

1 **Supplemental files**

2 **Document S0.**

3 A list of captions for the Supplemental files.

4 **Document S1.**

5 A complete and unabridged description of the materials and methods used in this study.

6 **Document S2.**

7 A stand-alone, HTML webpage reporting the conservation and synteny of the HiVir gene cluster
8 within the genomes of this study. These results were produced by CBLASTER [Gilchrist2021a] v
9 1.3.19, using as the input query the HiVir gene cluster from *P. ananatis* LMG 20103, as
10 described in Asselin et al. (2018) [38].

11 **Figure S1.**

12 The correlation plot, based on ANI (average nucleotide identity) values, categorizes *P.*
13 *agglomerans* (*Pagg*) strains according to their phylogroups. Pairwise ANI values for 187 *Pagg*
14 strains were calculated using FastANI (v. 1.3.3) and organized through a hierarchical clustering
15 method employing the Pearson correlation matrix. The strains are partitioned into two groups,
16 denoted by red and blue, aligning with phylogroup distinctions (blue corresponds to Phylogroup
17 I, while red corresponds to Phylogroup II).

18 **Figure S2. The phylogeny of *Pantoea agglomerans*, based on the *gyrB* gene is** 19 **largely congruent with the core genome phylogeny.**

20 (A) An approximately-maximum-likelihood phylogenetic tree, constructed using FastTree
21 v.2.1.11, relies on the 3111 core genes shared among 181 *P. agglomerans* strains (Figure 1).
22 Maximum-likelihood phylogenetic trees, constructed based on partial *gyrB* gene sequences
23 targeted by primers designed by (B) [45] and (C) [46], indicate that the *Bonasera gyrB* region

24 reveals a phylogeny congruent with that of the core genome. Strains not assigned to
25 phylogroups in the core genome tree (CFBP8785, T2, T3, T4, T5, and P5) form part of
26 phylogroups in *gyrB* trees. Phylogroups are highlighted according to the phylogroup designation
27 in the core genome tree, with a yellow box representing Phylogroup I and red representing
28 Phylogroup II. The association of strains with onions is indicated by the purple coloring of strain
29 names.

30 **Figure S3.**

31 Extended phylogroup identification of 226 *Pantoea agglomerans* strains using maximum-
32 likelihood phylogenetic analysis of partial *gyrB* gene as targeted by primers designed by [45].
33 *Tatumella pyseos* ATCC33301^T was used as an outgroup, and branch lengths were hidden.
34 Phylogroups are highlighted: Phylogroup I is marked with a yellow box, and Phylogroup II with
35 pink.

36 **Figure S4.**

37 A heatmap visualization of the presence (pink) and absence (yellow) genes (n= 27,204)
38 observed across 181 *Pagg* strains. The dendrogram drawn based on the gene presence and
39 absence categorizes *Pagg* strains into two groups whose members are consistent with those of
40 the Phylogroup I (indicated by black line) and Phylogroup II (red line).

41 **Figure S5.**

42 A combined hierarchical clustering and heatmap showing the relationship among all of the
43 replicons (chromosomes and plasmids) studied. The distances between replicons were
44 computed using Mashtree [59]. Values shown are 100×(1-MASH), which is analogous to
45 percentage identity.

46 **Figure S6.**

47 A combined hierarchical clustering and heatmap showing the relationship among the
48 **chromosomal** replicons studied. The distances between replicons were computed using
49 Mashtree [59]. Values shown are $100 \times (1 - \text{MASH})$, which is analogous to percentage identity.

50 **Figure S7.**

51 A combined hierarchical clustering and heatmap showing the relationship among the **LPP-1**
52 replicons studied. The distances between replicons were computed using Mashtree [59]. Values
53 shown are $100 \times (1 - \text{MASH})$, which is analogous to percentage identity.

54 **Figure S8.**

55 A combined hierarchical clustering and heatmap showing the relationship among the **pAggl**
56 replicons studied. The distances between replicons were computed using Mashtree [59]. Values
57 shown are $100 \times (1 - \text{MASH})$, which is analogous to percentage identity.

58 **Figure S9.**

59 A combined hierarchical clustering and heatmap showing the relationship among the **pOnion**
60 replicons studied. The distances between replicons were computed using Mashtree [59]. Values
61 shown are $100 \times (1 - \text{MASH})$, which is analogous to percentage identity.

62 **Figure S10.**

63 A combined hierarchical clustering and heatmap showing the relationship among the **pENY**
64 replicons studied. The distances between replicons were computed using Mashtree [59]. Values
65 shown are $100 \times (1 - \text{MASH})$, which is analogous to percentage identity.

66 **Figure S11.**

67 This figure shows pAggl plasmid plasmids pAR5_B and pAR1aB with core pAggl genes labeled
68 with green boxes. Core genes for pAggl were determined using Roary based on Prokka
69 annotations. Arrows denote CDSs from RefSeq annotations. Some genes appeared to occur in

70 clusters where genes had similar function or seemed to act in a single pathway. Where their
71 annotations allowed reasonable guesses of function, gene clusters are labeled with yellow bars
72 and expected functions.

73 **Figure S12.**

74 An alignment of *P. agglomerans* (*Pagg*) AR1a's plasmids C and D against *Pagg* T88c's plasmid C.
75 All three plasmids are members of the pOnion plaster. These results demonstrate the sequence
76 of T88c plasmid C can be reconstructed using sequence from AR1'a plasmids C and D. T88c
77 plasmid C CDSs with possible roles in partition and replication are filled in black, genes with
78 annotations suggesting involvement in recombination (e.g. transposases, integrases) are filled
79 in purple, genes annotated as *umuC* or *umuD* are filled in green, and CDSs involved in type IV
80 secretion are filled in pink. The figure was generated using the NCBI BLAST website, the blastn
81 tool, the "Align two or more sequences" option, optimized for somewhat similar sequences.
82 T88c plasmid C sequence was used as the query, and AR1a plasmids C and D as subject
83 sequences. The graphical representation of the alignment was pasted below a linear
84 representation of the plasmid generated in SeqBuilder 17.3 (DNASTar, Madison, WI) using pgap
85 annotations.

86 **Table S1.**

87 A comprehensive list of the primers utilized in this study for gene detection, sequencing, and
88 genetic manipulations.

89 **Table S2.**

90 The 618 *Pantoea* strains that are used in this manuscript and the tables in which they are
91 mentioned.

92 **Table S3.**

93 A list of the 100 public *Pantoea* genomes employed in this study for genome analysis, detailing
94 strain information including source, year, and place of isolation, and their respective GenBank
95 accession numbers.

96 **Table S4.**

97 A list of 501 *Pantoea* strains tested for red onion scale necrosis assay (RSN), and *P*
98 *agglomerans*-specific and *Pantoea*-HiVir PCR assays.

99 **Table S5.**

100 QCAST quality assessment results for the 87 *Pantoea agglomerans* genomes sequenced in this
101 study.

102 **Table S6.**

103 A table displaying the BUSCO genome assembly and annotation completeness scores for 81
104 genomes sequenced in this study.

105 **Table S7.**

106 **Pan-replicon Analysis Results of Chromosome Replicons**

107 A list of all of the gene groups identified by Roary [48] for replicons belonging to the
108 “chromosome” plaster. Each row reports a single gene group. The “Gene” and “Annotation”
109 values report the Prokka [47] predicted gene name and gene product for the genes in the
110 group. In the case that multiple different gene names or products are predicted, all predictions
111 are listed. The “Count” and “Pers” columns report the number and percentage of replicons in
112 the plaster group containing members of the genome group. The remaining columns indicate
113 the presence of the gene group in individual replicons.

114 **Table S8.**

115 **Pan-replicon Analysis Results of LPP-1 Replicons**

116 A list of all of the gene groups identified by Roary [48] for replicons belonging to the “LPP-1”
117 plaster. Each row reports a single gene group. The “Gene” and “Annotation” values report the
118 Prokka [47] predicted gene name and gene product for the genes in the group. In the case that
119 multiple different gene names or products are predicted, all predictions are listed. The “Count”
120 and “Pers” columns report the number and percentage of replicons in the plaster group
121 containing members of the genome group. The remaining columns indicate the presence of the
122 gene group in individual replicons.

123 **Table S9.**

124 **Pan-replicon Analysis Results of pAggl Replicons**

125 A list of all of the gene groups identified by Roary [48] for replicons belonging to the “pAggl”
126 plaster. Each row reports a single gene group. The “Gene” and “Annotation” values report the
127 Prokka [47] predicted gene name and gene product for the genes in the group. In the case that
128 multiple different gene names or products are predicted, all predictions are listed. The “Count”
129 and “Pers” columns report the number and percentage of replicons in the plaster group
130 containing members of the genome group. The remaining columns indicate the presence of the
131 gene group in individual replicons. The ortholog groups whose members include genes from the
132 HiVir gene cluster are shared blue.

133 **Table S10.**

134 **Pan-replicon Analysis Results of pOnion Replicons**

135 A list of all of the gene groups identified by Roary [48] for replicons belonging to the “pOnion”
136 plaster. Each row reports a single gene group. The “Gene” and “Annotation” values report the
137 Prokka [47] predicted gene name and gene product for the genes in the group. In the case that
138 multiple different gene names or products are predicted, all predictions are listed. The “Count”

139 and “Pers” columns report the number and percentage of replicons in the plaster group
140 containing members of the genome group. The remaining columns indicate the presence of the
141 gene group in individual replicons. The ortholog groups whose members include genes from the
142 cop, HiVir, and alt gene clusters are shared red, blue, and green, respectively.

143 **Table S11.**

144 **Pan-replicon Analysis Results of pENY Replicons**

145 A list of all of the gene groups identified by Roary [48] for replicons belonging to the “pENY”
146 plaster. Each row reports a single gene group. The “Gene” and “Annotation” values report the
147 Prokka [47] predicted gene name and gene product for the genes in the group. In the case that
148 multiple different gene names or products are predicted, all predictions are listed. The “Count”
149 and “Pers” columns report the number and percentage of replicons in the plaster group
150 containing members of the genome group. The remaining columns indicate the presence of the
151 gene group in individual replicons.

152 **Table S12.**

153 A pairwise average nucleotide identity (ANI) comparison table of 187 *Pantoea agglomerans*
154 genomes.

155 **Table S13.**

156 Type strains of *Pantoea* species and their respective GenBank accession numbers from which
157 the *gyrB* gene was extracted for phylogenetic analysis.

158 **Table S14.**

159 Metadata information utilized for annotating the 181 *Pantoea agglomerans* strains that appear
160 in Figure 1.

161 **Table S15.**

162 A pairwise nucleotide identity comparison of HiVir gene clusters among *Pantoea* species.

163 **Table S16.**

164 A compilation of 125 *Pantoea agglomerans* strains subjected to whole genome sequencing
165 (WGS) and *gyrB* gene sequencing in this study. Closed, complete genomes were generated for
166 19 strains, draft genomes were made for 68 strains, and *gyrB* gene sequences were made for 38
167 strains.

168 **Table S17.**

169 This table lists the strains whose genomes were used in the plaster analysis.

170 Initially, 38 closed, complete genomes were collected into a pool of *P. agglomerans* (*Pagg*)
171 genomes. The table lists, for each strain, its “phylogroup” and GenBank “accession” identifier.

172 In order to confirm these taxonomic identifications, we computed the ANI scores of each strain
173 compared to the *Pagg* type strain (here denoted strain FDAARGOS1447). We found that the
174 ANI scores of 3 strains (33.1, FL1, and HJS002), as found in the “fastani” column, fell below the
175 typical species cutoff of 95% and were removed from the pool. We used BUSCO to evaluate the
176 integrity of the remaining 35 assemblies; we found that that “busco_C” and “busco_D” scores
177 ranged from 98.7%-99.5% and 0.2%-2%, respectively, indicating the quality of these assemblies
178 ranged from good to excellent.

179 The “new” column records whether the genome assembly is new to this study. The “used”
180 column records whether or not the genome was ultimately used as input for the plaster
181 analysis.

182 **Table S18.**

183 A list of replicons used in the plaster analysis and their statistics. The pool of 35 genomes listed
184 in Table S17 consisted of 35 chromosomes plus 125 plasmids for a total of 160 replicons. The
185 chromosomes ranged in size from 3,978,822bp to 4,410,564; the plasmids ranged in size from
186 2,550bp to 613,013bp.

187 **Table S19.**

188 A list of all the plaster assignments by our plaster analysis to the 160 replicons listed in Table
189 S18. Our analysis identified 31 distinct plasmid clusters, or “plasters”, 9 of which were non-
190 trivial (>1 member). We assigned previously published names to 3 of the non-trivial plaster
191 plasters (“chromosome”, “LPP-1”, and “pPATH”) and new names to 3 others (“pAggl”,
192 “pOnion”, and “pENY”).

193 **Table S20.**

194 List of genes in the pAggl core genome in the context of the *Pantoea agglomerans* AR1a pAggl
195 plasmid. Core pAggl genes were called based on Prokka annotations using Roary. Where the
196 coordinates for the annotations did not match exactly (e.g. Prokka and RefSeq listed the gene
197 with different start sites), the coordinates for RefSeq are listed to be consistent with locus
198 name, GO terms and EC numbers.

199 **Table S21.**

200 List of genes in the pAggl core genome in the context of the *Pantoea agglomerans* AR5 pAggl
201 plasmid. Core pAggl genes were called based on Prokka annotations using Roary. Where the
202 coordinates for the annotations did not match exactly (e.g. Prokka and RefSeq listed the gene
203 with different start sites), the coordinates for RefSeq are listed to be consistent with locus
204 name, GO terms and EC numbers.

205 **Table S22.**

206 A list of 123 *Pantoea agglomerans* strains was utilized in the copper tolerance assays. The
207 presence of the *copC* gene was screened in these strains using the *copC* detection primers
208 designed in this study. Copper tolerance was assessed based on the level of colony growth on
209 CYE agar supplemented with varying concentrations of copper, with scores assigned as follows:
210 0 for no growth, 1 for non-confluent growth, and 2 for confluent growth.

211 **Table S23.**

212 **Predicted mobilizable and conjugative elements**

213 Mobilizable and conjugative elements were predicted by MacSyFinder v. 2.1.2 [56, 57] using the
214 CONJScan` `/Plasmids v. 2.0.1 models [58]. For each replicon, the number of predicted
215 instances of each submodel is shown. The submodels with predicted include, MOB
216 (mobilizable elements), T4SS_typeF and T4SS_typeG (functional conjugative T4SS's), and
217 dCONJ_typeG and dCONJ_typeT (decayed conjugative T4SS's).

218 **Table S24.**

219 Details of the *Pantoea agglomerans* strains used in conjugation experiments. The replicons
220 found in each strain are listed with some relevant features. TXXScan and CONJScan results using
221 MacSyFinder are listed for the strains. Strain derivatives (rifampicin resistant or altered to carry
222 a kanamycin resistance cassette) are presumed to have identical results to wild-type strains.

223 **Table S25.**

224 Results from preliminary conjugation experiments using *Pantoea agglomerans* (Pagg) strain
225 CB1 harboring a version of the pCB1C plasmid with a kanamycin-resistance cassette as donor
226 and four non-pathogenic strains of Pagg as potential recipients. For each mating pair, three
227 colonies, or fewer if three were not available, were screened via PCR with primers designed to
228 amplify the donor strain (CB1 3125-F/3125-R), the recipient strains (ImpA F1/R2), or from
229 within the HiVir cluster (Pag 3283 F1/R1), which is present on pCB1C but not in any of the wild-
230 type recipient strains. The results are color coded. Green means the colonies gave the expected
231 PCR result. Red means PCR indicated rifampicin-resistant donor colonies breaking through.
232 Yellow means the PCR result was ambiguous. Putative transconjugants were selected for
233 genome sequencing from conjugations with Pagg CB1 pCBC-AKan colony 30 as the donor.
234 Putative transconjugants were recovered from conjugations with AR8b Rp^r, MMD61212-C Rp^r,
235 and FC61912-B Rp^r as recipients, but not with J22c Rp^r.

236 **Dataset S1.**

237 The consolidated and reconciled annotations from all of the WGS strains used in this study. In
238 this study, Prokka was used, independently, by both the UGA and USDA authors to reannotate
239 genomes for input to Roary [48]. For genomes available from NCBI, both GenBank and RefSeq
240 annotations are available. Hence, for each genome, multiple sets of gene annotations are
241 available. For each WGS strain, we have provided a table listing which locus tags are equivalent
242 in the available annotation sets. That is, each row of each table lists the equivalent locus tags
243 from a genome's available annotation sets. Locus tags are considered equivalent if they refer to
244 CDS features whose amino acid sequences are equal.