

1 Predicting stop codon reassignment improves functional 2 annotation of bacteriophages

3

4 **Authors**

5 Ryan Cook^{*1}, Andrea Telatin¹, George Bouras^{2,3}, Antonio Pedro Camargo⁴, Martin Larralde⁵,
6 Robert A. Edwards⁶, and Evelien M. Adriaenssens¹

7

8 * *Denotes corresponding author:* Ryan.Cook@quadram.ac.uk

9

10 **Affiliations**

11 1: Food, Microbiome and Health Research Programme, Quadram Institute Bioscience, Norwich, NR4
12 7UQ, UK

13 2: Adelaide Medical School, Faculty of Health and Medical Sciences, The University of Adelaide,
14 Adelaide, SA 5070, Australia

15 3: Department of Surgery—Otolaryngology Head and Neck Surgery, University of Adelaide and the
16 Basil Hetzel Institute for Translational Health Research, Central Adelaide Local Health Network,
17 Adelaide, SA 5070, Australia

18 4: Department of Energy Joint Genome Institute, Lawrence Berkeley National Laboratory, Berkeley,
19 CA 94720, USA

20 5: Structural and Computational Biology Unit, European Molecular Biology Laboratory (EMBL),
21 Meyerhofstraße 1, 69117 Heidelberg, Germany

22 6: Flinders Accelerator for Microbiome Exploration, College of Science and Engineering, Flinders
23 University, Bedford Park, Adelaide, SA, 5042, Australia

24 **Abstract**

25 The majority of bacteriophage diversity remains uncharacterised, and new intriguing
26 mechanisms of their biology are being continually described. Members of some phage
27 lineages, such as the *Crassvirales*, repurpose stop codons to encode an amino acid by using
28 alternate genetic codes. Here, we investigated the prevalence of stop codon reassignment in
29 phage genomes and subsequent impacts on functional annotation. We predicted 76
30 genomes within INPHARED and 712 vOTUs from the Unified Human Gut Virome catalogue
31 (UHGVS) that repurpose a stop codon to encode an amino acid. We re-annotated these
32 sequences with modified versions of Pharokka and Prokka, called Pharokka-gv and Prokka-
33 gv, to automatically predict stop codon reassignment prior to annotation. Both tools
34 significantly improved the quality of annotations, with Pharokka-gv performing best. For
35 sequences predicted to repurpose TAG to glutamine (translation table 15), Pharokka-gv
36 increased the median gene length (median of per genome medians) from 287 to 481 bp for
37 UHGVS sequences (67.8% increase) and from 318 to 550 bp for INPHARED sequences (72.9%
38 increase). The re-annotation increased mean coding density from 66.8% to 90.0%, and from
39 69.0% to 89.8% for UHGVS and INPHARED sequences. Furthermore, the proportion of genes
40 that could be assigned functional annotation increased, including an increase in the number
41 of major capsid proteins that could be identified. We propose that automatic prediction of
42 stop codon reassignment before annotation is beneficial to downstream viral genomic and
43 metagenomic analyses.

44 Main Body

45 Bacteriophages, hereafter phages, are increasingly recognised as a vital component of
46 microbial communities in all environments where they have been studied in detail. Phages
47 are known to drive bacterial evolution and community composition through predator-prey
48 dynamics and their potential as agents of horizontal gene transfer. The use of viral
49 metagenomics, or viromics, has massively expanded our understanding of global viral
50 diversity and shed light on the ecological roles that phages play.

51

52 Much of the study into viral communities has been conducted on the human gut. Here,
53 viromics has uncovered ecologically important viruses that are difficult to bring into culture
54 using standard laboratory techniques¹, shown potential roles of viruses in disease states²,
55 and allowed for the recovery of enormous phage genomes larger than any brought into
56 culture³. As the majority of phage diversity remains uncharacterised, new and enigmatic
57 diversification mechanisms are being described continually, including the potential use of
58 alternative translation tables.

59

60 Lineage-specific stop codon reassignment has been described previously in
61 bacteriophages^{4,5}, whereby a stop codon is repurposed to encode an amino acid. Notably,
62 annotations of Lak “megaphages” assembled from metagenomes were observed to exhibit
63 unusually low coding density (~70%) when genes are predicted using the standard bacterial,
64 archaeal and plant plastid genetic code (translation table 11)³, much lower than the value
65 observed for most cultured phages of ~90%⁶. The Lak megaphages were predicted to
66 repurpose the TAG stop codon into an as-of-yet unknown amino acid³. More recently,
67 uncultured members of *Crassvirales* have been predicted to repurpose TAG to glutamine
68 (translation table 15), and TGA to tryptophan (translation table 4)⁵, and since then the use of
69 translation table 15 has been experimentally validated in two phages belonging to
70 *Crassvirales*⁷. As this feature may be widespread in human gut viruses, we trained a fork of
71 Prodigal⁸, named prodigal-gv, to predict stop codon reassignment in phages⁹ and
72 implemented in the pyrodigal-gv library to provide efficient Cython bindings to Prodigal-gv
73 with pyrodigal¹⁰. Additionally, the virus discovery tool geNomad incorporates pyrodigal-gv to
74 predict stop codon reassignment for viral sequences identified in metagenomes and

75 viromes⁹. However, the detection of translation table 15 still has limited support in many
76 tools, and the impacts of stop codon reassignment are rarely considered in viral genomics
77 and metagenomics.

78

79 To assess the extent of stop codon reassignment in studied phage genomes and the impacts
80 on functional annotation, we extracted phage genomes from INPHARED⁶ and predicted
81 those using alternative stop codons. We also added high-quality and complete vOTUs from
82 the Unified Human Gut Virome Catalog (UHGV; <https://github.com/snayfach/UHGV>)
83 predicted to use alternative codons. The viral genomes were re-annotated using modified
84 versions of the commonly used annotation pipelines Prokka¹¹, and Pharokka¹² implementing
85 prodigal-gv/pyrodigal-gv for gene prediction (Supplementary Methods). Hereafter, the
86 modified versions are referred to Prokka-gv and Pharokka-gv.

87

88 From INPHARED, 49 genomes (0.24%) were predicted to use translation table 15, and 27
89 (0.13%) were predicted to use translation table 4. From the UHGV, 666 vOTUs (1.2%) were
90 predicted to use translation table 15 and 46 (0.08%) were predicted to use translation table
91 4. These genomes and vOTUs were not constrained to one particular clade of viruses, being
92 predicted to occur on both dsDNA viruses of the realm *Duplodnaviria* and ssDNA viruses of
93 the realm *Monodnaviria*, suggesting it is a phenomenon that has arisen on at least two
94 occasions (Supplementary Table 1). The lower frequency of these genomes in cultured
95 isolates (INPHARED) versus human viromes (UHGV) may be due to culturing and sequencing
96 biases, perhaps including modifications to DNA that are known to be recalcitrant to
97 sequencing.

98

99 Although the mechanism for stop codon reassignment in phages is not fully understood,
100 suppressor tRNAs are suggested to play a role^{4,13}. Consistent with previous findings, we
101 found 375/715 (52.4%) phages predicted to use translation table 15 encoded at least one
102 suppressor tRNA corresponding to the *amber* stop codon (Sup-CTA tRNA), and 11/73 (15.1%)
103 of those predicted to use translation table 4 encoded at least one suppressor tRNA
104 corresponding to the opal stop codon (Sup-TCA tRNA)^{4,13,14}. Although fewer of those
105 predicted to use translation table 4 encoded the relevant suppressor tRNA, 22/27 (81%) of
106 the INPHARED phages predicted to use translation table 4 were viruses of *Mycoplasma* or

107 *Spiroplasma*. As *Mycoplasma* and *Spiroplasma* are known to use translation table 4, many of
108 the viruses predicted to use translation table 4 may be simply using the same translation
109 table as their host.

110

111 Prediction of stop codon reassignment led to improved annotations for both Prokka and
112 Pharokka, although the extent of this varied with the two datasets, translation tables, and
113 annotation pipelines tested. As Pharokka-gv outperformed Prokka-gv on all metrics tested,
114 only Pharokka-gv is discussed further, and the equivalent results for Prokka-gv can be found
115 in Supplementary Results.

116

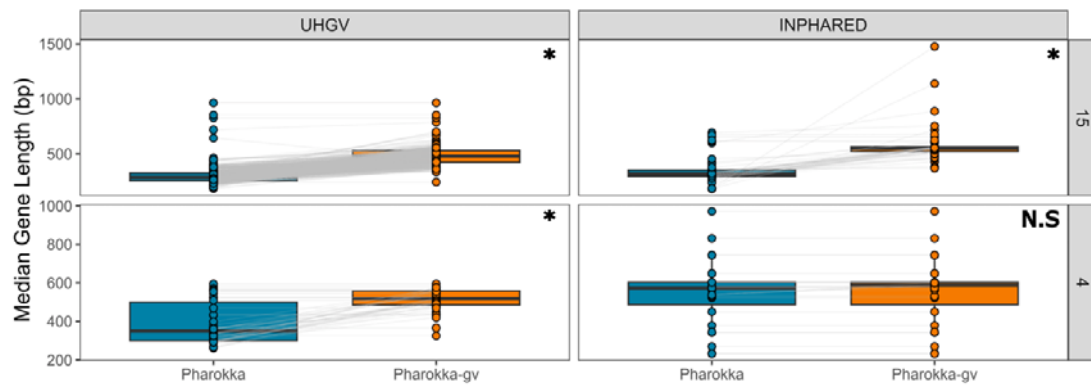
117 The largest differences were observed for sequences predicted to use translation table 15,
118 for which Pharokka-gv increased the median gene length (median of per genome medians)
119 from 287 to 481 bp for UHGV sequences (67.8% increase) and from 318 to 550 bp for
120 INPHARED sequences (72.9% increase; Figure 1A). This was also reflected in an increase of
121 median coding capacity from 66.8% to 90.0% for UHGV, and 69.0% to 89.8% for INPHARED
122 (Figure 1B). Overall, these improved gene calls led to an increased gene length, and a
123 reduction in the number of predicted genes per kb and the number of genes that could not
124 be assigned functional annotations (Supplementary Figure 2; Supplementary Table 2). As it is
125 commonly used as a phylogenetic marker for bacteriophages, we investigated how
126 commonly the major capsid protein (MCP) could be identified with and without predicted
127 stop codon reassignment¹⁵. For those viruses we predicted to use translation table 15,
128 annotation using the default translation table 11 only resulted in the MCP being identified in
129 407/715 (56.9%) of the genomes. In contrast, using translation table 15 with Pharokka-gv,
130 we could identify the MCP in 475/715 (66.4%).

131

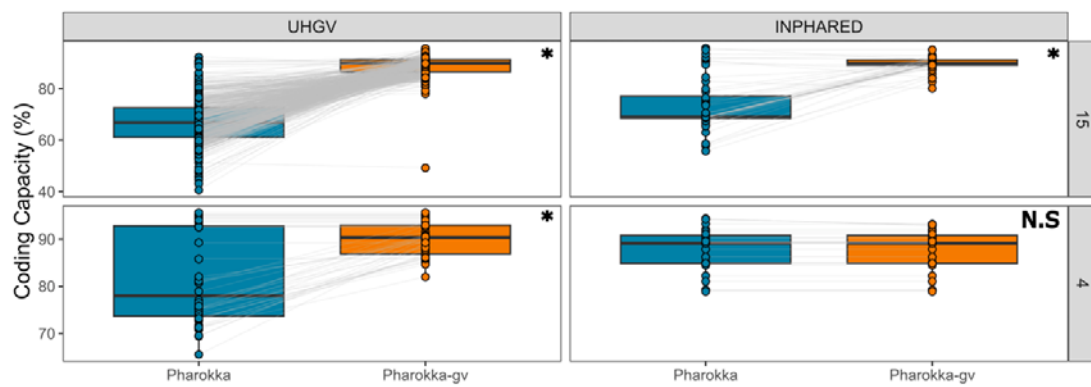
132 When investigating the sequences for which translation table 4 was predicted to be optimal,
133 a substantial increase was also observed for UHGV sequences, with Pharokka-gv increasing
134 median gene length (median of per genome medians) from 350 to 518 bp (a 48.0% increase
135 in length; Figure 1A), resulting in an increase of coding capacity from 78.0% to 90.4% (Figure
136 1B). However, the same was not observed for the 27 INPHARED genomes predicted to use
137 translation table 4. Reannotation resulted in a modest increase in median gene length
138 (median of per genome medians) from 573 to 588 bp (a 2.6% increase in length; Figure 1A).

139 Median coding capacity was not increased, with both Pharokka and Pharokka-gv obtaining
140 89.1% (Figure 1B). As the median gene length and coding capacity for INPHARED sequences
141 predicted to use translation table 4 are in line with expected values, their prediction may be
142 a false positive. Reassuringly, the prediction of translation table 4 has not hindered the
143 quality of annotations where it may be a false positive.

A



B



144

145

146 **Figure 1.** Re-annotating with predicted stop codon reassignment increases the quality of
147 annotations. Comparison of (A) median predicted gene length (bp) and (B) coding capacity
148 (%) for INPHARED genomes and UHGV vOTUs annotated with Pharokka (translation table 11
149 only) and Pharokka-gv (prediction of stop codon reassignment), grouped by dataset and
150 predicted stop codon reassignment. Asterisk indicates significance at $P \leq 10e-10$ with P
151 determined by a simple T test and adjusted with the Benjamini-Hochberg procedure.

152 The analysis of viral (meta)genomes relies on accurate protein predictions, with predicted
153 ORFs being used in common analyses, including (pro)phage prediction, functional
154 annotation, and phylogenetic analyses. The clear differences in protein predictions
155 with/without predicted stop codon reassignment will likely have downstream impacts upon
156 these analyses. However, this phenomenon is not yet widely considered in viral
157 (meta)genomics. We have demonstrated the impacts of stop codon reassignment in the
158 functional annotation of phages, and provide tools for the automatic prediction and
159 annotation of viral genomes that repurpose stop codons. Our analysis highlights the need for
160 accurate viral ORF prediction, and further experimental validation to elucidate the
161 mechanisms of stop codon reassignment.

162 **Data Availability**

163 The genomes used in this analysis are from two publicly available datasets; INPHARED
164 (<https://github.com/RyanCook94/inphared>) and the Unified Human Gut Virome (UHG
165 <https://github.com/snayfach/UHGV>). The details of included sequences are shown in
166 Supplementary Table 1. The code for Prokka-gv is available on GitHub
167 (<https://github.com/telatin/metaprokka>). The code for Pharokka is available on GitHub
168 (<https://github.com/gbouras13/pharokka>). The code for Prodigal-gv is available on GitHub
169 (<https://github.com/apcamargo/prodigal-gv>). The code for Pyrodigal-gv is available on
170 GitHub (<https://github.com/althonos/pyrodigal-gv>).

171

172 **Competing Interests**

173 The authors have nothing to declare.

174

175 **Funding**

176 This research was supported by the BBSRC Institute Strategic Programme Food Microbiome
177 and Health BB/X011054/1 and its constituent projects BBS/E/F/000PR13631 and
178 BBS/E/F/000PR13633; and by the BBSRC Institute Strategic Programme Microbes and Food
179 Safety BB/X011011/1 and its constituent projects BBS/E/F/000PR13634,
180 BBS/E/F/000PR13635 and BBS/E/F/000PR13636. R.C and E.M.A were supported by the
181 BBSRC grant Bacteriophages in Gut Health BB/W015706/1. This research was supported in
182 part by the NBI Research Computing through the High-Performance Computing cluster. We
183 gratefully acknowledge CLIMB-BIG-DATA infrastructure (MR/T030062/1) support for the
184 provision of cloud resources. RAE was supported by an award from the
185 NIH NIDDK RC2DK116713 and an award from the Australian Research
186 Council DP220102915. The work conducted by the US Department of Energy Joint Genome
187 Institute (<https://ror.org/04xm1d337>) and the National Energy Research Scientific
188 Computing Center (<https://ror.org/05v3mvq14>) is supported by the US Department of
189 Energy Office of Science user facilities, operated under contract no. DE-AC02-05CH11231.

190 References

- 191 1 Dutilh, B. E. *et al.* in *Nature Communications* Vol. 5 4498 (2014).
- 192 2 Clooney, A. G. *et al.* in *Cell Host & Microbe* Vol. 26 764-778.e765 (2019).
- 193 3 Devoto, A. E. *et al.* in *Nature Microbiology* (2019).
- 194 4 Ivanova, N. N. *et al.* Stop codon reassignments in the wild. *Science* **344**, 909-913 (2014).
- 195 <https://doi.org/10.1126/science.1250691>
- 196 5 Yutin, N. *et al.* Analysis of metagenome-assembled viral genomes from the human gut reveals
- 197 diverse putative CrAss-like phages with unique genomic features. *Nat Commun* **12**, 1044 (2021).
- 198 <https://doi.org/10.1038/s41467-021-21350-w>
- 199 6 Cook, R. *et al.* in *Phage* Vol. 2 214-223 (Cold Spring Harbor Laboratory, 2021).
- 200 7 Peters, S. L. *et al.* Experimental validation that human microbiome phages use alternative genetic
- 201 coding. *Nature Communications* **13**, 5710 (2022). <https://doi.org/10.1038/s41467-022-32979-6>
- 202 8 Hyatt, D. *et al.* in *BMC Bioinformatics* Vol. 11 1-11 (BioMed Central, 2010).
- 203 9 Camargo, A. P. *et al.* Identification of mobile genetic elements with geNomad. *Nat Biotechnol*
- 204 (2023). <https://doi.org/10.1038/s41587-023-01953-y>
- 205 10 Larralde, M. Pyrodigal: Python bindings and interface to Prodigal, an efficient method for gene
- 206 prediction in prokaryotes. *Journal of Open Source Software* **7**, 4296 (2022).
- 207 <https://doi.org/10.21105/joss.04296>
- 208 11 Seemann, T. in *Bioinformatics* Vol. 30 2068-2069 (2014).
- 209 12 Bouras, G. *et al.* Pharokka: a fast scalable bacteriophage annotation tool. *Bioinformatics* **39**
- 210 (2022). <https://doi.org/10.1093/bioinformatics/btac776>
- 211 13 Pfennig, A., Lomsadze, A. & Borodovsky, M. Annotation of Phage Genomes with Multiple Genetic
- 212 Codes. *bioRxiv*, 2022.2006.2029.495998 (2022). <https://doi.org/10.1101/2022.06.29.495998>
- 213 14 Chan, P. P. & Lowe, T. M. tRNAscan-SE: Searching for tRNA Genes in Genomic Sequences.
- 214 *Methods Mol Biol* **1962**, 1-14 (2019). https://doi.org/10.1007/978-1-4939-9173-0_1
- 215 15 Simmonds, P. *et al.* Four principles to establish a universal virus taxonomy. *PLOS Biology* **21**,
- 216 e3001922 (2023). <https://doi.org/10.1371/journal.pbio.3001922>
- 217 16 Telatin, A., Fariselli, P. & Birolo, G. SeqFu: A Suite of Utilities for the Robust and Reproducible
- 218 Manipulation of Sequence Files. *Bioengineering* **8**, 59 (2021).
- 219 17 Terzian, P. *et al.* in *NAR Genomics and Bioinformatics* Vol. 3 (Oxford Academic, 2021).
- 220 18 Team, R. C. *R: A language and environment for statistical computing.* (R Foundation for Statistical
- 221 Computing, 2018).
- 222 19 Benjamini, Y. & Hochberg, Y. in *Journal of the Royal Statistical Society: Series B (Methodological)*
- 223 Vol. 57 289-300 (John Wiley & Sons, Ltd, 1995).
- 224 20 Wickham, H. *Ggplot2: Elegant graphics for data analysis.* 2 edn, (Springer International
- 225 Publishing, 2016).

227 **Supplementary Methods**

228

229 **Datasets**

230 A multifasta file of phage genomes was downloaded from INPHARED
231 (<https://github.com/RyanCook94/inphared>; September 2023)⁶. Stop codon reassignment of
232 INPHARED genomes was predicted using Prodigal-gv v2.11.0
233 (<https://github.com/apcamargo/prodigal-gv>), a fork of Prodigal written to improve viral gene
234 calling⁸. Those predicted to use translation table 4 or 15 were retained for downstream
235 analysis.

236

237 The Unified Human Gut Virome Catalog (UHGV) was filtered for high quality and complete
238 vOTUs deemed to be a “high confidence” virus and predicted to use either translation table
239 4 or 15 (<https://github.com/snayfach/UHGV>). Stop codon reassignment had already been
240 predicted for UHGV vOTUs using Prodigal-gv and is available in the UHGV metadata.

241

242 **Prokka**

243 A fork of Prokka v1.14.5¹¹ was written that incorporates an initial stage of ORF prediction
244 using Prodigal-gv v2.11.0 (<https://github.com/apcamargo/prodigal-gv>)⁸. A first gene calling
245 step is used to infer the genetic code most likely adopted by the genome, then the predicted
246 genetic code is used to perform the translation FASTX::Seq, which we updated to accept
247 code 15 (metacpan.org/pod/FASTX::Seq)¹⁶. The code for this is available at
248 (github.com/telatin/metaprokka). We included publicly available HMMs of the PHROGs
249 database in our Prokka-gv annotations
250 (http://s3.climb.ac.uk/ADM_share/all_phrogs.hmm.gz)¹⁷. The fork is installable from
251 Bioconda as ‘metaprokka’.

252

253 **Pharokka**

254 Pharokka v1.5.0¹² was updated to include support for pyrodigal-gv implementing pyrodigal-
255 gv as a gene predictor. This is specified by using ‘-g prodigal-gv’ when running Pharokka. The
256 updated code is available on GitHub (<https://github.com/gbouras13/pharokka>). Pharokka
257 uses tRNAscan-SE for predicting tRNAs¹⁴.

258

259 **Statistical Analyses and Data Visualisation**

260 To test for significance in differences of results, a simple paired T test was performed in R
261 v4.2.2¹⁸ and P-values were adjusted using the Benjamini-Hochberg procedure¹⁹. Figure 1 was
262 produced using ggplot2 v3.4.2²⁰.

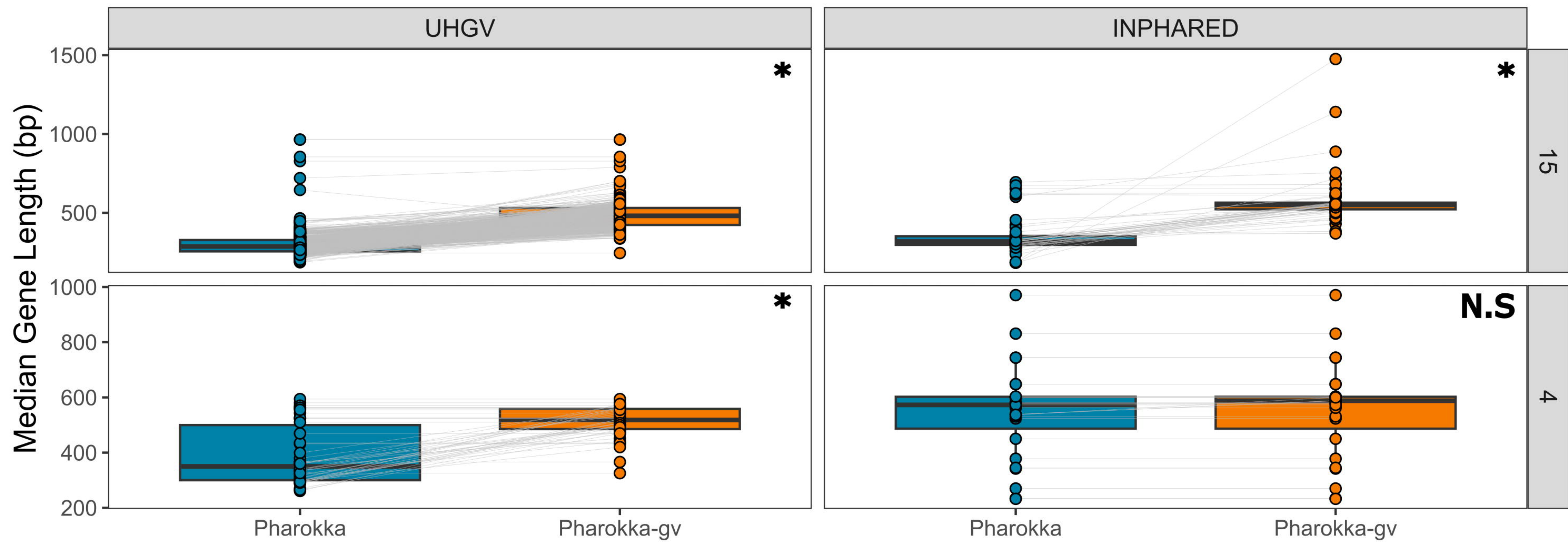
263 **Supplementary Results**

264 **Prokka-gv Annotations**

265 For Prokka-gv, the largest differences were observed for sequences predicted to use
266 translation table 15, for which Prokka-gv increased the median gene length (median of per
267 genome medians) from 276 to 396 bp for UHGV sequences (43.5% increase), and from 309
268 to 483 bp for INPHARED sequences (56.3% increase). This was also reflected in an increase
269 of median coding capacity from 66.6% to 86.7% for UHGV, and from 69.2% to 87.3% for
270 INPHARED. As it is commonly used as a phylogenetic marker for bacteriophages, we
271 investigated how commonly the major capsid protein (MCP) could be identified with and
272 without predicted stop codon reassignment¹⁵. For sequences predicted to use translation
273 table 15, the MCP could be identified on 382/715 (53.4%) sequences with Prokka and this
274 was marginally increased to 386/715 (53.9%) with Prokka-gv.

275

276 When investigating the sequences for which translation table 4 was predicted, a substantial
277 increase was also observed for UHGV sequences, with Prokka-gv increasing median median
278 gene length from 319 to 460 bp (44.2%), resulting in an increase of coding capacity from
279 78.4% to 91.4%. However, the same was not observed for INPHARED sequences predicted to
280 use translation table 4. These sequences observed a modest increase in median median
281 gene length from 573 to 584 bp (1.8%) for Prokka-gv. Median coding capacity was not
282 increased with Prokka and Prokka-gv both obtaining 86.2%.

A**B**