

Specific Codons Control Cellular Resources and Fitness

Aaron Love

Tufts University

Nikhil Nair (✉ nikhil.nair@tufts.edu)

Tufts University <https://orcid.org/0000-0001-7737-1385>

Article

Keywords: Codon bias, resource competition, tRNA, protein expression, genetic burden, translation

Posted Date: November 3rd, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-2198914/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: There is **NO** Competing Interest.

1 **Specific Codons Control Cellular Resources and Fitness**

2 Aaron M. Love^{a,b}, Nikhil U. Nair^{b,*}

3 ^a Manus Bio, Cambridge, MA 02138

4 ^b Department of Chemical & Biological Engineering, Tufts University, Medford, MA 02155

5 * Corresponding author: nikhil.nair@tufts.edu; @nair_lab

6

7

8 **KEYWORDS:** Codon bias; resource competition; tRNA; protein expression; genetic burden; translation

9

10

11

12

GLOSSARY

- 13 ▪ **RSCU:** Relative synonymous codon usage
- 14 ▪ **W_{ij}:** Relative adaptiveness (weight)
- 15 ▪ **CAI:** Codon adaptation index
- 16 ▪ **ENC:** Effective number of codons
- 17 ▪ **CUB:** Codon usage bias
- 18 ▪ **TAI:** tRNA adaptation index
- 19 ▪ **sTAI:** Species-specific tRNA adaptation index
- 20 ▪ **nTE:** Normalized translational efficiency
- 21 ▪ **RFM:** Ribosome flow model
- 22 ▪ **CFP:** Cyan fluorescent protein
- 23 ▪ **YFP:** Yellow fluorescent protein
- 24 ▪ **TxTL:** in vitro transcription-translation
- 25 ▪ **UTR:** Untranslated region
- 26 ▪ **AUC:** Area under the curve
- 27 ▪ **Fitness:** Performance of induced culture ÷ Performance of uninduced culture
- 28 ▪ **Growth Fitness:** AUC of growth curve (induced) ÷ AUC of growth curve (uninduced)
- 29 ▪ **Co-Expression Fitness:** AUC of YFP fluorescence (with induced CFP or mCherry) ÷ AUC of YFP fluorescence
- 30 (with uninduced CFP or mCherry)
- 31 ▪ **Expression Level:** AUC of fluorescence from induced over-expressed protein (CFP or mCherry)
- 32 ▪ **CHI (χ):** Codon harmony index
- 33 ▪ **MFE:** Mean free energy

34

35 **ABSTRACT:**

36

37 There is a degeneracy in codons – but they are not equivalent. While there is an understanding that codon use is
38 unequal in native genes, there is less knowledge of how this usage bias modulates the supply and demand of
39 protein translation resources. Here we investigate how the partitioning of microbial translational resources,
40 specifically through allocation of tRNA by incorporating dissimilar codon usage bias, can drastically alter expression
41 of proteins and reduce the burden on the host resources. By isolating individual codons experimentally, we find
42 heterologous gene expression can *trans*-regulate fitness of the host and other heterologous genes. Interestingly,
43 specific codons drive profitable or catastrophic phenotypic outcomes. We correlate codon usage patterns with
44 genetic fitness and empirically derive a novel coding scheme for multi-gene expression called Codon Harmony
45 Index (CHI, χ). CHI enables the design of harmonious multi-gene expression systems while avoiding catastrophic
46 cellular burden.

47

48 **INTRODUCTION:**

49

50 The genetic code is degenerate with 61 codons and only 20 amino acids, creating an astronomically high level of
51 mRNA sequence space for most protein coding genes. However, it is well accepted that synonymous codons are
52 not equivalent^{1,2}, as numerous reports of *cis* and *trans* effects have been documented³⁻¹¹ – from mRNA structure
53 and co-translational protein folding¹²⁻¹⁴ to tRNA and ribosome competition¹⁵⁻¹⁷. Re-coding proteins typically
54 proceeds through use of a codon adaptation index (CAI), which enables a gene to assume the codon usage bias
55 (CUB) of a reference set, often a set of highly expressed genes¹⁸. This strategy may generally correlate CUB with
56 protein expression, but it ignores the role CUB can play in partitioning translational resources such as tRNA and
57 ribosomes. Several recent studies have demonstrated the ability of heterologous genetic CUB to *trans*-regulate
58 host gene expression through translational resource competition^{19,20}, but there is little understanding of how
59 specific CUB alters host fitness given that cellular resources are invariably limited. Re-coding strategies such as the
60 tRNA adaptation index (tAI)^{7,21} and normalized translational efficiency (nTE)⁶ are attempts to address tRNA related
61 translational supply-demand constraints, but they are limited by how predictive natural CUB and/or tRNA levels
62 are for recombinant protein expression.

63 It is particularly important to consider translational resource competition in the context of multi-gene expression
64 (e.g., in the case of metabolic engineering and synthetic biology), where the objective is often for global organism
65 fitness in addition to high protein expression, and tradeoffs in protein expression can be highly consequential for
66 pathway or genetic circuit function and robustness²². This area is currently underexplored, as most studies to date
67 focus on feedback control mechanisms^{23,24}, resource partitioning^{25,26}, or attempt to draw inferences about
68 elongation in larger genes from libraries limited to the 5' sequence of a reporter^{27,28}, and experiments that do not
69 isolate translation elongation from initiation effects¹⁰. As cellular engineering becomes increasingly complex,
70 genetic resource competition can unravel designs and lead to unpredictable and undesirable phenotypes. While a
71 role for CUB in the partitioning of cellular resources has been reported²⁹, identification of specific codons that
72 present excess translational capacity could provide a novel avenue for harnessing underutilized resources that are
73 insofar ignored.

74 In this study, we systematically isolate the role of codon choices during translational elongation and identify
75 supply-demand constraints imposed on tRNA and ribosomal resources in *E. coli*. We demonstrate that tRNA
76 limitations lead to competition between overexpressed genes as well as with the host's demands. Select codons
77 over-represented in native highly expressed genes are found to cause severe fitness costs when present in
78 overexpressed protein sequences. While the traditional method of codon-optimization through maximizing CAI
79 may promote use of these codons, our data reveal their demand and supply are delicately balanced. We define a
80 new metric called "Codon Harmony Index" (CHI, χ) that quantitatively ranks codons by their capacity to remain
81 orthogonal to host demands. We also posit using this metric as a new codon optimization scheme to mitigate
82 competition with host demands and avoid growth defects. Genes characterized by high scores on this metric
83 scheme demonstrate relatively high expression while minimizing the burden on the host cells, allowing effective
84 multigene expression and cellular growth.

85

86

87 **RESULTS:**

88

89 **Fitness costs are incurred due to translation elongation limitation.**

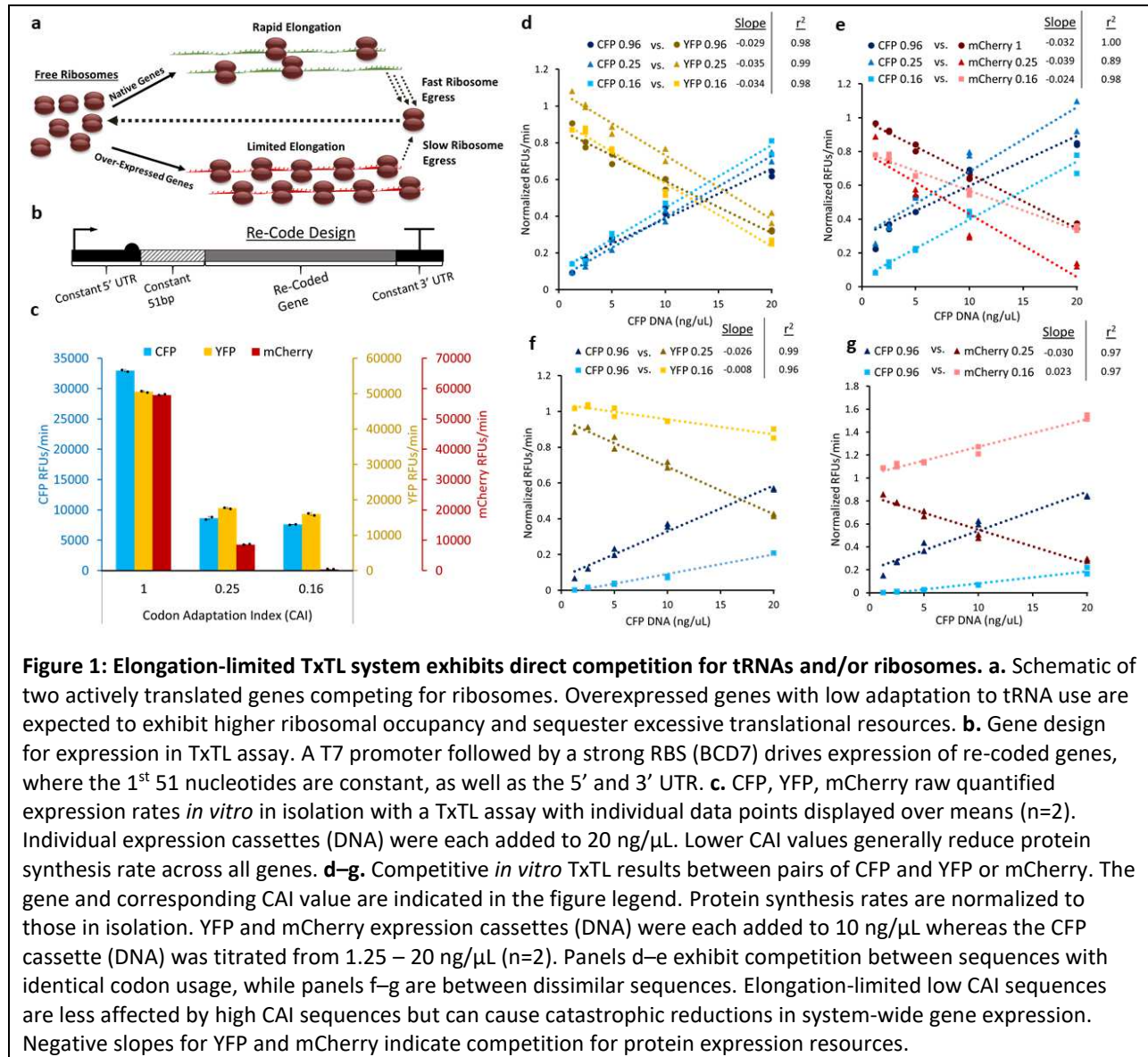
90

91 Genetic burden is frequently observed in microbial systems as a growth defect upon the overexpression of
92 recombinant proteins²⁴. While the cause of this effect varies, it is often attributed to resource competition at the
93 level of mRNA translation³⁰. In a fast-growing culture of *E. coli*, the availability of free ribosomes can limit mRNA
94 translation, especially in a system with overexpressed protein³¹ (**Figure 1a**). Elongation speed determines the rate
95 at which free ribosomes are made available, hence sub-optimal mRNA transcripts that are poorly translated have
96 higher ribosome occupancy. Such elongation limited mRNA sequences will sequester more ribosomes and return
97 them to the free pool at a slower rate, thus reducing ribosome availability. Translational resource competition has
98 been modeled in several ways³², including the ribosome flow model (RFM)³³, which can be useful in examining
99 translation rate as a function of elongation time that varies depending on the supply and demand of tRNAs in the
100 cell. Applying a previously developed RFM³⁴ to the model gene cyan and yellow fluorescent proteins (CFP and YFP
101 respectively) with high or low CAI values (where CAI is in reference to highly expressed *E. coli* genes) illustrates the
102 increase in mRNA ribosome occupancy that occurs when codons with longer elongation times³⁵ are used, and
103 indicates that elongation-limited sequences are less sensitive to changes in the rate of translation initiation (**Figure**
104 **S1**).

105 We first sought to investigate the impact of translation elongation resource competition using an *in vitro*
106 transcription-translation (TxTL) model. A significant challenge to investigating translational resource competition is
107 the difficulty in isolating any single sequence parameter experimentally, as any synonymous mutation can have a
108 multitude of effects on initiation, elongation, and mRNA structure². A TxTL system allows for better physical
109 control over the genetic expression environment by holding available resources (e.g., ribosomes, tRNAs, aminoacyl
110 tRNA synthetases, RNA polymerase etc.) constant, and allowing precise titration of genes of interest in the
111 reaction. We developed an assay for elongation limitation by leveraging the unique amino acid sequence similarity
112 between CFP and YFP derived from a super-folder green fluorescent protein³⁶, which only differ by 2 amino acids³⁷,
113 thus eliminating variability in protein structure and amino acid demand. The CFP-YFP pair permits the interrogation
114 of competition between various sequence designs using effectively identical proteins, which should also be less
115 susceptible to variation in co-translational protein folding due to their high stability. We also include mCherry in
116 the study, which is <30% identical to CFP-YFP and serves as a comparison point to find trends independent of
117 amino acid sequence (see supplementary data for sequences). The TxTL kit is based off the *E. coli* MRE600 strain,
118 which has a nearly identical CUB as K12 MG1655 and is therefore assumed to be a good proxy for the tRNA profile
119 in a K12 strain used subsequently (**Figure S2**). Reactions were driven by a T7 promoter using a bicistronic domain
120 (BCD) in place of a traditional ribosome binding site to minimize interactions between the 5' untranslated region
121 and gene of interest that could lead to differential expression³⁸. To further isolate translation elongation as the
122 primary variable in sequence design, we chose to keep the 5' and 3' untranslated regions (UTRs) as well as the first
123 51 base pairs (17 codons) constant to mitigate any effect sequence changes may have on translation initiation
124 (**Figure 1b**).

125 Utilizing the idealized TxTL competition assay, we evaluated baseline expression rates from CFP, YFP, and mCherry
126 re-codes with extreme CAI values (0.96, 0.25, or 0.16) (**Figure 1c**). We find that identical sequence pairs for CFP-
127 YFP behave very similarly in terms of relative expression, and that protein expression rates for CFP, YFP, and
128 mCherry correspond well with CAI value. This supports that TxTL recapitulates translation elongation limitation –
129 i.e., genes with lower CAI that use lower abundance tRNAs show lower protein synthesis rates. Next, we examined
130 competition between different pairs of genes. As in the RFM, we expected elongation-limited sequences with
131 lower CAI to disrupt expression of other genes through the sequestration of free ribosomes. We titrated CFP
132 template DNA against constant YFP or mCherry DNA using re-codes with either very high or very low CAI (**Figure**
133 **1d–g**). For instances of two identically re-coded sequences with any CAI tested, YFP and mCherry synthesis rates
134 are inversely correlated with CFP DNA concentration (**Figure 1d–e**), irrespective of their baseline expression,
135 indicating strong competition for limiting resources (i.e., tRNA). This indicates that while an excess protein
136 synthesis capacity exists in the TxTL system, sequences with lower CAI are still resource-limited, likely due to lower
137 availability of tRNA.

138



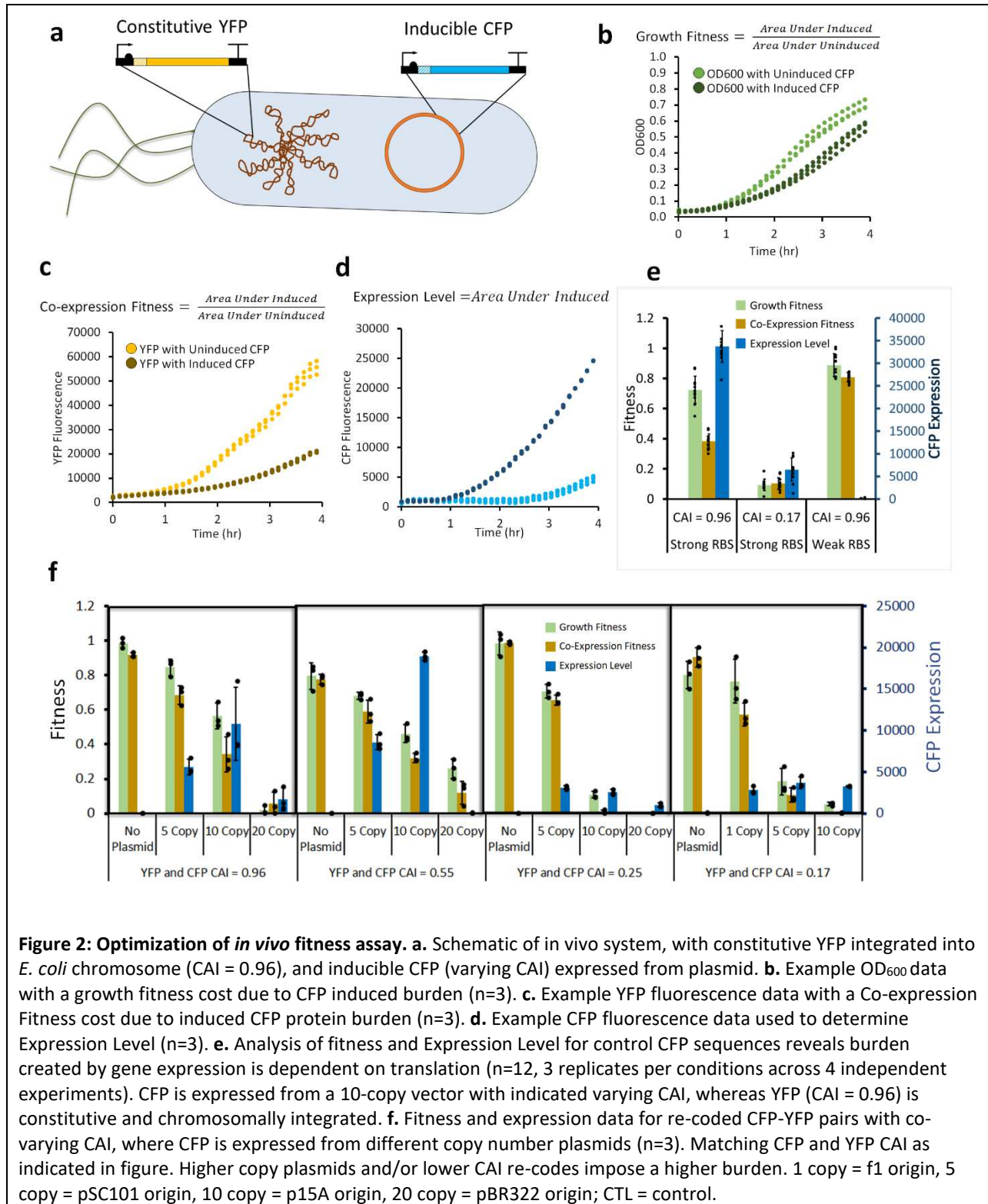
140

141 **Figure 1: Elongation-limited TxTL system exhibits direct competition for tRNAs and/or ribosomes.** **a.** Schematic of
 142 two actively translated genes competing for ribosomes. Overexpressed genes with low adaptation to tRNA use are
 143 expected to exhibit higher ribosomal occupancy and sequester excessive translational resources. **b.** Gene design
 144 for expression in TxTL assay. A T7 promoter followed by a strong RBS (BCD7) drives expression of re-coded genes,
 145 where the 1st 51 nucleotides are constant, as well as the 5' and 3' UTR. **c.** CFP, YFP, mCherry raw quantified
 146 expression rates *in vitro* in isolation with a TxTL assay with individual data points displayed over means (n=2).
 147 Individual expression cassettes (DNA) were each added to 20 ng/ μ L. Lower CAI values generally reduce protein
 148 synthesis rate across all genes. **d–g.** Competitive *in vitro* TxTL results between pairs of CFP and YFP or mCherry. The
 149 gene and corresponding CAI value are indicated in the figure legend. Protein synthesis rates are normalized to
 150 those in isolation. YFP and mCherry expression cassettes (DNA) were each added to 10 ng/ μ L whereas the CFP
 151 cassette (DNA) was titrated from 1.25 – 20 ng/ μ L (n=2). Panels d–e exhibit competition between sequences with
 152 identical codon usage, while panels f–g are between dissimilar sequences. Elongation-limited low CAI sequences
 153 are less affected by high CAI sequences but can cause catastrophic reductions in system-wide gene expression.
 154 Negative slopes for YFP and mCherry indicate competition for protein expression resources.

155

156 More interesting observations are seen when dissimilar CAI re-codes are under competition upon co-expression
 157 (**Figure 1e–f**). Low CAI YFP and mCherry synthesis rates are not very sensitive to increasing resource demand by
 158 high CAI CFP synthesis. Conversely, the relative CFP expression is much lower than we observed either in isolation
 159 or when competing with a high CAI sequence. The observed results appear to be consistent across different
 160 sequence pairs, indicating that this phenomenon is independent of protein sequence. When examined in the
 161 context of an RFM, we deduce that the rare codon enriched YFP and mCherry sequences sequester ribosomes to
 162 such a degree that even excess CFP template DNA does not yield high synthesis rates. On the other hand, YFP and
 163 mCherry are not affected due to severe elongation limitation. This model is further supported by our observation
 164 that YFP and mCherry rates are reduced when competing with similarly re-coded low CAI CFP sequences, which is a
 165 likely consequence of competition for scarce tRNAs. Overall, our data indicates that proteins coded with similar CAI
 166 (high or low) are strongly competitive due to demand for the same tRNA pool. Conversely, genes coded under
 167 distinct CAI regimes are constrained by the availability free ribosomes, which are in turn limited due to
 168 slow/stalled translation from scarce tRNA resources. Our TxTL data strongly support the argument that translation

169 elongation limitation could play an important role in cellular resource competition and highlights the impact to
 170 global translational resources (e.g., free ribosomes, tRNA) in multigene expression environments.
 171



172 **Figure 2: Optimization of *in vivo* fitness assay.** **a.** Schematic of *in vivo* system, with constitutive YFP integrated into
 173 *E. coli* chromosome (CAI = 0.96), and inducible CFP (varying CAI) expressed from plasmid. **b.** Example OD₆₀₀ data
 174 with a growth fitness cost due to CFP induced burden (n=3). **c.** Example YFP fluorescence data with a Co-expression
 175 Fitness cost due to induced CFP protein burden (n=3). **d.** Example CFP fluorescence data used to determine
 176 Expression Level (n=3). **e.** Analysis of fitness and Expression Level for control CFP sequences reveals burden
 177 created by gene expression is dependent on translation (n=12, 3 replicates per conditions across 4 independent
 178 experiments). CFP is expressed from a 10-copy vector with indicated varying CAI, whereas YFP (CAI = 0.96) is
 179 constitutive and chromosomally integrated. **f.** Fitness and expression data for re-coded CFP-YFP pairs with co-
 180 varying CAI, where CFP is expressed from different copy number plasmids (n=3). Matching CFP and YFP CAI as
 181 indicated in figure. Higher copy plasmids and/or lower CAI re-codes impose a higher burden. 1 copy = f1 origin, 5
 182 copy = pSC101 origin, 10 copy = p15A origin, 20 copy = pBR322 origin; CTL = control.

185
186 We next set out to optimize an in vivo system for *E. coli* expression to efficiently interrogate the effect alternative
187 recoding designs have on gene expression and host fitness. Our system generally consists of a strong constitutively
188 expressed YFP reporter gene (CAI = 0.96) integrated into the *E. coli* chromosome paired with an inducible CFP on a
189 plasmid driven by the inducible promoter P_{trc} with a strong RBS (**Figure 2a**). As before, we held the 1st 51 bases and
190 the 5' and 3' UTRs constant for all recodes. Cells grown in rich medium with a common pre-culture were passaged
191 under inducing or non-inducing conditions. The area under the curve (AUC) is used to measure each of the 3
192 signals (growth and 2 fluorescent proteins), which captures the aggregate effects of different lag phases and
193 expression rates (**Figure S3**). We define fitness as the ratio of AUC induced vs. uninduced, which ranges from 0 to 1
194 for low and high fitness, respectively (or conversely, high and low burden). Fitness can be in terms of Growth
195 Fitness based on OD₆₀₀, or Co-expression Fitness based on YFP fluorescence (chromosomal reporter), while
196 Expression Level is based on CFP or mCherry fluorescence (i.e., the overexpressed protein) (**Figure 2b-d**). We
197 generally observed a reduction in both growth and YFP fluorescence upon CFP induction. Examining several
198 controls expressed from a p15A origin in **Figure 2e**, a “codon-optimized” high CAI CFP gene expresses well but
199 elicits a significant fitness cost in terms of Growth Fitness and Co-expression Fitness. For a CFP recoded with rare
200 codons, the result is catastrophic, and cultures are unable to grow at all. The effect also seems mediated by
201 translation (not transcription) since the codon-optimized CFP with a very weak RBS, but intact promoter neither
202 synthesizes protein nor demonstrates much fitness cost. Upon varying plasmid copy number with several pairs of
203 CFP and YFP with different CAI levels, we found that fitness costs (Co-expression and Growth) were strongly
204 dependent on copy number that is further exacerbated by low CAI (**Figure 2f**). Interestingly, the CFP Expression
205 Level was not very correlated with CAI nor copy number. Based on these results, we picked the 10-copy vector
206 (p15A origin) with the YFP CAI = 0.96 reporter as the platform for further studies to investigate re-coding schemes
207 that may reduce fitness costs.

208

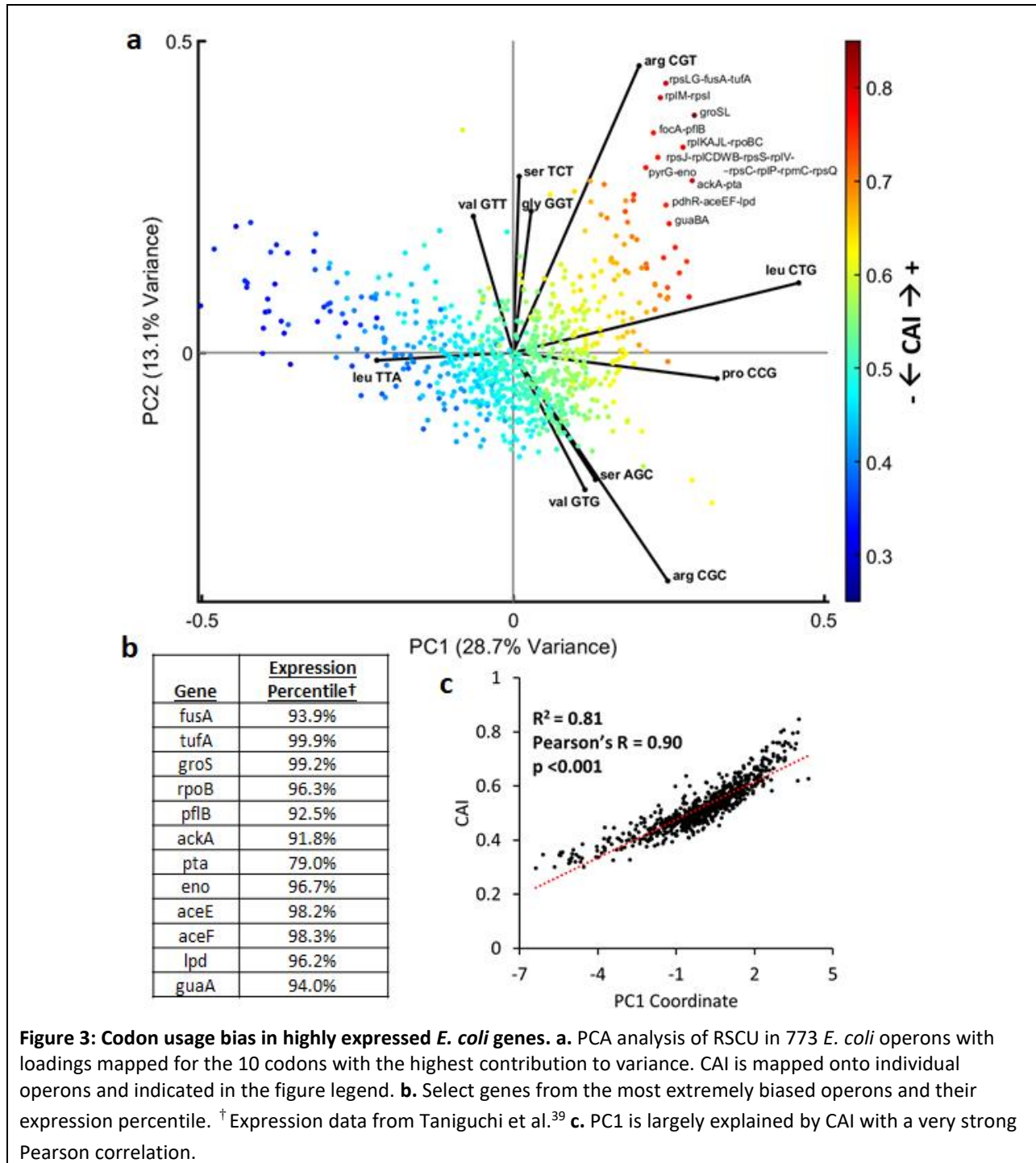
209 **Systematic analysis of codon use reveals supply and demand constraints in tRNA resources.**

210

211 Prior to designing novel re-coded genes that moderate translation elongation resources, we first investigated CUB
212 in the *E. coli* transcriptome. CAI calculations are typically based on the natural CUB in highly expressed genes. CUB
213 can be represented as a 64-dimensional space (total number of codons) using RSCU values (observed vs. expected
214 frequency) for each protein coding gene. Initial analysis revealed that groups of genes within the *E. coli*
215 transcriptome cluster according to distinct CUB schemes (**Figure S4**). We focused on a consolidated set of this
216 sequence space by analyzing all operons with at least 2 protein coding genes, given that functionally related genes
217 that naturally cluster have similar CUB (**Figure S5**). The resulting 64 dimensions of codon usage across 773 operons
218 can be represented in 2 dimensions accounting for 41.2% of total variance (**Figure S6**) using principal component
219 analysis (PCA) as shown in **Figure 3a**. The loading vectors mapped onto the plot represent the 10 codons that
220 contribute most significantly to codon bias across the 773 operons.

221 This analysis captures the CUB naturally observed in the *E. coli* transcriptome and highlights a positive correlation
222 between CAI and expression. This is expected because here CAI is calculated by optimizing towards CUB in highly
223 expressed genes¹⁸ (see methods) (**Figure S7**). Consistent with previous studies, we corroborate that genes in the
224 most extremely biased CUB space are some of the most highly expressed genes in the *E. coli* proteome that often
225 serve essential functions (**Figure 3b**). The natural bias leading to the CAI scale is very well explained by PC1 (**Figure**
226 **3c**). Despite the apparent correlation between CAI and expression, studies have reported that CAI often does not
227 predict higher gene expression¹⁰. Importantly, the CAI paradigm of re-coding proteins to match the CUB of highly
228 expressed genes ignores potential resource competition that can occur at the tRNA level. For 18 of 20 amino acids,
229 multiple codons exist, and 10 of 18 of those can be coded to use different tRNAs in *E. coli* K12 MG1655 (**Figure S8**).
230 Upon examining the PCA loadings, there are clearly particular codons that are very overrepresented in highly
231 expressed proteins (e.g., arg CGT, leu CTG, and pro CCG). For such high-demand codons, using alternative
232 codon/tRNA pairs, or even codons that recruit tRNAs with weaker affinity, have the potential to reduce translation
233 elongation-based resource competition between overexpressed proteins and native essential and/or highly
234 expressed genes.

235



236
 237 **Figure 3: Codon usage bias in highly expressed *E. coli* genes.** **a.** PCA analysis of RSCU in 773 *E. coli* operons with
 238 loadings mapped for the 10 codons with the highest contribution to variance. CAI is mapped onto individual
 239 operons and indicated in the figure legend. **b.** Select genes from the most extremely biased operons and their
 240 expression percentile. † Expression data from Taniguchi et al.³⁹ **c.** PC1 is largely explained by CAI with a very strong
 241 Pearson correlation.

242
 243
 244 Using our optimized in vivo assay, we sought to experimentally determine the contribution of individual codons to
 245 gene Expression Level and Co-Expression Fitness. The synonymous codon sequence space that could be explored in
 246 even a small gene such as CFP is experimentally intractable. Holding the first 51 bp constant and co-varying all
 247 possible synonymous codons would produce a massive library size of 1.8×10^{104} . While a more constrained codon
 248 library is possible, we chose a focused experimental approach by interrogating individual codon contribution to
 249 gene Expression Level and Co-Expression Fitness. Starting with a CFP or mCherry sequence having a high CAI (0.96

250 – 1.0) and using a single codon for each amino acid where the effective number of codons (ENC) = 20 (for details
251 on ENC, see methods), for each amino acid we re-coded every instance to another synonymous codon, resulting in
252 a total of 41 possible re-coded sequences (64 possible codons – 20 high CAI codons already in use – 3 stop codons
253 not changed) (**Figure 4a**). Results were normalized in terms of both Expression Level and Co-expression Fitness
254 (defined in **Figure 2b**) relative to the high CAI parent control (**Figure 4b**) and indicate wide ranging benefits or
255 costs. In several instances, alternative codons provide a significant improvement in Co-Expression Fitness across
256 both mCherry and CFP. Variations in phenotypes could in part be due to different amino acid composition between
257 mCherry and CFP, as the number of re-coded amino acids was not held constant between genes (**Figure S9**). We
258 chose to re-code all instances of each amino acid so as not to limit the number of altered codons to the amino acid
259 with the fewest instances. Most of the re-codes do not improve expression (**Figure 4c**), which is expected since
260 they were derived from (and normalized to) high CAI sequences that emulate highly expressed genes. CFP and
261 mCherry re-codes are also less consistent in Expression Level than Co-Expression Fitness, reflecting a higher degree
262 of variability between genes in *cis* compared *trans* effects. Notably, there are several alternative codons for
263 leucine, proline, and one for arginine, which robustly improve Co-expression Fitness, suggesting that dissimilar
264 codon use could be a means to generally reducing heterologous gene burden. Expression Level and Co-expression
265 Fitness do not correlate well for mCherry or CFP re-codes (**Figure 4d–e**), indicating that while there may be general
266 tradeoffs between expression and fitness, there are many instances where specific codon/tRNA pairs possess
267 excess translational capacity.

268
269

270 **Novel recoding scheme yields genes with robustly improved fitness.**

271

272 Next, we developed a new recoding index derived from Co-expression Fitness values for individual codons in
273 **Figure 4b**. We chose to focus on fitness rather than expression since our primary aim was to investigate how re-
274 coding schemes can modulate resource competition during translation elongation. To convert the Co-expression
275 Fitness data for CFP and mCherry re-codes into generalized codon weights, we took the Euclidean distance from
276 the origin to the coordinates of each data point shown in **Figure 4b** as a raw score for each sequence, where each
277 parent codon held a normalized coordinate value of (1,1). Similar to calculating CAI, relative adaptiveness (W_i)
278 scores were then determined by normalizing the raw weights from each amino acid codon set to the codon with
279 the highest fitness (see methods and **Data S1**). We refer to this new metric as the Codon Harmony Index (CHI or χ).

280

281 A comparative analysis between CUB in the overall *E. coli* genome, CAI (using highly expressed genes as a
282 reference), and χ reveals that χ favors very different codon use than CAI and discourages use of codons enriched in
283 highly expressed genes (**Figure 5a**), notably for Arg CGT, Leu CTG, and Pro CCG. There are instances where χ and
284 CAI do correspond well (e.g., Gly GGA, GGC, GGG), but many codons show inverse trends between the two scales.
285 Generally, amino acids with multiple available tRNAs (including Arg, Leu, and Pro) correspond with larger
286 differences between expected RSCU values calculated for CAI and χ (and shown in **Figure 5a**), suggesting that
287 recruitment of different tRNAs is playing a role in determining Co-Expression Fitness (**Figure S10**). Interestingly, χ
288 favored codons do not always correspond to amino acids with multiple available tRNAs, indicating tRNA
289 abundance may not alone account for the observed effect, which could also be in part due to different translation
290 efficiencies created by favorable interactions of tRNA codon-anticodon pairs.

291

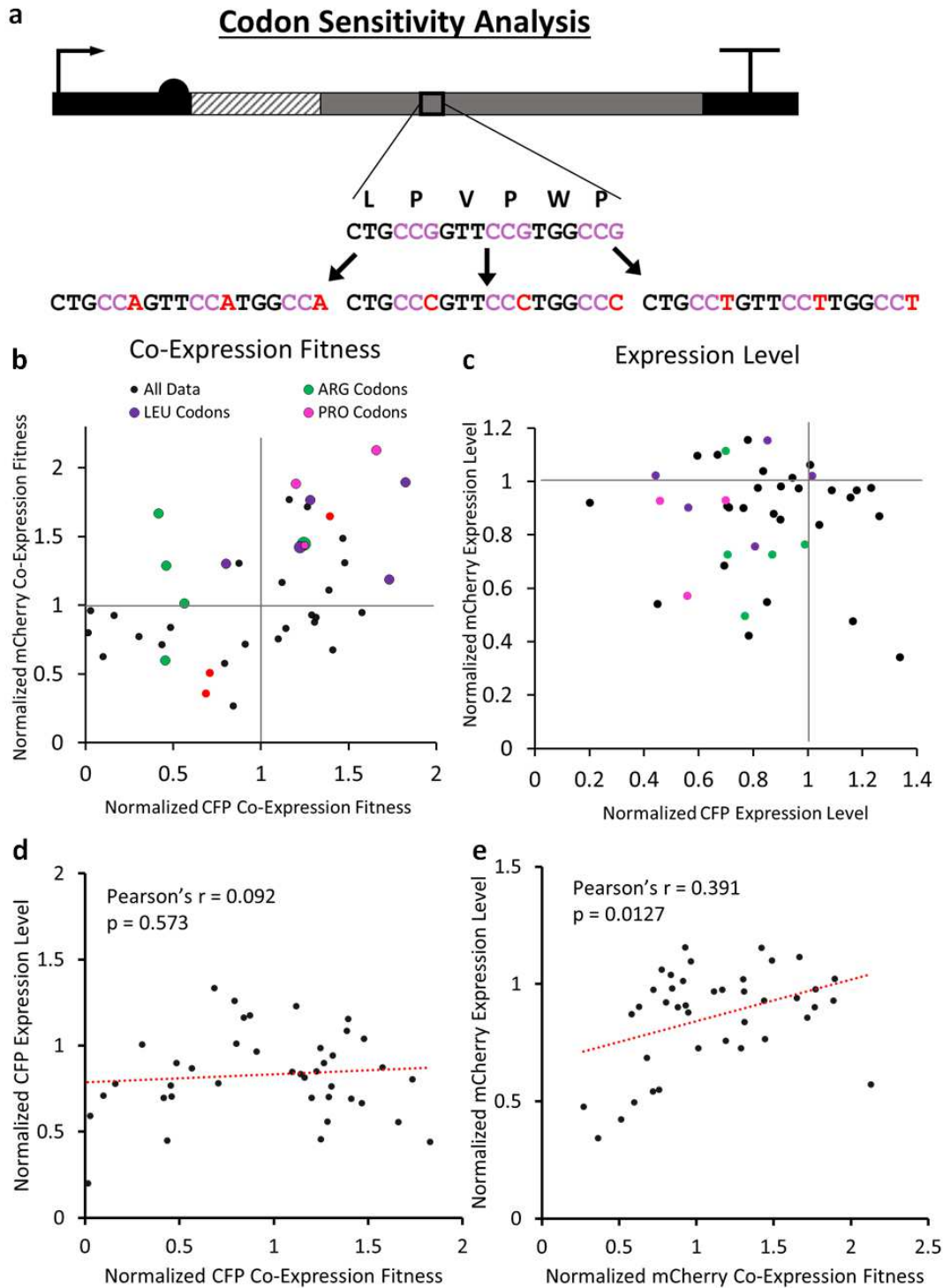


Figure 4: Systematic codon sensitivity analysis. **a.** Schematic of how genes are recoded for every amino acid. Starting with the highest CAI weighted codon for every instance of each amino acid, they are recoded to alternative synonymous codons. Example shown is for Proline. **b.** Mean fold change in (YFP) Co-expression Fitness upon CFP or mCherry co-expression, normalized to the parental (high CAI) control. **c.** Mean fold change in CFP or mCherry Expression Level relative to parental control. **d-e.** Poor correlations (Pearson's r) between fold change in Expression Level of CFP or mCherry recodes with Co-expression Fitness.

292
293
294
295
296
297
298
299

300 Utilizing the new χ weights, we next created several CFP and mCherry sequences that were optimized to varying
301 degrees on the new χ scale (**Figure 5b**). Specifically, we created a $\chi = 1$, ENC = 20 sequence, along with 4 sets of 3
302 different sequences each holding χ constant at 0.95, 0.85, 0.75, and 0.65 for both CFP and mCherry by using a
303 greedy algorithm (**Figure S11**). The lower end of the χ scale for the CFP/mCherry genes was approximately 0.6,
304 which is dictated by the protein sequence, and lowest W_{ij} values for each set of codons (see methods). When the χ
305 recoded sequences were assayed for fitness and expression (**Figure 5c–e**), there was a very strong positive
306 correlation between CFP and mCherry analogous re-codes for fitness and expression, indicating that these
307 synonymous coding schemes are a primary determinant for how a gene performs regardless of amino acid
308 sequence. Remarkably, we also observe a strong positive correlation between χ and, both, Growth Fitness and Co-
309 Expression Fitness—indicating that the weights derived from the individual codon assay are additive to improve
310 the fitness of various globally-recoded sequences (**Figure S12**). High χ sequences clearly provide reduced
311 competition for host resources and improved fitness. The χ scale is less predictive of expression, which is expected
312 as it was not part of the criteria used to create the codon weights. Despite this, there is a good correlation
313 between CFP and mCherry re-coded sequences in terms of Expression Level, indicating that codon usage bias does
314 generally predict expression. Importantly, there are several sequences with reduced burden that retain relatively
315 high expression, which represents an excess translational capacity for sequences re-coded using high χ values.
316

317 To investigate which codon usage bias patterns have the greatest contribution to Co-expression Fitness, we
318 analyzed RSCU across all variable 59 codon dimensions (excluding stop, Trp, and Met codons) for each of the CFP
319 and mCherry re-coded sequences (as seen in **Figure 5b**) using PCA (**Figure 6**). We were able to represent 46.7% of
320 the total sequence variation in the first 3 dimensions (**Figure S13**) when analyzing the CFP and mCherry recodes'
321 RSCU along with 773 *E. coli* operons. Here again PC1 and PC2 primarily explain variation across *E. coli* sequences,
322 but intriguingly we see a new highly orthogonal dimension in PC3 that explains variation in the χ sequences, and
323 PC1 vs. PC3 best differentiate the χ re-coded sequences from natural *E. coli* operons. The χ sequences generally
324 have intermediate to low values on the CAI scale with low overall CAI variation, meaning they would not have been
325 predicted to express well using CAI (**Figure 6a**). This is somewhat surprising given that many of the re-codes with
326 moderate to high χ (0.8–0.95) still exhibit relatively high expression compared with the high CAI control as
327 demonstrated in **Figure 5e**. When mapping χ values to the data, we see that χ describes variation along PC3 very
328 well (**Figure 6b, Figure S14**). *E. coli* operon sequences do not vary significantly on the χ scale, implying that the re-
329 coded sequences explore novel coding schemes orthogonal to natural sequence space. Examining the loadings for
330 the 3 most biased natural codons, we find that the high χ sequences are using synonymous variations for Arg, Leu,
331 and Pro that differ as expected from highly expressed genes. We conclude that competition for tRNA isoacceptors
332 in high demand by highly expressed essential genes primarily drives competition for translation elongation
333 resources and avoiding specific codons that are over-represented in such native genes provides a novel strategy to
334 improve the Co-Expression Fitness of heterologous genes.
335

336 Given the breadth of existing knowledge regarding codon optimization, we also evaluated how χ compares with
337 other reported CUB strategies such as the tRNA adaptation index (tAI)⁷ and normalized translation efficiency
338 (nTE)⁶. These approaches weight codons based on their co-adaptation to the tRNA pool or the tRNA supply vs.
339 codon demand respectively. We calculated the expected RSCU of a perfectly adapted gene sequence using these
340 various scales to assess their degree of similarity (**Figure S15**), and found that stAI (species specific TAI using *E. coli*
341 specific weights)²¹ correlates the closest with χ (Pearson's $r = 0.393$, $p = 0.002$), but does not provide as much
342 differentiation between codons available for each amino acid. We suspect the primary differentiator of the χ re-
343 coding strategy relative to tAI or nTE is that it provides empirical insight into which specific codons have excess
344 capacity for translation as opposed to an approach relying solely on genomic statistics and approximations. Further
345 analysis of the χ re-coded sequences did not reveal any consistent correlation with secondary structure or GC
346 content between CFP and mCherry re-codes, supporting the notion that specific codon use is likely driving
347 sequence behavior (**Figure S16**). We also re-coded 10 random genes with 3 free commercial re-coding algorithms
348 to analyze whether any of them exhibit exploration of χ related CUB strategies and found that they generally vary
349 along classical *E. coli* CUB and seek to adapt to host codon use without optimizing in the χ sequence space (**Figure**
350 **S17**).

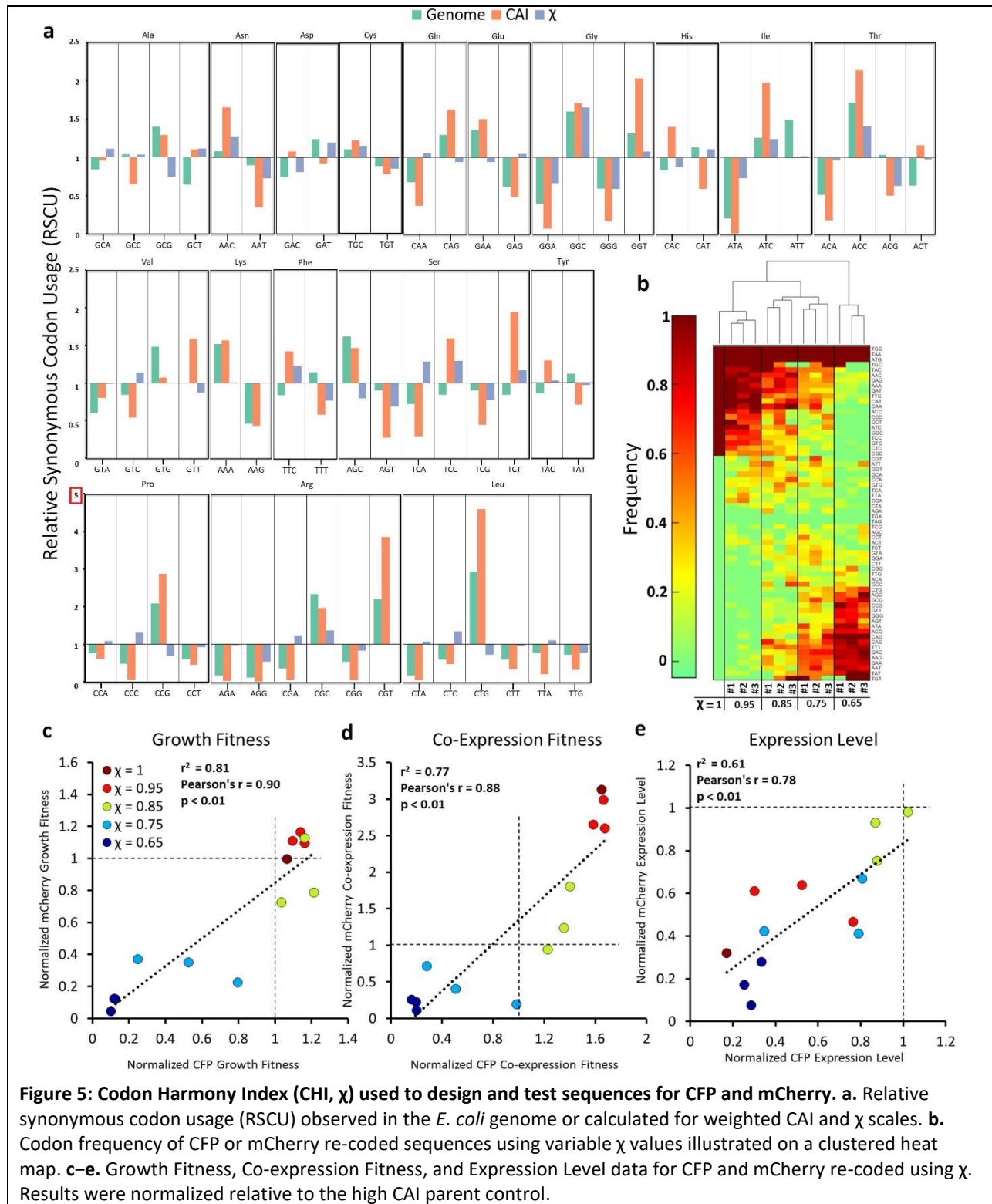
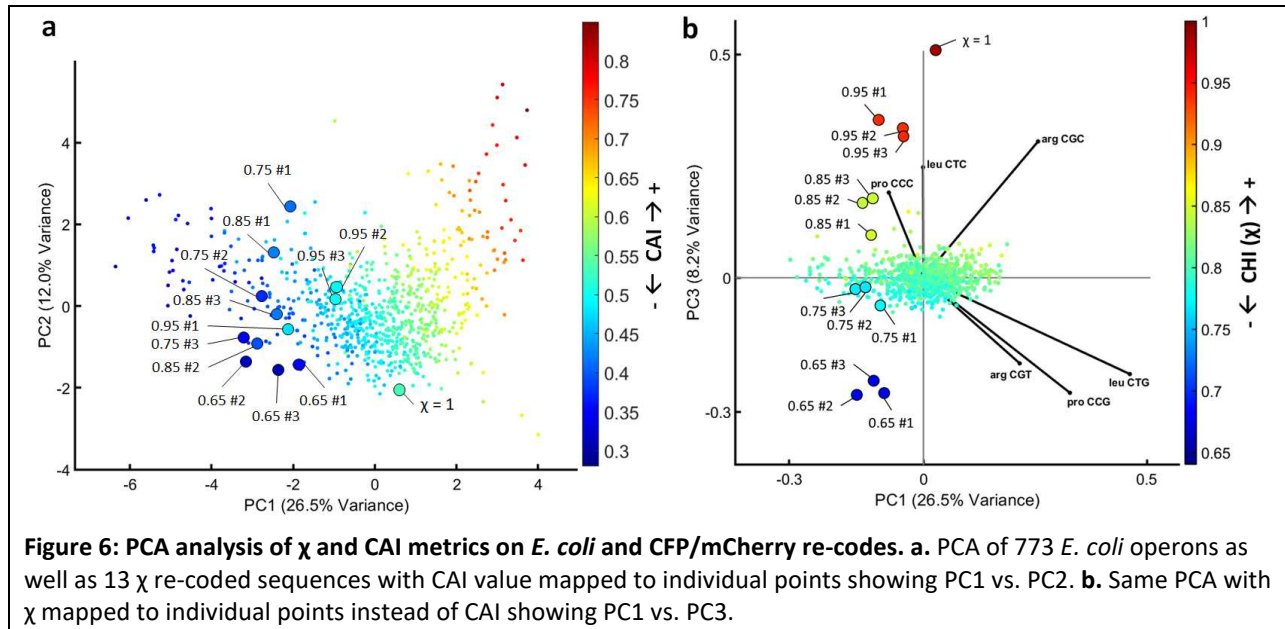


Figure 5: Codon Harmony Index (CHI, χ) used to design and test sequences for CFP and mCherry. a. Relative synonymous codon usage (RSCU) observed in the *E. coli* genome or calculated for weighted CAI and χ scales. **b.** Codon frequency of CFP or mCherry re-coded sequences using variable χ values illustrated on a clustered heatmap. **c–e.** Growth Fitness, Co-expression Fitness, and Expression Level data for CFP and mCherry re-coded using χ . Results were normalized relative to the high CAI parent control.

351
352
353
354
355
356
357
358
359
360
361

In theory, χ could also correlate with CUB in phages that infect *E. coli* and have co-adapted to maximize gene expression without overwhelming host resources. There have been reports of not only co-adaptation to tRNA pools^{40,41}, but also translational selection for CUB dissimilarity between viruses and hosts to avoid excessive

362 competition for tRNAs⁴². We examined codon usage in 12 common coliphages known to infect *E. coli* to examine
 363 whether CUB in such parasitic viruses may have evolved to harmonize with bacterial hosts as a means to allow
 364 better co-utilization of shared translational resources (**Figure S18**). Our analysis indicates that phage genes
 365 generally tend to avoid CUB at high values of CAI (>0.7) and exhibit a slightly higher mean χ than *E. coli* genes. This
 366 suggests that it may be more productive in the phage life cycle to avoid excessive similarity and competition with
 367 their host, but there is another unique aspect of the CUB in χ that was not strongly selected for in phages. It is
 368 possible that the translational resource demand from an overexpressed protein on a multi-copy vector is higher
 369 than natural genes have encountered and is thus under a higher level of translational selection resulting in novel
 370 types of advantageous CUB reflected by χ that cannot be inferred from natural sequence space.
 371



372
 373 **Figure 6: PCA analysis of χ and CAI metrics on *E. coli* and CFP/mCherry re-codes. a.** PCA of 773 *E. coli* operons as
 374 well as 13 χ re-coded sequences with CAI value mapped to individual points showing PC1 vs. PC2. **b.** Same PCA with
 375 χ mapped to individual points instead of CAI showing PC1 vs. PC3.

376
 377
 378 **DISCUSSION:**

379
 380 Protein translation is one of the most resource intensive cellular processes, which has yielded significant CUB
 381 observed in nature, especially in single cellular microorganisms often used as expression hosts⁴³. Most
 382 conventional codon optimization strategies operate under the key assumption that translational selection in
 383 naturally evolved systems provides CUB that is relevant for the overexpression of heterologous genes. This may be
 384 partially true, but realistically, the overexpression of genes can push host resource demand beyond levels required
 385 for native gene expression⁴⁴, resulting in translational selective pressures that organisms haven't evolved with.
 386 Protein expression must also be considered in the context of increasingly complicated engineered systems, and
 387 often in synthetic biology and metabolic engineering efforts, overexpression is not nearly as important as reliable
 388 and predictable gene expression and host fitness⁴⁵. Here we have revealed both in vitro and in an *E. coli* model
 389 that translation elongation can limit protein expression, and often has profitable or catastrophic consequences on
 390 system-wide resource availability.

391 In our TxTL assay, we found that proteins coded with similar CAI compete for the same tRNA supply, and re-coded
 392 genes can reduce such competition. Consequently, high CAI sequences are ribosome-limited, demonstrating
 393 reduced synthesis rates that are also highly sensitive to competition. In certain cases, low CAI genes are
 394 monopolistic or anti-competitive with free ribosomes and are thus insensitive to increased demand from high CAI
 395 sequences, albeit at the expense of overall resources. Theoretical frameworks have been well established to
 396 explain how resource limited translation can lead to the sequestration of ribosomes, but these studies generally
 397 rely on ribosome footprinting data³⁵ and tRNA copy number^{6,7} to infer codon elongation times, which are indirect
 398 measurements of ribosome flux on a given mRNA.

399 Our novel experimental approach using an *E. coli* model demonstrates the sensitivity of system resources at
400 individual codon resolution and reveals key differences between the optimal CUB for highly expressed native genes
401 vs. overexpressed proteins. Several previous studies have investigated CUB using randomized libraries that fail to
402 thoroughly explore the vast sequence space available when re-coding a gene⁴⁶. Such randomized sequences will
403 generally regress to intermediate RSCU values for each codon, and rarely sample the extremities of the sequence
404 space available (**Figure S19**). By systematically re-coding individual amino acids to each alternate codon in multiple
405 proteins, we have methodically investigated how individual codons contribute to gene Expression Level and Co-
406 Expression Fitness at further extremities of the theoretical design space than have been previously explored. The
407 avoidance of codons with very high CUB in native essential genes (e.g., for Arg/Leu/Pro) is a novel driver of
408 reduced genetic burden.

409 We used individual codon sensitivity data to create a new re-coding strategy that optimizes for fitness (CHI or χ)
410 and demonstrate how the new codon weighting method enables the creation of unique CUB strategies that are
411 not represented naturally in *E. coli*. Using PCA for dimensional reduction, our methodology reveals how sequences
412 with identical CAI scores can still exhibit distinct variations in CUB that result in different phenotypes, namely
413 improvements in Co-Expression Fitness. Remarkably, globally re-coded sequences were found to have predictable
414 phenotypes informed from the additive effects of individual codon use, allowing us to leverage a relatively small
415 dataset to predict phenotypes in a vast sequence space. While global sequence characteristics including GC
416 content, structure, and a variety of sequence motifs are all known to contribute to protein expression², our results
417 suggest that codon bias is a strong predictor of both protein expression and fitness and can be optimized
418 independently of the UTRs or 5' coding sequence. An analysis of *E. coli* phage CUB reveals that while parasitic
419 organisms may avoid over-use of preferred host codons, a concept that has been recently suggested⁴², the
420 demands of heterologous gene over-expression and resulting selective pressures are likely to have different
421 resource demands than those of viruses, and thus may have overlapping yet still largely distinct CUB fitness
422 landscapes.

423
424 The data-informed strategy in this study represents an approach that could be extended to other microbes
425 including eukaryotic systems, where ongoing controversy over the impact CUB has on host-gene fitness has been
426 unresolved^{47–51}. While our study included 2 proteins (CFP and mCherry) with very different amino acid sequences,
427 measuring Expression Level and Co-Expression Fitness for additional proteins could further refine χ , and provide
428 additional insight for maximizing expression and fitness together. The new χ metric is more predictive of *trans*
429 effects (Co-expression Fitness) than *cis* effects (Expression Level), thus further optimization of translation initiation
430 and CUB that maximizes both expression and fitness is an interesting future objective. The observation that there
431 are several sequences with relatively high expression and high fitness illustrates there are solutions to co-optimize
432 both genetic traits. In practice, re-coding genes with high CAI will often lead to higher expression with low overall
433 fitness, but re-coding with high χ values (between 0.9–0.95) should provide reasonably high expression with more
434 orthogonal resource demands. Similar data sets could also be collected for any organism where protein expression
435 is feasible, which could also provide insights into how species differ in the role CUB plays regarding resource
436 allocation. It is possible that with more inter-species data, organism specific χ weights could be predicted *a priori*
437 based on the avoidance of codons overrepresented in host genes. Practically, this study should improve the
438 predictability and robustness of genetic engineering by enabling the co-optimization of gene expression and
439 fitness, especially for multi-gene expression systems.

440

441 **MATERIALS AND METHODS:**

442

443 **Equations used to assess codon usage bias.**

444 We calculated codon adaptation following the classical method reported originally by Sharp and Li¹⁸. This method
 445 relies on first calculating relative synonymous codon usage (**RSCU**) in a genetic sequence, which is defined by

446 **Equation 1:**

$$RSCU_{ij} = \frac{X_{ij}}{\frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij}} \quad (1)$$

447 RSCU calculates the observed frequency of codon **j** belonging to amino acid **i** divided by expected frequency,
 448 where **X** is the number of occurrences for codon **j** in a given sequence. The expected frequency is simply the
 449 number of occurrences for any codon belonging to amino acid **i**, divided by the number of codons (**n**) available for
 450 that particular amino acid. RSCU is used instead of raw frequency values to normalize observed codon frequency
 451 based on the total codons available. An RSCU value < 1 indicates bias against the codon, while an RSCU value > 1
 452 indicates a bias toward the codon, and RSCU = 1 indicates no bias. The RSCU values for each codon can be used to
 453 calculate relative adaptiveness (**W**), which is defined by **Equation 2:**

$$W_{ij} = RSCU_{ij} / RSCU_{imax} \quad (2)$$

454 Relative adaptiveness is the RSCU for a codon **j** belonging to amino acid **i** divided by the RSCU for the codon in the
 455 set for amino acid **i** with the highest RSCU value (*imax*). In other words, **W** gives a value of 1 for codons in a target
 456 sequence that match the frequency of the most common codon in a reference sequence. **W** values are used in
 457 calculating the codon adaptation index (**CAI**) defined by **Equation 3:**

$$CAI = \left(\prod_{k=1}^L w_k \right)^{1/L} \quad (3)$$

458

459 CAI is the geometric mean of the **W** values for each codon in a given sequence containing **L** codons. Importantly,
 460 the reference sequence(s) and calculated RSCU values that **W** values are derived from can be from any source.
 461 Unless otherwise indicated, in this study, CAI refers to **W** values for a set of highly expressed set of *E. coli* genes.
 462 Alternatively, CAI can be computed based on **W** values for CUB across the entire genome, sTAI weights²¹, or χ
 463 weights (**See Data S2 for W values used in various calculations**). Normalized translational efficiency (nTE) was
 464 calculated as previously described⁶ by taking the ratio of species specific TAI weights for *E. coli*²¹ (supply) vs. the
 465 codon use across the *E. coli* transcriptome (demand) defined by **Equation 4:**

$$nTE_{ij} = sTAI_{ij} / Frequency_{ij} \quad (4)$$

466 The nTE_{ij} values are analogous to W_{ij} values for the calculation of nTE, which proceeds the same as for CAI by taking
 467 the geometric mean across a sequence (as in equation 3). In this study, nTE was calculated using genomic codon
 468 frequency as opposed to codon use (originally defined as codon occurrence multiplied by RNA transcript
 469 abundance), as the two were found to be highly correlated (**Figure S20**). Lastly, the effective number of codons
 470 (**ENC**) is often used as a measure of codon bias in a sequence, and is calculated using **Equation 5:**

$$ENC = 2 + \frac{9}{F_2} + \frac{1}{F_3} + \frac{5}{F_4} + \frac{3}{F_6} \quad (5)$$

471

472 ENC can take a value from 20, in the case of extreme bias where one codon is exclusively used for each amino acid,
473 to 61 when the use of alternative synonymous codons is equally likely. The value F is the average probability that
474 two randomly selected codons for an amino acid with n number of synonymous codons will be identical⁵².

475

476 **Data sources used in analysis.**

477 Genomic codon usage for *E. coli* K12 MG1655 and *E. coli* MRE600 were assessed by analyzing codon bias from
478 published annotated genomes obtained from NCBI under the accession numbers NC_000913.3 and CP014197.1
479 respectively using MATLAB. Phage analysis was done with annotated phage genomes from NCBI, and accession
480 numbers are listed in **Figure S18**. Exact codon frequencies and relative adaptiveness values (W) used in this study
481 for calculating CAI in reference to highly expressed genes CUB, entire genome CUB, sTAI, or nTE, can be found in
482 **Data S2**. The W values for χ and associated information from the study can be found in **Data S1**. W values for
483 highly expressed genes were originally downloaded online from GenScript, and were cross referenced to published
484 values⁵³. The sTAI codon weights were downloaded online from a publically available database ([http://tau-
485 tai.azurewebsites.net/](http://tau-tai.azurewebsites.net/))²¹. The tRNA copy numbers referenced in this study (**Figure S8**) were downloaded from the
486 Genomic tRNA Database (<http://gtrnadb.ucsc.edu/>)⁵⁴.

487

488 **Ribosome flow model.**

489 The implemented ribosome flow model (RFM) (**Figure S1**) was adapted from Zur et al. using open source Matlab®
490 code³⁴. In this model, an mRNA is divided into n number of chunks, where each chunk is 9 codons (27 bases),
491 approximately the footprint of an *E. coli* ribosome. Translation time of each chunk is based on local λ , which is a
492 sum of the individual times it takes to translate each codon in a chunk. Codon times used are available in **Data S3**.
493 Ribosome collisions are also accounted for in the model as a function of the ribosome density in adjacent
494 positions. In this model, the protein production rate is the rate of translation of the final position on the mRNA.
495 For this application, steady state ribosome densities were computed for CFP and YFP re-coded to use preferred
496 (high CAI) or rare (low CAI) codons. To demonstrate the relationship between initiation rate and translation rate
497 for different sequences, steady state protein production rates are calculated for different initiation rates.

498

499 **Gene design and re-coding.**

500 All genetic re-coding designs and analysis were executed in Matlab® using custom functions. Code is made
501 available online at <https://github.com/nair-lab>. A full list of amino acid and DNA sequences used in this study can
502 be found in **Data S4**. CFP and YFP were initially cloned through site directed mutagenesis of an existing super-
503 folder GFP protein based on previously reported sequences.^{37,55} For the systematic analysis of codon use design,
504 CFP or mCherry were re-coded starting from highly biased sequences using the most preferred codon for each
505 amino acid (CAI = 1 and ENC = 20), not taking into account the first 17 codons. The first 17 codons were held
506 constant for all re-codes and were based on previously used sequences that functionally expressed well. A Matlab®
507 script was then used to systematically design sequences where every instance of an amino acid was mutated to a
508 single alternate synonymous codon. In the design of sequences with novel re-coding schemes, a greedy algorithm
509 was used (**Figure S11**), that functions by randomly mutating a codon to a synonymous alternative, then evaluating
510 whether the new sequence is closer to the target CAI (or in this specific instance χ value). To re-code CFP and
511 mCherry to a desired χ value, a starting sequence was first randomized to ensure there was no initial bias, and then
512 the algorithm was followed to the target χ value. We generated several unique output sequences with the same χ
513 value but different coding sequences, then selected 3 sequences for each value of χ tested making sure they were
514 substantially different from each other based on hierarchal clustering done in Matlab®.

515

516 **Plasmids and strain construction.**

517 All plasmids were cloned from existing vectors with restriction enzyme sites already present (**figure S21, S23, Data**
518 **S4**), which also contained 5' and 3' UTRs. Genes were all custom ordered synthesized as full length double

519 stranded DNA fragments with AarI restriction sites on the 5' and 3' termini. A type IIS restriction enzyme cloning
520 approach with AarI was used to insert synthesized double stranded DNA gene fragments into the desired vector.
521 All constructs were sequence verified from clonally pure DNA using Sanger sequencing across the gene and UTRs.
522 The screening strain used to assess Co-Expression fitness was engineered from *E. coli* K12 MG1655 (CGSC#: 6300).
523 The YFP reporter was integrated in an intergenic region (~3,938,000 bp) between the *rsmG-atpI* genes using λ -
524 Red based homologous recombination of the YFP CAI = 0.96 sequence, which was under the control of a strong
525 constitutive promoter (FAB46) and RBS (BCD7) based on a previous study,³⁸ and a 5' insulator and 3' terminator
526 (**Figure S22, Data S4**). The method of integration and marker excision method has been previously reported
527 (Datsenko and Wanner).⁵⁶ Briefly, a linear cassette consisting of the gene, UTRs, and an attached kanamycin
528 resistance marker was amplified by PCR with ~500bp of homology to the desired locus on either end.
529 Chromosomally integrated clones were identified by colony PCR and sequence verified via Sanger sequencing of
530 the PCR product including several hundred bases of chromosomal DNA and the entire integrated heterologous
531 expression cassette. Sequence verified clones had the integrated kanamycin marker removed through the
532 previously described FLP-FRT site specific recombinase method and were again Sanger sequenced for final
533 verification.

534

535 **in vitro transcription-translation (TxTL) assay.**

536 The TxTL assay was carried out using the NEB PURExpress[®] kit (E6800). This assay relies on T7 polymerase, and
537 consists of purified reconstituted components. Accordingly, CFP, YFP, and mCherry expression cassettes were first
538 cloned into a pBAC vector with a T7 promoter and strong RBS (BCD7) (**Figure S23 a–b, Data S4**). The genes were
539 also flanked by an insulator and terminator sequence on the 5' and 3' UTR respectively. Once clonally pure and
540 sequence verified, expression cassettes were amplified by PCR (from the beginning of the insulator to end of the
541 terminator) and normalized in concentration using UV-vis spectroscopy at $\lambda = 260\text{nm}$. A master mix was first
542 prepared according to the PURExpress[®] published protocol, which was kept on ice until use. Reactions were scaled
543 down to 5 μL final volume and carried out in Corning[®] low volume 384-well white flat bottom polystyrene TC-
544 treated microplates (part # 3826). Reactions were initiated by the addition of DNA using a multi-channel pipette
545 ($n=2$ per condition), followed by immediate transfer to a Tecan Infinite[®] M1000 microplate reader. A DNA
546 concentration of 20ng/ μL each was found to generally maximize competition between two genetic cassettes
547 (**Figure S23 c-d**). Assays were run for 6hr. at 37°C with fluorescent reads every 5 minutes of each protein being
548 analyzed (CFP: Ex. 435nm, Em. 470nm, YFP: Ex. 510nm, Em. 530nm, mCherry: Ex. 585nm, Em. 612nm). Reported
549 reaction rates reflect the maximum rate observed for each individual replicate.

550

551 **in vivo fitness and expression assay.**

552 To assess Co-Expression Fitness, Growth Fitness, and Expression Level, sequence verified plasmid constructs were
553 transformed into *E. coli* K12 MG1655 with the chromosomally integrated YFP reporter. Unless noted otherwise,
554 overexpressed proteins were under control of the Trc promoter with a strong RBS (BCD7) (**Data S4**). 3 individual
555 transformants were isolated and grown overnight in 400 μL LB broth (BD Difco[™]) with selective antibiotic at 37°C in
556 96 deep well plates (Greiner Bio-One MASTERBLOCK[®], 96 Well, 2 ML Item: 780270) for 24 hr. Cultures were then
557 split and diluted 1:40 into LB broth with selective antibiotic and with or without 500 μM inducer (IPTG) in black 96
558 well clear bottom micro-titer plates (Thermo product: 165305). Plates were incubated for 8 hours with shaking at
559 37°C in a Tecan Infinite[®] M1000 microplate reader with monitoring every 5 minutes for OD600, as well as
560 fluorescence (CFP: Ex. 435nm, Em. 470nm, YFP: Ex. 510nm, Em. 530nm, mCherry: Ex. 585nm, Em. 612nm). Data
561 were analyzed by comparing independent induced vs. uninduced cultures in terms of fluorescence and growth. To
562 account for lag phase and differences in rates within a single term, the background subtracted area under the
563 curve (AUC) was used for each respective signal using a Matlab[®] numerical integrator. The timespan evaluated was
564 bounded by the time it took any sample to reach the upper limit of detection for fluorescence, which often took
565 between 4-6 hours. In most cases, the mean of 3 replicates was compared (fold change) relative to a control
566 sequence (e.g. the high CAI starting sequence).

567

568

569 **Additional data analysis.**

570 Except in the case of measured reaction rates, all data were collected from distinct samples. Mean, standard
571 deviation, linear regression, correlation analysis, dimensional reduction, and associated statistics were calculating
572 using built in functions in Matlab® or Microsoft Excel. Error bars in all plots represent standard deviation. Principal
573 component analysis and hierarchal clustering were always carried out on an m x n matrix of RSCU values with
574 codons in 61 rows and n number of gene sequences in columns. For RNA folding calculations, the minimum free
575 energy was calculated for sequences using the Vienna RNAfold Version 2.5.1 software.⁵⁷

576

577 **COMPETING INTERESTS:**

578 The authors declare no competing or conflicting interests.

579

580

581 **ACKNOWLEDGMENT AND FUNDING:**

582 ManusBio supported and funded this work. We thank members of the Nair lab for thoughtful discussions and
583 advice.

584

585

586 **ASSOCIATED CONTENT:**

587 There are 24 figures included in the Supplemental Information, and there are 4 supplementary data files.

588

589

590 **DATA AVAILABILITY STATEMENT:**

591 Additional data are available upon request. Additional supplementary Matlab® code can be found at
592 <https://github.com/nair-lab/CHI>.

593

594

595 **AUTHOR CONTRIBUTIONS:**

596 A.M.L. performed the experimental work and data analysis. A.M.L. and N.U.N. conceived the study, planned the
597 experiments, and wrote/edited the manuscript.

598

599 **REFERENCES:**

600

- 601 1. Plotkin, J. B. & Kudla, G. Synonymous but not the same: the causes and consequences of codon bias. *Nat. Rev. Genet.* **12**, 32–42 (2011).
- 602
- 603 2. Quax, T. E. F., Claassens, N. J., Söll, D. & van der Oost, J. Codon Bias as a Means to Fine-Tune Gene Expression. *Mol. Cell* **59**, 149–161 (2015).
- 604
- 605 3. Frumkin, I. *et al.* Gene Architectures that Minimize Cost of Gene Expression. *Mol. Cell* **65**, 142–153 (2017).
- 606 4. 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. *J. Mol. Biol.* **151**, 389–409 (1981).
- 607
- 608
- 609 5. Boël, G. *et al.* Codon influence on protein expression in E. coli correlates with mRNA levels. *Nature* **529**, 358–363 (2016).
- 610
- 611 6. Pechmann, S. & Frydman, J. Evolutionary conservation of codon optimality reveals hidden signatures of cotranslational folding. *Nat. Struct. Mol. Biol.* **20**, 237–243 (2013).
- 612
- 613 7. Reis, M. d., Savva, R. & Wernisch, L. Solving the riddle of codon usage preferences: a test for translational selection. *Nucleic Acids Res.* **32**, 5036–5044 (2004).
- 614
- 615 8. Tuller, T. *et al.* An Evolutionarily Conserved Mechanism for Controlling the Efficiency of Protein Translation. *Cell* **141**, 344–354 (2010).
- 616
- 617 9. Goodman, D. B., Church, G. M. & Kosuri, S. Causes and Effects of N-Terminal Codon Bias in Bacterial Genes. *Science (80-.)*. **342**, 475–479 (2013).
- 618
- 619 10. Kudla, G., Murray, A. W., Tollervey, D. & Plotkin, J. B. Coding-Sequence Determinants of Gene Expression in Escherichia coli. *Science (80-.)*. **324**, 255–258 (2009).
- 620
- 621 11. Presnyak, V. *et al.* Codon optimality is a major determinant of mRNA stability. *Cell* **160**, 1111–1124 (2015).
- 622 12. Saunders, R. & Deane, C. M. Synonymous codon usage influences the local protein structure observed. *Nucleic Acids Res.* **38**, 6719–6728 (2010).
- 623
- 624 13. Ciryam, P., Morimoto, R. I., Vendruscolo, M., Dobson, C. M. & O’Brien, E. P. In vivo translation rates can substantially delay the cotranslational folding of the Escherichia coli cytosolic proteome. *Proc. Natl. Acad. Sci.* **110**, E132–E140 (2013).
- 625
- 626
- 627 14. Yu, C.-H. *et al.* Codon Usage Influences the Local Rate of Translation Elongation to Regulate Co-translational Protein Folding. *Mol. Cell* **59**, 744–754 (2015).
- 628
- 629 15. Kafri, M., Metzler-Raz, E., Jona, G. & Barkai, N. The Cost of Protein Production. *Cell Rep.* **14**, 22–31 (2016).
- 630 16. Borkowski, O. *et al.* Cell-free prediction of protein expression costs for growing cells. *Nat. Commun.* **9**, 1457 (2018).
- 631
- 632 17. Dong, H., Nilsson, L. & Kurland, C. G. Co-variation of tRNA Abundance and Codon Usage in Escherichia coli at Different Growth Rates. *J. Mol. Biol.* **260**, 649–663 (1996).
- 633
- 634 18. Sharp, P. M. & Li, W. The codon adaptation index—a measure of directional synonymous codon usage bias, and its potential applications. *Nucleic Acids Res.* **15**, 1281–1295 (1987).
- 635
- 636 19. Lipinszki, Z. *et al.* Enhancing the Translational Capacity of E. coli by Resolving the Codon Bias. *ACS Synth. Biol.* **7**, 2656–2664 (2018).
- 637
- 638 20. Lyu, X., Yang, Q., Zhao, F. & Liu, Y. Codon usage and protein length-dependent feedback from translation elongation regulates translation initiation and elongation speed. *Nucleic Acids Res.* **49**, 9404–9423 (2021).
- 639
- 640 21. Sabi, R., Volvovitch Daniel, R. & Tuller, T. stAI calc : tRNA adaptation index calculator based on species-

- 641 specific weights. *Bioinformatics* **33**, btw647 (2016).
- 642 22. Boo, A., Ellis, T. & Stan, G. B. Host-aware synthetic biology. *Curr. Opin. Syst. Biol.* **14**, 66–72 (2019).
- 643 23. Shopera, T., He, L., Oyetunde, T., Tang, Y. J. & Moon, T. S. Decoupling Resource-Coupled Gene Expression in
644 Living Cells. *ACS Synth. Biol.* **6**, 1596–1604 (2017).
- 645 24. Ceroni, F. *et al.* Burden-driven feedback control of gene expression. *Nat. Methods* **15**, 387–393 (2018).
- 646 25. Huang, H.-H., Qian, Y. & Del Vecchio, D. A quasi-integral controller for adaptation of genetic modules to
647 variable ribosome demand. *Nat. Commun.* **9**, 5415 (2018).
- 648 26. Darlington, A. P. S., Kim, J., Jiménez, J. I. & Bates, D. G. Dynamic allocation of orthogonal ribosomes
649 facilitates uncoupling of co-expressed genes. *Nat. Commun.* **9**, 695 (2018).
- 650 27. Frumkin, I. *et al.* Gene Architectures that Minimize Cost of Gene Expression. *Mol. Cell* **65**, 142–153 (2017).
- 651 28. Cambray, G., Guimaraes, J. C. & Arkin, A. P. Evaluation of 244,000 synthetic sequences reveals design
652 principles to optimize translation in *Escherichia coli*. *Nat. Biotechnol.* **36**, 1005–1015 (2018).
- 653 29. Frumkin, I. *et al.* Codon usage of highly expressed genes affects proteome-wide translation efficiency. *Proc.*
654 *Natl. Acad. Sci.* **115**, E4940–E4949 (2018).
- 655 30. Nieß, A., Siemann-Herzberg, M. & Takors, R. Protein production in *Escherichia coli* is guided by the trade-
656 off between intracellular substrate availability and energy cost. *Microb. Cell Fact.* **18**, 8 (2019).
- 657 31. Vind, J., Sørensen, M. A., Rasmussen, M. D. & Pedersen, S. Synthesis of Proteins in *Escherichia coli* is
658 Limited by the Concentration of Free Ribosomes. *Journal of Molecular Biology* vol. 231 678–688 (1993).
- 659 32. Gorochowski, T. E., Avcilar-Kucukgoze, I., Bovenberg, R. A. L., Roubos, J. A. & Ignatova, Z. A Minimal Model
660 of Ribosome Allocation Dynamics Captures Trade-offs in Expression between Endogenous and Synthetic
661 Genes. *ACS Synth. Biol.* **5**, 710–720 (2016).
- 662 33. Reuveni, S., Meilijson, I., Kupiec, M., Ruppín, E. & Tuller, T. Genome-Scale Analysis of Translation
663 Elongation with a Ribosome Flow Model. *PLoS Comput. Biol.* **7**, e1002127 (2011).
- 664 34. Zur, H., Cohen-Kupiec, R., Vinokour, S. & Tuller, T. Algorithms for ribosome traffic engineering and their
665 potential in improving host cells' titer and growth rate. *Sci. Rep.* **10**, 21202 (2020).
- 666 35. Dana, A. & Tuller, T. The effect of tRNA levels on decoding times of mRNA codons. *Nucleic Acids Res.* **42**,
667 9171–9181 (2014).
- 668 36. Pédelacq, J.-D., Cabantous, S., Tran, T., Terwilliger, T. C. & Waldo, G. S. Engineering and characterization of
669 a superfolder green fluorescent protein. *Nat. Biotechnol.* **24**, 79–88 (2006).
- 670 37. Day, R. N. & Davidson, M. W. The fluorescent protein palette: tools for cellular imaging. *Chem. Soc. Rev.* **38**,
671 2887 (2009).
- 672 38. Mutalik, V. K. *et al.* Precise and reliable gene expression via standard transcription and translation initiation
673 elements. *Nat. Methods* **10**, 354–360 (2013).
- 674 39. Taniguchi, Y. *et al.* Quantifying *E. coli* Proteome and Transcriptome with Single-Molecule Sensitivity in
675 Single Cells. *Science (80-.)*. **329**, 533–538 (2010).
- 676 40. Chithambaram, S., Prabhakaran, R. & Xia, X. The Effect of Mutation and Selection on Codon Adaptation in
677 *Escherichia coli* Bacteriophage. *Genetics* **197**, 301–315 (2014).
- 678 41. Lucks, J. B., Nelson, D. R., Kudla, G. R. & Plotkin, J. B. Genome Landscapes and Bacteriophage Codon Usage.
679 *PLoS Comput. Biol.* **4**, e1000001 (2008).
- 680 42. Chen, F. *et al.* Dissimilation of synonymous codon usage bias in virus–host coevolution due to translational
681 selection. *Nat. Ecol. Evol.* **4**, 589–600 (2020).
- 682 43. Nieß, A., Siemann-Herzberg, M. & Takors, R. Protein production in *Escherichia coli* is guided by the trade-

- 683 off between intracellular substrate availability and energy cost. *Microb. Cell Fact.* **18**, 8 (2019).
- 684 44. Ceroni, F., Algar, R., Stan, G.-B. & Ellis, T. Quantifying cellular capacity identifies gene expression designs
685 with reduced burden. *Nat. Methods* **12**, 415–418 (2015).
- 686 45. McBride, C. D., Grunberg, T. W. & Del Vecchio, D. Design of genetic circuits that are robust to resource
687 competition. *Curr. Opin. Syst. Biol.* **28**, 100357 (2021).
- 688 46. Schmitz, A. & Zhang, F. Massively parallel gene expression variation measurement of a synonymous codon
689 library. *BMC Genomics* **22**, 149 (2021).
- 690 47. Shen, X., Song, S., Li, C. & Zhang, J. Synonymous mutations in representative yeast genes are mostly
691 strongly non-neutral. *Nature* **606**, 725–731 (2022).
- 692 48. Rodríguez-Beltrán, J. *et al.* Translational demand is not a major source of plasmid-associated fitness costs.
693 *Philos. Trans. R. Soc. B Biol. Sci.* **377**, (2022).
- 694 49. Torrent, M., Chalancon, G., De Groot, N. S., Wuster, A. & Madan Babu, M. Cells alter their tRNA abundance
695 to selectively regulate protein synthesis during stress conditions. *Sci. Signal.* **11**, 1–10 (2018).
- 696 50. Zhou, Z. *et al.* Codon usage is an important determinant of gene expression levels largely through its
697 effects on transcription. *Proc. Natl. Acad. Sci.* **113**, E6117–E6125 (2016).
- 698 51. Xia, X. A Major Controversy in Codon-Anticodon Adaptation Resolved by a New Codon Usage Index.
699 *Genetics* **199**, 573–579 (2015).
- 700 52. Wright, F. The ‘effective number of codons’ used in a gene. *Gene* **87**, 23–29 (1990).
- 701 53. Taniguchi, Y. *et al.* Quantifying E. coli Proteome and Transcriptome with Single-Molecule Sensitivity in
702 Single Cells. *Science (80-.)*. **329**, 533–538 (2010).
- 703 54. Chan, P. P. & Lowe, T. M. GtRNADB 2.0: an expanded database of transfer RNA genes identified in
704 complete and draft genomes. *Nucleic Acids Res.* **44**, D184–D189 (2016).
- 705 55. Pédelacq, J.-D., Cabantous, S., Tran, T., Terwilliger, T. C. & Waldo, G. S. SI_Engineering and characterization
706 of a superfolder green fluorescent protein. SF DsRed FR FR FR FR DsRed. *Nat. Biotechnol.* **37** (2006).
- 707 56. Datsenko, K. A. & Wanner, B. L. One-step inactivation of chromosomal genes in Escherichia coli K-12 using
708 PCR products. *Proc. Natl. Acad. Sci.* **97**, 6640–6645 (2000).
- 709 57. Lorenz, R. *et al.* ViennaRNA Package 2.0. *Algorithms Mol. Biol.* **6**, 26 (2011).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [DataS1CoexpressionFitnessWeightsforCHI.xlsx](#)
- [DataS2CodonweightsfrequenciesandRSCUvaluesused.xlsx](#)
- [DataS3CodonElongationTimesusedinRFM.xlsx](#)
- [DataS4Sequencesusedinthestudy.xlsx](#)
- [220821AMLResourceCompetitionManuscriptSlv5.docx](#)
- [221003nreditorialpolicychecklist.pdf](#)
- [221003nrreportingsummary.pdf](#)