

Additional results for CGA / CGG reassessments

For each CGA and/or CGG reassignment, we provide the following data files: an Excel spreadsheet with detailed information on all genomes belonging to the reassigned clades and close outgroups (including information on inferred CGA/CGG amino acid, CheckM genome completeness estimate from GTDB, presence/absence of specific tRNA genes, codon usage in aligned Pfam domains, and genomic GC content), full alignments of BUSCO genes featured in figures, alignments of tRNA sequences from the reassigned clade and outgroups, and a consensus phylogenetic tree of the tRNA sequences.

Reassignment of CGA and CGG to tryptophan in Absconditabacteria

Genetic code inference

Two clades in the Candidate Phyla Radiation– Absconditabacteria (SR1) and Gracilibacteria (BD1-5)– were previously described to use an alternative genetic code where the canonical stop codon UGA is translated as glycine. We confirmed that genomes annotated as Absconditabacteria and Gracilibacteria on NCBI were predicted to translate UGA as glycine by our genetic code inference method (paper Table 1). Some of these genomes were additionally predicted to have reassessments of the canonical arginine codons CGA and/or CGG as tryptophan.

Species phylogenetic tree

Since many genomes derived from the uncultivated Candidate Phyla Radiation have only broad phylogenetic annotation in NCBI, we investigated the phylogenetic distribution of this reassignment on the Genome Taxonomy Database (GTDB). In Figure 1, we show a part of the GTDB phylogenetic tree that includes all members of the Absconditabacteria (GTDB order o_Absconditabacterales, 10 species clusters), the sister group Gracilibacteria (GTDB order o_BD1-5, 15 species clusters), and all outgroup species out to Peregrinibacteria and Peribacteria (46 species clusters).

In order to help gauge whether a tRNA gene is not present in the genome or whether it is missing due to an incomplete assembly, we marked the most complete genomes on the tree. We believe that the CheckM estimates of genome completeness (based on presence of conserved protein sets and are provided on GTDB for each genome) tend to underestimate the completeness of Candidate Phyla Radiation genomes due to lineage-specific gene loss. Therefore, we additionally tabulated the presence of a minimal set of 22 required tRNA genes (excluding CGN-decoding tRNAs) found by tRNAscan-SE 2.0 (see Methods). In Figure 1, we marked the genomes which contained all 22 required tRNA genes and had CheckM completeness estimate of 75% or greater (maximum observed was 90%).

Reassignment of UGA to glycine in both Absconditabacteria and Gracilibacteria

On the GTDB phylogenetic tree (Figure 1), all members of the Absconditabacteria and Gracilibacteria were inferred to use the known reassignment of UGA as glycine, while none of the outgroup genomes were predicted to have an amino acid translation of the stop codon UGA. Consistent with that prediction, all members of the UGA-reassigned clades do not have an release factor 2 (RF2) gene (which terminates translation at UAA and UGA codons) and most in-

stead contain a glycine-type tRNA_{UCA}. The absence of a tRNA_{UCA} in some genomes may be due to the incomplete nature of most of these genomes, as they were either assembled from metagenomic data or derived from single-cell genome sequencing. None of the outgroup genomes encoded a tRNA_{UCA} and all instead had an RF2 gene.

Reassignment of CGA and CGG to tryptophan in Absconditabacteria

In the Absconditabacteria clade on the tree (Figure 1), all 10 species were inferred to translate CGA as tryptophan and CGG as either tryptophan, uninferred, or arginine (in one species).

Figure 2 shows four example alignments of conserved single-copy bacterial genes across Absconditabacteria, Gracilibacteria, and outgroup species. The four aligned genes show that multiple CGA and CGG positions across all 10 Absconditabacteria genomes appear to be used interchangeably with the canonical tryptophan codon UGG within the Absconditabacteria and align to positions often conserved for tryptophan in the Gracilibacteria and outgroup species. In contrast, outgroup CGAs and CGGs occur at positions broadly conserved for arginine.

The tRNAs present in Absconditabacteria genomes are consistent overall with the inferred codon translations (Figure 1). The tRNA_{CCG} and tRNA_{UCG} genes, which decode CGA and CGG codons, were predominantly tryptophan-type (classified by G73 discriminator base and the absence of A20), supporting the inferred translation for most Absconditabacteria genomes. Some genomes were missing tRNAs, however this is likely due to the incomplete nature of genomes from uncultivated bacteria.

One exception was the tRNA_{CCG} from GCA_002414185.1 (Patescibacteria group bacterium UBA5124; #9 on tree), which is the only Absconditabacteria genome with an arginine translation inferred for CGG. The tRNA_{CCG} could not be classified as either arginine- or tryptophan-type because it has a A73 discriminator base, which supports arginine identity, but has a C20 in the D-loop instead of the A20 that would be expected of arginine tRNAs. It is possible that this

unusual tRNA supports translation of CGG as arginine to some extent in this species.

In a phylogenetic tree built from arginine and tryptophan tRNA from across the Absconditabacteria, Gracilibacteria, and outgroup species (Figure 3), all tRNA_{CCG} and tRNA_{UCG} genes from Absconditabacteria (including GCA_002414185.1, #9 on species tree) cluster within the clade of tRNA_{CCA}^{Trp} genes from across the three groups. We compared the likelihood of two phylogenetic models, one where the Absconditabacteria tRNA_{CCG} and tRNA_{UCG} sequences are constrained to cluster with the tryptophan tRNA_{CCA} genes against another model where they are constrained to cluster with the arginine tRNAs. The log2 ratio of the likelihoods is 43, strongly favoring the tryptophan model. This indicates that the Absconditabacteria tRNA_{CCG} and tRNA_{UCG} genes are more similar in sequence to tRNA_{CCA}^{Trp} genes and not arginine tRNAs.

In Absconditabacteria, both CGA and CGG tend to be rare with codon usages of 14-37 per 10,000 codons in regions aligned to Pfam domains for CGA and 1-24 per 10,000 codons for CGG (Figure 1). At the extreme, Absconditabacteria genomes where CGG was either uninferred or inferred to code for arginine had the lowest CGG codon usage (<3 per 10,000 codons) resulting in fewer than 10 Pfam positions aligning to CGG. In contrast, other Absconditabacteria genomes such as GCA_002791215.1 (candidate division SR1 bacterium CG_4_9_14_3_um_filter_40_9; #8 on tree) had CGA and CGG codon usages (30 and 24 per 10,000) approaching that of the standard tryptophan codon UGG (35 per 10,000). The overall low usage of CGN codons in Absconditabacteria may be tied to the low GC content of the clade, which ranges between 0.29-0.38. GC content-driven substitutions towards AT-rich arginine codons, such as AGA and AGG, may have facilitated reassignment of CGA and CGG by reducing the number of compensatory substitutions needed to adapt to the new translation. Low GC content has been previously suggested to have played a role in the reassignment of UGA from stop to glycine in this clade by disfavoring UGA stop codon usage in favor of UAA and by creating a less GC rich codon for glycine (typically GGN) (1).

Translation of CGA and CGG in Gracilibacteria

The Gracilibacteria genomes on the tree (Figure 1) split into two clades (dubbed here as “clade 1” and “clade 2”) that differ in codon usage, genomic GC content, inferred CGA/CGG translation, and tRNA decoding capability. The 8 species in Gracilibacteria clade 1 have very low genomic GC content (ranging between 0.21-0.29) and extremely low usage of CGA and CGG codons (<7 per 10,000 for all genomes). Consequently, CGA and CGG codon meaning are uninferred for many of these species due to lack of aligned Pfam positions. When inferred, CGA is predicted to be an arginine codon while CGG is predicted to code for tryptophan. In contrast, the 7 species in Gracilibacteria clade 2 have higher genomic GC content (ranging between 0.35-0.53), CGA and CGG are abundantly used codons (33-254 per 10,000). Both CGA and CGG are inferred to be arginine codons in all Gracilibacteria clade 2 species, as in all outgroup genomes on the tree.

In the example alignments of conserved single-copy genes (Figure 2), CGA occurs at conserved arginine positions in members of Gracilibacteria clade 1 and clade 2; only one aligned CGA position in Gracilibacteria clade 1 is shown. In Gracilibacteria clade 1, three aligned CGG positions appear in the alignments; two of them coming from GCA_002435385.1 (Patescibacteria group bacterium UBA6489; #13 on tree). These CGG positions are encoded by the tryptophan codon UGG in other Gracilibacteria clade 1 genomes, and are sometimes more broadly conserved for tryptophan in Absconditabacteria and outgroup genomes. This supports possible translation of CGG as tryptophan in some members of Gracilibacteria clade 1.

The tRNAs to decode CGA and CGG in Gracilibacteria support the divide between the clade 1 and clade 2 (Figure 1). Clade 2 genomes all contain an arginine-type tRNA_{UCG}, which would translate both CGA and CGG codons as arginine. In contrast, none of the clade 1 genomes contain a tRNA_{UCG}. It is possible that we failed to find this tRNA due to a long intron or incomplete genome assembly, but if it is indeed missing, then there does not exist any tRNA

that recognizes CGA by conventional anticodon pairing rules in Gracilibacteria clade 1.

In Gracilibacteria clade 1, 6 of the 8 genomes contain a CGG-decoding tRNA_{CCG}. We could not confidently assign either arginine or tryptophan isotype to these tRNAs because they have an unusual D-loop sequence that makes it unclear which nucleotide is at position 20, contain the minor tryptophan identity element A1:U72 and a G73 discriminator which is compatible with either arginine or tryptophan isotype. We constructed a phylogenetic tree from an alignment of arginine and tryptophan tRNAs in Absconditabacteria, Gracilibacteria, and outgroups (Figure 3). In the consensus tree, the unusual tRNA_{CCG} sequences from Gracilibacteria clade 1 form a clade within the cluster of tryptophan tRNA_{CCA} from all groups and the reassigned tRNA_{UCG} and tRNA_{CCG} sequences from Absconditabacteria, and not with the arginine tRNAs from all groups. We compared the likelihoods of two phylogenetic models, one where the Gracilibacteria clade 1 tRNA_{CCG} sequences are constrained to form a clade with all tryptophan tRNAs and the reassigned tRNA_{CCG} and tRNA_{UCG} tRNAs from Absconditabacteria against another model where the Gracilibacteria clade 1 tRNA_{CCG} sequences are constrained to group with all arginine tRNAs (while the Absconditabacteria tRNA_{UCG} and tRNA_{CCG} cluster with tryptophan tRNAs). The log2 ratio of these two likelihoods is 3.9, indicating support for Gracilibacteria tRNA_{CCG} grouping with tryptophan tRNA sequences. All of this suggests that the clade 1 tRNA_{CCG} is more similar in sequence to tryptophan tRNAs and not arginine tRNAs, and was most likely derived from a tryptophan tRNA; however, this does not necessarily mean that this tRNA is charged with tryptophan *in vivo*.

All together, the evidence paints a picture where Gracilibacteria clade 2 species decode CGA and CGG as arginine, while clade 1 genomes might treat CGA as an arginine codon and CGG as a tryptophan codon but many questions remain regarding how these codons are decoded in living cells. It is unclear from tRNA gene sequence alone how the CGG-decoding tRNA_{CCG} in group 1 is aminoacylated. Charging with tryptophan would be consistent with the occurrence

of CGG in conserved tryptophan residues in conserved single-copy gene alignments in a few species and placement of tRNA_{CCG} with tryptophan tRNAs on a tRNA phylogeny. It is also possible that this unusual tRNA may be charged with arginine or even recognized by multiple aminoacyl-tRNA synthetases and ambiguously charged. If group 1 Gracilibacteria indeed do not have tRNA capable of reading CGA, CGA might be inefficiently decoded by non-cognate tRNAs, possibly by the arginine tRNA_{GCG} or perhaps the unusual tRNA_{CCG}.

Reassignment of CGG to glutamine in a clade of Clostridia

Genetic code inference

In our analysis of all sequenced bacterial genomes, we noticed that 6 genomes were predicted to translate the canonical arginine codon CGG as glutamine. These six genomes included uncultivated Clostridia assembled from metagenomic sequences and culturable species such as *Clostridium hiranonis*. Four of these genomes are included in the Genome Taxonomy Database (GTDB) phylogeny, where they are assigned to three species clusters within the GTDB genus Clostridium_U (family Peptostreptococcaceae, order Peptostreptococcales, class Clostridia), which additionally includes *Peptacetobacter hominis* (our program did not infer an amino acid translation for CGG in this species).

Species phylogenetic tree

In Figure 4, we show a part of the GTDB phylogenetic tree that includes all members of the GTDB genus Clostridium_U (4 species clusters) and outgroup species including the remaining species in the family Peptostreptococcaceae and the families Peptoclostridiaceae and Filifactoraceae (52 species clusters).

In order to help gauge whether a tRNA gene is not present in the genome or whether it is missing due to an incomplete genome, we marked the most complete genomes on the tree

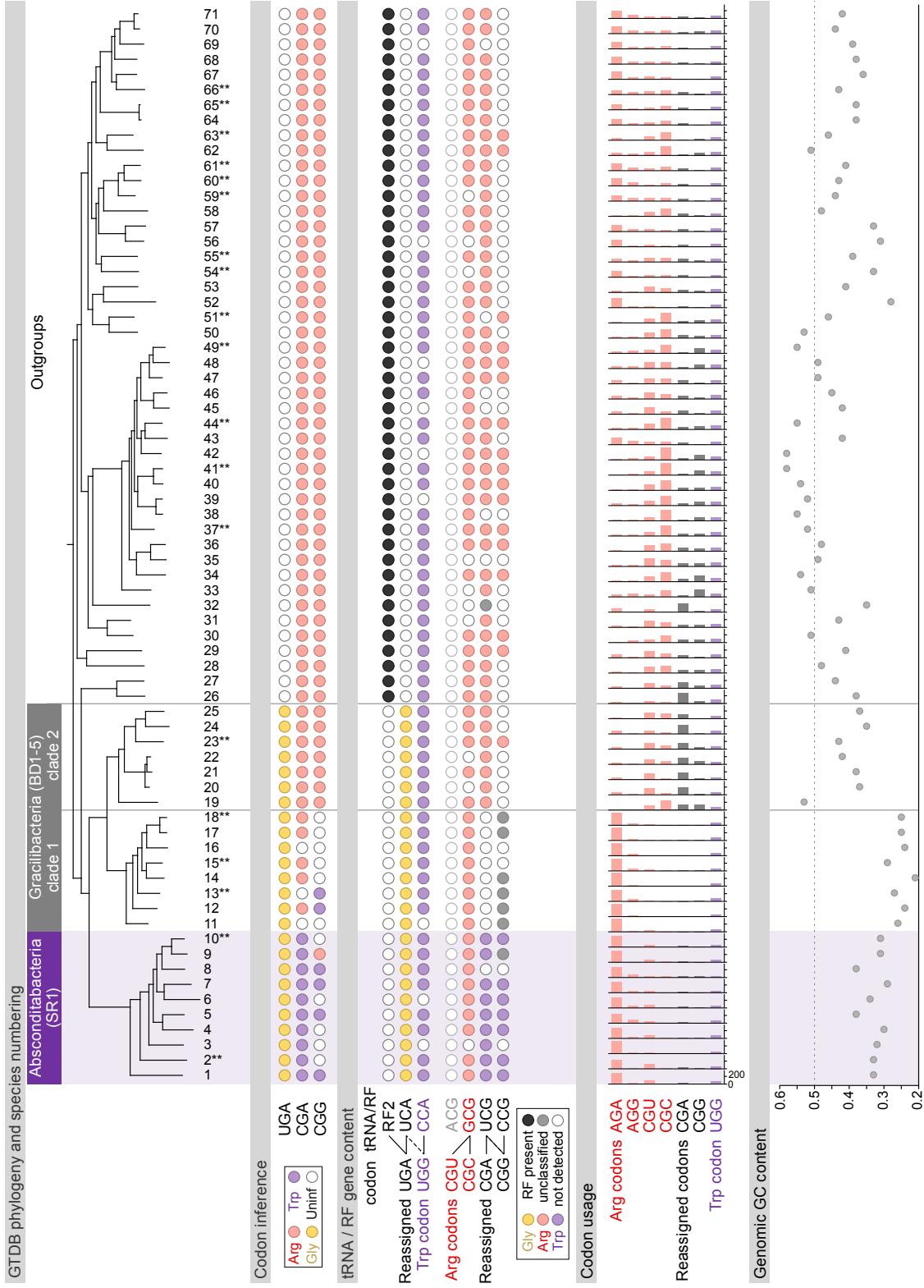


Figure 1: Phylogenetic tree from GTDB showing the Absconditabacteria, Gracilibacteria, and closest outgroup genomes (Peregrinibacteria and Peribacteria), each species indicated by a number which can be cross-referenced with the summary spreadsheet. Asterisks indicate genomes the most complete genomes, which have CheckM estimated genome completeness from GTDB $>75\%$ and the entire minimal set of 22 required tRNAs (excluding CGN-decoding tRNAs). For each species, the inferred translation of the three reassigned codons (UGA, CGA, and CGG) by our method is indicated by colored circles (red: arginine, purple: tryptophan, yellow: glycine, white: uninferred). The presence of release factor and tRNA genes that recognize the UGR- and CGN-codons is also indicated by filled circles, colored black if present and white is absent for release factor 2 and according to the predicted amino acid charging based on identity elements for tRNAs (see Methods). Anticodons in gray font are typically not found in the Candidate Phyla Radiation. The lines connecting codons and tRNA anticodons represent the likely wobble decoding capabilities, with dashed lines representing weaker interactions. For the anticodon UCG, U34 is presumed to be modified in a way that restricts wobble to CGA and CGG, at least in the Absconditabacteria, but could potentially recognize CGU and/or CGC depending on the true modification state. The remaining anticodons are not expected to be modified in a way that alters which codons are recognized. The codon usage for the reassigned codons CGA and CGG, the arginine codons AGA, AGG, CGU, and CGC, and the tryptophan codon UGG is the frequency per 10,000 codons aligned to Pfam positions. Genomic GC content is calculated over the entire genome.

(Figure 4), which had a CheckM completeness estimate of 98% or greater and contained the entire minimal set of 22 tRNA genes (excluding CGN-decoding tRNAs) that are required in all bacteria (see Methods). We noticed that some members of the reassigned clade and outgroup had very high CheckM completeness scores ($>98\%$) but were missing more than half of the required tRNA set. We attribute this due to the fact that tRNA genes are known to be highly clustered in Firmicutes (2) so it is possible for an assembly to be missing a tRNA cluster (and therefore a large fraction of tRNAs) without affecting the CheckM completeness significantly (which is based on the presence of non-clustered protein coding genes). As an example of the highly clustered nature of tRNA genes in this clade, the genome of *Peptostreptococcus stomatis* DSM 17678 (GCA_000147675.2, #40 on tree) was found to have 86 tRNA genes by tRNAscan-SE 2.0, 71 of which are located in one of two clusters (one 4kb cluster of 42 tRNAs and one 8kb cluster of 29 tRNAs).

Undecaprenyl diphosphate synthase (POG091H00BZ)

Absconditabacteria	2 . . . DGN [*] R [*] T [*] Q [*] A [*] Q [*] E [*] R [*] L [*] P [*] S [*] I [*] F [*] G [*] H [*] . . . ITL [*] α [*] G [*] A [*] S [*] T [*] E [*] N [*] I [*] Q [*] E [*] R [*] S [*] . . . FMT [*] α [*] W [*] I [*] S [*] Y [*] . . . KALE [*] α [*] Y [*] D [*] S [*] I [*] V [*] K [*] Y [*] R [*] N [*] F [*] G [*] K [*] 4 . . . DGN [*] R [*] T [*] W [*] A [*] K [*] E [*] L [*] G [*] K [*] T [*] S [*] L [*] G [*] E [*] H [*] . . . FTL [*] α [*] GL [*] S [*] T [*] E [*] N [*] L [*] K [*] N [*] R [*] T [*] . . . FML [*] α [*] W [*] I [*] G [*] Y [*] . . . KAIE [*] α [*] α [*] N [*] S [*] L [*] D [*] S [*] Q [*] N [*] F [*] G [*] K [*] 7 . . . DGN [*] R [*] T [*] W [*] A [*] K [*] L [*] N [*] Q [*] T [*] I [*] P [*] E [*] A [*] Y [*] . . . FTL [*] α [*] GL [*] S [*] T [*] E [*] N [*] A [*] N [*] K [*] R [*] P [*] . . . FMS [*] W [*] α [*] V [*] G [*] Y [*] . . . EALK [*] α [*] F [*] D [*] K [*] M [*] A [*] L [*] R [*] N [*] Y [*] G [*] K [*] 8 . . . DGN [*] R [*] T [*] Q [*] A [*] K [*] A [*] G [*] K [*] D [*] L [*] P [*] Q [*] A [*] Y [*] . . . FTL [*] α [*] GL [*] S [*] T [*] E [*] N [*] T [*] K [*] N [*] R [*] P [*] . . . FMS [*] α [*] I [*] G [*] Y [*] . . . ESLK [*] α [*] F [*] N [*] A [*] M [*] E [*] K [*] R [*] N [*] F [*] G [*] K [*] 9 . . . DGN [*] R [*] T [*] Q [*] A [*] R [*] E [*] S [*] N [*] R [*] T [*] V [*] α [*] E [*] A [*] Y [*] . . . ITL [*] W [*] G [*] L [*] S [*] T [*] E [*] N [*] T [*] K [*] R [*] P [*] . . . FMS [*] W [*] α [*] I [*] G [*] Y [*] . . . EALN [*] W [*] F [*] N [*] N [*] I [*] E [*] K [*] R [*] N [*] F [*] G [*] K [*] 10 . . . DGN [*] R [*] T [*] Q [*] A [*] R [*] E [*] S [*] N [*] R [*] T [*] V [*] α [*] E [*] A [*] Y [*] . . . FTL [*] α [*] GL [*] S [*] T [*] E [*] N [*] T [*] K [*] R [*] P [*] . . . FMS [*] W [*] α [*] I [*] G [*] Y [*] . . . EALAW [*] F [*] D [*] T [*] M [*] A [*] E [*] K [*] R [*] N [*] F [*] G [*] K [*] 13 . . . DGN [*] R [*] W [*] A [*] E [*] S [*] K [*] M [*] L [*] P [*] K [*] V [*] A [*] G [*] . . . ITL [*] W [*] A [*] L [*] STEN [*] L [*] I [*] K [*] R [*] D [*] . . . FLLFD [*] SAY [*] . . . EAIDT [*] FNK [*] --AK [*] R [*] N [*] F [*] G [*] K [*] 18 . . . DGN [*] R [*] W [*] A [*] E [*] K [*] G [*] F [*] P [*] K [*] F [*] V [*] G [*] H [*] . . . LTI [*] W [*] A [*] L [*] S [*] V [*] D [*] E [*] K [*] R [*] E [*] . . . FLLFD [*] SEY [*] . . . KAIDS [*] FAN [*] --SK [*] R [*] N [*] F [*] G [*] K [*] 19 . . . DGN [*] R [*] Q [*] A [*] W [*] A [*] K [*] N [*] G [*] L [*] V [*] K [*] T [*] I [*] G [*] H [*] . . . ASA [*] W [*] A [*] L [*] A [*] K [*] N [*] V [*] E [*] Y [*] D [*] . . . FLLYASEY [*] . . . QALAW [*] YDG [*] --CQ [*] R [*] N [*] F [*] G [*] Y [*] 23 . . . DGN [*] R [*] Q [*] A [*] W [*] A [*] K [*] L [*] G [*] N [*] L [*] A [*] L [*] A [*] G [*] H [*] . . . VSM [*] W [*] A [*] S [*] K [*] E [*] N [*] K [*] R [*] S [*] . . . YFLYQSAY [*] . . . EAIMS [*] FEG [*] --TK [*] R [*] N [*] F [*] G [*] K [*] 36 . . . DGN [*] R [*] W [*] A [*] R [*] A [*] Q [*] G [*] W [*] H [*] P [*] W [*] D [*] G [*] H [*] . . . LTI [*] W [*] C [*] F [*] S [*] TEN [*] W [*] K [*] R [*] D [*] . . . FFL [*] W [*] Q [*] S [*] V [*] . . . QILDKYHQ [*] --R [*] Y [*] α [*] FGG 42 . . . DGN [*] R [*] W [*] A [*] R [*] A [*] Q [*] G [*] L [*] Q [*] P [*] W [*] K [*] G [*] H [*] . . . LTV [*] W [*] C [*] F [*] S [*] TEN [*] W [*] K [*] R [*] E [*] . . . FAL [*] W [*] Q [*] S [*] V [*] . . . RAVDAFTL [*] --RT [*] α [*] FGA 54 . . . DGN [*] R [*] W [*] A [*] L [*] I [*] N [*] K [*] T [*] K [*] M [*] E [*] G [*] H [*] . . . LTI [*] W [*] G [*] L [*] S [*] E [*] N [*] L [*] K [*] E [*] Y [*] . . . FLP [*] W [*] Q [*] S [*] V [*] . . . KAI [*] YYNG [*] --AK [*] R [*] N [*] F [*] G [*] R [*] * * * * *
--------------------	---

tRNA (guanine-N1)-methyltransferase (POG091H01WE)

Absconditabacteria	1 . . . EATV [*] R [*] L [*] L [*] P [*] G [*] V [*] I [*] G [*] E [*] A [*] S [*] W [*] Q [*] Y [*] E [*] S [*] Y [*] . . . NHAAIEQ [*] α [*] R [*] K [*] D [*] N [*] . . . 2 . . . EAVV [*] R [*] L [*] L [*] P [*] G [*] V [*] I [*] T [*] S [*] W [*] I [*] E [*] S [*] Y [*] . . . HHKKIAE [*] W [*] K [*] K [*] D [*] N [*] . . . 3 . . . EAIT [*] R [*] L [*] L [*] P [*] G [*] V [*] I [*] K [*] A [*] Q [*] S [*] E [*] D [*] Y [*] . . . DHKKIEE [*] α [*] K [*] K [*] E [*] Q [*] . . . 4 . . . ESIT [*] R [*] L [*] I [*] P [*] G [*] V [*] I [*] E [*] S [*] W [*] Q [*] N [*] E [*] Y [*] . . . NTEEILK [*] α [*] R [*] K [*] N [*] N [*] . . . 6 . . . EAIS [*] R [*] L [*] V [*] P [*] G [*] V [*] I [*] E [*] S [*] G [*] E [*] S [*] Y [*] . . . DQKKIEE [*] W [*] R [*] K [*] E [*] . . . 7 . . . EAIT [*] R [*] L [*] V [*] P [*] G [*] V [*] I [*] E [*] S [*] G [*] E [*] S [*] Y [*] . . . HTKKIEA [*] W [*] K [*] D [*] K [*] N [*] . . . 8 . . . ESIV [*] R [*] L [*] V [*] P [*] G [*] V [*] I [*] E [*] S [*] G [*] E [*] Y [*] . . . HHKNIEE [*] α [*] R [*] R [*] K [*] K [*] . . . 9 . . . ESVV [*] R [*] L [*] I [*] P [*] G [*] V [*] I [*] E [*] S [*] W [*] Q [*] N [*] E [*] . . . HHKNIEK [*] W [*] K [*] K [*] N [*] . . . 10 . . . ESIV [*] R [*] L [*] I [*] P [*] N [*] V [*] I [*] E [*] S [*] W [*] Q [*] N [*] E [*] . . . HTKKIEE [*] W [*] K [*] K [*] D [*] N [*] . . . 13 . . . DALI [*] α [*] H [*] I [*] P [*] G [*] V [*] L [*] N [*] E [*] K [*] S [*] L [*] E [*] E [*] S [*] F [*] . . . NHKKIED [*] W [*] K [*] R [*] D [*] N [*] . . . 18 . . . DSFVR [*] N [*] I [*] S [*] G [*] V [*] L [*] G [*] N [*] K [*] L [*] S [*] L [*] E [*] D [*] S [*] F [*] . . . NHAEIEN [*] W [*] K [*] K [*] N [*] . . . 23 . . . DAVV [*] R [*] L [*] L [*] P [*] G [*] V [*] I [*] Q [*] S [*] -DS [*] R [*] E [*] E [*] S [*] F [*] . . . DMKAIE [*] T [*] W [*] K [*] N [*] H [*] . . . 41 . . . DAI [*] α [*] Q [*] I [*] P [*] G [*] V [*] L [*] G [*] D [*] E [*] S [*] A [*] T [*] E [*] S [*] F [*] . . . HHKEIEK [*] W [*] R [*] K [*] A [*] N [*] . . .
--------------------	--

Peptidase M50 (POG091H0131)

Absconditabacteria	2 . . . IPPKVATL [*] W [*] K [*] D [*] G [*] S [*] T [*] E [*] Y [*] . . . SFITASFLS [*] K [*] T [*] L [*] I [*] L [*] L [*] . . . 4 . . . LPPKICNL [*] W [*] K [*] D [*] K [*] G [*] T [*] Q [*] Y [*] . . . TLFTAPL [*] W [*] K [*] R [*] L [*] I [*] V [*] F [*] . . . 5 . . . MPPIATL [*] α [*] T [*] D [*] K [*] G [*] T [*] K [*] Y [*] . . . SFVKAKL [*] α [*] K [*] L [*] I [*] I [*] S [*] . . . 6 . . . IPPKAFKI [*] W [*] T [*] D [*] K [*] G [*] T [*] E [*] Y [*] . . . SFIKAKV [*] W [*] K [*] I [*] I [*] L [*] L [*] . . . 7 . . . IPPKACKL [*] G [*] K [*] D [*] K [*] G [*] T [*] Q [*] Y [*] . . . SLIKAPI [*] H [*] K [*] I [*] I [*] M [*] L [*] . . . 8 . . . IPPKVMTL [*] Y [*] K [*] D [*] K [*] G [*] T [*] E [*] Y [*] . . . SFIKAKVL [*] W [*] K [*] I [*] I [*] L [*] L [*] . . . 9 . . . IPPKICKI [*] W [*] T [*] D [*] K [*] G [*] T [*] E [*] Y [*] . . . SFIKAKVL [*] P [*] K [*] I [*] I [*] L [*] L [*] . . . 13 . . . IPPRAKKIG [*] K [*] D [*] K [*] H [*] G [*] T [*] I [*] Y [*] . . . NLSNKPA [*] Y [*] Q [*] S [*] I [*] I [*] V [*] V [*] . . . 18 . . . IPPRAKKL [*] F [*] T [*] D [*] K [*] G [*] T [*] I [*] F [*] . . . NLTNKPA [*] W [*] Q [*] S [*] I [*] I [*] I [*] L [*] . . . 23 . . . IPPRAKTL [*] F [*] Q [*] D [*] K [*] H [*] G [*] T [*] I [*] Y [*] . . . SFATKSWL [*] A [*] Q [*] S [*] A [*] V [*] L [*] L [*] . . . 25 . . . IPPKIKK [*] I [*] F [*] T [*] D [*] K [*] G [*] T [*] D [*] F [*] . . . AFMSKSLPK [*] α [*] L [*] L [*] V [*] L [*] V [*] . . . 49 . . . LPPKV [*] K [*] V [*] L [*] F [*] Y [*] -K [*] R [*] G [*] T [*] E [*] F [*] . . . SFGAATI [*] W [*] Q [*] Y [*] V [*] M [*] I [*] L [*] S [*] . . . 55 . . . LPPRIFGI [*] -K [*] R [*] G [*] T [*] E [*] Y [*] . . . SFMSKS [*] I [*] G [*] V [*] R [*] T [*] K [*] V [*] L [*] . . .
--------------------	---

DNA-directed RNA polymerase subunit beta (POG091H02K5)

Absconditabacteria	1 . . . DDEKHRLI [*] I [*] S [*] I [*] W [*] S [*] D [*] A [*] K [*] T [*] . . . SGA [*] R [*] G [*] T [*] α [*] Q [*] G [*] M [*] T [*] Q [*] M [*] A [*] . . . 2 . . . DDEKHRO [*] Q [*] I [*] V [*] O [*] I [*] W [*] T [*] D [*] I [*] K [*] N [*] . . . SGA [*] R [*] G [*] S [*] Y [*] N [*] S [*] T [*] Q [*] I [*] L [*] . . . 5 . . . EAEKHRLI [*] V [*] K [*] I [*] W [*] T [*] A [*] V [*] K [*] K [*] . . . SGA [*] R [*] G [*] S [*] Q [*] T [*] H [*] L [*] T [*] Q [*] I [*] S [*] . . . 6 . . . EQEKHRLI [*] I [*] N [*] V [*] α [*] S [*] E [*] V [*] K [*] T [*] . . . SGA [*] R [*] G [*] S [*] T [*] N [*] T [*] V [*] Q [*] I [*] S [*] . . . 7 . . . EEEKHRLI [*] I [*] K [*] V [*] W [*] T [*] D [*] V [*] K [*] S [*] . . . SKA [*] R [*] G [*] S [*] Q [*] T [*] H [*] L [*] T [*] Q [*] I [*] S [*] . . . 8 . . . DQEKHRN [*] V [*] E [*] I [*] W [*] S [*] K [*] V [*] G [*] K [*] G [*] . . . SGA [*] R [*] G [*] S [*] Q [*] T [*] H [*] I [*] T [*] Q [*] I [*] S [*] . . . 10 . . . DDEKHRS [*] I [*] K [*] I [*] α [*] T [*] E [*] V [*] K [*] T [*] . . . SGA [*] R [*] G [*] S [*] Q [*] T [*] H [*] M [*] T [*] Q [*] I [*] S [*] . . . 11 . . . ENEKYSQ [*] S [*] I [*] L [*] W [*] A [*] D [*] V [*] K [*] . . . SGA [*] R [*] G [*] N [*] Y [*] G [*] N [*] V [*] T [*] Q [*] L [*] C [*] . . . 13 . . . EEEKYQ [*] A [*] S [*] I [*] A [*] Y [*] E [*] T [*] K [*] . . . SGA [*] R [*] G [*] N [*] W [*] G [*] N [*] V [*] T [*] Q [*] L [*] C [*] . . . 18 . . . EDEKYNQ [*] S [*] I [*] K [*] W [*] A [*] Q [*] V [*] K [*] N [*] . . . SGA [*] R [*] G [*] N [*] W [*] G [*] N [*] V [*] T [*] Q [*] L [*] C [*] . . . 23 . . . EGE [*] α [*] Y [*] F [*] Q [*] S [*] L [*] N [*] V [*] W [*] H [*] A [*] T [*] K [*] N [*] . . . SGA [*] R [*] G [*] S [*] W [*] G [*] N [*] V [*] T [*] Q [*] L [*] C [*] . . . 41 . . . EDEY [*] Y [*] T [*] H [*] A [*] I [*] T [*] W [*] S [*] K [*] T [*] N [*] . . . SGA [*] R [*] G [*] N [*] W [*] G [*] Q [*] V [*] A [*] Q [*] L [*] C [*] . . . 66 . . . DDE [*] Y [*] L [*] H [*] T [*] I [*] K [*] V [*] W [*] S [*] E [*] A [*] K [*] S [*] E [*] K [*] . . . SGA [*] R [*] G [*] N [*] W [*] G [*] Q [*] I [*] T [*] Q [*] L [*] C [*] . . .
--------------------	--

Reassigned codon symbols	
α	CGA
Y	CGG

Figure 2: Multiple sequence alignments of undecaprenyl diphosphate synthase (BUSCO POG091H00BZ), tRNA (guanine-N1)-methyltransferase (BUSCO POG091H01WE), Peptidase M50 (BUSCO POG091H0131), and DNA-directed RNA polymerase subunit beta (BUSCO POG091H02K5) from Absconditabacteria, Gracilibacteria clades 1 and 2, and selected outgroup species. Alignment regions containing nearby CGA (α) or CGG (γ) positions are shown, with columns containing CGA or CGG in Absconditabacteria sequences or CGG in Gracilibacteria clade 1 sequences highlighted with an asterisk on the top, and columns containing CGA or CGG in Gracilibacteria clade 2 and outgroup sequences or CGA in Gracilibacteria clade 1 sequences highlighted with an asterisk below.

Reassignment of CGG to glutamine in Clostridium_U

In the Clostridium_U clade on the GTDB tree (Figure 4), three species were inferred to translate CGG as glutamine and one species did not have an inferred meaning for CGG (*Peptacetobacter hominis*, GCA_006861675.1). One non-representative genome for *Clostridium hiranonis* (species #1 on tree) was inferred to translate CGG as arginine; however, we suspect the signal may be coming from contaminating contigs (CheckM contamination estimate is 5%). All of the outgroup species were inferred to translate CGG as either arginine or the meaning was left uninferred in a few cases where CGG is very rare (<14 aligned Pfam positions).

Figure 5 shows four example alignments of conserved single-copy bacterial genes across Clostridium_U and outgroup species. The four aligned genes contain several positions encoded by CGG in Clostridium_U genomes and these residues are primarily conserved for glutamine in other species. The occurrence of CGG at conserved glutamine positions supports CGG being used as a glutamine codon in the reassigned clade.

In all Clostridium_U genomes, CGG is decoded by a tRNA_{CCG}. We did not classify this tRNA as an arginine tRNA due to lack of A20 in the D-loop. This tRNA also lacks many glutamine identity elements such as weak 1:72 base pair in the acceptor stem. Therefore, we could not determine which aminoacyl-tRNA synthetase recognizes the tRNA_{CCG} from the tRNA sequence alone. It is possible that the glutamyl-tRNA synthetase (GlnRS) still recognizes the

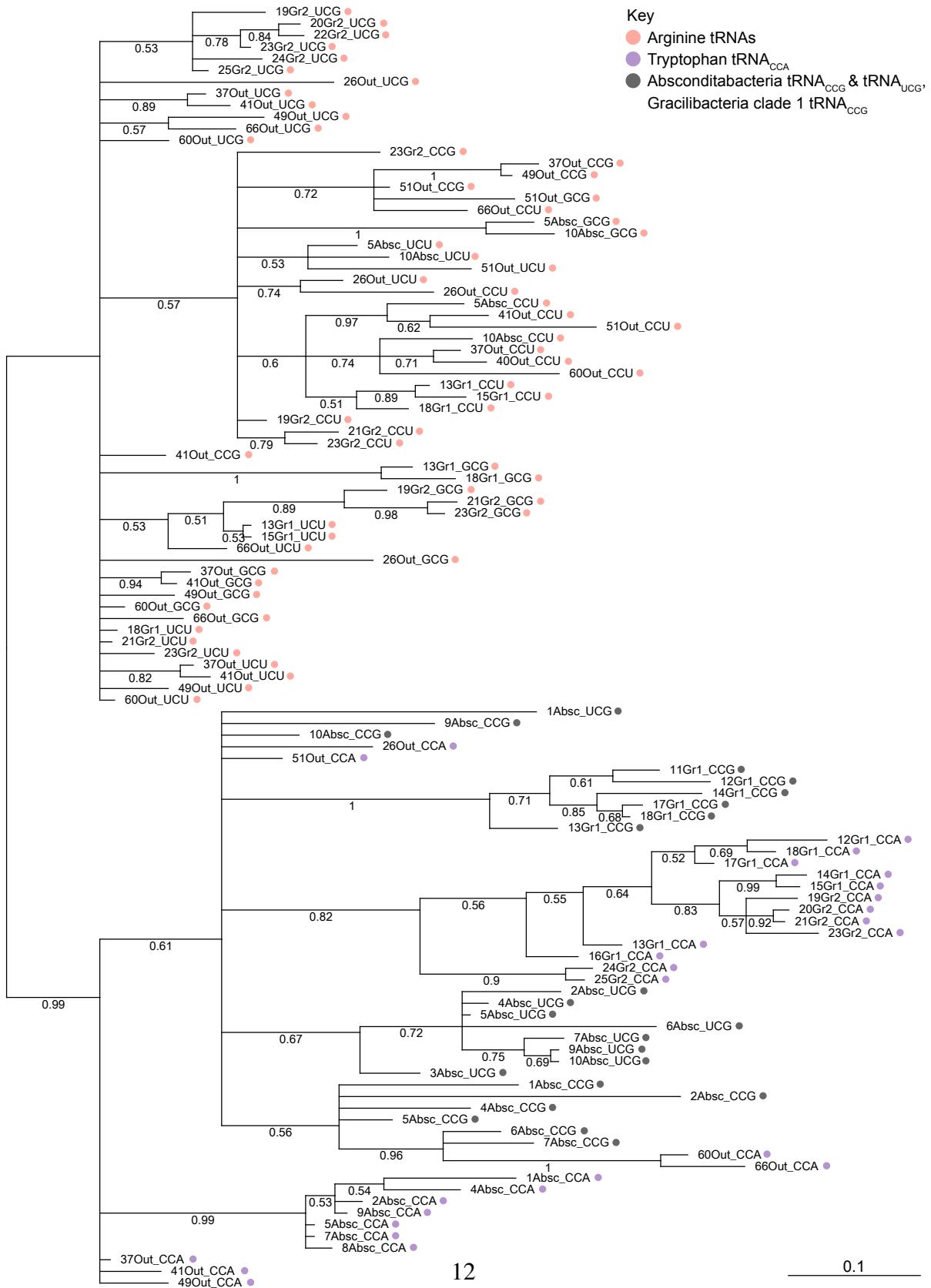


Figure 3: An unrooted phylogenetic tree of arginine, tryptophan, and reassigned tRNA sequences from Absconditabacteria, Gracilibacteria, and outgroup genomes, generated by a Bayesian approach implemented in MrBayes 3.2.7a. Values below each internal branch indicate the posterior probability of each clade. Since this is an unrooted tree, this value represents the probability of the sequences on either side of the branch form two distinct clades. Branch length scale in the bottom right is in expected substitutions per site. Sequence labels follow the format “species number - clade _ tRNA anticodon”. This tree includes sequences from tRNA_{CCG} and tRNA_{UCG} in the Absconditabacteria and tRNA_{CCG} in Gracilibacteria clade 1 (potentially involved in codon reassessments and indicated with a gray circle), the arginine tRNA_{UCU}, tRNA_{CCU}, and tRNA_{GCG} (along with tRNA_{UCG} in all Gracilibacteria and outgroups and tRNA_{CCG} in Gracilibacteria clade 2 and outgroups) indicated by a red circle, and the tryptophan tRNA_{CCA} from all groups indicated by a purple circle.

tRNA_{CCG} despite missing many of the known identity elements because the relaxed specificity of GlnRS towards non-cognate tRNAs. For example, in *E. coli*, many amber (UAG) suppressor tRNAs derived from non-glutamine tRNAs insert glutamine at UAG, even if they lack glutamine identity elements (3). The charging of tRNA_{CCG} could be experimentally confirmed in the future for culturable species such as *Clostridium hiranonis*.

In a tRNA phylogeny built from arginine and glutamine tRNA sequences from Clostridium_U species and outgroups, the reassigned tRNA_{CCG} branches outside of either the cluster of arginine or glutamine tRNAs (Figure 6). We compared the likelihood of two phylogenetic models, one where the Clostridium_U tRNA_{CCG} sequences are constrained to cluster with the glutamine tRNAs against another model where they are constrained to cluster with the arginine tRNAs. The log2 ratio of the likelihoods is 1.5, indicating more-or-less even support for both models.

In some Clostridium_U species, the tRNA genes needed to read the arginine codons CGA, CGC, and CGU (tRNA_{ACG} and/or tRNA_{UCG}) and the glutamine codon CAA (tRNA_{UUG}) cannot be found. This may be due to incomplete genome assembly because some members of Clostridium_U do have these tRNA genes, and all of the four Clostridium_U species are missing at least a few other required tRNA genes.

CGG is a rare codon in all members of Clostridium_U, with usage of 4-6 per 10,000 codons in aligned Pfam domains. The low genomic GC content in Clostridium_U species, ranging between 0.31-0.33, may have contributed to the low codon usage of CGG and the GC-rich arginine codons CGC, CGA, and AGG (each of these three codons has usage less than 22 per 10,000 codons) in favor of AGA (which is used at 314-355 per 10,000 codons). If low genomic GC content favored substitutions away from CGG prior to reassignment of CGG, this would have facilitated reassignment by minimizing the number of substitutions needed to adapt to the new translation.

In summary, the appearance of CGG at conserved glutamine positions in proteins has led to CGG being inferred as a glutamine codon in Clostridium_U species by the genetic code inference method, and is supported by examining multiple sequence alignments of conserved single-copy genes. CGG is decoded in this clade by an unusual tRNA_{CCG} whose isotype cannot be predicted by examination of the tRNA sequence. *Peptacetobacter hominis* (GCA_006861675.1, species #4 on tree) did not have an inferred CGG meaning by the genetic code inference program despite 145 Pfam positions aligned to CGG; however, the top amino acid model for CGG is glutamine (model probability of 0.984, below the threshold for reporting). This might possibly reflect a low level of assembly contamination, occurrence of CGG at weakly conserved positions, or possibly ambiguous translation of CGG as more than one amino acid.

CGG decoding in the outgroup species

In Clostridium_U, CGG is a rare codon, with codon usage of 4-6 per 10,000 codons in regions aligned to Pfam domains. In other members of the Peptostreptococcaceae (species #5-41 on tree) CGG is even rarer, occurring <2 times per 10,000 codons in all but two species. The overall low usage of CGG codons in the Peptostreptococcaceae may be tied to the low GC content of the clade, which ranges between 0.27-0.38 (Figure 4).

In the sister group to Clostridium_U within the Peptostreptococcaceae (species #5-34 on tree, genera *Clostridioides*, *Asaccharospora*, *Intestinibacter*, *Terrisporobacter*, *Paeniclostridium*, *Paraclostridium*, and *Romboutsia*), we could not find any tRNA capable of decoding CGG (Figure 4, gray box outline). We wanted to confirm that this is not due to incomplete genome assembly, so we analyzed the tRNA gene content of all 2,114 genome assemblies belonging to the 29 species in this clade (predominantly consisting of 2,014 *Clostridioides difficile* assemblies, species #8-11 on tree). None of the 2,114 genome assemblies contained a tRNA_{CCG} and only a single assembly contained a tRNA_{UCG} gene.

If a tRNA to decode CGG is truly missing in this clade, then the few CGGs in coding regions might be decoded inefficiently by the tRNA_{ACG}^{Arg} or by a non-cognate tRNA. The lack of a tRNA to decode CGG in some species, in conjunction with low GC content, may explain the extremely low usage of CGG in Peptostreptococcaceae genomes.

Reassignment of CGG to tryptophan in a clade of Bacilli

Genetic code inference

In our analysis of all sequenced bacterial genomes, 5 genomes broadly annotated as “uncultured Mollicutes” and “Bacillales sp” were predicted to translate the canonical arginine codon CGG as tryptophan. All of these genomes are included in the Genome Taxonomy Database (GTDB) phylogeny, where they belong to 4 species clusters that form a monophyletic clade within the GTDB genus g_UBA4855 (family f_CAG-826, order o_RFN20, class c_Bacilli). The clade using the predicted reassignment of AGG to methionine also belongs to the order o_RFN20, but to a different family.

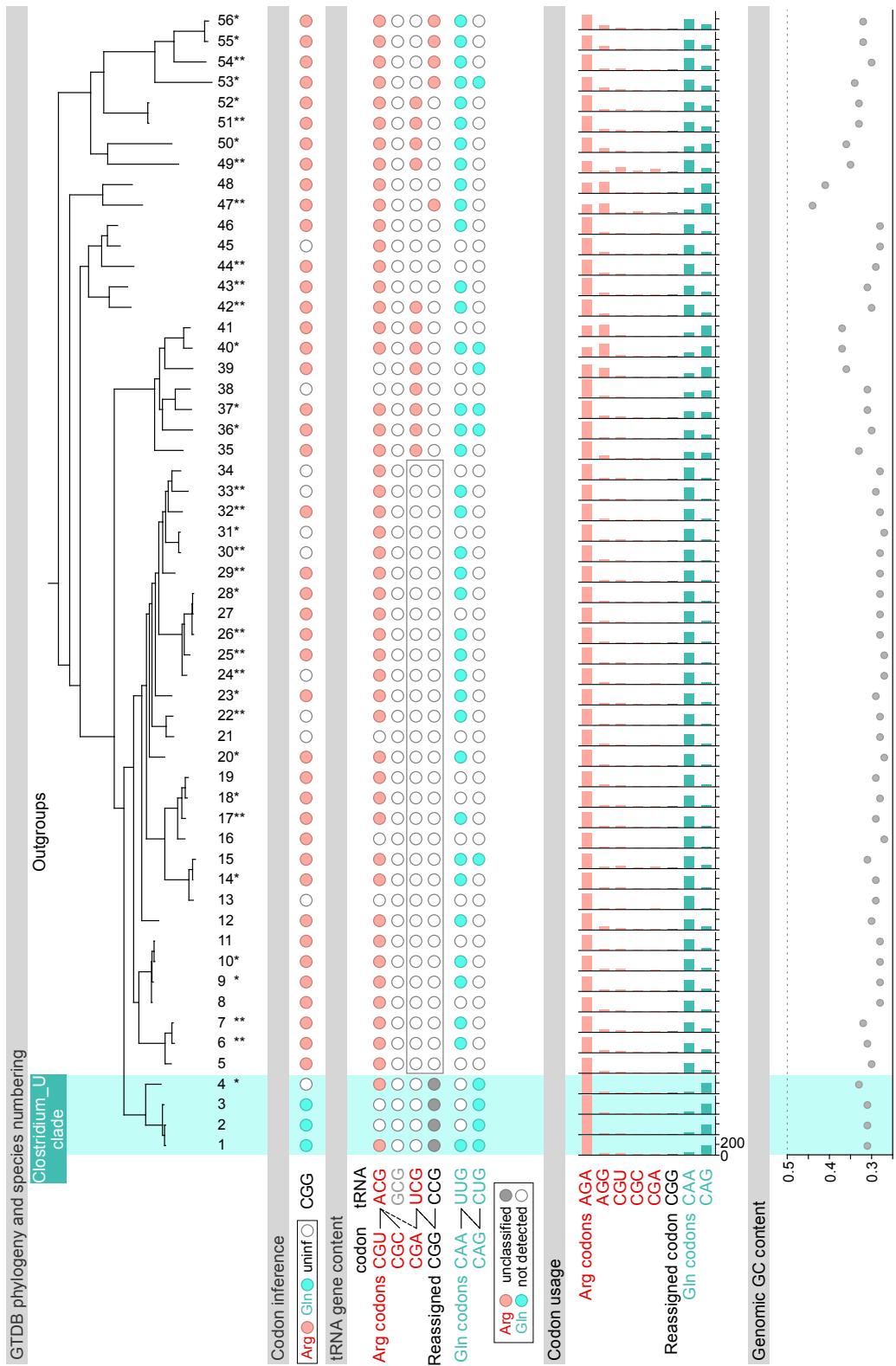


Figure 4: Phylogenetic tree from GTDB showing the Clostridium_U clade and closest outgroup genomes, each species indicated by a number which can be cross-referenced with the summary spreadsheet. Double asterisks indicate genomes the most complete genomes, which have CheckM estimated genome completeness from GTDB >98% and the entire minimal set of 22 required tRNAs (excluding CGN-decoding tRNAs), while single asterisks indicate genomes which have CheckM estimated genome completeness from GTDB >98% and >18 out of 22 required tRNAs. For each species, the inferred translation of the reassigned codon CGG by our method is indicated by colored circles (red: arginine, light blue: glutamine, white: uninferred). The presence of tRNA genes that recognize the CAR- and CGN-codons is also indicated by filled circles, colored according to the predicted amino acid charging based on identity elements for tRNAs (see Methods). Gray box outline highlights the inability to detect any CGG-decoding tRNA in the Peptostreptococcaceae (species #5-34). Anticodons in gray font are typically not found in the Firmicutes. The lines connecting codons and tRNA anticodons represent the likely wobble decoding capabilities, with dashed lines representing weaker interactions. The anticodon ACG is presumed to be modified to ICG, and UUG is presumed to be modified in a way that restricts wobble to CGA and CGG. The U34 of anticodon UCG is presumed to be modified in a way that restricts decoding to CGA and CGG, but could potentially recognize CGU and/or CGC depending on the true modification state. The remaining anticodons are not expected to be modified in a way that alters which codons are recognized. The codon usage for the reassigned codon CGG, the arginine codons AGA, AGG, CGU, CGC, and CGA, and the glutamine codons CAA and CAG is the frequency per 10,000 codons aligned to Pfam positions. Genomic GC content is calculated over the entire genome.

Transcription termination factor NusA (POG091H0124)

Clostridium_U	1 . . . KNFGQA α NVVEVFD . . . EGVMP α SEKID . . .	
	2 . . . KNFGQA α NVVEVFD . . . EGVMP α SEKID . . .	
	3 . . . KNFGQA α NVVEVFD . . . EGVMP α SEKID . . .	
	4 . . . KNFGQA α NVVEVFN . . . EGIIMT α NEKIA . . .	
Outgroup	6 . . . KNFGSA α NVRVEFD . . . EGVMT α SEQIP . . .	Reassigned codon symbols
	9 . . . KNFGSA α NVRVEFD . . . EGVMT α SEQIP . . .	α CGG
	14 . . . KNFGSA α NVRVSVD . . . EASMPENEQIP . . .	
	29 . . . KNFGSA α NVRVEFD . . . EAVMT α TEQM . . .	
	40 . . . KNFGSA α NVRIDMN . . . EGVLL α SEQIR . . .	
	44 . . . RNFGSC α NVRTEID . . . EGVLN α TEQIP . . .	
	47 . . . KNFGSS α NVRRIKID . . . EGVLNATEQIP . . .	
	52 . . . KNFGTSSNVRVEMD . . . EAVLPPSEQIP . . .	

ATP synthase F1, gamma subunit (POG091H01H6)

Clostridium_U	1 . . . IFKSVMEL α DPEKDMVV . . . RARQSSIT α EITEIAG . . .	
	2 . . . IFKSVMEL α NPKKDMVV . . . RARQSSIT α EITEIAG . . .	
	3 . . . IFKSVMEL α NLEKDMVV . . . RARQSSIT α EITEIAG . . .	
	4 . . . IFKREAESLMNKDTDMVI . . . RARQASIT α EITEIAG . . .	
Outgroup	6 . . . VLKETVNHMHDGKETVM . . . RARQSAVT α EITEIVG . . .	
	9 . . . VLKESVSHMEGKKEV . . . RARQSAVT α EITEIVG . . .	
	14 . . . IFLKLTAAHMDGKQE . . . RARQTAVT α EITEIVG . . .	
	29 . . . ILKAVLAHNSKNTAKVI . . . RARQAAVT α EISEIVA . . .	
	40 . . . AFKEAEKYMLVEDYVI . . . RARQGAVT α EITEISG . . .	
	44 . . . VLKTAINHMEHKQESII . . . RARQASIT α EISEIVA . . .	
	47 . . . SLKTAVALAHMEGKKEV . . . RARQATIT α EISEIVA . . .	
	52 . . . VLKMAVSHMDNQKYPVI . . . RARQATIT α EISEIVA . . .	

SsrA-binding protein (POG091H022D)

Clostridium_U	1 . . . LKGTEVKSI α GRVNLKEG . . .	
	2 . . . LKGTEVKSI α GRVNLKEG . . .	
	3 . . . LKGTEVKSI α GRVNLKEG . . .	
	4 . . . LKGTEVKSI α GKVNLSDG . . .	
Outgroup	6 . . . LKGTEVKSLR α GKANLSDG . . .	
	9 . . . LKGTEVKSI α GKVNLSDG . . .	
	14 . . . LKGTEVKSLRMGRVNLKD . . .	
	29 . . . LKGTEVKSI α GKLNLSDG . . .	
	40 . . . LKGTEVKSI α GRVNLKEG . . .	
	44 . . . LKGTEVKSI α GKLNLAEG . . .	
	47 . . . LKGTEVKSI α AGRINLKEG . . .	
	52 . . . LKGTEVKSI α GKVNLSGEG . . .	

16S rRNA methyltransferase GidB (POG091H00IF)

Clostridium_U	1 . . . ELLAEWN α KMNLTGIDDEKGT . . .	
	2 . . . ELLAEWN α KMNLTGIDDEKGT . . .	
	3 . . . ELLAEWN α KMNLTGIDDEKGT . . .	
	4 . . . EILVEWN α KMNLTGIEDEKEV . . .	
Outgroup	6 . . . EILVEWN α KMNLTGIEDEKEV . . .	
	9 . . . EILVDWNKKMNLTGIEDEKEV . . .	
	14 . . . EILVEWN α KMNLTGIEDEKEV . . .	
	29 . . . DMLADWN α HMNLTGIVEEKEV . . .	
	40 . . . QILVEYN α HMNLTGITEQREV . . .	
	47 . . . RLLEWNEKMNLTAITQE α EIY . . .	
	52 . . . DTLEYNKVMNLTAIEDPEEII . . .	

Figure 5: Multiple sequence alignments of transcription termination factor NusA (BUSCO POG091H0124), ATP synthase F1, gamma subunit (BUSCO POG091H01H6), SsrA-binding protein (BUSCO POG091H022D), and 16S rRNA methyltransferase GidB (BUSCO POG091H00IF) from the Clostridium_U clade and selected outgroup species. Alignment regions containing nearby CGG (α) positions are shown, with columns with CGG in Clostridium_U species sequences highlighted.

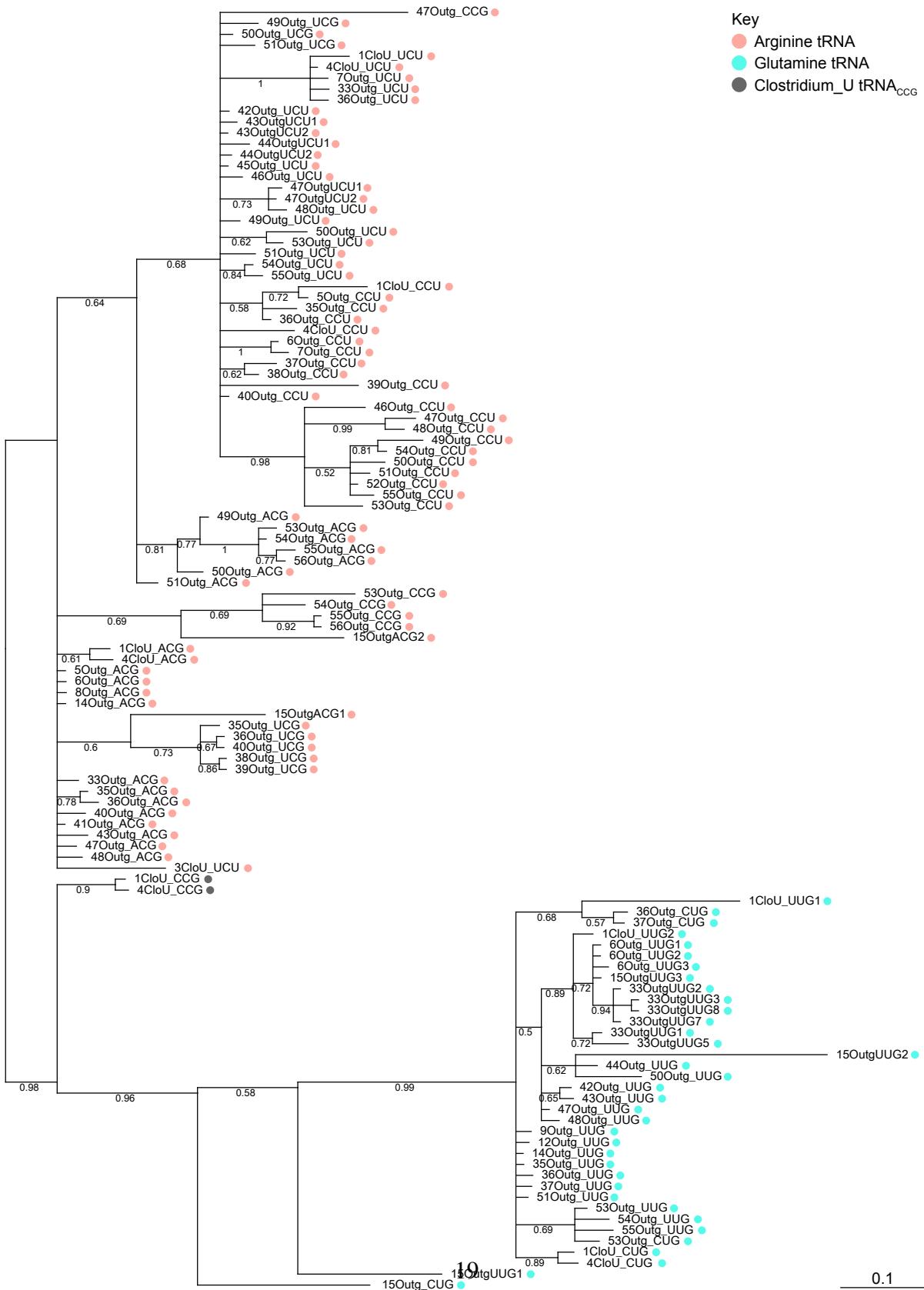


Figure 6: An unrooted phylogenetic tree of arginine, glutamine, and reassigned tRNA sequences from Clostridium_U and outgroup genomes, generated by a Bayesian approach implemented in MrBayes 3.2.7a. Values below each internal branch indicate the posterior probability of each clade. Since this is an unrooted tree, this value represents the probability of the sequences on either side of the branch form two distinct clades. Branch length scale in the bottom right is in expected substitutions per site. Sequence labels follow the format “species number - clade - tRNA anticodon”. This tree includes sequences from tRNA_{CCG} from Clostridium_U species (involved in the CGG reassignment and indicated with a gray circle), the arginine tRNA_{UCU}, tRNA_{CCU}, tRNA_{ACG}, and tRNA_{UCG} (and tRNA_{CCG} from outgroup species) indicated by a red circle, and the glutamine tRNA_{UUG} and tRNA_{CUG} from all groups indicated by a light blue circle.

Species phylogenetic tree

In Figure 7, we show a part of the GTDB phylogenetic tree that includes all members of the GTDB family f_CAG-826 which includes the reassigned clade (5 species) and outgroup species (45 species). We designated the reassigned clade as a branch of the tree containing all 4 species inferred to translate CGG as Trp plus one additional species that has CGG meaning uninferred (species s_UBA4855 sp002438785).

In order to help assess whether a tRNA gene is used by members of a clade or whether it is missing due to an incomplete genome, we marked the most complete genomes on the tree (Figure 7), which had a CheckM completeness estimate of 95% or greater and contained a minimal set of 22 tRNA genes (excluding CGN-decoding tRNAs) that are required in all bacteria.

Reassignment of CGG to tryptophan in the reassigned clade

In the reassigned clade (Figure 7, species #1-5), four species were inferred to translate CGG as tryptophan and one species did not have an inferred meaning for CGG (species s_UBA4855 sp002438785) due to one or two aligned Pfam positions at CGG codons. All of the outgroup species were inferred to translate CGG as either arginine or the meaning was left uninferred in a few cases where CGG is very rare (<10 aligned Pfam positions).

Figure 8 shows five example alignments of conserved single-copy bacterial genes across the reassigned clade and outgroup species. The five aligned genes show individual CGG positions in reassigned species #1,3,4 at residues that are primarily conserved for tryptophan in other members of the reassigned clade and outgroups. This suggests that CGG can be used interchangeably with the tryptophan codon UGG at conserved tryptophan residues in these species. These alignments do not help determine how CGG is translated in the uninferrered species (s_UBA4855 sp002438785, #5 on tree), as we could not find instances of CGG in aligned BUSCO genes. For the nearest outgroup species (s_UBA4855 sp002451465, #6 on tree), in the CTP synthase alignment, there is a single CGG position at a conserved arginine residue, supporting canonical arginine translation of CGG in this species.

In the reassigned clade, the four species that were inferred to translate CGG as tryptophan contain CGG-decoding tRNA_{CCG} gene. This tRNA does not appear to be an arginine tRNA due to lack of A20 in the D-loop, and contains elements supportive of tryptophan identity such as a G73 discriminator base and A/G1:U72 in the acceptor stem. The uninferrered member of the reassigned clade (s_UBA4855 sp002438785, #5 on tree) does not contain a CGG-decoding tRNA (either tRNA_{CCG} or tRNA_{UCG}) in any of the 3 genomes assigned to the species, indicating either an incomplete genome assembly, tRNA gene detection failure, or a bona fide inability to translate CGG.

In tRNA phylogenetic tree built from arginine and tryptophan tRNAs from the reassigned clade and outgroup species, the tRNA_{CCG} sequences from the reassigned clade do not cluster with the outgroup tRNA_{CCG} genes among the other arginine tRNAs, but instead branch right outside of the tryptophan tRNA_{CCA} sequences (Figure 9). We compared the likelihood of two phylogenetic models, one where the reassigned tRNA_{CCG} sequences are constrained to cluster with the tryptophan tRNAs against another model where they are constrained to cluster with the arginine tRNAs. The log2 ratio of the likelihoods is 4.5, in support of the tryptophan model and

consistent with the presence of tryptophan identity elements.

None of the reassigned species have a tRNA_{UCG} and instead presumably decode the arginine codons CGU, CGC, and CGA with a tRNA_{ACG}^{Arg} (A presumably modified to I). This is in contrast to many of the outgroup species (focusing on the more complete genomes as representatives), which appear to rely on a tRNA_{ACG}^{Arg} gene to decode CGU, CGC, and CGA and a tRNA_{UCG}^{Arg} to decode CGA and CGG, with an optional CGG-decoding tRNA_{CCG}^{Arg} present in some species.

CGG is a very rare codon in all members of the GTDB genus g_-UBA4855 (reassigned species #1-5 and outgroup species #6 on tree), with usage no greater than 2.5 per 10,000 codons in aligned Pfam domains. The very low genomic GC content in the clade, ranging between 0.26-0.30, may have contributed to the low codon usage of CGG and the GC-rich arginine codons CGC, CGA, and AGG (for each, usage <31 per 10,000) in favor of AGA (which is used at 220-292 per 10,000 codons). If low genomic GC content favored substitutions away from CGG prior to reassignment of CGG, this would have facilitated reassignment by minimizing the number of substitutions needed to adapt to the new translation.

In summary, the appearance of CGG at conserved tryptophan positions in proteins has led to CGG being inferred as a tryptophan codon in reassigned species by the genetic code inference method, and is supported by examining multiple sequence alignments of conserved single-copy genes. CGG is decoded in this clade by a tRNA_{CCG} that is consistent with a tryptophan isotype. s_-UBA4855 sp002438785 (species #5 on tree) did not have an inferred CGG meaning by the genetic code inference program and does not have a tRNA_{CCG}, but is included in the reassigned clade because it is placed in the same clade as the reassigned species according to the GTDB tree. The entire reassigned clade plus the closest outgroup species s_-UBA4855 sp002451465 (species #6 on tree) form the GTDB genus g_-UBA4855, which is characterized by low genomic GC content and extremely low codon usage of CGG.

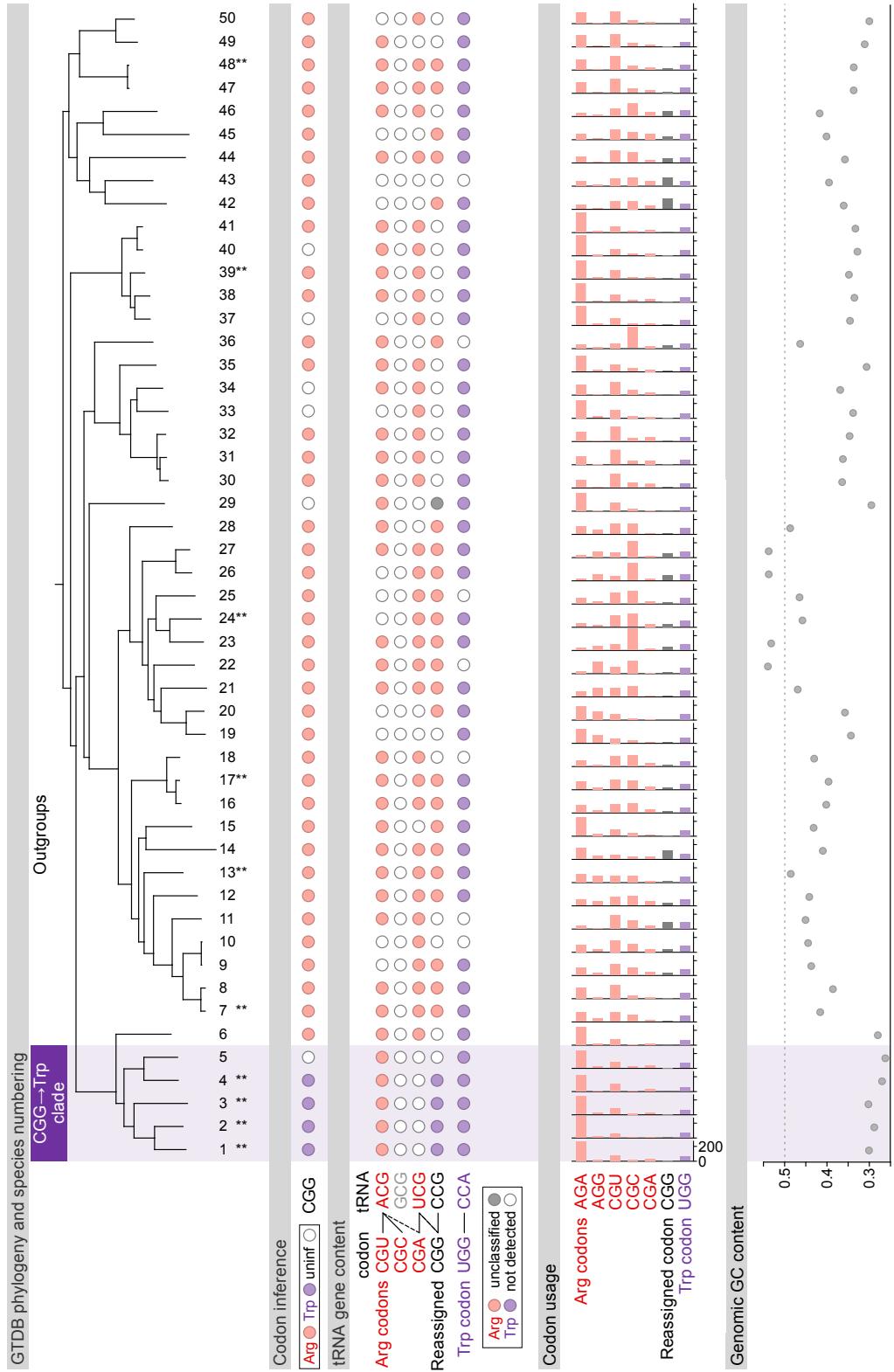


Figure 7: Phylogenetic tree from GTDB showing the Bacilli CGG→Trp clade and closest out-group genomes, each species indicated by a number which can be cross-referenced with the summary spreadsheet. Double asterisks indicate genomes the most complete genomes, which have CheckM estimated genome completeness from GTDB >95% and the entire minimal set of 22 required tRNAs (excluding CGN-decoding tRNAs). For each species, the inferred translation of the reassigned codon CGG by our method is indicated by colored circles (red: arginine, purple: tryptophan, white: uninferred). The presence of tRNA genes that recognize the CCA and CGN-codons is also indicated by filled circles, colored according to the predicted amino acid charging based on identity elements for tRNAs (see Methods). Anticodons in gray font are typically not found in the Firmicutes. The lines connecting codons and tRNA anticodons represent the likely wobble decoding capabilities, with dashed lines representing weaker interactions. The anticodon ACG is presumed to be modified to ICG. The U34 of anticodon UCG is presumed to be modified in a way that restricts decoding to CGA and CGG, but could potentially recognize CGU and/or CGC depending on the true modification state. The remaining anticodons are not expected to be modified in a way that alters which codons are recognized. The codon usage for the reassigned codon CGG, the arginine codons AGA, AGG, CGU, CGC, and CGA, and the tryptophan codon CCA is the frequency per 10,000 codons aligned to Pfam positions. Genomic GC content is calculated over the entire genome.

Reassignment of CGG to tryptophan in *Anaerococcus*

Genetic code inference

In our analysis of all sequenced bacterial genomes, four genomes annotated as belonging to the genus *Anaerococcus* were predicted to translate the canonical arginine codon CGG as tryptophan. Three of these genomes are included in the Genome Taxonomy Database (GTDB) phylogeny, where they belong to two species clusters within the genus g__*Anaerococcus* (family f__Helcoccaceae, order o__Tissierellales, class c__Clostridia) which is comprised of 23 species clusters in total.

Species phylogenetic tree

In Figure 10A, we show a part of the GTDB phylogenetic tree that includes all members of the GTDB family f__Helcoccaceae, including the genus *Anaerococcus* (23 species) and outgroup species (22 species). In order to help assess whether a tRNA gene is used by members of a clade

Adenylosuccinate synthetase (POG091H01G9)

Bacilli CGG→Trp clade	1 . . . VVVGSQ WG DEGK . . . YVTLPG W KEDISK . . . 2 . . . VVLGSQ WG DEGK . . . YKTFKG W TEDISK . . . 3 . . . VIQGTQ α GDEGK . . . YKDFCG α DEDISN . . . 4 . . . LVLGSQ WG DEGK . . . YKTFKG α DEDISK . . . 5 . . . VIEGSQ WG DEGK . . . YKEFKKFSFS-DK . . . 6 . . . LVIGAQ WG DEGK . . . YKTFAS W KEDISQ . . . 7 . . . AIQGMQ WG DEGK . . . YIELPS W KEDISS . . . 13 . . . AIQGSQ WG DEGK . . . YAIFKS W KEDISG . . . 15 . . . AIEGMQ WG DEGK . . . YISLPS W KEDISH . . . 28 . . . AIEGMQ WG DEGK . . . YKTLPG W KEDISN . . . 43 . . . AIVGVN WG DEGK . . . YEYLPGFNFEDISK . . . 48 . . . VLEGSQ WG DEGK . . . YITMPT W KEDITH . . .	Reassigned codon symbols α CGG
Outgroup		

CTP synthase (POG091H02IX)

Bacilli CGG→Trp clade	1 . . . DMSD WVK LIE . . . GFGK R GVEGK . . . 2 . . . NMDD WIDL IS . . . GFGN R GIEGK . . . 3 . . . NMDD αIKL ID . . . GFGN R GIEGK . . . 4 . . . EMSD WNELIR . . . GFGN R GVEGK . . . 5 . . . DMDD WQALIK . . . GFGT R GTEGK . . . 6 . . . DMDD WRHFVH . . . GFGK α GIDGM . . . 8 . . . DLHE WRSWCD . . . GFGE R GSKGK . . . 13 . . . DLRN WQKWCD . . . GFGE R GSEGK . . . 17 . . . DLHN WEKWVD . . . GFGE R GTEGK . . . 21 . . . QLVDYEEFVS . . . GFGR R GTDGM . . . 35 . . . DLSDFK QOLIK . . . GFGK R GIEGK . . . 45 . . . TIEFW QOLIQ . . . GFG α ATEGK . . .	
Outgroup		

Protein translocase subunit SecA (POG091H01RS)

Bacilli CGG→Trp clade	1 . . . YQIKRRE α DKETADQ . . . 2 . . . YLARKKE W DKDVASQ . . . 3 . . . YNSLKKNW P KEEIDK . . . 5 . . . YSQQKK W EPEIAEK . . . 6 . . . YLRRKK W DKEFAEK . . . 9 . . . YLEKRKT W PAADADK . . . 17 . . . YVNKRKE W PNEVADQ . . . 21 . . . YLKKRK W PKEVADR . . . 29 . . . YLVRRKD W -KELADQ . . . 36 . . . YVN α RKE W GEEIANQ . . . 38 . . . YVERRKD W PEELQGG . . . 48 . . . YLERRKE W GDEVAEN . . .	
Outgroup		

DNA ligase (POG091H024G)

Bacilli CGG→Trp clade	1 . . . KIMEMDG WSNK SVDK . . . 2 . . . EIVEIDG WSHK SIDK . . . 3 . . . DIINSDG WSYR STDN . . . 4 . . . NIIQIEG α SYKSTES . . . 6 . . . ELMNIDG WSIK SVTN . . . 10 . . . EIKNIDG WSDK SMNS . . . 13 . . . MIKELDG WSDK SINS . . . 21 . . . EIKALDG WSDK SISS . . . 29 . . . QILNLEG WSHK SFNN . . . 34 . . . EIINIEG WSTK SIDN . . . 38 . . . EIIALDG WEKS SIDN . . . 48 . . . DLLMIDG FSDK SVDK . . .	
Outgroup		

Alanyl-tRNA synthetase (POG091H01PM)

Bacilli CGG→Trp clade	1 . . . FFDRGEK WDP KHLGV . . . 2 . . . FFDRGEK W DKDNIGVD . . . 3 . . . FFDRGEK α DPEHIGIK . . . 4 . . . FYDRGEKYDPNHLGV . . . 5 . . . FYDRGEK W DKDHLGVK . . . 6 . . . FFDRGEKYDPNHIGID . . . 7 . . . FYDRGE W DPKHLGIK . . . 13 . . . FFDRGEK W DPKHLGVK . . . 21 . . . FFDRGEKYDPKHLGVK . . . 34 . . . FFDRGEKYDEKHLGIK . . . 40 . . . HFDRGEKFDPHEVGVK . . . 48 . . . FFDRGEKYDPDHLGIR . . .	
Outgroup		

Figure 8: Multiple sequence alignments of adenylosuccinate synthetase (BUSCO POG091H01G9), CTP synthase (BUSCO POG091H02IX), protein translocase subunit SecA (BUSCO POG091H01RS), DNA ligase (BUSCO POG091H024G), and alanyl-tRNA synthetase (BUSCO POG091H01PM) from the Bacilli CGG→Trp clade and selected outgroup species. Alignment regions containing CGG (α) at conserved positions are shown, with columns with CGG in Bacilli CGG→Trp clade and the closest outgroup (species #6) sequences highlighted.

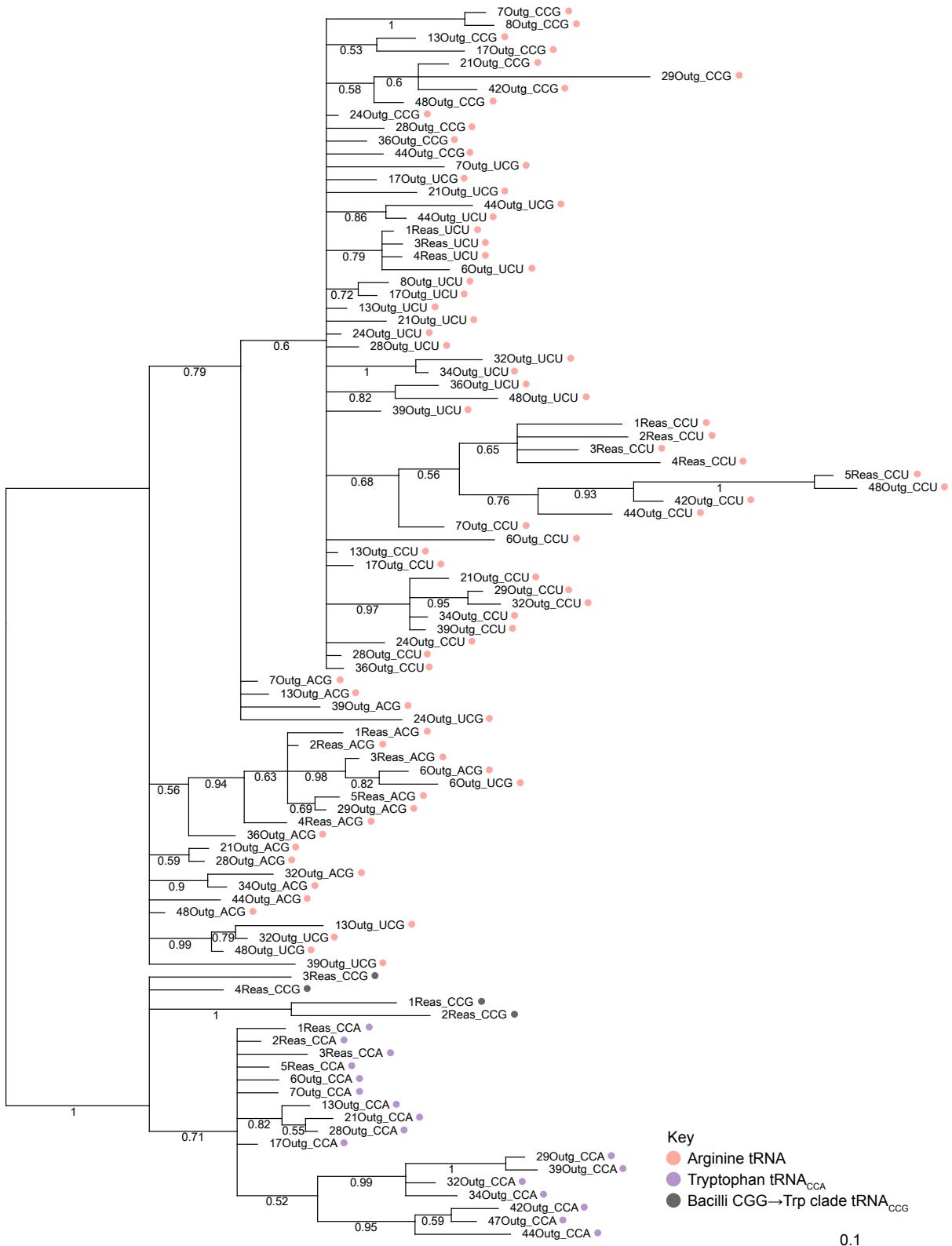


Figure 9: An unrooted phylogenetic tree of arginine, tryptophan, and reassigned tRNA sequences from Bacilli CGG→Trp clade and outgroup genomes, generated by a Bayesian approach implemented in MrBayes 3.2.7a. Values below each internal branch indicate the posterior probability of each clade. Since this is an unrooted tree, this value represents the probability of the sequences on either side of the branch form two distinct clades. Branch length scale in the bottom right is in expected substitutions per site. Sequence labels follow the format “species number - clade _ tRNA anticodon”. This tree includes sequences from tRNA_{CCG} from Bacilli CGG→Trp clade (involved in the CGG reassignment and indicated with a gray circle), the arginine tRNA_{UCU}, tRNA_{CCU}, tRNA_{ACG}, and tRNA_{UCG} (and tRNA_{CCG} from outgroup species) indicated by a red circle, and the tryptophan tRNA_{CCA} from all groups indicated by a purple circle.

or whether it missing due to an incomplete genome, we marked the most complete genomes on the tree (Figure 10A), based on a CheckM completeness estimate of 98% or greater and the presence of a minimal set of 22 tRNA genes (excluding CGN-decoding tRNAs) that are required in all bacteria.

Reassignment of CGG to tryptophan in the reassigned clade

In *Anaerococcus* (Figure 10A, species #1-23), two species were inferred to translate CGG as tryptophan, two species were inferred to translate CGG was arginine, and the remaining 19 species did not have an inferred meaning for CGG due to fewer than 18 aligned Pfam positions at CGG codons. We refer to all of *Anaerococcus* as the reassigned clade based on the presence of the tryptophan-type tRNA_{CCG} (see below). All of the outgroup species were inferred to translate CGG as either arginine or the meaning was left uninferred in species where CGG is very rare (<9 aligned Pfam positions).

We constructed multiple sequence alignments of conserved single-copy bacterial genes from the BUSCO dataset across *Anaerococcus* species and the outgroup species to find examples of CGG occurring at conserved positions for some amino acid. Unfortunately, CGG is so rare in *Anaerococcus* that only a single CGG position could be found in a well-aligned region of a

BUSCO gene alignment. In *Anaerococcus nagyae* (#14 on tree, GCA_003433955.1) a single CGG occurs at position conserved for tryptophan in other *Anaerococcus* and outgroup species in an alignment of DNA ligase (Figure 10B).

In *Anaerococcus*, all 23 species contain CGG-decoding tRNA_{CCG} which does not appear to be an arginine tRNA due to lack of A20 in the D-loop, and contains elements supportive of tryptophan identity such as a G73 discriminator base and G1:U72 in the acceptor stem. None of the *Anaerococcus* genomes have a tRNA_{UCG} and instead presumably decode the arginine codons CGU, CGC, and CGA with a tRNA_{ACG}^{Arg} (A presumably modified to I).

To further explore the origin of the tRNA_{CCG} in *Anaerococcus*, we built a tRNA phylogenetic tree from sequences of arginine and tryptophan tRNAs from *Anaerococcus* and outgroup species (Figure 11). In this phylogenetic tree, the *Anaerococcus* tRNA_{CCG} sequences form a clade that is not placed within the cluster of arginine tRNAs, but instead falls within a three-way multifurcation at the base of the cluster of tryptophan tRNAs. We compared the likelihood of two phylogenetic models, one where the *Anaerococcus* tRNA_{CCG} sequences are constrained to cluster with the tryptophan tRNA_{CCA} sequences against another model where they are constrained to cluster with the arginine tRNAs. The log2 ratio of the likelihoods is 2.3, favoring the tryptophan model. This indicates that the *Anaerococcus* tRNA_{CCG} genes are more similar in sequences to tRNA_{CCA}^{Trp} than arginine tRNAs.

CGG is a very rare codon in all *Anaerococcus* species, with usage no greater than 1.5 per 10,000 codons in aligned Pfam domains. The very low genomic GC content in the clade, ranging between 0.28-0.36, may have contributed to the low codon usage of CGG and the GC-rich arginine codons CGC, CGA, and AGG in favor of AGA (which is used at 196-278 per 10,000 codons). This would have facilitated the reassignment of CGG by minimizing the number of substitutions needed to adapt to the new translation. This codon usage pattern is repeated in other low GC content genomes in the outgroup, while outgroup genomes with higher GC

content (>0.40) favor other arginine codons such as AGG, CGU, CGG, or CGG.

In summary, two *Anaerococcus* species were predicted to translate CGG as tryptophan by the genetic code inference method, while most species did not have an inferred amino acid for CGG and two species were predicted to translate CGG as arginine. Overall, CGG is very rare in this genus (possibly due to the low genomic GC content), and many species had CGG uninferred because there were too few aligned Pfam positions to make a confident amino acid prediction. Despite the sparse distribution of CGG translation as tryptophan, all *Anaerococcus* species encode a tRNA_{CCG} gene that is consistent with a tryptophan isotype.

CGG translation in the outgroup

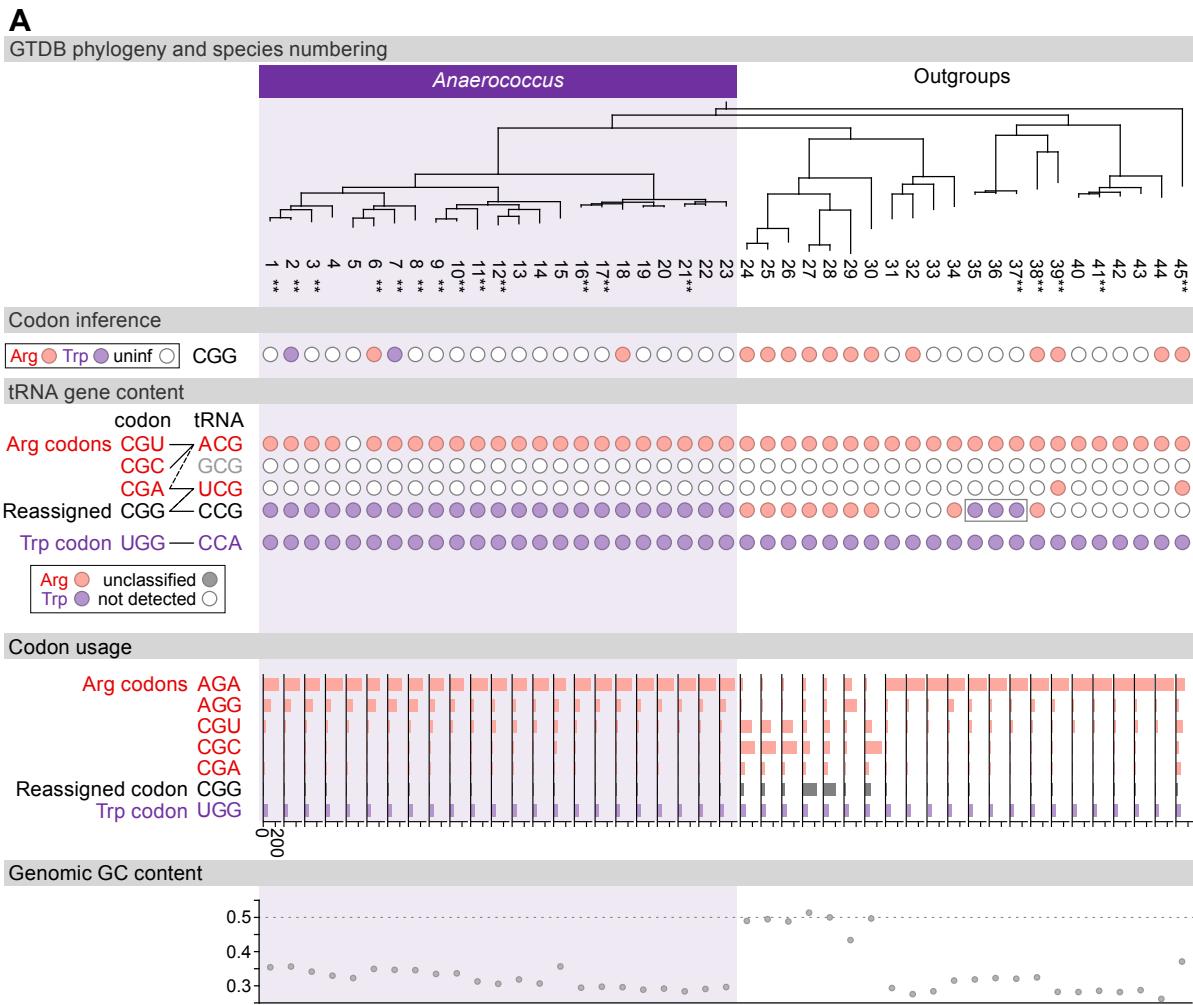
The outgroup consists of all remaining 22 species in the family f_Helcoccaceae. All species were predicted to translate CGG as either arginine or no amino acid meaning was predicted. All outgroup species presumably use a tRNA_{ACG}^{Arg} (A modified to I) to decode the arginine codons CGU, CGC, and CGA. In some species, we could not find a CGG-decoding tRNA, which may be due to incomplete genome assembly, tRNA gene detection failure, or a true lack of CGG-decoding capability.

Outgroup species in the genus *Finegoldia* (species #35-37 on tree) have a CGG-decoding tRNA_{CCG} that lacks the arginine identity element A20 and contains elements consistent with tryptophan decoding including a G73 discriminator base and A1:U72 in the acceptor stem (Figure 11, gray box outline). On the tRNA phylogenetic tree, *Finegoldia* tRNA_{CCG} sequences form a clade within other tryptophan tRNAs (Figure 11). We compared the likelihood of two phylogenetic models, one where the *Finegoldia* tRNA_{CCG} sequences are constrained to cluster with the tryptophan tRNAs and *Anaerococcus* tRNA_{CCG} against a model where the *Finegoldia* tRNA_{CCG} sequences are constrained to cluster within arginine tRNAs. The log₂ ratio of the likelihoods was 24.3, indicated strong support for the tryptophan model. This suggests that

the *Finegoldia* tRNA_{CCG} is more similar in sequence to and may have been derived from a tryptophan tRNA. CGG is very rare in this genus (fewer than 0.5 per 10,000 codons in Pfam alignments) and there were insufficient aligned Pfam positions for the genetic code inference method to predict an amino acid inference, and no CGGs could be found in well-aligned regions of BUSCO alignments. It is possible that this represents an early stage in a reassignment of CGG to tryptophan. It is unclear whether the *Finegoldia* and *Anaerococcus* tRNA_{CCG} genes are related or independently evolved.

References

1. C. Rinke, *et al.*, *Nature* **499**, 431 (2013).
2. T. T. T. Tran, H. Belahbib, V. Bonnefoy, E. Talla, *Genome Biology and Evolution* **8**, 282 (2015).
3. J. Normanly, L. G. Kleina, J. M. Masson, J. Abelson, J. H. Miller, *Journal of Molecular Biology* **213**, 719 (1990).



B DNA ligase (POG091H024G)

Anaerococcus	1 ... KFEAEYEYTTLREVVWNVGRSGKVTPSAILDP... 3 ... KFEAEYEYTTLRKVWWNVGRTGKVTPSAILDP... 4 ... KYEAEETTTLKEVWWNVGRTGKVTPSAILEP... 5 ... KFEPEEFTTKLIDVWWNVGRTGKVTPSAILEP... 11 ... KFEAEYEYTTLLDVWWNVGRTGKVTPSAILEP... 14 ... KFEAEYEYTTLLDVWWNVGRTGKVTPSALLEP... 15 ... KFEAEYEYTTLLDVWWNVGRTGKVTPSAILEP... 21 ... KYDPEEYTTKLLDVWWNVGRTGKVTPSAILEP... 27 ... KFEAEEVTTIL α SWEWNVG α TGKVTPSAILDP... 30 ... KFEAEEVTTILQAVEWNVGRTGKVTPTAQLDP... 39 ... KFEPEEVTTILKEVIWNVGRTGKVPTAILEP... 41 ... KFEAEYEYSTILKEVWWNVGRTGKVTPTAILEP...
Outgroup	

Reassigned codon symbols

α	CGG
----------	-----

Figure 10: (A) Phylogenetic tree from GTDB showing *Anaerococcus* and closest outgroup genomes, each species indicated by a number which can be cross-referenced with the summary spreadsheet. Double asterisks indicate genomes the most complete genomes, which have CheckM estimated genome completeness from GTDB >98% and the entire minimal set of 22 required tRNAs (excluding CGN-decoding tRNAs). For each species, the inferred translation of the reassigned codon CGG by our method is indicated by colored circles (red: arginine, purple: tryptophan, white: uninferred). The presence of tRNA genes that recognize the CCA and CGN-codons is also indicated by filled circles, colored according to the predicted amino acid charging based on identity elements for tRNAs (see Methods). Gray box outline highlights the tRNA_{CCG} in *Finegoldia*. Anticodons in gray font are typically not found in the Firmicutes. The lines connecting codons and tRNA anticodons represent the likely wobble decoding capabilities, with dashed lines representing weaker interactions. The anticodon ACG is presumed to be modified to ICG. The U34 of anticodon UCG is presumed to be modified in a way that restricts decoding to CGA and CGG, but could potentially recognize CGU and/or CGC depending on the true modification state. The remaining anticodons are not expected to be modified in a way that alters which codons are recognized. The codon usage for the reassigned codon CGG, the arginine codons AGA, AGG, CGU, CGC, and CGA, and the tryptophan codon CCA is the frequency per 10,000 codons aligned to Pfam positions. Genomic GC content is calculated over the entire genome. (B) Multiple sequence alignment of DNA ligase (BUSCO POG091H024G) from the *Anaerococcus* species and selected outgroup species. An alignment region containing CGG (α) at conserved positions is shown, with a column with CGG in a single *Anaerococcus* species highlighted.

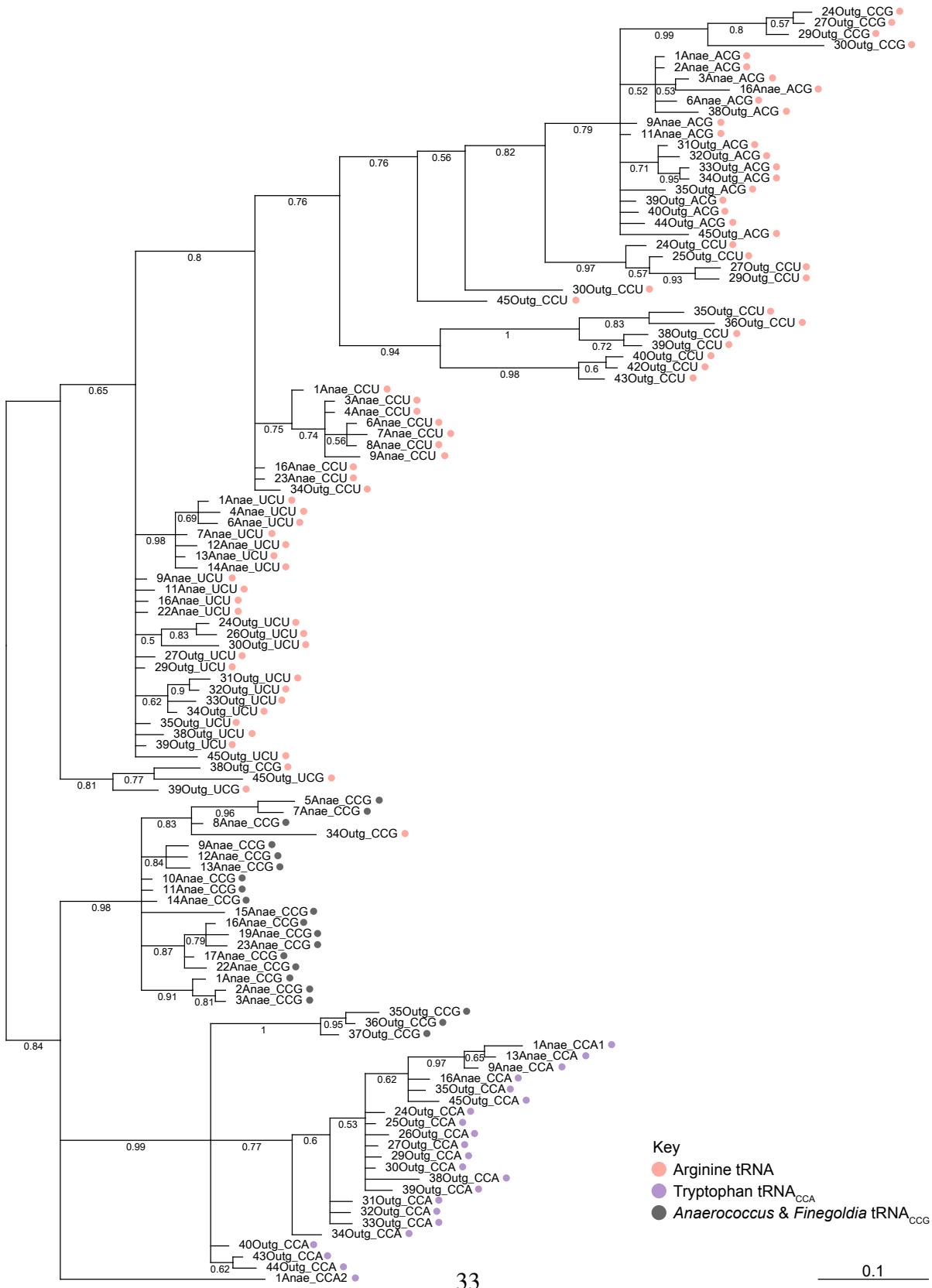


Figure 11: An unrooted phylogenetic tree of arginine, tryptophan, and reassigned tRNA sequences from *Anaerococcus* and outgroup genomes, generated by a Bayesian approach implemented in MrBayes 3.2.7a. Values below each internal branch indicate the posterior probability of each clade. Since this is an unrooted tree, this value represents the probability of the sequences on either side of the branch form two distinct clades. Branch length scale in the bottom right is in expected substitutions per site. Sequence labels follow the format “species number - clade _ tRNA anticodon”. This tree includes sequences from tRNA_{CCG} from *Anaerococcus* and *Finegoldia* (potentially involved in the CGG reassignment and indicated with a gray circle), the arginine tRNA_{UCU}, tRNA_{CCU}, tRNA_{ACG}, and tRNA_{UCG} (and tRNA_{CCG} from outgroup species) indicated by a red circle, and the tryptophan tRNA_{CCA} from all groups indicated by a purple circle.