

1
2
3

1. Extended Data

Figure #	Figure title One sentence only	Filename This should be the name the file is saved as when it is uploaded to our system. Please include the file extension. i.e.: <i>Smith_ED_Fig1.jpg</i>	Figure Legend If you are citing a reference for the first time in these legends, please include all new references in the main text Methods References section, and carry on the numbering from the main References section of the paper. If your paper does not have a Methods section, include all new references at the end of the main Reference list.
Extended Data Fig. 1	Identification of recoded phages in the Giant Tortoise microbiome.	Extended_Data_1.jpg	A. Dereplicated complete or near complete ($\geq 90\%$) phage genomes from the Giant Tortoise gut microbiome. Phages are plotted by size and coding density (CD) in standard code (Code11) B. Replotting of phage genomes from panel A, but with coding density of the alternatively coded phage calculated with the predicted genetic code instead of standard code. In all plots, phages that have recoded the TGA stop codon are indicated in green, and phages that have recoded the TAG stop codon are indicated in orange.
Extended Data Fig. 2	True coding density of standard and alternatively coded phages.	Extended_Data_2.jpg	A-F. Replotting of phage genomes from Figure 1 in the main text, but this time the coding density of alternatively coded phage was calculated with their predicted genetic code, not standard code. In all plots, symbol color represents genetic code (TGA recoding = green, TAG recoding = orange, standard code = grey).
Extended Data Fig. 3	Evolution of alternative coding.	Extended_Data_3.jpg	A. Global alignment of an 80 kilobase (kb) partial TGA recoded Agate genome (Cattle_ERR2019405_scaffold_1063) and a close standard code relative (pig_ID_3053_F60_scaffold_12). Homologous collinear sequences are shown with colored blocks (red and green here), where color corresponds to nucleotide alignment between the two genomes and lack

			of color represents lack of alignment. Genome structure for each phage is shown under the alignment graph, with DNA replication machinery represented as yellow bars and structural and lysis genes with pink bars. TGA stop codons have predominantly arisen in structural and lysis genes (individual recoded genes below in green).
Extended Data Fig. 4	Genomic maps of Jade, Sapphire, Agate and Topaz phages.	Extended_Data_4.jpg	A-D. TGA recoded genomes (A) contain genes with in-frame TGA codons (green) while TAG recoded genomes (B-D) have genes with in-frame TAG codons (orange). Suppressor tRNAs (tRNA TGA or tRNA TAG, red) are predicted to suppress translation termination at TGA and TAG stop codons, respectively. Regions of the genome encoding structural and lysis genes (pink) coincide with high use of alternative code. Contrastingly, genes involved in DNA replication (yellow) are variably encoded in alternative code. Genomes with a GC skew patterns indicative of bidirectional replication and clear origins and termini (C) have unique replichores marked in alternating shades of blue. Genomes with GC skew patterns most consistent with unidirectional replication (A-B,D) have no replication-related annotation. In some cases, unique or interesting genes have been noted with text. Clade representatives: Jade = JS_HF2_S141_scaffold_159238, Sapphire= SRR1747018_scaffold_13, Agate = Cattle_ERR2019359_scaffold_1067472, Topaz = pig_ID_1851_F40_2_B1_scaffold_1589
Extended Data Fig. 5	Genomic maps of Lak, Garnet, and Amethyst phages.	Extended_Data_5.jpg	A-C. TAG recoded genomes have genes with in-frame TAG codons (orange). Suppressor tRNAs (tRNA TAG, red) are predicted to suppress translation termination at TAG stop codons. Regions of the genome encoding structural and lysis genes (pink) coincide with high use of alternative code. In Lak phage (A), genes involved in DNA replication (yellow) are mostly encoded in alternative code. Origins and termini are unmarked in these genomes as we were unable to define clear replichores for Lak (A) and Garnet and Amethyst (B-C) appear to utilize unidirectional

			genome replication based on GC skew patterns. In some cases, unique or interesting genes have been noted with text. Clade representatives: Lak = C1--CH_A02_001D1_final, Garnet = pig_ID_3640_F65_scaffold_1252, Amethyst = pig_EL5596_F5_scaffold_275.
Extended Data Fig. 6	Code change machinery in two TGA-recoded Jade phages.	Extended_Data_6.jpg	A. An operon implicated in changing the genetic code from standard code (TGA = Stop) to code 4 (TGA = W) is directly upstream of the lysis cassette. The code change genes themselves are encoded in standard code, while some genes in the lysis cassette have in frame TGA codons (green). TrpRS = Tryptophanyl tRNA synthetase, RF1 = Release Factor 1, TM-domain = Transmembrane domain.
Extended Data Fig. 7	Read mapping to Garnet and Topaz lysogens.	Extended_Data_7.jpg	A. Reads were mapped against a manually curated Garnet lysogen. Read coverage for the <i>Prevotella</i> DNA is ~2x higher than the read coverage of the Garnet prophage, indicating that the bacterial population in this sample is incompletely lysogenized. Supporting this conclusion are paired reads that span the prophage (not shown), as well as some individual reads which show imperfect mapping to the lysogen consensus sequence (marked with asterisk), which represent the contiguous bacterial sequence. Identical sequence blocks are indicated with color. B. Reads were mapped against a manually curated Topaz lysogen. Read coverage for the integrated Topaz phage genome is ~50x higher than the neighboring <i>Oscillospiraceae</i> sequence. This indicates that the phage is actively replicating in this sample. Supporting this conclusion are paired reads that span the length of the prophage (not shown), as well as individual reads which show imperfect mapping to the lysogen consensus sequence at the 5' end of the prophage (light blue) and the 3' end of the prophage (dark blue). The reads correspond to circularized sequences. Identical sequence blocks are indicated with color.

Extended Data Fig. 8	Alignments of free and integrated phage genomes.	Extended_Data_8.jpg	<p>A. A 25kb circular TAG-recoded Garnet phage aligned to a prophage integrated in a <i>Prevotella</i> genome (<i>Prevotella</i> genes = brown). The prophage boundaries are marked by the phage integrase (pink) and the host tRNA Met. B. A 24kb circular TAG-recoded Topaz phage aligned to a prophage integrated into a <i>Oscillospiraceae</i> genome (<i>Oscillospiraceae</i> genes = blue). The prophage boundaries are marked by the phage integrase (pink) and the host tRNA Thr.</p>
----------------------	--	---------------------	--

4

2. Supplementary Information:

5

6

A. Flat Files

7

8

9

Item	Present?	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. The extension must be .pdf	A brief, numerical description of file contents. i.e.: <i>Supplementary Figures 1-4, Supplementary Discussion, and Supplementary Tables 1-4.</i>
Supplementary Information	Yes	SI_combined.pdf	<i>Supplementary Figures 1-2</i>
Reporting Summary	Yes	nr-reporting-summary.pdf	
Peer Review Information	Yes.	<i>BanfieldTPRFile.pdf</i>	

B. Additional Supplementary Files

10

11

12

Type	Number If there are multiple files of the same type this should be the numerical indicator. i.e. "1" for Video 1, "2" for Video 2, etc.	Filename This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. i.e.: <i>Smith_Supplementary_Video_1.mov</i>	Legend or Descriptive Caption Describe the contents of the file
Supplementary Tables	1-5	combined_supp_tables.xlsx	<p>Table S1: Table of source metagenomes used in this study.</p> <p>Table S2: Table of clades of alternatively coded phages found in this study.</p> <p>Table S3: Table of alternatively coded phage genomes and relatives from this study.</p> <p>Table S4: Table of tRNAs for all alternatively coded phage genomes and relatives from this study.</p> <p>Table S5: Table of release factor and tRNA synthetase counts in alternatively coded phage clades.</p>

13

14 **3. Source Data**

15

16

Parent Figure or Table	Filename	Data description
-------------------------------	-----------------	-------------------------

	This should be the name the file is saved as when it is uploaded to our system, and should include the file extension. i.e.: <i>Smith_SourceData_Fig1.xls</i> , or <i>Smith_Unmodified_Gels_Fig1.pdf</i>	i.e.: Unprocessed Western Blots and/or gels, Statistical Source Data, etc.
Source Data Fig. 1	Fig1_source.xlsx	Statistical Source Data
Source Data Fig. 2	Fig2_source.xlsx	Statistical Source Data
Source Data Fig. 3	Fig3_source.xlsx	Statistical Source Data
Source Data Fig. 4	Fig4_source.xlsx	Statistical Source Data
Source Data Extended Data Fig. 1	ExD_Fig1_source.xlsx	Statistical Source Data
Source Data Extended Data Fig. 2	ExD_Fig2_source.xlsx	Statistical Source Data
Source Data Extended Data Fig. 3	ExD_Fig3_source.xlsx	Statistical Source Data
Source Data Extended Data Fig. 4	ExD_Fig4_source.xlsx	Statistical Source Data
Source Data Extended Data Fig. 5	ExD_Fig5_source.xlsx	Statistical Source Data

18

19 **Widespread stop-codon recoding in bacteriophages may regulate translation of lytic genes**

20

21 Adair L. Borges^{1,2}, Yue Clare Lou^{1,3}, Rohan Sachdeva^{1,4}, Basem Al-Shayeb^{1,3}, Petar I. Penev⁴,
22 Alexander L. Jaffe³, Shufei Lei⁴, Joanne M. Santini⁵, Jillian F. Banfield^{1,2,4,6,7*}

23

24 ¹ Innovative Genomics Institute, University of California, Berkeley, CA, USA

25 ² Environmental Science, Policy and Management, University of California, Berkeley, CA, USA

26 ³ Department of Plant and Microbial Biology, University of California, Berkeley, CA, USA

27 ⁴ Earth and Planetary Science, University of California, Berkeley, CA, USA

28 ⁵ Department of Structural and Molecular Biology, Division of Biosciences, University College
29 London, London, UK

30 ⁶ Lawrence Berkeley National Laboratory, Berkeley, CA, USA

31 ⁷ The University of Melbourne, Australia

32

33 * Corresponding author: jbanfield@berkeley.edu

34

35

36 **Abstract**

37 Bacteriophages (phages) are obligate parasites that use host bacterial translation machinery to
38 produce viral proteins. However, some phages have alternative genetic codes with reassigned
39 stop codons that are predicted to be incompatible with bacterial translation systems. We
40 analysed 9422 phage genomes and found that stop-codon recoding has evolved in diverse
41 clades of phages that infect bacteria present in both human and animal gut microbiota. Recoded
42 stop codons are particularly over-represented in phage structural and lysis genes. We propose
43 that recoded stop-codons might function to prevent premature production of late-stage proteins.
44 Stop-codon recoding has evolved several times in closely related lineages, which suggests that
45 adaptive recoding can occur over very short evolutionary timescales.

46

47 **Main Text**

48

49 **Introduction**

50 The genetic code is highly conserved and considered to be evolutionarily static¹. However, some
51 organisms have alternate genetic codes that reassign one or more codons². Alternative genetic

52 codes are seen in the nuclear genomes of some ciliates³⁻⁵, diplomonads⁶, green algae⁷ and
53 yeasts^{8,9}, as well as genomes of some endosymbionts and mitochondria². Among bacteria,
54 *Mycoplasma*^{10,11} and *Spiroplasma*¹² have reassigned the TGA stop codon to tryptophan (genetic
55 code 4), and members of the Candidate Phyla Radiation (CPR) Gracilibacteria and
56 Absconditabacteria have reassigned the TGA stop codon to glycine¹³⁻¹⁵ (genetic code 25).
57 Furthermore, a computational screen revealed that some bacterial lineages have also reassigned
58 codons that encode amino acids¹⁶.

59 Alternate genetic codes and genetic code expansion can be beneficial. Programmed
60 incorporation of selenocysteine (the 21st amino acid) into selenoproteins is directed at specific in-
61 frame TGA codons in both prokaryotes¹⁷ and eukaryotes¹⁸, and pyrrolysine (the 22nd amino acid)
62 is inserted at in-frame TAG codons in some Archaea, where it boosts enzyme activity^{19,20}. In an
63 unusual case, the pathogenic yeast *Candida albicans* has almost completely reassigned a serine
64 CTG codon to leucine, but still decodes CTG as serine at low levels. This codon-level ambiguity
65 expands the yeast proteome and generates phenotypic diversity, potentially increasing
66 adaptability²¹.

67 Some large, uncultivated phages of the gut microbiome - predominantly Lak phages and
68 crAssphages - have recoded the TAG or TGA stop codon²²⁻²⁶ (genetic codes 15 and 4). A 2014
69 analysis of a single recoded phage genome²² (now classified as a crAssphage²⁵) proposed
70 alternative coding could be a manifestation of phage-host antagonism. In this model, the TAG-
71 recoded crAssphage was hypothesized to infect TGA-recoded bacteria, with both phage and host
72 disrupting translation of the others' genes. However, recent analyses²⁵ revealed that recoded
73 crAssphages infect standard code hosts (genetic code 11), which has cast doubt on that model.
74 Thus, it remains unclear why some phages have evolved genetic codes that are incompatible with
75 host translation systems. In-frame stop codons should induce phage lethality by preventing
76 translation of full length gene products, motivating our study of phages that employ alternative
77 genetic codes with recoded stop codons.

78 Here, we carry out an analysis of stop-codon recoding in 9422 phage genomes recovered
79 from human and animal gut metagenomes. We identify diverse lineages of phages with recoded
80 stop codons that are predicted to infect bacteria that use standard code, and use gene and
81 genome level analyses to propose a regulatory role for stop-codon recoding in the phage life
82 cycle.

83

84 **Results**

85

86

87 **Genome recoding in gut microbiome phage**

88 We recovered 9422 complete or nearly complete ($\geq 90\%$ complete) dereplicated phage genomes²⁷
89 from 726 human and animal gut metagenomes (**Supplementary Table 1**). To broadly sample
90 phage diversity within the human gut, we analyzed gut microbiomes from individuals inferred to
91 consume westernized^{28–31} and non-westernized diets^{23,29,30,32}, based on diet and location-related
92 metadata provided in the original studies. To sample phage diversity beyond the human gut, we
93 recovered phages from gut microbiomes of baboons³³, pigs^{34,35}, cattle³⁶, horses, and giant
94 tortoises. To identify instances of TAG or TGA stop-codon recoding, we predicted phage genes
95 in standard code (code 11), or alternative genetic codes with TAG or TGA stop codons recoded
96 (code 15 or code 4, respectively) and calculated coding density for each phage genome in each
97 code. As stop codon recoding leads to gene fragmentation in standard code, we identified
98 genomes that underwent a 5-10% coding density increase when genes were predicted with an
99 alternative code (**Fig. 1A, Supplementary Information Fig. 1A-C**). We then manually verified
100 these putative alternatively coded genomes (**Methods**), arriving at a final set of 473 recoded
101 double-stranded DNA phage genomes.

102

103 Previously stop codon recoding had only been found in phages with large genomes:
104 crAssphages^{22,25} (95-190kb), jumbophages²⁴ (200-500 kb), and megaphages^{23,26} (>500kb to 660
105 kb). We identified complete recoded phage genomes across a very wide diversity of sizes, ranging
106 down to 14.7 kb (**Fig. 1B**). We observed that TAG recoding is more common than TGA recoding
107 (75% TAG recoded, 25% TGA recoded, **Fig. 1B,C**). While each gut microbiome type has recoded
108 phages present, recoding was least common in phages recovered from humans inferred to
109 consume a westernized diet, and was most common in baboon phages (**Fig. 1C**). We conclude
110 that alternative coding is a common feature of phage populations in the human and animal gut,
111 and occurs in phages of diverse genome sizes (**Fig. 1D-I, Extended Data Fig. 1A**).

112

113 **Diversity and evolution of recoded phages**

114 We constructed a phylogenetic tree of large terminase subunits from recoded phages and their
115 standard code relatives, finding many sequences form clades with high bootstrap support
116 ($\geq 95\%$)(**Fig. 2A**). Inspired by the historical designation of TAG and TGA as the amber and opal
117 stop codons, we chose to name the six newly discovered clades of TAG and TGA recoded phages
118 after other gemstones (Garnet, Amethyst, Jade, Sapphire, Agate, Topaz). These clade
119 designations are not intended as taxonomic names. Including previously discovered Lak²³ and

120 crAss-like families^{22,25}, we describe 8 independent phage clades that use recoded stop codons in
121 human and animal gut microbiomes (see **Supplementary Table 2** for clade-level data,
122 **Supplementary Table 3** for genome-level data).

123
124 To identify the genetic codes used by the recovered recoded genomes, we analyzed the
125 alignments of terminase sequences with in-frame recoded stop codons translated to X, and found
126 that in most cases TGA aligned with tryptophan (genetic code 4) and TAG aligned with glutamine
127 (genetic code 15). Use of genetic code 15 has since been confirmed in crAss-like phages via
128 metaproteomics of human samples³⁷. Many clades encompassed multiple genetic codes (**Fig.**
129 **2A, Supplementary Table 2**). Some recoded phages used code 25, where TGA is reassigned
130 to glycine. We predict these phages infect *Candidatus Absconditabacteria*, which also uses code
131 25¹⁴(**Supplementary Table 3**). In all other cases, the recoded phage clades are predicted to
132 infect bacteria from common standard code gut phyla, Firmicutes and Bacteroidetes
133 (**Supplementary Table 2, Supplementary Table 3**).

134
135 Lak, crAss, Jade, Sapphire and Agate phages have larger than average genome length ($562 \pm$
136 44kb , $210 \pm 35\text{ kb}$, $201 \pm 22\text{ kb}$, $154 \pm 27\text{ kb}$, $\pm\text{SD}$) and Garnet, Amethyst, and Topaz phages
137 have smaller than average genomes ($34 \pm 5\text{kb}$, $34 \pm 6\text{kb}$, $22 \pm 3\text{ kb}$, $\pm\text{SD}$) (**Supplementary Table**
138 **2**). These eight clades have uneven distributions across the environments analyzed here (**Fig.**
139 **2B**). Notably, the recoded phage present in westernized-diet microbiomes mainly comprises
140 crAss-like phages, whereas other gut microbiomes have higher diversity of recoded phages (**Fig.**
141 **2B**).

142 143 **Mechanisms of phage recoding**

144 In bacteria that use the standard genetic code, the TAG, TGA, and TAA stop codons are
145 recognized by specific release factors (RF1 or RF2), which trigger translation termination.
146 Suppressor tRNAs recognize TAG, TGA, or TAA codons, and have been previously identified in
147 recoded phage genomes^{22–26}, where they presumably mediate code-change. We predicted tRNAs
148 in all phage genomes, and calculated the frequency at which genomes of each code encoded
149 suppressor tRNAs for the TAG, TGA, or TAA stop codons. We found a strong relationship
150 between stop codon recoding and suppressor tRNA usage, detecting TAG suppressor tRNAs in
151 40% of TAG recoded genomes, and TGA suppressor tRNAs in 35% of TGA recoded (code 4)
152 genomes (**Fig. 2C, Supplementary Table 4**) Surprisingly, when we analyzed suppressor tRNA
153 occurrence across phage phylogeny, we found that suppressor tRNAs were strongly partitioned

154 between phage clades (**Fig. 2A**). Specifically, almost all the suppressor tRNAs detected were
155 found in the Lak, crAss-like, Sapphire, and Agate phages which have large genomes. In contrast,
156 the small-genome Garnet, Amethyst, and Topaz clades rarely encoded suppressor tRNAs.

157

158 We also searched for phage-encoded release factors (RF), which terminate translation at “true”
159 stop codons and have been previously observed in Lak²³ and crAss-like²² phages. We identified
160 RF2 (terminates translation at TAA and TGA) in six TAG recoded Lak phages and two TAG
161 recoded Agate phages (**Supplementary Table 5**). RF1, which terminates translation at TAG and
162 TAA stop codons was identified in two TGA recoded Jade phages (**Supplementary Table 5**). We
163 also identified tryptophanyl tRNA synthetases in the same two Jade genomes (**Supplementary**
164 **Table 5**), which we predict could ligate the amino acid tryptophan to the TGA suppressor tRNA,
165 thus mediating the TGA → W code change.

166

167

168 **Relationships between recoded and standard-code phages**

169 We next calculated the average nucleotide identity (ANI) between all phage in our dataset to
170 identify examples of very closely related genomes that use different genetic codes. We identified
171 a set of Agate clade genomes with greater than 80% ANI that includes TGA-recoded code 4
172 genomes and standard code genomes (**Fig. 3A**). One standard code phage had acquired a TGA
173 suppressor tRNA (+), potentially preceding the code change (**Fig. 3A**). We also found an example
174 where a TGA recoded Agate phage Cattle_ERR2019405_scaffold_1063 and a standard code
175 Agate phage pig_ID_3053_F60_scaffold_12 share greater than 90% ANI (**Fig. 3A, Extended**
176 **Data Fig. 3A**). This indicates that genetic code can change over short evolutionary timescales.

177

178 The Agate phage pig_ID_3053_F60_scaffold_12 uses standard code and only 4 out of 146 genes
179 (2.7%) use TGA stop codons. In contrast, 34 genes use TAG and 108 genes use TAA. This
180 suggests divestment in TGA as a stop codon may precede its reassignment in Agate phages,
181 consistent with the codon capture hypothesis of genetic code evolution³⁸. To test for TAG or TGA
182 stop codon loss at a wider scale, we surveyed stop codon usage across all standard code phages
183 that are closely related to recoded phages. These relatives are more likely to use the TAA stop
184 codon than the TAG and TGA stop codons (**Fig. 3B**, TAG vs. TAA: $Z = -19.71$, $p = 1.63e-86$, TGA
185 vs. TAA: $Z = -19.65$, $p = 5.98e-86$, two-sided Wilcoxon Rank-Sum Test). We also observed that
186 TAG is rarer than TGA (**Fig. 3B**, $Z = -6.43$, $p = 1.24e-10$, two-sided Wilcoxon Rank-Sum Test).
187 This depletion of TAG and TGA is likely due to the reduced GC content (**Fig. 3C**, $Z = -12.59$, $p =$

188 2.33e-36, two-sided Wilcoxon Rank-Sum Test) in these phages compared to standard code
189 phages that are not close relatives of recoded phages. Thus, stop codon loss driven by low GC
190 content may be an evolutionary precursor to stop codon recoding.

191

192 **Recoding may regulate cell lysis**

193 We analyzed functional predictions for genes with in-frame recoded stop codons in the genomes
194 of representatives of each recoded phage clade (**Fig. 4A-B, Extended Data Fig. 4A-D,**
195 **Extended Data Fig. 5A-C**). Consistent with previous observations, we saw that both Lak²³ and
196 crAss-like^{22,25} genomes use alternative code for their “late” structural and lysis genes.
197 Furthermore, we observed that Garnet, Amethyst, Jade, Sapphire, Agate, and Topaz phages also
198 use alternative code for structural and lysis genes. In contrast, use of alternative code was
199 variable in the DNA replication machinery. In crAss-like, Garnet, Amethyst, and Topaz phages,
200 all the structural and lysis genes are encoded together a single alternatively-coded genomic unit
201 (**Fig. 4A, Extended Data Fig. 4D, Extended Data Fig. 5B-C**). In Jade, Sapphire, Agate, and Lak
202 phages, the structural and lysis genes are in multiple alternatively-coded modules that are spread
203 across the genome (**Fig. 4B, Extended Data Fig. 4A-C, Extended Data Fig. 5A**). As structural
204 and lysis proteins encoded with recoded stop codons cannot be expressed before the code
205 change is manifested, stop codon recoding could effectively regulate the timing of protein
206 expression from related gene modules.

207

208 We next identified the gene families most biased towards use of recoded stop codons, as they would be
209 most impacted by this proposed form of gene regulation. We measured the codon preference in two
210 phage types that were represented by a sufficiently large set of related genomes to enable gene-
211 level statistics: ~105 kb TAG-recoded crAss-like phages and ~127 kb TGA-recoded Agate
212 phages. While many genes in these phages have at least one in-frame recoded stop codon (**Fig.**
213 **4A-B**), only a few gene families preferentially use recoded stop codons over standard code
214 encodings of glutamine (crAss-like phages, TAG → Q) or tryptophan (Agate phages, TGA → W).

215

216 In the crAss-like phage genomes analyzed, only four gene families preferentially use TAG over
217 CAG or CAA to encode Q (two-sided Wilcoxon Rank-Sum Test, corrected for multiple
218 comparisons) (**Fig. 4C**). Two of these four families are essential components of the lysis cassette:
219 a lysozyme type amidase ($Z = 2.91$, $p = 6.82e-3$) and a spanin, which is a critical regulator of lysis
220 of gram-negative bacteria^{39,40} ($Z = 4.82$, $p = 9.00e-6$). A tail tube gene family that is encoded two
221 genes downstream (1.1 kb) of the spanin gene is also preferentially recoded ($Z = 3.56$, $p = 7.59e-$

222 4). Having multiple strongly alternatively coded genes in the same transcript may amplify the stop
223 codon mediated translation block. The fourth recoded gene family is of unknown function.

224

225 In the Agate genomes analyzed, three gene families preferentially use TGA instead of TGG to
226 encode W (two-sided Wilcoxon Rank-Sum Test, corrected for multiple comparisons) (**Fig. 4D**).
227 One of these is a group I intron endonuclease ($Z = 3.88$, $p = 1.32e-3$) inserted in the DNA
228 replication module, which predominately uses standard code. This self-splicing intron is expected
229 to excise itself from the mRNA, but then in-frame recoded stop codons should prevent homing
230 endonuclease production until late in the infection cycle. A tail gene directly upstream of the lysis
231 cassette ($Z = 2.69$, $p = 4.18e-2$) is also preferentially recoded, analogous to the tail tube gene in
232 the crAss-like phages. A preferentially recoded transmembrane domain protein ($Z = 2.63$, $p =$
233 $4.64e-2$) in the lysis cassette belongs to a family of transmembrane proteins that are assigned
234 various lysis and lysis regulation related functions (holin, spanin, lysis regulatory protein, and ATP
235 synthetase B chain precursor). We hypothesize a putative role for this protein in controlling lysis,
236 potentially by depolarizing the cell membrane⁴¹⁻⁴⁴.

237

238 We also identified a “code change module” composed of a suppressor tRNA, a tRNA synthetase,
239 and a release factor directly upstream of the lysis cassette in Jade phages (**Extended Data Table**
240 **2, Extended Data Fig. 6A-B**). These code-change related genes are all encoded in standard
241 code, whereas the lysis genes directly downstream use alternative code. We anticipate that
242 expression of these code change genes would drive expression of the lysis program. Overall, we
243 propose that by changing the genetic code of the infected cell over time, these phages can use
244 stop codon recoding to coordinate protein expression from related late genes and also to
245 suppress misexpression of critical lytic gene products.

246

247 **Is recoding in prophages a lysogeny switch?**

248 Many phages integrate into their bacterial host chromosome as prophages. Excitingly, we found
249 recoded Garnet and Topaz prophages integrated into standard code bacterial contigs, two of
250 which we analyzed in depth (**Fig. 5A-B**).

251

252 The Garnet prophage is part of a 94 kb *Prevotella* contig (SRR1747048_scaffold_47) assembled
253 from a baboon metagenome. (**Fig. 5A**). When we mapped reads to this prophage we observed
254 that the sequencing read depth of the bacterial region was twice that of the integrated prophage
255 (**Extended Data Fig. 7A**). Some reads spanned the prophage, corresponding to *Prevotella*

256 genomes that lack the integrated prophage. Thus, the exact prophage 24,371 bp genome could
257 be defined.

258

259 The Topaz genome is part of a 36.9 kb *Oscillospiraceae* contig (SRR1747065_scaffold_956)
260 assembled from a baboon metagenome (**Fig. 5B**). In this case, sequencing reads coverage over
261 the prophage region is ~50 times higher than the flanking genome (**Extended Data Fig. 7B**). We
262 infer the vast majority of phages in the sample were replicating and only a subset remained
263 integrated at the time of sampling. Based on the sequence margins, we determined that the length
264 of this prophage genome is 23,706 bp.

265

266 We also identified circular free phage genomes in related baboon samples that were nearly 100%
267 identical to the Garnet and Topaz prophages analyzed here. This supports our conclusion that
268 these prophages represent actively-replicating viable phages (**Extended Data Fig. 8A-B**) and
269 verifies the lengths determined from the read mapping analysis.

270

271 We noticed that while almost all of the prophage genes were extremely fragmented in standard
272 code, the integrase genes did not contain recoded stop codons (**Fig. 5A-B, Extended Data Fig.**
273 **4D, Extended Data Fig. 5B**). When we measured codon preference across all gene families
274 encoded by alternatively coded Garnet and Topaz phages, we found the integrase gene families
275 strongly avoided use of recoded stop codons. (Garnet: $Z = -3.97$, $p = 5.12 \times 10^{-3}$, Topaz: $Z = -8.87$, p
276 $= 1.23 \times 10^{-16}$, two-sided Wilcoxon Rank-Sum Test, corrected for multiple comparisons). We
277 hypothesize that these phages are using stop codon recoding as a regulator of the lytic-lysogenic
278 switch. In this scenario, the standard code translation environment of the host promotes
279 expression of the integrase and establishment of lysogeny, with strong suppression of lytic genes.
280 Likewise, a switch to alternative code would promote expression of lysis-related proteins during
281 initial infection or prophage induction. Thus, genetic code may function as a mechanism to
282 partition two distinct arms of the phage life cycle.

283

284

285 **Discussion**

286 Using a computational analysis, we detected widespread use of recoded stop codons in eight
287 families of phage and prophage present in human and animal gut microbiomes. We hypothesize
288 that evolution of alternate code involves ancestral depletion of TAG and TGA stop codons, and

289 propose a model in which stop codon recoding is a post-transcriptional regulator of protein
290 expression in phages and prophages (**Fig. 6**).

291 We propose an evolutionary route to recoding that begins with depletion of TAG or TGA
292 stop codons in standard code phages with low GC content. Via acquisition of a suppressor tRNA,
293 in-frame stop codons can accumulate in positions that would previously have been lethal for the
294 phage. We identified TGA suppressor tRNA acquisition by standard code close relatives of TGA
295 recoded phages, which supports this model. We also found that TAG stop codons are more rare
296 than TGA stop codons in standard code relatives of recoded phages, potentially explaining the
297 higher prevalence of TAG recoding compared with TGA recoding.

298 After in-frame stop codons are “detoxified” by the acquisition of suppressor mechanisms
299 such as tRNAs, selection enriches or depletes recoded stop codons across specific gene families
300 to create patterns of codon use that can be harnessed as a form of gene regulation. Clades of
301 recoded phages have independently converged upon using recoded stop codons to encode lysis
302 and structural proteins. This is consistent with more limited observations of structural gene
303 recoding seen in Lak²³ and crAsslike^{22,25} phages, and supports a model where the genetic code
304 of the infected cell changes throughout the phage infection cycle. Dynamic codon use throughout
305 the infection cycle has been demonstrated in T4-like phages that encode large tRNA arrays,
306 where late-expressed genes have codon use aligned with the phage tRNA repertoire^{45,46}. This
307 may represent a mechanism to toggle translation efficiency of late genes throughout the phage
308 life cycle.

309 Stop codons have low to no translation efficiency, so we hypothesize that use of recoded
310 stop codons in late expressed genes is an extreme form of codon based regulation in phages.
311 We found two distinct lineages of phages with preferential recoding of the lysis cassette, for which
312 precisely timed expression is crucial. Premature lysis aborts the phage life cycle and limits phage
313 production, and some anti-phage immune systems even exploit this by forcing early lysis^{44,47}. By
314 encoding lysis regulators with in-frame recoded stop codons, these phages block both accidental
315 or host-forced premature expression of these proteins.

316 We also identified prophages with recoded stop codons that were integrated into standard
317 code hosts. The decision to enter lytic growth or lysogeny is a crucial point in the temperate phage
318 life cycle, and phages have evolved elaborate regulatory mechanisms to precisely control this
319 decision^{48–50}. We hypothesize that alternate coding may function in the lysis-lysogeny switch in
320 these recoded temperate phages.

321 Most suppressor tRNAs identified here were encoded by phages with large genomes,
322 consistent with previous reports that tRNAs number increases with phage genome size^{24,51}.

323 However, we predict that all recoded phages that infect standard code hosts would require
324 suppressor tRNAs to decode recoded stop codons. One possibility is that small phages “piggy-
325 back” on large phages of the same code, to use larger-phage suppressor tRNAs during
326 coinfection. Some huge phages have been shown to carry CRISPR-Cas systems that target
327 small phages²⁴, consistent with the hypothesis that small phages may parasitize large phages.

328 Stop codon recoding could allow phages of any size to sense the presence of co-resident
329 phages that use the same genetic code via activity of translation-related molecules such as
330 suppressor tRNAs. This would be beneficial to prophages that are induced in response to a
331 superinfecting lytic phage, or for coinfecting phages to coordinate the timing of their lytic program.

332

333 **Conclusion**

334 Stop codon recoding may have an important but previously unappreciated role in the phage life
335 cycle. Further, understanding alternative genetic code use in phage is crucial to our ability to
336 detect and classify phage sequences. Broadening our view of genetic code diversity in phages
337 has the potential to augment our understanding of basic phage biology and bacterial translation,
338 as well as improving synthetic biology strategies to design new genetic codes.

339

340 **Online Methods**

341

342 **Phage prediction**

343 Phage prediction tools Seeker⁵² (predict-metagenome) and VIBRANT⁵³ were run on assembled
344 metagenomes (contigs > 5kb) using default settings. CheckV⁵⁴ (end-to-end) was run on predicted
345 phages and trimmed proviruses to evaluate completeness and quality. Contigs evaluated as low
346 quality by both CheckV and VIBRANT were removed from analysis. Contigs < 100 kb with viral
347 genes > host genes and contigs > 100 kb with < 20% host genes were maintained as high
348 confidence phages. All deposited phage genomes are compliant with MIUVIG standards²⁷.

349

350

351 **Phage dereplication**

352 Phage scaffolds for each ecosystem were dereplicated at 99% ANI using the dRep⁵⁵ dereplicate
353 module (-sa 0.99 --ignoreGenomeQuality -l 5000 -nc 0.5 --clusterAlg single -N50W 0 -sizeW 1).

354

355 **Identification of phage genomes with recoded stop codons**

356 Prodigal⁵⁶ (single mode) was used to predict genes on dereplicated $\geq 90\%$ complete phage
357 genomes using genetic codes 4, 11, and 15. Coding density was calculated by summing the
358 length of genes for each contig and dividing by the total contig length. Contigs 5-100 kb that had
359 an increase of greater than 10% coding density in code 4 or code 15 relative to code 11 were
360 tentatively assigned that genetic code, as were contigs ≥ 100 kb with a coding density increase
361 $>5\%$. All code assignments were confirmed by manual analysis of each contig. If the alternative
362 genetic code resulted in more contiguous operon structure, reduced strand switching, correct-
363 length genes (as checked by blastp⁵⁷ against NCBI database), and did not result in gene fusions
364 (as checked by blastp against NCBI database) the phage was confirmed as alternatively coded.

365

366 **Structural and functional annotations**

367 Coding sequences predicted by prodigal using genetic code 4 for TGA recoded phages and code
368 15 for TAG recoded phages. HMMER⁵⁸ (hmmsearch) was used to annotate the resulting
369 sequences with the PFAM, pVOG, VOG, and TIGRFAM HMM libraries. tRNAs were predicted
370 using tRNAscan-s.e. V.2.0 in general mode⁵⁹.

371

372 **Host prediction**

373 A combination of CRISPR spacer analysis and taxonomic classification were used to predict
374 putative host phyla for recoded phages and their standard code relatives. Contigs with a minimum
375 length of 5 kb from the human and animal metagenomes analyzed in this study were searched
376 for CRISPR spacers using minCED⁶⁰. blastn short was used to identify matches between phage
377 and spacer of $>90\%$ identity and $>90\%$ spacer coverage. Taxonomic profiling was performed by
378 using DIAMOND⁶¹ (fast mode, $e = 0.0001$) to search all phage proteins against a custom version
379 of the UNIREF100 database that retained NCBI taxonomic identifiers. tRep⁶² was then used to
380 profile the taxonomy of each phage contig. For each contig, the bacterial phylum with most hits
381 was considered to be the putative host, but only if that phylum had more than 3x hits than the
382 second most common phylum²⁴. In almost every case, the CRISPR spacer analysis and the
383 taxonomic profiling agreed on the phage host phyla. In the rare cases that these analyses were
384 not in agreement, the host phyla was considered unknown.

385

386

387 **Phage genome clustering by average nucleotide identity (ANI)**

388 Our total dataset of 9422 non-dereplicated phage scaffolds from all ecosystems was augmented
389 with 1428 phage genomes from other animal/human microbiomes from ggkbase, and the

390 genomes clustered using dRep⁵⁵ compare module (-sa 0.8 -pa 0.8 -nc .1 --clusterAlg single).
391 Whole genome alignment was visualized using Mauve⁶³ (progressiveMauve algorithm)
392 implemented in Geneious Prime 2021.0.3 (<https://www.geneious.com>).

393

394

395 **Phage clustering with Vcontact2**

396 Phages scaffold from the dereplicated dataset of $\geq 90\%$ complete phage scaffolds for each
397 ecosystem were clustered into viral clusters with Refseq viruses using Vcontact2⁶⁴ (--rel-mode
398 'Diamond' --db ProkaryoticViralRefSeq201-Merged --pcs-mode MCL --vcs-mode ClusterONE).
399 Standard code phages that were in the same viral cluster (VC) as at least one alternatively coded
400 phage were considered to be close relatives of alternatively coded phages.

401

402 **Phylogenetic analysis of large terminase subunit of recoded phages and standard code** 403 **relatives**

404 Terminases were found using two rounds of HMM-based classification. Proteins were initially
405 annotated using PFAM, pVOG, VOG, and TIGRfam HMMs. This did not result in complete
406 recovery of terminases for all phages of interest. To increase sensitivity, we clustered proteins
407 into subfamilies using MMseqs⁶⁵ (-s 7.5, -c 0.5, -e 0.001), and used HHblits⁶⁶ to generate hmms
408 of each subfamily based on alignments generated with the MMseqs result2msa parameter. We
409 used HHSearch⁶⁷ (-p 50 -E 0.001) to perform an HMM-HMM comparison with the PFAM
410 database. We then identified subfamilies with a best hit to large terminase HMMs with a $>95\%$
411 probability. Putative terminase subfamilies with a low number of primary terminase annotations
412 were confirmed by blastp against the NCBI database. If subfamily members had hits to terminases
413 in known phages, we considered the subfamily to be a true terminase subfamily. In rare cases,
414 the terminase gene was fragmented due to assembly error or mobile intron insertion. In these
415 cases we chose the larger of the gene fragments for downstream analysis. Terminases from
416 recoded phages and these standard code relatives (from vContact2⁶⁴) were searched against
417 the Refseq protein database using blastp, retaining the top 10 hits per protein. The recovered
418 Refseq proteins were dereplicated at 90% using CD-HIT⁶⁸. Recoded phage, standard code
419 relative, and dereplicated Refseq terminases were combined and aligned using MAFFT⁶⁹, and
420 the alignment trimmed with trimAL⁷⁰ (-gt 0.5). IQ-TREE⁷¹ was used to build a tree using the
421 VT+F+R10 model and ultrafast bootstrap with 1000 iterations. Tree was visualized using iTol⁷².

422

423 **Codon preference analysis**

424 **TAG-recoded crAss and TGA-recoded Agate analysis:** ANI-based genome clustering showed
425 high representation of a lineage of TGA recoded ~127 kb Agate phages as well as a lineage of
426 TAG recoded ~105 kb crAss-like phages, which were chosen for further analysis. For each phage
427 lineage, proteins were clustered into families created using a two step protein clustering method.
428 First, proteins were clustered into subfamilies using MMseqs⁶⁵ (-s 7.5, -c 0.5, -e 0.001), and
429 HHBlits⁶⁶ was used to generate HMMs of each subfamily based on alignments generated with
430 the MMseqs result2msa parameter. These HMMs were then compared to one another using
431 HHBlits (-p 50 -E 0.001). MCLclustering (--coverage 0.70 -l 2.0 --probs 0.95) was used to
432 generate families from the HMM-HMM comparisons. Two-sided Wilcoxon rank sum test was used
433 to evaluate protein families that preferred the in-frame recoded stop codon to the standard coding
434 for the recoded amino acid. The Benjamini-Hochberg p-value correction was used to correct for
435 multiple hypothesis testing with a false discovery rate of 5%. For TGA → W recoded phages, TGA
436 occurrence was compared to the occurrence of the standard codon for Tryptophan (TGG). For
437 TAG → Q recoded phages, TAG occurrence was compared to the occurrence of the standard
438 codons for Glutamine (CAG, CAA). Proteins were annotated by PFAM, pVOG, VOG, and
439 TIGRFAM as well as BLAST searches against the NCBI database. In some cases, the HHPred
440 webserver⁷³ and the Phyre2 webserver⁷⁴ were used to augment initial annotations. Gene
441 neighborhoods were visualized using Clinker⁷⁵.

442
443 **Garnet and Topaz integrase analysis:** Garnet and Topaz proteins were clustered into families
444 using the two step method detailed above. We identified the integrase families for each phage
445 clade using PFAM, pVOG, VOG, and TIGRFAM HMM annotations. We observed that the majority
446 of the integrase genes had zero in-frame recoded stop codons. A few genes had one in-frame
447 stop, and when we examined alignments of the integrase families we found that in all cases the
448 in-frame recoded stop was in a N or C terminal extension of the protein. We believe that this
449 corresponds to incorrect start codon prediction (N terminal extensions) or legitimate use of the
450 codon to terminate the integrase gene (C terminal extensions). We used a two-sided Wilcoxon
451 rank sum test to evaluate all protein families in each phage clade for avoidance of in-frame
452 recoded stop codons relative to the rates at which they use the standard codons for Glutamine
453 (for TAG → Q recoded phages) or Tryptophan (for TGA → W recoded phages). The Benjamini-
454 Hochberg p-value correction was used to correct for multiple hypothesis testing with a false
455 discovery rate of 5%. We found that for both Garnet and Topaz phages, the integrase gene
456 families strongly avoided in-frame recoded stop codons relative to the rate at which they used

457 standard code encodings for glutamine (TAG → Q recoded phages) or tryptophan (TGA → W
458 recoded phages)

459

460 **Origin and terminus determination via GC Skew**

461 GC skew (G-C/G+C) and cumulative GC skew were calculated across the phage genome using
462 the iRep package (gc_skew.py)⁷⁶. This allowed us to predict origins of replication, replication
463 termini, and define individual replichores. We observed a variety of replication styles: double origin
464 bi-directional replication, single origin bi-directional replication, and unidirectional replication. We
465 also observed GC skew patterns of unknown significance. See **Supplementary Figure 2A-I** for
466 cumulative GC skew plots from the representatives of each phage clade.

467

468 **Lysogen read mapping**

469 Reads from the source metagenome were mapped against lysogenic contigs with Bowtie 2⁷⁷
470 using default settings. Contigs and mapped reads were visualized in Geneious Prime 2021.0.3
471 (<https://www.geneious.com>).

472

473 **Statistics** **and** **Reproducibility**

474

475 This study was designed to capture a broad range of gut microbiome phage diversity. We
476 recovered phage from 7 gut microbiome ecosystem types that we and others had sampled
477 sufficiently to allow high recovery of near-complete phage genomes. Only high confidence phage
478 genomes were used in this study. Phage-like contigs that were evaluated as low quality by both
479 CheckV and VIBRANT were excluded from this study. Phage-like contigs < 100 kilobases with
480 more host genes than viral genes were excluded from this study. Phage-like contigs > 100
481 kilobases with > 20% host genes were excluded from this study. We validated these cutoffs by
482 manually inspecting contigs with high host gene content, and found that they often represented
483 plasmids or chromosomal fragments. These cutoffs were employed to ensure we only had phage
484 genomes in our dataset. We also excluded phage genomes that were less than 90% complete
485 from our survey. Since stop codon recoding is often only present in part of the genome, the
486 recoded region of the genome may be greatly reduced or even entirely missing from an
487 incomplete genome. This means that use of genome fragments to determine phage genetic code
488 is unreliable. All phage genomes in our study were dereplicated, to ensure we were measuring
489 independent phage genomes, and were not measuring the “same” phage across multiple different
490 samples. We used a two-sided Wilcoxon Rank-Sum Test to compare differences between groups

491 of genomes (GC content, stop codon use) or gene families (alternative coding bias). When
492 comparing large numbers of gene families, we used Benjamini-Hochberg p-value correction to
493 correct for multiple hypothesis testing with a false discovery rate of 5%. No statistical method was
494 used to predetermine sample size for any analyses. The experiments were not randomized. The
495 Investigators were not blinded to allocation during experiments and outcome assessment.

496

497 **Data Availability**

498

499 Accessions for MIUVIG-compliant genomes²⁷ and associated reads for alternatively coded
500 phages and relatives are provided in **Supplementary Table 3**. Genomes and predicted proteins
501 for alternatively coded phages and relatives, the terminase phylogenetic tree file, closely related
502 Agate and crAss-like genomes, and untrimmed lysogenic contigs are available through Zenodo
503 (10.5281/zenodo.6410225). The UniRef100 database is available through
504 <ftp.uniprot.org/pub/databases/uniprot/uniref/uniref100>.

505

506

507 **Code Availability**

508

509 Python script used to analyze coding density and predict genetic code is available on Github:
510 https://github.com/borgesadair1/AC_phage_analysis/releases/tag/v1.0.0

511

512 **Acknowledgments**

513 We thank Yun Song, Jamie Cate, Kimberly Seed, Grayson Chadwick, Lin-Xing Chen, Jacob
514 West-Roberts, and Spencer Diamond for helpful discussions. We thank Ka Ki Lily Law and Jordan
515 Hoff for technical support. This work was supported by a Miller Basic Research Fellowship to
516 A.L.B, an NSF Graduate Research Fellowship to B.A-S (No. DGE 1752814), and NIH award
517 RAI092531A, a Chan Zuckerberg Biohub award, and Innovative Genome Institute funding to
518 J.F.B.

519

520 **Author Contributions**

521 A.L.B and J.F.B. developed the project, led analyses, and wrote the manuscript with input from
522 all authors. A.L.B, J.F.B, Y.C.L, R.S., and S.L. compiled the phage dataset. B.A-S assembled
523 public metagenome data and provided support for phage genome analyses. Phage genomes
524 were manually curated by J.F.B. P.I.P contributed to phage tRNA analyses. A.L.J. and Y.C.L

525 contributed to design of statistical analyses. J.M.S. contributed DNA samples from animal and
526 arsenic-exposed human gut microbiomes.

527

528 **Competing Interests**

529 J.F.B. is a founder of Metagenomi. The other authors declare no competing interests.

530

531 **Main Text Figure Legends**

532

533 **Fig. 1: Identification of recoded phage in human and animal microbiomes.**

534 **A.** A 5-10% coding density increase between standard code and alternative code was used to
535 identify putative recoded phages, followed by manual confirmation of code.

536 **B.** Recoded phage genomes spanned a wide size range from 14.7 kilobases (kb) to 660 kb.

537 **C.** Abundance of recoded phages varied from ~2-6 % of the total phage population in the gut
538 microbiome types surveyed in this study. WD = westernized diet, NWD = non-westernized diet.

539 **D-I.** Phage genomes recovered from the indicated human or animal microbiome. The number of
540 phage genomes (n) recovered after dereplication from each environment is indicated in the title
541 of each plot. Individual phage genomes are represented by single points and plotted by genome
542 size and coding density (CD) in standard code (code 11). In all plots, phage genomes have been
543 dereplicated and are complete or near complete ($\geq 90\%$). Symbol color represents genetic code
544 (TGA recoding = green, TAG recoding = orange, standard code = grey). See **Extended Data Fig.**
545 **1A-F** for plots with coding density re-calculated using the predicted genetic code.

546

547 **Fig. 2: Phylogeny of recoded phages and suppressor tRNA usage.**

548 **A.** The phylogeny of recoded phages was reconstructed using large terminase sequences from
549 a dereplicated set of complete or near-complete ($\geq 90\%$) recoded phages (n=444) and their close
550 standard code relatives (n =258), as well as related proteins from Refseq r205 (n=410).
551 Terminase sequences from eukaryotic herpesviruses (n=8) were used to root the tree. The inner
552 to outer ring shows phage clade ($\geq 95\%$ bootstrap support), genetic code for phages from this
553 study, suppressor tRNA presence, host phylum as predicted by taxonomic profiling and CRISPR
554 spacer matches, and genome size with a grey line at 100 kilobases (kb) for scale. Genetic code,
555 suppressor tRNA presence, and genome size were not included for Refseq proteins since some
556 proteins were derived from prophages and/or incomplete phage genomes.

557 **B.** Distribution of recoded phages by clade across the 7 types of gut microbiomes evaluated in
558 this study. WD = westernized diet, NWD = non-westernized diet.

559 **C.** Heatmap of the percent of genomes of each genetic code that have tRNAs predicted to
560 suppress translation termination at the TAG, TGA, or TAA stop codons.

561

562 **Fig. 3: Evolutionary relationships among phages according to genetic code**

563 **A.** Dendrogram of average nucleotide identity (ANI) across a set of Agate phage genomes.
564 Standard code (grey) and TGA recoded (green) phages share >80% ANI, and one cluster of
565 >90% ANI genomes (orange) has both standard and TGA recoded genomes, indicating an
566 extremely close evolutionary relationship.

567 **B.** Close relatives of recoded phages (n = 260 biologically independent phage genomes) use the
568 TAA stop codon at a higher rate than the TAG and TGA stop codons (TAG vs. TAA: Z = -
569 19.71, p = 1.63e-86, TGA vs. TAA: Z = -19.65, p = 5.98e-86, two-sided Wilcoxon Rank-Sum
570 Test). The TAG stop codon is depleted relative to TGA (Z = -6.43, p = 1.24e-10, two-sided
571 Wilcoxon Rank-Sum Test) in these phages. TAG frequencies: Minima = 0.0, Maxima = 0.38,
572 Median = 0.09, IQR = 0.14, Q1 = 0.04, Q3 = 0.18. TGA frequencies: Minima = 0.0, Maxima =
573 0.44, Median = 0.16, IQR = 0.13, Q1 = 0.11, Q3 = 0.25. TAA frequencies: Minima = 0.40, Maxima
574 = 0.94, Median = 0.68, IQR = 0.16, Q1 = 0.59, Q3 = 0.76.

575 **C.** Close relatives of alternatively coded phages have a lower mean GC content relative to all other standard
576 code phages (Z = -12.59, p = 2.33e-36, two-sided Wilcoxon Rank-Sum Test). Close relatives: n = 260
577 biologically independent phage genomes, Minima = 27.97, Maxima = 45.73, Median = 35.10, IQR =
578 5.30, Q1 = 33.27, Q3 = 38.57. Unrelated standard code phages: n = 8689 biologically independent phage
579 genomes, Minima = 19.60, Maxima = 67.07, Median = 41.827, IQR = 12.67, Q1 = 35.95, Q3 = 48.62.
580 **** p ≤ 0.0001, two-sided Wilcoxon Rank-Sum test.

581

582 **Fig. 4: Preferential recoding of lysis-related genes in recoded phages**

583 **A-B.** Genomic maps of manually-curated representatives of crAss-like phages (js4906-23-
584 2_S13_scaffold_20) and Agate phages (GiantTortoise_AD_1_scaffold_344). TAG recoded
585 genomes (**A**) contain genes with in-frame TAG codons (orange) while TGA recoded genomes (**B**)
586 have genes with in-frame TGA codons (green). Suppressor tRNAs (red labels) are predicted to
587 suppress translation termination at recoded stop codons. Regions of the genome encoding
588 structural and lysis genes (pink) coincide with high use of alternative code. In these phages, DNA
589 replication machinery (yellow) is encoded in standard code. Origins and termini were identified
590 based on GC skew patterns indicative of bidirectional replication, and unique replichoes are
591 marked in alternating shades of blue.

592 **C-D.** Genomic maps of highly-recoded lysis cassette neighborhoods from representative TAG-
593 recoded crAss-like phages (**C**) and TGA-recoded Agate phages (**D**). Lysis genes (pink) as well as
594 structural genes (purple) that were significantly biased towards use of in-frame recoded stop codons are
595 marked with black striping. In crAss-like phage (C), lytic amidase ($p = 6.82e-3$), spanin ($p=9.00e-6$)
596 and tail tube ($p = 7.59e-4$) gene families preferentially used TAG to encode glutamine (Q) . In
597 Agate phages, a tail gene family ($p = 4.18e-2$) and a transmembrane domain protein family (TM-
598 domain, $p = 4.64e-2$) preferentially use TGA to encode tryptophan (W). **** $p \leq 0.0001$, *** $p \leq$
599 0.001 , ** $p \leq 0.01$, * $p \leq 0.05$, Benjamini-Hochberg p-value corrected two-sided Wilcoxon Rank-
600 Sum Test. This statistical test was used to analyze rates of TAG use relative to standard code
601 encoding of glutamine (TAG \rightarrow Q recoded phage in C) or TGA use relative to the standard code
602 encoding of tryptophan (TGA \rightarrow W recoded phage in D).

603

604 **Fig. 5: Recoded prophages integrated into bacterial genomes.**

605 **A.** A manually curated 24,371 bp TAG-recoded Garnet prophage integrated in a *Prevotella sp.*
606 genome.

607 **B.** A manually curated 23,706 bp TAG-recoded Topaz prophage integrated in a *Oscillospiraceae*
608 *sp.* genome. In both A and B, the bacterial hosts use standard code (black gene predictions).
609 Standard code results in highly fragmented gene predictions in the prophages, due to the high
610 number of genes with in-frame TAGs (orange). In both A and B, the integrase is one of the few
611 prophage genes that does not have in-frame TAG codons (grey). An increase in GC content (blue
612 line) and transition from phage to bacterial gene content marks prophage boundaries. LS = Large
613 subunit, SS = Small Subunit, TMP = Tape Measure Protein.

614

615 **Fig. 6: A model for recoding in the phage life cycle.**

616 Infection of a standard code host begins with the production of proteins from standard code
617 compatible genes. In some phage, this is a route to integrase production and establishment of
618 lysogeny. In other phage, this early phase involves the production of molecules involved in
619 switching from standard to alternative code such as suppressor tRNAs (Sup tRNA), amino acyl
620 tRNA synthetases (aaRS) and release factors (RF1/2). As infection proceeds, recoded gene
621 products initially suppressed by in-frame recoded stop codons code can be produced. This allows
622 for expression of phage structural proteins and ultimately triggers lysis.

623

624

625

626 **References**

- 627 1. Crick, F. H. The origin of the genetic code. *J. Mol. Biol.* **38**, 367–379 (1968).
- 628 2. Knight, R. D., Freeland, S. J. & Landweber, L. F. Rewiring the keyboard: evolvability of the
629 genetic code. *Nat. Rev. Genet.* **2**, 49–58 (2001).
- 630 3. Horowitz, S. & Gorovsky, M. A. An unusual genetic code in nuclear genes of Tetrahymena.
631 *Proc. Natl. Acad. Sci. U. S. A.* **82**, 2452–2455 (1985).
- 632 4. Caron, F. & Meyer, E. Does Paramecium primaurelia use a different genetic code in its
633 macronucleus? *Nature* **314**, 185–188 (1985).
- 634 5. Preer, J. R., Jr, Preer, L. B., Rudman, B. M. & Barnett, A. J. Deviation from the universal
635 code shown by the gene for surface protein 51A in Paramecium. *Nature* **314**, 188–190
636 (1985).
- 637 6. Keeling, P. J. & Doolittle, W. F. A non-canonical genetic code in an early diverging
638 eukaryotic lineage. *EMBO J.* **15**, 2285–2290 (1996).
- 639 7. Schneider, S. U., Leible, M. B. & Yang, X. P. Strong homology between the small subunit of
640 ribulose-1,5-bisphosphate carboxylase/oxygenase of two species of Acetabularia and the
641 occurrence of unusual codon usage. *Mol. Gen. Genet.* **218**, 445–452 (1989).
- 642 8. Santos, M. A., Keith, G. & Tuite, M. F. Non-standard translational events in Candida
643 albicans mediated by an unusual seryl-tRNA with a 5'-CAG-3' (leucine) anticodon. *EMBO J.*
644 **12**, 607–616 (1993).
- 645 9. Ohama, T. *et al.* Non-universal decoding of the leucine codon CUG in several Candida
646 species. *Nucleic Acids Res.* **21**, 4039–4045 (1993).
- 647 10. Inamine, J. M., Ho, K. C., Loechel, S. & Hu, P. C. Evidence that UGA is read as a
648 tryptophan codon rather than as a stop codon by Mycoplasma pneumoniae, Mycoplasma
649 genitalium, and Mycoplasma gallisepticum. *J. Bacteriol.* **172**, 504–506 (1990).
- 650 11. Yamao, F. *et al.* UGA is read as tryptophan in Mycoplasma capricolum. *Proc. Natl. Acad.*

- 651 *Sci. U. S. A.* **82**, 2306–2309 (1985).
- 652 12. Stamburski, C., Renaudin, J. & Bové, J. M. Mutagenesis of a tryptophan codon from TGG
653 to TGA in the *cat* gene does not prevent its expression in the helical mollicute *Spiroplasma*
654 *citri*. *Gene* **110**, 133–134 (1992).
- 655 13. Wrighton, K. C. *et al.* Fermentation, hydrogen, and sulfur metabolism in multiple
656 uncultivated bacterial phyla. *Science* **337**, 1661–1665 (2012).
- 657 14. Campbell, J. H. *et al.* UGA is an additional glycine codon in uncultured SR1 bacteria from
658 the human microbiota. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 5540–5545 (2013).
- 659 15. Hanke, A. *et al.* Recoding of the stop codon UGA to glycine by a BD1-5/SN-2 bacterium
660 and niche partitioning between Alpha- and Gammaproteobacteria in a tidal sediment
661 microbial community naturally selected in a laboratory chemostat. *Front. Microbiol.* **5**, 231
662 (2014).
- 663 16. Shulgina, Y. & Eddy, S. R. A computational screen for alternative genetic codes in over
664 250,000 genomes. *Elife* **10**, e71402 (2021).
- 665 17. Zinoni, F., Birkmann, A., Leinfelder, W. & Böck, A. Cotranslational insertion of
666 selenocysteine into formate dehydrogenase from *Escherichia coli* directed by a UGA
667 codon. *Proc. Natl. Acad. Sci. U. S. A.* **84**, 3156–3160 (1987).
- 668 18. Berry, M. J. *et al.* Recognition of UGA as a selenocysteine codon in type I deiodinase
669 requires sequences in the 3' untranslated region. *Nature* **353**, 273–276 (1991).
- 670 19. Hao, B. *et al.* A new UAG-encoded residue in the structure of a methanogen
671 methyltransferase. *Science* **296**, 1462–1466 (2002).
- 672 20. Sun, J. *et al.* Recoding of stop codons expands the metabolic potential of two novel
673 Asgardarchaeota lineages. *ISME Communications* **1**, 1–14 (2021).
- 674 21. Gomes, A. C. *et al.* A genetic code alteration generates a proteome of high diversity in the
675 human pathogen *Candida albicans*. *Genome Biol.* **8**, R206 (2007).
- 676 22. Ivanova, N. N. *et al.* Stop codon reassignments in the wild. *Science* **344**, 909–913 (2014).

- 677 23. Devoto, A. E. *et al.* Megaphages infect *Prevotella* and variants are widespread in gut
678 microbiomes. *Nat Microbiol* **4**, 693–700 (2019).
- 679 24. Al-Shayeb, B. *et al.* Clades of huge phages from across Earth's ecosystems. *Nature* **578**,
680 425–431 (2020).
- 681 25. Yutin, N. *et al.* Analysis of metagenome-assembled viral genomes from the human gut
682 reveals diverse putative CrAss-like phages with unique genomic features. *Nat. Commun.*
683 **12**, 1–11 (2021).
- 684 26. Crisci, M. A. *et al.* Closely related Lak megaphages replicate in the microbiomes of diverse
685 animals. *iScience* **24**, 102875 (2021).
- 686 27. Roux, S. *et al.* Minimum Information about an Uncultivated Virus Genome (MIUViG). *Nat.*
687 *Biotechnol.* **37**, 29–37 (2019).
- 688 28. Goltsman, D. S. A. *et al.* Metagenomic analysis with strain-level resolution reveals fine-
689 scale variation in the human pregnancy microbiome. *Genome Res.* **28**, 1467–1480 (2018).
- 690 29. Obregon-Tito, A. J. *et al.* Subsistence strategies in traditional societies distinguish gut
691 microbiomes. *Nat. Commun.* **6**, 6505 (2015).
- 692 30. Rampelli, S. *et al.* Metagenome Sequencing of the Hadza Hunter-Gatherer Gut Microbiota.
693 *Curr. Biol.* **25**, 1682–1693 (2015).
- 694 31. Lou, Y. C. *et al.* Infant gut strain persistence is associated with maternal origin, phylogeny,
695 and traits including surface adhesion and iron acquisition. *Cell Rep Med* **2**, 100393 (2021).
- 696 32. David, L. A. *et al.* Gut microbial succession follows acute secretory diarrhea in humans.
697 *MBio* **6**, e00381–15 (2015).
- 698 33. Tung, J. *et al.* Social networks predict gut microbiome composition in wild baboons. *Elife* **4**,
699 e05224 (2015).
- 700 34. Munk, P. *et al.* A sampling and metagenomic sequencing-based methodology for
701 monitoring antimicrobial resistance in swine herds. *J. Antimicrob. Chemother.* **72**, 385–392
702 (2017).

- 703 35. Andersen, V. D. *et al.* Predicting effects of changed antimicrobial usage on the abundance
704 of antimicrobial resistance genes in finisher' gut microbiomes. *Prev. Vet. Med.* **174**, 104853
705 (2020).
- 706 36. Wallace, R. J. *et al.* A heritable subset of the core rumen microbiome dictates dairy cow
707 productivity and emissions. *Sci Adv* **5**, eaav8391 (2019).
- 708 37. Peters, S. L. *et al.* Validation that human microbiome phages use alternative genetic coding
709 with TAG stop read as Q. *bioRxiv* 2022.01.06.475225 (2022)
710 doi:10.1101/2022.01.06.475225.
- 711 38. Osawa, S. & Jukes, T. H. Codon reassignment (codon capture) in evolution. *J. Mol. Evol.*
712 **28**, 271–278 (1989).
- 713 39. Berry, J., Rajaure, M., Pang, T. & Young, R. The spanin complex is essential for lambda
714 lysis. *J. Bacteriol.* **194**, 5667–5674 (2012).
- 715 40. Young, R. Phage lysis: three steps, three choices, one outcome. *J. Microbiol.* **52**, 243–258
716 (2014).
- 717 41. Doermann, A. H. The intracellular growth of bacteriophages. I. Liberation of intracellular
718 bacteriophage T4 by premature lysis with another phage or with cyanide. *J. Gen. Physiol.*
719 **35**, 645–656 (1952).
- 720 42. Heagy, F. C. The effect of 2,4-dinitrophenol and phage T2 on Escherichia coli B. *J.*
721 *Bacteriol.* **59**, 367–373 (1950).
- 722 43. Park, T., Struck, D. K., Dankenbring, C. A. & Young, R. The pinholin of lambdoid phage 21:
723 control of lysis by membrane depolarization. *J. Bacteriol.* **189**, 9135–9139 (2007).
- 724 44. Hays, S. G. & Seed, K. D. Dominant *Vibrio cholerae* phage exhibits lysis inhibition sensitive
725 to disruption by a defensive phage satellite. *Elife* **9**, (2020).
- 726 45. Cowe, E. & Sharp, P. M. Molecular evolution of bacteriophages: Discrete patterns of codon
727 usage in T4 genes are related to the time of gene expression. *J. Mol. Evol.* **33**, 13–22
728 (1991).

- 729 46. Yang, J. Y. *et al.* Degradation of host translational machinery drives tRNA acquisition in
730 viruses. *Cell Syst* (2021) doi:10.1016/j.cels.2021.05.019.
- 731 47. Durmaz, E. & Klaenhammer, T. R. Abortive phage resistance mechanism AbiZ speeds the
732 lysis clock to cause premature lysis of phage-infected *Lactococcus lactis*. *J. Bacteriol.* **189**,
733 1417–1425 (2007).
- 734 48. Zeng, L. *et al.* Decision making at a subcellular level determines the outcome of
735 bacteriophage infection. *Cell* **141**, 682–691 (2010).
- 736 49. Erez, Z. *et al.* Communication between viruses guides lysis-lysogeny decisions. *Nature*
737 **541**, 488–493 (2017).
- 738 50. Silpe, J. E. & Bassler, B. L. A Host-Produced Quorum-Sensing Autoinducer Controls a
739 Phage Lysis-Lysogeny Decision. *Cell* **176**, 268–280.e13 (2019).
- 740 51. Bailly-Bechet, M., Vergassola, M. & Rocha, E. Causes for the intriguing presence of tRNAs
741 in phages. *Genome Res.* **17**, 1486–1495 (2007).
- 742 52. Auslander, N., Gussow, A. B., Benler, S., Wolf, Y. I. & Koonin, E. V. Seeker: alignment-free
743 identification of bacteriophage genomes by deep learning. *Nucleic Acids Res.* **48**, e121–
744 e121 (2020).
- 745 53. Kieft, K., Zhou, Z. & Anantharaman, K. VIBRANT: automated recovery, annotation and
746 curation of microbial viruses, and evaluation of viral community function from genomic
747 sequences. *Microbiome* **8**, 90 (2020).
- 748 54. Nayfach, S. *et al.* CheckV assesses the quality and completeness of metagenome-
749 assembled viral genomes. *Nat. Biotechnol.* **39**, 578–585 (2020).
- 750 55. Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate
751 genomic comparisons that enables improved genome recovery from metagenomes through
752 de-replication. *ISME J.* **11**, 2864–2868 (2017).
- 753 56. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
754 identification. *BMC Bioinformatics* **11**, 119 (2010).

- 755 57. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment
756 search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 757 58. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
- 758 59. Chan, P. P. & Lowe, T. M. tRNAscan-SE: Searching for tRNA Genes in Genomic
759 Sequences. *Methods Mol. Biol.* **1962**, 1–14 (2019).
- 760 60. Skennerton, C. T. *minced: Mining CRISPRs in Environmental Datasets*. (Github).
- 761 61. Buchfink, B., Reuter, K. & Drost, H.-G. Sensitive protein alignments at tree-of-life scale
762 using DIAMOND. *Nat. Methods* **18**, 366–368 (2021).
- 763 62. Olm, M. *tRep: Quick get the taxonomy of a genome*. (Github).
- 764 63. Darling, A. C. E., Mau, B., Blattner, F. R. & Perna, N. T. Mauve: multiple alignment of
765 conserved genomic sequence with rearrangements. *Genome Res.* **14**, 1394–1403 (2004).
- 766 64. Bin Jang, H. *et al.* Taxonomic assignment of uncultivated prokaryotic virus genomes is
767 enabled by gene-sharing networks. *Nat. Biotechnol.* **37**, 632–639 (2019).
- 768 65. Hauser, M., Steinegger, M. & Söding, J. MMseqs software suite for fast and deep clustering
769 and searching of large protein sequence sets. *Bioinformatics* **32**, 1323–1330 (2016).
- 770 66. Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: lightning-fast iterative protein
771 sequence searching by HMM-HMM alignment. *Nat. Methods* **9**, 173–175 (2011).
- 772 67. Söding, J. Protein homology detection by HMM-HMM comparison. *Bioinformatics* **21**, 951–
773 960 (2005).
- 774 68. Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein
775 or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).
- 776 69. Katoh, K., Misawa, K., Kuma, K.-I. & Miyata, T. MAFFT: a novel method for rapid multiple
777 sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066
778 (2002).
- 779 70. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated
780 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973

- 781 (2009).
- 782 71. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and effective
783 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**,
784 268–274 (2015).
- 785 72. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree
786 display and annotation. *Nucleic Acids Res.* **49**, W293–W296 (2021).
- 787 73. Zimmermann, L. *et al.* A Completely Reimplemented MPI Bioinformatics Toolkit with a New
788 HHpred Server at its Core. *J. Mol. Biol.* **430**, 2237–2243 (2018).
- 789 74. Kelley, L. A., Mezulis, S., Yates, C. M., Wass, M. N. & Sternberg, M. J. E. The Phyre2 web
790 portal for protein modeling, prediction and analysis. *Nat. Protoc.* **10**, 845–858 (2015).
- 791 75. Gilchrist, C. L. M. & Chooi, Y.-H. clinker & clustermap.js: Automatic generation of gene
792 cluster comparison figures. *bioRxiv* 2020.11.08.370650 (2020)
793 doi:10.1101/2020.11.08.370650.
- 794 76. Brown, C. T., Olm, M. R., Thomas, B. C. & Banfield, J. F. Measurement of bacterial
795 replication rates in microbial communities. *Nat. Biotechnol.* **34**, 1256–1263 (2016).
- 796 77. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*
797 **9**, 357–359 (2012).

798

Figure 1

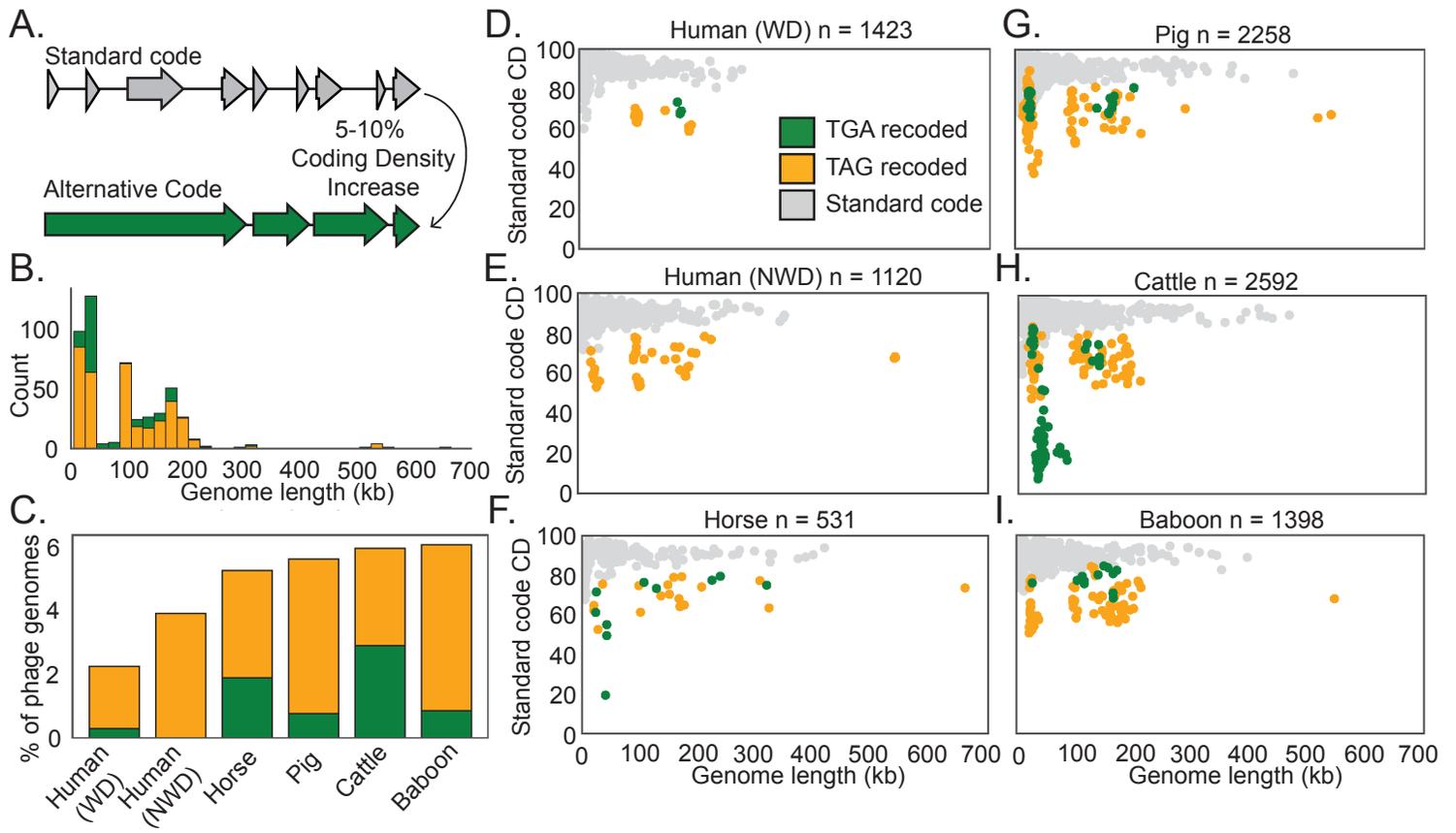


Figure 2

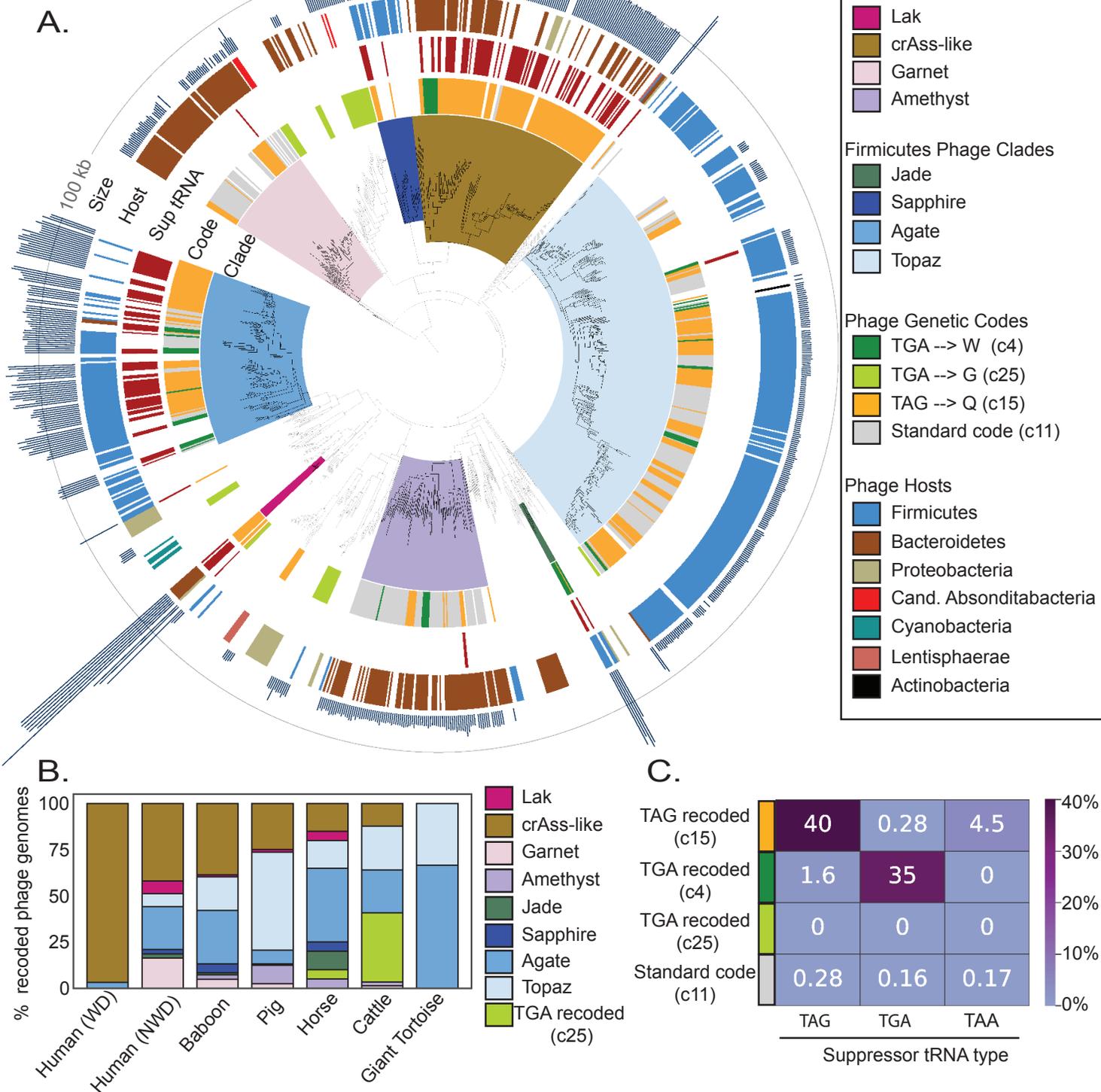
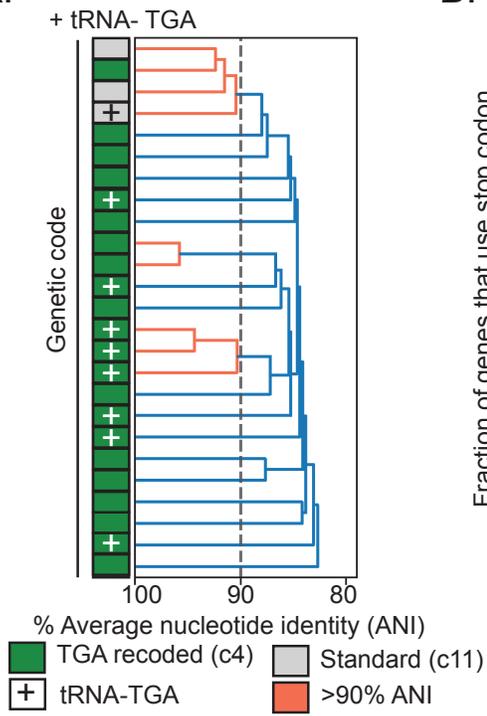
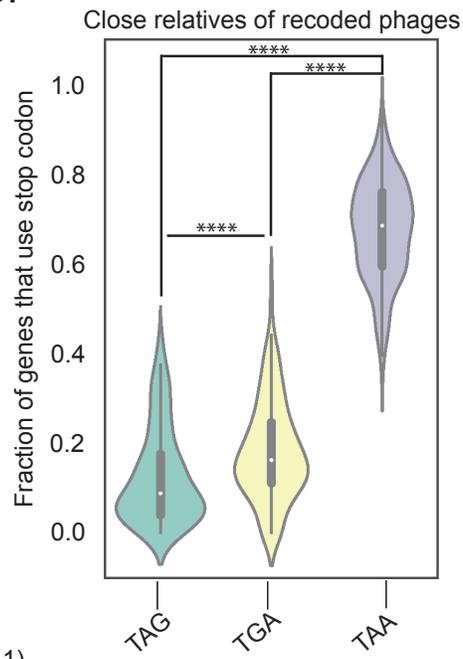


Figure 3

A.



B.



C.

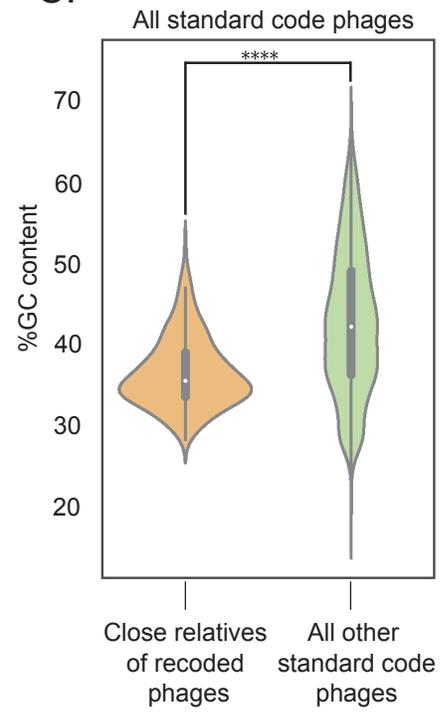


Figure 4

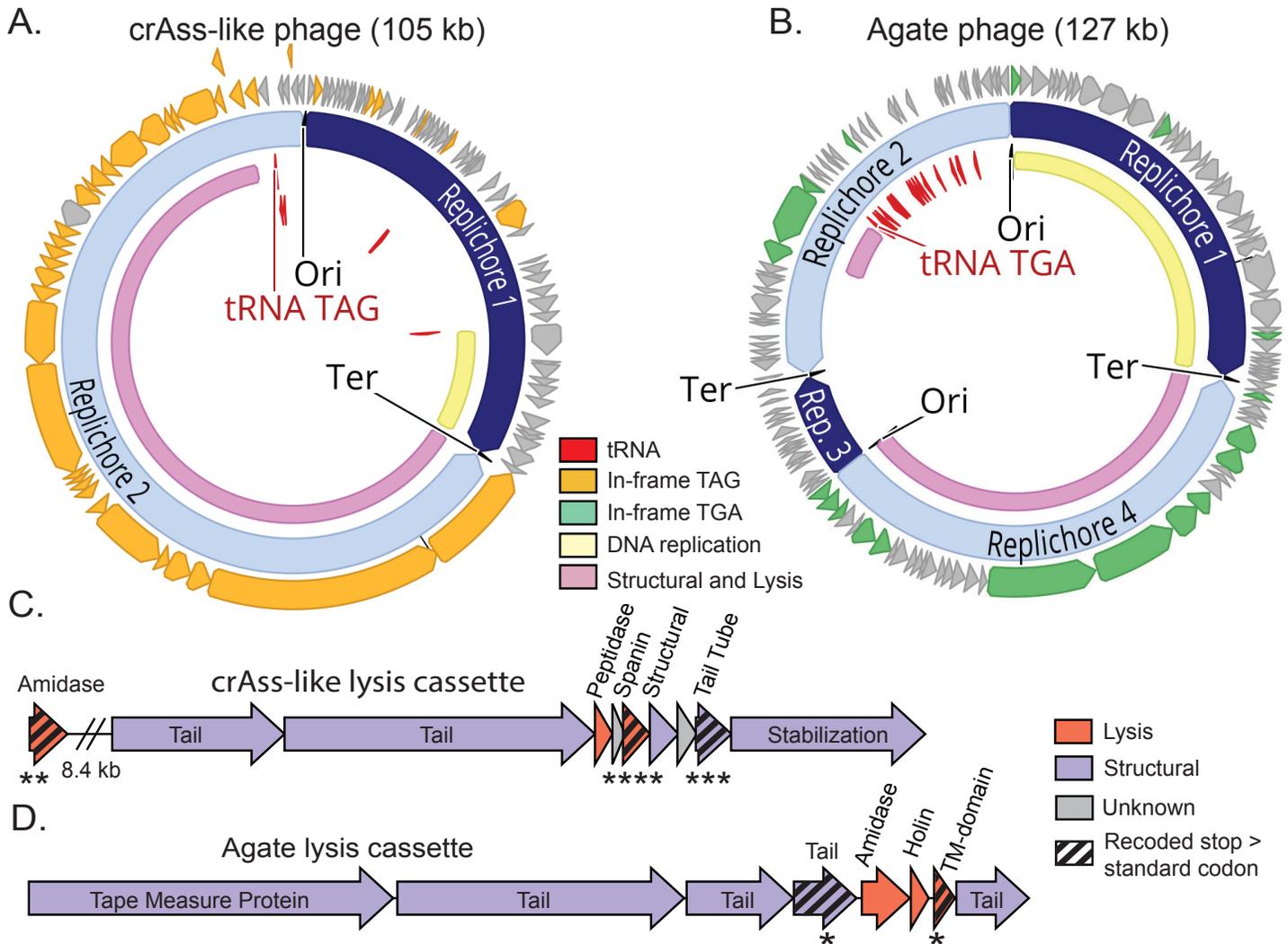


Figure 5.

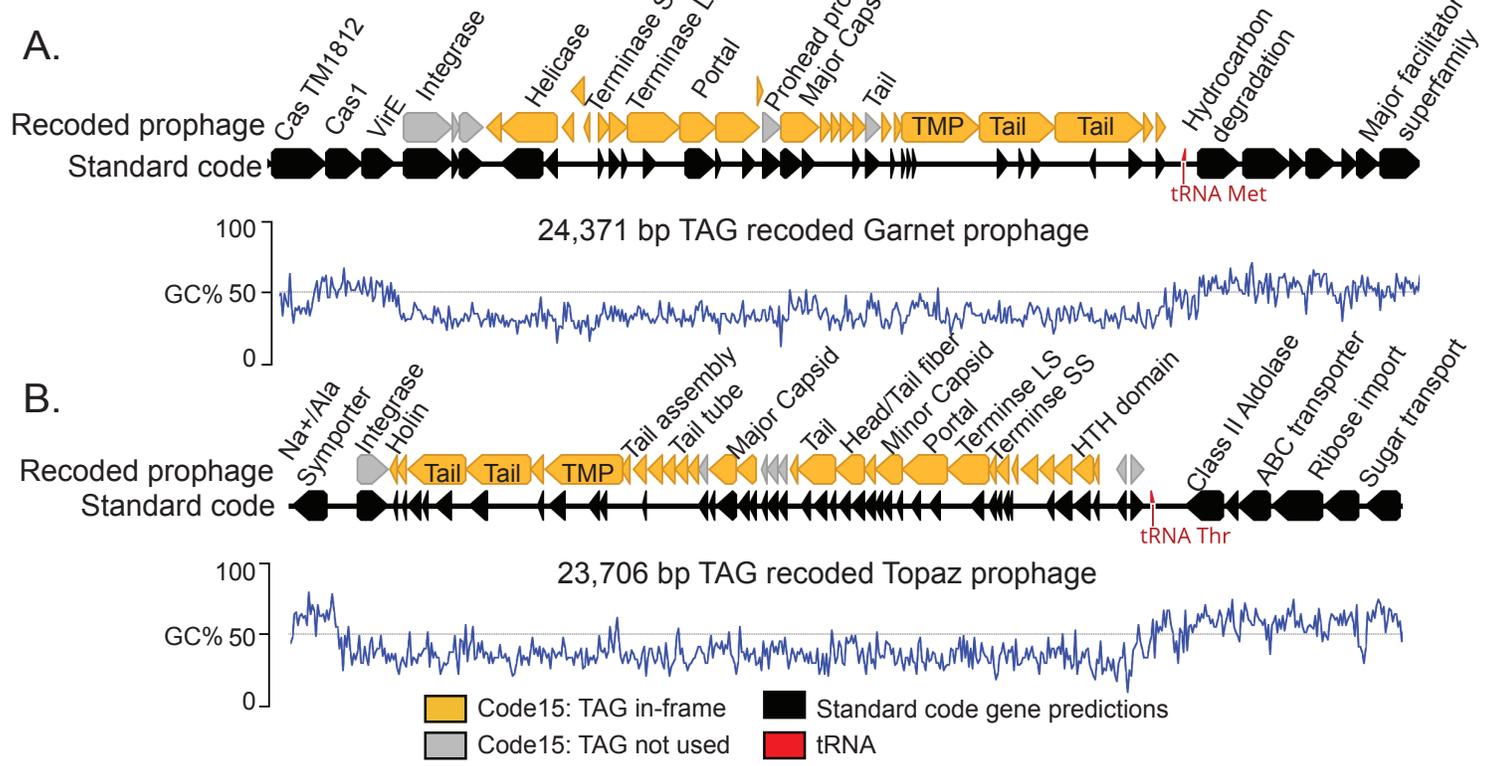
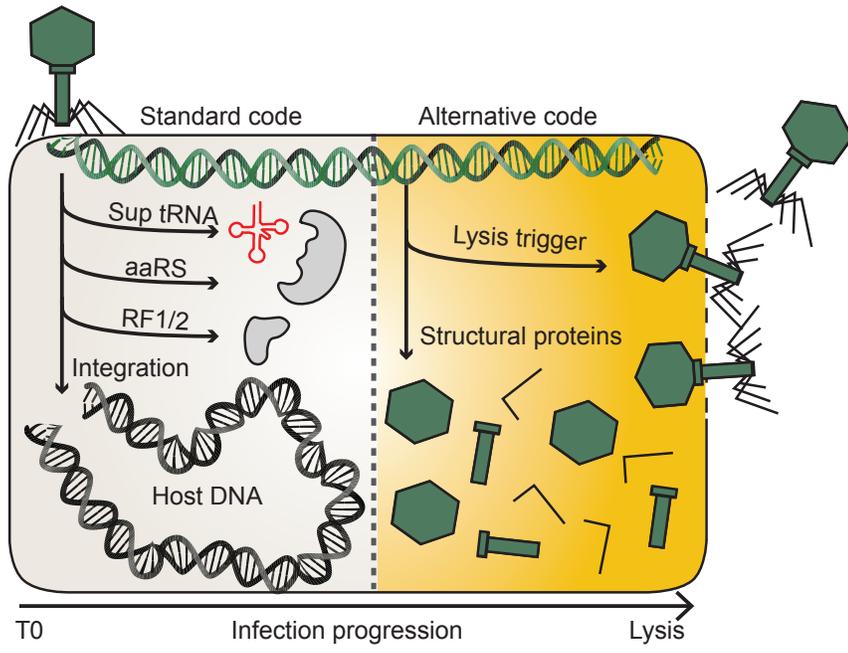
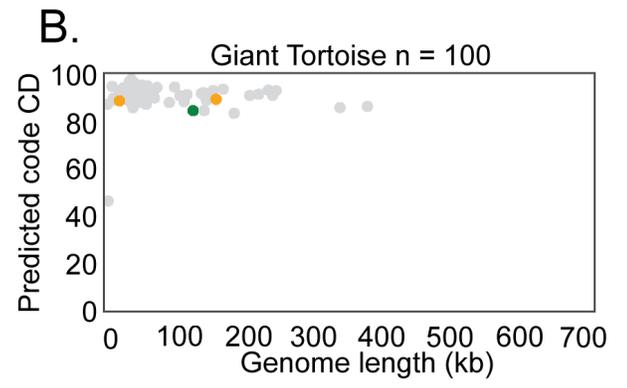
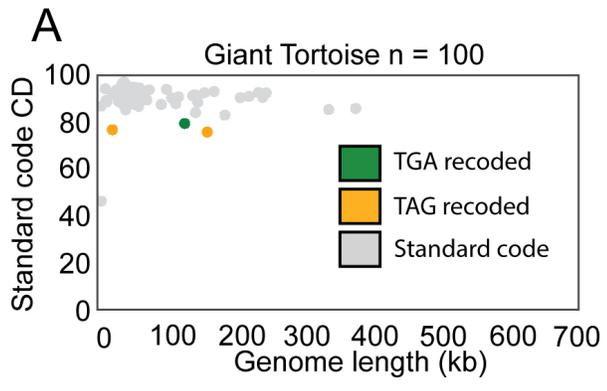


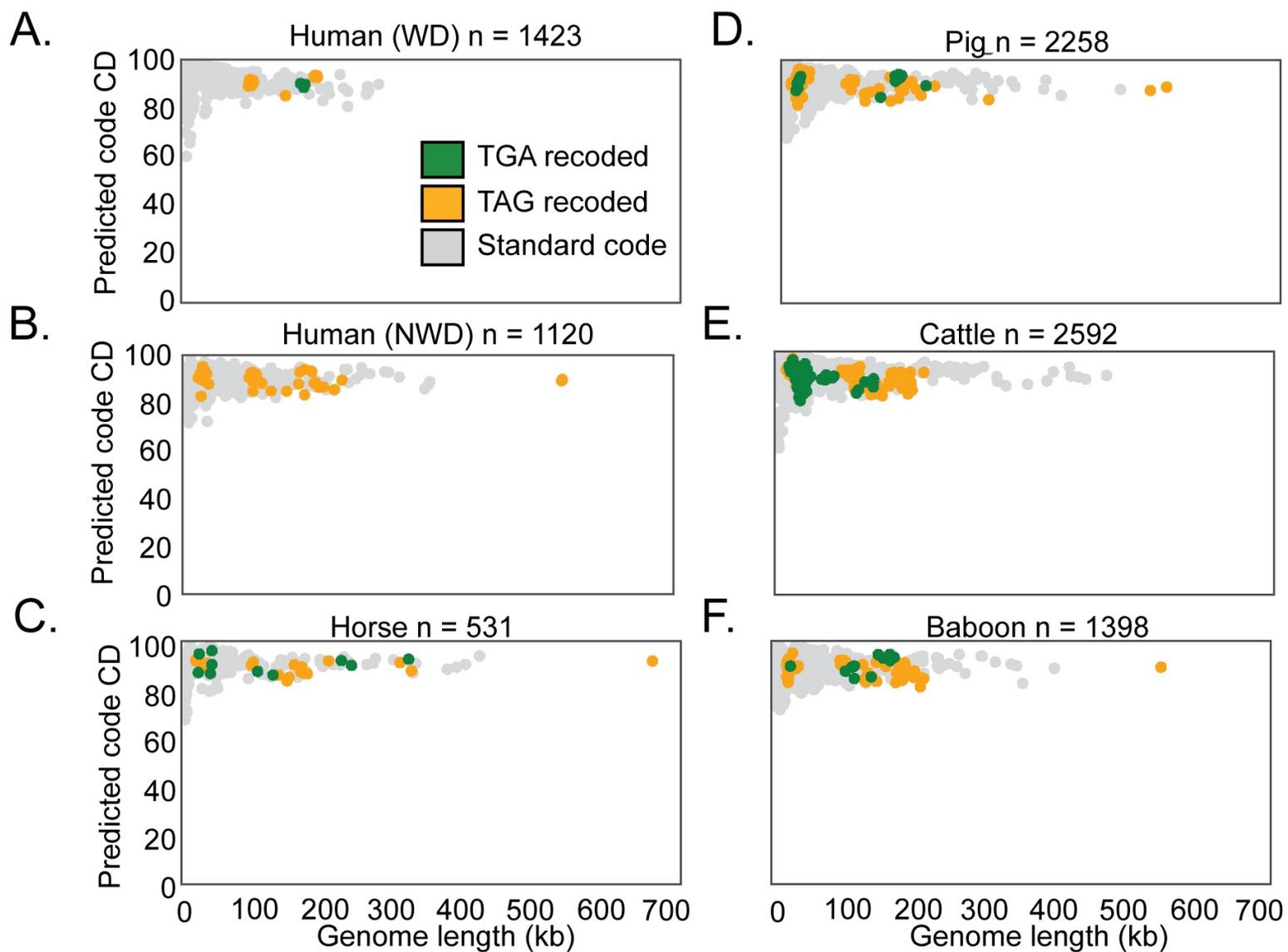
Figure 6



Extended Data Figure 1



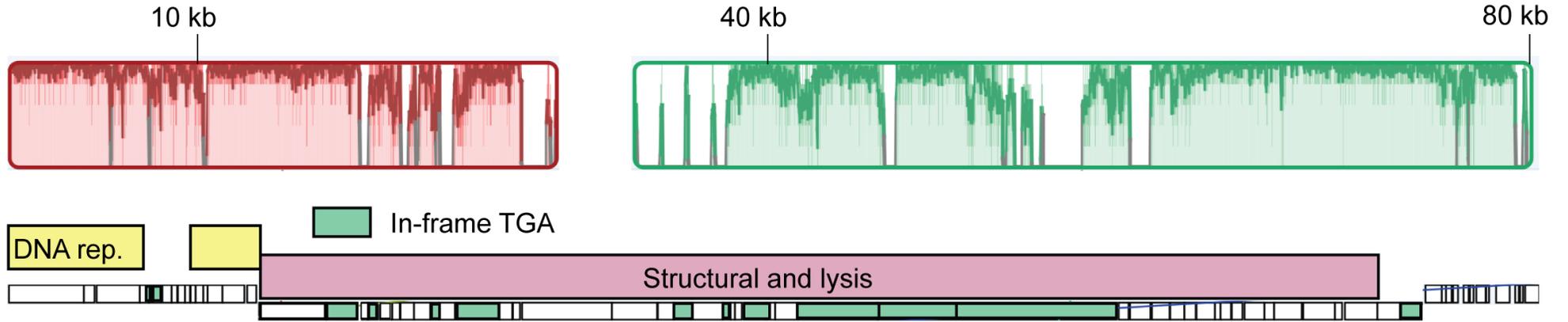
Extended Data Figure 2



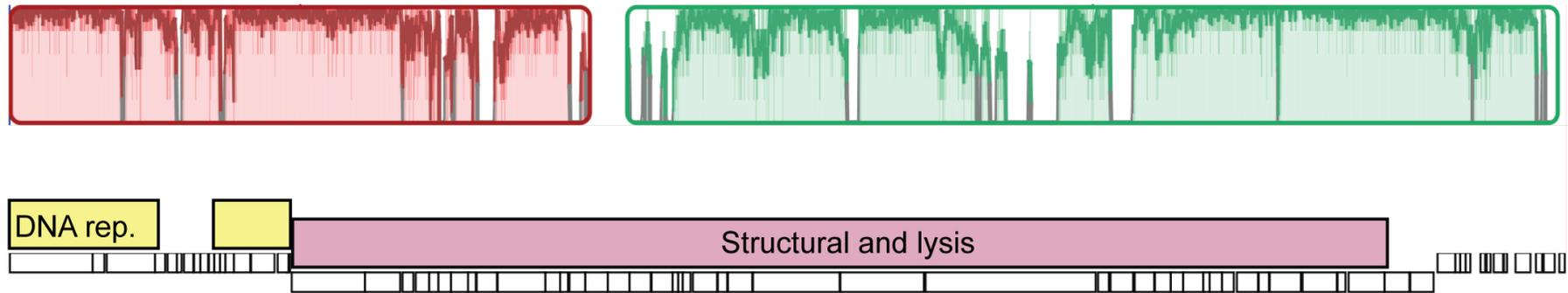
Extended Data Figure 3

A.

TGA recoded Agate phage



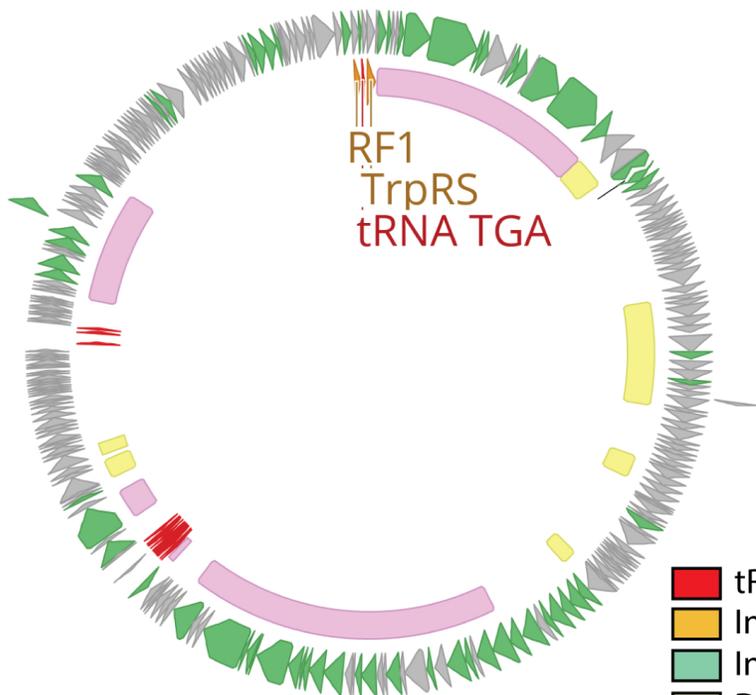
Standard code Agate phage



Extended Data Figure 4

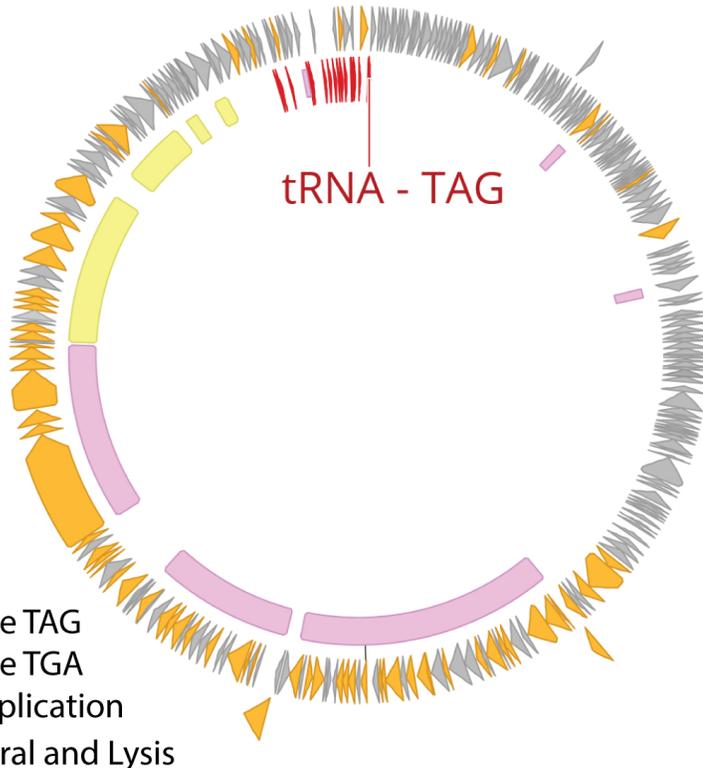
A.

Jade (224 kb)



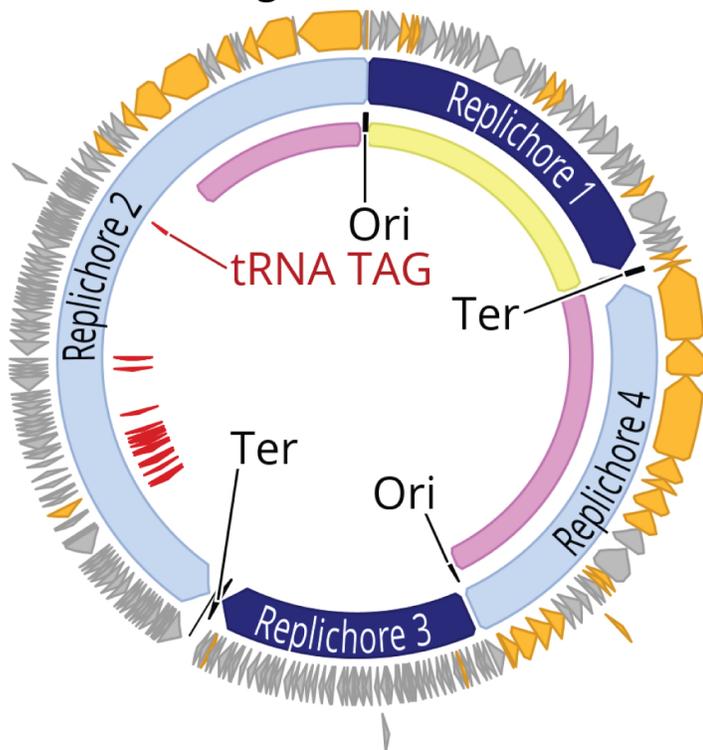
B.

Sapphire (210 kb)



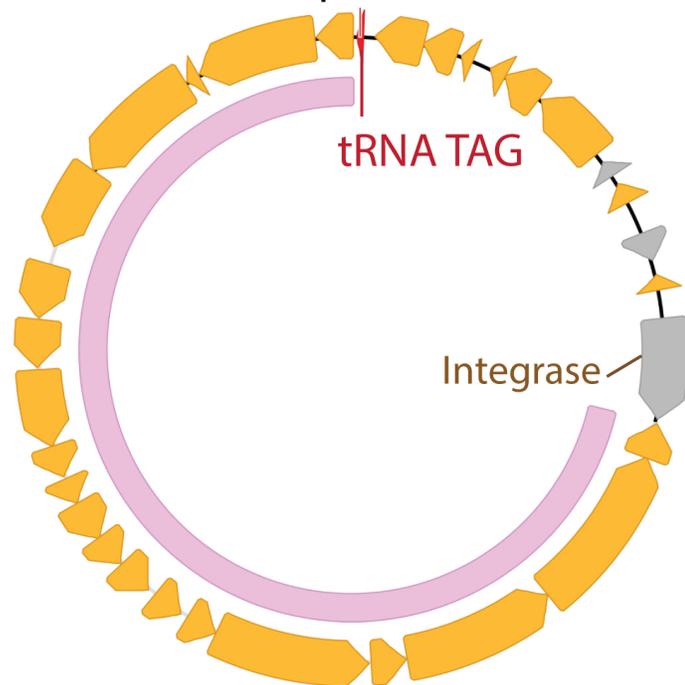
C.

Agate (180kb)



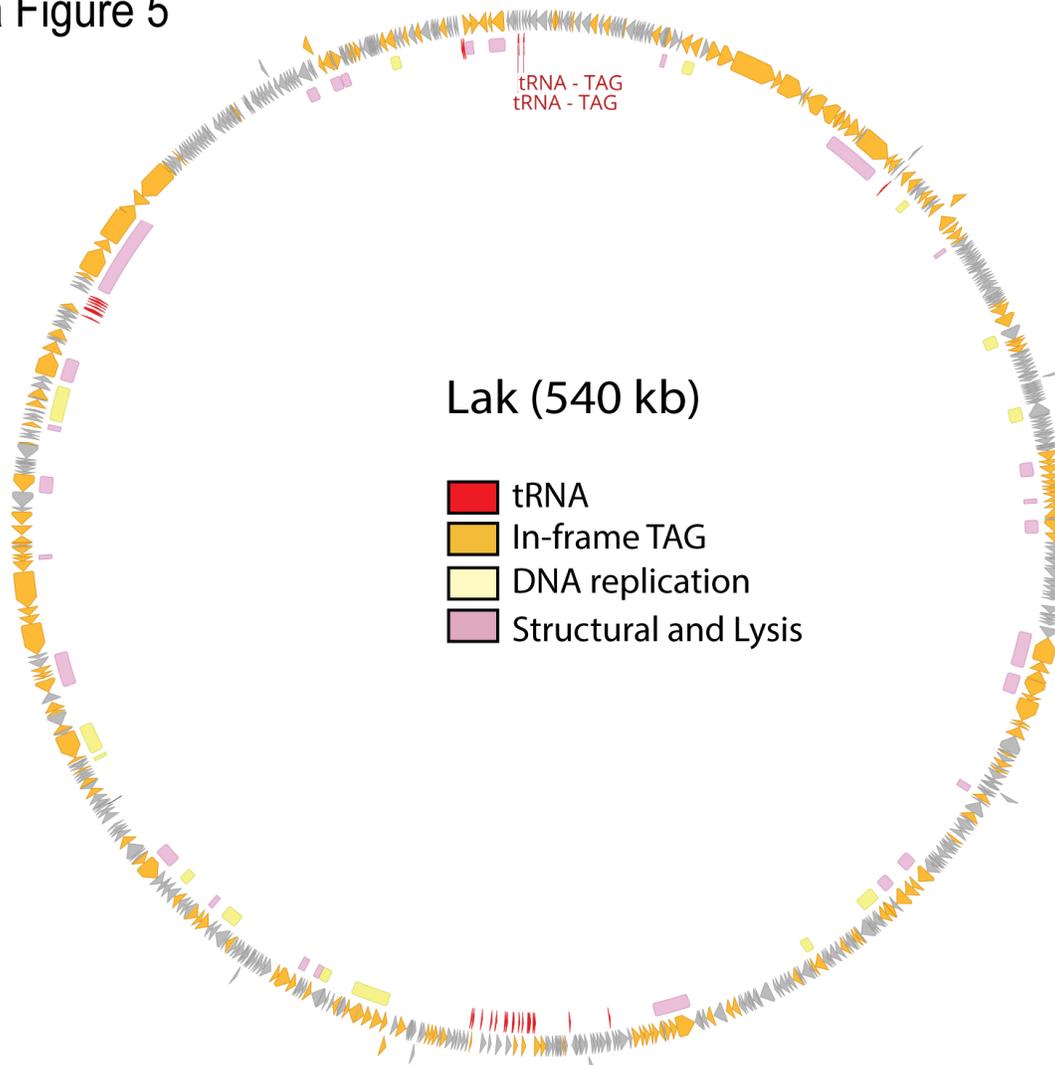
D.

Topaz (20kb)



Extended Data Figure 5

A.

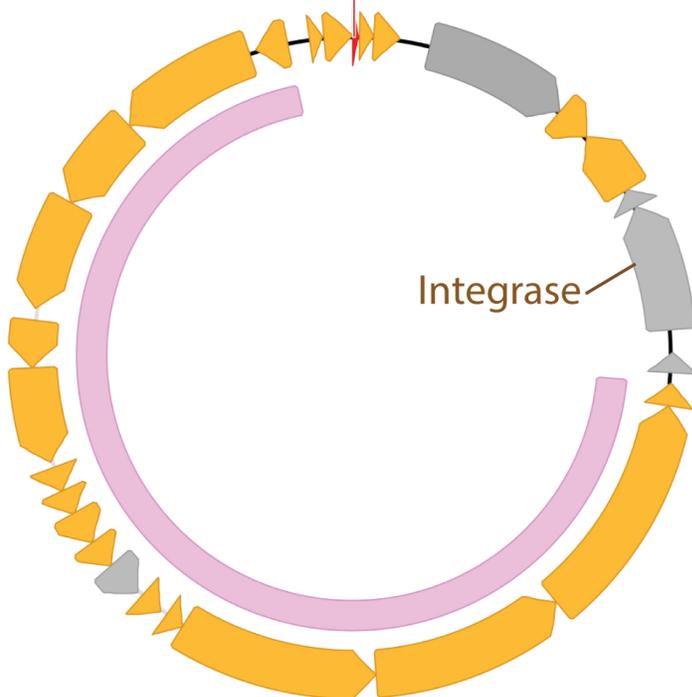


B.

Garnet (25kb)

tRNA TAG

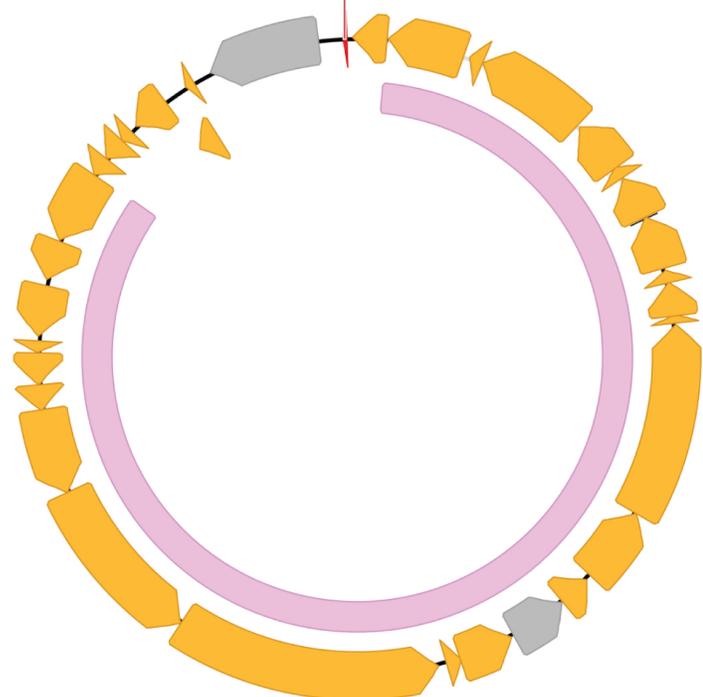
Integrase



C.

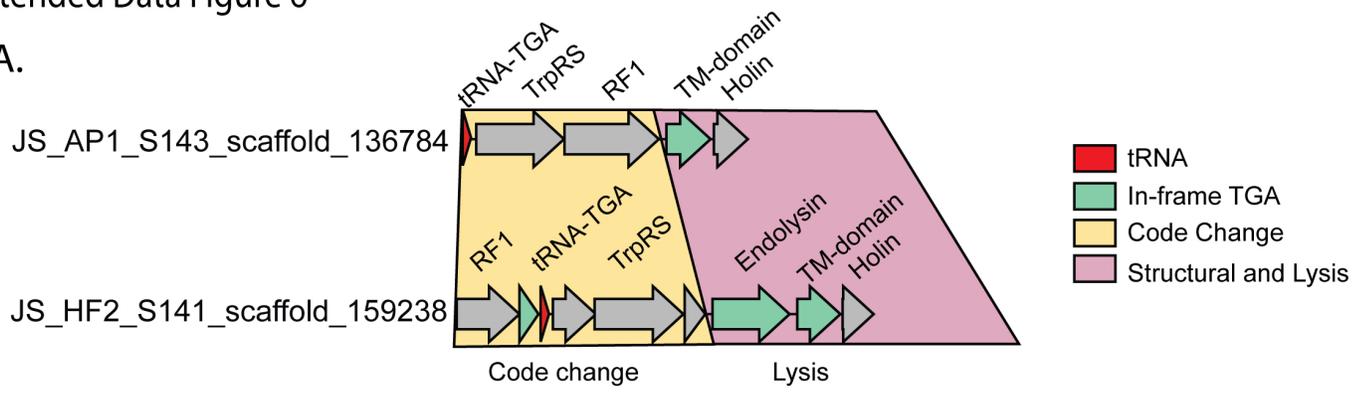
Amethyst (32kb)

tRNA TAG



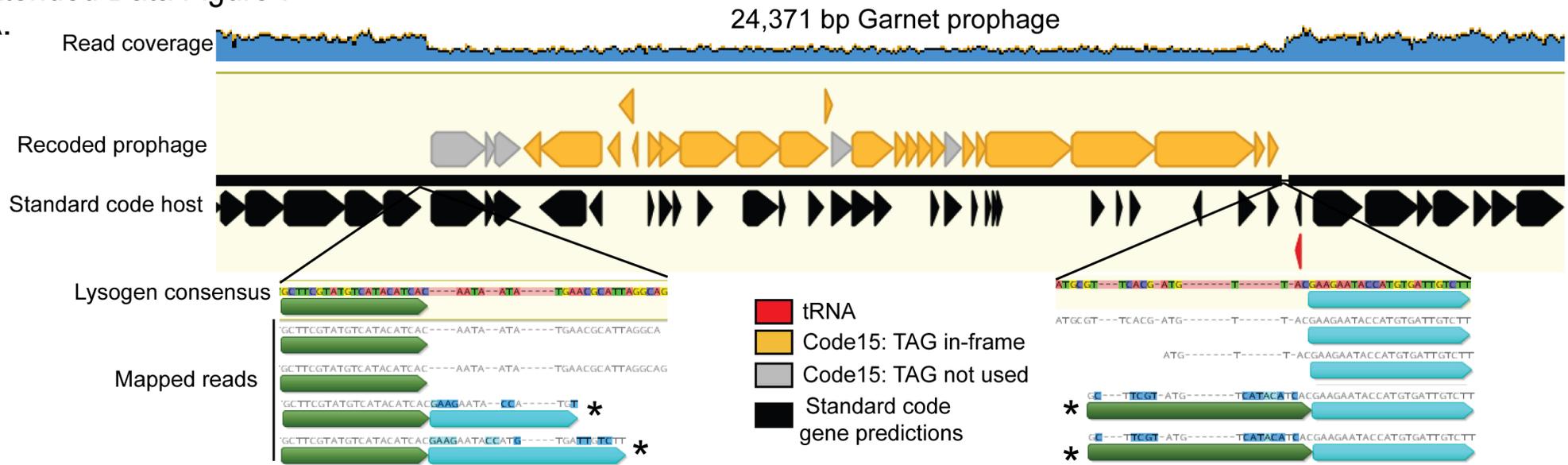
Extended Data Figure 6

A.

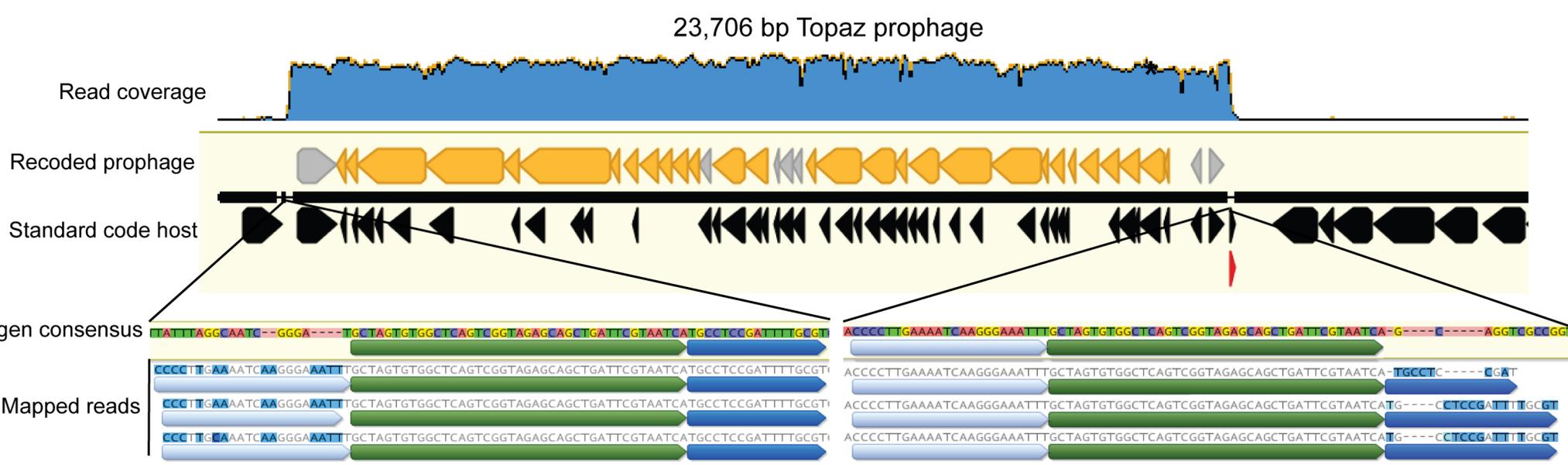


Extended Data Figure 7

A.



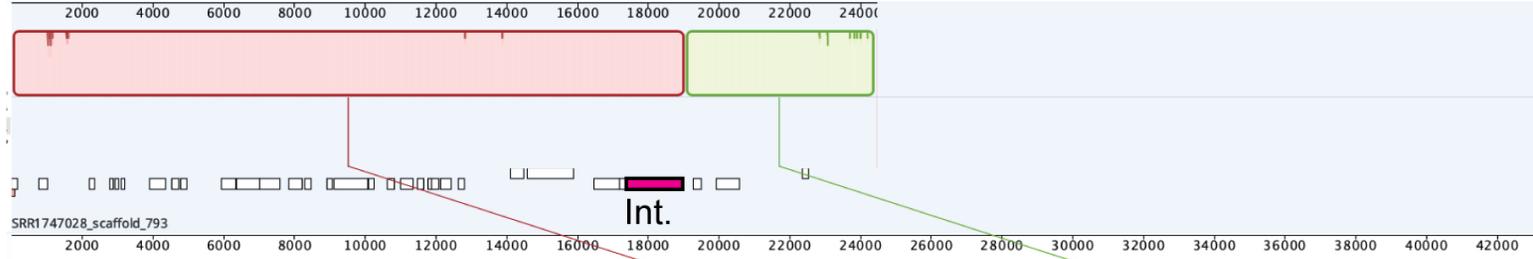
B.



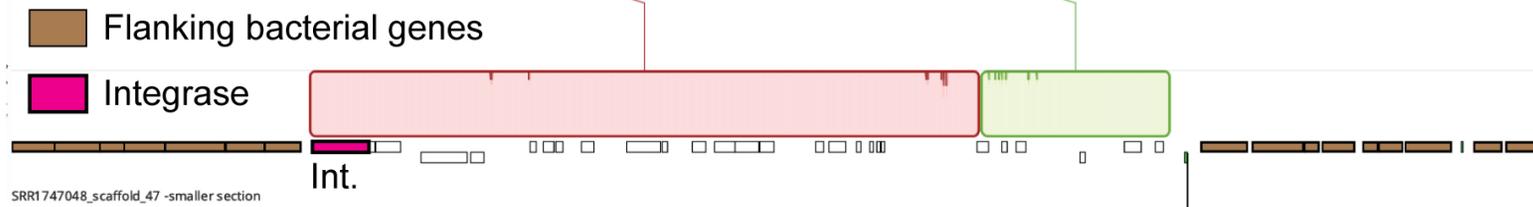
Extended Data Figure 8

A.

25 kb circular Garnet phage
SRR1747028_scaffold_793

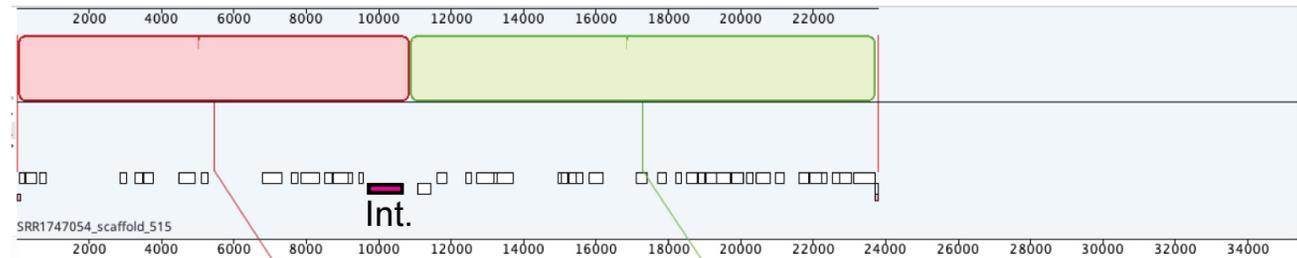


Garnet prophage
SRR1747048_scaffold_47

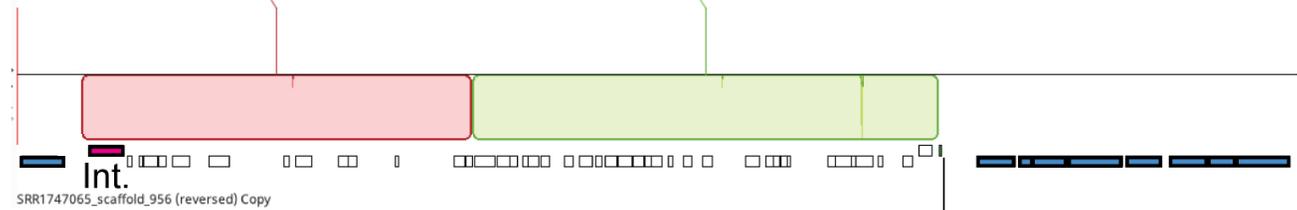


B.

24 kb circular Topaz phage
SRR1747054_scaffold_515



Topaz prophage
SRR1747065_scaffold_956



Flanking bacterial genes

Integrase

tRNA
Met

tRNA
Thr