



Xenologs show low expression levels in the cyanobacterium *Synechococcus elongatus*

Journal:	<i>Genome Biology and Evolution</i>
Manuscript ID:	Draft
Manuscript Type:	Research Article
Date Submitted by the Author:	n/a
Complete List of Authors:	<p>Delaye, Luis; Cinvestav Irapuato, Departamento de Ingenieria Genetica González-Domenech, Carmen; Universidad de Granada, Facultad de Farmacia De la Cruz, Fernando; Universidad de Cantabria-CSIC-SODERCAN, Departamento de Biología Molecular e Instituto de Biomedicina y Biotecnología de Cantabria (IBBTEC) Moya, Andres; Universitat de València, Instituto Cavanilles de Biodiversidad y Biología Evolutiva</p>
Keywords:	Correspondence analysis, codon usage, horizontal gene transfer, translation efficiency

SCHOLARONE™
Manuscripts

view

1
2
3 **Xenologs show low expression levels in the cyanobacterium**
4
5
6 ***Synechococcus elongatus***
7
8
9

10
11 Luis Delaye^{1*}, Carmen M González-Domenech², Fernando de la Cruz³ and Andrés
12 Moya⁴
13
14
15

16
17 ^{1*}Author for Correspondence: Luis Delaye, Departamento de Ingeniería Genética
18 CINVESTAV-Irapuato, Km. 9.6 Libramiento Norte, Carretera Irapuato-León, 36821
19 Irapuato, Guanajuato, México. Tel: +52 (462) 6239669. Fax +52 (462) 6245846. LD:
20 ldelaye@ira.cinvestav.mx.
21
22
23
24
25

26
27 ²Facultad de Farmacia, Universidad de Granada, Campus de Cartuja s/n. 18071
28 Granada, Spain
29
30
31

32 ³Departamento de Biología Molecular e Instituto de Biomedicina y Biotecnología de
33 Cantabria (IBBTEC), Universidad de Cantabria-CSIC-SODERCAN, Santander, Spain.
34
35
36

37 ⁴Instituto Cavanilles de Biodiversidad y Biología Evolutiva, Universitat de València,
38 Spain.
39
40
41
42
43
44

45 Email addresses:
46
47

48 LD: ldelaye@ira.cinvestav.mx
49

50
51 CG: cmgodo@ugr.es
52

53
54 FC: fernando.cruz@unican.es
55

56
57 AM: andres.moya@uv.es
58
59
60

Abstract

Horizontal gene transfer is a central process in prokaryotic evolution. However, little is known about the factors which determine successful gene transfers. In particular, few studies have tackled the importance of gene expression levels with respect to their role in the success of a horizontally acquired gene. To answer this question, we studied the expression level of genes that clearly arrived by horizontal gene transfer into the genome of the cyanobacterium *Synechococcus elongatus* PCC 7942. We have found that xenologs show lower expression levels than the rest of the genes in the genome. We have further found that there is selection for translation efficiency in *S. elongatus*, and that the overall expression levels of xenologs, as suggested by CAI values, are low. We interpret this pattern as an indication that most successful horizontal gene transfers are those that least affect the fitness of the recipient cell. Finally, we discuss the case of the *ccoN* gene coding for a cytochrome oxidase as an example of a successful xenolog co-opted for (Fe) starvation in the context of selection of synonymous codon usage.

Keywords

Correspondence analysis, codon usage, horizontal gene transfer, translation efficiency

Background

Horizontal gene transfer (HGT) is recognized as one of the major processes in prokaryotic evolution (Zhaxybayeva et al. 2011). Starting on the micro-evolutionary scale, recent comparisons of 61 sequenced *Escherichia coli* genomes show that the core genome, i.e., the set of genes shared by all compared genomes, is comprised of only 993 gene families, while the pangenome, i.e., the set of genes shared by 61 or more genomes, is composed of 15,741 gene families (Lukjancenko, et al. 2010). This suggests that if the ancestral *E. coli* genome had approximately 4,000 genes, extant strains would have acquired about two-thirds of their genomes by HGT.

The role of HGT in the evolution of bacterial adaptation is clearly exemplified by higher levels of resistance to antibiotics among pathogenic bacteria (Maclean et al. 2010). Other examples of adaptation by HGT include the acquisition of several genes involved in salt tolerance in the halophilic bacterium *Salinibacter ruber* (Mongodin et al. 2005) and the adaptation to low phosphorous levels by the acquisition of genes for the sulfoquinovose synthesis in *Bacillus coahuilensis* (Alcaraz et al. 2008). Furthermore, the evolution by HGT of new metabolic pathways which degrade xenobiotics exemplifies the rapidity and efficacy of this adaptive process (Springael et al. 2004). Xenobiotics are defined as anthropogenic compounds which are believed to have been alien to extant microbial communities. For example, the *pcpB* gene, which initiates the degradation of pentachlorophenol (considered a 'true' xenobiotic), was acquired by the pentachlorophenol-degrading bacterium *Sphingomonas chlorophenolicum* (Tirola et al. 2002) via HGT. More dramatic examples include the evolution of a pathway to degrade chlorobenzene in the *Ralstonia* spp. strain JS705. In this case, the two critical genes for chlorobenzene degradation originated from different bacterial species, one from

1
2
3 *Ralstonia* spp. strain JS745 and the other from *Pseudomonas* sp. strain B13 (Müller et
4
5 al. 2003).

6
7
8 Finally, on the macro-evolutionary scale, the structure of the universal tree of life, as
9
10 revealed by the 16/18S rRNA molecule, is being challenged by the impact of HGT in
11
12 prokaryotic evolution (Doolittle, 1999; Doolittle, 2007). How to represent prokaryotic
13
14 evolution in light of HGT is, as of today, a work in progress (Lazcano et al. 2011).

15
16
17 Despite the importance of HGT in prokaryotic evolution, the factors determining the
18
19 success of acquired genes are not totally understood. Once a protein coding gene is
20
21 acquired by a genome via HGT, its protein product interacts with the rest of the
22
23 molecules present in the intracellular environment. The fate of the new xenolog depends
24
25 upon what effects these molecular interactions have on the fitness of the recipient cell.
26
27 For example, an *in silico* analysis of the metabolic network of *E. coli* has shown that the
28
29 chance of acquiring a gene by HGT is up to six times higher if an enzyme that catalyses
30
31 a coupled metabolite flux is already encoded in the genome (Pál et al. 2005). Another
32
33 factor that seems to determine the success of a transferred gene is codon usage. In
34
35 species where there has been selection on codon usage for translation efficiency, there
36
37 will be a bias in the codon composition of highly expressed genes when compared to the
38
39 codon usage of the rest of genes in the genome (Henry et al. 2007 and references
40
41 therein). For instance, the codon adaptation index (CAI) uses this bias to identify highly
42
43 expressed genes (Sharp et al. 1987). When a gene is acquired by a genome via HGT, it
44
45 is very likely that its codon usage differs from the average codon usage (or the optimal
46
47 codon usage) of the genes in the recipient genome. It has been proposed that with the
48
49 passing of time, xenologs acquire the codon composition of the recipient genome in a
50
51 process known as amelioration (Lawrence et al. 1997). This is because xenologs suffer
52
53 the same mutation pressure as the rest of the genes in the genome. However, a different
54
55
56
57
58
59
60

1
2
3 interpretation has been given to the observation that xenologs (at least some of them)
4 resemble the codon usage of recipient genes. Accordingly, codon usage compatibility
5 between foreign genes and recipient genomes increases the fixation probability of HGT
6 events (Medrano-Soto et al. 2004).
7
8
9

10
11
12 Recently, it has been shown that expression level hampers HGT in *Escherichia coli*
13 (Park et al. 2012). It was proposed that xenologs deteriorate the fitness of the recipient
14 cells due to: a) energy expenditure in transcription and translation; b) cytotoxic protein
15 misfolding; c) reduction in cellular transcriptional efficiency; d) detrimental protein
16 interaction; and e) disturbance of the optimal protein concentration or cell physiology
17 (Park et al. 2012). Here we show that the xenologs that have been transferred into
18 *Synechococcus elongatus* PCC 7942 from non-cyanobacterial species also show low
19 expression levels, thus corroborating and expanding previous results for *E. coli* (Park et
20 al. 2012).
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

37 Results

38
39 We found that xenologs in *Synechococcus elongatus* PCC 7942 (hereafter *S. elongatus*)
40 have, on average, lower expression levels than the rest of the genes in the genome
41 (Figure 1). The difference is statistically significant at the $P < 0.001$ level (one-sided
42 Wilcoxon test). This difference is statistically significant even when genes belonging to
43 a large phage island are removed from the dataset. A summary of xenologs here
44 identified is shown in Figure 2, and a more complete description of them is found in
45 Supplementary Table 1.
46
47
48
49
50
51
52
53
54

55
56 It has been suggested that the Codon Adaptation Index (CAI) is a better indication of
57 the overall expression level of a gene than the direct measurement of its expression on a
58
59
60

1
2
3 single laboratory condition (Fraser et al. 2004). To determine if xenologs also show low
4
5 CAI values, we first studied to what extent bias in codon usage among *S. elongatus*
6
7 genes is determined by G+C mutational bias. This is because CAI values are only
8
9 meaningful in species where there has been selection of codon usage for translation
10
11 efficiency (Henry et al. 2007; Puigbò et al. 2008).
12

13
14 One way to evaluate the dependence of bias in codon usage to G+C content is by using
15
16 an effective number of codons (ENc) plot (Wright, 1990). Genes whose codon choice is
17
18 constrained only by a G+C mutational bias will lie on or just below their predicted ENc
19
20 values (see methods). The ENc plot for *S. elongatus* is shown in Figure 3. The
21
22 correlation between the observed and “expected” ENc values is moderate (correlation
23
24 coefficient $r \sim 0.51$). This result shows that variation in codon usage in *S. elongatus* is
25
26 not entirely due to G+C content.
27
28

29
30 We have further studied the degree of determination of G+C content on the bias in
31
32 codon usage among *S. elongatus* genes by performing a Correspondence Analysis
33
34 (COA). In essence, COA creates a series of orthogonal axes in order to identify trends
35
36 in data variation, with each subsequent axis explaining a decreasing amount of the
37
38 variation (Peden, 1999 and references therein). When applied to the study of the
39
40 variation in Relative Synonymous Codon Usage (RSCU) among *S. elongatus* genes,
41
42 COA shows that the first axis strongly correlates with G+C content at third codon
43
44 positions (GC_{3s}) (correlation coefficient $r \sim -0.86$, $p < 0.001$ Spearman rank correlation)
45
46 and explains ~12% of all variation (Figure S1).
47
48
49

50
51 Although this result shows that G+C mutational bias impacts codon usage, several other
52
53 aspects of the first axis are better explained by selection of codon usage for translation
54
55 efficiency. These are: a) that there is an enrichment of ribosomal protein coding genes at
56
57
58
59
60

1
2
3 one extreme of the values of the first axis, and of hypothetical genes at the other
4
5 extreme; b) codons predicted to be optimal by COA in *S. elongatus* are also optimal in
6
7 other species known to be under selection for translation efficiency; c) the principal axis
8
9 also correlates with observed ENc values (correlation coefficient $r \sim 0.43$, $p < 0.001$
10
11 Spearman rank correlation), which is an independent measure of codon bias; and d)
12
13 although slightly, axis 1 also correlates with gene expression (correlation coefficient $r \sim$
14
15 -0.06 , $p < 0.001$ Spearman rank correlation). (See supplementary information for a
16
17 complete description of COA results, Figures S1 to S4 and Table S2.)
18
19

20
21 Previous results clearly indicate that, although bias in G+C mutation does affect bias in
22
23 codon usage among *S. elongatus* genes, selection for translation efficiency has been
24
25 strong enough to leave a pattern in RSCU variation among genes, rendering CAI values
26
27 meaningful, i.e., there has been selection for a sub-set of codons in highly expressed
28
29 genes. This is supported by the correlation between CAI values and expression level
30
31 (Figure 4) (correlation coefficient $r \sim 0.1$, $p < 0.001$ Spearman rank correlation).
32
33

34
35 With the exception of ORPhans, CAI values of xenologs are on average lower than the
36
37 CAI values of the rest of the genes in the genome (Figure 5). The difference is
38
39 statistically significant at the $P < 0.001$ level (one-sided Wilcoxon test).
40
41
42
43
44

45 46 **Discussion**

47
48 Here we show that xenologs in *S. elongatus* have, on average, lower expression levels
49
50 than the rest of the genes in the genome. This is true for direct measurements of gene
51
52 expression levels by RNA sequencing (Vijayan et al. 2011) as well as for CAI values,
53
54 which are an indirect measure of gene expression. Differing from the study by (Park et
55
56 al. 2012), we have further validated the use of CAI by showing that selection for
57
58
59
60

1
2
3 translation efficiency has been strong enough to leave a mark in codon bias in *S.*
4
5 *elongatus*. Our results confirm and expand previous observations showing that
6
7 expression level is the principal factor determining the success of an HGT (Park et al.
8
9 2012).

10
11
12 It has been argued that most successful xenologs are neutral or nearly neutral (Park et al.
13
14 2012; Gogarten and Townsend 2005). This is surprising given the efficacy of natural
15
16 selection on prokaryotes due to their large population sizes. However, the fact that
17
18 xenologs tend to show low expression levels supports previous claims about their
19
20 relative neutrality, i.e., genes which least affect fitness are the ones that are the most
21
22 persistent once they are acquired by HGT. If this is the case, then it is possible to picture
23
24 an image where the xenome (the set of genes frequently and successfully transferred
25
26 between species) is running in the background in terms of fitness, waiting for natural
27
28 selection to act and pick a lucky genetic variant when there is a favorable change in
29
30 environmental conditions. In those cases where the presence of a xenolog is adaptive,
31
32 natural selection can further modulate its expression level to better fit its product to the
33
34 inner workings of the cell.
35
36
37
38
39

40
41 It is likely that the above scenario explains the presence of some of the xenologs in *S.*
42
43 *elongatus*. For instance, *ccoN* (Synpcc7942_0202) and *catD* (Synpcc7942_2603) are
44
45 two genes for cytochrome oxidases coded in *S. elongatus*. The first one, *ccoN*, was
46
47 acquired by HGT, while the other, *catD*, shows a phylogeny consistent with a history of
48
49 vertical inheritance (Supplementary material). Interestingly, it has been shown that the
50
51 expression level of the gene *ccoN* significantly increases under (Fe) starvation (Nodop
52
53 et al. 2008). This is not the case for *catD*. And, unexpectedly for a xenolog, the
54
55 expression level of *ccoN* as measured by RNA sequencing (Vijayan et al. 2011) and its
56
57 CAI value are larger than those of the indigenous cytochrome oxidase *ctaD*. One could
58
59
60

1
2
3 argue that the expression level of *ccoN* has been adapted by modifying its codon usage
4 as an adaptation to (Fe) starvation. Consistent with this hypothesis is the fact that in
5 phylogeny, the branch leading to *ccoN*, the number of d_S substitutions has sensibly
6 increased with respect to the number of d_N changes (Supplementary material).
7
8
9

10
11
12 Other xenologs that may have been co-opted for cellular functions are
13 Synpcc7942_0458 and Synpcc7942_0459, which are predicted to code for an S-
14 formylglutathione hydrolase and a S-(hydroxymethyl)glutathione dehydrogenase. These
15 genes are predicted to be in a single transcription unit according to bioinformatic
16 predictions (Caspi et al. 2012) as well as experimental determinations (Vijayan et al.
17 2012). And, according to the KEGG database (Kanehisa et al. 2012), these enzymes
18 catalyze two successive steps to synthesize formate from S-hydroxy-methyl-glutathione.
19 However, they could also participate in the synthesis of cofactors. In another example,
20 there is experimental evidence that the proteins coded by the xenologs *moaA* and *moaC*
21 (Synpcc7942_1282 and Synpcc7942_1285 respectively) together with other genes in
22 the *narA* locus, participate in molybdenum cofactor biosynthesis for nitrate reduction
23 (Rubio et al. 1998). More studies are needed to find out if these enzymes were initially
24 neutral and subsequently retained by natural selection.
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

42 Xenologs here identified are certainly a subset of all genes acquired by HGT in *S.*
43 *elongatus*. This is because we only identified genes acquired from non-cyanobacterial
44 origins. It remains to be seen whether this pattern also applies to xenologs acquired
45 from other cyanobacteria.
46
47
48
49
50

51 52 53 54 55 **Material and methods** 56 57 58 59 60

1
2
3 *Xenologs identification.* The protein sequences from *S. elongatus* PCC 7942 were
4
5 downloaded from Genbank database (Benson et al. 2011). The genome sequence data
6
7 for this organism was produced by the US Department of Energy Joint Genome Institute
8
9 (www.jgi.doe.gov). BLAST (Camacho et al. 2009) searches were conducted locally for
10
11 all *S. elongatus* PCC 7942 protein sequences in the non-redundant (nr) database from
12
13 NCBI (ftp.ncbi.nlm.nih.gov/blast/db/FASTA/). Only homologs showing an e-value
14
15 smaller than 0.001 and a minimum of 70% coverage of the query relative to the
16
17 sequence of the nr database were retrieved. For each *S. elongatus* protein, up to 50
18
19 homologs were considered. The above criteria were used to retrieve only the closest
20
21 homologs to *S. elongatus* PCC 7942 protein sequences.
22
23

24
25
26 We then reconstructed a phylogeny for each set of homologs as follows. Identified
27
28 homologs were aligned with MUSCLE v3.8 (Edgar, 2004), and for each multiple
29
30 alignment, the best-fit model of evolution was identified with PROTTEST v2.4
31
32 according to the Akaike Information Criterion (Abascal et al. 2005). Phylogenies for
33
34 each set of homologs were reconstructed with PHYML v3.0 following the best model of
35
36 evolution detected with PROTTEST (search parameter BEST for NNI and SPR, and six
37
38 rate categories) (Guindon et al. 2010). For statistical support of branches, we used
39
40 aLRT. And finally, by using Perl scripts and the Python package E.T.E. v2.1 (Huerta-
41
42 Cepas et al. 2010), we classified as *xenologs* all *S. elongatus* PCC 7942 proteins having
43
44 homologs only from species other than cyanobacteria.
45
46
47

48
49 *Effective Number of Codons (ENc) plot.* We studied the relationship between codon and
50
51 mutational bias with a formula that predicts the expected codon usage solely as a
52
53 function of G+C in the third codon position (Wright, 1990).
54
55

$$ENc = 2 + S + (29/(S^2 + (1 - S)^2))$$

1
2
3 Where S is the frequency of G+C in the third codon position. Genes whose codon
4 choice is constrained only by a G+C mutational bias will lie on or just below the curve
5 predicted by the above equation.
6
7
8

9
10 *Correspondence Analysis (COA)*. The software CodonW was used for COA (Peden,
11 1999). Genes with fewer than 50 codons were initially excluded from the analysis to
12 reduce signal noise. Then, by using the Relative Synonymous Codon Usage (RSCU),
13 the vectors of the COA were generated from those genes with more than 50 codons.
14 After vectors were generated, those genes with less than 50 codons were added to the
15 COA. Finally, to identify optimal codons, the RSCU of the 5% of genes in the most
16 extreme values of Axis 1 were contrasted.
17
18
19
20
21
22
23
24

25
26 *Codon Adaptation Index (CAI)*. The CAI values were calculated by using CodonW
27 software (Peden, 1999).
28
29
30
31

32 *Rates of synonymous (d_S) and non synonymous (d_N) substitutions*. To infer the number
33 of synonymous and non-synonymous substitutions in the branch leading to the *ccoN*
34 gene, we used CODEML from the PAML v4.4 package (Yang, 2007). We contrasted
35 two models by using a Likelihood Ratio Test (LRT). The first model had one parameter
36 ($\omega = d_N/d_S$) for all branches; and the second model had two parameters, one for the
37 branch leading to *ccoN* ($\omega_1 = d_N/d_S$) and the other for the rest of the branches ($\omega_0 =$
38 d_N/d_S).
39
40
41
42
43
44
45
46
47

48 *Statistical analysis*. All statistical analyses were conducted in R (www.r-project.org/).
49
50
51
52
53
54

55 **List of abbreviations**

56
57 COA, RSCU, LRT, CAI, ENc, HGT.
58
59
60

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

LD identified homologs, carried the statistical analysis and drafted the manuscript. CG classified genes and drafted the manuscript. FC conceived the study and made important contributions to draft the manuscript. AM conceived the study and made important contributions to draft the manuscript. All authors read and approved the final manuscript.

Description of additional data files

Supplementary Fig. 1. G+C content at third codon position versus Axis 1 of COA.

Supplementary Fig. 2. Occurrence of genes coding for ribosomal proteins (red), for proteins participating in the photosystem (green) or for hypothetical proteins (gray) along Axis 1 of COA. Each bar denotes the number of occurrences along 50 genes.

Supplementary Fig. 3. Observed ENc values versus Axis 1 of COA.

Supplementary Fig. 4. Expression level versus Axis 1. Expression level is measured as Mean RNA Sequencing over ORF (MRSO).

Supplementary Fig. 5. The number of silent substitutions (d_s) has sensibly increased in the branch leading to *ccoN*. a) Maximum-Likelihood tree of *ccoN*. The branch leading to *ccoN* is indicated in green and with a hash tag (#). *ccoN* is indicated by the gi number

1
2
3 of the coded protein (81299013). The number of d_N and d_S substitutions is shown below
4
5 the branch. Only bootstrap values ≥ 50 are shown; b) the distribution of d_N and d_S
6
7 substitutions. The number of d_S substitutions in the branch leading to *ccoN* is indicated
8
9 with an arrow. The tree with two parameters (omegas), one for the branch leading to
10
11 *ccoN* and other for the rest of the branches, explains better the variation in d_N and d_S
12
13 than the tree with only one parameter ($P = 0.0045$).
14
15

16 17 **Supplementary table 1.**

18
19 Xenologs in *S. elongatus* acquired from non-cyanobacterial sources. Contiguous genes
20
21 in the genome are indicated in blue. MRSO: expression level in “Mean Sequences over
22
23 ORF” according to (Vijayan et al. 2011). In violet we show genes belonging to a large
24
25 phage related island.
26
27

28
29 **Supplementary table 2.** Position of each codon by Axis 1 in COA. Those codons that
30
31 occur significantly more often ($p < 0.01$) are indicated in yellow, and those codons with
32
33 $p < 0.05$ in green. Typically optimal and rare codons are from *E. coli*, *B. subtilis*, *S.*
34
35 *cerevisiae*, *S. pombe*, and *D. melanogaster* (Peden, 1999 and references therein).
36
37

38
39 **Supplementary file of xenologs.** The phylogenies of all xenologs in newick format.
40
41

42 43 44 45 **Authors' information**

46
47
48 LD: Full professor of Genetics; CG: postdoctoral specialist in bioinformatics; FC: Full
49
50 professor of Genetics; AM: Full professor of Genetics.
51
52

53 54 55 56 **Acknowledgements**

1
2
3 This work was supported by Consejo Nacional de Ciencia y Tecnología CONACYT
4
5 CB-2010-01 [grant number 157220]. This work was further funded by the grants
6
7 SAF2009-13032-C02-01, and SAF2012-31187 from the Spanish Ministry of Economy
8
9 and Competitiveness, Prometeo/2009/092 from Generalitat Valenciana (Spain), and ST-
10
11 FLOW (EU).
12
13

14 15 16 17 18 **References**

19
20
21 Abascal F, Zardoya R, Posada D. 2005. ProtTest: selection of best-fit models of protein
22
23 evolution. *Bioinformatics*. 21:2104-2105.
24

25
26 Alcaraz LD, et al. 2008. The genome of *Bacillus coahuilensis* reveals adaptations
27
28 essential for survival in the relic of an ancient marine environment. *Proc Natl Acad Sci*
29
30 U S A. 105:5803-5808.
31

32
33 Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. 2011.
34
35 *Nucleic Acids Res*. 39(Database issue):D32-7.
36

37
38 Camacho C, et al. 2009. BLAST+: architecture and applications. *BMC Bioinformatics*.
39
40 10:421.
41

42
43 Caspi R, et al. 2012. The MetaCyc database of metabolic pathways and enzymes and
44
45 the BioCyc collection of pathway/genome databases. *Nucleic Acids Res*. 40(Database
46
47 issue):D742-53.
48

49
50
51 Doolittle WF. 1999. Phylogenetic classification and the universal tree. *Science*.
52
53 284:2124-2129.
54
55
56
57
58
59
60

1
2
3 Doolittle WF, Baptiste E. 2007. Pattern pluralism and the Tree of Life hypothesis. Proc
4
5 Natl Acad Sci U S A. 104:2043-2049.
6

7
8 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high
9
10 throughput. Nucleic Acids Res. 32:1792-1797.
11

12
13 Fraser HB, Hirsh AE, Wall DP, Eisen MB. 2004. Coevolution of gene expression
14
15 among interacting proteins. Proc Natl Acad Sci U S A. 101:9033-9038.
16

17
18 Gogarten JP, Townsend JP. 2005. Horizontal gene transfer, genome innovation and
19
20 evolution. Nat Rev Microbiol. 3:679-687.
21

22
23 Guindon S, et al. 2010. New algorithms and methods to estimate maximum-likelihood
24
25 phylogenies: assessing the performance of PhyML 3.0. Syst Biol. 59:307-321.
26

27
28 Henry I, Sharp PM. 2007. Predicting gene expression level from codon usage bias. Mol
29
30 Biol Evol. 24:10-12.
31

32
33 Huerta-Cepas J, Dopazo J, Gabaldón T. 2010. ETE: a python environment for tree
34
35 exploration. BMC Bioinformatics. 11:24.
36

37
38 Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M. 2012. KEGG for integration and
39
40 interpretation of large-scale molecular data sets. Nucleic Acids Res. 40(Database
41
42 issue):D109-14.
43

44
45 Lawrence JG, Ochman H. 1997. Amelioration of bacterial genomes: rates of change and
46
47 exchange. J Mol Evol. 44:383-397.
48

49
50 Lazcano A. 2011. Natural history, microbes and sequences: shouldn't we look back
51
52 again to organisms? PLoS One. 6(8):e21334.
53
54
55
56
57
58
59
60

1
2
3 Lukjancenko O, Wassenaar TM, Ussery DW. 2010. Comparison of 61 sequenced
4
5 *Escherichia coli* genomes. Microb Ecol. 60:708-720.
6
7

8 Maclean RC, Hall AR, Perron GG, Buckling A. 2010. The evolution of antibiotic
9
10 resistance: insight into the roles of molecular mechanisms of resistance and treatment
11
12 context. Discov Med. 10:112-118.
13
14

15 Medrano-Soto A, Moreno-Hagelsieb G, Vinuesa P, Christen JA, Collado-Vides J. 2004.
16
17 Successful lateral transfer requires codon usage compatibility between foreign genes
18
19 and recipient genomes. Mol Biol Evol. 21:1884-1894.
20
21
22

23 Mongodin EF, et al. 2005. The genome of *Salinibacter ruber*: convergence and gene
24
25 exchange among hyperhalophilic bacteria and archaea. Proc Natl Acad Sci U S A. 102:
26
27 18147-18152.
28
29

30 Müller TA, Werlen C, Spain J, Van Der Meer JR. 2003. Evolution of a chlorobenzene
31
32 degradative pathway among bacteria in a contaminated groundwater mediated by a
33
34 genomic island in *Ralstonia*. Environ Microbiol. 5:163-173.
35
36
37

38 Nodop A, et al. 2008. Transcript profiling reveals new insights into the acclimation of
39
40 the mesophilic fresh-water cyanobacterium *Synechococcus elongatus* PCC 7942 to iron
41
42 starvation. Plant Physiol. 147:747-763.
43
44

45 Pál C, Papp B, Lercher MJ. 2005. Horizontal gene transfer depends on gene content of
46
47 the host. Bioinformatics. 21 Suppl 2:ii222-3.
48
49

50 Park Ch, Zhang J. 2012. High expression hampers horizontal gene transfer. Genome
51
52 Biol. Evol. 4: 523-532.
53
54

55 Peden J. 1999. (<http://codonw.sourceforge.net>).
56
57
58
59
60

1
2
3 Puigbò P, Romeu A, García-Vallvé S. 2008. HEG-DB: a database of predicted highly
4 expressed genes in prokaryotic complete genomes under translational selection. *Nucleic*
5 *Acids Res.* vol. 36, Database issue: D524-D527.
6
7

8
9
10 Rubio LM, Flores E, Herrero A. 1998. The *narA* locus of *Synechococcus* sp. strain PCC
11 7942 consists of a cluster of molybdopterin biosynthesis genes. *J Bacteriol.* 180:1200-
12 1206.
13
14

15
16
17 Sharp PM, Li WH. 1987. The codon Adaptation Index--a measure of directional
18 synonymous codon usage bias, and its potential applications. *Nucleic Acids Res.*
19 15:1281-1295.
20
21

22
23
24 Springael D, Top EM. 2004. Horizontal gene transfer and microbial adaptation to
25 xenobiotics: new types of mobile genetic elements and lessons from ecological studies.
26 *Trends Microbiol.* 12:53-58.
27
28

29
30
31 Tiirola MA, Wang H, Paulin L, Kulomaa MS. 2002. Evidence for natural horizontal
32 transfer of the *pcpB* gene in the evolution of polychlorophenol-degrading
33 sphingomonads. *Appl Environ Microbiol.* 68:4495-4501.
34
35

36
37
38 Vijayan V, Jain IH, O'Shea EK. 2012. A high resolution map of a cyanobacterial
39 transcriptome. *Genome Biol.* 12:R47.
40
41

42
43
44 Wright F. 1990. The effective number of codons in a gene. *Gene.* 87:23-29.
45
46

47
48
49 Yang Z. 2007. PAML 4: a program package for phylogenetic analysis by maximum
50 likelihood. *Mol Biol Evol.* 24: 1586-1591
51
52

53
54
55 Zhaxybayeva O, Doolittle WF. 2011. Lateral gene transfer. *Curr Biol.* 21:R242-6.
56
57
58
59
60

Figure legends

Fig. 1. Xenologs have low expression levels. The expression level of xenologs is on average lower than the expression level of the rest of the genes of the genome (p value < 0.001 , one sided Wilcoxon-Mann-Whitney test). Each box denotes the median, upper and lower quartiles, the inter quartile range (denoted by the whiskers), and the outliers (denoted by dots), of the gene expression level distribution for each gene category. The number of genes in each category is as follow: (o) ORPhans, 136; (r) rest of the genes, 2273; (ps) genes from the photosystem, 50; (rp) ribosomal protein coding genes, 52; (x) xenologs, 101.

Fig. 2. Xenologs classified according to function. All genes related to phage shown here are found in a single locus in *S. elongatus*.

Fig. 3. ENc plot showing the relation between predicted and observed bias in codon usage according to G+C content. Red dots indicate the expected codon usage according to G+C content at third codon positions. Gray dots show observed codon usage for each gene.

Fig. 4. Expression level versus CAI. The number of genes in each category is as follow: (o) ORPhans in blue, 136; (r) rest of the genes in gray, 2273; (ps) genes from the photosystem in green, 50; (rp) ribosomal protein coding genes in red, 52; and (x) xenologs in orange, 101.

Fig. 5. Xenologs have low CAI values. The CAI values of xenologs is on average lower than the expression level of the rest of the genes of the genome (p value < 0.001 , one sided Wilcoxon-Mann-Whitney test). Each box denotes the median, upper and lower quartiles, the inter quartile range (denoted by the whiskers), and the outliers (denoted by dots), of the gene expression level distribution for each gene category. The

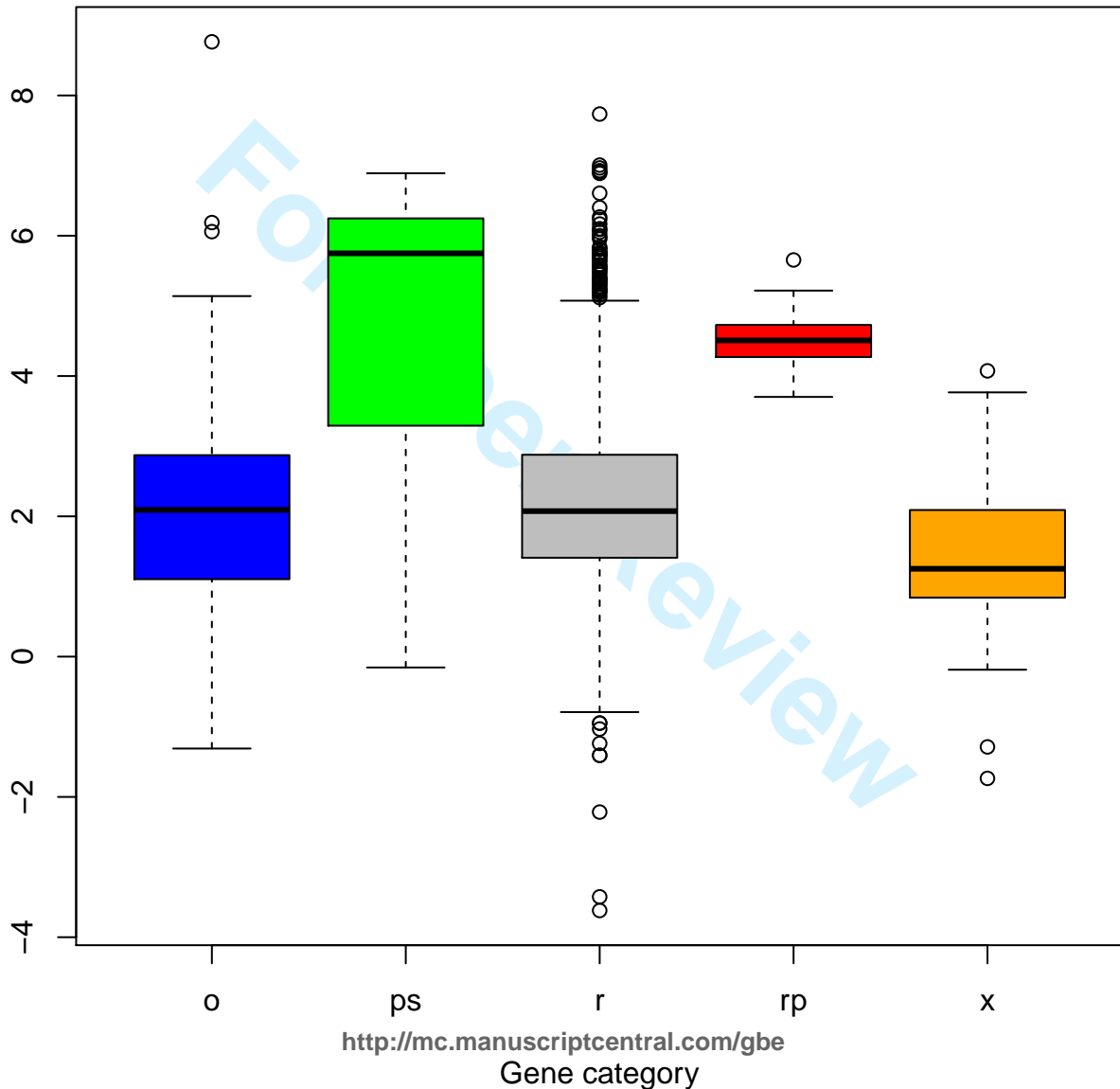
1
2
3 number of genes in each category is as follow: (o) ORPhans, 136; (r) rest of the genes,
4
5 2273; (ps) genes from the photosystem, 50; (rp) ribosomal protein coding genes, 52; (x)
6
7 xenologs, 101.
8
9

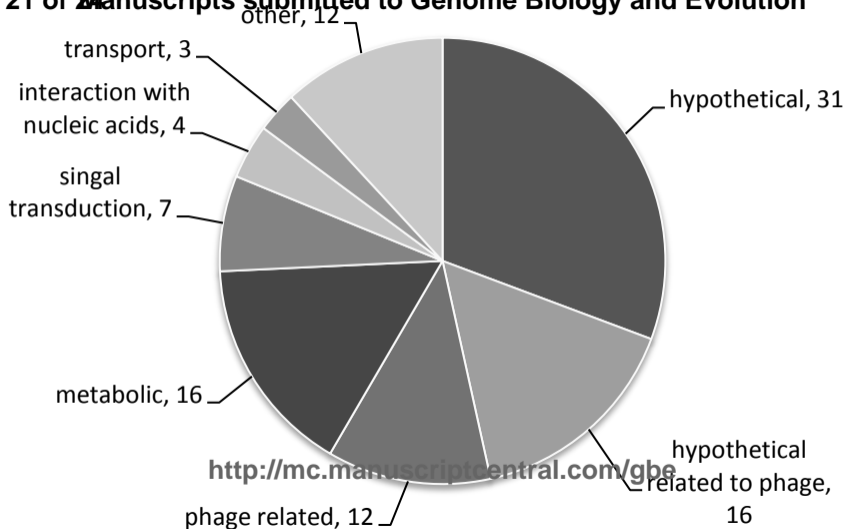
10 **Tables**

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

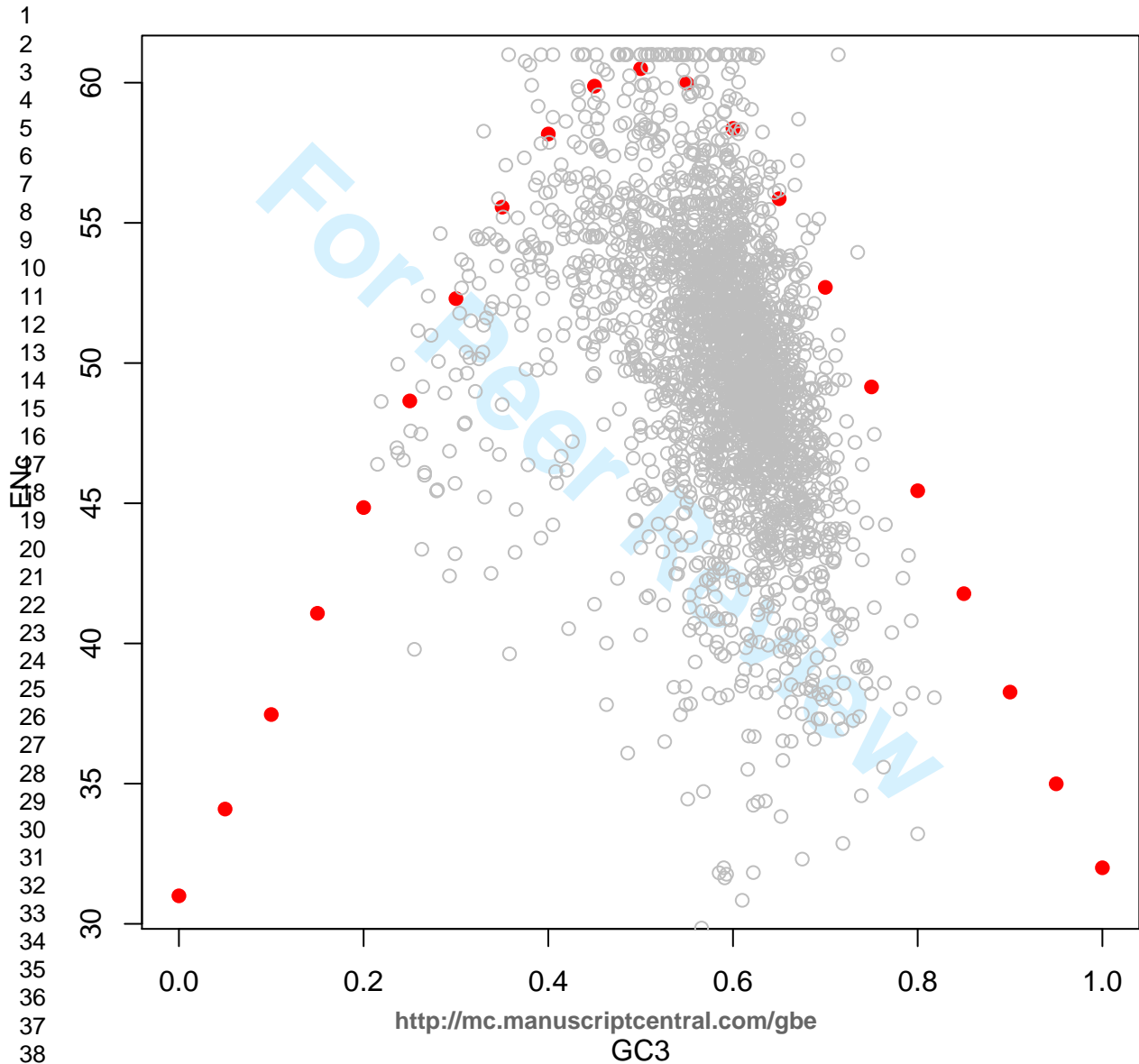
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

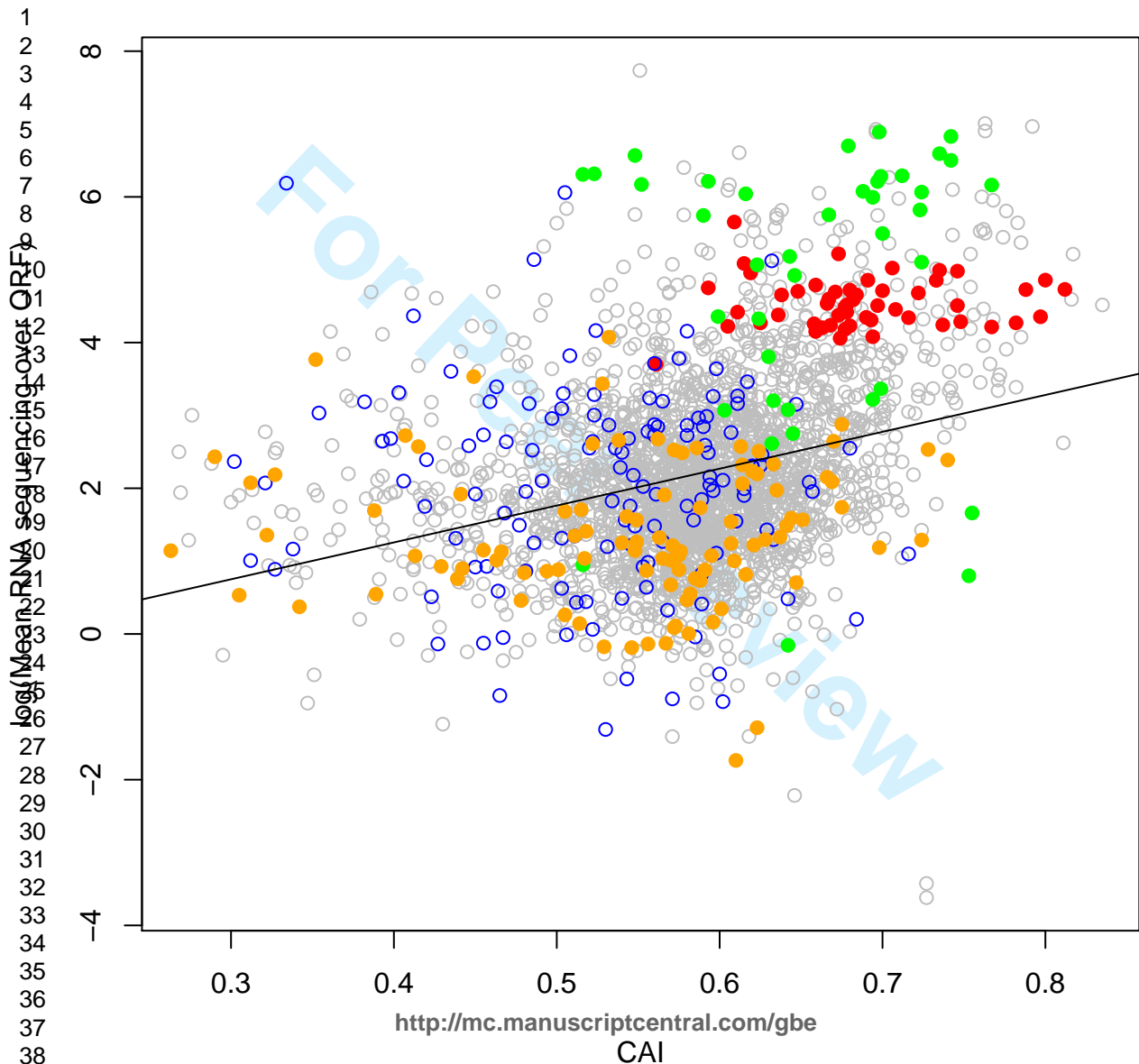




1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16

<http://mc.manuscriptcentral.com/gbe>





1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

