

1 **Just the Two of Us? A Family of *Pseudomonas***
2 **Megaplasmiids Offers a Rare Glimpse Into the Evolution of**
3 **Large Mobile Elements**

4

5

6 Brian A. Smith^{1*}, Courtney Leligdon¹, and David A. Baltrus^{1,2}

7

8

9 ¹School of Plant Sciences, University of Arizona, Tucson, AZ, USA

10 ²School of Animal and Comparative Biomedical Sciences, University of Arizona,

11 Tucson, AZ, USA

12

13 *corresponding author

14 basmith@email.arizona.edu

15

16

17

18 **Abstract**

19 Pseudomonads are ubiquitous group of environmental proteobacteria, well known
20 for their roles in biogeochemical cycling, in the breakdown of xenobiotic materials,
21 as plant growth promoters, and as pathogens of a variety of host organisms. We
22 have previously identified a large megaplasmid present within one isolate the plant
23 pathogen *Pseudomonas syringae*, and here we report that a second member of this
24 megaplasmid family is found within an environmental Pseudomonad isolate most
25 closely related to *P. putida*. Many of the shared genes are involved in critical cellular
26 processes like replication, transcription, translation, and DNA repair. We argue that
27 presence of these shared pathways sheds new light on discussions about the types
28 of genes that undergo horizontal gene transfer (i.e. the complexity hypothesis) as
29 well as the evolution of pangenomes. Furthermore, although both megaplasms
30 display a high level of synteny, genes that are shared differ by over 30% on average
31 at the amino acid level. This combination of conservation in gene order despite
32 divergence in gene sequence suggests that this Pseudomonad megaplasmid family is
33 relatively old, that gene order is under strong selection within this family, and that
34 there are likely many more members of this megaplasmid family waiting to be found
35 in nature.

36

37 **Introduction**

38 Horizontal Gene Transfer (HGT) of megaplastids can rapidly create dramatic
39 phenotypic differences between otherwise closely related bacterial strains, with
40 potential for over a thousand genes to be gained by a strain in a single event.
41 Although there have been numerous attempts to identify overarching themes for the
42 evolutionary effects of HGT based on types of genes and pathways transferred, such
43 efforts have often neglected to incorporate intrinsic characteristics of megaplastids
44 ¹⁻⁶. Furthermore, since secondary replicons are prone to rapid reshuffling of gene
45 order as well as extensive gains and losses of loci, it has traditionally been
46 challenging to analyze past evolutionary dynamics to understand overall historical
47 pressures acting on this class of mobile elements⁷⁻¹⁴. More thorough investigation of
48 evolutionary dynamics within relatively large plasmid families could therefore
49 provide new viewpoints into the evolutionary effects of gene transfer and may also
50 enable broader generalizations about selective forces driving the composition and
51 overall structure of megaplastids, chromids, and second chromosomes.

52

53 Megaplastids are generally characterized as low copy extrachromosomal replicons
54 >350kb in size and which are dispensable to the bacterial cell under a subset of
55 conditions¹⁴. As with many plasmid families, they have often been identified because
56 they impart beneficial phenotypes such as resistance to antimicrobial compounds or
57 introduce novel catabolic pathways into host cells¹⁴. Given their size and gene
58 content, it is possible that megaplastids possess greater potential for generating
59 evolutionary costs than smaller plasmids when transferred to naive hosts^{15,16}.

60 However, efforts to identify shared genes and pathways across megaplasmiids and to
61 use this information to make predictions about potential systems level conflicts
62 have been hampered by poor sampling across novel megaplasmiid families¹⁴.

63 Identification of additional examples can help to fill this gap in current knowledge
64 and may uncover new evolutionary trends that govern megaplasmiid-chromosomal
65 interactions.

66

67 Consideration of megaplasmiids may uniquely inform general discussions about
68 evolutionary effects HGT in ways that have been overlooked by analyses focusing
69 simply on distributions of genes and pathways maintained after transfers and
70 without considering timing of HGT or linkage between loci. For instance, many of
71 the earliest discussions concerning evolutionary constraints of HGT, the so called
72 “complexity hypothesis”, found that loci associated with “more complex” cellular
73 processes undergo lower rates of transfer than other genes^{4,17}. Interpretations of
74 these patterns have changed through time with examples of horizontally acquired
75 informational genes, suggesting that it is actually the shape of protein interaction
76 networks that are critical for maintenance after gene transfer.^{4,18-22} However,
77 larger mobile elements like megaplasmiids can contain genes encoding proteins and
78 pathways that could be classified as “complex” (i.e. proteins involved in translation)
79 and which have the potential to interact with numerous chromosomally encoded
80 pathways^{14,23}. Likewise, a variety of recent papers have focused on selective forces
81 (or lack thereof) governing microbial pangenomes^{1-3,24}. These discussions have
82 largely focused on population level distributions for single genes that compose a

83 pangenome, but by their nature intrinsically fail to consider linkage of genes on
84 plasmids that are frequently acquired and lost. While many genes within the
85 pangenome may indeed be ‘adaptive’, such a viewpoint overlooks the idea that no
86 single gene need be adaptive for the bacterial cell if selection acts at the level of
87 plasmid transfer and hundreds of genes that could be linked to that process. More
88 thorough characterization of multiple megaplasmid families and identification of
89 new megaplasmids will enable identification of the types of genes and pathways
90 canonically associated with these large vectors and patterns that emerge can be
91 incorporated into greater discussions of the general role of HGT across bacterial
92 species.

93

94 Previously we have described a megaplasmid, pMPPla0107, found within one isolate
95 of the phytopathogen *Pseudomonas syringae*²³. pMPPla107 is self-transmissible
96 across the *Pseudomonas* phylogeny, harbors numerous loci that could be annotated
97 as “housekeeping” genes, and it is stably maintained within recipient cells^{23,25,26}. We
98 have also demonstrated that acquisition of this megaplasmid through HGT also
99 imparts significant phenotypic costs to recipient cells, likely mediated by
100 detrimental interactions between chromosomal and plasmid encoded proteins²⁶. It
101 is unclear if pMPPla107 is the only megaplasmid of its kind and size with
102 conjugation abilities across an entire genus of bacteria or if it is a member of a larger
103 family of secondary replicons.

104

105 Here we identify a new megaplasmid, pBASL58, related to pMPPla107 and use
106 molecular and computational approaches to characterize this megaplasmid family
107 more broadly. We show that, while these megaplasmids are similar in size, genetic
108 structure, nucleotide bias, and functionality, there is a high level of divergence
109 across shared orthologous gene groups, a dissimilar cargo region, and differing
110 CRISPR loci. This overall level of divergence suggests that both members of this
111 megaplasmid family have been independently evolving for a relatively long period
112 of time. Even more, we find that these divergent orthologous pathways demonstrate
113 high levels of synteny in the context of overall plasmid structure; suggesting that
114 conservation of gene orientation and order is under relatively strong selective
115 pressures. Lastly, characterization of these plasmids allows for a chance to
116 emphasize that pathways found on both megaplasmids are likely involved in
117 important cellular processes like nucleotide synthesis and DNA replication, which
118 highlights new discussion points to add to the complexity hypothesis as well as the
119 adaptive nature of pangenomes.

120

121 **Methods**

122

123 *Identification of pBASL58*

124 Initial BLASTp searches coupled with inspection of contigs from a draft genome
125 assembly of *Pseudomonas* sp. Leaf58²⁷ suggested that this strain could contain a
126 megaplasmid related to pMPPla017, represented within contig
127 Ga0102293_111 within the draft genome assembly containing 18 contigs total

128 (Genbank accession GCA_001422615.1). An isolate of *Pseudomonas* sp. Leaf58 was
129 obtained from DSMZ (DSM-102683), and a single colony was picked from a culture
130 from the freeze-dried ampule plated on unsupplemented LB media. To test for
131 circularization of contig Ga0102293_111, primers (BAS 17-31) were designed to
132 amplify off of the edges of the contig and overlap each other. An approximate 1.5kb
133 size PCR product was amplified from an overnight culture of this strain grown in KB
134 media, demonstrating circularization of this contig. Sanger sequencing of this PCR
135 fragment demonstrated that the initial draft contig contained a missassembly, and
136 so this sequence was corrected by hand, and the contig as reoriented and used in all
137 analyses in this manuscript. Sequence of this contig can be found at Figshare
138 (doi:[10.6084/m9.figshare.6914033](https://doi.org/10.6084/m9.figshare.6914033)). For consistency of analyses throughout the
139 manuscript we reannotated the megaplasmid sequence with Prokka v1.12 using
140 default parameters, and this annotation can also be found at Figshare
141 (doi:[10.6084/m9.figshare.6914033](https://doi.org/10.6084/m9.figshare.6914033)).

142

143 *Genome Sequencing, assembly, and annotations of Pseudomonas sp. Leaf58*

144 As further evidence of the existence of a megaplasmid in Leaf58, we generated a
145 complete genome assembly for this strain (currently found at Figshare
146 (doi:[10.6084/m9.figshare.6914033](https://doi.org/10.6084/m9.figshare.6914033)), Genbank accession TBD). After revival from
147 the Baltrus lab stock, a single colony *Pseudomonas* sp. Leaf58 was picked to an
148 overnight culture in KB media, and grown in a shaking incubator at 27°C. After
149 approximately 24 hours, DNA was extracted from this culture using a Promega
150 Wizard kit. A rapid sequencing library was created using this DNA, and 169,316

151 reads (933,937,907 total bp, 5,515bp average read size) were generated on an R9.4
152 flowcell using a Rapid sequencing kit (SQK-RAD004). Additionally, 100bp paired
153 end Illumina reads used to generate the original draft genome of this strain were
154 downloaded from the SRA (Accession ERR1103815)²⁷. A complete genome
155 sequence for *Pseudomonas* sp. Leaf58 was generated by combining these short and
156 long reads in Unicycler (version 0.4.4)²⁸. This sequence consists of a single
157 chromosome (5,432,868 bp) and the pBASL58 megaplasmid (904,253bp), both of
158 which were circular according to Unicycler.

159

160 *Genome Sequencing, assembly, and annotations of Pla107*

161 A single colony of the Baltrus lab stock of *Pseudomonas syringae* pv. *lachrymans* 107
162 (MAFF31015) was picked to an overnight culture in KB media, and grown in a
163 shaking incubator at 27°C. After approximately 24 hours, DNA was extracted from
164 this culture using a Promega Wizard kit. Illumina sequencing of was performed by
165 MicrobesNG, and generated 2,771,213 250bp paired end reads (231 median read
166 length after trimming, ~166x coverage of the genome) on an Illumina MiSeq.
167 Assembly was performed using SPAdes v3.10.1 with default parameters as well as
168 through MicrobeNG's bioinformatics pipeline, which matches the reads to the best
169 reference using Kraken and maps reads back to that reference using BWA-MEM²⁹.
170 MicrobeNG also uses *de novo* assembly with SPAdes. pMPPla107 assembled
171 completely from these reads, and this version of the megaplasmid sequence can be
172 found Figshare (doi:[10.6084/m9.figshare.6914033](https://doi.org/10.6084/m9.figshare.6914033)) and was used for all analyses
173 throughout this manuscript. Gene annotation of this version of the megaplasmid

174 sequence was performed with Prokka v1.12 using default parameters. This gene
175 model used for all coding sequence analyses within the manuscript and can be
176 found at Figshare (doi:[10.6084/m9.figshare.6914033](https://doi.org/10.6084/m9.figshare.6914033)). We additionally generated
177 long read sequences for Pla107 using a MinION from Oxford Nanopore. A rapid
178 sequencing library was created from an independent genomic isolation of a
179 derivative of Pla107, DBL328, which contains an integrated version of the
180 pMTN1907 marker plasmid and which has been selected to for kanamycin
181 resistance from this marker plasmid. As above, a single colony of this strain was
182 picked to an overnight culture in KB media and DNA was extracted with a Promega
183 Wizard kit. 15,461 reads (139,041,576 total bp, 8,993 average read size) were
184 generated on an R9.4 flowcell using a Rapid sequencing kit (SQK-RAD004). A whole
185 genome assembly was created by combining both MiSeq and MinION reads using
186 Unicycler (version 0.4.4)²⁸ with default parameters. This whole genome sequence
187 consists of a circular chromosome (6,075,120 bp), pMPPla107 (971,889 bp, and
188 sequence identical to the assembly from SPADES alone), and two other plasmids,
189 pPla107-1 (62,136 bp) and pPla107-2 (40,720 bp). Three of these sequences
190 (except pPla017-1) were complete and circular contigs according to Unicycler
191 assembly. This assembly was used to update the Genbank version of this genome,
192 and is found at accessions (CP031225, CP031226 CP031227, CP031228). Gene
193 annotations in this Genbank file were generated by NCBI's PGAAP pipeline³⁰.

194

195 *Identifying Origins of Replication*

196 To identify putative origins of replication for both megaplasmiids, we used a
197 modified GC skew script³¹ to scan the entirety of pMPPla107 and pBASL58 and
198 combined this information with characterization of repetitive motifs that could
199 represent *oriV* sites. GC skew and repetitive motifs suggest pMPPla107 and
200 pBASL58 have predicted origins of replication within a similar genomic region near
201 partitioning genes (Figure 2). Based on this information we oriented the sequences
202 of pMPPla107 and pBASL58 to begin at the start codon of shared *parA*-like loci. We
203 chose the *parA*-like locus as the starting point because it is shared by both
204 sequences, is near the predicted origin of replication, and is predicted to be an
205 important gene for plasmid partitioning.

206

207 *CRISPR Identification*

208 CRISPR-Cas and repeat structure annotations were identified using both Prokka
209 annotations and the web tool CRISPRCasFinder³²⁻³⁴

210

211 *Plasmid Comparisons With BLASTp and MAUVE*

212 Amino acid sequence names were changed to numbers in an increasing order using
213 the `mod_protein_id.py` script. We then used the BLAST 2.6.0+ package³⁵. BLASTp
214 parameters were altered to ensure only the top hit was returned and that there
215 were zero overlapping hits. The BLAST command used was:

216

```
217 blastp -db [blastdb] -query [query_file] -culling_limit 1 -max_target_seqs 1 -
```

```
218 max_hsps 1 -out [out_file] -outfmt 6
```

219

220 Data was extracted from the BLAST output at 40, 50, 60, and 70% identity cutoffs
221 and plotted in R using ggplot2.

222

223 pMPPla107 and pBASL58 sequences were input into Progressive Mauve 2.3.1 to
224 compare megaplasmid sequences within the software Geneious ³⁶.

225

226 *Gene Mapping Visualization with Circa*

227 The BLASTp output data mentioned above was altered in a format to comply with
228 input to Circa using gff_info_extract.py followed by geneid_match.py. The
229 parameters used to generate the Circa map and the Python scripts used to generate
230 the data can be found at the https://github.com/basmith89/megaplasmid_compare.

231

232 *Tetranucleotide frequency Comparisons*

233 We performed pairwise comparisons of tetranucleotide frequencies between
234 chromosome sequences and secondary replicon sequences in an all by all method.
235 Tetranucleotide frequencies were calculated with the calc.kmerfreq.pl script created
236 by Mads Albertsen³⁷ found at [https://github.com/MadsAlbertsen/multi-](https://github.com/MadsAlbertsen/multi-metagenome)
237 [metagenome](https://github.com/MadsAlbertsen/multi-metagenome). Output of this script was plotted using ggplot2 and R² values were
238 calculated in R.

239

240 *Functional Comparisons With KEGG, KASS, and UProC*

241 We carried out two analyses utilizing the Kyoto Encyclopedia of Genes and Genomes
242 (KEGG) database³⁸. Amino acid sequences of coding regions predicted by Prokka
243 were input into the protein sequence classification software, UProC³⁹. UProC's
244 output is a list of KEGG IDs and counts. We designed a perl script,
245 `kegg_path_counter.pl`, to extract these ID's and counts and associated them with
246 KEGG functional pathways. The script and ID key can be found at
247 https://github.com/basmith89/megaplasmid_compare. These data were then
248 plotted with the Plotly package in R.

249

250 Amino acid sequences output by Prokka were also run through a Python script to
251 produce a list of gene annotations that both megaplasמידs have in common
252 https://github.com/basmith89/megaplasmid_compare. Amino acid sequences
253 from genes on this shared list were then run through KASS (KEGG Automatic
254 Annotation Server) to determine what pathways are shared by the megaplasמידs⁴⁰.
255 Pathways maps were then condensed into one figure by hand.

256

257 **Results**

258 *A new member of the pMPPla107 megaplasמיד family*

259 pMPPla107 was originally identified from an assembly using both 454 and 30bp
260 Illumina sequencing reads²³. However, due limitations of these early technologies,
261 this assembly of pMPPla107 remained incomplete and consisted of linked scaffolds.
262 We therefore utilized updated sequencing and assembly technologies to sequence
263 the *P. syringae* genome containing pMPPla107, yielding a complete circular

264 sequence for this megaplasmid (971,889bp compared to 963,598bp in original
265 sequence) (Table 1). Additionally, multiple searches using protein sequences from
266 pMPPla107 consistently yielded high quality matches to the scaffold
267 Ga0102293_111 (referred to as pBASL58 from here on) from a public genome
268 assembly of *Pseudomonas* sp. Leaf58. This strain was originally isolated as part of a
269 project to thoroughly sample cultureable strains from the phyllosphere of
270 Arabidopsis and is most closely related to *P. putida* strains⁴¹. We independently
271 confirmed circularization (Figure 1) of this contig from Leaf 58, using both PCR and
272 long read nanopore sequencing, definitively showing this contig was indeed a large
273 megaplasmid separate from the chromosome.

274

275 *Both megaplasmsids contain numerous tRNA loci*

276 The size, number of predicted genes, number of tRNAs, and GC content are highly
277 similar between pMPPla107 and the pBASL58 (Table 1). Overall GC content was
278 similar in Leaf58 and *P. syringae lac107*, and the GC content in both pMPPla107 and
279 pBASL58 were lower than their respective chromosomal partners. pBASL58 and
280 pMPPla107 contained 54 and 44 regions annotated as tRNA loci, respectively.
281 pBASL58 encodes 20 unique tRNAs and pMPPla107 encodes 10, some of which
282 were repetitive like tRNA-Glu(ttc) in pBASL58 occurring six times. When observing
283 tRNA amino acid products, pMPPla107 encodes for 16/20 possible amino acids and
284 pBASL58 encodes for 19/20 possible amino acids possibly indicating pBASL58 is
285 less dependent on host tRNAs. In addition to the 16 amino acids produced by
286 pMPPla107, pBASL58 is predicted to code for the ability to charge tRNAs with

287 tryptophan, glutamate, and aspartate and both plasmids are missing any anticodons
288 to produce histidine. These differences could suggest an amino acid preference for
289 the maintenance or protein production of the plasmids.

290

291 *Identifying genomic similarities of pMPPla107 and the Leaf58 plasmid*

292 Both megaplasms within this new family are highly syntenic (Figure 3A, with
293 Mauve alignment showing that 72.8% (707,677bp out of 971,871bp) of pMPPla107
294 aligns well with 71.6% (646,763bp out of 903,765bp) of pBASL58 (Supplemental
295 Figure 1. The regions of highest similarity occur near the origin of replication.

296 Despite overall high levels of synteny, there is a highly dissimilar region
297 (approximately 300kb in size) occurring within the first half of the sequences and a
298 \approx 50kb inversion in the last half indicating these megaplasms have also undergone
299 structural diversification.

300

301 Even though both megaplasms display high levels of synteny, preliminary
302 comparisons of protein sequences suggested a relatively high level of divergence
303 between orthologues shared by both megaplasms (Figures 3B and C). The highest
304 levels of average amino acid similarity (48.6%) occur near the predicted origin of
305 replication where genes for plasmid replication, partitioning, and conjugation are
306 common. Areas near the terminus still demonstrate strong synteny but have higher
307 divergence in amino acid identity (\approx 38.2% similarity). These data suggest
308 pMPPla107 and pBASL58 are structurally related to each other and share a common
309 plasmid ancestor, but have experienced independent evolutionary pressures for

310 long enough time for significant diversification to occur within shared protein
311 sequences.
312
313 To further gauge relationships between both megaplasmiids and the chromosomes
314 of their host strains, we compared tetranucleotide frequencies for each of these
315 replicons⁴²⁻⁴⁴. Pairwise comparisons demonstrated that pMPPla107/pBASL58 (R^2
316 = 0.878) and the *P. syringae*/Leaf58 (R^2 = 0.889) chromosomes are most similar in
317 frequencies (Figure 4). All remaining pairwise comparisons reported R^2 values less
318 than 0.780. pMPPla107 shows the greatest differences in tetranucleotide
319 frequencies when compared to both the *P. syringae* and the Leaf58 chromosomes
320 with R^2 values of 0.524 and 0.393 respectively. pBASL58 shares slightly more
321 similar frequency preferences indicative of R^2 values of 0.780 and 0.695, to *P.*
322 *syringae* and Leaf58 chromosomes respectively. This data suggest that mutational
323 biases affecting these secondary replicons are most similar to each other, which
324 suggests that they have not been replicating within these host strains long enough
325 to be subject to amelioration.

326

327 *Housekeeping gene functionality is shared by pBASL58 and pMPPla107*

328 Based on the structural similarities established, we hypothesized that pMPPla107
329 and pBASL58 would share similar functional pathways. UProC called 9% (85) and
330 10% (111) of the predicted coding regions for pBASL58 and pMPPla107,
331 respectively, indicating the majority of predicted gene functionality is unknown.
332 Annotation with Prokka returned similar results (13% of genes with annotated

333 functions). The pathways and functions most frequently annotated were replication
334 and repair at 2.3% (22 genes) for pBASL58 and 2.4% (26) for pMPPla107, global
335 and overview maps at 2.1% (20) for pBASL58 and 2.2% (24) times for pMPPla107,
336 and nucleotide metabolism at 1.3% (12) for pBASL58 and 1.8% (19) for pMPPla107
337 (Figure 5). KEGG KASS also predicated that the two megaplasmids share 57.6%
338 (99/172) of annotated genes. Therefore pBASL58 and pMPPla107 carry 31 and 42
339 unique genes respectively. Again, the overall distribution of gene products present
340 on both megaplasmids tends towards DNA synthesis, DNA repair, and synthesis of
341 deoxyribonucleotide-triphosphates (Supplemental Table 1 and Supplemental Figure
342 3). These shared groups include DNA polymerase III subunits, helicases, primase,
343 ligases, recombination proteins, and exonucleases indicating these megaplasmids
344 encode for pathways associated with their maintenance. Other gene products on
345 these megaplasmids are involved in metabolic pathways such as fatty acid
346 biosynthesis, RNA degradation, Aminoacyl tRNA biosynthesis, and NOD-like
347 receptor signaling pathways. Interestingly, both plasmids also encode for several
348 membrane and multidrug efflux pump genes. Both shared efflux genes belong to the
349 Resistance-Nodulation-Division (RND) family of transporters and are known for
350 their multidrug resistance efflux capabilities indicating potential selective factors
351 enabling maintenance in host cells.

352

353 *Differences of pMPPla107 and the Leaf58 plasmid*

354 pBASL58 is predicted to encode a complete CRISPR system from 229-241kb,
355 including two *cas*, three *csy* genes, and a repeat region that includes 36 repeats and

356 spacers(Figure 6). This CRISPR is located in the region of dissimilarity between
357 pMPPla107 and pBASL58 and is not found in pMPPla107. pBASL58 and pMPPla107
358 do share a (presumably) incomplete CRISPR systems at 436kb and 576kb
359 respectively (Figure 6). These regions include *cas3*, *csy3*, and *csy4* but lack *csy1*,
360 *csy2*. pMPPla107 lacks a repeat region altogether associated with this locus while
361 pBASL58 has a repeat region at 720kb encoding 9 repeats and spacers. To our
362 knowledge these are the first complete CRISPR systems located on plasmids found
363 within Proteobacteria.

364

365 There exists a region of dissimilarity across both megaplasms, occurring after
366 approximately 170kb (Figure 3), which could be classified as a cargo region. In
367 pMPPla107 this region consists of 468 predicted genes, of which 27 are annotated.
368 18 of these 27 annotated genes can be found in pBASL58 and again encode for genes
369 associated with DNA replication, repair, and metabolism. These genes also include
370 membranous proteins like FtsH, which is known to degrade unnecessary or
371 damaged membrane proteins ^{45,46}. We have also found that this region can largely
372 be deleted from pMPPla107 during lab adaptation (unpublished) even though the
373 rest of the plasmid is maintained. These data suggest that although this large region
374 may be expendable in some strains, pBASL58 has maintained many of the annotated
375 genes perhaps pointing to their importance in megaplasms stability or
376 maintenance.

377

378 **Discussion**

379 We report a family of divergent, yet syntenic megaplasmids found in single isolates
380 across distinct *Pseudomonas* species. High levels of synteny are matched by shared
381 signals in both tetranucleotide bias and protein pathway functionality. However,
382 these plasmids hosted by strains that are phylogenetically and geographically
383 separated; *Pla107* (containing pMPPla107) was found within a *P. syringae* isolate as
384 a causative agent of cucumber disease in Japan, while Leaf58 was found as an
385 epiphyte of *Arabidopsis* in Switzerland in a strain most closely related to *P. putida*⁴¹.
386 Furthermore, despite high levels of synteny and shared protein functionality,
387 consistently high levels of divergence across shared proteins ($\approx 30\%$) suggest both
388 plasmids have been independently evolving for a relatively long period of time.
389 From this data we infer that multiple additional members of a family of relatively
390 large ($\approx 1\text{Mb}$) “cryptic” megaplasmids likely persist within *Pseudomonas* strains.
391
392 That there have been no signs of these megaplasmids in the numerous sequences of
393 pseudomonads closely related to each of these isolates is strong indication that
394 these megaplasmids have been relatively recently acquired by their host strains.
395 This pattern, coupled with high levels of divergence between members of this
396 megaplasmid family, suggest that these replicons likely have a high turnover rate
397 within strains over evolutionary time and may persist within communities through
398 frequent horizontal transfer. In other words, presence of this megaplasmid family
399 may be transient in any given genome, but has likely been maintained within
400 Pseudomonads for a long time. Such a lifestyle is consistent with high levels of
401 conjugation as observed in pMPPla107 under laboratory conditions²⁵.

402

403 Replication, transcription, and translation of horizontally transferred genes are
404 known to incur costs on host cell resources with protein production likely having
405 the greatest effect on fitness^{15,47-49}. Previous work on pMPPla107 suggests that
406 acquisition of the megaplasmid results in lowered fitness and other phenotypic
407 changes which could be costly in some environments, yet it still transfers readily
408 and is maintained within host cells^{25,26}. Such costs could likely be the reason
409 pMPPla107 and pBASL58 encode a large number of genes involved in critical
410 functions regarding plasmid maintenance and transmission as well as potential
411 addiction systems and could enable long-term survival despite a transient lifestyle.
412 In particular, there are various proteins found in pMPPla107 and pBASL58 involved
413 in synthesizing precursors for nucleotides such as: thymidylate synthase, guanylate
414 kinase, ribonucleoside diphosphate reductase, deoxycytidine triphosphate
415 deaminase, and glutamate synthase (Supplemental Table 1 and Supplemental Figure
416 3). The megaplasmids may carry these proteins in order to increase flux to
417 nucleotide synthesis and drive replication and transcription processes to alleviate
418 any physiological costs an additional ≈ 1 Mb of newly acquired DNA may bring. Many
419 of these genes do not encode for complete pathways, indicating possible parasitic
420 behavior of host resources while ensuring the necessary building blocks for plasmid
421 maintenance are available.

422

423 Plasmid usage of host tRNA pools has been shown to deplete tRNAs resulting in
424 reduced growth and fitness^{20,50-52}. The large number of tRNAs and presence of a

425 handful of annotated tRNA ligases encoded on the megaplasמידs may serve the
426 purpose of avoiding translational costs due to tRNA depletion or may accommodate
427 codon usage bias between chromosome and megaplasמיד. Both megaplasמידs are
428 also predicted to encode Mfd, Rep, DnaB, and RecA all known to resolve replication
429 and transcription complex conflicts ensuring successful replicon duplication and
430 transcription^{53,54}. We hypothesize the megaplasמידs maximize their ability to
431 persist by eliminating or compensating for these potential costs by encoding a
432 variety of housekeeping genes coupled with high levels of horizontal transfer
433 through conjugation.

434

435 Evolutionary relationships between pBASL58 and pMPPla107, their relatively large
436 size and contribution to gene content of single strains, coupled with maintenance of
437 “housekeeping” genes, and high levels of transfer across pseudomonads suggest that
438 this megaplasמיד will provide unique insights into an evolutionary argument
439 concerning horizontal transfer referred to as the complexity hypothesis¹⁷. The
440 complexity hypothesis has been through multiple revisions, but is currently
441 interpreted as a trend where horizontally transferred genes are less likely to be
442 involved with complex processes (like translation) and maintain a lower number of
443 protein-protein interactions than vertically inherited loci⁴. One current limitation of
444 the complexity hypothesis, as highlighted by these megaplasמיד families, is that it
445 fails to reconcile gene conservation in the context of highly mobile selfish DNA like
446 plasmids. Both pBASL58 and pMPPla107 contain numerous “complex” genes,
447 including those involved in nucleotide synthesis, DNA replication, and translation

448 and yet these genes are clearly horizontally transferred across strains. Therefore,
449 the presented family of megaplastids potentially necessitates a caveat to the
450 complexity hypothesis in which “complex” genes can be horizontally transferred
451 frequently but aren’t maintained over time, because they are linked together on
452 megaplastids that require these pathways to ameliorate physiological costs.
453
454 Likewise, there have been numerous recent discussions about whether bacterial
455 pangenomes are adaptive or neutral. Similar to the complexity hypothesis, these
456 discussions tend to focus on the presence/absence of single genes across a variety of
457 closely related genomes rather than the linked gain/loss of genes that compose a
458 pangenome¹⁻³. To put this in perspective, recent findings suggest that the *P.*
459 *syringae* pangenome is composed of 77,728 genes, meaning that 1.5% of these are
460 solely present on pMPPI107⁵⁵. Since megaplastids have the potential to add
461 thousands of genes to a pangenome linked together in a single transfer event⁵⁶, one
462 has to consider that evolutionary pressures may act differentially on subsets of the
463 pangenome. Our data suggest that a majority of genes on these megaplastids may
464 be either neutral or costly to the host when selection is considered in the context of
465 the host genome. However, a majority of genes linked on the megaplastid may be
466 selectively beneficial for megaplastid maintenance and/or transfer regardless of
467 fitness of the host cell. Thus, presence of a majority of genes on the megaplastid
468 (and which are part of the pangenome) are under selection at some level, but only a
469 minority of these may be beneficial at the level of bacterial strains or populations.
470

471 CRISPR-Cas systems have become popularized recently because of their utility in
472 genome editing, however, these systems likely originated in bacteria as defense
473 mechanisms against invasion of foreign genetic material⁵⁷⁻⁶¹. CRISPR arrays are
474 often carried and transferred by larger plasmids in bacteria and archaea, yet *cas*
475 genes are rarely found on plasmids^{62,63}. Here we characterize a potentially shared
476 CRISPR-Cas system bound to the bacterial megaplasmids pMPPla107 and pBASL58.
477 Although pBASL58 encodes a fully intact CRISPR-Cas3 system with a region
478 containing 36 spacers and repeats, this repeat and spacer region are not present
479 within pMPPla107 leading us to believe pMPPla107's system is nonfunctional.
480 Regardless of functionality, it is quite interesting that at least one of these
481 megaplasmids contains an intact CRISPR locus given the widespread idea that these
482 systems are used by bacteria to defend against parasites and mobile elements.
483 Perhaps the presence of a CRISPR system is a beneficial and selective trait for
484 retention of pBASL58 in host cells in that it provides a transferable immune
485 pathway. However, the recent description of CRISPR spacers that target sites on
486 bacterial chromosomes also suggest that these loci may also function in gene
487 regulation⁶⁴⁻⁶⁷.

488

489 Using comparative computational and molecular approaches we have characterized
490 pBASL58, the second member of a family of large megaplasmids found in
491 Pseudomonads. Conservation of pathway presence and megaplasmid structure
492 strongly suggests that a majority of the sequences on pBASL58 and pMPPla107 have
493 diverged from a common ancestral plasmid. However, the consistent levels of

494 divergence between proteins shared by both plasmids suggest that this common
495 ancestral plasmid did not recently exist. Finding two related plasmids with such
496 high level of divergence also highlights the likelihood that other members of this
497 megaplasmid family exist in nature and are waiting to be found. Our work serves as
498 a guide to discover megaplasmid families as well as a foundation of understanding
499 the forces that structure megaplasmid evolution, maintenance, and transfer.

500

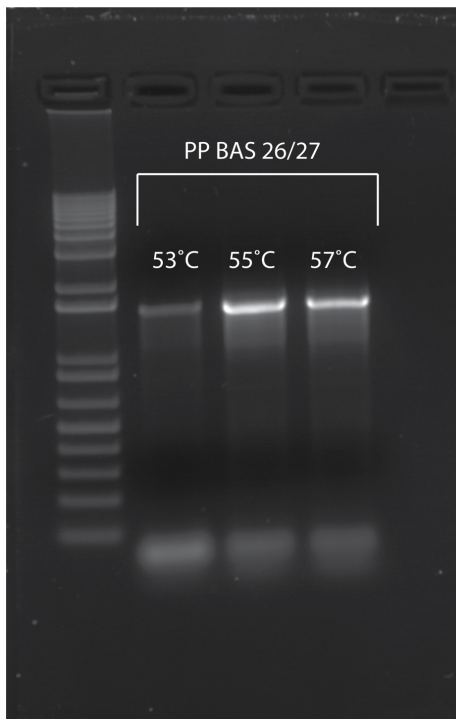
501

502

503

504

505



506

507 **Figure 1: Confirmation of circular DNA molecule of Leaf58 megaplasmid.**

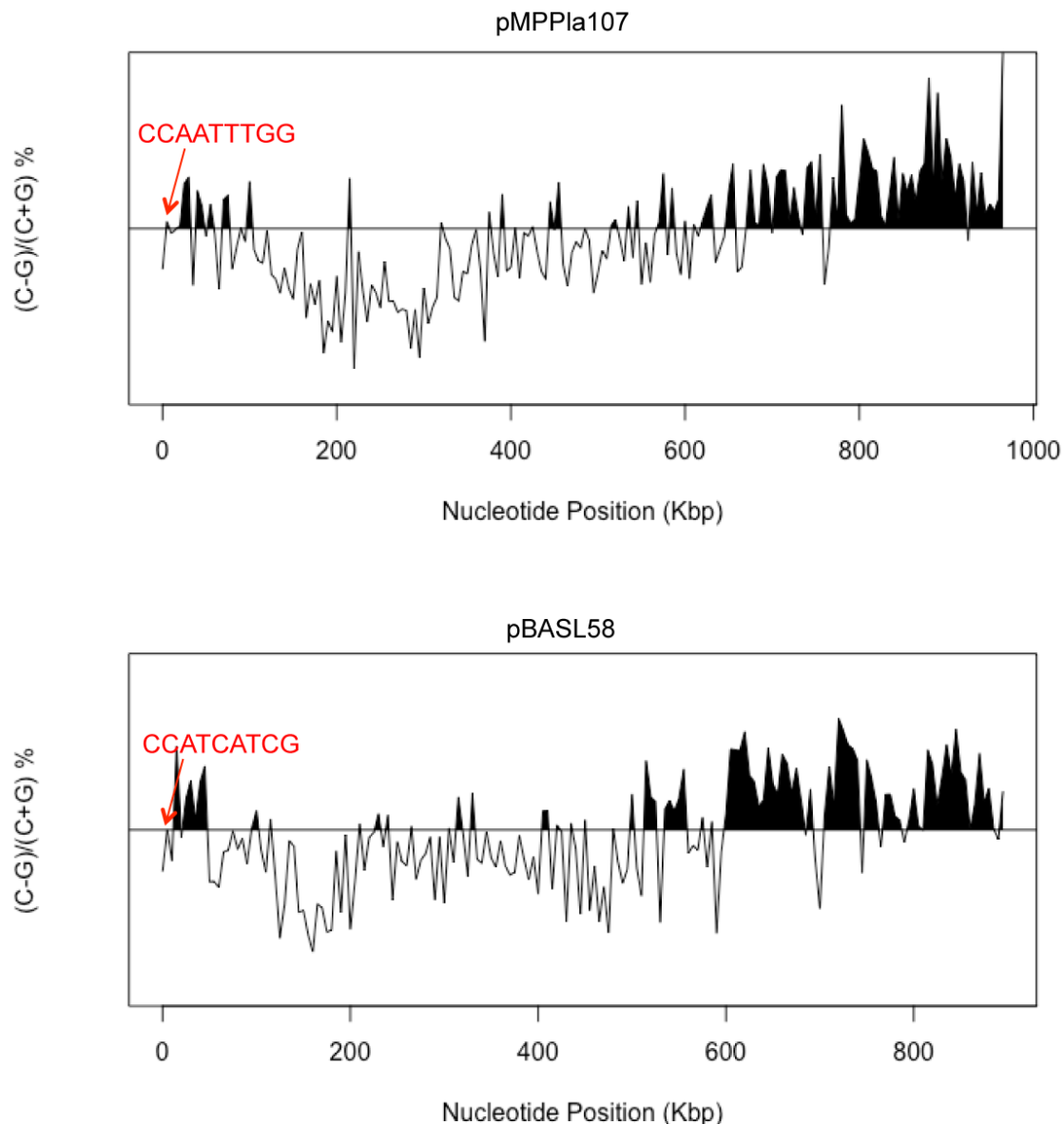
508 Primers designed to amplify the ends of the Leaf58 contig and disregarding the

509 misassembled repeat region successfully amplified products of an expected size.

510 Three annealing temperatures (53, 55, and 57°C) were used due to difficulties

511 amplifying this region.

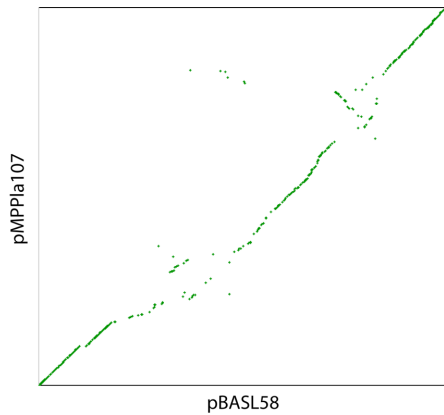
512



513

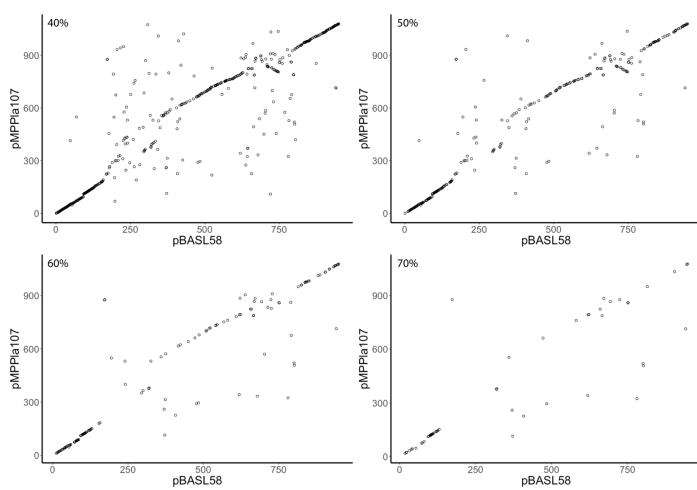
514 **Figure 2: Predicted origins of replication occur within a similar region of**
515 **pMPPla107 and pBASL58.** GC skew was calculated and is a known predictor of
516 origins of replication by a dramatic shift in GC content. Repetitive motifs were also
517 calculated for areas near the predicted origin of replication as repetitive binding
518 regions occur near replication sites. The most common motif is indicated in red.

519 **A)**



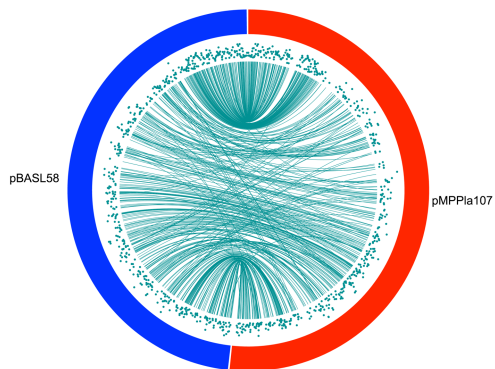
520

521 **B)**



522

523 **C)**



524

525 **Figure 3: pMPPla107 and pBASL58 share synteny and demonstrate divergence**

526 **on the amino acid level. A)** SynMap output of pMPPla107 vs. pBASL58 sequences

527 suggests highly syntenic megaplasmids. *X*-axis is pBASL58 gene order where $x_{1...N} =$

528 $gene_{1...N}$, and the *y*-axis is pMPPla107 gene order where $y_{1...N} = gene_{1...N}$. Completely

529 syntenic sequences would be represented by $y = 1x + b$. **B-C)** BLAST data was used to

530 plot pMPPla107 vs. pBASL58 synteny and amino acid divergence data together. **B)**

531 The majority of syntenic genes have $\geq 50\%$ sequence identity. BLAST data was

532 plotted in gene order to mimic SynMap's plot with amino acid sequence identity

533 cutoffs at 40%, 50%, 60%, and 70%. The best hit for each pBASL58 gene against

534 pMPPla107 is plotted. Each axis indicates gene position within the corresponding

535 sequence. **C)** Circa plot using BLAST data indicates higher synteny near the origin,

536 while areas near the terminus are less syntenic and experience more noise. Teal lines

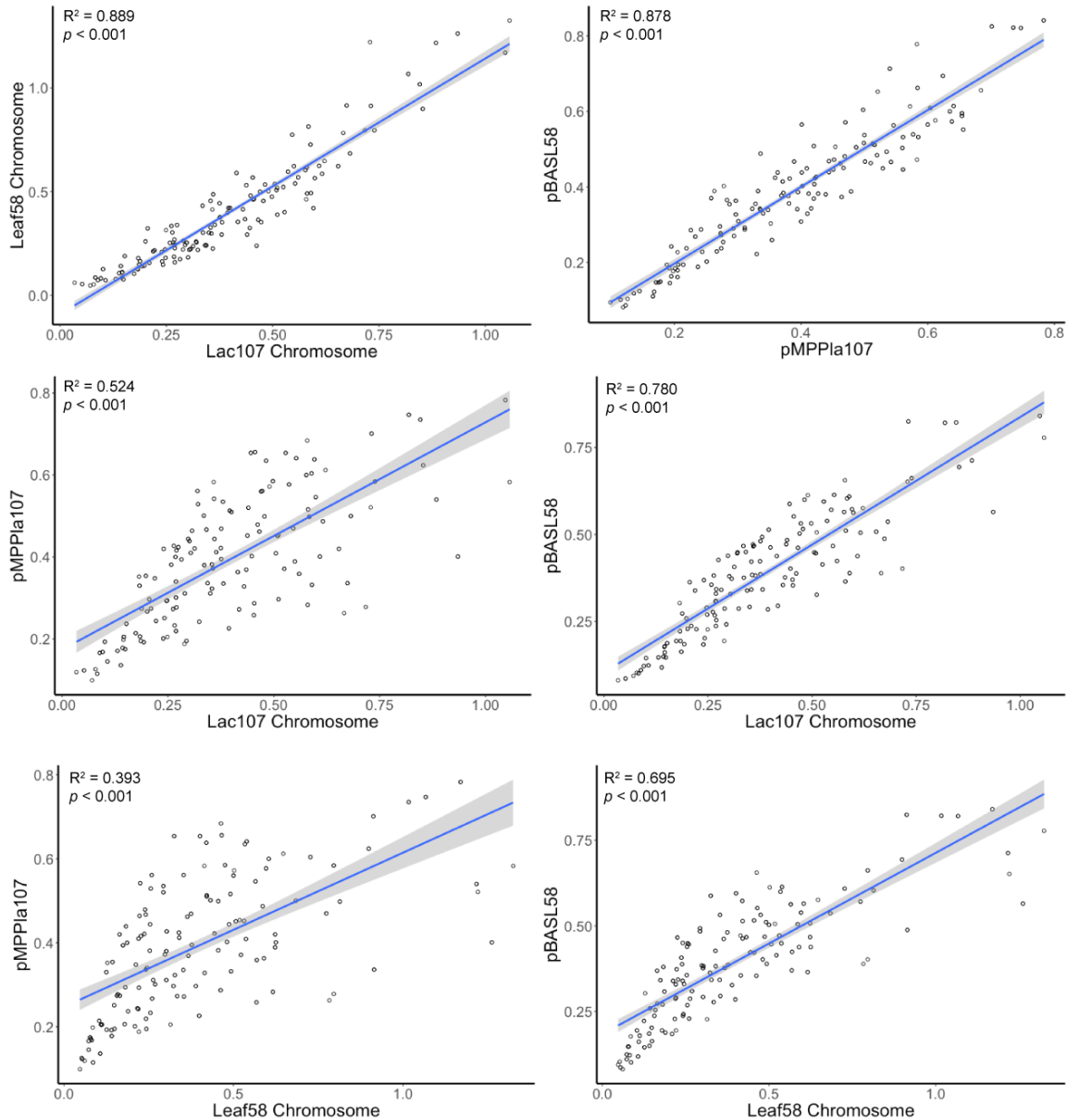
537 connect gene start position on pBASL58 to gene start position on pMPPla107. Teal

538 scatter plots are amino acid sequence identity with 40% = 0 (bottom) and 100% =

539 100 (top)

540

541



542

543 **Figure 4: Tetranucleotide frequencies between pMPPla107 and pBASL58**

544 **suggest an evolutionary relationship.** Nucleotide biases were determined to

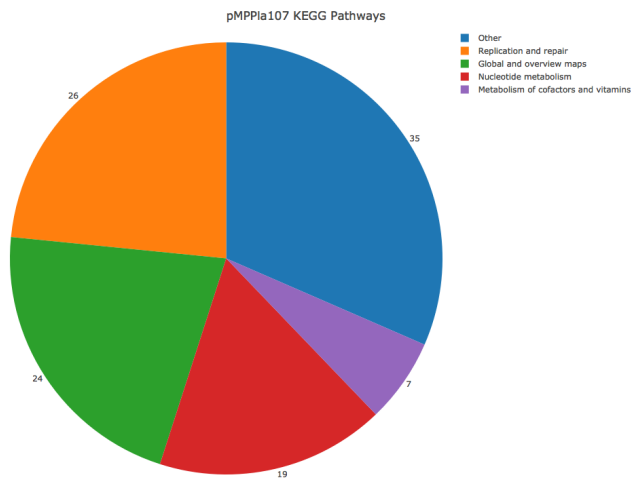
545 demonstrate relatedness of megaplasmid and chromosomal sequences in a pairwise

546 fashion. The blue line represents the linear regression model with the surrounding

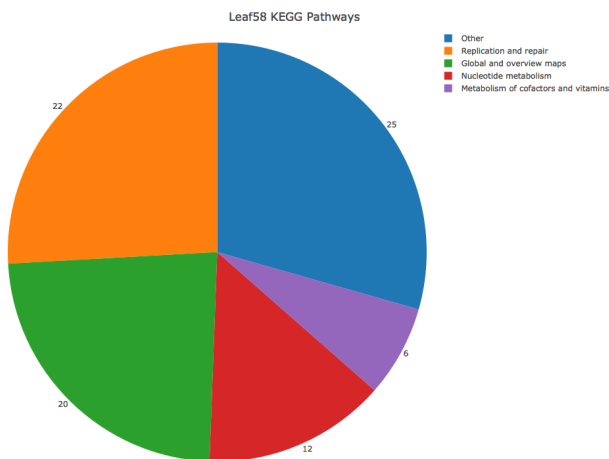
547 shaded grey area indicating a 95% confidence interval. R^2 and p values are listed for

548 each comparison.

549



550

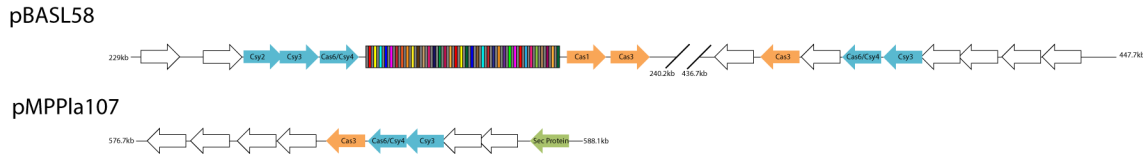


551

552 **Figure 5: pMPPla107 and pBASL58 share similar functional profiles.**

553 pMPPla107 and pBASL58 were input into UProC which counts the number of amino
554 acid sequences that are predicted to belong to a KEGG pathway IDs. KEGG IDs were
555 then matched with the correct pathway. All functional groups with < 5 counts were
556 grouped into an “other” category.

557



559

560 **Figure 6: CRISPR systems on pBASL58 and pMPPla107.** pBASL58 encodes two

561 CRISPR loci, one of which contains a repeat-spacer regions of 36 repeats.

562 pMPPla107 contains a CRISPR locus without any repeat-spacer regions. Direction of

563 arrows indicates gene orientation. Arrows are colored as: blue) *cys* genes, orange)

564 *cas* genes, green) secretion genes, and white) hypothetical genes. The multicolored

565 boxes indicate the repeat-spacer region, where grey boxes are spacers and colored

566 boxes are repeats.

567

568

569

Name	Size	Genes (CDS)	tRNA	GC Content
pMPPla107	971871	1082	54	52.84
pBASL58	903765	996	44	55.4
Leaf58 Concatenated	5378738	4847	80	62.35
Lac107 Concatenated	5936302	5436	58	58.26

570

571 **Table 1: General features of *P. syringae* and *Pseudomonas* Leaf58 replicons**

572 **have similarities.** Size, coding regions, tRNAs, and GC content were calculated to

573 understand the relationship between the four sequences on a broad scale.

574 “Concatenated” indicates that all sequences from the assemblies except either the

575 pMPPla107 or pBASL58 replicons to their respective genome were concatenated

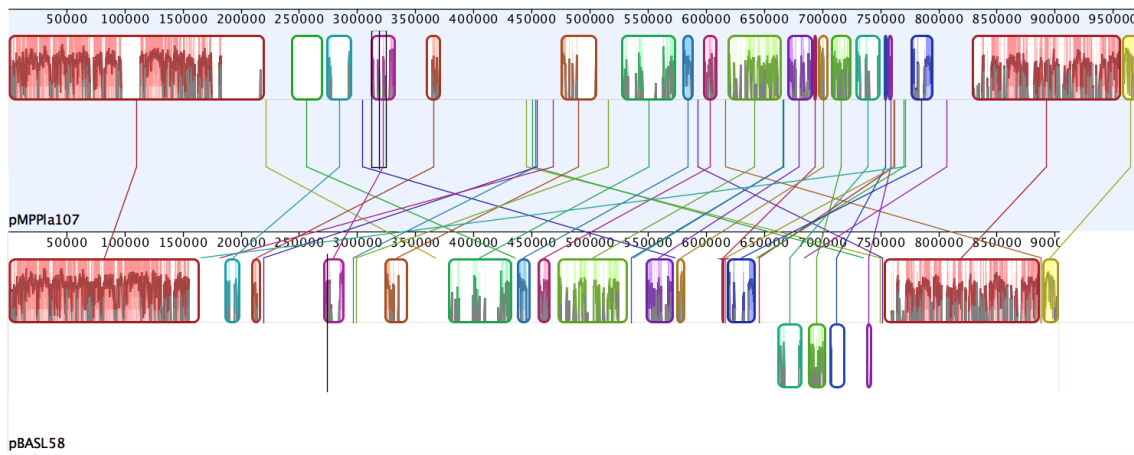
576 together.

577

578

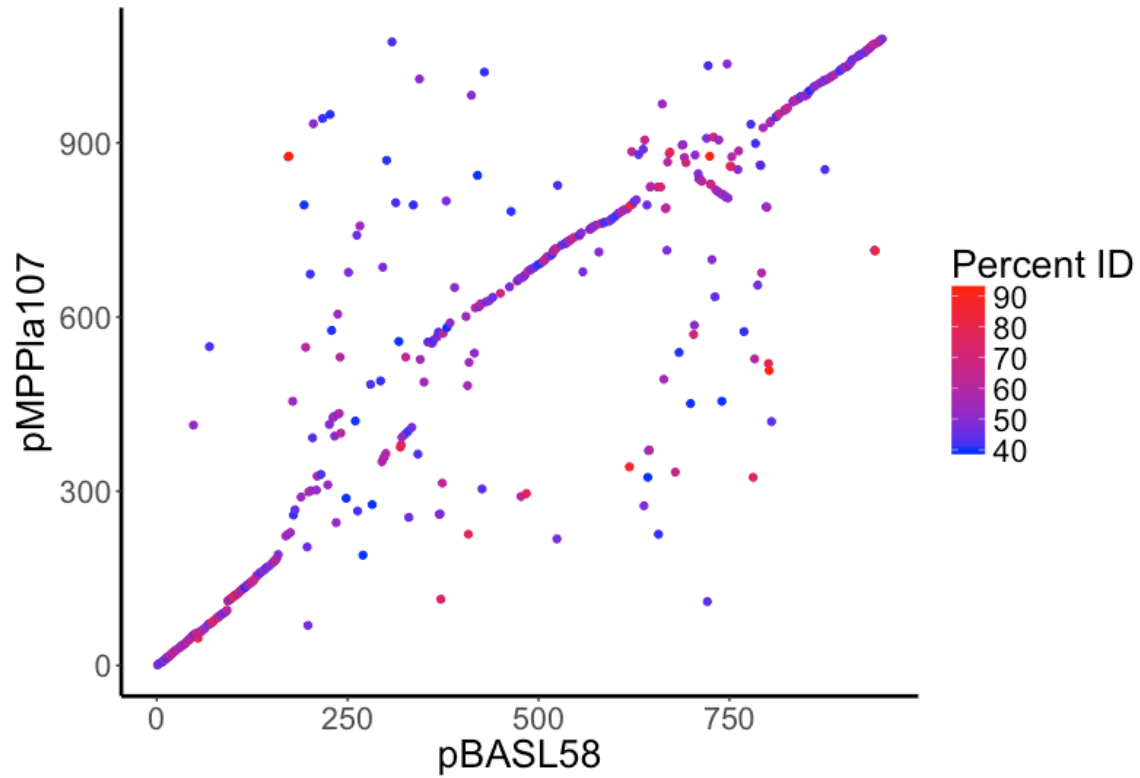
579

580 **Supplemental Figures**



581

582 **Supplemental Figure 1:** MAUVE analysis of pMPPla107 and pBASL58 demonstrate
583 local collinear blocks (LCBs) and areas of synteny. Lines connect LCBs with each
584 other between megaplasmids. Blocks below the midline for each sequence indicate
585 inverted regions. Colored areas within LCBs indicate higher levels of homology
586 between sequences.



587

588 **Supplemental Figure 2:** BLASTp results of pMPPla107 and pBASL58 indicate
589 syntenic and divergent sequences. Axes indicate gene position order. The color
590 gradient is set to BLASTp percent identity results for each comparison. A percent
591 identity cutoff of $\geq 40\%$ was used.

592

602

603 **References**

- 604 1. Vos, M. & Eyre-Walker, A. Are pangenomes adaptive or not? *Nature*
605 *Microbiology* **2**, 1576–1576 (2017).
- 606 2. McInerney, J. O., McNally, A. & O'Connell, M. J. Why prokaryotes have
607 pangenomes. *Nature Microbiology* **2**, 17040 (2017).
- 608 3. Shapiro, B. J. The population genetics of pangenomes. *Nature Microbiology* **2**,
609 1574–1574 (2017).
- 610 4. Cohen, O., Gophna, U. & Pupko, T. The complexity hypothesis revisited:
611 connectivity rather than function constitutes a barrier to horizontal gene
612 transfer. *Mol Biol Evol* **28**, 1481–1489 (2011).
- 613 5. Lercher, M. J. & Pál, C. Integration of horizontally transferred genes into
614 regulatory interaction networks takes many million years. *Mol Biol Evol* **25**,
615 559–567 (2008).
- 616 6. Wellner, A. & Gophna, U. Neutrality of foreign complex subunits in an
617 experimental model of lateral gene transfer. *Mol Biol Evol* **25**, 1835–1840
618 (2008).
- 619 7. Cooper, V. S., Vohr, S. H., Wrocklage, S. C. & Hatcher, P. J. Why genes evolve
620 faster on secondary chromosomes in bacteria. *PLOS Comput Biol* **6**,
621 e1000732 (2010).
- 622 8. Choudhary, M., Zanhua, X., Fu, Y. X. & Kaplan, S. Genome analyses of three
623 strains of *Rhodobacter sphaeroides*: evidence of rapid evolution of
624 chromosome II. *Journal of Bacteriology* **189**, 1914–1921 (2007).
- 625 9. Guo, H., Sun, S., Eardly, B., Finan, T. & Xu, J. Genome variation in the
626 symbiotic nitrogen-fixing bacterium *Sinorhizobium meliloti*. *Genome* **52**,
627 862–875 (2009).
- 628 10. Epstein, B. *et al.* Population genomics of the facultatively mutualistic
629 bacteria *Sinorhizobium meliloti* and *S. medicae*. *PLoS Genet.* **8**, e1002868
630 (2012).
- 631 11. Holden, M. T. G. *et al.* Genomic plasticity of the causative agent of
632 melioidosis, *Burkholderia pseudomallei*. *PNAS* **101**, 14240–14245 (2004).
- 633 12. Holden, M. T. G. *et al.* The genome of *Burkholderia cenocepacia* J2315, an
634 epidemic pathogen of cystic fibrosis patients. *Journal of Bacteriology* **191**,
635 261–277 (2009).
- 636 13. Janssen, P. J. *et al.* The complete genome sequence of *Cupriavidus*
637 *metallidurans* strain CH34, a master survivalist in harsh and anthropogenic
638 environments. *PLoS ONE* **5**, e10433 (2010).
- 639 14. diCenzo, G. C. & Finan, T. M. The Divided Bacterial Genome: Structure,
640 Function, and Evolution. *Microbiol. Mol. Biol. Rev.* **81**, e00019–17 (2017).
- 641 15. Baltrus, D. A. Exploring the costs of horizontal gene transfer. *Trends in*
642 *Ecology & Evolution* **28**, 489–495 (2013).
- 643 16. San Millan, A. & MacLean, R. C. Fitness Costs of Plasmids: a Limit to Plasmid
644 Transmission. *Microbiol Spectr* **5**, (2017).
- 645 17. Jain, R., Rivera, M. C. & Lake, J. A. Horizontal gene transfer among genomes:

- 646 the complexity hypothesis. *PNAS* **96**, 3801–3806 (1999).
- 647 18. MacLean, R. C. & San Millan, A. Microbial Evolution: Towards Resolving the
648 Plasmid Paradox. *Curr. Biol.* **25**, R764–7 (2015).
- 649 19. Kacar, B., Garmendia, E., Tuncbag, N., Andersson, D. I. & Hughes, D.
650 Functional Constraints on Replacing an Essential Gene with Its Ancient and
651 Modern Homologs. *MBio* **8**, e01276–17 (2017).
- 652 20. Harrison, E., Guymier, D., Spiers, A. J., Paterson, S. & Brockhurst, M. A. Parallel
653 compensatory evolution stabilizes plasmids across the parasitism-
654 mutualism continuum. *Curr. Biol.* **25**, 2034–2039 (2015).
- 655 21. Tett, A. *et al.* Sequence-based analysis of pQBR103; a representative of a
656 unique, transfer-proficient mega plasmid resident in the microbial
657 community of sugar beet. *ISME J* **1**, 331–340 (2007).
- 658 22. Hall, J. P. J. *et al.* Environmentally co-occurring mercury resistance plasmids
659 are genetically and phenotypically diverse and confer variable context-
660 dependent fitness effects. *Environ. Microbiol.* **17**, 5008–5022 (2015).
- 661 23. Baltrus, D. A. *et al.* Dynamic Evolution of Pathogenicity Revealed by
662 Sequencing and Comparative Genomics of 19 *Pseudomonas syringae* Isolates.
663 *PLoS Pathog* **7**, e1002132 (2011).
- 664 24. Andreani, N. A., Hesse, E. & Vos, M. Prokaryote genome fluidity is dependent
665 on effective population size. *ISME J* **11**, 1719–1721 (2017).
- 666 25. Romanchuk, A. *et al.* Bigger is not always better: transmission and fitness
667 burden of ~1MB *Pseudomonas syringae* megaplasmid pMPPla107. *Plasmid*
668 **73**, 16–25 (2014).
- 669 26. Dougherty, K. *et al.* Multiple phenotypic changes associated with large-scale
670 horizontal gene transfer. *PLoS ONE* **9**, e102170 (2014).
- 671 27. Bai, Y. *et al.* Functional overlap of the Arabidopsis leaf and root microbiota.
672 *Nature* **528**, 364–369 (2015).
- 673 28. Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving
674 bacterial genome assemblies from short and long sequencing reads. *PLoS*
675 *Comput Biol* **13**, e1005595 (2017).
- 676 29. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence
677 classification using exact alignments. *Genome Biol.* **15**, R46 (2014).
- 678 30. Tatusova, T., DiCuccio, M., acids, A. B. N.2016. NCBI prokaryotic genome
679 annotation pipeline. *academic.oup.com*
- 680
- 681 31. Charif, D. & Lobry, J. R. in *Structural Approaches to Sequence Evolution* 207–
682 232 (Springer, Berlin, Heidelberg, 2007). doi:10.1007/978-3-540-35306-
683 5_10
- 684 32. Abby, S. S., Néron, B., Ménager, H., Touchon, M. & Rocha, E. P. C.
685 MacSyFinder: a program to mine genomes for molecular systems with an
686 application to CRISPR-Cas systems. *PLoS ONE* **9**, e110726 (2014).
- 687 33. Grissa, I., Vergnaud, G. & Pourcel, C. CRISPRFinder: a web tool to identify
688 clustered regularly interspaced short palindromic repeats. *Nucleic Acids Res.*
689 **35**, W52–7 (2007).
- 690 34. Couvin, D. *et al.* CRISPRCasFinder, an update of CRISPRFinder, includes a
691 portable version, enhanced performance and integrates search for Cas

- 692 proteins. *Nucleic Acids Res.* **99**, 7536 (2018).
- 693 35. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
694 alignment search tool. *Journal of Molecular Biology* **215**, 403–410 (1990).
- 695 36. Darling, A. E., Mau, B. & Perna, N. T. progressiveMauve: multiple genome
696 alignment with gene gain, loss and rearrangement. *PLoS ONE* **5**, e11147
697 (2010).
- 698 37. Albertsen, M. *et al.* Genome sequences of rare, uncultured bacteria obtained
699 by differential coverage binning of multiple metagenomes. *Nature*
700 *Biotechnology* **31**, 533–538 (2013).
- 701 38. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes.
702 *Nucleic Acids Res.* **28**, 27–30 (2000).
- 703 39. Meinicke, P. UProC: tools for ultra-fast protein domain classification.
704 *Bioinformatics* **31**, 1382–1388 (2015).
- 705 40. Moriya, Y., Itoh, M., Okuda, S., Yoshizawa, A. C. & Kanehisa, M. KAAS: an
706 automatic genome annotation and pathway reconstruction server. *Nucleic*
707 *Acids Res.* **35**, W182–5 (2007).
- 708 41. Hesse, C. *et al.* Genome-based evolutionary history of *Pseudomonas* spp.
709 *Environ. Microbiol.* **7**, e1002132 (2018).
- 710 42. Richter, M. & Rosselló-Móra, R. Shifting the genomic gold standard for the
711 prokaryotic species definition. *Proc. Natl. Acad. Sci. U.S.A.* **106**, 19126–19131
712 (2009).
- 713 43. Nishida, H., Abe, R., Nagayama, T. & Yano, K. Genome Signature Difference
714 between *Deinococcus radiodurans* and *Thermus thermophilus*. *Int J Evol Biol*
715 **2012**, 205274–6 (2012).
- 716 44. Teeling, H., Meyerdierks, A., Bauer, M., Amann, R. & Glöckner, F. O.
717 Application of tetranucleotide frequencies for the assignment of genomic
718 fragments. *Environ. Microbiol.* **6**, 938–947 (2004).
- 719 45. Tomoyasu, T. *et al.* Escherichia coli FtsH is a membrane-bound, ATP-
720 dependent protease which degrades the heat-shock transcription factor
721 sigma 32. *EMBO J.* **14**, 2551–2560 (1995).
- 722 46. Ito, K. & Akiyama, Y. Cellular functions, Mechanism of action, and regulation
723 of *ftsH* protease. *Ann. Rev. Microbiol.* **59**, 211–231 (2005).
- 724 47. Hall, J. P. J., Brockhurst, M. A. & Harrison, E. Sampling the mobile gene pool:
725 innovation via horizontal gene transfer in bacteria. *Philos. Trans. R. Soc.*
726 *Lond., B, Biol. Sci.* **372**, 20160424 (2017).
- 727 48. Bragg, J. G. & Wagner, A. Protein material costs: single atoms can make an
728 evolutionary difference. *Trends Genet.* **25**, 5–8 (2009).
- 729 49. Shachrai, I., Zaslaver, A., Alon, U. & Dekel, E. Cost of unneeded proteins in *E.*
730 *coli* is reduced after several generations in exponential growth. *Molecular*
731 *Cell* **38**, 758–767 (2010).
- 732 50. Dittmar, K. A., Sørensen, M. A., Elf, J., Ehrenberg, M. & Pan, T. Selective
733 charging of tRNA isoacceptors induced by amino-acid starvation. *EMBO Rep.*
734 **6**, 151–157 (2005).
- 735 51. Elf, J., Nilsson, D., Tenson, T. & Ehrenberg, M. Selective charging of tRNA
736 isoacceptors explains patterns of codon usage. *Science* **300**, 1718–1722
737 (2003).

- 738 52. Bonomo, J. & Gill, R. T. Amino acid content of recombinant proteins
739 influences the metabolic burden response. *Biotechnol. Bioeng.* **90**, 116–126
740 (2005).
- 741 53. Hamperl, S. & Cimprich, K. A. Conflict Resolution in the Genome: How
742 Transcription and Replication Make It Work. *Cell* **167**, 1455–1467 (2016).
- 743 54. McGlynn, P., Savery, N. J. & Dillingham, M. S. The conflict between DNA
744 replication and transcription. *Mol. Microbiol.* **85**, 12–20 (2012).
- 745 55. Dillion, M. M., Thakur, S., Almeida, R. N. & Guttman, D. S. Recombination of
746 ecologically and evolutionarily significant loci maintains genetic cohesion in
747 the *Pseudomonas syringae* species complex. *bioRxiv* 227413 (2017).
748 doi:10.1101/227413
- 749 56. Nowell, R. W., Green, S., Laue, B. E. & Sharp, P. M. The extent of genome flux
750 and its role in the differentiation of bacterial lineages. *Genome Biol Evol* **6**,
751 1514–1529 (2014).
- 752 57. Doudna, J. A. & Charpentier, E. Genome editing. The new frontier of genome
753 engineering with CRISPR-Cas9. *Science* **346**, 1258096–1258096 (2014).
- 754 58. Hsu, P. D., Lander, E. S. & Zhang, F. Development and applications of CRISPR-
755 Cas9 for genome engineering. *Cell* **157**, 1262–1278 (2014).
- 756 59. Marraffini, L. A. & Sontheimer, E. J. CRISPR interference: RNA-directed
757 adaptive immunity in bacteria and archaea. *Nat. Rev. Genet.* **11**, 181–190
758 (2010).
- 759 60. Deveau, H., Garneau, J. E. & Moineau, S. CRISPR/Cas system and its role in
760 phage-bacteria interactions. *Annual Review of Microbiology* **64**, 475–493
761 (2010).
- 762 61. Makarova, K. S. *et al.* An updated evolutionary classification of CRISPR-Cas
763 systems. *Nat Rev Micro* **13**, 722–736 (2015).
- 764 62. Lillestøl, R. K. *et al.* CRISPR families of the crenarchaeal genus *Sulfolobus*:
765 bidirectional transcription and dynamic properties. *Mol. Microbiol.* **72**, 259–
766 272 (2009).
- 767 63. Godde, J. S. & Bickerton, A. The repetitive DNA elements called CRISPRs and
768 their associated genes: evidence of horizontal transfer among prokaryotes. *J.*
769 *Mol. Evol.* **62**, 718–729 (2006).
- 770 64. Guan, J., Wang, W., Sun, B. & Fey, P. D. Chromosomal Targeting by the Type
771 III-A CRISPR-Cas System Can Reshape Genomes in *Staphylococcus aureus*.
772 *mSphere* **2**, e00403–17 (2017).
- 773 65. Vercoe, R. B. *et al.* Cytotoxic chromosomal targeting by CRISPR/Cas systems
774 can reshape bacterial genomes and expel or remodel pathogenicity islands.
775 *PLoS Genet.* **9**, e1003454 (2013).
- 776 66. Briner, A. E. *et al.* Occurrence and Diversity of CRISPR-Cas Systems in the
777 Genus *Bifidobacterium*. *PLoS ONE* **10**, e0133661 (2015).
- 778 67. Stern, A., Keren, L., Wurtzel, O., Amitai, G. & Sorek, R. Self-targeting by
779 CRISPR: gene regulation or autoimmunity? *Trends Genet.* **26**, 335–340
780 (2010).
- 781