# Codon-usage optimization in the prokaryotic tree of life:

# How synonymous codons are differentially selected in sequence

# domains with different expression levels and degrees of conservation.

**José Luis López [1], Mauricio Javier Lozano [1], María Laura Fabre [1], and Antonio Lagares [1] ***

[1]IBBM - Instituto de Biotecnología y Biología Molecular, CONICET, CCT-La Plata, Departamento de Ciencias Biológicas, Facultad de Ciencias Exactas, Universidad Nacional de La Plata, calles 47 y 115, 1900-La Plata, Argentina.

J.L.L. and M.J.L. contributed equally to this work.

**Key words:** codon usage selection, mutational bias, genome evolution, core genes, singletons, translation efficiency, translation accuracy

**\*Corresponding author**

Phone: +54-221-422-9777
Fax: +54-221-422-3409 ext. 56
E-mail: lagares@biol.unlp.edu.ar

1                                            **ABSTRACT**

2    Prokaryote genomes exhibit a wide range of GC contents and codon usages, both

3    resulting from an interaction between mutational bias and natural selection. In order to

4    investigate the basis underlying specific codon changes, we performed a comprehensive

5    analysis of 29-different prokaryote families. The analysis of core-gene sets with

6    increasing ancestries in each family lineage revealed that the codon usages became

7    progressively more adapted to the tRNA pools. While, as previously reported, highly-

8    expressed genes presented the more optimized codon usage, the singletons contained

9    the less selectively-favored codons. Results showed that usually codons with the highest

10    translational adaptation were preferentially enriched. In agreement with previous reports,

11    a C-bias in 2- to 3-fold codons, and a U-bias in 4-fold codons occurred in all families,

12    irrespective of the global genomic-GC content. Furthermore, the U-biases suggested that

13    $U_3$-mRNA–$U_{34}$-tRNA interactions were responsible for a prominent codon optimization in

14    both the more ancestral core and the highly expressed genes. A comparative analysis of

15    sequences that encode conserved-(*cr*) or variable-(v*r*) translated products, with each one

16    being under high- (HEP) and low- (LEP) expression levels, demonstrated that the

17    efficiency was more relevant (by a factor of 2) than accuracy to modelling codon usage.

18    Finally, analysis of the third position of codons (GC3) revealed that, in genomes of global-

19    GC contents higher than 35-40%, selection favored a GC3 increase; whereas in

20    genomes with very low-GC contents, a decrease in GC3 occurred. A comprehensive final

21    model is presented where all patterns of codon usage variations are condensed in five-

22    distinct behavioral groups.

23 **IMPORTANCE**

24

25   The prokaryotic genomes—the current heritage of the more ancient life forms on earth—

26   are comprised of diverse gene sets; all characterized by varied origins, ancestries, and

27   spatial-temporal–expression patterns. Such genetic diversity has for a long time raised

28   the question of how cells shape their coding strategies to optimize protein demands (*i.e.*,

29   product abundance) and accuracy (*i.e.*, translation fidelity) through the use of the same

30   genetic code in genomes with GC-contents that range from less than 20 to over 80%. In

31   this work, we present evidence on how codon usage is adjusted in the prokaryote tree of

32   life, and on how specific biases have operated to improve translation. Through the use of

33   proteome data, we characterized conserved and variable sequence domains in genes of

34   either high-  or low-expression level, and quantitated the relative weight of efficiency and

35   accuracy—as well as their interaction—in shaping codon usage in prokaryotes.

## INTRODUCTION

The wide range of GC contents exhibited by prokaryote genomes—*i. e.*, from less than 20 to 80%—are believed to be primarily caused by interspecies differences in mutational processes that operate on both the coding and the noncoding regions (1–6). Since prokaryote genomes consist mainly of coding regions that tightly reflect the genomic GC content, mutational bias is a main force that shapes the codon usage of the majority of the genes (7, 8). Thus, understanding how selection is coupled to mutational processes to model codon usage under such diverse GC contents is an essential issue (9–11). Recent evidence suggests that prokaryotic genomes with intermediate-to-high GC contents are affected by mutations that are universally biased in favor of AT replacements (12, 13). That process is counterbalanced by selection-based constraints that, in turn, increase the GC content and fine-tune codon usage—*i. e.*, the so-called mutation-selection-drift model (14–16). Intragenomic codon-usage heterogeneities, however, are always present among different gene sets—*i. e.*, between core genes that are shared throughout a given lineage, and singletons (unique accessory genes) that are taxa- and/or strain-specific (17, 18). Furthermore, in a multipartite genome, the linkage between the physical patterns of heterogeneity in codon usage and the replicon location of the different core genes has also been recently demonstrated (19). The analysis of intragenomic codon usage heterogeneities by different authors (20, 21) has served to identify at least the following three distinctive gene groups: The first comprises the majority of the coding sequences that are associated with the so-called typical codon usage, while the second consists of the putative highly expressed (PHE) genes involving codon usages that are the best adapted to the translational machinery (20, 22–26). The third contains genes that encode the accessory information including the singletons (unique genes) that are present in mobile genetic elements as well as in the more stable replicons (27–31). The intracellular variations in codon usage can be explained on the basis of selective pressures that operate with different strengths depending on gene

64    function and the resulting impact on cellular fitness (32). A search for the biochemical

65    basis associated with the heterogeneity in codon usage among different gene sets has

66    been the focus of numerous studies. Several lines of evidence have indicated that the

67    biased codon usage in PHE genes correlates with the copy number of the specific tRNA

68    species that decode the preferred codons (23, 33, 34) and with an optimal codon-

69    anticodon interaction (35). The latter includes both the classical Watson-and-Crick

70    interactions (WCIs) and a wobble-base pairing with the corresponding cognate tRNAs. All

71    these interactions have been taken into consideration in order to define different

72    numerical indices (36, 37) as estimators of the codon adaptation to the existing tRNA

73    pool. Though not considered in currently used translation-adaptation indices, evidence

74    has also been found for other nonstandard codon-anticodon interactions which, by

75    improving the decoding capacity, are also relevant to codon-usage evolution (38–41).

76        The analysis of an extensive number of genes with different functions,

77    ubiquitousness, and degrees of phylogenetic conservation has demonstrated that codon

78    usage is related to gene-expression level (33, 42, 43), the degree of conservation (18,

79    32, 44, 45), the genomic location—*i. e.*, chromosome, chromid, plasmidome (19, 46,

80    47)—and different features such as codon ramps, and mRNA secondary structure,

81    among others (48–50). Current evidence indicates that accessory genes involve atypical

82    codon usages (21, 47, 51, 52) compared to the more conserved (ancestral) core genes in

83    a given lineage. These latter genes, for their part, exhibit adaptational variations in codon

84    usages ranging from typical to more biased as the one observed in genes that

85    correspond to highly abundant proteins which are coded by PHE (53). Moreover, that

86    core genes may also exhibit remarkable codon-usage heterogeneities has been recently

87    demonstrated (19).

88        In the work reported here after examining 29 different prokaryote families, we

89    performed a consolidated analysis aimed at characterizing the specific intragenomic

90    codon variations that lead to differences in codon usage between gene sets with diverse

91    expression levels and degrees of conservation in a given lineage. The evaluation of

92    intragenic regions with different coding characteristics—compared to strategies based on

93    the global analyses of complete genes—enabled the recognition of different patterns of

94    codon usages within a message to be translated. Thus the questions emerged of (i)

95    whether the codon-usage patterns associated with highly expressed amino-acid

96    sequences (*i. e.*, affecting efficiency) were the same as those associated with genes

97    encoding highly conserved sequences (*i. e.*, affecting accuracy), and (ii) whether the

98    requirements for translation efficiency and accuracy were fully independent or whether

99    those two types of demands interacted. The results have indicated how, even in

100   organisms with quite different GC contents, alterations in specific codons are associated

101   with a selective adaptation of the more ancestral genes compared to the adaptation of

102   those genes that are newer in the phylogeny. Through an independent analysis of

103   sequences associated with variable or conserved regions having different expression

104   levels (*i. e.*, low *versus* high), we were able to identify the specific codon usages

105   associated with translation efficiency and accuracy as well as quantitatively estimate their

106   relative relevance to codon usage.

107 **MATERIALS AND METHODS**

108

109 **Prokaryote families selected for analysis in this work and identification of core**

110 **genes and singletons**

111 We screened the EDGAR public-projects database (54, 55) available at

112 https://edgar.computational.bio.uni-giessen.de/cgi-bin/edgar_login.cgi, chose several

113 prokaryote families that included at least 20 complete genomes each, and finally selected

114 27 bacterial and 2 archaeal families (Table S1_a, tab 1). A specific core-gene set was

115 defined as a group of genes whose orthologs are present in a given set of species under

116 investigation. For each of the families selected sequential core-gene sets with increasing

117 ancestry (C1 through Cn) were calculated. To that end, first the phylogenetic tree for

118 each family was extracted from EDGAR and one species per family chosen as a

119 reference. Next, the different core-gene sets were obtained by incorporating into the

120 analysis new species having sequentially increasing phylogenetic distances from the

121 reference strain (accordingly, by following the tree from the branches to the root). Table

122 S1_a-c, lists the phylogenetic trees used for these calculations as well as the particular

123 species that were included in each core-gene set (C1 to Cn) for the different prokaryote

124 families. The phylogenetic trees were edited with the Figtree (56) and Inkskape programs

125 (TEAM-Inkscape). At least six core-gene sets differing in size from *ca*. 50 to 100 genes

126 each were calculated per family. In each prokaryote family, the most ancestral core-gene

127 set (Cn) consisted of 100 to 500 orthologs. Table S2, tabs 2 to 30 lists the singletons—

128 those corresponding to genes that were specific to the reference strains with no orthologs

129 within the family—as calculated with EDGAR.

130

131 **PHE genes.**

132 For each of the selected reference genomes, we obtained a set of genes

133 encoding ribosomal proteins and tRNA synthetases (24, 57). Table S2, tab 1 itemizes the

134 PHE genes whose orthologs were obtained and analyzed in each reference genome.

135

**Highly (HEP) and lowly (LEP) expressed proteins within the same core-gene set**

136

137     Integrated expression data from the Protein Abundance Database (PaxDB; (58))

138 were retrieved for the bacterial strains *Yersinia pestis* CO92, *Streptococcus pyogenes* M1

139 GAS, *Campylobacter jejuni* subsp. jejuni NCTC 11168, *Bacillus subtilis* subsp. subtilis str.

140 168, *Bacteroides thetaiotaomicron* VPI-5482, and *Mycobacterium tuberculosis* H37Rv.

141 Assuming that orthologs have comparable expression levels within the same—or closely

142 related—species and using the PaxDB data from the above indicated 6 strains, we

143 inferred putative expression data for the proteomes of the microorganisms presented in

144 Figs. 4 to 6, and listed in Table S3. Then, for selected core fractions, we obtained one

145 subset of genes encoding HEP plus another subset codifying LEP. For 23 out of the 29

146 prokaryotic genomes that we analyzed here, no proteome data were available, nor were

147 any in phylogenetically related microorganisms.

148

**Analysis of codon usage in gene-sequence regions that encode either conserved**

149

**(*cr*) or variable (*vr*) amino-acid positions in the HEP and LEP subsets**

150

151     Individual genes that belonged to the HEP and LEP groups were aligned with the

152 corresponding orthologs. Then codons corresponding to conserved and variable amino-

153 acid positions in the HEP genes were separated and each concatenated to generate the

154 HEP_*cr*, HEP_*vr* sequence groups. Through the use of a similar procedure with the LEP

155 genes, the LEP_*cr*, and LEP_*vr* sequences were also generated. Codons categorized as

156 belonging to the *cr* group were those associated with positions with fully conserved amino

157 acids throughout the alignment. Codons categorized as belonging to the *vr* group were

158 those associated with positions where none of the amino acids aligned (at that specific

159 point) reached a proportion higher than 0.5. The modal codon usage (47) of each

160 collection of *cr* and *vr* sequences were calculated and used for further analysis.

161

**Raw codon counts (RCC)-based Correspondence analyses (CAs).**

162

163    The RCC-based CAs were performed using bash and R-software homemade scripts

164    which can be found at CUBES software page (this work, available at

165    https://github.com/maurijlozano/CUBES). Briefly, G. Olsen codon usage software was

166    used to count codons on coding sequences (available at

167    http://www.life.illinois.edu/gary/programs.html), data were loaded on R, and the

168    correspondence analyses were run using FactoMiner (59) and Factoextra

169    (https://CRAN.R-project.org/package=factoextra) packages. Plots were made using

170    ggplot2, ggrepel, ggthemes and gtools R packages. For each core-gene set the CA

171    coordinates were calculated as the arithmetic mean of the first and second dimensions of

172    all the genes present in that set (centroids). Then, a plot was generated containing all the

173    coding sequences, together with the projections of the core-gene sets (C1 to Cn), the

174    singletons and PHE genes.

175

176    **Relative synonymous-codon usages (RSCUs)-based CAs, and calculation of modal**

177    **codon usages.**

178          The RSCUs (60) of all individual genes from a given genome were calculated by

179    CodonW with DNA sequences as input data (61) and then used to perform the 59-

180    variable–based correspondence analysis (CA)—*i. e.*, with all the codons except those for

181    Met (AUG), Trp (UGG), and the three stop codons (UAA, UAG, and UGA). The modal

182    codon usages (47) were calculated for the core genes, singletons, and PHE genes.

183    Artificial sequences representing modal codon usages (*i. e.*, modal sequences) and the

184    amino-acid composition present in each core fraction (Cn) were generated through the

185    use of a homemade Perl script (calculate_modals2.pl) from the CUBES package. In order

186    to accurately represent the modal codon usage, particularly for synonymous codons from

187    low-abundance amino acids, modal sequences were designed with a length of at least

188    ten thousand codons. These modal sequences were used as an additional input in their

189    respective CAs. CA plots were generated through the use of Ggplot2 program (62) and

190    edited with Inkskape (TEAM-Inkscape).

191

**tRNA-gene–copy number and modal species-specific tRNA-adaptation index (m-tAI)**

The gene-copy number of each tRNA in the different prokaryote species studied here were downloaded from the GtRNAdb server (http://gtrnadb.ucsc.edu), which website uses predictions made by the program tRNAscan-SE (63). For each reference genome, the copy number for the tRNAs and the sequences of all the open reading frames were used to train the $S_{ij}$ weights as previously reported, with that parameter estimating the efficiency of the interaction between the $i$th codon and the $j$th anticodon in a given organism (36, 37). The procedure stated in brief: With a given $n$, and randomly generated $S_{ij}$ starting points—$i. e.$, having values that range between 0 and 1—the tAI was calculated for each coding sequence by means of the tAI package ((36), https://github.com/mariodosreis/tai). Next, the directional codon-bias score (DCBS; (37)) was calculated through the use of the script seq2DCBS.pl (CUBES package). Finally, the Nelder-Mead optimization algorithm from R project was used (instead of the hill-climbing algorithm) to search for the $S_{ij}$ value that maximized the Spearman rank correlation between the DCBS and the tAI. A script for bulk $S_{ij}$ estimation is available in the CUBES package (https://github.com/maurijlozano/CUBES, calculate_sopt_DCBS_GNM_f.sh). The estimated sets of $S_{ij}$ weights were used to calculate the modal tRNA-adaptation index (m-tAI) for different species and gene sets ($i. e.$, core and PHE genes plus singletons) as a measure of their efficiency in being recognized by the intracellular tRNA pool. The m-tAIs were calculated from the previously generated modal sequences by means of the tAI_Modal_g.sh script from the CUBES package.

214    **RESULTS**

215

216    **Ancestry-dependent codon-usage bias as revealed by the analysis of core genes**

217    **from diverse prokaryotic families**

218         López et al. (19) have recently demonstrated that, in a model proteobacterium,

219    the more ancestral the core genes were the better adapted their codon usages were to

220    the translational machinery. In order to investigate if such correlation was associated with

221    a general phenomenon in different prokaryote taxa, we assembled different core-gene

222    sets that progressed deeper into the phylogenies of 27 Gram-negative and -positive

223    eubacterial families spanning the phyla Proteobacteria, Actinobacteria, Firmicutes, and

224    Bacteroidetes along with 2 archaeal families from the phylum Euryarchaeota. Table S1_a

225    (tab 1) itemizes for each taxon the number of genes in each gene set from the most

226    recent core 1 (C1), to the most ancestral core n (Cn). The codon-usage variation with

227    gene ancestry within a given prokaryote family was evaluated through a correspondence

228    analysis (CA) that used as variables the raw codon counts (RCC) of the individual genes

229    in each genome analyzed (see Materials and Methods). The average values of the first

230    two components for the core-gene sets C1 to Cn were projected on the CA plots. Fig. 1

231    (left panels) depicts the CAs for four different genomes specifically selected to represent

232    groups of organisms with different types of CA plots and GC contents, namely Groups A

233    to D. CAs were also calculated using RSCUs as input variables instead of RCC as

234    presented in Fig S1A. In agreement with a recent study in *Sinorhizobium meliloti* (19), in

235    all instances a directional shift in the codon usage positions was evident from the most

236    recent (C1) towards the most ancestral (Cn) core-gene set. That this ancestry-dependent

237    pattern of codon-usage variation had been observed in even quite distant prokaryote

238    families among those analyzed here was remarkable (*cf*. the CA plots for all other

239    species in Fig. S1B, left panels). In the evolution of core codon usages, however, the

240    extent of the observed shifts and the type of synonymous codons enriched in each taxon

241  (*i. e.*, the direction of change) varied markedly among different families (Fig. 1, Fig S1A

242  and Fig S1B, right panels).

243      The general features that characterized the bias in codon usages can be

244  summarized as follows. First, a general pattern indicated that in bacteria from Groups B,

245  C and D the PHE genes are enriched in codons with higher GC3 when compared with

246  singletons (Fig. 1 and Fig. S1, right panels). Conversely, an AU enrichment in the third

247  position of codons was observed in the ancestral core fractions of organisms from Group

248  A which have extremely low GC contents. Second, from C1 to Cn in the CA plot, the

249  codon usages gradually shifted away from the position of the singletons (the unique

250  genes) to approach the region where the PHE genes were located (Fig. 1 and Fig S1, left

251  panels). Similar results were obtained when PHE genes were subtracted from the

252  different Cn cores (see Fig. S2). Thus, the overall evidence suggested that gene ancestry

253  correlated with a codon-usage optimization that resembled the one observed in the PHE

254  genes. Nonetheless, the more ancestral core genes (*i. e.*, the Cn gene sets) never

255  overlapped with the position of the PHE genes in the CA plots. In most prokaryote

256  species, the order of positions in the CA plot followed the sequence singletons-C1-Cn,

257  which series was associated with an enrichment in some of the C-ending 2-/3-fold

258  degenerate codon families (*i. e.*, the 2-/3-fold C-bias; *cf*. the distribution of red circles in

259  Fig. 1 and Fig. S1, right panels); whereas the position of the PHE genes compared to that

260  of Cn was characterized by an additional enrichment in U-ending 4-fold degenerate

261  codon families (*i. e.* the 4-fold U bias; *cf*. the distribution of light-blue circles). Each of the

262  previous effects varied in relative intensity among the different prokaryote families, where

263  other specific codon changes (brown circles in Fig. 1 and Fig. S1) also occurred from C1

264  to Cn to PHE and accompanied the above-mentioned 2-/3-fold C, and 4-fold U biases.

265  Wald et al. (41) have previously reported that the C and the U biases are associated with

266  an improved codon-usage correspondence to the anticodons of the tRNA pool. The

267  combined effects of the C and U biases are the basis for the "rabbit head" distribution of

268  genes that we observe in most of the CA plots (gray dots), an effect that was originally

269    described in *E. coli* (21). Contrasting with the codon usage of core and PHE genes, the

270    singleton genes tend to be enriched in A/U-ending codons.

271

272    **Indication from m-tAI values that the codon usages of more ancestral genes are**

273    **better adapted to the cellular translational machinery**

274        In order to explore how extensive the correlations between codon usage, gene

275    ancestry, and translation efficiency were, we calculated the m-tAI values for the C1 to the

276    Cn core genes for a given strain and used those indices to estimate the adaptation of

277    each gene set to the tRNA pool. Each m-tAI takes into consideration both the copy

278    number of each tRNA structural gene as an estimation of that tRNA's cellular

279    concentration and the codon-anticodon interactions including the classical Watson-and-

280    Crick interactions (WCIs) along with the wobble rules (see Materials and Methods).

281    Unfortunately, nonstandard forms of base pairing, such as U:U interactions and others,

282    are not included in the m-tAI calculations. In Fig. 2, the left panels illustrate how with

283    progressive gene ancestry the m-tAI generally increases to often approach that of the

284    PHE genes, thus evidencing that the more ancestral cores are enriched in genes that

285    displayed adaptive—*i. e.*, selection-dependent—changes in their codon usage. That such

286    m-tAI increases with progressive ancestry had been observed in strains from 18

287    prokaryote families (17 eubacteria, 1 archaea) was indeed remarkable (*cf.* Fig. 2 and Fig.

288    S3, left panels a1, a3, a5, a6, a8, a10, and a12 to a19). In the reference strains from

289    these prokaryote families, the PHE genes (red dashed lines) were always associated with

290    higher m-tAI values than those of the core-gene sets from the same genome.

291    Conversely, singletons (blue-dashed lines) were always the gene sets with the lower m-

292    tAIs, thus suggesting that accessory genes (*i. e.*, those present in plasmids, phages, and

293    the unique genes in chromosomes) involve codon usages that—most likely due to their

294    non-essential character—is far from being optimized with respect to the host-translation

295    machinery. Strains with the characteristics described above have genomes with quite

296    diverse GC contents, ranging from *ca.* 30% to over 70%. Exceptions to the general

297    increase in the m-tAI values with ancestry are likely due to m-tAI deficiencies to

298    quantitate non-standard codon-anticodon interactions (i.e. those different from WCIs, and

299    wobble base pairing) (36).

300

301    **Effect of codon optimization on the GC content**

302    An analysis of the prokaryote genomes with different GC contents enabled us to

303    explore how the GC composition at the third base of codons (*i. e.*, the GC3) changed in

304    the core-gene sets over ancestry, and to compare the results with the GC3 in PHE genes

305    and singletons. Since the first two positions in codons are constrained by the protein-

306    coding information, most of the GC changes result in variations in synonymous codons

307    (2). As we have seen in the two previous sections, core genes adjust their codon usages

308    in the direction of the PHE genes (Fig. 1 and Fig. S1, left panels) in order to improve

309    translation (Fig. 2 and Fig. S3, left panels). The question thus became raised as to how

310    bacteria with different GC contents changed their GC3 composition in the process of

311    adapting their codon usage. The results presented in Fig. 2 and Fig. S3 (right panels)

312    show that changes in GC3 in genomes from Groups A to D each follows a distinctive

313    pattern from singletons-to-Ci-to-PHE. Whereas in genomes that belong to Group A

314    (overall GC content lower than ca. 35%) the GC3 decreases from singletons to Ci to PHE

315    (*cf.* Fig. 2, panel b1), in the genomes included in Group C the GC3 either increases from

316    singletons to Ci-to PHE (*cf.* Fig. 2, panel b3) or plateaus in Ci to PHE at a high level (*cf.*

317    Fig. S3, panel b17). In contrast, genomes pertaining to Group B; exhibited a biphasic

318    pattern with an initial GC3 increase from the level of the singletons up to the contents of

319    the Ci series (with i varying from 1 to n) followed by a later decrease from the Cn values

320    down to those of the PHE genes (*cf.* Fig. 2, panel b2). Those changes in the Group-B

321    genomes were reflected in pronounced forward and backward movements in the position

322    of the core genes in the CA plots, first from singletons to Ci and then from Cn to the PHE

323    genes (*cf.* Fig. S1, organisms in Group B). A similar biphasic pattern in the CA plots could

324    also be recognized, though softened, in certain species that were included in Group C or

325　even Group D where the PHE genes did not evidence a decrease in GC3 levels when

326　compared to those of the core genes. The genomes in Group D had extremely high

327　global GC contents and had GC3 values in all their core-gene sets (C1 to Cn) that were

328　comparable—though slightly higher—than the corresponding values in their PHE genes.

329　In the next section we will describe how individual codons for a given amino acid are

330　selected in the more ancestral core-gene sets.

331

332　**Characterization of codons that improve adaptation to the tRNA pool**

333　　　　The variations in the use of individual codons when progressing from the C1 to

334　the Cn gene sets were analyzed in the different prokaryote genomes, together with the

335　tRNA-gene–copy numbers and the codon-adaptation indices (*Wi*(s); *cf*. Materials and

336　Methods). Figs. 3 and S4 illustrate the CUFs (codon-usage frequencies, *cf*. Materials and

337　Methods) for singletons, PHE genes, and core genes with increasing ancestry together

338　with the tRNA-gene copies and the *Wis* (Fig. S5 summarizes the *Wis*, ΔCn-C1, and

339　ΔPHE-Cn in the different genomes studied). In agreement with previous reports (10), our

340　results demonstrated that the CUF values among synonymous codons were strongly

341　influenced by the global GC content in each genome—*i. e.*, codons with G and C in the 3'

342　position ($N_3$) were the more abundant synonymous codons in the GC-rich genomes,

343　whereas A and U become predominant in that position in the genomes with low GC

344　content (Figs. 3 and S4). An inspection of the proportion of codon usage for each amino

345　acid in ancestral cores compared to the more recent ones (curves in Figs.3, S4, and S5)

346　revealed that in most genomes a C-bias enrichment occurred with increased ancestry at

347　the 3' position of the 2-fold pyrimidine-ending codons—for Asp (GAU, GAC), Phe (UUU,

348　UUC), His (CAU, CAC), Asn (AAU, AAC), and Tyr (UAU, UAC)—as well as in the unique

349　3-fold codons for Ile (AUU, AUC, AUA), which three included the pyrimidine-ending pair

350　AUY (Figs. 3, S4, and S5). Corresponding to the observed C bias, in all these examples

351　high *Wi* values (shown in parenthesis in the figure) were observed for the C-ending

352　codons, which triplets were decoded through exact WCIs with the cognate tRNA species

353    (*i. e.*, with the anticodon $G_{34}N_{35}N_{36}$). Because of the absence of tRNA species bearing

354    anticodons $A_{34}N_{35}N_{36}$ for these six amino acids, lower *Wi* values were obtained for the U-

355    ending codons as the consequence of a weaker wobble codon-anticodon non-WCI

356    recognition. Especially noteworthy was the observation that, though to a lesser extent,

357    the bacteria with extremely low GC contents likewise exhibited a C bias in the 2- to 3-

358    fold–codon family, irrespective of a global decrease in the GC3 value, as in the example

359    of *Clostridium beijerinckii* (*cf.* Figs. 2 and 3).

360         In the instance of the 2-fold purine-ending codons—that is GAA and GAG for Glu,

361    AAA and AAG for Lys, and CAA and CAG for Gln—we observed that the codons with G

362    or A in the 3' position were enriched from C1 to Cn and from Cn to PHE (*i. e.,* ΔCn-C1

363    and/or ΔPHE-Cn in Fig. S5, respectively) depending upon which tRNA species

364    (anticodons) were present. In those examples where only the tRNAs bearing the

365    $U_{34}N_{35}N_{36}$ anticodons were present, the cognate A-ending codons recognized by WCIs

366    were the ones that became enriched in the more ancestral core and/or PHE genes (*cf.* in

367    Fig. S5, the GA<u>A</u> triplet for Glu in *C. violaceum*, *P. graminis*, *Bacillus subtilis*, *Bordetella*

368    *holmesii*, and *Leisingera methylohalidivorans*; the AA<u>A</u> for Lys in *M. smithii* and *Bacillus*

369    *subtilis*; and the CA<u>A</u> for Gln in *M. smithii, Streptococcus equii*, and *B. subtilis*).

370    Accordingly, these 3'-A–ending codons were associated with higher *Wi* values than the

371    corresponding codons ending in G, as the latter were recognized only by wobble-base

372    pairing (*i. e.*, $G_3$-$U_{34}$ interaction). In a second circumstance, where both tRNA species for

373    the same amino acid (*i. e.*, those bearing anticodons $U_{34}N_{35}N_{36}$ or $C_{34}N_{35}N_{36}$) were

374    present, a more frequent enrichment in G–ending codons was observed (with few

375    exceptions) since such codons can be decoded by either Watson-Crick or wobble

376    interactions with the tRNA anticodons $C_{34}N_{35}N_{36}$ or $U_{34}N_{35}N_{36}$, respectively. In those few

377    examples where the A-ending were more enriched than the G-ending codons, a higher

378    copy number of the tRNA genes was always observed with anticodons $U_{34}N_{35}N_{36}$ than

379    that obtained with the anticodons $C_{34}N_{35}N_{36}$ (*cf.* in Figs. 3 and S4, the GA<u>A</u> triplets for Glu

380    in *Bacteroides vulgatus* and *C. beijerinckii*; the AA<u>A</u> triplets for Lys in *S. multivorans*; and

381    the CA<u>A</u> triplets for Gln in *C. beijerinckii* and *S. multivorans*).

382          A different codon usage bias—in a pattern not found in the 2-/3-fold–degenerate

383    amino acids—was observed in codons encoded by 4-fold–degenerate amino acids (Val,

384    Thr, Pro, Gly, Ala) or by the 4-fold boxes of the 6-fold degenerate amino acids (Ser, Leu,

385    Arg). In these 4-fold groups a U-bias enrichment (*i. e.*, a NNU-codon enrichment) was

386    observed in the PHE genes from most of the genomes irrespective of their GC content

387    (Figs. 3, S4, and Fig. S5). This enrichment in U-ending codons, previously reported as a

388    U bias (41), could not be explained by WCIs with $A_{34}N_{35}N_{36}$ tRNAs because these latter

389    species were not present in prokaryotes, except in the example of Arg. The observed U

390    bias likely occurred through the previously proposed nonconventional codon-

391    $U_3$:anticodon-$U_{34}$ interaction that was known to exist in bacteria (64). The presence of

392    $U_{34}N_{35}N_{36}$ tRNA species might then lead to an increase in both NNA and NNU codons as

393    a consequence of positive WCIs and $U_3$-$U_{34}$ interactions, respectively.

394          All the codon adaptations that we have described in this section referring to core

395    genes proved to be more prominent in the PHE genes, whose triplets were even better

396    adapted to the translational machinery. Contrasting with such a strong pattern of

397    selection-associated codon bias, the singletons displayed codon usages that were in

398    general the most distant from those observed in the PHE genes (as exemplified in the

399    CUFs in Figs. 3, S4, and in the CA plots from Figs. 1 and S1). These observations are

400    also in agreement with variations in the m-tAIs for the different gene sets presented in the

401    previous section.

402

403    **Search for coding signatures for translation efficiency and accuracy: Codon-usage**

404    **profiles associated with sequences encoding highly-expressed–variable (HEP_*vr*)**

405    **and -conserved (HEP_*cr*) translated domains**

406          Expression level and amino-acid–sequence conservation are both parameters

407    that positively correlate with codon-usage optimization (65). Nevertheless, the relative

408    relevance of efficiency and accuracy to translation plus the way in which either one of

409    those parameters affects the other have not yet been investigated in detail. A central

410    limitation that made such studies difficult was associated with the natural genomic

411    heterogeneity in gene ancestry along with the expression level and the sequence

412    conservation (structural constraints) in the translated products. In order to reduce the

413    degrees of freedom in the analysis; for each of six different bacterial species, we created

414    two distinct gene sets based on the experimental proteome data. One of those gene sets

415    consisted of genes encoding proteins with the highest expression levels in the cell (*i. e.*,

416    the HEP), while the other was associated with proteins with low cellular abundance (*i. e.*,

417    the LEP). Then, the conserved (*cr*) and variable (*vr*) sequences among the orthologs

418    were collected from each individual gene (*cf*. Materials and Methods), the corresponding

419    HEP_*cr*, HEP_*vr*, LEP_*cr*, and LEP_*vr* modal codon usages were used to calculate the

420    relative distances illustrated in the neighbor-joining tree presented in Fig. 4. In five out of

421    the six species present in the trees (Fig. 4A to 4E), the HEP_*cr* and HEP_*vr* sequences

422    separated from those of the singletons, the core genes, and all the LEPs as the result of

423    a strong codon-usage adaptation (also reflected in the low effective number of codons

424    (Ncs) associated with the HEPs, indicated in parentheses following labels in the tree).

425    Furthermore, the large distance in the tree between HEP_*cr* and LEP_*cr* (where both

426    sequences encode regions with conserved amino acids, but with different expression

427    levels) compared to the much shorter distance between HEP_*cr* and HEP_*vr* (where both

428    encode highly expressed products with different degrees of conservation) pointed to the

429    quantitatively stronger effect of efficiency over accuracy in shaping codon-usage bias.

430    However, codons that were optimized as a result of accuracy under high and under low

431    expression—*i. e.*, [HEP_*cr*–HEP_*vr*] and [LEP_*cr*–LEP_*vr*], respectively, labelled ***A*** for

432    accuracy at the bottom of Fig. 5—were highly coincident with the codons that were

433    optimized through efficiency—*i. e.*, [HEP_*cr*–LEP_*cr*] and [HEP_*vr*–LEP_*vr*], labelled ***E***

434    for efficiency. In some organisms, the greater distance between HEP_*cr* and HEP_*vr*

435    than between LEP_*cr* and LEP_*vr* (Fig. 4) indicates a stronger influence of accuracy in

436   codon-usage optimization when operating under high-expression conditions, thus

437   pointing to an interaction between the simultaneous requirements of high fidelity and

438   efficiency. The more relevant contributions to the global difference in codon usage

439   between HEP and LEP resulted to be efficiency (both in conserved and in variable

440   regions) (*E* columns in Fig. 5) followed by accuracy under high expression (first *A* column

441   in Fig 5)(the stronger the contribution of each factor, either *E* or *A*, the shorter the

442   distance in brackets to HEP-LEP in the figure). The heat maps display the complete

443   profiles of preferred codons for sequences requiring high translational accuracy and/or

444   efficiency (protein demands). As expected, the preferred codon for each amino was in

445   agreement with the C and U bias and the tRNA-copy number described in the previous

446   sections. In light of these results, the highly-expressed variable and conserved domains

447   (HEP-*vr*/*cr*) constitute the basis for explaining the observed codon-usage optimization in

448   the more ancestral core-gene sets (Cn), which concentrate HEPs (Table S3). Fig. 6

449   illustrates that HEP sequences (red dots) are those under the highest selective pressure

450   to optimize codon usage because of both their expression level and their degree of

451   conservation.

452　**DISCUSSION**

453　　　　Since gene adaptation to a host cell is expected to be associated with an

454　improved codon selection for translation efficiency and accuracy (43, 66), we investigated

455　correlations between core-gene ancestry and their modal codon usage within a given

456　prokaryote family. In order to ascertain if the adaptation of the more ancestral core genes

457　was an extensive phenomenon among prokaryotes, we analyzed core modal codon

458　usages in 27 different species of Bacteria and 2 of Archaea. That in the CA plots the

459　more ancestral core genes had been the ones with the closest location to the PHE genes

460　in all families was remarkable and strongly indicated a core codon-usage adaptation that

461　likely operated to improve translation. In agreement with the position of the different gene

462　sets in the CA plots, the m-tAI values served to confirm that the PHE genes were the best

463　adapted gene set, followed by the Cn to C1 core genes, in that order, and finally by the

464　singletons, with those being the least adapted genes with the lowest m-tAIs in the

465　genome. Studies made by others have previously demonstrated that the level of gene

466　expression together with the need to preserve accuracy during the translation of

467　conserved amino-acid regions are both among the main parameters that govern codon-

468　usage selection (65). The bioinformatics isolation of conserved (*cr*) and variable (*vr*)

469　coding-sequence domains from genes under high- (HEP) and low-expression (LEP)

470　regimes served in this work to ascertain quantitatively the relative contribution of

471　efficiency (expression level) *versus* accuracy during the selection-based codon-usage

472　optimization. According to the observed neighbor-joining distances (Table S3 worksheet

473　"distances", and tree in Fig. 4), changes in codon usage derived from differences in gene-

474　expression levels—*i. e.*, the efficiency in terms of the distance from the LEP to the HEP—

475　were between 1.25 to 2.35 times greater than the changes in codon usage resulting in

476　increased accuracy—*i. e.*, the distance from *vr* to *cr*—. The increasing proportion of

477　highly expressed-variable and specially -conserved sequences (*i. e.*, HEP_*vr* and

478　HEP_*cr*) in the more ancestral gene sets constituted the basis for explaining the

479  corresponding high degree of codon-usage optimization that gradually increased from C1

480  to Cn.

481      The central question therefore was how adaptive changes in codon usage—which

482  alterations become reflected in m-tAI values—occurred in prokaryotes with quite diverse

483  GC contents (10). Because of the small amount of intergenic DNA in prokaryotes,

484  genomic differences in base composition must mainly derive from changes in the coding

485  regions. Within the alterations in the open reading frames, changes in GC are

486  preferentially associated with modifications in the GC3, and only to a lower extent with

487  alterations in the GC content of the first two codon positions (2, 4). How mutational bias

488  (12) competes with selection (15) to drive all these changes is not yet fully understood.

489  The codon-usage biases described here were associated with the four different patterns

490  of GC3 changes summarized in the schemes presented in Fig. 7 (*i. e.* the genome

491  Groups A, B, C, and D depicted in the figure). The Group-A genomes, those having an

492  extremely low GC content and with their GC3 frequency decreasing from C1 to Cn,

493  proved to have only the tRNA-$U_{34}$ to recognize 4-fold synonymous codons in one or more

494  amino acids. In such instances, the observed core-gene AT enrichment over ancestry

495  appeared to be directly affected by selection (as with the PHE genes), where codons

496  NNA (via WCIs with the tRNA-$U_{34}$) and NNU (via nonconventional U-U interactions) were

497  preferentially enriched over NNC and NNG codons. Though both of those changes were

498  probably related to improvements in translation efficiency, such increases are not always

499  reflected in the m-tAIs since, as stated earlier, U-U interactions are not considered in the

500  calculation of that index. Unfortunately, when we (not shown) and others (37) have

501  attempted to improve the tRNA-adaptation index by including additional nonstandard

502  base pairings, we obtained no better results. Nonetheless, under the assumption that the

503  PHE genes are the best adapted to the translational machinery, in genomes with

504  extremely low GC content—such as those belonging to Group-A—the observed AT3

505  enrichment from C1 to Cn to PHE (Fig. 7, right panel) should mainly result from selection.

506  According to Hildebrand et al. (15), the mutational processes in very low-GC organisms

507   favor a GC3 enrichment. That the core and PHE genes in bacteria that belonged to

508   Group A had been selected to bear lower GC3 values than singletons in order to improve

509   translation in view of the previous pattern of increasing GC content was remarkable, with

510   this circumstance being a result of the above-mentioned enrichment in NNA and NNU

511   triplets compared to their proportion in the synonymous codons (Fig. 7, right panel). In

512   Group-B genomes, the biphasic pattern observed from singletons to PHE genes could be

513   explained by an initial increase in GC3-rich codons from singletons to core genes,

514   followed by a later U bias from core genes to PHE genes. That initial GC3 enrichment

515   followed by a U3 increase was sufficient to explain the basis of the previously reported

516   "rabbit head" distribution of codon usages that characterizes most prokaryote genomes

517   (21, 67). What should be also especially noted is that the PHE genes separated from the

518   Cn (in both the CA, and the GC3 plots) because of a much more intense U bias likely

519   associated with the difference in expression levels between the two gene sets. In the

520   type-C genomes, in which the GC3 always increased, the absence of a strong U bias

521   from the Cn to the PHE genes led to a less pronounced—*i. e.*, more linear—"rabbit-head"

522   distribution of genes in the CA plot. In addition to that general trend, *Yersinia*

523   *enterocolítica*, *Methanolacinia petrolearia*, and *Sphingomonas parapaucimobilis* could be

524   considered as having an intermediate behavior between the bacteria in Group C and

525   those in Group B. Finally, the Group-D genomes, which had extremely-high GC contents,

526   were the most restricted with respect to GC3 variations. The quite small compositional

527   variations in that group of genomes became apparent in the compacted location of the

528   different core and PHE genes in the CA plots. What was remarkable is that in Group-D

529   genomes a U bias (though much less intense than in the genomes of Groups A, B, and

530   C) was still a visible variable that differentiated codon usages between the core and the

531   PHE genes. As stated above, the noninclusion of U:U interactions in the m-tAI calculation

532   limited the use of this parameter to expresses the translational adaptation of those gene

533   sets in which a U bias was dominant. Pouyet et al (11) present a model to predict and

534   separate the relative contribution of mutational bias (N-layer), codon selection (C-layer),

535    and amino acid composition (A-layer) on the global GC and the GC3 content. Our

536    analysis is fully consistent with the results reported by Pouyet et al (11) where the C-layer

537    (codon selection / translational selection) has a stronger influence on the GC3 of genes

538    with low effective number of codons (Nc)(such as Cn and PHE) compared to the

539    influence on genes with the highest Nc (such as the C1).

540         The results presented here together with previous evidence from other authors

541    have enabled a comprehensive analysis of the principal basis underlying the changes

542    associated with the optimization of codon usage in prokaryotes in different gene sets and

543    in organisms with different GC contents. As stated previously, the overall codon usage is

544    known to be constrained by genome-wide mutational processes (7, 8, 10) that are

545    considered as a main force in shaping the global GC content. The intragenomic codon

546    usage will concurrently become accommodated through selection-driven processes, as

547    has also been extensively reported (35, 42, 48, 68). In order to further our knowledge of

548    the relevance of those factors/forces generating intragenomic variations, we investigated

549    the different nucleotide-base changes underlying the selection of preferred codons in the

550    core and PHE genes of representative prokaryote species. The analysis of gene sets with

551    different expression levels and degrees of conservation in organisms with diverse global

552    GC contents enabled a description of how core codon usage approaches that existing in

553    the PHE genes and how nucleotide changes correlate with an improved compatibility

554    between the genes and the coexisting tRNA pool. That C- and U-ending codons in 2-/3-

555    fold- and in 4-fold–degenerate amino acids, respectively, were specifically enriched as a

556    result of selection to improve translation has been previously reported for different

557    prokaryotic genomes (41). Using separate analyses focused on different gene sets, we

558    demonstrated here that similar selection-driven adaptations in codon usage has taken

559    place from singletons to core genes to PHE genes. The intensity and relevance of the C

560    and U bias was dependent on the particular genome—and especially on the genomic GC

561    content—as well as on the gene fractions under consideration. In contrast to the codon-

562    usage variations occasioned by selection in the core and PHE genes, the singletons

563     constituted the gene set characterized by both the lowest GC3 content as a result of the

564     AT mutation that is universally biased in prokaryotes (12) and a much more relaxed

565     selection than that of the more ancestral genes, with the sole exception of the extremely

566     low-GC–containing genomes of Group A. In addition to a description of the basic

567     changes that together conform the intracellular-codon–usage variations, further

568     investigation should be focused on the analysis of the time course required by the newly

569     acquired information to be properly incorporated into the genetic language of the host cell

570     (codon usage tuning).

571

572     **ACKNOWLEDGEMENTS**

573

584 **LEGENDS TO FIGURES**

585

586 **Fig. 1. Raw codon counts-based Correspondence Analysis (CA) plots of core-gene**

587 **sets with different degrees of conservation throughout the phylogeny of selected**

588 **prokaryote families (Groups A to D). Panels a1 to a4:** In 4 reference strains with

589 different GC contents, individual genes (in gray) are represented in the space of the first-

590 two CA components, with the percent variation of components 1 and 2 being indicated on

591 the axes. CAs were computed using raw codon counts (RCC) as the input variables.

592 Average coordinates (centroids) for different gene sets (i.e. singletons in blue, C1 to Cn

593 in a gradient from blue to red, and PHE in red) were projected on the CA space. In C1 to

594 Cn the higher number denote a more ancestral core-gene set within the phylogeny. Table

595 S1_a (tab 1) lists the prokaryote species that were used to construct each Ci gene set by

596 means of the EDGAR software (54, 55). **Panel b1 to b4**: Plots describing codon relative

597 weight in the first two principal-component positions of the CA. Codons with the highest

598 CUF enrichment for each amino acid from C1 to PHE (i.e. those codons that better

599 represent translational adaptation) were colored in brown, except when those same

600 codons corresponded also to a C- or to a U-bias in which cases they were colored in red

601 and light blue, respectively.

602

603 **Fig. 2. Codon-usage adaptations to the cellular tRNA-pool, and changes in the GC3**

604 **content of different prokaryote core genes.** The reference strains represented here

605 are the same five as in Fig. 1. **Panels a1 to a4:** In each panel, the modal tRNA-

606 adaptation index (m-tAI) calculated for each of the Ci gene-sets as described in Materials

607 and Methods is plotted on the *ordinate* as a function of the evolutionary distance

608 indicated on the *abscissa* (Table S1_a, tab 2) as inferred from the corresponding

609 phylogenetic trees included in Table S1_a-c. Higher values of m-tAI indicate an

610 enrichment in the codon usage frequencies of those synonymous codons better adapted

611 to the host-cell tRNA pool. The C1 to Cn gene sets plotted here are the same as those

612    presented in Fig. 1. The red and blue horizontal dashed lines correspond to the

613    respective m-tAI values calculated for the PHE genes and the singletons. **Panels b1 to**

614    **b4:** In each panel, the average GC3 content in each core-gene set of increasing ancestry

615    is plotted on the *ordinate* as a function of the evolutionary distance indicated on the

616    *abscissa* as in panels a1 to a4. The PHE genes and the singletons are represented as

617    red and blue horizontal dashed lines, respectively.

618

619    **Fig. 3. Codon usage frequencies and adaptation indices (*Wi*s) of the gene sets**

620    **analyzed in this work, together with the tRNA-gene–copy numbers for strains of**

621    **the four reference Groups A to D.** For the amino acid denoted by the corresponding

622    single-letter identification code located above each panel, the change in the modal

623    codon-usage frequencies (CUFs; see Materials and Methods) of the core-gene sets with

624    increasing ancestries (left to right, the C1 to Cn), the PHE genes, and the singletons are

625    plotted in the upper panels as solid horizontal curves for each of the indicated codon

626    triplets between the two vertical broken lines, for the singletons to the left of the first of

627    those lines, and for the PHE genes to the right of the second (with singletons and PHE

628    genes being located at the beginning and the end of the curves, respectively). The CUFs

629    are represented by different colors with the associated absolute codon-adaptation index

630    (*Wi*, (36)) being denoted within parentheses beside each triplet. Finally, the presence and

631    gene-copy number (N tRNA) of the cognate tRNA species of a given synonymous codon

632    bearing the exact complementary anticodon is depicted with a number and a bar of

633    proportional height in the lower panel in the same color as the corresponding triplet and

634    curve in the upper panel.

635

636    **Fig. 4. Neighbor-joining–distance trees of different gene sets encoding** HEP, LEP,

637    **and their associated conserved (*cr*) and variable (*vr*) regions based on the**

638    **corresponding modal codon usage.** Modal codon usage-based neighbor-joining trees

639    were constructed for the indicated gene sets and intragenic regions (*cr* and *vr*) following

640    the method described by Karberg et al. (17) along with the neighbor-joining program of

641    the Phylip package (69). Phylogenetic trees were drawn through the use of the Figtree

642    application (70). The figure abbreviations are as follows: C1 to Ci, core-gene sets with

643    increasing ancestry; single, singletons; HEP, genes encoding proteins with the highest

644    expression level; LEP, genes encoding proteins with the lowest expression level;

645    HEP_cr, conserved HEP sequences (dark red); HEP_vr, variable HEP sequences (light

646    red); LEP_cr, conserved LEP sequences (dark blue); and LEP_vr, variable LEP

647    sequences (light blue). HEP and LEP cr and vr subfractions were recovered as indicated

648    in Materials and Methods through the use of the polypeptide sequences included in C13

649    for *Yersinia enterocolitica* subsp. *palearctica* Y11, C10 for *Streptococcus equi* ATCC

650    33398, C8 for *Sulfurospirillum multivorans* DSM 12446, C9 for *Bacillus subtilis* subsp.

651    *spizizenii* TU B 10, C6 for *Bacteroides vulgatus* ATCC 8482, and C12 for *Mycobacterium*

652    *fortuitum* subsp. *fortuitum* DSM 46621 (ATCC 6841). The effective number of codons

653    (Nc$_s$) as previously defined by Wright (71) are indicated in brackets for the *cr* and *vr*

654    subset of sequences.

655

656    **Fig. 5. Heat-map representation expressing differences in modal-codon-usage**

657    **profiles between the indicated gene sets.** The color scale from red to blue indicates

658    the relative level of use of each particular codon in a gene set compared to that of

659    another (*i. e.*, gene set 1 *versus* gene set 2). The blue color corresponds to the dominant

660    use of a particular codon in Gene set 1 over the use of the same codon in Gene set 2

661    (and *vice versa* for the red color). Amino acids are indicated in the standard three-letter

662    code. The heat map was constructed through the use of the phytools R package (72).

663    The abbreviations are as follows: HEP, genes encoding proteins with the highest

664    expression level; LEP, genes encoding proteins with the lowest expression level;

665    HEP_cr, conserved HEP sequences; HEP_vr, variable HEP sequences; LEP_cr,

666    conserved LEP sequences; and LEP_vr, variable LEP sequences. The [HEP (gene set 1)

667    – LEP (gene set 2)] column represents the profile of the optimized codons when

668    comparing the coding strategies in high- *versus* low-expression genes (*i. e.*, reflecting

669    differences in their modal codon usages). The columns indicated by ***A*** correspond to the

670    profiles of codons optimized as a result of accuracy (*i. e.*, differences between [HEP_*cr* –

671    HEP_*vr*] and [LEP_*cr* – LEP_*vr*]). The columns indicated by ***E*** correspond to the profiles

672    of optimized codons through high expression (*i. e.*, reflecting differences in efficiency

673    between [HEP_cr – LEP_*cr*] and [HEP_*vr* – LEP_*vr*]). The numbers in brackets indicate

674    the extent to which changes induced by either efficiency or accuracy approach the

675    differences in codon usage between HEP and LEP (*i. e.*, distances from each column to

676    the column HEP – LEP). The shorter the distance in brackets the stronger the

677    contribution of the indicated factor (*i. e.*, A = accuracy or E = expression level) to codon

678    optimization in the HEP.

679

680    **Fig. 6. Amino-acid–sequence conservation in proteins with different cellular**

681    **abundances.** The amino-acid–sequence conservation calculated for proteins of the

682    indicated bacterial species and core fractions (Materials and Methods) are plotted on the

683    *ordinate* as a function of the logarithm of the corresponding protein abundance (logPA)

684    on the *abscissa*. The red and blue dots correspond to HEP and LEP, respectively, with all

685    the other proteins of the same core represented in gray. The linear regressions and

686    graphs were all performed with the ggplot2 library from the R package.

687

688    **Fig. 7. Schematic representation of general codon-usage patterns observed in**

689    **different prokaryote families.** For the prokaryote strains whose genomes were

690    classified as belonging to Groups A, B, C, and D and listed to the right of each panel,

691    cartoons with the associated correspondence-analysis (CA) and GC3-variation pattern

692    among the core-gene sets of increasing ancestry (light gray) are presented, along with

693    the corresponding singletons (blue) and PHE genes (red). The light-blue arrow indicates

694    the direction of the U bias and the red arrow that of the C bias. The right panel is a plot of

695    the GC3 content on the *ordinate* as a function of increasing evolutionary distance on the

696    *abscissa* with the red horizontal dashed line denoting the PHE genes and the blue the

697    singletons.

**REFERENCES**

698

699   1.   Sueoka N. 1962. On the genetic basis of variation and heterogeneity of DNA base

700        composition. Proc Natl Acad Sci.

701   2.   Sueoka N. 1988. Directional mutation pressure and neutral molecular evolution.

702        Proc Natl Acad Sci.

703   3.   Osawa S, Ohama T, Yamao F, Muto A, Jukes TH, Ozeki H, Umesono K. 1988.

704        Directional mutation pressure and transfer RNA in choice of the third nucleotide of

705        synonymous two-codon sets. Proc Natl Acad Sci.

706   4.   Sueoka N. 1992. Directional mutation pressure, selective constraints, and genetic

707        equilibria. J Mol Evol.

708   5.   Sueoka N. 1995. Intrastrand parity rules of DNA base composition and usage

709        biases of synonymous codons. J Mol Evol.

710   6.   Sueoka N. 1999. Two aspects of DNA base composition: G+C content and

711        translation- coupled deviation from intra-strand rule of A = T and G = C. J Mol Evol.

712   7.   Knight RD, Stephen J, Landweber L. 2001. A simple model based on mutation and

713        selection explains trends in codon and amino-acid usage and GC composition

714        within and across genomes 1–13.

715   8.   Chen SL, Lee W, Hottes AK, Shapiro L, Mcadams HH. 2004. Codon usage

716        between genomes is constrained by genome-wide mutational processes.

717   9.   Karimpour I, Cutler M, Shih D, Smith J, Kleene K, Hill KE, Lloyd RS, Yang JG,

718        Read R, Hurk RF, Burk RF, Lawrence RA, Lane JM, Rurk RF, Rellew T, Morrison-

719        Plummer J, Bellew T, Palmer IS, Berry MJ, Banu L, Chen Y, Mandel SJ, Kieffer

720        JD, Harney JW, Larsen PR. 1993. Codon usage: mutational bias, translational

721        selection, or both?Proc. Natl. Acad. Sci. U.S.A.

722   10.  Hershberg R, Petrov DA. 2009. General rules for optimal codon choice. PLoS

723        Genet 5.

724   11.  Pouyet F, Bailly-Bechet M, Mouchiroud D, Guéguen L. 2016. SENCA: A

725        multilayered codon model to study the origins and dynamics of codon usage.

726      Genome Biol Evol.

727   12.   Hershberg R, Petrov DA. 2010. Evidence that mutation is universally biased

728      towards AT in bacteria. PLoS Genet 6.

729   13.   Bohlin J, Eldholm V, Brynildsrud O, Petterson JH-O, Alfsnes K. 2018. Modeling of

730      the GC content of the substituted bases in bacterial core genomes. BMC

731      Genomics 19:589.

732   14.   Bulmer M. 1991. The selection-mutation-drift theory of synonymous codon usage.

733      Genetics.

734   15.   Hildebrand F, Meyer A, Eyre-Walker A. 2010. Evidence of selection upon genomic

735      GC-content in bacteria. PLoS Genet 6.

736   16.   Raghavan R, Kelkar YD, Ochman H. 2012. A selective force favoring increased

737      G+C content in bacterial genes. Proc Natl Acad Sci 109:14504–14507.

738   17.   Karberg KA, Olsen GJ, Davis JJ. 2011. Similarity of genes horizontally acquired by

739      *Escherichia coli* and *Salmonella enterica* is evidence of a supraspecies

740      pangenome. Proc Natl Acad Sci 108:20154–20159.

741   18.   Bohlin J, Eldholm V, Pettersson JHO, Brynildsrud O, Snipen L. 2017. The

742      nucleotide composition of microbial genomes indicates differential patterns of

743      selection on core and accessory genomes. BMC Genomics 18:1–11.

744   19.   López JL, Lozano MJ, Lagares A, Fabre ML, Draghi WO, Del Papa MF, Pistorio M,

745      Becker A, Wibberg D, Schlüter A, Pühler A, Blom J, Goesmann A, Lagares A.

746      2019. Codon usage heterogeneity in the multipartite prokaryote genome:

747      Selection-based coding bias associated with gene location, expression level, and

748      ancestry. MBio.

749   20.   Grantham R, Gautier C, Gouy M, Jacobzone M, Mercier R, Evolution G, Biometrie

750      L De, I UL. 1981. Nucleic Acids Research Codon catalog usage is a genome

751      strategy modulated for gene expressivity Nucleic Acids Research 9.

752   21.   Médigue C, Rouxel T, Vigier P, Hénaut A, Danchin A. 1991. Evidence for

753      horizontal gene transfer in *Escherichia coli* speciation. J Mol Biol 222:851–856.

754    22.    Grosjean H, Fiers W. 1982. Preferential codon usage in prokaryotic genes: the

755           optimal codon-anticodon interaction energy and the selective codon usage in

756           efficiently expressed genes. Gene.

757    23.    Kanaya S, Yamada Y, Kudo Y, Ikemura T. 1999. Studies of codon usage and

758           tRNA genes of 18 unicellular organisms and quantification of *Bacillus subtilis*

759           tRNAs: Gene expression level and species-specific diversity of codon usage based

760           on multivariate analysis. Gene 238:143–155.

761    24.    Karlin S, Mrazek J. 2000. Predicted highly expressed genes of diverse prokaryotic

762           genomes. J Bacteriol.

763    25.    dos Reis M, Wernisch L, Savva R. 2003. Unexpected correlations between gene

764           expression and codon usage bias from microarray data for the whole Escherichia

765           coli K-12 genome. Nucleic Acids Res.

766    26.    Supek F, Škunca N, Repar J, Vlahoviček K, Šmuc T. 2010. Translational selection

767           is ubiquitous in prokaryotes. PLoS Genet 6:1–13.

768    27.    Médigue C, Rouxel T, Vigier P, Hénaut A, Danchin A. 1991. Evidence for

769           horizontal gene transfer in Escherichia coli speciation. J Mol Biol.

770    28.    Mrázek J, Karlin S. 1999. Detecting alien genes in bacterial genomesAnnals of the

771           New York Academy of Sciences.

772    29.    Daubin V, Lerat E, Perrière G. 2003. The source of laterally transferred genes in

773           bacterial genomes. Genome Biol.

774    30.    Daubin V, Ochman H, Daubin V, Ochman H. 2004. Bacterial Genomes as New

775           Gene Homes : The Genealogy of ORFans in Bacterial Genomes as New Gene

776           Homes : The Genealogy of ORFans in E . coli 1036–1042.

777    31.    Ochman H, Lerat E, Daubin V. 2005. Examining bacterial species under the

778           specter of gene transfer and exchange. Proc Natl Acad Sci 102:6595–6599.

779    32.    Yannai A, Katz S, Hershberg R. 2018. The codon usage of lowly expressed genes

780           is subject to natural selection. Genome Biol Evol 10:1237–1246.

781    33.    Ikemura T. 1981. Correlation between the abundance of Escherichia coli transfer

782    RNAs and the occurrence of the respective codons in its protein genes: A proposal

783    for a synonymous codon choice that is optimal for the E. coli translational system.

784    J Mol Biol.

785  34.  Ikemura T. 1982. Correlation between the abundance of yeast transfer RNAs and

786    the occurrence of the respective codons in protein genes. J Mol Biol.

787  35.  Bulmer M. 1987. Coevolution of codon usage and transfer RNA abundance.

788    Nature.

789  36.  dos Reis M, Savva R, Wernisch L. 2004. Solving the riddle of codon usage

790    preferences: A test for translational selection. Nucleic Acids Res 32:5036–5044.

791  37.  Sabi R, Tuller T. 2014. Modelling the Efficiency of Codon-tRNA Interactions Based

792    on Codon Usage Bias. DNA Res 21:1–15.

793  38.  Gerber AP, Keller W. 2001. RNA editing by base deamination: More enzymes,

794    more targets, new mysteries. Trends Biochem Sci.

795  39.  Agris PF, Vendeix FAP, Graham WD. 2007. tRNA's Wobble Decoding of the

796    Genome: 40 Years of Modification. J Mol Biol.

797  40.  Novoa EM, Pavon-Eternod M, Pan T, Ribas de Pouplana L. 2012. A Role for tRNA

798    Modifications in Genome Structure and Codon Usage. Cell 149:202–213.

799  41.  Wald N, Alroy M, Botzman M, Margali H. 2012. Codon usage bias in prokaryotic

800    pyrimidine-ending codons is associated with the degeneracy of the encoded amino

801    acids. Nucleic Acids Res 40:7074–7083.

802  42.  Sharp PM, Bailes E, Grocock RJ, Peden JF, Sockett RE. 2005. Variation in the

803    strength of selected codon usage bias among bacteria. Nucleic Acids Res

804    33:1141–1153.

805  43.  Stoletzki N, Eyre-Walker A. 2007. Synonymous codon usage in *Escherichia coli*:

806    Selection for translational accuracy. Mol Biol Evol 24:374–381.

807  44.  Ran W, Kristensen DM, Koonin E V. 2014. Coupling between protein level

808    selection and codon usage. MBio 5:1–11.

809  45.  Jara E, Morel MA, Lamolle G, Castro-Sowinski S, Simón D, Iriarte A, Musto H.

810      2018. The complex pattern of codon usage evolution in the family

811      Comamonadaceae. Ecol Genet Genomics 6:1–8.

812   46.  McInerney JO. 1998. Replicational and transcriptional selection on codon usage in

813      Borrelia burgdorferi. Proc Natl Acad Sci USA 95:10698–10703.

814   47.  Davis JJ, Olsen GJ. 2010. Modal codon usage: Assessing the typical codon usage

815      of a genome. Mol Biol Evol 27:800–810.

816   48.  Sharp PM, Emery LR, Zeng K. 2010. Forces that influence the evolution of codon

817      bias. Philos Trans R Soc B Biol Sci 365:1203–1212.

818   49.  Shabalina SA, Spiridonov NA, Kashina A. 2013. Sounds of silence: Synonymous

819      nucleotides as a key to biological regulation and complexity. Nucleic Acids Res

820      41:2073–2094.

821   50.  Quax TEF, Claassens NJ, Söll D, van der Oost J. 2015. Codon Bias as a Means to

822      Fine-Tune Gene Expression. Mol Cell 59:149–161.

823   51.  Daubin V, Lerat E, Perrière G. 2003. The source of laterally transferred genes in

824      bacterial genomes. Genome Biol 4:R57.

825   52.  Lerat E, Daubin V, Ochman H, Moran NA. 2005. Evolutionary origins of genomic

826      repertoires in bacteria. PLoS Biol 3:0807–0814.

827   53.  Davis JJ, Olsen GJ. 2011. Characterizing the native codon usages of a genome:

828      An axis projection approach. Mol Biol Evol 28:211–221.

829   54.  Blom J, Albaum SP, Doppmeier D, Pühler A, Vorhölter FJ, Zakrzewski M,

830      Goesmann A. 2009. EDGAR: A software framework for the comparative analysis

831      of prokaryotic genomes. BMC Bioinformatics.

832   55.  Yu J, Blom J, Glaeser SP, Jaenicke S, Juhre T, Rupp O, Schwengers O, Spänig S,

833      Goesmann A. 2017. A review of bioinformatics platforms for comparative

834      genomics. Recent developments of the EDGAR 2.0 platform and its utility for

835      taxonomic and phylogenetic studies. J Biotechnol.

836   56.  Rambaut A, Drummond A. 2009. FigTree v1. 3.1. Website http//tree bio ed ac

837      uk/software/figtree.

838    57.    Karlin S, Barnett MJ, Campbell AM, Fisher RF, Mrazek J, Mrázek J. 2003.

839            Predicting gene expression levels from codon biases in alpha-proteobacterial

840            genomes. Proc Natl Acad Sci U S A.

841    58.    Wang M, Herrmann CJ, Simonovic M, Szklarczyk D, von Mering C. 2015. Version

842            4.0 of PaxDb: Protein abundance data, integrated across model organisms,

843            tissues, and cell-lines. Proteomics.

844    59.    Lê S, Josse J, Husson F. 2008. FactoMineR: An R Package for Multivariate

845            Analysis. J Stat Software, Artic 25:1–18.

846    60.    Sharp PM, Tuohy TMF, Mosurski KR. 1986. Codon usage in yeast: Cluster

847            analysis clearly differentiates highly and lowly expressed genes. Nucleic Acids

848            Res.

849    61.    Peden J. 1999. Analysis of codon usage [PhD dissertation]. University of

850            Nottingham.

851    62.    Wickham H. 2009. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag

852            New York.

853    63.    Lowe TM, Chan PP. 2016. tRNAscan-SE On-line: integrating search and context

854            for analysis of transfer RNA genes. Nucleic Acids Res 44:W54–W57.

855    64.    Ran W, Higgs PG. 2010. The influence of anticodon-codon interactions and

856            modified bases on codon usage bias in bacteria. Mol Biol Evol.

857    65.    Gingold H, Pilpel Y. 2011. Determinants of translation efficiency and accuracy. Mol

858            Syst Biol.

859    66.    Lawrence JG, Ochman H. 1997. Amelioration of bacterial genomes: Rates of

860            change and exchange. J Mol Evol 44:383–397.

861    67.    McInerney JO. 1997. Prokaryotic Genome Evolution as Assessed by Multivariate

862            Analysis of Codon Usage Patterns. Microb Comp Genomics 2:89–97.

863    68.    Ikemura T. 1981. Correlation between the abundance of *Escherichia coli* transfer

864            RNAs and the occurrence of the respective codons in its protein genes: A proposal

865        for a synonymous codon choice that is optimal for the *E. coli* translational system.

866        J Mol Biol 151:389–409.

867    69.  Felsenstein J. 2005. PHYLIP (Phylogeny Inference Package) version 3.6. Seattle.

868    70.  Rambaut A. 2009. FigTree, a graphical viewer of phylogenetic trees. Inst Evol Biol

869        Univ Edinburgh Ashworth Lab Edinburgh EH9 3JT.

870    71.  Wright F. 1990. The "effective number of codons" used in a gene. Gene.

871    72.  Revell LJ. 2012. phytools: An R package for phylogenetic comparative biology

872        (and other things). Methods Ecol Evol.

**FIGURE 1**

**FIGURE 2**

**FIGURE 3**

**FIGURE 4**

**FIGURE 5**

# FIGURE 6

# FIGURE 7

**Group A** — *Clostridium beijerinckii* NCIMB8052 — 29.9 %GC

**Group B** — *Porphyromonas gingivalis* ATCC 33277 — 48.4 %GC

**Group C** — *Corynebacterium resistens* DSM45100 — 57.1 %GC

**Group D** — *Actinomyces gerencseriae* DSM 6844 — 70.7 %GC

A  *Bacillus subtilis subsp. spizenii TU B10*

B  *Bacteroides vulgatus ATCC 8482*

C  *Sulfurospirillum multivorans DSM 12446*

D  *Streptococcus equi ATCC 33398*

E  *Yersinia enterocolitica subsp. palearctica Y11*

F  *Mycobacterium fortuitum subsp. fortuitum DSM 46621*

|  | CA | GC3 | Species |
|---|---|---|---|
| **Group A** | | | *C. beijerinckii* NCIMB_8052<br>*M. smithii* ATCC35061<br>*T. halophilus* NBRC12172 |

**Legend:** U-bias, C-bias; PHE; C1→Cn; Singletons

**Group B species:**
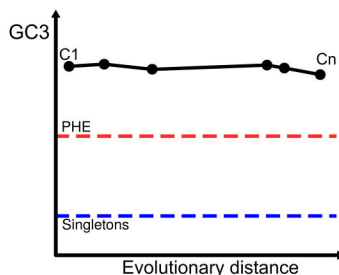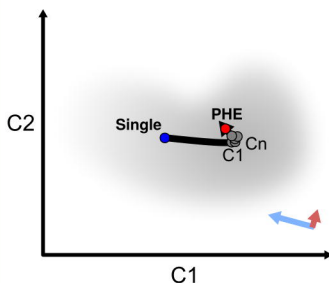*S. multivorans* DSM12446
*S. equi* ATCC33398
*B. vulgatus* ATCC8482
*B. subtilis subsp. spizenii* TUB10
*M. bovis* CCUG2133
*C. violaceum* ATCC12472
*P. gingivalis* ATCC33277
*P. graminis* DSM15220
*P. gaetbulicola* Gung47

**Group C species:**
*T. succinifaciens* DSM2489
*P. melaninogenica* ATCC25845
*A. parvulum* DSM20469
*Y. enterocolitica subsp. palearctica* Y11
*M. petrolearia* DSM11571
*C. resistens* DSM45100
*B. longum subsp. longum* JCM1217
*B. holmesii* ATCC51541
*M. fortuitum subsp. fortuitum* DSM46621
*S. parapaucimobilis* NBRC15100

**Group D species:**
*L. methylohalidivorans* DSM14336
*O. pacificum* MCCC1A02656
*A. enclensis* NIO1008
*R. denitrificans* 2APBS1
*M. aurum* KACC15219
*R. mucosa* ATCCBAA692
*A. gerencseriae* DSM6844