1    **Horizontal gene transfer to a defensive symbiont with a reduced genome**

2    **amongst a multipartite beetle microbiome**

3    Samantha C. Waterworth[a], Laura V. Flórez[b], Evan R. Rees[a], Christian Hertweck[c,d],

4    Martin Kaltenpoth[b] and Jason C. Kwan[a]#

5

6    Division of Pharmaceutical Sciences, School of Pharmacy, University of Wisconsin-

7    Madison, Madison, Wisconsin, USA[a]

8    Department of Evolutionary Ecology, Institute of Organismic and Molecular Evolution,

9    Johannes Gutenburg University, Mainz, Germany[b]

10    Department of Biomolecular Chemistry, Leibniz Institute for Natural Products Research

11    and Infection Biology, Jena, Germany[c]

12    Department of Natural Product Chemistry, Friedrich Schiller University, Jena, Germany[d]

13

14    #Address correspondence to Jason C. Kwan, jason.kwan@wisc.edu

15

16

17

18

**ABSTRACT**

The loss of functions required for independent life when living within a host gives rise to reduced genomes in obligate bacterial symbionts. Although this phenomenon can be explained by existing evolutionary models, its initiation is not well understood. Here, we describe the microbiome associated with eggs of the beetle *Lagria villosa*, containing multiple bacterial symbionts related to *Burkholderia gladioli* including a reduced-genome symbiont thought to produce the defensive compound lagriamide. We find that the putative lagriamide producer is the only symbiont undergoing genome reduction, and that it has already lost most primary metabolism and DNA repair pathways. The horizontal acquisition of the lagriamide biosynthetic gene cluster likely preceded genome reduction, and unexpectedly we found that the symbiont accepted additional genes horizontally during genome reduction, even though it lacks the capacity for homologous recombination. These horizontal gene transfers suggest that absolute genetic isolation is not a requirement for genome reduction.

**INTRODUCTION**

Mutualistic symbioses between animals and bacteria, widespread in nature, serve a variety of functions such as biosynthesis of nutrients not found in the host's diet (Akman et al., 2002; Shigenobu et al., 2000), and protection from predation (Lopera et al., 2017; Miller et al., 2016a; Piel, 2002) or infection (Currie et al., 1999; Flórez et al., 2018; Kroiss et al., 2010). Such relationships exist on a continuous spectrum of dependency and exclusivity from the perspective of both the host and symbiont. Symbionts generally become obligate after a prolonged period of exclusive association with the host, and the

2

42  symbionts that become obligate tend to carry out highly important functions for the host

43  (Latorre and Manzano-Marín, 2017; Lo et al., 2016; McCutcheon and Moran, 2012). For

44  example, mitochondria and chloroplasts, organelles that are required for energy

45  production and carbon fixation in eukaryotic and plant cells, originated from

46  endosymbiotic capture of alphaproteobacteria and cyanobacteria ~1.2 Bya and ~900

47  Mya, respectively (Shih and Matzke, 2013). The acquisition of these organelles allowed

48  the diversification of eukaryotic species (López-García et al., 2017). More recently,

49  aphids evolved to feed on plant sap depleted of several essential amino acids only

50  through capture of an endosymbiont, *Buchnera aphidicola*, that can synthesize these

51  nutrients, ~160–280 Mya (Moran et al., 1993; Munson et al., 1991). In these cases, the

52  microbial symbiont has lost the ability to live independently of the host, and the hosts

53  are also dependent on their symbionts.

54

55  The mechanism by which symbionts become obligate is through loss of genes required

56  for independent (but not host-associated) life, leading to an overall reduction in genome

57  size (Latorre and Manzano-Marín, 2017; Lo et al., 2016; McCutcheon and Moran,

58  2012). This gene loss is the result of relaxed selection on genes for functions provided

59  by the host, and also increased genetic drift as a result of small effective populations in

60  strict vertical transmission (Latorre and Manzano-Marín, 2017). A general mutational

61  bias towards deletion in bacteria (Mira et al., 2001) combined with many successive

62  population bottlenecks that allow the fixation of slightly deleterious mutations (Latorre

63  and Manzano-Marín, 2017) mediates general gene degradation and genome reduction.

64  While these processes are largely thought to be nonadaptive, there is some evidence

65    that increase in AT and reduction in genome size could be selected for to reduce

66    metabolic burden on the host (Dietel et al., 2019, 2018). The early stages of this

67    process are manifested by a proliferation of nonfunctional pseudogenes and a decrease

68    in coding density in the genome (Lo et al., 2016), before the intergenic sequences are

69    lost to eventually give tiny <1 Mbp genomes (Latorre and Manzano-Marín, 2017;

70    McCutcheon and Moran, 2012). While there is a robust model for the evolutionary

71    forces that drive this process once a symbiont becomes host-restricted, it is unknown

72    how bacteria first become obligate and start on the road to genome reduction (Latorre

73    and Manzano-Marín, 2017).

74

75    The beetle subfamily Lagriinae offers an opportunity to examine this question. Various

76    Lagriinae have evolved special symbiont-bearing structures that serve to facilitate the

77    vertical transmission of bacteria (Flórez et al., 2017). Beetles are typically co-infected

78    with multiple symbiont strains related to the plant pathogen *Burkholderia gladioli*, that

79    are secreted onto the surface of eggs as they are laid (Flórez et al., 2017). In the

80    species *Lagria villosa*, a South American soybean pest, at least one symbiotic *B.*

81    *gladioli* strain (Lv-StA) has been cultured and is still capable of infecting plants (Flórez

82    et al., 2017). The same strain produces antibacterial and antifungal compounds that can

83    protect the beetle's eggs from infection (Dose et al., 2018; Flórez et al., 2017). This is

84    consistent with the hypothesis that the *B. gladioli* symbionts evolved from plant

85    pathogens to become beetle mutualists. However, in field collections of *L. villosa*, Lv-

86    StA is only found sporadically, and is never highly abundant (Flórez and Kaltenpoth,

87    2017). Instead, the most abundant strain is often the uncultured Lv-StB (Flórez and

4

88     Kaltenpoth, 2017), which has been implicated in the production of the antifungal

89     lagriamide, a defensive compound found in field egg collections (Flórez et al., 2018).

90     We previously found through metagenomic sequencing that the genome of Lv-StB was

91     much smaller than that of Lv-StA, suggesting it has undergone genome reduction

92     (Flórez et al., 2018). It would seem that while *L. villosa* has multiple options for

93     symbionts that produce potential chemical defenses, only a subset have specialized as

94     obligate mutualists. The presence of multiple related strains in this system, with

95     selective genome reduction of a single strain, could potentially shed light on why the

96     genomes of some symbionts become reduced.

97

98     Here, we show that in the *L. villosa* microbiome, Lv-StB is uniquely undergoing genome

99     reduction, despite other community members possessing biosynthetic pathways for

100     potentially defensive molecules. We also suggest that this process was likely driven not

101     only by horizontal acquisition of the putative lagriamide pathway, but also by loss of

102     genes that limit cell division and translation, and gain of *zot*, a toxin also found in *Vibrio*

103     *cholerae* that aids invasion of host membranes. Further, we present evidence that these

104     horizontal gene transfers occurred concurrently with genome reduction, suggesting that

105     complete genetic isolation is not a main driving force for the reduction process.

106

107     **RESULTS AND DISCUSSION**

108     **Selective genome reduction of strain Lv-StB in *Lagria villosa***

109     We previously analyzed the metagenome of eight *L. villosa* egg clutches (Flórez et al.,

110     2018), using our binning pipeline Autometa (Miller et al., 2019). This method has the

5

111    advantage that it can separate noncharacterized eukaryotic contamination from

112    metagenomes, and it uses multiple factors (nucleotide composition, sequence

113    homology, the presence of single-copy marker genes and coverage) to accurately

114    produce bins from individual datasets. Because we had implemented several bugfixes

115    and small improvements to the pipeline since our original analysis, we re-ran Autometa

116    on the same metagenomic assembly. Despite some minor differences, the new bins

117    were broadly similar to our previous results (Dataset S1A), with 19 bins. As before, the

118    Lv-StB bin had the highest coverage, at 1,977×, such that the constituent contigs are

119    unlikely to be repeats from lower-coverage bins. We classified the bins according to a

120    new standardized bacterial taxonomy utilized by the Genome Taxonomy Database

121    (GTDB) that minimizes polyphyletic taxa and standardizes divergence between taxa of

122    the same rank (Parks et al., 2018). Notably, the GTDB taxonomy reclassifies

123    betaproteobacteria as being under class gammaproteobacteira. By this classification, all

124    bins were in class Gammaproteobacteria, in three different orders:

125    Betaproteobacteriales, Pseudomonadales and Xanthomonadales (Dataset S1B). The

126    most abundant bins were all in the family Burkholderiaceae, with the highest abundance

127    corresponding to the Lv-StB strain previously found to harbor the putative lagriamide

128    biosynthetic gene cluster (BGC) (Flórez et al., 2018). Interestingly, the average

129    nucleotide identity (ANI) of Lv-StB to the reference *B. gladioli* genome in GTDB (strain

130    ATCC 10248) is 85.7%, much lower than the 95% cutoff suggested for species

131    identifications (Goris et al., 2007) (Dataset S1B). This divergence suggests that Lv-StB

132    is a novel species in the genus *Burkholderia*, even though we previously classified it as

133    *B. gladioli* on the basis of 16S rRNA gene sequence (Flórez et al., 2018; Flórez and

134    Kaltenpoth, 2017), and therefore we refer to the strain here as "*Burkholderia* Lv-StB".

135    Likewise, most bins were found to be novel species, with one (DBSCAN_round2_3)

136    being divergent enough to be a representative of a novel genus in the family

137    Burkholderiaceae. Notably, the cultured *B. gladioli* strain that we have previously

138    isolated from *L. villosa* eggs, Lv-StA (Flórez et al., 2018, 2017), was not found to be

139    present in this metagenome.

140

141    The bins obtained had a range of different sizes (Table 1), which could be due to either

142    genome reduction or poor assembly and/or binning of a larger genome, which is often

143    observed if there are many related strains in a metagenome (Miller et al., 2019). As part

144    of the binning procedure, genome completeness was estimated based on the presence

145    of 139 single-copy marker genes (Rinke et al., 2013) (Dataset S1C). However, as some

146    complete genomes of genome-reduced symbionts have low apparent completeness by

147    this measure (Miller et al., 2017, 2016b), this figure cannot be used alone to determine

148    the size of an incompletely assembled genome. Conversely, even the drastically-

149    reduced genomes of intracellular obligate insect symbionts have been found to almost

150    universally maintain certain genes that we refer to here as 'core genes', involved in

151    replication, transcription, protein folding/stability, tRNA modification, sulfur metabolism,

152    RNA modification and translation (McCutcheon and Moran, 2012). We would expect,

153    therefore, that a well-assembled reduced genome would contain a near complete core

154    gene set, but not necessarily the whole set of 139 single-copy marker genes.

155    Conversely, incompletely assembled genomes are likely to be missing a significant

156    number of core genes that are required even in symbionts with highly reduced

157    genomes. We examined the presence of core genes in all metagenomic bins, as well as

158    *B. gladioli* Lv-StA for comparison (Dataset S1D and Table 1). Nine bins were close in

159    size to the genome of their respective closest relative, while maintaining most core

160    genes, and are classified as "nonreduced", and ten bins were small but also lacked a

161    significant fraction of core genes, and are classified as "incomplete". Only the Lv-StB

162    genome can be classified as reduced, on the basis of reduced size compared to its

163    close relative (2.07 Mbp, 23.5%) and maintenance of most core genes (85.7%).  This

164    bin exhibited additional features of genomes undergoing reduction, namely reduced

165    GC% compared to the *B. gladioli* reference genome (58.7% vs. 67.9%), and a

166    proliferation of pseudogenes accounting for 45.29% of the annotated ORFs (Dataset

167    S1C, S1E and Figure 1). Because of this, the Lv-StB genome exhibits a low coding

168    density (59.04%), and it also possesses a large number of ORFs containing

169    transposases (159). Both of these characteristics are hallmarks of symbionts in the

170    early stages of genome reduction, where there are high numbers of pseudogenes and

171    genome rearrangements (McCutcheon and Moran, 2012). The Lv-StB genome also

172    contains a low number of genes compared with its free-living relative, with 744 ORFs

173    that are not pseudogenes, transposases or hypothetical, versus 4,778 such genes in *B.*

174    *gladioli* Lv-StA.

175

176    **Diversity of biosynthetic gene clusters in the *L. villosa* microbiome**

177    Because we had previously isolated the non-reduced *B. gladioli* Lv-StA strain from *L.*

178    *villosa* (Flórez et al., 2018), and found it to produce protective compounds despite

179    having sporadic distribution and low abundance in field-collected beetles, we asked

8

180    whether BGCs were a unique feature of Lv-StB in the metagenome, or whether other

181    community members have the biosynthetic machinery for potential chemical defenses.

182    AntiSMASH (Blin et al., 2017) searches revealed a total of 105 BGCs in the

183    metagenome (Fig. 2), with variable BGC content in the bins, from zero to 566 kbp (0 to

184    14 BGCs, with 16 BGCs in the unclustered bin), while the *B. gladioli* Lv-StA genome

185    contained 1006 kbp BGCs (21 BGCs). This indicates that Lv-StB is not the only strain in

186    the egg microbiome with the potential to produce complex natural products, and many

187    strains harbor multiple BGCs. In the metagenome, bins in the family Burkholderiaceae

188    collectively contained the most BGCs by length (963 kbp), followed by a single bin in

189    family Rhodanobacteraceae (DBSCAN_round1_2, 566 kbp, Fig. 2). This distribution

190    suggests that although Burkholderiaceae appear to be an important reservoir of BGCs

191    in the *L. villosa* egg microbiome, other groups have significant biosynthetic potential.

192    Out of the 126 BGCs detected in the metagenome and in the *B. gladioli* Lv-StA genome,

193    there were 17 > 50 kbp in length predicted to produce complex nonribosomal peptides

194    or polyketides (Fig. 3). Two of these have been putatively assigned to production of the

195    antibiotic lagriene in Lv-StA (Flórez et al., 2017), and the antifungal lagriamide in Lv-StB

196    (Flórez et al., 2018), whereas five other small molecules known to be produced by Lv-

197    StA have been assigned to shorter BGCs (Dose et al., 2018; Flórez et al., 2017). That

198    15 out of the 17 largest assembled BGCs remain without characterized products

199    suggests that the majority of biosynthetic pathways in the *L. villosa* egg microbiome

200    likely codes for novel small molecule products. We compared the 126 identified BGCs

201    using BIG-SCAPE (Navarro-Muñoz et al., 2018), and found only 7 examples of BGCs

202    occurring in multiple strains, indicating that the biosynthetic potential in the metagenome

9

203    and *B. gladioli* Lv-StA is largely nonredundant. Taken together, this suggests that there

204    is a large amount of undefined biosynthetic potential for small molecule production in *L.*

205    *villosa* symbionts, beyond *B. gladioli* Lv-StA and *Burkholderia* Lv-StB.

206

**Divergence between Lv-StB and the closest free-living relative**

208    We constructed a phylogenetic tree of metagenomic bins assigned to the genus

209    *Burkholderia* as well as *B. gladioli* Lv-StA, based on 120 marker genes (Fig. 4A). This

210    showed that *Burkholderia* Lv-StB is most highly related to the *B. gladioli* clade but is

211    divergent from it. We calculated genome-wide ANI values for pairs of *Burkholderia*

212    genomes, and found *B. gladioli* strains shared between 97-100% ANI, while at most

213    *Burkholderia* Lv-StB shared 85.79% ANI with *B. gladioli* A1 (Dataset S1F). During

214    genome reduction, symbionts are known to undergo rapid evolution due to the loss of

215    DNA repair pathways (as found in the Lv-StB genome, see below) and the relaxation of

216    selection (McCutcheon and Moran, 2012), and so the divergence of Lv-StB from *B.*

217    *gladioli* may have been accelerated relative to free-living lineages. Genome-reduced

218    symbionts have often been vertically transmitted for evolutionary timescales and across

219    host speciation events, and therefore it is possible to calculate evolution rates where

220    related symbionts occur in hosts with known divergence times inferred from fossil

221    records. Such estimates in insect symbionts vary widely over three orders of magnitude,

222    but more recent ant and sharp-shooter symbiont lineages (established < 50 Mya for

223    "*Candidatus* Baumannia cicadellinicola", BAU; *Blochmannia obliquus*, BOB; *Bl.*

224    *pennsylvanicus*, BPN and *Bl. floridanus*, BFL) show high rates of divergence per

225    synonymous site per year (dS/t) between $1.1 \times 10^{-8}$ and $8.9 \times 10^{-8}$ (Silva and Santos-

226    Garcia, 2015) (the divergence rates used here are found in Table S1). Because of the

227    large number of pseudogenes in the Lv-StB genome, we reasoned that it is likely to be

228    a recent symbiont, and therefore used these rates to estimate divergence times

229    between Lv-StB and *B. gladioli* A1. We found a dS of 0.5486 per site between these

230    genomes, and calculated divergence times of 6.15, 8.55, 6.93 and 49.76 My for rates

231    BFL, BPN, BOB and BAU, respectively (Table S1). We should note here that these

232    figures are very approximate, and are possibly overestimates as symbiont evolution

233    rates are likely not constant, with particularly rapid evolution occurring during lifestyle

234    transitions (Lo et al., 2016). The range of these estimated divergence times suggests

235    that the common ancestor of *Burkholderia* Lv-StB and *B. gladioli* existed after the

236    evolution of symbiont bearing structures in Lagriinae beetles (see below).

237

238    We then sought to quantify the conservation of genes in Lv-StB compared to 13 related

239    *B. gladioli* strains, by identifying homologous gene groups among the entire set of non-

240    pseudogenes in these strains with OMA (Altenhoff et al., 2018). This pipeline aims to

241    identify orthologous groups while discounting paralogs among the genes in a given set

242    of genomes (Dessimoz et al., 2005). Of the 1,388 genes in Lv-StB that are not

243    pseudogenes (Dataset S1C), 492 were not included in any OMA orthologous group

244    (see below). A crude analysis of the OMA groups in Lv-StB, *B. gladioli* Lv-StA and *B.*

245    *gladioli* A1 revealed that Lv-StB retains a small subset of groups found in both Lv-StA

246    and A1, and has few unique groups (Fig. 4B), suggesting that Lv-StB has lost many of

247    the genes conserved in *B. gladioli*. Consistent with this notion, we visualized the

248    pangenome of *B. gladioli* and Lv-StB with Roary (Page et al., 2015) (Fig. 4C) and found

249    a large number of gene clusters that are conserved in *B. gladioli*, but not Lv-StB. The

250    gene clusters that are more variable amongst *B. gladioli* are also generally not found in

251    Lv-StB. Conversely, there were gene clusters found in Lv-StB that are not present in *B.*

252    *gladioli* strains, and these clusters may have been obtained by horizontal transfer after

253    the divergence of Lv-StB, or alternatively were lost in *B. gladioli*.

254

255    Remarkably, out of the 1,149 pseudogenes detected in the Lv-StB genome, 976 were

256    hypothetical and 129 were transposases, leaving only 44 that were recognizable

257    (Dataset S1G). This set of pseudogenes included a number of important genes in the

258    categories noted to be depleted below. For instance, the DNA polymerase I gene (*polA*)

259    appears to have been disrupted by a transposase, which is now flanked by two DNA pol

260    I fragments (E5299_1120 and E5299_01122). Likewise, the *uvrC* gene (E5299_00503),

261    a component of the nucleotide excision repair system (Lin and Sancar, 1992), is also

262    present as a truncated gene adjacent to a transposase. There were also pseudogenes

263    involved in the Entner-Doudoroff and glycolysis energy-producing pathways

264    (phosphogluconate dehydratase (Carter et al., 1993) and glucokinase (Lunin et al.,

265    2004)), as well as purine biosynthesis (phosphoribosylglycinamide formyltransferase

266    (Almassy et al., 1992), phosphoribosylformylglycinamidine synthase (Schendel et al.,

267    1989)). Interestingly, we found two pseudogenes that negatively regulate cell division

268    and translation. Septum protein Maf is a nucleotide pyrophosphatase that has been

269    shown to arrest cell division, especially after transformation or DNA damage

270    (Tchigvintsev et al., 2013). Deletion of the *E. coli* gene for homolog YhdE increased

271    growth rate, while overexpression decreased growth rate (Jin et al., 2015). Therefore

272    the loss of Maf in Lv-StB would be expected to increase the rate of cell division and

273    reduce the conversion of nucleotides, which it probably obtains from the host, to the

274    monophosphates. The gene for the energy-dependent translational throttle protein EttA

275    was also found to be truncated. This protein slows translation through interacting with

276    the ribosome in both the ATP- and ADP-bound forms (Boël et al., 2014; Chen et al.,

277    2014). Under energy-depleted conditions (i.e. high ADP), EttA was found to stabilize

278    ribosomes and prevent commitment of metabolic resources, and thus the deletion

279    mutant displayed reduced fitness during extended stationary phase (Boël et al., 2014).

280    However, under circumstances where the host supplies ample nucleotides, the loss of

281    *ettA* would be expected to increase translation rates.

282

283    **Degradation of primary metabolic pathways in *Burkholderia* Lv-StB**

284    Lv-StB is deficient in many metabolic pathways that are complete in related *B. gladioli*

285    strains (Fig. 5), including the glyoxylate shunt (Dolan and Welch, 2018), various carbon

286    degradation pathways, mixed acid fermentation, as well as sulfur and nitrogen

287    metabolism. Although the extent of gene losses could be overestimated due to the draft

288    status of the Lv-StB genome, the pervasiveness of metabolic gaps found combined with

289    the high coverage of the genome (Dataset S1A) suggest generalized gene loss in many

290    functional categories. Lv-StB appears incapable of making any of the following

291    compounds due to the absence of several biosynthetic genes: Thiamine, riboflavin,

292    nicotinate, pantothenate, vitamin B12 and biotin. Likewise, there were deficiencies in

293    amino acid biosynthesis (Fig. S1). We predict that Lv-StB would be able to make

294    chorismate, isoleucine, leucine, ornithine, proline and threonine, but likely lacks the

13

295     ability to make aromatic amino acids, serine, methionine, lysine, histidine, cysteine,

296     glutamine and arginine due to the absence of several key genes in these pathways. The

297     genome of Lv-StB also lacks genes involved in chemotaxis and flagella, suggesting that

298     after the symbiont mixture is spread on eggs, the colonization of the dorsal cuticular

299     structures in the embryo (Flórez et al., 2017) does not require symbiont motility.

300     Interestingly, the Lv-StB genome includes a trimeric autotransporter adhesin (TAA)

301     related to SadA (Raghunathan et al., 2011), which is involved in the pathogenicity of

302     *Salmonella typhimurium* by aiding cell aggregation, biofilm formation, and adhesion to

303     human intestinal epithelial cells. TAAs are found in Proteobacteria and consist of

304     anchor, stalk and head domains, of which the head forms the adhesive component

305     (Linke et al., 2006). The bacterial honey-bee symbiont, *Snodgrassella alvi* is

306     hypothesized to utilize TAAs in combination with other extracellular components during

307     colonization of the host gut (Powell et al., 2016), and similar genes were identified in *S.*

308     *alvi* symbionts in bumble bees and are predicted to perform a similar role (Kwong et al.,

309     2014). Therefore, this gene may play a role in the adhesion of Lv-StB cells to *L. villosa*

310     eggs.

311

312     The Lv-StB genome is also missing several enzymes within glycolysis, most notably

313     phosphoglycerate kinase, phosphoenolpyruvate carboxykinase and others, which would

314     suggest that Lv-StB has lost the ability to perform glycolysis. The loss of glycolysis is

315     often substituted by an alternative pathway, such as the pentose phosphate or Entner-

316     Doudoroff pathway (Chen et al., 2016). This is not the case for Lv-StB, where both the

317     glucose-6-phosphate dehydrogenase and 6-phosphogluconolactonase genes appear to

318     be missing from the genome. The citrate cycle is largely complete, except that it is

319     missing pyruvate carboxylase, the enzyme that converts pyruvate to oxaloacetate.

320     However, phosphoenolpyruvate carboxylase is present in Lv-StB (E5299_00983) and

321     may alternatively be used for the production of oxaloacetate from exogenous

322     phosphoenol-pyruvate (Takeya et al., 2017) as alternative pathways for the supply of

323     oxaloacetate are also incomplete. Lv-StB is also missing all genes required to form

324     cytochrome c oxidase and cytochrome b-c complexes. However, all genes required to

325     encode NADH:quinone oxidoreductase (Complex I), succinate dehydrogenase

326     (Complex II) and cytochrome o ubiquinol oxidase are present, along with all genes

327     required for the F-type ATPase. The lack of Complex III would likely result in a

328     decreased rate of ATP production in Lv-StB as observed in fungi with alternative

329     oxidases that bypass Complexes III and IV (Duarte and Videira, 2009). Lv-StB may be

330     similar to the psyllid endosymbiont "*Candidatus* Liberibacter asiaticus", which has lost

331     both key glycolysis and glyoxalase genes and instead relies on the scavenging of ATP

332     from the host (Jain et al., 2017).

333

334     We also found many deficiencies in both the *de novo* and salvage nucleotide pathways

335     (Fig. S2). In the pyrimidine biosynthetic pathway, genes for dihydrooratase (*pyrC*)

336     (Porter et al., 2004) and orotate phosphoribosyl transferase (*pyrE*) (Aghajari et al.,

337     1994) were missing, suggesting that Lv-StB cannot produce orotate from *N*-

338     carbamoylaspartate, and cannot create nucleotides from free pyrimidine bases.

339     Ribonucleotide reductase (*nrdAB*) (Brignole et al., 2012) and thymidylate synthase

340     (*thyA*) (Carreras and Santi, 1995) are present, suggesting that deoxypyrimidine

15

341    nucleotides can be made from CTP. The deficiencies in purine synthesis were more

342    profound. The majority of the *de novo* pathway (Nelson et al., 2008) was missing

343    (*purCDEFHLMNT*, IMP dehydrogenase/*guaB*), except for adenylosuccinate lyase

344    (*purB*), adenylosuccinate synthase (*purA*) and GMP synthase (*guaA*). Lv-StB should

345    therefore be able to make AMP from IMP, and GMP from XMP (plus their deoxy

346    analogs through ribonucleotide reductase), but cannot make purines *de novo*. We were

347    also not able to find adenine phosphoribosyltransferase or hypoxanthine-guanine

348    phosphoribosyltransferase (Nelson et al., 2008), meaning that purine bases cannot be

349    salvaged to make nucleotides.

350

351    **Degradation of DNA repair pathways in *Burkholderia* Lv-StB**

352    The genome of Lv-StB is missing many genes involved in DNA repair (Fig. 5B), similar

353    to other examples of genome-reduced symbionts (McCutcheon and Moran, 2012).

354    Compared to closely-related *B. gladioli* strains, Lv-StB lacks genes in every repair

355    pathway. In particular DNA polymerase I (*polA*), used in homologous recombination,

356    nucleotide excision repair and base excision repair, is only present as two truncated

357    pseudogenes (see above). Even though *polA* is involved in many different DNA repair

358    pathways, it has been found to be nonessential in *Escherichia coli* (Gerdes et al., 2003;

359    Goodall et al., 2018), *B. pseudomallei* (Moule et al., 2014) and *B. cenocepacia* (Higgins

360    et al., 2017). In the homologous recombination pathway, Lv-StB lacks *recA*, *polA*, *ruvA*,

361    *ruvB*, *ruvC* and *recG*, all of which have been found to be essential for homologous

362    recombination in *E. coli* (Kowalczykowski et al., 1994). Likewise, Lv-StB is also missing

363    the two components of the nonhomologous end-joining pathway, *ku* and *ligD* (Pitcher et

16

364    al., 2007), suggesting that it cannot recover from double-strand breaks. In the base

365    excision repair pathway, Lv-StB lacks several DNA glycosylases which are responsible

366    for removing chemically modified bases from double-stranded DNA (McCullough et al.,

367    1999). Some of these losses simply reduce redundancy, but it has also lost

368    nonredundant glycosylases *mutM* and *mug*. The former recognizes 2,6-diamino-4-

369    hydroxy-5-*N*-methylformamidopyrimidine (Fapy) and 8-hydroxyguanine (Boiteux et al.,

370    1992), while the latter recognizes G:U and G:T mismatches (Barrett et al., 1998) as well

371    as epsilonC (Saparbaev and Laval, 1998), 8-HM-epsilonC (Hang et al., 2002), 1,N(2)-

372    epsilonG (Saparbaev et al., 2002) and 5-formyluracil (Liu et al., 2003). Finally, in the

373    mismatch repair system, Lv-StB is missing *mutS*, which is required for the recognition of

374    mismatches in methyl-directed repair (Schofield and Hsieh, 2003). In summary, Lv-StB

375    is likely to be completely incapable of nonhomologous end-joining, homologous

376    recombination, and mismatch repair, while being impaired in nucleotide excision repair

377    and base excision repair due to the loss of DNA polymerase I and several DNA

378    glycosylases.

379

380    **Timing of horizontal acquisition of defensive and other genes in the Lv-StB**

381    **genome**

382    We then attempted to identify genes in the *Burkholderia* Lv-StB genome that are likely

383    to have been acquired by horizontal gene transfer (HGT). A total of 497 non-

384    pseudogenes were identified as unique to Lv-StB, and following removal of genes with

385    no matches against the BLAST NR database, and removal of genes that were

386    homologous to other *B. gladioli* genomes not included in this study, there were 148

17

387    genes that appeared to be more closely related to species other than *B. gladioli*, that

388    have potentially been acquired through horizontal transfer (Figure 6A). Most genes

389    appear to have been obtained from gammaproteobacteria and alphaproteobacteria, with

390    a small number from firmicutes, cyanobacteria and phages (see Dataset S1H). The

391    distribution is consistent with the notion that horizontal transfer occurs most frequently

392    between closely-related species (Gillings, 2017). In particular, Burkholderiaceae was

393    the most frequent apparent donor amongst gammaproteobacterial proteins.

394    Interestingly, the alphaproteobacterial genus *Ochrobactrum* (family Brucellaceae in the

395    NCBI taxonomy, Rhizobiaceae in GTDB) was a major putative gene donor. This genus

396    includes several symbionts of termites (Mathew et al., 2012), army worms (Jones et al.,

397    2019), weevils (Montagna et al., 2015) and leeches (McClure et al., 2019; Rio et al.,

398    2009). In previous 16S amplicon investigations of *Lagria* beetles, *Ochrobactrum* were

399    often found (Flórez et al., 2017; Flórez and Kaltenpoth, 2017) (Fig. S3), suggesting that

400    the donors of these genes could have also been associates of *L. villosa*. *Ochrobactrum*

401    OTUs account for 5–20% of 16S rRNA gene reads, but this genus was not observed in

402    the shotgun metagenome. However, disparities between 16S and shotgun metagenome

403    abundances are not uncommon due to variable 16S copy number, primer and

404    sequencing biases (Chen et al., 2017; Delforno et al., 2017). Based on the evidence for

405    putatively horizontally transferred genes, we asked whether these could have

406    contributed to the dominance of Lv-StB in *L. villosa*, and set out to estimate the timings

407    of horizontal transfer events.

408

409    Horizontally transferred genes are often detected on the basis of nucleotide composition

410    differing from other genes in the genome (Becq et al., 2010). Such genes initially exhibit

411    nucleotide composition consistent with the donor genome, which over time will

412    eventually normalize to the composition of the recipient genome (Lawrence and

413    Ochman, 1997). The rate of this "amelioration" process ($\Delta GC^{HT}$) has been modeled by

414    Lawrence and Ochman (Lawrence and Ochman, 1997), based on the substitution rate

415    ($S$), the transition/transversion ratio ($\kappa$) and GC content of both the recipient genome

416    ($GC^{EQ}$) and putatively horizontally transferred genes ($GC^{HT}$), according to equation 1.

417    By iterating this equation repeatedly until $GC^{HT}$ equals $GC^{EQ}$, the time required from the

418    present day to complete amelioration can be estimated. If the GC content of the donor

419    genome is known, then equation 1 can be used in reverse to estimate the time since

420    introgression. However, if the donor GC content is not known, then the differing

421    selection pressures on the first, second and third codon positions can be exploited to

422    estimate the introgression time. Because these positions have different degrees of

423    amino acid degeneracy, they are subject to different degrees of selection, and therefore

424    they ameliorate at different rates. As a consequence, Lawrence and Ochman (Lawrence

425    and Ochman, 1997) found that for genes in the process of amelioration, the relationship

426    between overall GC content and the GC content at individual codon positions seen in

427    genes at equilibrium (Lawrence and Ochman, 1997; Muto and Osawa, 1987) (Equations

428    2, 3 and 4) does not hold. So if equation 1 is applied in reverse separately for each

429    codon position, the time since introgression can be inferred at the iteration yielding the

430    minimum square difference from equations 2–4. Application of equation 1 also yields an

431    estimate for the original donor GC.

19

432

433
$$\Delta GC^{HT} = S \times \frac{\kappa + \frac{1}{2}}{\kappa + 1} \times [GC^{EQ} - GC^{HT}] \qquad (1)$$

434

435
$$GC_{1st} = 0.615 \times GC_{Genome} + 26.9 \qquad (2)$$

436

437
$$GC_{2nd} = 0.270 \times GC_{Genome} + 26.7 \qquad (3)$$

438

439
$$GC_{3rd} = 1.692 \times GC_{Genome} - 32.3 \qquad (4)$$

440

441 We identified groups of consecutive genes in the putative HGT set that could have been

442 acquired together, and used the above method to estimate their introgression time

443 (Dataset S1H). Out of the 18 identified gene groups, 7 were found to have atypical GC

444 content (defined by Lawrence and Ochman (Lawrence and Ochman, 1997) as either

445 >10% lower or >8% higher GC% in first or third codon positions compared to the

446 genome as a whole). The method above was used to estimate time of introgression for

447 each gene group, using the BFL, BPN, BOB and BAU divergence rates (see above, Fig.

448 6B, Dataset S1H). The oldest horizontal transfer was found to be group "hypo4" (1.79–

449 14.36 Mya), representing two phage proteins, and the next oldest is the lagriamide BGC

450 and neighboring genes (*lga*, 0.8–6.42 Mya). The *lga* BGC is predicted to have come

451 from a high-GC organism, with an original GC content of 72%. The closest relatives of

452 many of the *lga* genes are found in *Pseudomonas* strains, which typically do not have

453 GC contents this high. However, as BGCs are often thought to be horizontally

454 transferred (Jensen, 2016), *Pseudomonas* may not be the direct source of *lga* in Lv-StB.

455

456    Several other HGT gene groups were predicted to be involved in transport

457    (sugar_transp, bmp, ext, dinj, tonb, tonb2), with predicted substrates including purine

458    nucleosides, trehalose and vitamin B12. Of these, the sugar_transp and bmp groups

459    predicted to be involved in trehalose and purine import, respectively, are relatively old

460    (0.41–3.31 My and 0.27–2.16 My), and the groups likely involved in B12 import (dinj,

461    tonb, tonb2) are estimated to be <5,000 y old. Trehalose is the most abundant

462    component of insect hemolymph, and in a leaf beetle system was found to be

463    provisioned by the host to its genome-reduced symbiont (Bauer et al., 2019). We also

464    found HGT gene groups putatively involved in heat shock response (hsp20, 0.05–0.4

465    Mya), and DNA repair (ura, rep, ochro). In the latter category was a uracil-DNA

466    glycosylase (in ura, 0.21–1.68 Mya) used in base excision repair in case of

467    deamination, RecB used in homologous recombination (in rep, 0.01–0.055 Mya)

468    (Nelson et al., 2008) and YedK (in ochro, 0.01-0.035 Mya), a protein used in the SOS

469    response that binds to abasic sites in single strand DNA (Mohni et al., 2019). The

470    transferred transport functions and DNA repair proteins match functional categories that

471    are currently lacking genes due to genome reduction (see above), and therefore the

472    transfers could have been contemporaneous with the reduction process, acting as

473    compensatory mechanisms for lost functions.

474

475    One putative HGT toxin that is unique to Lv-StB in the metagenome and amongst *B.*

476    *gladioli* strains is zonular occludens toxin (zot, estimated as acquired <5,000 ya). The

477    *zot* gene is responsible for the production of the zonula occuludens toxin, a virulence

478   factor which was initially identified in *Vibrio cholera* and was found to lead to the

479   disassembly of intracellular tight junctions, leading to increased permeability of

480   mammalian epithelium (Di Pierro et al., 2001). Co-localized with the predicted zot gene

481   was a gene encoding DUF2523, which we found often accompanied *zot* in searches of

482   the STRING database (Szklarczyk et al., 2017). Zot proteins have been identified in

483   several strains of *Camplylobacter* and have been shown to elicit an inflammatory

484   response in intestinal epithelial cells (Liu et al., 2016; Mahendran et al., 2016).

485   Furthermore, a significant correlation was found between the presence of the Zot

486   protein and hyper-invasive strains of *Neisseria meningitides* (Joseph et al., 2011).

487   Potentially, *zot* may aid in the infection of the *L. villosa* embryonic structures through

488   increasing permeability across the outer layers of the egg, although this remains

489   speculative.

490

491   **CONCLUSIONS**

492   In the 1920s it was observed (Jürgen Stammer, 1929) that beetles in the *Lagria* and

493   *Cerogria* genera contained structures now known to harbor *Burkholderia* symbionts in *L.*

494   *villosa* and *L. hirta* (Flórez and Kaltenpoth, 2017). Other genera in the Lagriinae

495   subfamily, such as *Adynata* and *Arthromacra*, do not contain these structures.

496   According to the tree found at timetree.org (Kumar et al., 2017), *Lagria* and *Cerogria*

497   diverged 55 Mya, and the common ancestor of *Lagria*, *Cerogria* and *Adynata* existed 82

498   Mya (this region of the tree utilizes data from Kergoat *et al*. (Kergoat et al., 2014)).

499   Based on these estimates, the symbiont-bearing structures in *Lagria* and *Cerogria* likely

500   evolved between 82 and 55 Mya. Our analysis suggests that the divergence of

22

501    *Burkholderia* Lv-StB from *B. gladioli* occurred after that point (6.15–49.76 Mya). During

502    that time, the genome of *Burkholderia* Lv-StB became reduced, and it is likely

503    dependent on the host due to deficiencies in energy metabolism and nucleotide

504    biosynthesis. Notably, the profound metabolic insufficiencies and incomplete DNA repair

505    pathways in Lv-StB are typical of symbionts with smaller genomes, such as "*Ca.*

506    Endolissoclinum faulkneri", an intracellular tunicate symbiont with a 1.48 Mbp genome

507    and a similar number of genes (783) (Kwan et al., 2012), estimated to have been a

508    symbiont for at least 6-31 My (Kwan and Schmidt, 2013). While the presence of

509    *Burkholderia* Lv-StB and its defensive compound lagriamide has been shown to

510    decrease the rate of fungal egg infection (Flórez et al., 2018), the symbiont is not

511    essential for beetle reproduction (Flórez et al., 2017). Therefore, the relationship is

512    facultative from the perspective of the host, while Lv-StB is in the process of becoming

513    dependent on *L. villosa*. A central question we aimed to answer in this work was how

514    the genome of Lv-StB became reduced, when *L. villosa* appears to maintain multiple

515    other nonreduced *Burkholderia* and other symbionts.

516

517    It is clear from previous work that the *Lagria* symbionts related to *B. gladioli* evolved

518    from plant associated strains (Flórez and Kaltenpoth, 2017), likely transmitted to the

519    insects from the plant environment. The probable advantage of this early association for

520    the hosts was protection of eggs from infection, through small molecules made by its

521    microbiome. The strains characterized here, as well as the previously isolated Lv-StA,

522    were found to contain ample biosynthetic potential, and both Lv-StA and Lv-StB

523    produce antifungal compounds that protect eggs from fungal infection in lab

23

524     experiments (Flórez and Kaltenpoth, 2017). Yet Lv-StA is only found sporadically in the

525     field as a minor component of the microbiome (Flórez et al., 2018). It is probably

526     advantageous for *Lagria* beetles to maintain a pool of facultative symbionts with

527     different biosynthetic capability, to allow for fast adaptation to different environmental

528     infection pressures (Flórez et al., 2015). However, there may be less selection pressure

529     on a facultative symbiont to stay associated with its host if it can also survive in the

530     environment and infect plants.

531

532     The foundational event in the establishment of the symbiosis between Lv-StB and *L.*

533     *villosa* was likely the acquisition of the *lga* pathway, which putatively produces

534     lagriamide, in a non-reduced ancestor of Lv-StB (Flórez et al., 2018). We place this as

535     the first event for four reasons. First, the *lga* BGC is almost the oldest detectable

536     horizontal transfer that survives in the reduced genome of Lv-StB. Second, we found

537     little evidence that Lv-StB is capable of making metabolites of use to the host, indicating

538     that the symbiosis is likely not based on nutrition. *L. villosa*'s diet of plant leaves may be

539     nitrogen poor, with hard to digest plant cell wall components (Salem et al., 2017), but

540     we didn't find polysaccharide degrading pathways or extensive biosynthesis of essential

541     amino acids in the Lv-StB genome. Therefore, the *lga* BGC is the oldest remaining

542     feature that potentially increases host fitness. Third, the reduced coding density seen in

543     the Lv-StB genome may be indicative of a recent transitional event (Lo et al., 2016),

544     such as strict host association or a move to vertical transmission. Fourth, even though

545     we found genes missing from all DNA repair pathways, which is thought to be a driver

546     for increased AT content in symbiont genomes (McCutcheon and Moran, 2012), and

547    increase in AT may have an adaptive component that reduces the metabolic costs of

548    symbionts (Dietel et al., 2019), the GC content of the Lv-StB genome is not very

549    different from free-living *B. gladioli* strains, when compared to other "transitional"

550    symbionts. For instance, "*Candidatus* Pantoea carbekii" has a reduced genome and like

551    Lv-StB lacks full-length DNA polymerase I (Kenyon et al., 2015). This strain has a GC

552    content almost 30% lower than its closest freeliving relatives (Lo et al., 2016), while the

553    genome of Lv-StB has a GC content ~10% lower than its closest relatives. Therefore,

554    we propose that loss of DNA repair pathways and other genome degradation events in

555    Lv-StB occurred very recently, after the acquisition of *lga*.

556

557    It can be envisioned that *lga* provided a sustained survival advantage in an environment

558    where lagriamide consistently reduced egg fungal infections, and there was positive

559    selection on beetles that vertically transmitted *lga*-bearing symbionts. In *L. villosa*

560    symbionts are stored extracellularly, and they are spread onto the outside of eggs as

561    they are laid. According to observations in the congeneric species *L. hirta*, the

562    symbionts first enter through the egg micropyle to reach the embryonic organs where

563    they are housed throughout larval development (Jürgen Stammer, 1929). It is thus likely

564    that only a few of the cells are vertically transmitted by colonizing these structures,

565    potentially providing the population bottlenecks that could have caused initial

566    accumulation of deleterious mutations that started the process of genome reduction.

567    Meanwhile, loss of certain proteins limiting growth rate (see above) may have been

568    selected through increased Lv-StB populations and compound production. It is unknown

569    to what extent Lv-StB is genetically isolated in the larval or adult host, but we found

570    evidence of ongoing horizontal transfer events in the recent past, presumably through

571    contact with a complex microbiome associated with *L. villosa* egg surfaces. These

572    horizontal gene transfers likely happened concurrently with the ongoing genome

573    reduction process and may have been compensatory for gene losses (see above).

574    There is some precedent for extracellular symbionts with profoundly reduced genomes

575    (Kaiwa et al., 2014; Kikuchi et al., 2009; Nikoh et al., 2011; Salem et al., 2017). For

576    instance, the leaf beetle *Cassida rubiginosa* harbors a symbiont with the smallest

577    genome of any extracellular organism (0.27 Mbp), "*Candidatus* Stammera capleta",

578    which provides pectinolytic enzymes to help break down the host's leafy diet (Salem et

579    al., 2017). In many of these cases, symbionts are stored as isolated monocultures

580    within specialized structures in adult hosts, while vertical transmission is assisted by

581    packaging symbiont cells in protective "caplets" attached to eggs (Salem et al., 2017), a

582    "symbiont capsule" encased in chitin (Nikoh et al., 2011), or secreted in a galactan-

583    based jelly ingested by hatched larvae (Kaiwa et al., 2014), although reduced-genome

584    symbionts have also been known to be vertically transmitted by simple egg surface

585    contamination (Kikuchi et al., 2009), as in *L. villosa*. Because these examples are

586    advanced cases of genome reduction, it would be difficult to determine whether

587    horizontal transfer events occurred before or during genome reduction, and none have

588    been noted. The complexity of the *L. villosa* microbiome appears to be different, as it

589    afforded ample opportunity for horizontal gene transfer even while the genome of Lv-

590    StB was actively undergoing reduction.

591

26

592   Horizontal acquisition of genes has been observed in two types of reduced genome

593   symbiont, eukaryotic parasites in the genus *Encephalitozoon* (Pombert et al., 2012),

594   and Acetobacteraceae strains associated with the gut community of red carpenter ants

595   (Brown and Wernegreen, 2019). In both these cases symbiont genomes were similar in

596   size to Lv-StB (~2 Mbp), but with far greater coding density and fewer pseudogenes.

597   Furthermore, both *Encephalitozoon* and Acetobacteraceae strains were culturable,

598   suggesting that they are facultative symbionts in a less advanced state of genome

599   reduction compared to Lv-StB. The genome of Lv-StB appears to be different from

600   these examples, because there is evidence of recent horizontal transfers, even as

601   genes required for homologous recombination are currently missing. Either the loss of

602   homologous recombination was very recent, or such transfers could have occurred in a

603   RecA-independent manner. For example, plasmids could have been transferred into Lv-

604   StB cells, followed by the RecA-independent transposition of genes to the chromosome

605   (Harmer and Hall, 2016; Zupancic et al., 1983).

606

607   It is unclear whether Lv-StB will continue on the path of genome reduction to become

608   drastically reduced with a <1 Mbp genome. Where symbionts are required for host

609   survival and are genetically isolated within host cells or specialized structures, such a

610   process appears to be irreversible and unstoppable (Bennett and Moran, 2015; Moran,

611   1996). However, a number of alternate fates could be envisioned for Lv-StB. With a

612   complex microbiome, if ongoing gene losses in Lv-StB reduce its fitness past a certain

613   point, then it could be replaced by another strain, potentially accompanied by horizontal

614   transfer of the *Iga* pathway to a less reduced genomic chassis. Alternatively, horizontal

27

615    transfers of genes to Lv-StB could lead to an equilibrium of gene loss and gain.

616    Interestingly, it appears that up until the present time horizontal transfer has not

617    occurred fast enough to prevent widespread loss of metabolism and DNA repair in the

618    Lv-StB genome. The host could also evolve strategies to maintain an increasingly

619    genome-reduced Lv-StB, perhaps by selective extracellular partitioning and packaging

620    for vertical transfer similar to the examples outlined above. However, it is unclear

621    whether such an evolutionary path would be favorable, given that the co-infection of

622    multiple BGC-bearing symbiont strains could be advantageous in environments with

623    variable pathogen pressures.

624

625    In summary, evidence gathered here suggests that the introduction of the lagriamide

626    BGC initiated genome erosion of Lv-StB, potentially through selection of beetles that

627    transferred the symbiont vertically, leading to a population structure with frequent

628    bottlenecks. Simultaneous advantageous gene acquisitions may have enabled the

629    preferential survival of Lv-StB and its dominance in the adult host and the egg surface.

630

631

632    **MATERIALS AND METHODS**

633    **Sequencing and assembly of *Burkholderia gladioli* Lv-StA genome.** Genome

634    sequencing of the isolated *B. gladioli* Lv-StA strain was carried out using PacBio with

635    Single Molecule, Real-Time (SMRT) technology. For *de novo* assembly (carried out by

636    Eurofins Genomics), the HGAP pipeline was used (Heirarchical Genome Assembly

637    Process). Briefly, a preassembly of long and accurate sequences was generated by

638    mapping filtered subreads to so-called seed reads. Subsequently, the Celera assembler

639    was used to generate a draft assembly using multi-kb long reads, which in this case

640    rendered full genome closure. Finally, the Quiver algorithm was used to correct inDel

641    and substitution errors by considering the quality values from the bas.h5 files.

642

643    **Metagenomic binning and annotation.** Metagenomic assembly files were clustered

644    into putative genomic bins using Autometa (Master branch - commit bbcea30) (Miller et

645    al., 2019). Contigs with lengths smaller than 3,000bp were excluded from the binning

646    process and a taxonomy table was produced. Contigs classified as bacterial were

647    further binned into putative genomic bins using run_autometa.py. Unclustered contigs

648    were recruited into clusters using ML_recruitment.py. Results were summarized using

649    cluster_process.py. Resultant genome bins were compared to earlier versions (Flórez et

650    al., 2018) using Mash version 2.1.1 (Ondov et al., 2016) which hashes genomes to

651    patterns of $k$-mers (sketching) allowing for rapid distance calculations between two

652    sketches. All bins were sketched and distances were computed in a pairwise fashion.

653    Pairwise distances were visualized in R as a dendrogram and enabled the

654    determination of equivalent old and updated putative genome bins between analyses.

655    The updated putative genomic bins were annotated using Prokka version 1.13

656    (Seemann, 2014), with genbank compliance enabled. Reference genomes downloaded

657    from NCBI were similarly annotated with Prokka in order to maintain consistency

658    between datasets. Amino acid sequences of open-reading frames (ORFs) were further

659    annotated using DIAMOND blastp version 0.9.21.122 (Buchfink et al., 2015) against the

660    diamond formatted NR database. The search was limited to returning a maximum of 1

29

661    target sequence and the maximum number of high-scoring pairs per subject sequence

662    was set to 1. Results were summarized in BLAST tabular format with qseqid (Query

663    gene ID), stitle (aligned gene ID), pident (Percentage of identical matches), evalue

664    (Expected value), qlen (Query sequence length) and slen (aligned gene sequence

665    length) as desired parameter output. Pseudogenes were identified by finding Lv-StB

666    genes that were more than 20% shorter than their respective BLAST matches. This

667    criteria has been used previously to identify pseudogenes (Kwan and Schmidt, 2013;

668    Lerat and Ochman, 2005).

669

670    Coding density was calculated as the sum of all protein coding sequences (coding

671    sequence) as a percentage of the sum of all contigs (total sequence). In cases where

672    protein coding genes were found to overlap, the length of the overlap region was

673    counted only once. This calculation was performed for all binned genomes, on both

674    initial datasets (genbank files generated during Prokka annotations) and edited datasets

675    where pseudogenes had been removed. For the identification and count of genes

676    encoding transposases and hypothetical proteins, protein-coding gene amino acid files

677    (*.faa) containing Prokka annotations were parsed for gene descriptions containing

678    "transposase" and "hypothetical" strings.

679

680    **Taxonomic classification of genome bins.** Putative genome bins clustered from the

681    *L. villosa* metagenomic dataset were taxonomically classified using GTDB-Tk v0.2.2

682    (reference database gtdbtk.r86_v2) with default parameters (Parks et al., 2018). GTDB-

683    Tk identifies and aligns 120 bacterial marker genes per genome before calculating the

30

684    optimal placement of the respective alignments in the pre-computed GTDB-Tk

685    reference tree which consists of 94,759 genomes (Dataset S1B). A species was

686    assigned to a genome if it shared 95% or more ANI with a reference genome.

687

688    **Identification of "core" genes.** The set of "core" genes generally found in even the

689    most reduced symbiont genomes was taken from Table 2 of McCutcheon and Moran

690    2012 (McCutcheon and Moran, 2012). GFF files produced for *B. gladioli* Lv-StA and the

691    metagenomic bins by Prokka were searched for the following gene symbols: '*dnaE*',

692    '*dnaQ*', '*rpoA*', '*rpoB*', '*rpoC*', '*rpoD*', '*groL*', '*groS*', '*dnaK*', '*mnmA*', '*mnmE*', '*mnmG*',

693    '*sufS*', '*sufB*', '*sufC*', '*iscS*', '*iscA*', '*iscU*', '*rluA*', '*rluB*', '*rluC*', '*rluD*', '*rluE*', '*rluF*', '*infA*', '*infB*',

694    '*infC*', '*fusA*', '*tsf*', '*prfA*', '*prfB*', '*frr*', '*def*', '*alaS*', '*gltX*', '*glyQ*', '*ileS*', '*metG*', '*pheS*', '*trpS*',

695    '*valS*', '*rpsA*', '*rpsB*', '*rpsC*', '*rpsD*', '*rpsE*', '*rpsG*', '*rpsH*', '*rpsI*', '*rpsJ*', '*rpsK*', '*rpsL*', '*rpsM*',

696    '*rpsN*', '*rpsP*', '*rpsQ*', '*rpsR*', '*rpsS*', '*rplB*', '*rplC*', '*rplD*', '*rplE*', '*rplF*', '*rplK*', '*rplM*', '*rplN*',

697    '*rplO*', '*rplP*', '*rplT*', '*rplV*', '*rpmA*', '*rpmB*', '*rpmG*', '*rpmJ*', '*tRNA-Met*', '*tRNA-Gly*', '*tRNA-

698    Cys*', '*tRNA-Phe*', '*tRNA-Lys*', '*tRNA-Ala*', '*tRNA-Glu*', '*tRNA-Pro*', '*tRNA-Gln*', '*tRNA-Ile*'.

699    The presence of single or multiple examples of these genes per genome/bin was

700    tabulated in excel to produce Dataset S1D, and the percentage of the core gene set

701    found in a genome/bin was used for Table 1.

702

703    **Annotation and analysis of BGCs.** Putative biosynthetic gene clusters were identified

704    in all binned genomes using the AntiSMASH (Blin et al., 2019) docker image (Image ID:

705    8942d142d9ac). Entire genome HMMer analysis was enabled, and identified clusters

706    were compared to both antiSMASH-predicted clusters, the MIBiG database and

31

707　　secondary metabolite orthologous groups. Similarities between identified putative

708　　biosynthetic gene clusters were assessed using BiG-SCAPE version 20181005

709　　(Navarro-Muñoz et al., 2018) in "glocal" mode.

710

711　　**Construction of multilocus species tree.** Genomes of *Burkholderia* Lv-StB, *B. gladioli*

712　　Lv-StA and binned genomes taxonomically classified within the *Burkholderia* genus:

713　　DBSCAN_round6_14, DBSCAN_round6_18 and DBSCAN_round4_0, were uploaded to

714　　the AutoMLST website (Alanjary et al., 2019). A concatenated species tree was

715　　constructed in *de novo* mode, with default options as well as the IQ-TREE Ultrafast

716　　Bootstrap analysis and ModelFinder options enabled.

717

718　　**Calculation of average nucleotide identities.** The average nucleotide identities (ANIs)

719　　of *Burkholderia* Lv-StB and *B. gladioli* were calculated in a pairwise manner using

720　　FastANI (Jain et al., 2018) against 45 Burkholderia reference genomes downloaded

721　　from NCBI. A total of 13 genomes shared over 85% ANI (Dataset S1H). These

722　　genomes were all identified as *B. gladioli* species and were used in downstream

723　　analyses.

724

725　　**Quantification of divergence between *Burkholderia* Lv-StB and *B. gladioli* A1.**

726　　Orthologous protein sequences were identified in non-pseudogene sequence files of

727　　*Burkholderia* Lv-StB and 13 closely related genomes (identified through the ANI

728　　analyses: *B. gladioli* Lv-StA, *B. gladioli* A1, *B. gladioli* UCDUG, *B. gladioli*

729　　FDAARGOS_389, *B. gladioli* ATCC25417, *B. gladioli* Co14, *B. gladioli* SN82F6, *B.*

730    *gladioli* ATCC10248, *B. gladioli* NBRC13700, *B. gladioli* FDAARGOS_188, *B. gladioli*

731    MSMB1756, *B. gladioli* BSR3, *B. gladioli* KACC11889) using OMA version 2.2.0

732    (Altenhoff et al., 2018). A subset (797 groups) of the resultant orthologous groups (OGs)

733    was identified which included genes from all 14 genomes used in the analysis. Each set

734    of OG sequences were aligned using MUSCLE v3.8.31 (Edgar, 2004) and

735    corresponding nucleotide files were extracted and aligned against the amino acid

736    sequences using the PAL2NAL docker image (Image ID: ce3b1d7d83ab) (Suyama et

737    al., 2006) using codon table 11 and specifying no gaps with paml as the output format.

738    The resultant paml files were used to estimate pairwise dS (synonymous divergence

739    rate), dN (non-synonymous divergence rate) and kappa (transition/transversion ratio)

740    between individual genes per orthologous group with codeml (Yang, 2007) in the

741    PAL2NAL package. The following parameters were specified in the control file: runmode

742    = -2 (pairwise), model = 0 (one) fix_kappa = 0 (kappa to be estimated), fix_omega = 0

743    (estimate omega) where omega is the dN/dS ratio, with initial omega set to 0.2. Any

744    orthologous gene sets that included genes that gave a dS value over 3 were removed

745    from the analysis (Yang, 2014). Individual sequences from remaining OGs were then

746    gathered into genome-specific files (i.e all Lv-StB genes in all OGs were moved into an

747    ordered Lv-StB.faa/.ffn file). Stop codons were removed from each nucleotide

748    sequence. Sequences per genome were then concatenated to produce a single

749    sequence per genome. The concatenated amino acid sequences and corresponding

750    nucleotide sequences were aligned against one another using PAL2NAL as performed

751    for individual genes. Pairwise estimations of dS, dN and kappa were calculated as

752    before using codeml. Additionally, the concatenated genes were analysed a second

33

753    time using an alternative control file, in which the model was set to 2. The likelihood

754    ratio test value between pairwise null and alternative likelihood scores was calculated (

755    2x Alt_lnl - Null_lnl) for Lv-StB relative to the 13 reference genomes and found to be 0 in

756    all cases indicating that the omega (dN/dS) ratio was consistent between Lv-StB and

757    the reference genomes. Individual dS values were used to estimate divergence

758    between Lv-StB and the 13 reference genomes using divergence rates estimated by

759    Silva and Santos-Garcia (Silva and Santos-Garcia, 2015) (Table S1) in the equation:

760    Age of divergence (Mya) = dS ÷ divergence rate x 1,000,000. As Lv-StB shared the

761    greatest ANI with *B. gladioli* A1, the kappa value found between these two genomes

762    (7.03487) was used for amelioration estimates.

763

764    **Pangenome analysis.** To assess the pangenome of Lv-StB and other *B. gladioli*

765    genomes, GFF files generated by Prokka were analysed using Roary (Page et al.,

766    2015) which identifies core and accessory genes per genome. Concatenation and

767    alignment of orthologous genes was enabled in Roary and used to build a phylogenetic

768    tree with FastTree version 2.1.10 (Price et al., 2010). The resultant phylogenetic tree

769    and presence/absence matrix of genes in all genomes were visualized with the

770    roary_plots.py script. Additionally, non-pseudogenes of all genomes were annotated

771    against the KEGG database (Kanehisa et al., 2019; Kanehisa and Goto, 2000) using

772    kofamscan (Aramaki et al., 2019) with output in mapper format. An overview of the

773    completeness of general metabolic pathways was visualized using KEGG-Decoder

774    (Graham et al., 2018) with kofamscan annotations. For specific pathways of interest

775    (amino acids, DNA repair, nucleotide de novo biosynthesis), presence/absence

776    matrices of genes per KEGG pathway entry were visualized in R version 3.6.0 using the

777    tidyr, ggplot2 and viridis libraries.

778

779    **Identification of genes putatively acquired by horizontal transfer.** All amino acid

780    sequence files of non-pseudogenes of all genomes used in the ANI analysis (*B. gladioli*

781    Lv-StA and 45 *Burkholderia* reference genomes) were concatenated and converted into

782    a DIAMOND BLAST (Buchfink et al., 2015) database (build 125). The amino acid

783    sequence files of non-pseudogenes in Lv-StB were then searched against this

784    database. All genes that had no significant hit, or a significant hit but with a shared

785    percentage of less than 50% were considered "unique" to Lv-StB. Non-pseudogenes of

786    Lv-StB that were not found to have an ortholog in the OMA analysis were used to

787    validate this list. These "unique" genes were then compared to the NR database using

788    DIAMOND blastp (as described above) and any genes that had no significant hit were

789    removed from the "unique" set of genes. Manual inspection of the remaining genes

790    resulted in the removal of any genes that were closely related to *B. gladioli* genes (i.e.

791    found within *B. gladioli* genomes other than the ones investigated here). The remaining

792    unique genes were henceforth considered as gene potentially acquired via horizontal

793    transfer. This list was expanded with other genes within Lv-StB, for which homologs in

794    *Burkholderia* genomes could be found, but the closest match against the NR database

795    belonged to a different genus. For example, E5299_02249 of the "addic" group shared

796    53.1% sequence identity with a gene from *B. contaminans* strain LMG 23361 but shares

797    a higher sequence identity of 96.9% with *Ochrobactrum pituitosum*.

798

799  **De-amelioration of putatively horizontally transferred genes.** The method of

800  Lawrence and Ochman (Lawrence and Ochman, 1997) was implemented in Python and

801  is available at https://bitbucket.org/jason_c_kwan/age_horizontal_transfers.py. The

802  script takes as input 1. an in-frame nucleotide FASTA file containing the sequences of

803  putatively horizontally transferred genes, 2. an in-frame nucleotide FASTA file

804  containing a comparison set of gene sequences from the genome, 3. a synonymous

805  mutation rate in substitutions per 100 sites per million years, 4. a nonsynonymous

806  mutation rate in substitutions per 100 sites per million years, 5. a transition/transversion

807  ratio ($\kappa$), 6. a step time in millions of years, and 7. a maximum time to iterate to. The

808  script outputs GC content of each codon position (plus overall GC) at each timepoint,

809  and reports the estimated age of the gene cluster as corresponding to the iteration with

810  the smallest sum of squared deviations from equations 2–4. The substitution rates used

811  in our calculations were half of the divergence rates estimated by Silva and Santos-

812  Garcia (Silva and Santos-Garcia, 2015), BAU: synonymous 0.55, nonsynonymous 0.05;

813  BOB: synonymous 3.95; nonsynonymous 0.26; BPN: synonymous 3.2, nonsynonymous

814  0.28; BFL: synonymous 4.45, nonsynonymous 0.395. A value of 7.0348 was used for $\kappa$,

815  previously calculated for *Burkholderia* Lv-StB and *B. gladioli* A1 (see above). A step

816  time of 0.005 My and a maximum time of 50 My was used in all calculations. The

817  comparison gene set included only non-pseudogenes that were not identified as

818  putative horizontally transferred genes.

819

820  **Microbial community analysis.** 16S rRNA amplicon sequence datasets analysed via

821  oligotyping previously (Flórez et al., 2018) were re-analysed using Mothur v.1.40.3

822    (Schloss et al., 2009). Reads shorter than 200bp, or containing ambiguous bases or

823    homopolymeric runs longer than 7 bases were removed from the dataset. Chimeric

824    sequences were identified using VSEARCH (Rognes et al., 2016) and removed from

825    the dataset. Reads were taxonomically classified against the Silva database (version

826    132) and all reads classified as unknown, eukaryotic, mitochondrial or as chloroplasts

827    were removed from the dataset. Reads were aligned using the Silva database (v. 132)

828    as reference and clustered into operational taxonomic units (OTUs) at a distance of

829    0.03: an approximation to bacterial species. Counts of OTUs per sample were

830    generated and the top 10 most abundant OTUs were plotted (Figure S3). The top 50

831    most abundant OTUs were queried against the "nt" nucleotide database using blast for

832    taxonomic classification.

833

834    **Data availability.** The complete *Burkholderia gladioli* Lv-StA genome, the draft

835    assembly of the *Burkholderia* Lv-StB genome and other bacterial metagenomic bins will

836    be deposited in GenBank, and the respective accession numbers will be included in the

837    accepted version of this manuscript.

838

839    **ACKNOWLEDGMENTS**

845    Sciences. The CHTC is supported by UW-Madison, the Advanced Computing Initiative,

846    the Wisconsin Alumni Research Foundation, and Wisconsin Institutes for Discovery,

847    and the National Science Foundation and is an active member of the Open Science

848    Grid, which is supported by the National Science Foundation and the U.S. Department

849    of Energy's Office of Science.

850

851    **COMPETING INTERESTS**

852    The authors disclose no competing interests.

## REFERENCES

853

854    Aghajari N, Jensen KF, Gajhede M. 1994. Crystallization and preliminary X-ray

855    diffraction studies on the apo form of orotate phosphoribosyltransferase from

856    *Escherichia coli*. *J Mol Biol* **241**:292–294.

857    Akman L, Yamashita A, Watanabe H, Oshima K, Shiba T, Hattori M, Aksoy S. 2002.

858    Genome sequence of the endocellular obligate symbiont of tsetse flies,

859    *Wigglesworthia glossinidia*. *Nat Genet* **32**:402–407.

860    Alanjary M, Steinke K, Ziemert N. 2019. AutoMLST: An automated web server for

861    generating multi-locus species trees highlighting natural product potential. *Nucleic*

862    *Acids Res* **47**:W276–W282.

863    Almassy RJ, Janson CA, Kan CC, Hostomska Z. 1992. Structures of apo and

864    complexed *Escherichia coli* glycinamide ribonucleotide transformylase. *Proc Natl*

865    *Acad Sci U S A* **89**:6114–6118.

866    Altenhoff AM, Glover NM, Train C-M, Kaleb K, Warwick Vesztrocy A, Dylus D, de Farias

867    TM, Zile K, Stevenson C, Long J, Redestig H, Gonnet GH, Dessimoz C. 2018. The

868    OMA orthology database in 2018: Retrieving evolutionary relationships among all

869    domains of life through richer web and programmatic interfaces. *Nucleic Acids Res*

870    **46**:D477–D485.

871    Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, Ogata H. 2019.

872    KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive

873    score threshold. *bioRxiv*. doi:10.1101/602110

874    Barrett TE, Savva R, Panayotou G, Barlow T, Brown T, Jiricny J, Pearl LH. 1998.

875    Crystal structure of a G:T/U mismatch-specific DNA glycosylase: Mismatch

876    recognition by complementary-strand interactions. *Cell* **92**:117–129.

877    Bauer E, Kaltenpoth M, Salem H. 2019. Minimal fermentative metabolism fuels

878    extracellular symbiont in a leaf beetle In review.

879    Becq J, Churlaud C, Deschavanne P. 2010. A benchmark of parametric methods for

880    horizontal transfers detection. *PLoS One* **5**:e9989.

881    Bennett GM, Moran NA. 2015. Heritable symbiosis: The advantages and perils of an

882    evolutionary rabbit hole. *Proc Natl Acad Sci U S A* **112**:10169–10176.

883    Blin K, Shaw S, Steinke K, Villebro R, Ziemert N, Lee SY, Medema MH, Weber T. 2019.

884    antiSMASH 5.0: Updates to the secondary metabolite genome mining pipeline.

885    *Nucleic Acids Res* **47**:W81–W87.

886    Blin K, Wolf T, Chevrette MG, Lu X, Schwalen CJ, Kautsar SA, Suarez Duran HG, de

887    Los Santos ELC, Kim HU, Nave M, Dickschat JS, Mitchell DA, Shelest E, Breitling

888    R, Takano E, Lee SY, Weber T, Medema MH. 2017. antiSMASH 4.0-improvements

889    in chemistry prediction and gene cluster boundary identification. *Nucleic Acids Res*

890    **45**:W36–W41.

891    Boël G, Smith PC, Ning W, Englander MT, Chen B, Hashem Y, Testa AJ, Fischer JJ,

892    Wieden H-J, Frank J, Gonzalez RL Jr, Hunt JF. 2014. The ABC-F protein EttA

893    gates ribosome entry into the translation elongation cycle. *Nat Struct Mol Biol*

894    **21**:143–151.

895    Boiteux S, Gajewski E, Laval J, Dizdaroglu M. 1992. Substrate specificity of the

896    *Escherichia coli* Fpg protein formamidopyrimidine-DNA glycosylase: Excision of

897    purine lesions in DNA produced by ionizing radiation or photosensitization.

898    *Biochemistry*. doi:10.1021/bi00116a016

899    Brignole EJ, Ando N, Zimanyi CM, Drennan CL. 2012. The prototypic class Ia

900        ribonucleotide reductase from *Escherichia coli*: still surprising after all these years.

901        *Biochem Soc Trans* **40**:523–530.

902    Brown BP, Wernegreen JJ. 2019. Genomic erosion and extensive horizontal gene

903        transfer in gut-associated Acetobacteraceae. *BMC Genomics* **20**:472.

904    Buchfink B, Xie C, Huson DH. 2015. Fast and sensitive protein alignment using

905        DIAMOND. *Nat Methods* **12**:59–60.

906    Carreras CW, Santi DV. 1995. The catalytic mechanism and structure of thymidylate

907        synthase. *Annu Rev Biochem* **64**:721–762.

908    Carter AT, Pearson BM, Dickinson JR, Lancashire WE. 1993. Sequence of the

909        *Escherichia coli* K-12 *edd* and *eda* genes of the Entner-Doudoroff pathway. *Gene*

910        **130**:155–156.

911    Chen B, Boël G, Hashem Y, Ning W, Fei J, Wang C, Gonzalez RL Jr, Hunt JF, Frank J.

912        2014. EttA regulates translation by binding the ribosomal E site and restricting

913        ribosome-tRNA dynamics. *Nat Struct Mol Biol* **21**:152–159.

914    Chen EZ, Bushman FD, Li H. 2017. A Model-Based Approach For Species Abundance

915        Quantification Based On Shotgun Metagenomic Data. *Stat Biosci* **9**:13–27.

916    Chen X, Schreiber K, Appel J, Makowka A, Fähnrich B, Roettger M, Hajirezaei MR,

917        Sönnichsen FD, Schönheit P, Martin WF, Gutekunst K. 2016. The Entner-

918        Doudoroff pathway is an overlooked glycolytic route in cyanobacteria and plants.

919        *Proc Natl Acad Sci U S A* **113**:5441–5446.

920    Currie CR, Scott JA, Summerbell RC, Malloch D. 1999. Fungus-growing ants use

921        antibiotic-producing bacteria to control garden parasites. *Nature*.

41

922        doi:10.1038/19519

923    Delforno TP, Lacerda Júnior GV, Noronha MF, Sakamoto IK, Varesche MBA, Oliveira

924        VM. 2017. Microbial diversity of a full-scale UASB reactor applied to poultry

925        slaughterhouse wastewater treatment: integration of 16S rRNA gene amplicon and

926        shotgun metagenomic sequencing. *Microbiologyopen* **6**. doi:10.1002/mbo3.443

927    Dessimoz C, Cannarozzi G, Gil M, Margadant D, Roth A, Schneider A, Gonnet GH.

928        2005. OMA, a comprehensive, automated project for the identification of orthologs

929        from complete genome data: Introduction and first achievementsComparative

930        Genomics. Springer Berlin Heidelberg. pp. 61–72.

931    Dietel A-K, Kaltenpoth M, Kost C. 2018. Convergent evolution in intracellular elements:

932        Plasmids as model endosymbionts. *Trends Microbiol* **26**:755–768.

933    Dietel A-K, Merker H, Kaltenpoth M, Kost C. 2019. Selective advantages favour high

934        genomic AT-contents in intracellular elements. *PLoS Genet* **15**:e1007778.

935    Di Pierro M, Lu R, Uzzau S, Wang W, Margaretten K, Pazzani C, Maimone F, Fasano

936        A. 2001. Zonula occludens toxin structure-function analysis. Identification of the

937        fragment biologically active on tight junctions and of the zonulin receptor binding

938        domain. *J Biol Chem* **276**:19160–19165.

939    Dolan SK, Welch M. 2018. The glyoxylate shunt, 60 years on. *Annu Rev Microbiol*

940        **72**:309–330.

941    Dose B, Niehs SP, Scherlach K, Flórez LV, Kaltenpoth M, Hertweck C. 2018.

942        Unexpected bacterial origin of the antibiotic icosalide: Two-tailed depsipeptide

943        assembly in multifarious *Burkholderia* symbionts. *ACS Chem Biol* **13**:2414–2420.

944    Duarte M, Videira A. 2009. Effects of mitochondrial complex III disruption in the

945    respiratory chain of *Neurospora crassa*. *Mol Microbiol* **72**:246–258.

946    Edgar RC. 2004. MUSCLE: Multiple sequence alignment with high accuracy and high

947         throughput. *Nucleic Acids Res* **32**:1792–1797.

948    Flórez LV, Biedermann PHW, Engl T, Kaltenpoth M. 2015. Defensive symbioses of

949         animals with prokaryotic and eukaryotic microorganisms. *Nat Prod Rep* **32**:904–

950         936.

951    Flórez LV, Kaltenpoth M. 2017. Symbiont dynamics and strain diversity in the defensive

952         mutualism between *Lagria* beetles and *Burkholderia*. *Environ Microbiol* **19**:3674–

953         3688.

954    Flórez LV, Scherlach K, Gaube P, Ross C, Sitte E, Hermes C, Rodrigues A, Hertweck

955         C, Kaltenpoth M. 2017. Antibiotic-producing symbionts dynamically transition

956         between plant pathogenicity and insect-defensive mutualism. *Nat Commun*

957         **8**:15172.

958    Flórez LV, Scherlach K, Miller IJ, Rodrigues A, Kwan JC, Hertweck C, Kaltenpoth M.

959         2018. An antifungal polyketide associated with horizontally acquired genes

960         supports symbiont-mediated defense in *Lagria villosa* beetles. *Nat Commun*

961         **9**:2478.

962    Gerdes SY, Scholle MD, Campbell JW, Balázsi G, Ravasz E, Daugherty MD, Somera

963         AL, Kyrpides NC, Anderson I, Gelfand MS, Bhattacharya A, Kapatral V, D'Souza M,

964         Baev MV, Grechkin Y, Mseeh F, Fonstein MY, Overbeek R, Barabási A-L, Oltvai

965         ZN, Osterman AL. 2003. Experimental determination and system level analysis of

966         essential genes in *Escherichia coli* MG1655. *J Bacteriol* **185**:5673–5684.

967    Gillings MR. 2017. Lateral gene transfer, bacterial genome evolution, and the

968    Anthropocene. *Ann N Y Acad Sci* **1389**:20–36.

969    Goodall ECA, Robinson A, Johnston IG, Jabbari S, Turner KA, Cunningham AF, Lund

970    PA, Cole JA, Henderson IR. 2018. The essential genome of *Escherichia coli* K-12.

971    *MBio* **9**:e02096–17.

972    Goris J, Konstantinidis KT, Klappenbach JA, Coenye T, Vandamme P, Tiedje JM. 2007.

973    DNA–DNA hybridization values and their relationship to whole-genome sequence

974    similarities. *Int J Syst Evol Microbiol* **57**:81–91.

975    Graham ED, Heidelberg JF, Tully BJ. 2018. Potential for primary productivity in a

976    globally-distributed bacterial phototroph. *ISME J* **12**:1861–1866.

977    Hang B, Downing G, Guliaev AB, Singer B. 2002. Novel activity of *Escherichia coli*

978    mismatch uracil-DNA glycosylase (Mug) excising 8-(hydroxymethyl)-3,$N^4$-

979    ethenocytosine, a potential product resulting from glycidaldehyde reaction.

980    *Biochemistry* **41**:2158–2165.

981    Harmer CJ, Hall RM. 2016. IS*26*-mediated formation of transposons carrying antibiotic

982    resistance genes. *mSphere* **1**:e00038–16.

983    Higgins S, Sanchez-Contreras M, Gualdi S, Pinto-Carbó M, Carlier A, Eberl L. 2017.

984    The essential genome of *Burkholderia cenocepacia* H111. *J Bacteriol* **199**:e00260–

985    17.

986    Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. 2018. High

987    throughput ANI analysis of 90K prokaryotic genomes reveals clear species

988    boundaries. *Nat Commun* **9**:5114.

989    Jain M, Munoz-Bodnar A, Gabriel DW. 2017. Concomitant loss of the glyoxalase system

990    and glycolysis makes the uncultured pathogen "*Candidatus* Liberibacter asiaticus"

991    an energy scavenger. *Appl Environ Microbiol* **83**:e01670–17.

992    Jensen PR. 2016. Natural products and the gene cluster revolution. *Trends Microbiol*

993    **24**:968–977.

994    Jin J, Wu R, Zhu J, Yang S, Lei Z, Wang N, Singh VK, Zheng J, Jia Z. 2015. Insights

995    into the cellular function of YhdE, a nucleotide pyrophosphatase from *Escherichia*

996    *coli*. *PLoS One* **10**:e0117823.

997    Jones AG, Mason CJ, Felton GW, Hoover K. 2019. Host plant and population source

998    drive diversity of microbial gut communities in two polyphagous insects. *Sci Rep*

999    **9**:2792.

1000   Joseph B, Schwarz RF, Linke B, Blom J, Becker A, Claus H, Goesmann A, Frosch M,

1001   Müller T, Vogel U, Schoen C. 2011. Virulence evolution of the human pathogen

1002   *Neisseria meningitidis* by recombination in the core and accessory genome. *PLoS*

1003   *One* **6**:e18441.

1004   Jürgen Stammer H. 1929. Die Symbiose der Lagriiden (Coleoptera). *Zoomorphology*

1005   **15**:1–34.

1006   Kaiwa N, Hosokawa T, Nikoh N, Tanahashi M, Moriyama M, Meng X-Y, Maeda T,

1007   Yamaguchi K, Shigenobu S, Ito M, Fukatsu T. 2014. Symbiont-supplemented

1008   maternal investment underpinning host's ecological adaptation. *Curr Biol* **24**:2465–

1009   2470.

1010   Kanehisa M, Goto S. 2000. KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic*

1011   *Acids Res* **28**:27–30.

1012   Kanehisa M, Sato Y, Furumichi M, Morishima K, Tanabe M. 2019. New approach for

1013   understanding genome variations in KEGG. *Nucleic Acids Res* **47**:D590–D595.

1014    Kenyon LJ, Meulia T, Sabree ZL. 2015. Habitat visualization and genomic analysis of

1015    "*Candidatus* Pantoea carbekii," the primary symbiont of the brown marmorated

1016    stink bug. *Genome Biol Evol* **7**:620–635.

1017    Kergoat GJ, Bouchard P, Clamens A-L, Abbate JL, Jourdan H, Jabbour-Zahab R,

1018    Genson G, Soldati L, Condamine FL. 2014. Cretaceous environmental changes led

1019    to high extinction rates in a hyperdiverse beetle family. *BMC Evol Biol* **14**:220.

1020    Kikuchi Y, Hosokawa T, Nikoh N, Meng X-Y, Kamagata Y, Fukatsu T. 2009. Host-

1021    symbiont co-speciation and reductive genome evolution in gut symbiotic bacteria of

1022    acanthosomatid stinkbugs. *BMC Biol* **7**:2.

1023    Kowalczykowski SC, Dixon DA, Eggleston AK, Lauder SD, Rehrauer WM. 1994.

1024    Biochemistry of homologous recombination in *Escherichia coli*. *Microbiol Rev*

1025    **58**:401–465.

1026    Kroiss J, Kaltenpoth M, Schneider B, Schwinger M-G, Hertweck C, Maddula RK,

1027    Strohm E, Svatos A. 2010. Symbiotic *Streptomycetes* provide antibiotic

1028    combination prophylaxis for wasp offspring. *Nat Chem Biol* **6**:261–263.

1029    Kumar S, Stecher G, Suleski M, Hedges SB. 2017. TimeTree: A resource for timelines,

1030    timetrees, and divergence times. *Mol Biol Evol* **34**:1812–1819.

1031    Kwan JC, Donia MS, Han AW, Hirose E, Haygood MG, Schmidt EW. 2012. Genome

1032    streamlining and chemical defense in a coral reef symbiosis. *Proc Natl Acad Sci U*

1033    *S A* **109**:20655–20660.

1034    Kwan JC, Schmidt EW. 2013. Bacterial endosymbiosis in a chordate host: Long-term

1035    co-evolution and conservation of secondary metabolism. *PLoS One* **8**:e80822.

1036    Kwong WK, Engel P, Koch H, Moran NA. 2014. Genomics and host specialization of

1037    honey bee and bumble bee gut symbionts. *Proc Natl Acad Sci U S A* **111**:11509–

1038        11514.

1039    Latorre A, Manzano-Marín A. 2017. Dissecting genome reduction and trait loss in insect

1040        endosymbionts. *Ann N Y Acad Sci* **1389**:52–75.

1041    Lawrence JG, Ochman H. 1997. Amelioration of bacterial genomes: Rates of change

1042        and exchange. *J Mol Evol* **44**:383–397.

1043    Lerat E, Ochman H. 2005. Recognizing the pseudogenes in bacterial genomes. *Nucleic*

1044        *Acids Res* **33**:3125–3132.

1045    Lin J-J, Sancar A. 1992. Active site of (A) BC excinuclease. I. Evidence for 5'incision by

1046        UvrC through a catalytic site involving Asp399, Asp438, Asp466, and His538

1047        residues. *J Biol Chem* **267**:17688–17692.

1048    Linke D, Riess T, Autenrieth IB, Lupas A, Kempf VAJ. 2006. Trimeric autotransporter

1049        adhesins: Variable structure, common function. *Trends Microbiol* **14**:264–270.

1050    Liu F, Lee H, Lan R, Zhang L. 2016. Zonula occludens toxins and their prophages in

1051        Campylobacter species. *Gut Pathog* **8**:43.

1052    Liu P, Burdzy A, Sowers LC. 2003. Repair of the mutagenic DNA oxidation product, 5-

1053        formyluracil. *DNA Repair* **2**:199–210.

1054    Lopera J, Miller IJ, McPhail KL, Kwan JC. 2017. Increased Biosynthetic Gene Dosage in

1055        a Genome-Reduced Defensive Bacterial Symbiont. *mSystems* **2**.

1056        doi:10.1128/mSystems.00096-17

1057    López-García P, Eme L, Moreira D. 2017. Symbiosis in eukaryotic evolution. *J Theor*

1058        *Biol* **434**:20–33.

1059    Lo W-S, Huang Y-Y, Kuo C-H. 2016. Winding paths to simplicity: Genome evolution in

1060      facultative insect symbionts. *FEMS Microbiol Rev* **40**:855–874.

1061  Lunin VV, Li Y, Schrag JD, Iannuzzi P, Cygler M, Matte A. 2004. Crystal structures of

1062      *Escherichia coli* ATP-dependent glucokinase and its complex with glucose. *J*

1063      *Bacteriol* **186**:6915–6927.

1064  Mahendran V, Liu F, Riordan SM, Grimm MC, Tanaka MM, Zhang L. 2016. Examination

1065      of the effects of *Campylobacter concisus* zonula occludens toxin on intestinal

1066      epithelial cells and macrophages. *Gut Pathog* **8**:18.

1067  Mathew GM, Ju Y-M, Lai C-Y, Mathew DC, Huang CC. 2012. Microbial community

1068      analysis in the termite gut and fungus comb of *Odontotermes formosanus*: the

1069      implication of *Bacillus* as mutualists. *FEMS Microbiol Ecol* **79**:504–517.

1070  McClure EA, Nelson MC, Lin A, Graf J. 2019. *Macrobdella decora*: Old world leech gut

1071      microbial community structure conserved in a new world leech. *bioRxiv*.

1072  McCullough AK, Dodson ML, Lloyd RS. 1999. Initiation of base excision repair:

1073      Glycosylase mechanisms and structures. *Annu Rev Biochem* **68**:255–285.

1074  McCutcheon JP, Moran NA. 2012. Extreme genome reduction in symbiotic bacteria. *Nat*

1075      *Rev Microbiol* **10**:13–26.

1076  Miller IJ, Chevrette MG, Kwan JC. 2017. Interpreting Microbial Biosynthesis in the

1077      Genomic Age: Biological and Practical Considerations. *Mar Drugs* **15**:165.

1078  Miller IJ, Rees ER, Ross J, Miller I, Baxa J, Lopera J, Kerby RL, Rey FE, Kwan JC.

1079      2019. Autometa: Automated extraction of microbial genomes from individual

1080      shotgun metagenomes. *Nucleic Acids Res* **47**:e57.

1081  Miller IJ, Vanee N, Fong SS, Lim-Fong GE, Kwan JC. 2016a. Lack of overt genome

1082      reduction in the bryostatin-producing bryozoan symbiont, "*Candidatus* Endobugula

48

1083    sertula." *Appl Environ Microbiol* **82**:6573–6583.

1084    Miller IJ, Weyna TR, Fong SS, Lim-Fong GE, Kwan JC. 2016b. Single sample

1085         resolution of rare microbial dark matter in a marine invertebrate metagenome. *Sci*

1086         *Rep* **6**:34362.

1087    Mira A, Ochman H, Moran NA. 2001. Deletional bias and the evolution of bacterial

1088         genomes. *Trends Genet* **17**:589–596.

1089    Mohni KN, Wessel SR, Zhao R, Wojciechowski AC, Luzwick JW, Layden H, Eichman

1090         BF, Thompson PS, Mehta KPM, Cortez D. 2019. HMCES maintains genome

1091         integrity by shielding abasic sites in single-strand DNA. *Cell* **176**:144–153.e13.

1092    Montagna M, Chouaia B, Mazza G, Prosdocimi EM, Crotti E, Mereghetti V, Vacchini V,

1093         Giorgi A, De Biase A, Longo S, Cervo R, Lozzia GC, Alma A, Bandi C, Daffonchio

1094         D. 2015. Effects of the diet on the microbiota of the red palm weevil (Coleoptera:

1095         Dryophthoridae). *PLoS One* **10**:e0117439.

1096    Moran NA. 1996. Accelerated evolution and Muller's rachet in endosymbiotic bacteria.

1097         *Proc Natl Acad Sci U S A* **93**:2873–2878.

1098    Moran NA, Munson MA, Baumann P, Ishikawa H. 1993. A molecular clock in

1099         endosymbiotic bacteria is calibrated using the insect hosts. *Proc R Soc Lond B Biol*

1100         *Sci* **253**:167–171.

1101    Moule MG, Hemsley CM, Seet Q, Guerra-Assunção JA, Lim J, Sarkar-Tyson M, Clark

1102         TG, Tan PBO, Titball RW, Cuccui J, Wren BW. 2014. Genome-wide saturation

1103         mutagenesis of *Burkholderia pseudomallei* K96243 predicts essential genes and

1104         novel targets for antimicrobial development. *MBio* **5**:e00926–13.

1105    Munson MA, Baumann P, Clark MA, Baumann L, Moran NA, Voegtlin DJ, Campbell BC.

1106      1991. Evidence for the establishment of aphid-eubacterium endosymbiosis in an

1107      ancestor of four aphid families. *J Bacteriol* **173**:6321–6324.

1108    Muto A, Osawa S. 1987. The guanine and cytosine content of genomic DNA and

1109      bacterial evolution. *Proc Natl Acad Sci U S A* **84**:166–169.

1110    Navarro-Muñoz JC, Selem-Mojica N, Mullowney MW, Kautsar S, Tryon JH, Parkinson

1111      EI, De Los Santos ELC, Yeong M, Cruz-Morales P, Abubucker S, Roeters A,

1112      Lokhorst W, Fernandez-Guerra A, Cappelini LTD, Thomson RJ, Metcalf WW,

1113      Kelleher NL, Barona-Gomez F, Medema MH. 2018. A computational framework for

1114      systematic exploration of biosynthetic diversity from large-scale genomic data.

1115      *bioRxiv*. doi:10.1101/445270

1116    Nelson DL, Lehninger AL, Cox MM. 2008. Lehninger Principles of Biochemistry,

1117      Lehninger Principles of Biochemistry. W. H. Freeman.

1118    Nikoh N, Hosokawa T, Oshima K, Hattori M, Fukatsu T. 2011. Reductive evolution of

1119      bacterial genome in insect gut environment. *Genome Biol Evol* **3**:702–714.

1120    Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, Phillippy AM.

1121      2016. Mash: Fast genome and metagenome distance estimation using MinHash.

1122      *Genome Biol* **17**:132.

1123    Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, Fookes M, Falush

1124      D, Keane JA, Parkhill J. 2015. Roary: Rapid large-scale prokaryote pan genome

1125      analysis. *Bioinformatics* **31**:3691–3693.

1126    Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil P-A,

1127      Hugenholtz P. 2018. A standardized bacterial taxonomy based on genome

1128      phylogeny substantially revises the tree of life. *Nat Biotechnol* **36**:996–1004.

1129    Piel J. 2002. A polyketide synthase-peptide synthetase gene cluster from an uncultured

1130        bacterial symbiont of *Paederus* beetles. *Proc Natl Acad Sci U S A* **99**:14002–

1131        14007.

1132    Pitcher RS, Brissett NC, Doherty AJ. 2007. Nonhomologous end-joining in bacteria: A

1133        microbial perspective. *Annu Rev Microbiol* **61**:259–282.

1134    Pombert J-F, Selman M, Burki F, Bardell FT, Farinelli L, Solter LF, Whitman DW, Weiss

1135        LM, Corradi N, Keeling PJ. 2012. Gain and loss of multiple functionally related,

1136        horizontally transferred genes in the reduced genomes of two microsporidian

1137        parasites. *Proc Natl Acad Sci U S A* **109**:12638–12643.

1138    Porter TN, Li Y, Raushel FM. 2004. Mechanism of the dihydroorotase reaction.

1139        *Biochemistry* **43**:16285–16292.

1140    Powell JE, Leonard SP, Kwong WK, Engel P, Moran NA. 2016. Genome-wide screen

1141        identifies host colonization determinants in a bacterial gut symbiont. *Proc Natl Acad*

1142        *Sci U S A* **113**:13887–13892.

1143    Price MN, Dehal PS, Arkin AP. 2010. FastTree 2--approximately maximum-likelihood

1144        trees for large alignments. *PLoS One* **5**:e9490.

1145    Raghunathan D, Wells TJ, Morris FC, Shaw RK, Bobat S, Peters SE, Paterson GK,

1146        Jensen KT, Leyton DL, Blair JMA, Browning DF, Pravin J, Flores-Langarica A,

1147        Hitchcock JR, Moraes CTP, Piazza RMF, Maskell DJ, Webber MA, May RC,

1148        MacLennan CA, Piddock LJ, Cunningham AF, Henderson IR. 2011. SadA, a

1149        trimeric autotransporter from *Salmonella enterica* serovar Typhimurium, can

1150        promote biofilm formation and provides limited protection against infection. *Infect*

1151        *Immun* **79**:4342–4352.

1152 Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, Darling A,

1153     Malfatti S, Swan BK, Gies EA, Dodsworth JA, Hedlund BP, Tsiamis G, Sievert SM,

1154     Liu W-T, Eisen JA, Hallam SJ, Kyrpides NC, Stepanauskas R, Rubin EM,

1155     Hugenholtz P, Woyke T. 2013. Insights into the phylogeny and coding potential of

1156     microbial dark matter. *Nature* **499**:431–437.

1157 Rio RVM, Maltz M, McCormick B, Reiss A, Graf J. 2009. Symbiont succession during

1158     embryonic development of the European medicinal leech, *Hirudo verbana*. *Appl*

1159     *Environ Microbiol* **75**:6890–6895.

1160 Rognes T, Flouri T, Nichols B, Quince C, Mahé F. 2016. VSEARCH: A versatile open

1161     source tool for metagenomics. *PeerJ* **4**:e2584.

1162 Salem H, Bauer E, Kirsch R, Berasategui A, Cripps M, Weiss B, Koga R, Fukumori K,

1163     Vogel H, Fukatsu T, Kaltenpoth M. 2017. Drastic genome reduction in an

1164     herbivore's pectinolytic symbiont. *Cell* **171**:1520–1531.e13.

1165 Saparbaev M, Langouët S, Privezentzev CV, Guengerich FP, Cai H, Elder RH, Laval J.

1166     2002. 1,$N^2$-Ethenoguanine, a mutagenic DNA adduct, is a primary substrate of

1167     *Escherichia coli* mismatch-specific uracil-DNA glycosylase and human alkylpurine-

1168     DNA-*N*-glycosylase. *J Biol Chem* **277**:26987–26993.

1169 Saparbaev M, Laval J. 1998. 3,$N^4$-ethenocytosine, a highly mutagenic adduct, is a

1170     primary substrate for *Escherichia coli* double-stranded uracil-DNA glycosylase and

1171     human mismatch-specific thymine-DNA glycosylase. *Proc Natl Acad Sci U S A*

1172     **95**:8508–8513.

1173 Schendel FJ, Mueller E, Stubbe J, Shiau A, Smith JM. 1989. Formylglycinamide

1174     ribonucleotide synthetase from *Escherichia coli*: Cloning, sequencing,

1175    overproduction, isolation, and characterization. *Biochemistry* **28**:2459–2471.

1176   Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA,

1177    Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn

1178    DJ, Weber CF. 2009. Introducing Mothur: Open-source, platform-independent,

1179    community-supported software for describing and comparing microbial

1180    communities. *Appl Environ Microbiol* **75**:7537–7541.

1181   Schofield MJ, Hsieh P. 2003. DNA mismatch repair: Molecular mechanisms and

1182    biological function. *Annu Rev Microbiol* **57**:579–608.

1183   Seemann T. 2014. Prokka: Rapid prokaryotic genome annotation. *Bioinformatics*

1184    **30**:2068–2069.

1185   Shigenobu S, Watanabe H, Hattori M, Sakaki Y, Ishikawa H. 2000. Genome sequence

1186    of the endocellular bacterial symbiont of aphids *Buchnera* sp. APS. *Nature* **407**:81–

1187    86.

1188   Shih PM, Matzke NJ. 2013. Primary endosymbiosis events date to the later Proterozoic

1189    with cross-calibrated phylogenetic dating of duplicated ATPase proteins. *Proc Natl*

1190    *Acad Sci U S A* **110**:12355–12360.

1191   Silva FJ, Santos-Garcia D. 2015. Slow and fast evolving endosymbiont lineages:

1192    Positive correlation between the rates of synonymous and non-synonymous

1193    substitution. *Front Microbiol* **6**:1279.

1194   Suyama M, Torrents D, Bork P. 2006. PAL2NAL: Robust conversion of protein

1195    sequence alignments into the corresponding codon alignments. *Nucleic Acids Res*

1196    **34**:W609–12.

1197   Szklarczyk D, Morris JH, Cook H, Kuhn M, Wyder S, Simonovic M, Santos A, Doncheva

1198      NT, Roth A, Bork P, Jensen LJ, von Mering C. 2017. The STRING database in

1199      2017: Quality-controlled protein–protein association networks, made broadly

1200      accessible. *Nucleic Acids Res* **45**:D362–D368.

1201      Takeya M, Hirai MY, Osanai T. 2017. Allosteric inhibition of phosphoenolpyruvate

1202      carboxylases is determined by a single amino acid residue in cyanobacteria. *Sci*

1203      *Rep* **7**:41080.

1204      Tchigvintsev A, Tchigvintsev D, Flick R, Popovic A, Dong A, Xu X, Brown G, Lu W, Wu

1205      H, Cui H, Dombrowski L, Joo JC, Beloglazova N, Min J, Savchenko A, Caudy AA,

1206      Rabinowitz JD, Murzin AG, Yakunin AF. 2013. Biochemical and structural studies

1207      of conserved Maf proteins revealed nucleotide pyrophosphatases with a preference

1208      for modified nucleotides. *Chem Biol* **20**:1386–1398.
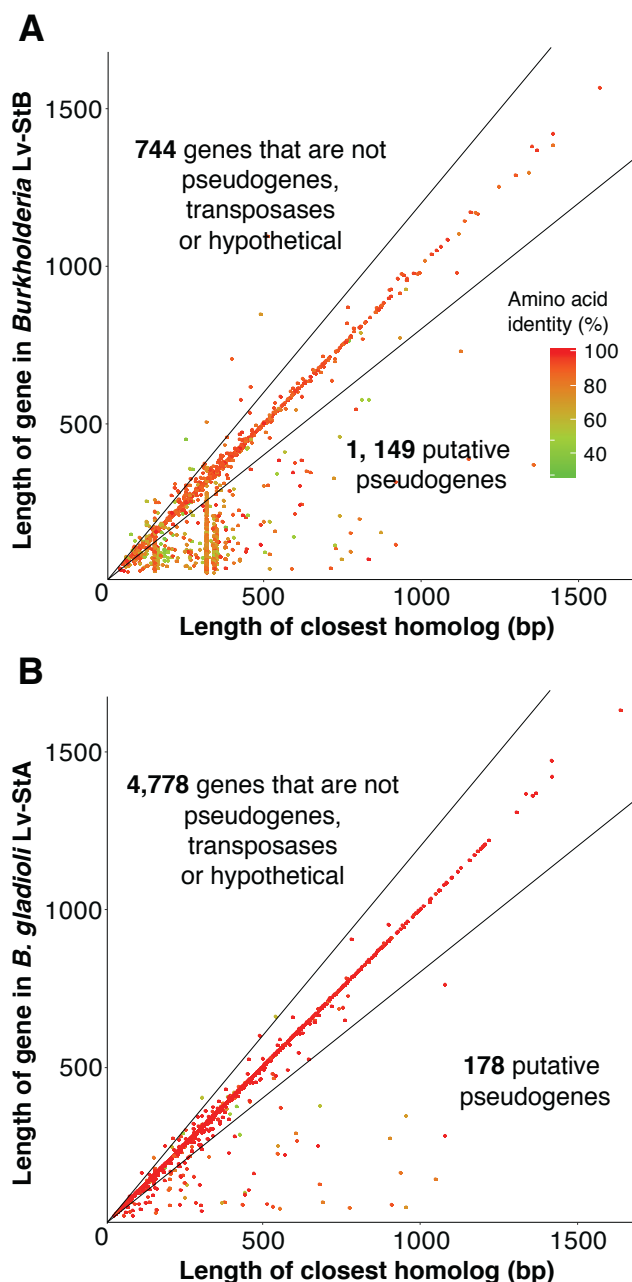
1209      Yang Z. 2014. Molecular Evolution: A Statistical Approach. OUP Oxford.

1210      Yang Z. 2007. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol*

1211      **24**:1586–1591.

1212      Zupancic TJ, Marvo SL, Chung JH, Peralta EG, Jaskunas SR. 1983. RecA-independent

1213      recombination between direct repeats of IS50. *Cell* **33**:629–637.
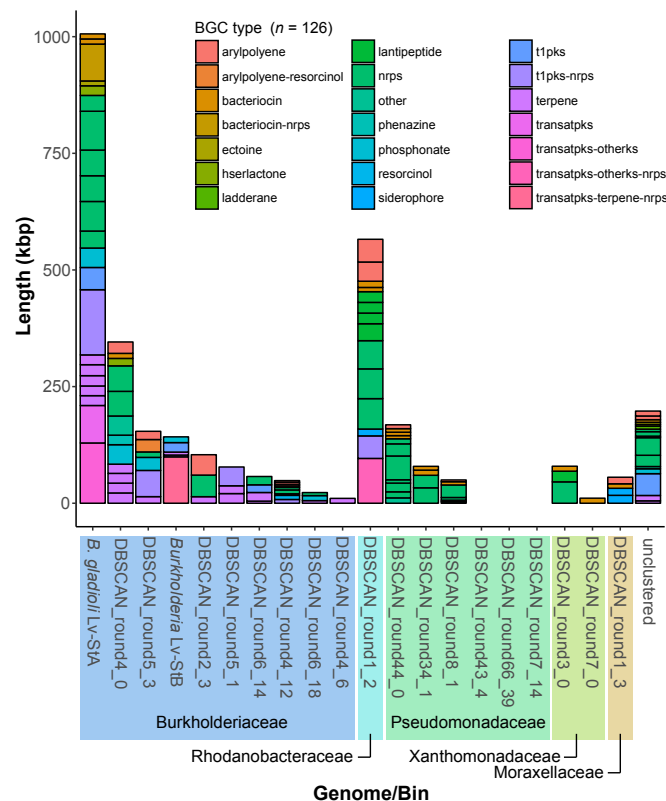
1214

1215

1216

**Figure 1.** Comparison of the lengths of genes in the *Burkholderia* Lv-StB genome (**A**)

and *B. gladioli* Lv-StA (**B**) with the closest homologs identified through BLAST searches

against the NR database (see Methods). Genes which are less than 80% of the length

of the closest relative (i.e. below the lower black line) are putatively assigned as

pseudogenes, as described previously (Lerat and Ochman, 2005; Lopera et al., 2017).

1222    Note: Two vertical groupings of pseudogenes in Lv-StB correspond to multiple copies of

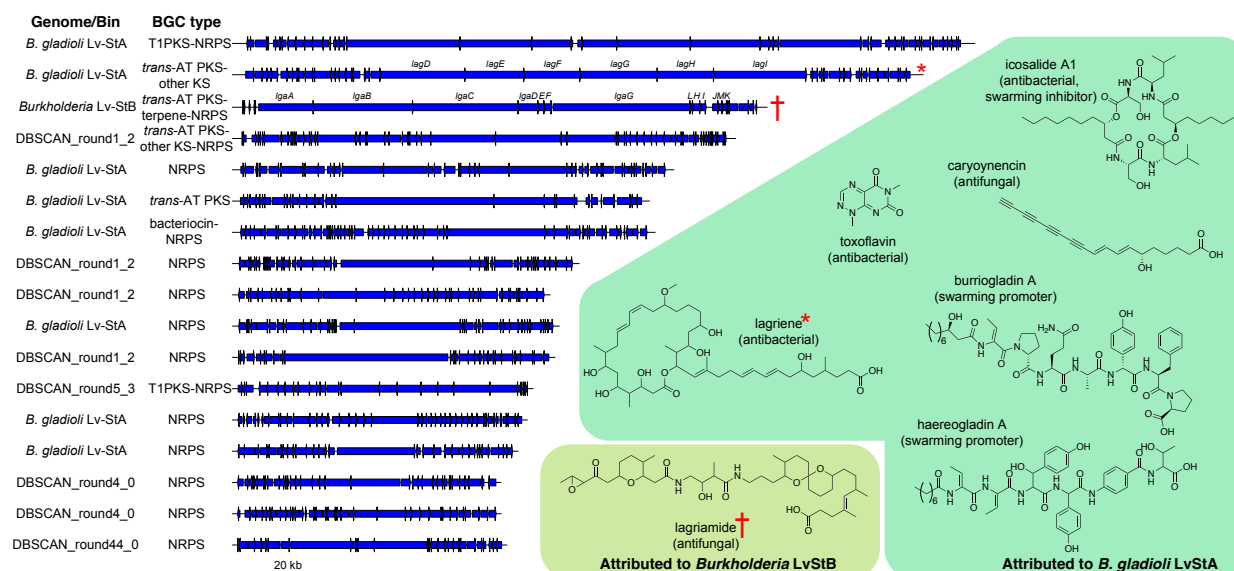1223    a 155 bp hypothetical gene and an IS5 family transposase.
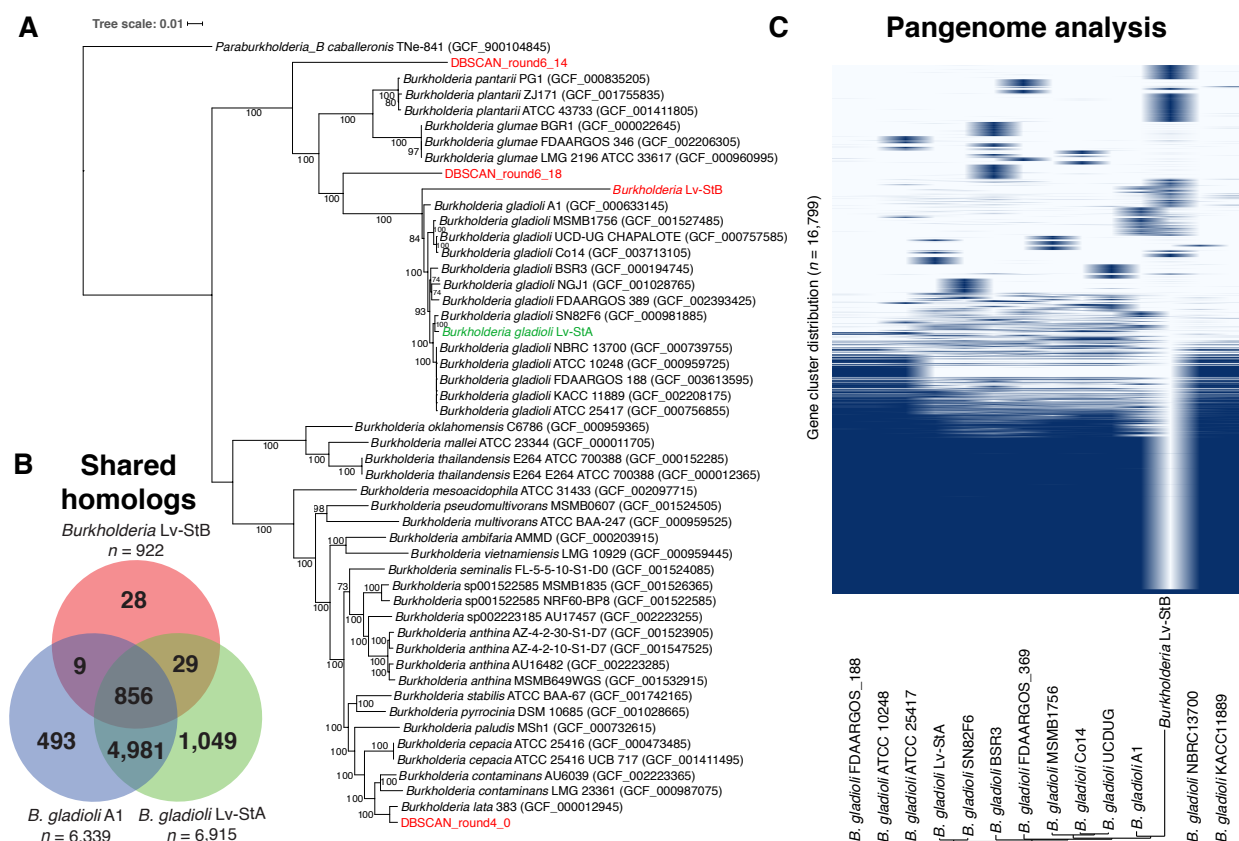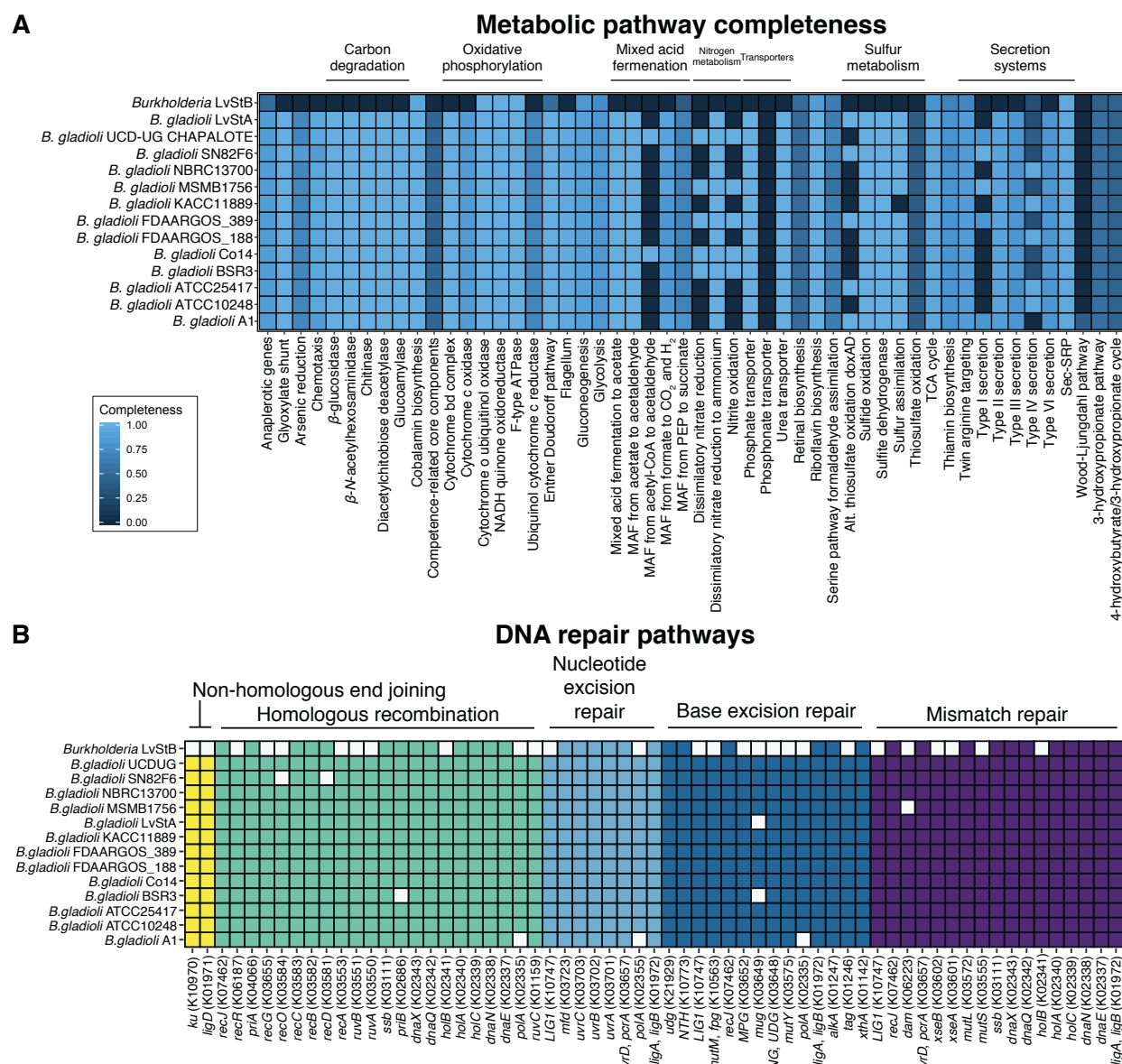
1224

1225

**Figure 2.** Distribution of biosynthetic gene clusters (BGCs) amongst the *L. villosa* metagenome bins and the genome of *B. gladioli* Lv-StA. Colors indicate the type of BGC annotated by antiSMASH (126 identified) (Blin et al., 2017).

57

**Figure 3.** Schematics of all BGCs with greater than 50 kbp length assembled from the

*L. villosa* metagenome and in the *B. gladioli* Lv-StA isolate genome, out of 126 identified

by antiSMASH (left). Shown on the right are structures of compounds putatively

assigned to BGCs in the *B. gladioli* Lv-StA genome (Dose et al., 2018; Flórez et al.,

2017), or putatively assigned to a BGC in the *Burkholderia* Lv-StB genome (Flórez et

al., 2018). Structures highlighted with red symbols have been attributed to the indicated

BGCs.

1242



1243 **Figure 4.** (**A**) Maximum-likelihood multilocus species tree of metagenomic bins

1244 classified in the genus *Burkholderia*, plus *B. gladioli* Lv-StB, using 120 marker gene

1245 protein sequences. Bootstrap proportions greater than 70% are expressed to the left of

1246 each node as a percentage of 1,000 replicates. Metagenomic bins are shown in red,

1247 while the *L. villosa*-associated isolate *B. gladioli* Lv-StA is shown in green. (**B**) Shared

1248 homologous gene groups in *Burkholderia* Lv-StA, *B. gladioli* A1 and *B. gladioli* Lv-StB,

1249 after discounting pseudogenes. Note: Homologous groups are counted only once per

1250 genome (i.e. collapsing paralogs), and therefore counts are lower than absolute gene

1251 counts. Also, 492 genes in the Lv-StB genome were singletons and not included in

1252 homologous groups. (**C**) Hierarchical clustering of homologous gene clusters, showing

1253 presence and absence in *Burkholderia* Lv-StB and closely-related strains of *B. gladioli*.

1254

59

**Figure 5.** Completeness of metabolic and DNA repair pathways in *Burkholderia* Lv-StB in comparison to closely-related strains of *B. gladioli*. (**A**) Completeness of various metabolic pathways as determined by KEGG-decoder. Note: Categories that were not found in any of the examined genomes have been removed. (**B**) The presence (colored squares) and absence (white squares) of genes in the different DNA repair pathways in Lv-StB and related *B. gladioli* strains.

**Figure 6. (A)** Putative sources of genes unique to the *Burkholderia* Lv-StB genome

(compared to *B. gladioli* strains) based on hits to BLASTP searches. Note: Three

proteins of putative phage origin (see Dataset S1I) are not included in the figure. **(B)**

Estimated introgression time for putative HGT gene sets, using the "BAU" substitution

rates (Silva and Santos-Garcia, 2015). Note: For clarity, the gene sets with estimated

introgression time of < 5,000 ya are not labeled. For these and ages estimated with

other substitution rates, see Dataset S1I.

1273 **Table 1.** Genome characteristics of *Burkholderia* symbiont Lv-StB, its relative *B. gladioli*

1274 Lv-StA and other bins obtained from the metagenome.

| Genome | Coverage | Size (% of closest relative)* | Core genes (%) | Category |
|---|---|---|---|---|
| *B. gladioli* Lv-StA | N/A | 96.2 | 95.2 | Nonreduced |
| *B. gladioli* Lv-StB | 1,977 | 23.5 | 85.7 | Reduced |
| DBSCAN_cluster_round5_1 | 355 | 105 | 92.9 | Nonreduced |
| DBSCAN_cluster_round2_3 | 298 | 74.7 | 95.2 | Nonreduced |
| DBSCAN_cluster_round6_18 | 207 | 18.2 | 36.9 | Incomplete |
| DBSCAN_cluster_round4_6 | 170 | 86.1 | 95.2 | Nonreduced |
| DBSCAN_cluster_round43_4 | 142 | 4.74* | 38.1 | Incomplete |
| DBSCAN_cluster_round8_1 | 94 | 39.6 | 25 | Incomplete |
| DBSCAN_cluster_round3_0 | 90 | 86.1 | 81 | Nonreduced |
| DBSCAN_cluster_round5_3 | 71 | 77.2 | 81 | Nonreduced |
| DBSCAN_cluster_round1_2 | 68 | 128 | 92.9 | Nonreduced |
| DBSCAN_cluster_round6_14 | 61 | 20.9 | 39.3 | Incomplete |
| DBSCAN_cluster_round66_39 | 47 | 0.351† | 25 | Incomplete |
| DBSCAN_cluster_round7_14 | 37 | 5.3 | 16.7 | Incomplete |
| DBSCAN_cluster_round34_1 | 34 | 40.9 | 20.2 | Incomplete |
| DBSCAN_cluster_round7_0 | 24 | 62.8 | 31 | Incomplete |
| DBSCAN_cluster_round44_0 | 18 | 58.4 | 34.5 | Incomplete |
| DBSCAN_cluster_round1_3 | 16 | 103 | 94 | Nonreduced |
| DBSCAN_cluster_round4_0 | 8 | 102 | 89.3 | Nonreduced |
| DBSCAN_cluster_round4_12 | 5 | 63.1 | 47.6 | Incomplete |

1275 *Calculated relative to the genome of the closest relative identified by GTDB-Tk (see Dataset
1276 S1B)
1277 †Calculated relative to the average size of 517 *Pseudomonas* genomes taken from the GTDB
1278 database
1279

1280

1281

1282

1283    **Table S1.** Divergence rates used in this study (taken from Silva and Santos-Garcia

1284    2015 (Silva and Santos-Garcia, 2015)).

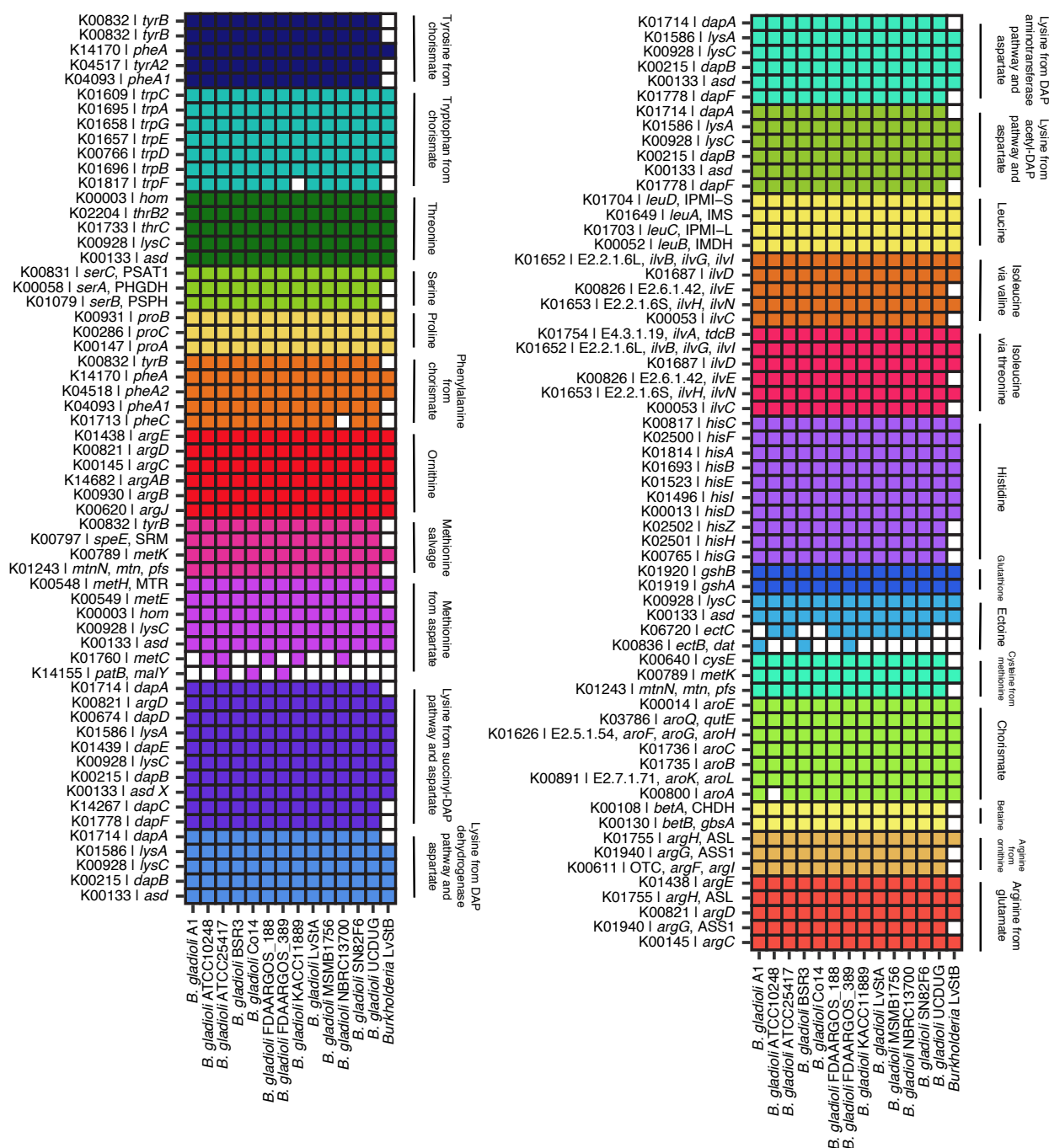| Abbreviation | Symbiont | Host | dS/t | dN/t |
|---|---|---|---|---|
| BFL | *Blochmannia floridans* | *Camponotus floridans* | $8.9 \times 10^{-8}$ | $7.9 \times 10^{-9}$ |
| BPN | *Blochmannia pennsylvanicus* | *Camponotus pennsylvanicus* | $6.4 \times 10^{-8}$ | $5.6 \times 10^{-9}$ |
| BOB | *Blochmannia obliquus* | *Colobopsis obliquus* | $7.9 \times 10^{-8}$ | $5.2 \times 10^{-9}$ |
| BAU | *Baumannia cicadellinicola* | *Graphocephala atropunctata, Homalodisca vitripennis* | $1.1 \times 10^{-8}$ | $1.0 \times 10^{-9}$ |

1285

1286

1287

1288

1289    **Dataset S1.** Comparative analysis of the *Burkholderia* Lv-StB to other genomes in the

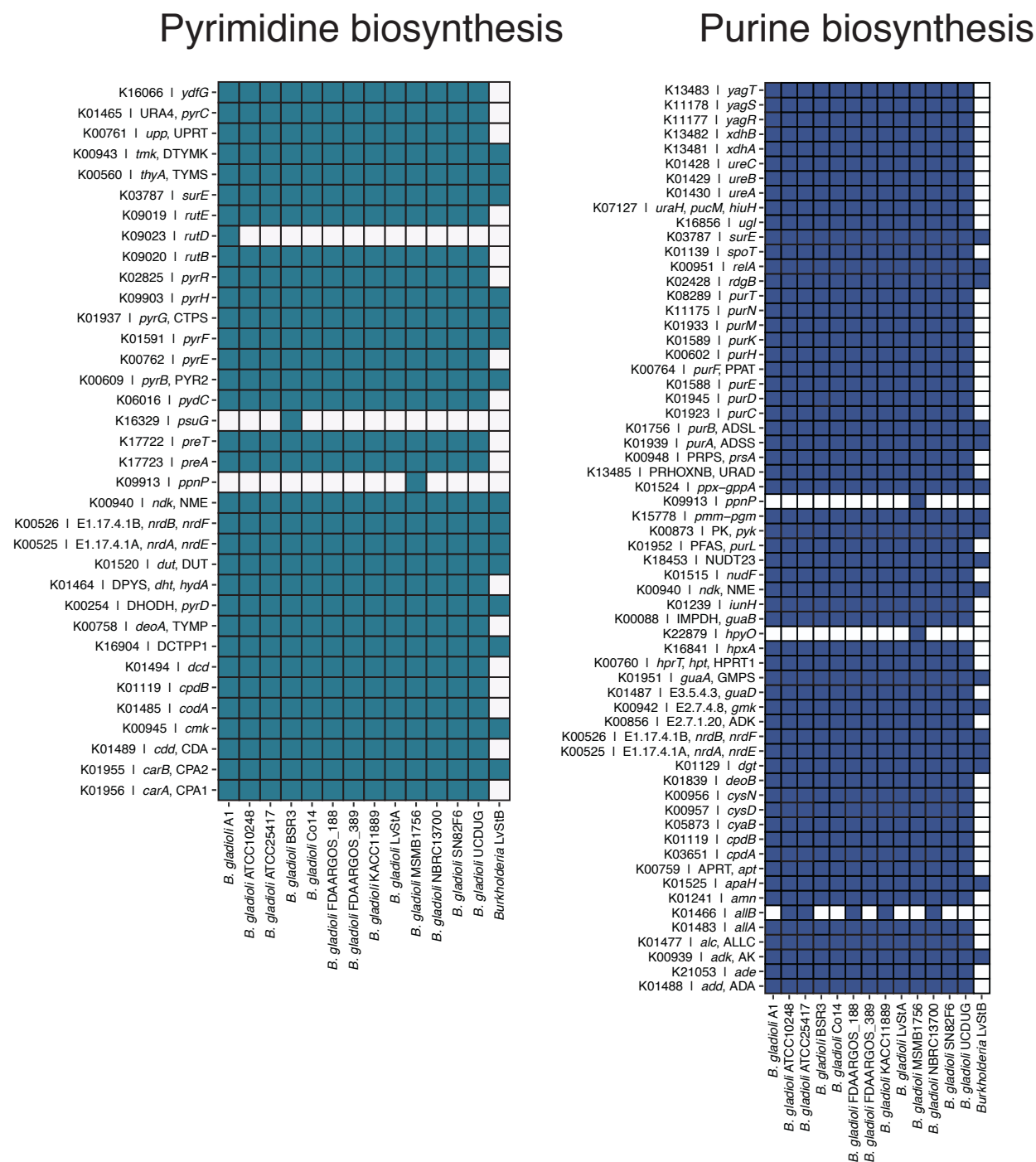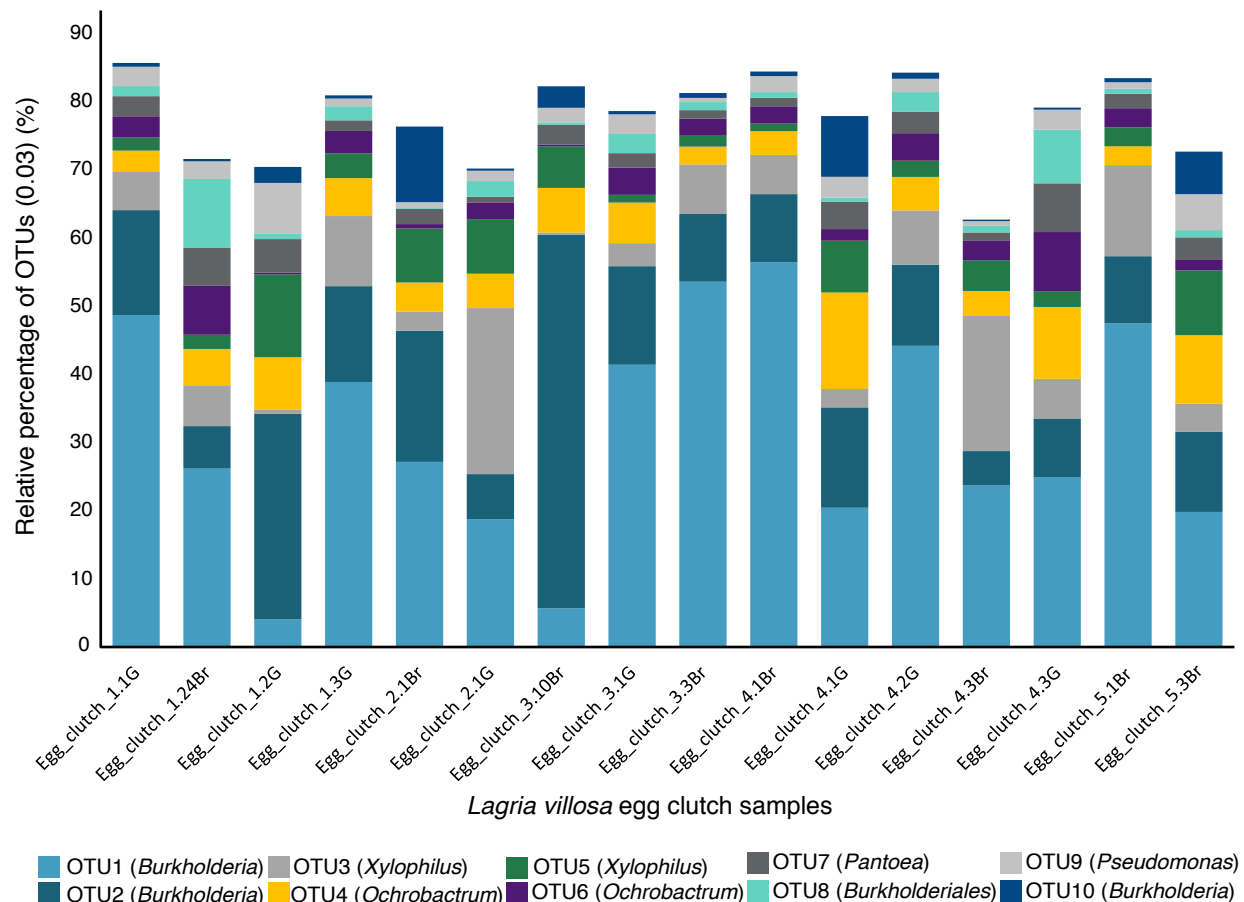1290    genus *Burkholderia*.

1291

1292

1293

**Figure S1.** Completeness of amino acid biosynthesis pathways in *Burkholderia* Lv-StB in comparison to closely-related strains of *B. gladioli*.

1294

1295

1296

1297

1298

**Figure S2.** Completeness of nucleotide biosynthesis pathways in *Burkholderia* Lv-StB

in comparison to closely-related strains of *B. gladioli*.

1301

1302

**Figure S3.** Reanalysis of 16S rRNA amplicon data used in Flórez *et al.* 2017 (Flórez et al., 2017), showing distribution of dominant microbial communities associated with *L. villosa* egg clusters. Amplicon 16S rRNA gene sequences were clustered into operational taxonomic units (OTUs) at a distance of 0.03 as an approximation to bacterial species. The putative taxonomic classification of each OTU is indicated with a colored key. Abundance of the top 10 most abundant OTUs is indicated as a relative percentage of the total reads for each cluster of eggs collected.