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

Alexander L. Cope[1], Robert L. Hettich[1,2], and Michael A. Gilchrist[1,3,4]

[1]Genome Science and Technology, University of Tennessee, Knoxville

[2]Chemical Sciences Division, Oak Ridge National Laboratory, Oak Ridge, TN

[3]Department of Ecology and Evolutionary Biology, University of Tennessee,Knoxville

[4]National Institute for Mathematical and Biological Synthesis, Knoxville, TN

Last compiled on Thursday 28th June, 2018 at 15:15.

1

# 1 Abstract

2 The Sec secretion pathway is found across all domains of life. A critical feature of Sec
3 secreted proteins is the signal peptide, a short peptide with distinct physicochemical prop-
4 erties located at the N-terminus of the protein. Previous work indicates signal peptides are
5 biased towards translationally inefficient codons, which is hypothesized to be an adaptation
6 driven by selection to improve the efficacy and efficiency of the protein secretion mechanisms.
7 We investigate codon usage in the signal peptides of *E. coli* using the Codon Adaptation
8 Index (CAI), the tRNA Adaptation Index (tAI), and the ribosomal overhead cost formu-
9 lation of the stochastic evolutionary model of protein production rates (ROC-SEMPPR).
10 Comparisons between signal peptides and 5'-end of cytoplasmic proteins using CAI and tAI
11 are consistent with a preference for inefficient codons in signal peptides. Simulations reveal
12 these differences are due to amino acid usage and gene expression – we find these differences
13 disappear when accounting for both factors. In contrast, ROC-SEMPPR, a mechanistic
14 population genetics model capable of separating the effects of selection and mutation bias,
15 shows codon usage bias (CUB) of the signal peptides is indistinguishable from the 5'-ends of
16 cytoplasmic proteins. Additionally, we find CUB at the 5'-ends is weaker than later segments
17 of the gene. Results illustrate the value in using models grounded in population genetics
18 to interpret genetic data. We show failure to account for mutation bias and the effects
19 of gene expression on the efficacy of selection against translation inefficiency can lead to a
20 misinterpretation of codon usage patterns.

21 **Keywords:** Codon usage bias; signal peptides; protein secretion

# 22 Introduction

23 A secreted protein can broadly be defined as any protein entering a secretory pathway for

24 transport through a cellular membrane. These proteins serve important cellular functions,

25 including metabolism and antibiotic resistance [15, 37]. Secreted proteins also play essential

26  roles in the virulence of pathogenic bacteria [15]. Numerous secretion systems exists and vary

27  between and within taxa [1, 15, 37]. Despite the diversity of secretion pathways, the general

28  secretion pathway, also commonly referred to as the Sec pathway, is found across all domains

29  of life [15, 26]. In brief, proteins are transported to the SecYEG translocon located in the

30  membrane in a chaperone-dependent (SecA/B and SRP) or chaperone-independent manner

31  [26, 43]. All SecA/B-dependent proteins and chaperone-independent, as well as some SRP-

32  dependent proteins, contain a short peptide chain located at the N-terminus of the protein

33  known as the signal peptide [15, 26, 43]. The signal peptide is an essential component of

34  the Sec pathway, serving as a binding site for the appropriate chaperones and/or helping

35  delay the folding of the protein [26, 43]. Although signal peptides do vary in their amino

36  acid sequences, signal peptides have distinct physicochemical properties which biases their

37  amino acid usage [26, 43, 49]. A signal peptide generally consists of 3 regions: a positively

38  charged N-terminus, a hydrophobic core, and a polar C-terminus, where the signal peptide

39  is cleaved from the rest of the protein, sometimes referred to as the "mature peptide."

40  The ability to accurately predict signal peptides is useful for identifying secreted proteins

41  in non-model organisms; this has led to the development of machine learning approaches

42  to predict signal peptides which take advantage of the distinct physicochemical properties

43  of signal peptides, such as SignalP [31]. Although the physicochemical properties of signal

44  peptides are consistent, altering the N-terminus has a range of effects on protein secretion:

45  from a decrease in the number of proteins secreted to no observable effect [18, 27, 34, 45].

46  The variability in the outcomes of neutralizing the N-terminal positive charge led to a search

47  for other mechanisms which also contribute to the efficacy of protein secretion [49, 50].

48  Numerous studies suggests codon usage bias (CUB) – the non-uniform usage of synony-

49  mous codons – contributes to effective protein secretion in *E. coli* [3, 32, 52, 51, 53, 55].

50  [32] found *E. coli* K12 MG1655 signal peptides are biased for translation inefficient codons,

51  which are predicted to be translated slower than their synonymous counterparts. This is

52  in stark contrast to the rest of the *E. coli* proteome, where *E. coli* is biased towards the

53  most efficient codons [17, 32]. [20, 21, 24] examined the usage of inefficient codons in sig-

54  nal peptides of *S. coelicolor*, *S. cerevisiae*, and various multicellular eukaryotes and came

55  to similar conclusions when applying codon usage indices such as the Codon Adaptation

56  Index (CAI) [41] and tRNA Adaptation Index (tAI) [7]. Consistent across this work is the

57  interpretation that selection is driving the apparent increase in inefficient codon usage in

58  signal peptides. Furthermore, [54] concluded an overabundance of the lysine codon AAA at

59  the second position in the signal peptide promoted efficient translation initiation.

60  [49] hypothesized an adaptive role for inefficient codons in the protein secretion process in

61  which the combination of efficient translation initiation and inefficient translation reduced the

62  distance between sequential ribosomes along the mRNA, leading to more efficient recycling

63  of the necessary chaperones. Other explanations for the observed increase in inefficient

64  codons include the inability of *E. coli* SRP to induce a translational pause following signal

65  peptide recognition [33, 49] and slowing down the co-translational folding of the protein, as

66  a folded protein cannot be translocated through the SecYEG translocon [32, 52, 51, 50]. If

67  signal peptides have a different CUB relative to the rest of the genome, then codon-level

68  information could be incorporated into signal peptide prediction tools.

69  In contrast [21] found no significant differences in the ribosome densities between the

70  signal peptides and the 5'-ends of nonsecretory genes in various eukaryotes. Ribosome den-

71  sities are expected to be higher in signal peptides relative to the 5'-end of nonsecretory genes

72  if selection is acting to increase translation inefficiency in the signal peptide. Additionally,

73  while both [24] and [21] examined codon usage in relation to secretion in *H. sapiens* using

74  a metric based on tAI, only [24] found results consistent with increased frequencies of in-

75  efficient codons in signal peptides. From a population genetics perspective, it is surprising

76  statistically significant results were obtained in a mammal, which usually have little adaptive

77  CUB due to their lower effective population sizes [5, 22]. More recently, [38] found codon

78  optimization of a signal peptide improved localization of the protein to the periplasm of *E.*

79  *coli*, seemingly contradicting a general role for inefficient codon usage in signal peptides. A

80  potential reason for these contradictions is the previous analyses of signal peptide codon

81  usage by [20, 21, 24, 32] did not adequately account for the effects of mutation bias and drift

82  in shaping codon usage [2, 13, 11, 12, 40, 46].

83  We re-examined CUB in signal peptides of *E. coli* using CAI, tAI, and ROC-SEMPPR

84  - a population genetics model which accounts for selection, mutation bias, and gene expres-

85  sion - to determine if selection on codon usage in signal peptides differs from the 5'-ends

86  of genes. Although we find significant differences in codon usage using CAI and tAI, we

87  present evidence these differences are due to signal peptide-specific amino acid biases and

88  differences in the gene expression distributions of genes with and without signal peptides.

89  When comparing signal peptides and the 5'-ends of genes not containing a signal peptide

90  with ROC-SEMPPR, we find signal peptide codon usage is consistent with the 5'-ends. We

91  find selection on codon usage favors the efficient codons, but the strength of selection is

92  weaker at the 5'-ends, corroborating previous analyses [9, 13, 11, 32, 35].

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

94  framework, as well as highlights potential limitations with using more common metrics such

95  as CAI for analyzing codon usage on relatively small regions of the genome. Failure to

96  account for variation in the strength of selection due to variation in gene expression can lead

97  to conflating mutation bias with selection, resulting in a misinterpretation of observed codon

98  usage patterns. Our work also illustrates the importance of considering non-adaptive forces

99  in shaping biological phenomenon before invoking adaptive explanations [14]. We believe this

100  is particularly important in the modern genomic-age when the combination of large datasets,

101  misinterpretation of p-values, and an inherent bias towards adaptationist interpretations can

102  lead to the proliferation of over-interpreted hypotheses within the biological community.

# Materials and Methods

## Signal Peptide Prediction

Signal peptides were predicted using SignalP 4.1 [31] using both the default cutoff D-score of 0.51 and a more conservative D-score of 0.75. In brief, SignalP consists of two neural networks, one for determining the amino acid sequence similarity to signal peptides and the other for identifying the most likely cleavage site. The results of both neural networks are combined into one value, called the D-score, which ranges between 0 and 1. Setting the cutoff D-score closer to 1 results in a lower false positive rate. A set of confirmed signal peptides for *E. coli* K12 MG1655 was taken from The Signal Peptide Website (http://www.signalpeptide.de/). All analyses in the main text will focus on the set of signal peptides with $D \geq 0.51$ as this set provides us with the most data; analyses of the $D > 0.75$ and set of confirmed signal peptides give similar results (see Supplementary Material).

## ROC-SEMPPR

Given a set of protein-coding genes, ROC-SEMPPR employs a Markov Chain Monte Carlo (MCMC) to estimate codon specific parameters for mutation bias $\Delta M$ and pausing times $\Delta \eta$ for each codon within a synonymous codon family (Table 1). In previous work, $\Delta \eta$ was scaled relative to the most efficient codon, which had $\Delta \eta$ and $\Delta M$ values fixed at 0. To avoid the choice of reference codon affecting our comparisons of CUB between regions, all $\Delta \eta$ values in this paper are re-scaled such that these values are centered around 0 for each amino acid. The $\Delta \eta$ values reflect the strength and direction of selection against translation inefficiency in a set of protein-coding regions (e.g. the signal peptides). A region with stronger selection against translation inefficiency will have higher $\Delta \eta$ values on average than a region with weaker selection. Similarly, a region which favors translation inefficiency would be expected to have $\Delta \eta$ values which negatively correlate with a region which favors translation efficiency.

ROC-SEMPPR also estimates an average protein production rate $\phi$ for each gene (Table

128 1). It is important to note ROC-SEMPPR is structured such that the average value of $\phi$

129 across the genome is 1. This choice of scaling means the pausing times $\Delta\eta$ represent the

130 average strength of selection relative to genetic drift for or against a given codon. We find

131 ROC-SEMPPR estimated $\phi$ values correlate well with empirical measurements of protein

132 production rates for *E. coli* (see Supplementary Methods: Assessing ROC-SEMPPR Model

133 Adequacy and Figures S1 - S2). If changes in synonymous codon usage alter the efficiency

134 at which a protein is translated, then such a change will have the largest impact on the

135 energetic costs of proteins with high production rates, making $\phi$ a more appropriate gene

136 expression metric than say, mRNA abundance or protein abundance. Thus, we use protein

137 production rates $\phi$ as our metric of gene expression. For more details on ROC-SEMPPR,

138 see [12]. Analysis of CUB with ROC-SEMPPR was performed using AnaCoDa [19].

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

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

## 139 CAI and tAI

140 Analysis of CUB was also performed using CAI [41] and tAI [7]. Both CAI and tAI quantify

141 CUB by assigning weights to the 61 sense codons. For CAI, each codon is assigned a weight

142 based on its relative frequency to its synonymous counterparts in a reference set of highly

143 expressed genes, such as ribosomal protein coding genes. The key assumption of CAI is the

144 most frequent codons in the reference set are the most efficient codons [41]. In contrast, tAI

145 assigns weights based on tRNA abundances corresponding to a codon, as well as accounting

146 for codon-anticodon interactions. The key assumption of tAI is the most efficient codons are

147 usually those with the most abundant tRNA [7].

148     CAI and tAI both range between 0 and 1. A CAI score closer to 1 represents a sequence

149 which more closely resembles the codon usage of the reference set of genes, while a tAI

150 closer to 1 indicates a sequence is more closely adapted to the genomic tRNA pool [7, 41].

151 Calculations for CAI were performed using the AnaCoDa [19], while tAI was calculated using

152 the R package tAI [6].

## Generating Datasets

154 Previous analysis of the *E. coli* genome found a set of genes with CAI values that had a

155 negative correlation with their gene expression estimates [8]. It is believed many of these

156 genes were the result of horizontal gene transfer and had not yet reached evolutionary equi-

157 librium with respect to their CUB. We repeated the analysis described in [8] on the current

158 *E. coli* K12 MG1655 genome (version 3, NC_000913.3). Briefly, correspondence analysis was

159 performed using CodonW [30], followed by clustering based on the principle axis scores using

160 the CLARA algorithm [23] in R. Our analysis was consistent with the findings of [8], reveal-

161 ing 782 genes with a CUB deviating significantly from the majority of the *E. coli* genome.

162 We will refer to this set of 782 genes as the "exogenous" component of the genome and the

163 rest of the *E. coli* genome as the "endogenous" for simplicity. All analyses presented will

164 consider only "endogenous" genes because the "exogenous" genes may violate the implicit

165 assumptions of CAI and tAI and the explicit assumptions of ROC-SEMPPR.

166 Proteins with a signal peptide were split into the signal peptide and the mature peptide

167 – the segment of the peptide chain after the signal peptide. On average, the signal peptides

168 were 23 codons long. For comparisons to the 5'-ends of nonsecretory genes – defined here

169 as those lacking a signal peptide – the first 23 codons of the nonsecretory genes were used.

170 We note the secretory genes have an average protein production rate $\phi$ approximately 10%

171 higher than that of the nonsecretory genes ($\bar{\phi} = 1.08$ and $\bar{\phi} = 0.992$,respectively, Figure S3).

172 As the strength of selection on CUB scales with protein production rate $\phi$, we created a

173 control group that eliminates differences in the distribution of $\phi$ for the nonsecretory genes

174 and signal peptide genes. Specifically, the nonsecretory genes were selected using acceptance-

175 rejection sampling to create the "pseudo-secreted proteins". In brief, acceptance-rejection

176 sampling is a procedure for sampling from a population such that its distribution of a metric

177 for one population mirrors the distribution of the same metric for another population. In

178 this case, the pseudo-secreted proteins were sampled such that the mean and variance of the

179 $\log(\phi)$ values reflected those of the genes with a signal peptide. The CUB signature of a

180 gene varies with protein production rate $\phi$; thus we can be more confident any differences

181 seen between genes with a signal peptide and pseudo-signal peptide genes are not due to

182 differences in their respective $\phi$ distributions. All pseudo-secreted proteins were split into two

183 regions we will refer to as the "pseudo-signal peptides" and the "pseudo-mature peptides"

184 (the first 23 codons and the remainder of the gene, respectively).

185 To assess the performance of CAI and tAI when comparing regions with differences in

186 the distributions of protein production rates $\phi$ and amino acid biases, simulated sequences

187 were used. Sequences based on the 5'-ends of nonsecretory genes, pseudo-signal peptides, and

188 signal peptides were simulated using the AnaCoDa package [19]. To normalize for amino acid

189 usage, sequences 23 amino acids in length were randomly generated to match the amino acid

190 frequencies of the signal peptides. The codon usage of these sequences was also simulated

191 in AnaCoDa, assuming either the $\phi$ distribution of the nonsecretory genes or the pseudo-

192 secreted proteins. All sequences were simulated using the pausing times $\Delta\eta$ and mutation

193 bias $\Delta M$ parameters estimated from the 5'-end of endogenous nonsecretory genes.

## Analysis of Codon Usage with CAI, tAI, and ROC-SEMPPR

195 We estimated protein production rates $\phi$ by fitting ROC-SEMPPR to the protein-coding

196 sequences in the *E. coli* K12 MG1655 genome. Analysis of intragenic (e.g. signal vs. mature

197 peptides) and intergenic (e.g. pseudo-signal peptides vs. real signal peptides) CUB was

198 carried out using the mixture distribution functionality available in the AnaCoDa imple-

199 mentation of ROC-SEMPPR [19]. We assumed mutation bias was consistent for the entire

200 genome; thus, we forced mutation bias $\Delta M$ parameters to be equal across the groups of

201  regions. Each group of regions (e.g. signal peptides, mature peptides, etc.) was assumed to

202  have an independent set of pausing time parameters, allowing pausing time $\Delta\eta$ estimates to

203  vary between them. $\phi$ was fixed for each region of a gene at the value estimated when the

204  model was fit to the entire protein-coding sequence. This is done for two reasons: (a) shorter

205  regions, such as the signal peptide, likely have insufficient information to accurately estimate

206  $\phi$ and (b) this guarantees our gene expression metric has the same impact on the estimates

207  of $\Delta\eta$ and $\Delta M$ for intragenic regions, such as a signal peptide and its corresponding mature

208  peptide. We note the use of empirical $\phi$ estimates in place of ROC-SEMPPR estimated $\phi$

209  did not impact our interpretations.

210      A Model-II regression was used to compare estimated pausing times $\Delta\eta$ between regions.

211  Unlike ordinary least squares, Model-II regression, or errors-in-variables regression, accounts

212  for errors in both the $x$ and $y$ variables [42]. When both variables are subject to error, which

213  is the case for the $\Delta\eta$ estimates, the use ordinary least squares leads to downwardly biased

214  parameter estimates. A Model-II regression slope $\beta = 1$ (or the $y = x$ line) will serve as

215  the null hypothesis, as this indicates both the strength and direction of selection between

216  two regions are the same. The intercept parameter was fixed at $\alpha = 0$ because the $\Delta\eta$

217  estimates are scaled such that the mean value of $\Delta\eta$ is 0. We note that when we allowed

218  the $\alpha$ parameter to vary, it was as expected, approximately 0. For more details on our use

219  of Model-II regression, see Supplementary Methods.

220      CAI and tAI were used to compare codon usage between signal peptides, 5'-ends, and

221  pseudo-signal peptides [8, 7, 41]. As recommended by [41], methionine and tryptophan were

222  not included when normalizing for the length of the gene in our calculations of CAI. Statisti-

223  cal significance was assessed using a one-tailed Welch's t-test in R [36]. R and Python scripts

224  used for this paper can be found at https://github.com/acope3/Signal_Peptide_Scripts.

# Results

Our analysis of CUB in signal peptides and the 5'-ends of nonsecretory genes using ROC-SEMPPR revealed these regions to be indistinguishable. Qualitatively, the expected codon frequencies for the 5'-ends of nonsecretory genes and the signal-peptides based on the pausing time $\Delta\eta$ and mutation bias $\Delta M$ values estimated from these regions are indistinguishable (Figure S4). Cysteine, aspartic acid, lysine, glutamine, and tyrosine are apparent exceptions, but only the 95% posterior probability intervals of cysteine and glutamine fail to overlap with $y = x$ line. When comparing the pausing times $\Delta\eta$ of signal peptides to the 5'-ends of nonsecretory genes using a Model-II regression, we find no significant difference from the $y = x$ line (slope $\beta$ 95% confidence interval: $0.923 - 1.128$, Figure 1a). To determine if differences were not detected due to underlying differences in the distributions of $\phi$, we compared $\Delta\eta$ estimated from signal peptides and pseudo-signal peptides. Again, no statistically significant difference from the $y = x$ line was found and the expected codon frequencies are similar ($\beta$ 95% confidence interval: $0.939 - 1.149$, Figure 1b and S5). Similar results are obtained using the signal peptides with a D-score greater than 0.75 or the confirmed signal peptides (Figures S6 - S7). We also see no significant result when using empirically estimated $\phi$ values ($\beta = 0.908$, 95% confidence interval: $0.671 - 1.168$, Figure S8), although these results show much more variability. The increased variability in the $\Delta\eta$ values and corresponding regression line is unsurprising given the empirically estimated $\phi$ values are subject to significant noise (Figure S2), but are, in this case, treated as error free estimates of a gene's true $\phi$ value.
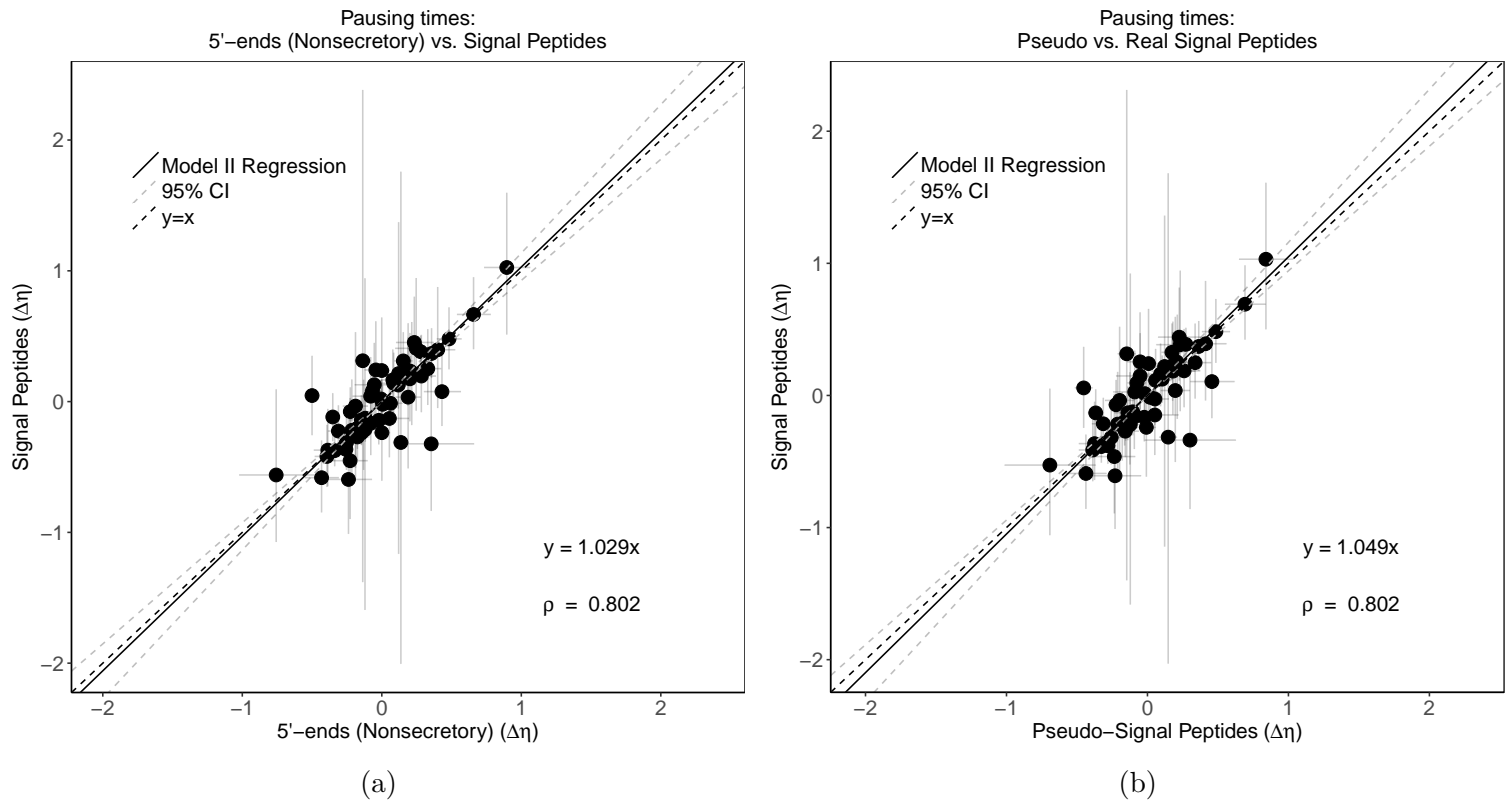
Figure 1: Comparing the pausing time estimates $\Delta\eta$ between (a) the 5%-ends of nonsecretory genes or (b) pseudo-signal peptides to signal peptides. Grey dashed lines represent the 95% confidence intervals of the regression line. Results clearly show a strong positive linear relationship ($\rho = 0.802$) between the regions and a regression line not significantly different from $y = x$.

The Model-II regression lines estimated from the mature vs. signal peptide comparison and the pseudo-mature vs. pseudo-signal peptide comparison are similar, providing further evidence the nature and magnitude of selection on codon usage in signal peptides and the 5'-ends of nonsecretory genes is indistinguishable (Figure 2). The mature vs. signal peptide comparison produces a regression line with slope $\beta = 0.480$ (95% confidence interval: 0.428 - 0.574), which is approximately 50% of the slope observed when comparing signal peptides to the 5'-ends of nonsecretory genes and pseudo-signal peptides. This indicates selection on codon usage in the mature peptides is stronger than it is in signal peptides, although the nature of selection is still *against* translation inefficiency. Similar behavior is observed when comparing the pseudo-mature vs. pseudo-signal peptide comparison ($\beta = 0.509$, 95%

255   confidence interval: 0.490 - 0.533). The slope estimate from the mature vs. signal peptide

256   comparison is not significantly different from $\beta = 0.509$ (Two-tailed Z-test, $p = 0.0682$).

257   Similar regression lines would not be expected if differences in selection on codon usage

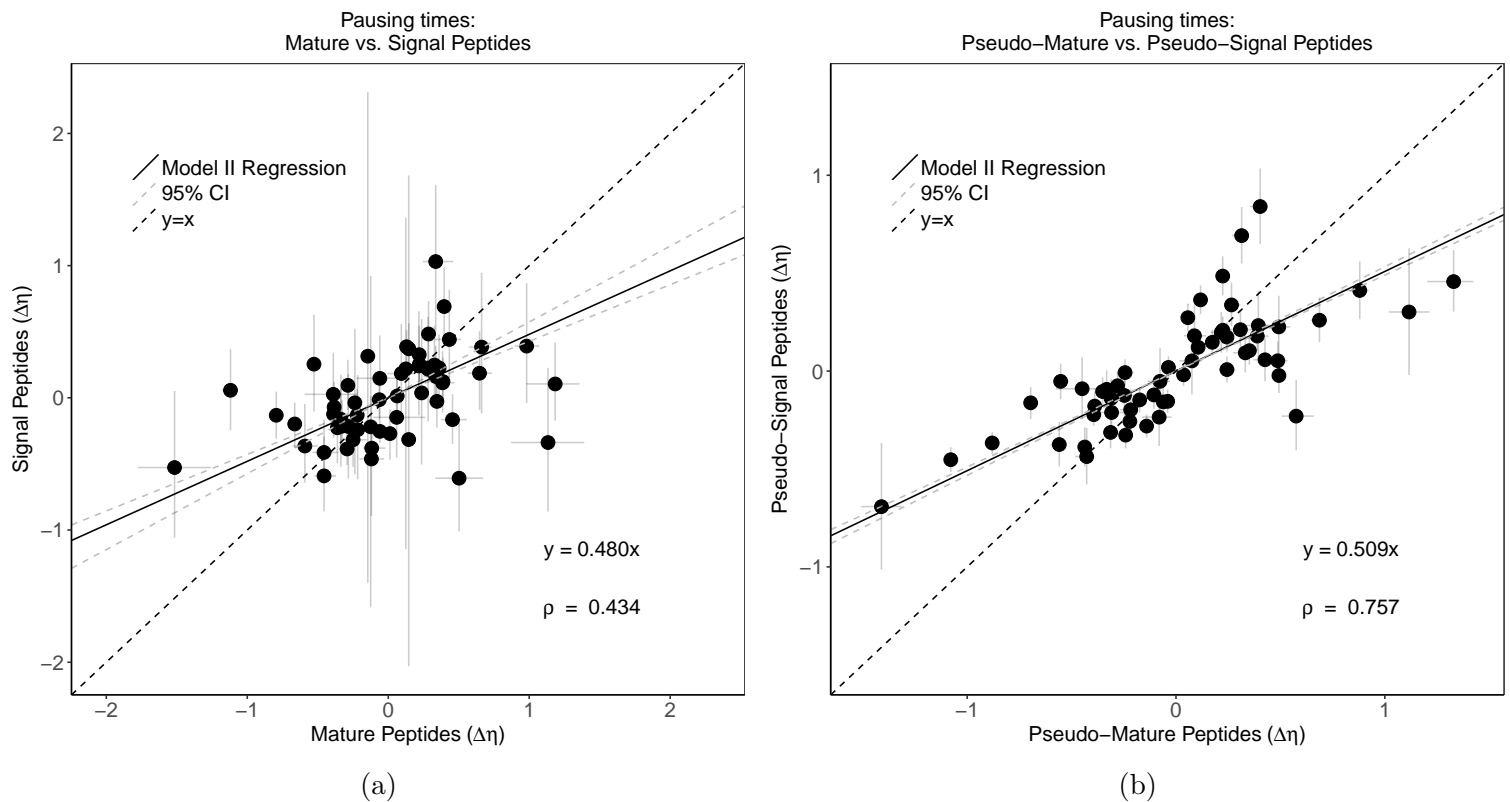258   existed between signal peptides and the pseudo-signal peptides.



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 indistinguishable from those estimated for the mature and signal peptide comparison (Likelihood Ratio test, $p = 0.562$).

259   Noting CAI and tAI do not account for the effects of gene expression, mutation bias, drift,

260   or amino acid biases, we found signal peptides have lower CAI and tAI values compared to

261   the first 23 codons of nonsecretory genes (one-tailed Welch's t-test, $p < 10^{-5}$). This was also

262   the case when looking at the pseudo-signal peptides, which normalizes for protein production

263   rates $\phi$. These results with CAI and tAI can potentially be explained by either the preferred

264 use of inefficient codons in signal peptides *or* as artifacts of amino acid biases. Signal peptides

265 have a different amino acid composition from the 5'-end due to the required physicochemical

266 properties of this region (Figure S9). We examined the robustness of tAI and CAI as a

267 means of quantifying differences in selection on codon usage when underlying differences

268 between amino acid composition and $\phi$ exists using data simulated under the same mutation

269 bias $\Delta M$ and pausing time $\Delta \eta$ parameters. When comparing simulated signal peptides to

270 simulated 5'-end of nonsecretory genes and simulated pseudo-signal peptides using CAI, the

271 simulated signal peptides are found to have a significantly lower mean CAI (Welch's t-test,

272 $p < 0.05$) 100% of the time (Figure 3A-B), despite the fact the $\Delta \eta$ and $\Delta M$ parameters used

273 to simulate these regions were the same. This suggests differences in amino acid usage and

274 not adaptation to novel selective forces, explains the lower CAI of the signal peptides.

275 When using simulated 5'-ends of nonsecretory genes which have amino acid composition

276 consistent with the signal peptides, the p-values were heavily skewed towards 1. (Figure

277 3C). This odd behavior is due to the differences in the $\phi$ distribution differences of the signal

278 peptide and nonsecretory genes. As the former has a higher mean $\phi$, the signal peptides on

279 average will have a stronger CUB after normalizing for the amino acid biases. A one-tailed

280 Welch's t-test with the alternative hypothesis being signal peptides have a lower mean CAI,

281 when in reality they likely have a larger mean CAI, would skew the p-value distribution

282 towards 1. Importantly, ROC-SEMPPR did not detect significant differences between signal

283 peptides and the 5'-ends of non-secretory genes, despite differences in the $\phi$ distributions

284 (Figure 1a). When normalizing for both amino acid usage and $\phi$, significant differences in

285 CAI are found approximately 4% of the time, which is close to the expected number of false

286 positives at the 0.05 significance level (Figure 3D). Similar results are seen when using tAI

287 (Figure S10). Our results indicate CAI and tAI are prone to inflating differences in CUB

288 between two regions when differences in $\phi$ and amino acid usage are not accounted for.

Comparing CAI of Simulated 5' Regions to Simulated Signal Peptides:
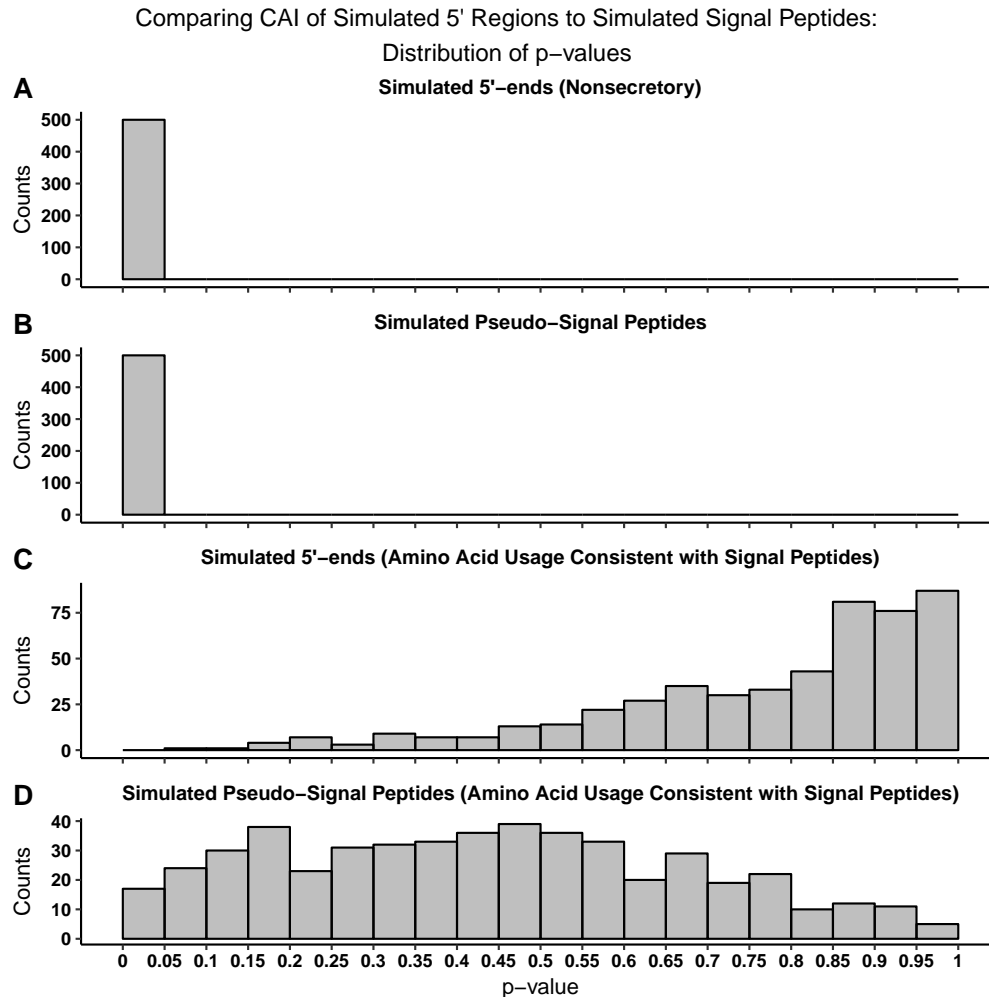Distribution of p−values



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 $\Delta M$ 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.

289 Notably, selection on codon usage near the N-terminus appears to be on average approxi-

290 mately 50% weaker than the remainder of the gene based on the slopes $\beta$. Previous analyses

291 using a variety of codon usage metrics found CUB near the 5'-end to be weaker than middle

292 sections of the gene, with these differences being attributed to selection against nonsense er-

293 rors and to maintain translation initiation efficiency by reducing mRNA secondary structure

[9, 13, 11, 16, 35, 32]. We confirm this trend using ROC-SEMPPR (Figure S11).

[54] proposed selection for translation initiation efficiency was shaping signal peptide codon usage, particularly the use of lysine codon AAA, at the second amino acid position. While AAA appears to be slightly favored in signal peptides, which is not the case in the pseudo-signal peptides, the 95% posterior probability interval overlaps with the $y = x$ line (Figure S12). If the insignificant increased usage of AAA is due to greater selection for translation initiation efficiency in signal peptides, then removing the first 3 codons when analyzing signal peptide codon usage should remove this effect. Doing so results in no change in the behavior of AAA, suggesting if there is any selection for increased AAA usage in signal peptides, it is not due to selection for increased translation initiation efficiency (Figure S13). Notably, AAA is both mutationally and selectively-favored for lysine in *E. coli*. Keeping in mind selection on CUB is weaker near the 5'-end of the genes in *E. coli*, the combination of weaker selection, mutational favorability, and a slight increase in the occurrence of lysine in signal peptides (Figure S9) likely drives up the frequency of codon AAA in signal peptides relative to the 5'-ends of nonsecretory genes.

# Discussion

In summary, we found no evidence to support the hypothesis that selection on codon usage in signal peptides and the 5'-ends of nonsecretory genes in *E. coli* using a mechanistic model of CUB which incorporates the effects of selection, mutation bias, gene expression, and amino acid usage. We find commonly employed codon usage metrics CAI and tAI produce spurious differences between signal peptides and 5'-ends of nonsecretory genes due to differences in amino acid usage and gene expression of signal peptide containing genes relative to the rest of the genome. Importantly, both amino acid usage and $\phi$ were significant confounding factors when analyzing CUB with CAI and tAI – only accounting for one of these factors still suggested significant differences between the simulated regions. Although we are not the

first to note potential issues with metrics like CAI or tAI for intragenic CUB analysis [16], our results demonstrate these metrics are insufficient for intragenic CUB analysis when these regions have drastically different amino acid usage or $\phi$ distributions, resulting in incorrect biological interpretation.

This is not to say CUB plays no role in the secretion of specific proteins. For example, experimental evidence demonstrates codon optimization of the *E. coli* maltose binding protein's (MBP) signal peptide results in a decrease in protein abundance. Evidence suggests this is due to increased targeting of the codon optimized MBP by proteases due to improper folding [52, 53]. However, CUB as a means to guide proper co-translational folding is not a phenomenon unique to proteins with a signal peptide [4, 29, 48]. Although inefficient codons might be crucial to the fold of certain secreted proteins, our results do not indicate this is any more or less so than nonsecretory genes.

Although we found no general difference in selection on codon usage between signal peptides and the 5'-ends, it is possible CUB differences exist between the chaperone-dependent and chaperone-independent mechanisms of the Sec pathway. Previous analyses revealed patterns consistent with a region of slower translation at the 5'-ends of transmembrane proteins, which are typically secreted via SRP in bacteria [26]. [10] found transmembrane proteins in *E. coli* have a higher frequency of "programmed pause sites," areas of high ribosomal density downstream from Shine-Dalgarno-like sequences, near the 5'-end. This region of higher ribosomal density was not observed in periplasmic proteins, which are normally secreted via SecA/B [26, 43]. Notably, [25] challenged the assertion that Shine-Dalgarno-like sequences are responsible for inducing translational pauses in bacteria, concluding signals previously seen were an artifact of the method for assigning ribosome occupancy along the transcript. [28] also found a consistent trend of inefficient codons 35-40 codons downstream of the SRP-binding site in various yeasts species using a modified form of the tAI. Ribosomal profiling data taken from *S. cerevisiae* provided experimental support for this hypothesis, but this analysis was limited to a small, closely-related phylogeny. Further work is needed

346  to determine the general mechanistic role, if any, of codon-induced inefficient translation in

347  SRP-dependent protein secretion, as well as to determine if any specific codon biases exists

348  for SecA/B-dependent or chaperone-independent secreted proteins.

349  We do find selection on CUB is weaker at the 5'-ends relative to later portions of the

350  gene, corroborating previous work [9, 13, 11, 16, 32, 35]. Weaker selection at the 5'-ends is

351  often attributed to selection against nonsense errors and selection against mRNA secondary

352  structure. Importantly, the advent of ribosome profiling suggested the presence of high

353  ribosomal density at the 5'-ends, often referred to as the "5'-ramp" [44]. The 5'-ramp

354  was originally thought to be the result of increased selection for slow translation at the 5'-

355  end to reduce ribosomal interference further down the transcript, but simulations suggest

356  the 5'-ramp is an artifact of short genes with high initiation rates [39]. Selection for co-

357  translational folding is also thought to shape intragenic CUB [4, 29, 48]. Further work is

358  needed to understand how these various selective forces are balanced to maintain translation

359  efficiency and efficacious protein biogenesis.

360  Although it may be tempting to explain statistically significant results in the context of

361  selection and adaptation, it is important to assert results cannot be explained by nonadap-

362  tive evolutionary forces (e.g. mutation bias and genetic drift) and/or as an artifact of some

363  other constraint on the trait of interest (e.g. amino acid biases). We are certainly not the

364  first to note the importance of considering nonadaptive explanations. Almost four decades

365  ago, [14] critiqued the propensity of evolutionary biologists to invoke natural selection and

366  adaptation without seriously considering possible nonadaptive explanations. The explosion

367  of genomic data means now, more than ever, biologists should be hesitant to adopt adapta-

368  tionist explanations to biological phenomenon without first investigating if such results could

369  be shaped by nonadaptive forces. The embrace of "big data" by biological researchers is a

370  double-edged sword: while we have the ability to investigate patterns and explore hypotheses

371  which would not have been possible 20 years ago, the indiscriminate analysis of large datasets

372  can lead to spurious, but statistically significant p-values, which are often misinterpreted as

373 both evidence of a strong effect and a small probability of the null hypothesis being true
374 [47]. The misinterpretation of p-values and a bias towards adaptationist explanations can be
375 a dangerous combination, leading to a misinterpretation of results and, in turn, misleading
376 other researchers.

377    The development of models incorporating both adaptive and nonadaptive evolutionary
378 forces will be important for understanding the selective forces shaping complex biological
379 data. In the case of the studying CUB, codon indices like CAI have long been employed,
380 but these metrics often are sensitive to and, thus, unable to disentangle the effects of amino
381 acid and mutation biases from selection. While often good proxies of gene expression, these
382 indices do not directly incorporate gene expression information into the weights estimated for
383 each codon. This could lead to further problems of conflating mutation bias with selection
384 when comparing CUB across regions. In contrast, because ROC-SEMPPR is grounded in
385 population genetics and thus, is able to decouple selection and mutation bias, it serves as
386 a more accurate and evolutionarily-grounded tool for the study of CUB. Ultimately, our
387 work further illustrates the value of employing population genetics models which include
388 nonadaptive evolutionary forces for analyzing genomic data.

# Acknowledgments

# References

[1] BENDTSEN, J. D., KIEMER, L., FAUSBØLL, A., AND BRUNAK, S. Non-classical protein secretion in bacteria. *BMC Microbiology 5*, 1 (2005), 58.

[2] BULMER, M. The effect of context on synonymous codon usage in genes with low codon usage bias. *Nucleic Acids Res. 18*, 10 (1990), 2869–2873.

[3] BURNS, D., AND BEACHAMN, I. Rare codons in *E. coli* and *S. typhimurium* signal sequences. *FEBS Letters 189* (1985), 318–324.

[4] CHANEY, J., AND CLARK, P. Roles for synonymous codon usage in protein biogenesis. *Annu. Rev. Biophysics 44* (2015), 143–166.

[5] CHARLESWORTH, B. Effective population size and patterns of molecular evolution and variation. *Nature Reviews Genetics 10* (2009), 195–205.

[6] DOS REIS, M. *tAI: The tRNA adaptation index*, 2016. R package version 0.2.

[7] DOS REIS, M., SAVVA, R., AND WERNISCH, L. Solving the riddle of codon usage preferences: a test for translational selection. *Nucleic Acids Research 32*, 17 (2004), 5036–5044.

[8] DOS REIS, M., WERNISCH, L., AND SAVVA, R. 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 (2003), 6976–6985.

[9] EYRE-WALKER, A. Synonymous codon bias is related to gene length in *Escherichia coli*: Selection for translational accuracy? *Mol. Biol. Evol. 13*, 6 (1996), 864–872.

[10] FLUMAN, N., NAVON, S., BIBI, E., AND PILPEL, Y. mrna-programmed translation pauses in the targeting of e. coli membrane proteins. *eLife 3* (2014), e03440.

[11] GILCHRIST, M. Combining models of protein translation and population genetics to predict protein production rates from codon usage patterns. *Mol. Biol. Evol. 24*, 11 (2007), 2362–2372.

[12] GILCHRIST, M., CHEN, W., SHAH, P., LANDERER, C., AND ZARETZKI, R. Estimating gene expression and codon-specific translational efficiencies, mutation biases, and selection coefficients from genomic data alone. *Genome Biology and Evolution 7* (2015), 1559–1579.

[13] GILCHRIST, M., AND WAGNER, A. A model of protein translation inducing codon bias, nonsense errors, and ribosome recyling. *Journal of Theoretical Biology 239* (2006), 417–434.

[14] GOULD, S., AND LEWONTIN, R. The spandrels of san marco and the panglossian paradigm: A critique of the adaptationist programme. *Proceedings of the Royal Society of London 205*, 1161 (1979), 581–598.

[15] GREEN, E., AND MECSAS, J. Bacterial secretion systems - an overview. *Microbiol Spectr. 4*, 1 (2016).

[16] HOCKENBERRY, A., SIRER, M., AMARAL, L., AND JEWETT, M. Quantifying position-dependent codon usage bias. *mol. Biol. Evol 31*, 7 (2014), 1880–1893.

[17] IKEMURA, T. 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* (1981), 389–409.

[18] INOUYE, S., SOBERON, X., FRANCESCHINI, T., NAKAMURA, K., ITAKURA, K., AND INOUYE, M. 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* (1982), 3438–3441.

[19] LANDERER, C., COPE, A., ZARETZKI, R., AND GILCHRIST, M. Anacoda: analyzing codon data with bayesian mixture models. *Bioinformatics* (2018), bty138.

[20] LI, Y., XIE, Z., DU, Y., ZHOU, Z., MAO, X., LV, L., AND LI, Y. The rapid evolution of signal peptides is mainly caused by relaxed selection on non-synonynous and synonymous sites. *Gene 436* (2009), 8–11.

[21] LIU, H., RAHMAN, S., MAO, Y., XU, X., AND TAO, S. Codon usage bias in 5' terminal coding sequences reveals distinct enrichment of gene functions. *Genomics 109* (2017), 506–513.

[22] LYNCH, M., ACKERMAN, M., GOUT, J., LONG, H., SUNG, W., THOMAS, W., AND FOSTER, P. Genetic drift, selection and the evolution of the mutation rate. *Nature Reviews Genetics 17* (2016), 704–714.

[23] MAECHLER, M., ROUSSEEUW, P., STRUYF, A., HUBERT, M., AND HORNIK, K. *cluster: Cluster Analysis Basics and Extensions*, 2018. R package version 2.0.7-1 — For new features, see the 'Changelog' file (in the package source).

[24] MAHLAB, S., AND LINIAL, M. Speed controls in translating secretory proteins in eukaryotes - an evolutionary perspective. *PLoS Computational Biology 10*, 1 (2014), e1003294.

[25] MOHAMMAD, F., WOOLSTENHULME, C., GREEN, R., AND BUSKIRK, A. Clarifying the translational pausing landscape in bacteria by ribosome profiling. *Cell Reports 14* (2016), 686–694.

[26] NATALE, P., BRUSER, T., AND DRIESSEN, A. Sec- and tat-mediated protein secretion across the bacterial cytoplasmic membrane—distinct translocases and mechanisms. *Biochimica et Biophysica Acta 1778* (2008), 1735–1756.

[27] Nesmeyanova, M., Karamyshev, A., Karamysheva, Z., Kalinin, A., Ksenzenko, V., and Kajava, A. 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* (1997), 203–207.

[28] Pechmann, S., Chartron, J., and Frydman, J. Local slowdown of translation by nonoptimal codons promotes nascent-chain recognition by srp *in vivo. Nature Structural and Molecular Biology 21*, 12 (2014), 1100–1105.

[29] Pechmann, S., and Frydman, J. Evolutionary conservation of codon optimality reveals hidden signatures of cotranslational folding. *Nature Structural and Molecular Biology 20*, 2 (2013), 237–243.

[30] Peden, J. *Analysis of Codon Usage.* PhD thesis, University of Nottingham, 8 1999.

[31] Petersen, T., Brunak, S., von Heijne, G., and Nielsen, H. Signalp 4.0: discriminating signal peptides from transmembrane regions. *Nature Methods 8*, 10 (2011), 785–786.

[32] Power, P., Jones, R., Beacham, I., Bucholtz, C., and Jennings, M. Whole genome analysis reveals a high incidence of non-optimal codons in secretory signal sequences of *Escherichia coli. Biochemical and Biophysical Research Communications 322* (2004), 1038–1044.

[33] Powers, T., and Walter, P. Co-translational protein targeting catalyzed by the *Escherichia coli* signal recognition particle and its receptor. *The EMBO Journal 16*, 16 (1997), 4880–4886.

[34] Puziss, J., Fikes, J., and Bassford, P. Analysis of mutational alterations in the hydrophilic segment of the maltose-binding protein signal peptide. *Journal of Bacteriology 171* (1989), 2303–2311.

[35] QIN, H., WU, W., KREITMAN, J. C. M., AND LI, W. Intragenic spatial patterns of codon usage bias in prokaryotic and eukaryotic genomes. *Genetics 168* (2004), 2245–2260.

[36] R CORE TEAM. *R: A Language and Environment for Statistical Computing.* R Foundation for Statistical Computing, Vienna, Austria, 2018.

[37] SAIER, M. Protein secretion systems in gram-negative bacteria. *Microbe 1*, 9 (2006), 414–419.

[38] SAMANT, S., GUPTA, G., KARTHIKEYAN, S., AMD A. NAIR, S. H., SAMBASIVAM, G., AND SUKUMARAN, S. Effect of codon-optimized *E. coli* signal peptides on recombinant *Bacillus stearothermophilus* maltogenic amylase periplasmic localization, yield and activity. *J. Ind. Microbial Biotechnol 41* (2014), 1435–1442.

[39] SHAH, P., DING, Y., NIEMCZYK, M., KUDLA, G., AND PLOTKIN, J. Rate-limiting steps in yeast protein translation. *Cell 153* (2013), 1589–1601.

[40] SHAH, P., AND GILCHRIST, M. Explaining complex codon usage patterns with selection for translational efficiency, mutation bias, and genetic drift. *PNAS 108*, 25 (2011), 10231–10236.

[41] SHARP, P., AND LI, W. The codon adaptation index - a measure of directional synonymous codon usage bias, and its potential applications. *Nucl. Acids Research 15*, 3 (1987), 1281–1295.

[42] SOKAL, R., AND ROHLF, F. *Biometry - The Principles and Practices of Statistics in Biological Research*, 3rd ed. W.H. Freeman, New York, 1995.

[43] TSIRIGOTAKI, A., GEYTER, J. D., SOSTARIC, N., ECONOMOU, A., AND KARAMANOU, S. Protein export through the bacterial sec pathway. *Nature Reviews: Microbiology 15* (2017), 21–36.

[44] TULLER, T., CARMI, A., VESTSIGIAN, K., NAVON, S., DORFAN, Y., ZABORSKE, J., PAN, T., DAHAN, O., FURMAN, I., AND PILPEP, Y. An evolutionarily conserved mechanism for controlling the efficiency of protein translation. *Cell 141* (2010), 344–354.

[45] VLASUK, G., INOUYE, S., ITO, H., ITAKURA, K., AND INOUYE, M. 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* (1983), 7141–7148.

[46] WALLACE, E., AIROLDI, E., AND DRUMMOND, D. Estimating selection on synonymus codon usage from noisy experimental data. *Molecular Biology and Evolution 30*, 6 (2013), 1438–1453.

[47] WASSERSTEIN, R., AND LAZAR, N. The asa's statement on p-values: Context, process, and purpose. *The American Statistician 70*, 2 (2016), 129–133.

[48] YU, C., DANG, Y., ZHOU, Z., WU, C., ZHAO, F., SACHS, M., AND LIU, Y. Codon usage influences the local rate of translation elongation to regulate co-translational protein folding. *Molecular Cell 59* (2015), 744–754.

[49] ZALUCKI, Y., BEACHAM, I., AND JENNINGS, M. Biased codon usage in signal peptides: a role in protein export. *Trends in Microbiology 17*, 4 (2009), 146–150.

[50] ZALUCKI, Y., BEACHAM, I., AND JENNINGS, M. Coupling between codon usage, translation and protein export in *Escherichia coli. Biotechnology Journal 6* (2011), 660–667.

[51] ZALUCKI, Y., GITTINS, K., AND JENNINGS, M. Secretory signal sequence non-optimal codons are required for expression and export of $\beta$-lactamase. *Biochemical and Biophysical Research Communications 366* (2008), 135–141.

[52] ZALUCKI, Y., AND JENNINGS, M. Experimental confirmation of a key role for non-optimal codons in protein export. *Biochemical and Biophysical Research Communications 355* (2007), 143–148.

[53] ZALUCKI, Y., JONES, C., NG, P., SCHULZ, B., AND JENNINGS, M. Signal sequence non-optimal codons are required for the correct folding of mature maltose binding protein. *Biochimica et Biophysica Acta 1798* (2010), 1244–1249.

[54] ZALUCKI, Y., POWER, P., AND JENNINGS, M. Selection for efficient translation initiation biases codon usage at the second amino acid position in secretory proteins. *Nucleic Acids Research* (2007), 1–7.

[55] ZALUCKI, Y., SHAFER, W., AND JENNINGS, M. Directed evolution of effeicient secretion in the srp-dependent export of tolb. *Biochemica et Biophysica Acta 1808* (2011), 2544–2550.