

1 **A hyperconserved protein in *Prochlorococcus* and marine *Synechococcus***

2

3

4

5 Olga Zhaxybayeva*¹, J Peter Gogarten² and W. Ford Doolittle¹

6

7 ¹Department of Biochemistry and Molecular Biology, Dalhousie University, 5850

8 College Street, Halifax, NS, B3H 1X5, Canada; olgazh@dal.ca, ford@dal.ca

9 ²Department of Molecular and Cell Biology, University of Connecticut, Storrs, CT,

10 USA; gogarten@uconn.edu

11

12 *corresponding author

13

14

15 The definitive version of this manuscript is available at [www.blackwell-](http://www.blackwell-synergy.com)

16 [synergy.com](http://www.blackwell-synergy.com). The full citation for the article is *FEMS Microbiol Lett*, 2007, **274**:

17 30-34

18

1 **Abstract**

2 We describe an open reading frame (encoding a hyperconserved protein, or
3 HCP) of unknown function in *Prochlorococcus*/marine *Synechococcus* genomes
4 that is 100% conserved on the amino acid level but lacks homologs outside this
5 group. Such conservation of an uncharacterized group-specific gene is unusual,
6 and unexpected given the extensive genomic divergence of this group at other
7 loci. Comparative analyses indicate that HCP is under stabilizing selection and
8 has resided in these genomes since the last common ancestor of the group.

9

10 **Introduction**

11 *Prochlorococcus* and marine *Synechococcus* are abundant unicellular marine
12 cyanobacteria (Waterbury *et al.*, 1986; Chisholm *et al.*, 1992), for which 13
13 completed genomes are now available (Table 1). Although members of the two
14 genera are more than 96% identical in their 16S rRNA sequences, their genomes
15 are surprisingly divergent, possibly due to high mutation rates and the lack of
16 several DNA repair enzymes in some of the genomes (Rocap *et al.*, 2003).
17 Pairwise Average Nucleotide Identity [ANI] values are as low as 64.5% and
18 Average Amino acid Identity [AAI] as low as 57.4%. Such divergence is found
19 even among the 933 “core” genes present in all 13 *Prochlorococcus* and marine
20 *Synechococcus* genomes. On average these genes show amino acid identity
21 across all 13 genomes at only 44% of their sites. It is noteworthy, then, that one
22 open reading frame (here designated HCP for Hyper-Conserved Protein), of 62

1 amino acids, is conserved at all positions in all genomes. In this article, we
2 characterize this protein *in silico*.

3

4 **Materials and Methods**

5 The ORF encoding the Hyper Conserved Protein (HCP) was used as query in
6 BLAST (Altschul *et al.*, 1997) searches against the nr database, 705 microbial
7 genomes and 21 environmental metagenomic databases that were available at
8 NCBI BLAST web server on January 30, 2007. BLASTP, TBLASTN and PSI-
9 BLAST searches, as well as CDD search with RPS-BLAST were performed. A
10 TBLASTN search of the Global Ocean Sampling Expedition database was
11 performed at the CAMERA version 0.7 website (<http://camera.calit2.net/>).
12 Prediction of transmembrane helices, coiled-coils and helix-turn-helix motifs, and
13 calculation of the predicted isoelectric point was performed using programs from
14 the EMBOSS suite (Rice *et al.*, 2000). The core of 13 completely sequenced
15 genomes was defined as a set of genes picking each other as top-scoring hits in
16 a fully transitive BLASTP search (E-value < 10⁻⁴). The core gene sets were
17 aligned in ClustalW v. 1.83 (Thompson *et al.*, 1994). The alignment conservation
18 was defined as proportion of alignment sites per gene that did not have amino
19 acid substitutions (gaps were treated as missing data) and calculated using an
20 in-house Perl script. Similarly, core genes and alignment conservation were
21 calculated for three other groups of genomes with several completely sequenced
22 genomes (genome lists are available upon request). The number of substitutions
23 per site in the most conserved 20 genes in the 13 genomes was calculated as a

1 total branch length of a phylogenetic tree reconstructed from the concatenated
2 alignment (the tree was reconstructed using NEIGHBOR from the PHYLIP
3 package (Felsenstein, 1993) with distances calculated in TREE-PUZZLE
4 (Schmidt *et al.*, 2002) under JTT+G substitution model). The probability of
5 observing no amino acid substitutions was calculated under a Poisson
6 distribution. Genome-wide ANI and AAI were calculated as described in
7 (Konstantinidis & Tiedje, 2005a; Konstantinidis & Tiedje, 2005b). Values of GC3
8 were calculated using in-house Perl script. 10 other cyanobacterial genomes
9 that were used in the genome context analyses are *Anabaena variabilis* ATCC
10 29413, *Gloeobacter violaceus* PCC 7421, *Nostoc* sp. PCC 7120, *Synechococcus*
11 *elongatus* PCC 6301, *Synechococcus elongatus* PCC 7942, *Synechococcus* sp.
12 JA-2-3B'a(2-13), *Synechococcus* sp. JA-3-3Ab, *Synechocystis* sp. PCC 6803,
13 *Thermosynechococcus elongatus* BP-1, *Trichodesmium erythraeum* IMS101 (all
14 available through the NCBI).

15

16 **Results and Discussion**

17 *Remarkable conservation of the HCP*

18 The HCP is conserved in all positions in all 13 *Prochlorococcus*/marine
19 *Synechococcus* genomes. Given that the twenty next most conserved core
20 genes have approximately 0.29 amino acid substitutions per site across all
21 genomes, 100% conservation is surprising. Assuming this substitution rate the
22 probability of observing a completely conserved protein with 62 amino acids is
23 2×10^{-8} .

1

2 *Misannotation of the HCP gene in three out of 13 genomes and on variation at*
3 *the C-terminal*

4 In the CC9902, MIT9515 and CC9311 genomes the HCP gene is annotated with
5 an alternative start codon. This appears to be an artifact of automated genome
6 annotations, since in the CC9902 genome this alternative start codon leaves only
7 three nucleotides between the upstream Trp-tRNA gene and the HCP, while in
8 the latter two genomes it overlaps with the Trp-tRNA gene by 29 nucleotides.
9 There appears to be a length variant of the HCP gene, resulting from mutation of
10 the stop codon TAA to CAA and hence causing two extra amino acids at the C
11 terminal (QK) present in two out of 13 genomes (MIT9312 and AS9601), as well
12 as in the Global Ocean Sampling Expedition database (see below).

13

14 *HCP is a group-specific protein*

15 This ORF is found in no other cyanobacterium (10 sequenced genomes), nor in
16 any other genome or database, with the exception of the “Global Ocean
17 Sampling Expedition” metagenome. Several samples of the Global Ocean
18 Sampling Expedition database had significant BLAST hits to the HCP (total of
19 119 hits with E-value $<10^{-27}$). The samples that return BLAST hits are from the
20 Sargasso Sea, Caribbean Sea, Eastern Tropical Pacific, and from the Pacific
21 Ocean near the Galapagos Islands. Somewhat unexpectedly, no significant
22 BLAST hits were found in samples from the North Pacific Subtropical Gyre
23 (DeLong *et al.*, 2006), but this might be due to this metagenome being not

1 sampled as deeply as the ones from the Global Ocean Sampling Expedition.
2 Also, no significant BLAST hits were found in the viral metagenomes that were
3 available for the BLAST at the CAMERA website (see Materials and Methods).
4
5 *HCP is not a recent transfer within the Prochlorococcus/marine Synechococcus*
6 *group*
7 At the nucleotide level there are 65 positions with substitutions within aligned
8 HCP genes, all synonymous, indicating that the ORF does indeed encode
9 protein, that it is under purifying selection, and that it has resided for some time in
10 these genomes. Indeed, the 13 analyzed genomes have a very broad range of
11 GC content, varying from 31% to 59% GC, and the GC content of the HCP gene
12 follows that of the host genome (although it is in general higher due to the fact
13 that 23 out its 62 amino acids require GC rich codons) (see Table 1). The
14 correlation between GC content of the HCP gene and the host genome is more
15 visible if only third codon positions are considered (see Fig. 1). A probable
16 scenario is that the HCP gene was inserted between the two tRNA genes in the
17 last common ancestor of the *Prochlorococcus/marine Synechococcus* group.
18 Based on microarray experiments, HCP is expressed in *Prochlorococcus*
19 *marinus* CCMP1986 (S.W. Chisholm, personal communication). The 5'-adjacent
20 region is moderately conserved and the Trp-tRNA gene immediately upstream
21 can be predicted to contain a promoter (Vogel *et al.*, 2003).

22

23 *Genomic context of the HCP coding gene*

1 Although dot plots (not shown) indicate that synteny is not widely conserved
2 among the 13 genomes, the HCP gene occurs in the same genomic context in
3 all, flanked by Glutamyl tRNA synthetase (gltX) and Asp-tRNA coding genes on
4 its 3' end and by the ribosomal protein L19 (rpl19) and the Trp-tRNA genes on its
5 5' end (Fig. 2). In the 10 other completely sequenced cyanobacterial genomes,
6 the gltX+Asp-tRNA and L19+Trp-tRNA gene clusters are usually found, but not in
7 close proximity to each other. In several genomes, the L19 and Trp-tRNA genes
8 are followed by a small ORF coding for SecE. However, no significant similarity
9 was found between the SecE and the HCP genes.

10

11 *HCP might be involved in DNA/RNA interactions*

12 Computational predictions of transmembrane helices, coiled-coils and helix-turn-
13 helix motifs did not return any significant results. However, the protein has a
14 predicted isoelectric point of 11.3, and a high proportion of positively charged
15 amino acids (16%), which suggests that it might be involved in DNA/RNA
16 interactions. For the MIT9313 genome only 47 (2%) of the proteins are more
17 alkaline and have a higher proportion of positively charged amino acids.

18

19 **Conclusions**

20 It is remarkable that HCP has been under such strong purifying selection that not
21 a single amino acid change has occurred during the divergence of the
22 *Prochlorococcus*/marine *Synechococcus* group (roughly estimated to be 100-200
23 millions years based on 16S rRNA divergence (Rappe & Giovannoni, 2003)) and

1 especially so since HCP is of unknown function. Examination of several other
2 groups for which multiple genomes are available shows that the highly conserved
3 core proteins (none were 100% conserved, but several families had 97-99%
4 conserved alignment sites) are usually well-characterized functionally, with
5 homologs outside the examined group (data not shown). The functional role of
6 the HCP in the *Prochlorococcus*/marine *Synechococcus* group therefore invites
7 further, experimental, investigation.

8

9 **Acknowledgements**

10 We would like to thank R. Thane Papke, Yuri Wolf and Kira Makarova for
11 insightful discussions and suggestions. OZ is supported through a CIHR
12 Postdoctoral Fellowship and is an honorary Killam Postdoctoral Fellow at
13 Dalhousie University.

14

15 **References:**

16 Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W & Lipman DJ
17 (1997) Gapped BLAST and PSI-BLAST: a new generation of protein
18 database search programs. *Nucleic Acids Res* **25**: 3389-3402.
19 Chisholm SW, Frankel SL, Goericke R, Olson RJ, Palenik B, Waterbury JB,
20 West-Johnsrud L & Zettler ER (1992) *Prochlorococcus marinus* nov. gen.
21 nov. sp.: an oxyphototrophic marine prokaryote containing divinyl
22 chlorophyll a and b. *Archives of Microbiology* **157**: 297.

- 1 DeLong EF, Preston CM, Mincer T *et al.* (2006) Community genomics among
2 stratified microbial assemblages in the ocean's interior. *Science* **311**: 496-
3 503.
- 4 Felsenstein J (1993) PHYLIP (Phylogeny Inference Package). Distributed by the
5 author. *Department of Genetics, University of Washington, Seattle.*
- 6 Konstantinidis KT & Tiedje JM (2005a) Genomic insights that advance the
7 species definition for prokaryotes. *Proc Natl Acad Sci U S A* **102**: 2567-
8 2572.
- 9 Konstantinidis KT & Tiedje JM (2005b) Towards a genome-based taxonomy for
10 prokaryotes. *J Bacteriol* **187**: 6258-6264.
- 11 Rappe MS & Giovannoni SJ (2003) The Uncultured Microbial Majority. *Annual*
12 *Review of Microbiology* **57**: 369-394.
- 13 Rice P, Longden I & Bleasby A (2000) EMBOSS: The European Molecular
14 Biology Open Software Suite. *Trends in Genetics* **16**: 276-277.
- 15 Rocap G, Larimer FW, Lamerdin J *et al.* (2003) Genome divergence in two
16 *Prochlorococcus* ecotypes reflects oceanic niche differentiation. *Nature*
17 **424**: 1042.
- 18 Schmidt HA, Strimmer K, Vingron M & von Haeseler A (2002) TREE-PUZZLE:
19 maximum likelihood phylogenetic analysis using quartets and parallel
20 computing. *Bioinformatics* **18**: 502-504.
- 21 Thompson JD, Higgins DG & Gibson TJ (1994) CLUSTAL W: improving the
22 sensitivity of progressive multiple sequence alignment through sequence

1 weighting, position-specific gap penalties and weight matrix choice.
2 *Nucleic Acids Res* **22**: 4673-4680.

3 Vogel J, Axmann IM, Herzel H & Hess WR (2003) Experimental and
4 computational analysis of transcriptional start sites in the cyanobacterium
5 *Prochlorococcus* MED4. *Nucleic Acids Res* **31**: 2890-2899.

6 Waterbury J, Watson S, Valois F & Franks D (1986) Biological and ecological
7 characterization of the marine unicellular cyanobacterium *Synechococcus*.
8 *Can Bull Fish Aquat Sci* **214**: 71-120.

9

1 **Table 1. GI numbers for the HCP in completely sequenced genomes*.**

Genome	GI number	GC content of HCP, %	GC content of the genome, %
<i>Prochlorococcus marinus</i> str. MIT 9312	78778859	42.6	31.2
<i>Prochlorococcus marinus</i> str. MIT 9313	33863575	55.0	50.7
<i>Prochlorococcus marinus</i> str. MIT 9303	124022390	55.0	50.0
<i>Prochlorococcus marinus</i> str. MIT 9515	123965773	45.5	30.8
<i>Prochlorococcus marinus</i> str. AS9601	123968066	43.1	31.3
<i>Prochlorococcus marinus</i> str. NATL1A	124025241	49.2	35.0
<i>Prochlorococcus marinus</i> str. NATL2A	72383643	48.7	35.1
<i>Prochlorococcus marinus</i> subsp. <i>marinus</i> str. CCMP1375	33239923	45.5	36.4
<i>Prochlorococcus marinus</i> subsp. <i>pastoris</i> str. CCMP1986	33861031	43.9	30.8
<i>Synechococcus</i> sp. WH 8102	33866357	60.8	59.4
<i>Synechococcus</i> sp. CC9605	78212192	58.7	59.2
<i>Synechococcus</i> sp. CC9902	78185284	58.7	54.2
<i>Synechococcus</i> sp. CC9311	113953642	52.4	52.4

2 *The HCP was also found in the following available unfinished genomes:

3 *Synechococcus* sp. WH7805, *Synechococcus* sp. WH5701, *Synechococcus* sp. RS9917,
4 *Synechococcus* sp. RS9916, *Synechococcus* sp. BL107, *Prochlorococcus marinus* str. MIT9211.

5

6

1 **Figure legends:**

2 **Fig. 1.** Comparison of GC content at the third codon positions in the HCP and in
3 each genome. GC3 values are given as a proportion of G+C nucleotides. The
4 error bars correspond to the 95% confidence interval calculated assuming a
5 binomial distribution for the GC3 values. The graph indicates that GC3 values of
6 the HCP and the host genome are equal within the sampling error.

7

8 **Fig. 2.** A. Genomic context of the HCP gene in 13 *Prochlorococcus* and marine
9 *Synechococcus* genomes. The unconserved region downstream from the HCP
10 gene varies from 149 to 302 nucleotides in length, and in CC9311 (the genome
11 with the largest, 302 nucleotide region) contains a hypothetical ORF that codes
12 for a 54 a.a. protein (locus tag sync_2165). The latter ORF does not have
13 matches anywhere in databases (including other marine
14 *Synechococcus/Prochlorococcus* genomes) and perhaps is decayed in other
15 genomes with smaller downstream regions. B. Genomic context of genes
16 adjacent to the HCP gene in three selected cyanobacterial genomes. See text
17 for more details.

18

19

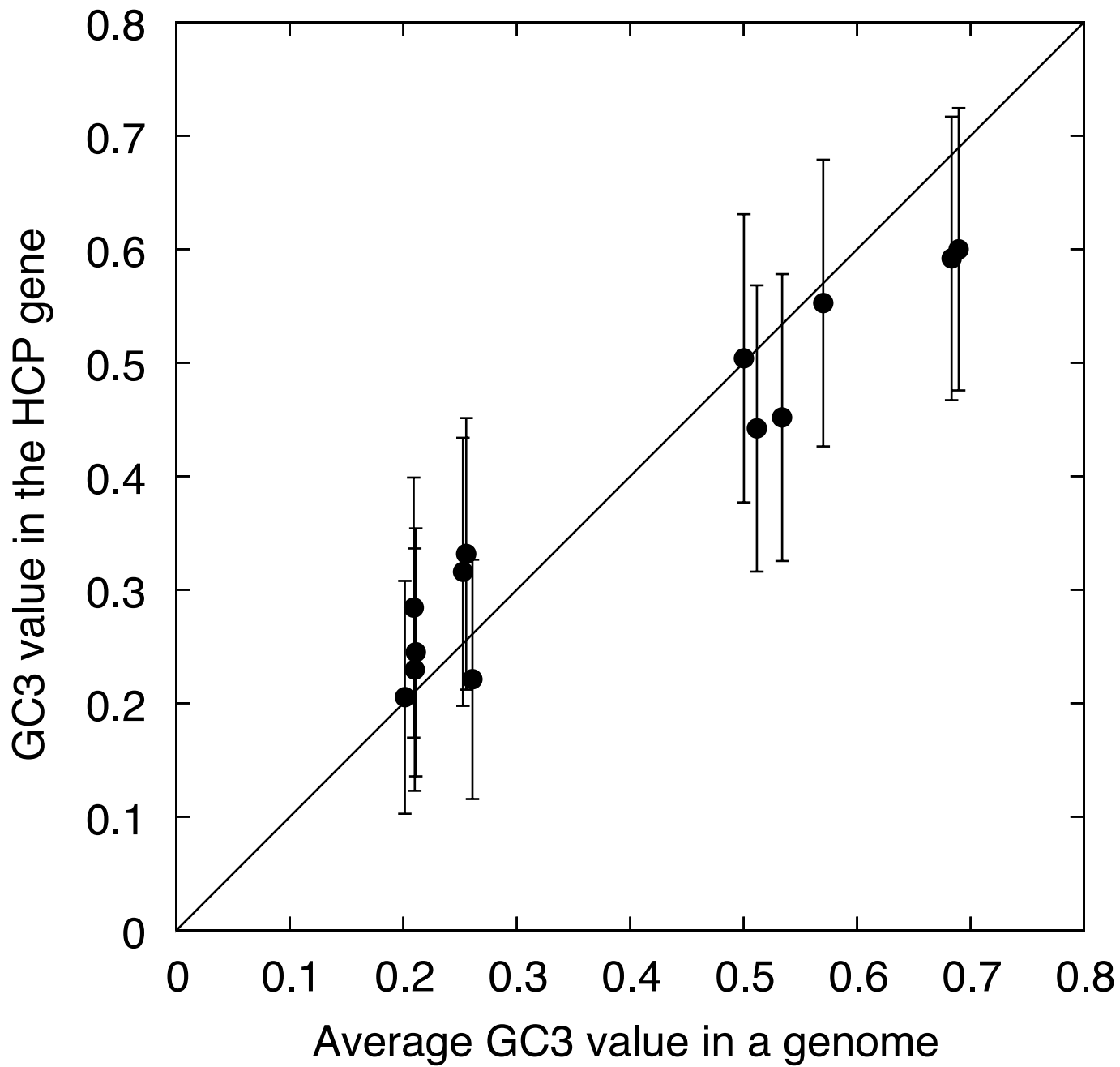


Figure 1

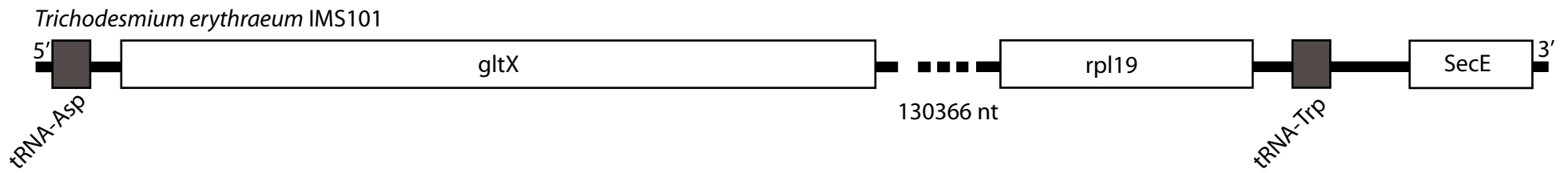
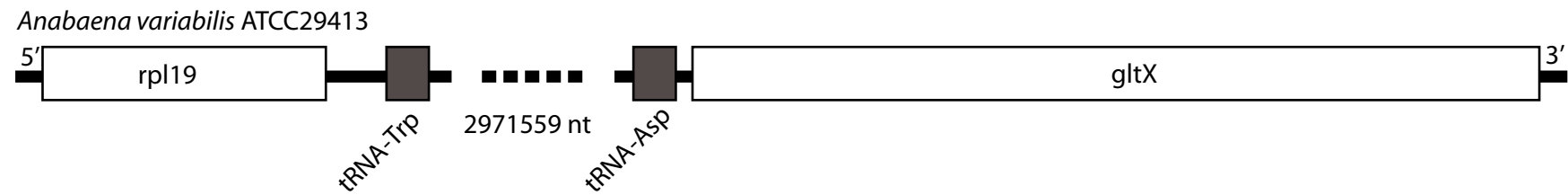
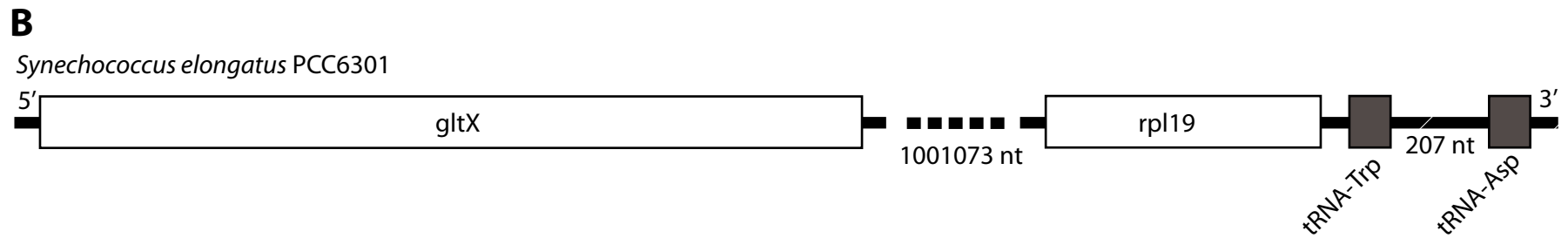
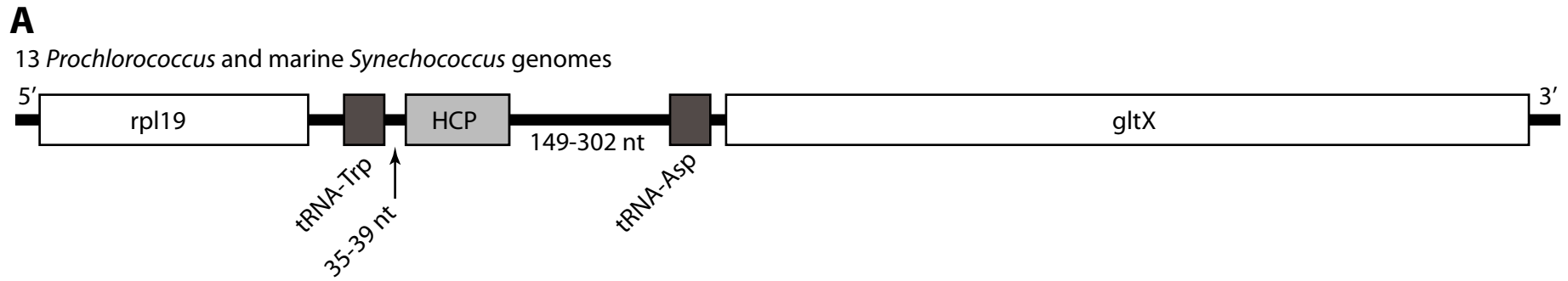


Figure 2