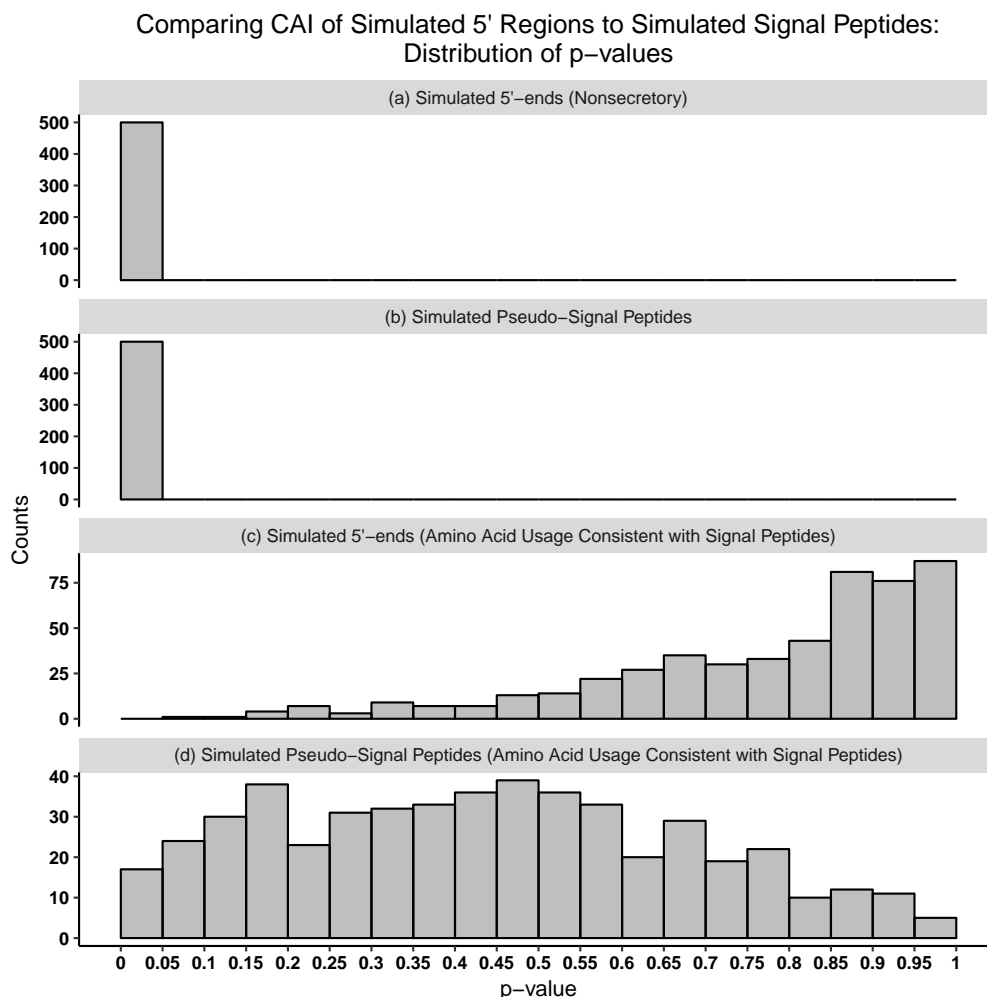


Figure 3: Distribution of p-values from a one-tailed Welch's t-test comparing CAI in simulated nonsecretory 5'-ends, pseudo-signal peptides, and signal peptides in which all regions were simulated using the same pausing time $\Delta\eta$ and ΔM parameters. (a-b) The CAI of simulated signal peptides was found to be significantly lower on average at a 100% false positive rate when compared to simulated 5'-ends of nonsecretory genes and simulated pseudo-signal peptides. (c) Adjusting the amino acid frequencies of the 5'-end of nonsecretory genes to match those of the signal peptides results in a heavily skewed distribution. (d) Adjusting the amino acid frequencies of the pseudo-signal peptides to match those of the signal peptides results in a more uniform distribution.

298 Notably, selection on codon usage near the N-terminus appears to be on average approx-

299 imately 50% weaker than the remainder of the gene based on the slope β . Previous analyses
300 using a variety of codon usage metrics found CUB near the 5'-end to be weaker than mid-
301 dle sections of the gene (Eyre-Walker, 1996; Gilchrist and Wagner, 2006; Gilchrist, 2007;
302 Hockenberry *et al.*, 2014; Qin *et al.*, 2004; Power *et al.*, 2004), with these differences being
303 attributed to selection against nonsense errors and to maintain translation initiation effi-
304 ciency by reducing mRNA secondary structure. We confirm this trend using ROC-SEMPPR
305 (Figure S11).

306 It was also proposed selection for translation initiation efficiency was shaping signal
307 peptide codon usage, particularly the use of lysine codon AAA, in signal peptides at position
308 2 of the peptide (Zalucki *et al.*, 2007). We do find AAA appears to be slightly favored in signal
309 peptides, which is not the case in the pseudo-signal peptides, although the 95% posterior
310 probability interval overlaps with the $y = x$ line (Figure S12). If the slight but statistically
311 insignificant favored usage of AAA is due to an increased selection for translation initiation
312 efficiency in signal peptides, then removing the first 3 codons when analyzing signal peptide
313 codon usage should remove this effect. Doing so results in no change in the behavior of AAA,
314 suggesting if there is any selection for increased AAA usage in signal peptides, it is not due
315 to selection for increased translation initiation efficiency (Figure S13). Notably, AAA is
316 both mutationally and selectively-favored for lysine by *E. coli*. Keeping in mind selection on
317 CUB is weaker near the 5'-end of the genes in *E. coli*, the combination of weaker selection,
318 mutational favorability, and a slight increase in the occurrence of lysine in signal peptides
319 (Figure S9) likely drives up the frequency of codon AAA in signal peptides relative to the
320 5'-ends of nonsecretory genes.

321 Discussion

322 In summary, we found no evidence suggesting a general significant difference between selec-
323 tion on codon usage in signal peptides and the 5'-ends of nonsecretory genes in *E. coli* using

324 a mechanistic model of CUB which incorporates the effects of selection, mutation bias, gene
325 expression, and amino acid usage. Instead, we find failures to account for amino acid usage
326 and protein production rate ϕ resulted in the commonly used codon metrics CAI and tAI
327 indicating significant differences between regions simulated under the same parameters, but
328 these differences disappear when accounting for both amino acid usage and ϕ . Importantly,
329 both amino acid usage and ϕ were significant confounding factors when analyzing CUB with
330 CAI and tAI – only accounting for one of these factors still suggested significant differences
331 between the simulated regions. Although we are not the first to note potential issues with
332 metrics like CAI or tAI for intragenic CUB analysis (Hockenberry *et al.*, 2014), our results
333 demonstrate these metrics are insufficient for intragenic CUB analysis when these regions
334 have drastically different amino acid usage or ϕ distributions, resulting in incorrect biological
335 interpretation.

336 This is not to say CUB plays no role in the secretion of specific proteins. For example,
337 experimental evidence demonstrates codon optimization of the *E. coli* maltose binding pro-
338 tein's (MBP) signal peptide results in a decrease in protein abundance. Evidence suggests
339 this is due to increased targeting of the codon optimized MBP by proteases due to improper
340 folding (Zalucki and Jennings, 2007; Zalucki *et al.*, 2008). However, CUB as a means to
341 guide proper co-translational folding is not a phenomenon unique to proteins with a signal
342 peptide (Chaney and Clark, 2015; Pechmann and Frydman, 2013; Yu *et al.*, 2015). Although
343 inefficient codons might be crucial to the fold of certain secreted proteins, our results do not
344 indicate this is any more or less so than nonsecretory genes.

345 Although we found no general difference in selection on codon usage between signal pep-
346 tides and the 5'-ends, it is possible CUB differences exist among the chaperone-dependent
347 and chaperone-independent mechanisms of the Sec pathway. We are unaware of any CUB
348 comparisons of these three groups, but researchers have noted a region of slower translation
349 downstream from the signal peptide of transmembrane proteins, which are typically secreted
350 via SRP in bacteria (Natale *et al.*, 2008). Using a modified form of the tAI, previous efforts

351 found a consistent trend of inefficient codons 35-40 codons downstream of the SRP-binding
352 site in various yeasts species (Pechmann *et al.*, 2014). Ribosomal profiling data taken from *S.*
353 *cerevisiae* provided experimental support for this hypothesis; however this analysis was lim-
354 ited to a small, closely-related phylogeny. Further work is needed to determine the generality
355 of this observation to bacteria and other eukaryotes. Similarly, SRP-dependent transmem-
356 brane proteins in *E. coli* have a higher frequency of "programmed pause sites," areas of high
357 ribosomal density downstream from Shine-Dalgarno-like sequences, at the beginning of the
358 gene (Fluman *et al.*, 2014). A higher frequency of programmed pause sites was not observed
359 in the region downstream from the signal peptides in periplasmic proteins. Notably, this
360 region of higher ribosome density downstream from the signal peptides was not observed in
361 periplasmic proteins, which are normally secreted via SecA/B (Natale *et al.*, 2008; Tsirig-
362 otaki *et al.*, 2017) However, recent work challenges the findings that Shine-Dalgarno-like
363 sequences are largely responsible for translational pause (Mohammad *et al.*, 2016).

364 Notably, we do find selection on CUB is weaker at the 5'-ends relative to later portions
365 of the gene, corroborating previous work (Eyre-Walker, 1996; Gilchrist and Wagner, 2006;
366 Gilchrist, 2007; Hockenberry *et al.*, 2014; Power *et al.*, 2004; Qin *et al.*, 2004). Weaker
367 selection at the 5'-ends is often attributed to selection against nonsense errors and selection
368 against mRNA secondary structure. Importantly, the advent of ribosome profiling revealed
369 the presence of high ribosomal density at the 5'-ends, often referred to as the "5'-ramp"
370 (Tuller *et al.*, 2010). The 5'-ramp was originally thought to be the result of increased
371 selection for slow translation at the 5'-end to reduce ribosomal interference further down
372 the transcript, but simulations suggest the 5'-ramp is an artifact of short genes with high
373 initiation rates (Shah *et al.*, 2013). Selection for co-translational folding is also thought to
374 shape intragenic CUB (Chaney and Clark, 2015; Pechmann and Frydman, 2013; Yu *et al.*,
375 2015). Further work is needed to understand how these various selective forces are balanced
376 to maintain translation efficiency and efficacious protein biogenesis.

377 Ultimately, our work further illustrates the value of population genetics models which

378 include nonadaptive evolutionary forces when analyzing genomic data. Biologists are often
379 tempted to explain statistically significant results in the context of selection and adaptation,
380 but researchers must first provide evidence these results cannot be explained by nonadap-
381 tive evolutionary forces (eg. mutation bias and genetic drift) and/or as an artifact of some
382 other constraint on the trait of interest (eg. amino acid biases). We are certainly not the
383 first to note the importance of considering nonadaptive explanations. Almost four decades
384 ago, Gould and Lewontin (1979) critiqued the propensity of evolutionary biologists to invoke
385 natural selection and adaptation without seriously considering possible nonadaptive expla-
386 nations. The explosion of genomic data means now, more than ever, biologists should be
387 hesitant to adopt adaptationists explanations to biological phenomenon without first inves-
388 tigating if such results could be shaped by nonadaptive forces. The embrace of "big data"
389 by biological researchers is a double-edged sword: while we have the ability to investigate
390 patterns and explore hypotheses which would not have been possible 20 years ago, the use
391 of large datasets can lead to incredibly small p-values, which are often misinterpreted as
392 both evidence of a strong effect and a small probability of the null hypothesis being true
393 (Wasserstein and Lazar, 2016). The misinterpretation of p-values and a bias towards adap-
394 tationist explanations can be a dangerous combination, with researchers over-interpreting
395 their results and misleading other researchers.

396 The development of models incorporating both adaptive and nonadaptive evolutionary
397 forces will be important for understanding the selective forces shaping complex biological
398 data. In the case of the studying CUB, codon indices like CAI have long been employed,
399 but these metrics often are unable to disentangle the effects of amino acid biases, mutation,
400 and selection. While often good proxies of gene expression, these indices do not directly
401 incorporate gene expression information into the weights estimated for each codon. This
402 could lead to further problems of conflating mutation bias with selection when comparing
403 CUB across regions. In contrast, because ROC-SEMPPR is grounded in population genetics
404 and thus, is able to decouple selection and mutation bias, it serves as a more accurate and

405 evolutionarily-grounded tool for researchers interested in studying CUB.

References

- Bendtsen, J. D., Kiemer, L., Fausbøll, A., and Brunak, S. (2005). Non-classical protein secretion in bacteria. *BMC Microbiology*, **5**(1), 58.
- Bulmer, M. (1990). The effect of context on synonymous codon usage in genes with low codon usage bias. *Nucleic Acids Res.*, **18**(10), 2869–2873.
- Burns, D. and Beacham, I. (1985). Rare codons in *E. coli* and *S. typhimurium* signal sequences. *FEBS Letters*, **189**, 318–324.
- Chaney, J. and Clark, P. (2015). Roles for synonymous codon usage in protein biogenesis. *Annu. Rev. Biophysics*, **44**, 143–166.
- Charlesworth, B. (2009). Effective population size and patterns of molecular evolution and variation. *Nature Reviews Genetics*, **10**, 195–205.
- dos Reis, M. (2016). *tAI: The tRNA adaptation index*. R package version 0.2.
- dos Reis, M., Wernisch, L., and Savva, R. (2003). Unexpected correlations between gene expression and codon usage bias from microarray data for the whole *Escherichia coli* k-12 genome. *Nucleic Acids Research*, **31**(23), 6976–6985.
- dos Reis, M., Savva, R., and Wernisch, L. (2004). Solving the riddle of codon usage preferences: a test for translational selection. *Nucleic Acids Research*, **32**(17), 5036–5044.
- Eyre-Walker, A. (1996). Synonymous codon bias is related to gene length in *Escherichia coli*: Selection for translational accuracy? *Mol. Biol. Evol.*, **13**(6), 864–872.
- Fluman, N., Navon, S., Bibi, E., and Pilpel, Y. (2014). mrna-programmed translation pauses in the targeting of e. coli membrane proteins. *eLife*, **3**, e03440.

- Gilchrist, M. (2007). Combining models of protein translation and population genetics to predict protein production rates from codon usage patterns. *Mol. Biol. Evol.*, **24**(11), 2362–2372.
- Gilchrist, M. and Wagner, A. (2006). A model of protein translation inducing codon bias, nonsense errors, and ribosome recycling. *Journal of Theoretical Biology*, **239**, 417–434.
- Gilchrist, M., Chen, W., Shah, P., Landerer, C., and Zaretzki, R. (2015). Estimating gene expression and codon-specific translational efficiencies, mutation biases, and selection coefficients from genomic data alone. *Genome Biology and Evolution*, **7**, 1559–1579.
- Gould, S. and Lewontin, R. (1979). The spandrels of san marco and the panglossian paradigm: A critique of the adaptationist programme. *Proceedings of the Royal Society of London*, **205**(1161), 581–598.
- Green, E. and Mecsas, J. (2016). Bacterial secretion systems - an overview. *Microbiol Spectr.*, **4**(1).
- Hockenberry, A., Sirer, M., Amaral, L., and Jewett, M. (2014). Quantifying position-dependent codon usage bias. *mol. Biol. Evol.*, **31**(7), 1880–1893.
- Ikemura, T. (1981). Correlation between the abundance of *Escherichia coli* transfer rnas and the occurrence of the respective codons in its protein genes: A proposal for a synonymous codon choice that is optimal for the *E. coli* translational system. *Journal of Molecular Biology*, **151**, 389–409.
- Inouye, S., Soberon, X., Franceschini, T., Nakamura, K., Itakura, K., and Inouye, M. (1982). Role of positive charge on the amino-terminal region of the signal peptide in protein secretion across the membrane. *Proc. Natl. Acad. Sci. USA.*, **79**, 3438–3441.
- Landerer, C., Cope, A., Zaretzki, R., and Gilchrist, M. (2018). Anacoda: analyzing codon data with bayesian mixture models. *Bioinformatics*, page bty138.

- Li, Y., Xie, Z., Du, Y., Zhou, Z., Mao, X., Lv, L., and Li, Y. (2009). The rapid evolution of signal peptides is mainly caused by relaxed selection on non-synonymous and synonymous sites. *Gene*, **436**, 8–11.
- Liu, H., Rahman, S., Mao, Y., Xu, X., and Tao, S. (2017). Codon usage bias in 5' terminal coding sequences reveals distinct enrichment of gene functions. *Genomics*, **109**, 506–513.
- Lynch, M., Ackerman, M., Gout, J., Long, H., Sung, W., Thomas, W., and Foster, P. (2016). Genetic drift, selection and the evolution of the mutation rate. *Nature Reviews Genetics*, **17**, 704–714.
- Maechler, M., Rousseeuw, P., Struyf, A., Hubert, M., and Hornik, K. (2018). *cluster: Cluster Analysis Basics and Extensions*. R package version 2.0.7-1 — For new features, see the 'Changelog' file (in the package source).
- Mahlab, S. and Linial, M. (2014). Speed controls in translating secretory proteins in eukaryotes - an evolutionary perspective. *PLoS Computational Biology*, **10**(1), e1003294.
- Mohammad, F., Woolstenhulme, C., Green, R., and Buskirk, A. (2016). Clarifying the translational pausing landscape in bacteria by ribosome profiling. *Cell Reports*, **14**, 686–694.
- Natale, P., Bruser, T., and Driessen, A. (2008). Sec- and tat-mediated protein secretion across the bacterial cytoplasmic membrane—distinct translocases and mechanisms. *Biochimica et Biophysica Acta*, **1778**, 1735–1756.
- Nesmeyanova, M., Karamyshev, A., Karamysheva, Z., Kalinin, A., Ksenzenko, V., and Kajava, A. (1997). Positively charged lysine at the n-terminus of the signal peptide of the *Escherichia coli* alkaline phosphatase provides the secretion efficiency and is involved in the interaction with anionic phospholipids. *FEBS Letters*, **403**, 203–207.

- Pechmann, S. and Frydman, J. (2013). Evolutionary conservation of codon optimality reveals hidden signatures of cotranslational folding. *Nature Structural and Molecular Biology*, **20**(2), 237–243.
- Pechmann, S., Chartron, J., and Frydman, J. (2014). Local slowdown of translation by nonoptimal codons promotes nascent-chain recognition by srp *in vivo*. *Nature Structural and Molecular Biology*, **21**(12), 1100–1105.
- Peden, J. (1999). *Analysis of Codon Usage*. Ph.D. thesis, University of Nottingham.
- Petersen, T., Brunak, S., von Heijne, G., and Nielsen, H. (2011). Signalp 4.0: discriminating signal peptides from transmembrane regions. *Nature Methods*, **8**(10), 785–786.
- Power, P., Jones, R., Beacham, I., Bucholtz, C., and Jennings, M. (2004). Whole genome analysis reveals a high incidence of non-optimal codons in secretory signal sequences of *Escherichia coli*. *Biochemical and Biophysical Research Communications*, **322**, 1038–1044.
- Powers, T. and Walter, P. (1997). Co-translational protein targeting catalyzed by the *Escherichia coli* signal recognition particle and its receptor. *The EMBO Journal*, **16**(16), 4880–4886.
- Puziss, J., Fikes, J., and Bassford, P. (1989). Analysis of mutational alterations in the hydrophilic segment of the maltose-binding protein signal peptide. *Journal of Bacteriology*, **171**, 2303–2311.
- Qin, H., Wu, W., Kreitman, J. C. M., and Li, W. (2004). Intra-genic spatial patterns of codon usage bias in prokaryotic and eukaryotic genomes. *Genetics*, **168**, 2245–2260.
- R Core Team (2018). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Saier, M. (2006). Protein secretion systems in gram-negative bacteria. *Microbe*, **1**(9), 414–419.

REFERENCES

27

- Samant, S., Gupta, G., Karthikeyan, S., and A. Nair, S. H., Sambasivam, G., and Suku-
maran, S. (2014). Effect of codon-optimized *E. coli* signal peptides on recombinant *Bacillus*
stearothermophilus maltogenic amylase periplasmic localization, yield and activity. *J. Ind.*
Microbial Biotechnol, **41**, 1435–1442.
- Shah, P. and Gilchrist, M. (2011). Explaining complex codon usage patterns with selection
for translational efficiency, mutation bias, and genetic drift. *PNAS*, **108**(25), 10231–10236.
- Shah, P., Ding, Y., Niemczyk, M., Kudla, G., and Plotkin, J. (2013). Rate-limiting steps in
yeast protein translation. *Cell*, **153**, 1589–1601.
- Sharp, P. and Li, W. (1987). The codon adaptation index - a measure of directional syn-
onymous codon usage bias, and its potential applications. *Nucl. Acids Research*, **15**(3),
1281–1295.
- Sokal, R. and Rohlf, F. (1995). *Biometry - The Principles and Practices of Statistics in*
Biological Research. W.H. Freeman, New York, 3rd edition.
- Tsirigotaki, A., Geyter, J. D., Sostaric, N., Economou, A., and Karamanou, S. (2017).
Protein export through the bacterial sec pathway. *Nature Reviews: Microbiology*, **15**,
21–36.
- Tuller, T., Carmi, A., Vestsigian, K., Navon, S., Dorfan, Y., Zaborske, J., Pan, T., Dahan, O.,
Furman, I., and Pilpep, Y. (2010). An evolutionarily conserved mechanism for controlling
the efficiency of protein translation. *Cell*, **141**, 344–354.
- Vlasuk, G., Inouye, S., Ito, H., Itakura, K., and Inouye, M. (1983). Effects of the complete
removal of basic amino acid residues from the signal peptide on secretion of lipoprotein in
Escherichia coli. *J. Biol. Chem.*, **258**, 7141–7148.
- Wallace, E., Airoidi, E., and Drummond, D. (2013). Estimating selection on synonymous

- codon usage from noisy experimental data. *Molecular Biology and Evolution*, **30**(6), 1438–1453.
- Wasserstein, R. and Lazar, N. (2016). The asa’s statement on p-values: Context, process, and purpose. *The American Statistician*, **70**(2), 129–133.
- Yu, C., Dang, Y., Zhou, Z., Wu, C., Zhao, F., Sachs, M., and Liu, Y. (2015). Codon usage influences the local rate of translation elongation to regulate co-translational protein folding. *Molecular Cell*, **59**, 744–754.
- Zalucki, Y. and Jennings, M. (2007). Experimental confirmation of a key role for non-optimal codons in protein export. *Biochemical and Biophysical Research Communications*, **355**, 143–148.
- Zalucki, Y., Power, P., and Jennings, M. (2007). Selection for efficient translation initiation biases codon usage at the second amino acid position in secretory proteins. *Nucleic Acids Research*, pages 1–7.
- Zalucki, Y., Gittins, K., and Jennings, M. (2008). Secretory signal sequence non-optimal codons are required for expression and export of β -lactamase. *Biochemical and Biophysical Research Communications*, **366**, 135–141.
- Zalucki, Y., Beacham, I., and Jennings, M. (2009). Biased codon usage in signal peptides: a role in protein export. *Trends in Microbiology*, **17**(4), 146–150.
- Zalucki, Y., Jones, C., Ng, P., Schulz, B., and Jennings, M. (2010). Signal sequence non-optimal codons are required for the correct folding of mature maltose binding protein. *Biochimica et Biophysica Acta*, **1798**, 1244–1249.
- Zalucki, Y., Beacham, I., and Jennings, M. (2011a). Coupling between codon usage, translation and protein export in *Escherichia coli*. *Biotechnology Journal*, **6**, 660–667.

REFERENCES

29

Zalucki, Y., Shafer, W., and Jennings, M. (2011b). Directed evolution of efficient secretion in the srp-dependent export of tolB. *Biochimica et Biophysica Acta*, **1808**, 2544–2550.