1 **Assessing biosynthetic gene cluster diversity in a multipartite**

2 **nutritional symbiosis between herbivorous turtle ants and**

3 **conserved gut symbionts**

4

5 **Running title:** Biosynthetic gene clusters of gut bacteria in ants

6
7 **Anaïs Chanson[1*], Corrie S. Moreau[2], Christophe Duplais[3]**

8 [1]Université de Guyane, UMR EcoFoG, AgroParisTech, CNRS, Cirad, INRAE, Université des

9 Antilles, Kourou, France

10 [2]Departments of Entomology and Ecology & Evolutionary Biology, Cornell University, Ithaca,

11 NY, USA

12 [3]CNRS, UMR EcoFoG, AgroParisTech, Cirad, INRAE, Université des Antilles, Université de

13 Guyane, Kourou, France

14

15 **\*Corresponding author**

16 Anaïs Chanson

17 UMR EcoFoG at the Pasteur Institute of French Guiana

18 23 avenue Pasteur, 97300 Cayenne, France

19 anais.chanson@ecofog.gf

20 Phone: +594 594 293 134

21

22    **Abstract**

23    In insect-microbe nutritional symbioses the symbiont supplements the low nutrient diet of the

24    host by producing amino acids and vitamins, and degrading lignin or polysaccharides. In

25    multipartite mutualisms composed of multiple symbionts from different taxonomical orders, it

26    has been suggested that in addition to the genes involved in the nutritional symbiosis the

27    symbionts maintain genes responsible for the production of metabolites putatively playing a

28    role in the maintenance and interaction of the bacterial communities living in close proximity.

29    To test this hypothesis we investigated the diversity of biosynthetic gene clusters (BGCs) in the

30    genomes and metagenomes of obligate gut symbionts associated with the herbivorous turtle

31    ants (genus: *Cephalotes*). We studied 17 *Cephalotes* species collected across several

32    geographical areas to reveal that (i) mining bacterial metagenomes and genomes provides

33    complementary results demonstrating the robustness of this approach with metagenomic data,

34    (ii) symbiotic gut bacteria have a high diversity of BGCs which is correlated with host

35    geography but not host phylogeny, (iii) the majority of the BGCs comes from the bacteria

36    involved in the nutritional symbiosis supporting conserved metabolic functions for

37    colonization, communication and competition in the gut environment, (iv) phylogenetic

38    analysis of arylpolyene, polyketide (PK), and siderophore shows high similarity between BGCs

39    of a single symbiont across different ant host species, while non-ribosomal peptide (NRP)

40    shows high similarity between BGCs from different bacterial orders within a single host species

41    suggesting multiple mechanisms for genome evolution of these obligate mutualistic gut

42    bacteria.

43

44    **Introduction**

45    Bacterial symbiosis is widespread among insects and has shaped the evolution of their hosts[1,2].

46    The symbiotic interactions between bacterial communities inside the host, and between bacteria

2

47   and hosts, relies massively on metabolites, metabolic pathways and on the enzymes that

48   regulate metabolic flux. Among mutualistic bacteria the nutritional symbionts supplement a

49   large diversity of nutrients to the diet of their insect hosts which often rely on nutritionally

50   limited food sources. In the case of strictly blood-sucking ticks, intracellular bacterial symbionts

51   provide vitamin B to the host[3,4]. Xylophagous termites rely on a diverse community of gut

52   bacteria to degrade the lignocellulose from wood and also fix nitrogen for their hosts.[5]

53   Herbivorous aphids and ants feeding on nitrogen-poor diets depend on gut bacteria to recycle

54   nitrogen from food waste and contribute to the biosynthesis of essential and non-essential amino

55   acids[6,7]. Bees have a conserved gut bacteria community able to degrade polysaccharides[8].

56   Herbivorous beetles have bacteriocytes hosting bacterial symbionts which enrich the host

57   metabolism with aromatic amino acid to support cuticle formation[9]. Herbivorous turtle ants also

58   need their gut bacterial symbionts to support normal cuticle formation (Duplais et al. in review).

59   Another type of microbial symbiosis involves defensive mechanisms where bacterial

60   symbionts enhance insect resistance to a variety of natural enemies including microbes, fungi

61   or nematodes[10,11]. Bacteria can produce a wide range of secondary metabolites which are not

62   required for the immediate survival of the bacteria, but serve important functions in microbial-

63   microbial interactions, such as defense, competition and communication[12,13]. For example, in

64   several insect groups, including beetles, wasps, bees and ants, bacterial symbionts produce a

65   cocktail of antimicrobial compounds to protect the nest[14–16], eggs[17] or a mutualistic fungal

66   strain[15,18].

67   Bacterial metabolites are produced via specific groups of genes, called biosynthetic gene

68   clusters (BGCs), which are located in close proximity to each other in the bacterial genome.

69   Together, the genes composing these clusters encode for specific enzymatic steps in the

70   biosynthetic pathway of metabolites[19]. Bioinformatics tools have been recently developed to

71   mine bacterial genomes[19,20] and retrieve the BGCs of different chemical families including

72   arylpolyene, lanthipeptide, non-ribosomal peptide, polyketide and terpene. Metabolites play

73   several roles in bacterial communities for colonizing tissue, communication through quorum

74   sensing, competition using antimicrobial molecules, and nutrient acquisition with siderophore.

75   BGC studies often focus on the genomes of culturable strains from environmental bacteria[21,22]

76   or host-associated bacteria[23,24] and can contribute to the discovery of bioactive compounds for

77   medicine. Mining of bacterial metagenomes has also been investigated from environmental

78   samples[25], however this approach is only starting to be applied to bacterial communities

79   associated within hosts[26].

80        The study of BGCs from symbiotic bacteria may help us predict the role of metabolites

81   encoded in symbiont genomes and provide a framework to understand which functional genes

82   are maintained to support bacterial colonization and maintenance within a host. In addition the

83   comparison between bacterial BGCs across different hosts[27] from different geographic areas[28]

84   allows us to test what environmental, ecological or phylogenetic factors shape BGC diversity

85   in host-associated bacterial communities.

86        To investigate the potential drivers of BGCs diversity in insect symbionts, we

87   investigated the association between gut bacteria and herbivorous *Cephalotes* ants. In this

88   nutritional symbiosis the nitrogen-poor diet of ants is supplemented by a core microbiome

89   which recycle urea food waste into amino acids beneficial to the ant[29]. Genomic approaches

90   have revealed the redundancy in functions related to the nitrogen flux in the genome of five

91   conserved bacterial families from the gut of 17 species of *Cephalotes*[7]. Unlike most intra- or

92   extracellular single nutritional symbionts, the *Cephalotes* multipartite gut bacteria mutualism

93   may retain metabolic functions selected not solely for direct benefit to hosts, but for sustaining

94   diverse bacterial community members. To assess if the gut bacteria associated with *Cephalotes*

95   possess BGCs, and if the environment and evolutionary history of the *Cephalotes* are correlated

96   with BGCs diversity, we focused on bacterial genomes and metagenomes from the gut of 17

97   *Cephalotes* species collected from North and South America[7]. We retrieved a high number of

98   BGCs from both genomes and metagenomes of bacteria associated with *Cephalotes* ants and

99   tested the correlation with geography and host phylogeny. To study the potential role of

100  metabolites in this multipartite mutualism we recorded the bacterial origin of each BGCs to

101  understand their occurrence in gut bacteria involved in this nutritional symbiosis with

102  *Cephalotes.* An analysis of core genes from BGCs (arylpolyene, non-ribosomal peptide (NRP),

103  polyketide (PK) and siderophore) was performed to identify architectural similarity patterns

104  across different symbionts and host species. Our work shows the relevance of genome and

105  metagenome mining in an insect-microbe symbiosis to study the evolution of BGCs in symbiont

106  genomes across the host phylogeny within an ecology and evolutionary framework.

107

108  **Material and Methods**

109  **Genomes and metagenomes analysis**

110  The 14 genomes of cultured gut bacteria and 18 metagenomes of *Cephalotes* gut bacteria were

111  obtained from JGI-IMG version 5.0[30] (Table S1 and S2 respectively) from the previous projects

112  Gs0085494 (“*Cephalotes varians* microbial communities from the Florida Keys, USA”),

113  Gs0114286 (“Symbiotic bacteria isolated from *Cephalotes varians*”), Gs0117930 (“*Cephalotes*

114  ants gut microbiomes”) and Gs0118097 (“Symbiotic bacteria isolated from *Cephalotes*

115  *rohweri*”)[7]. The metagenomic data used in this study have two metagenomes from the same *C.*

116  *varians* species. However, for the metagenome of *C. varians* PL010W the sequencing quality

117  is very different from all the other metagenomes (number of reads, GC content), therefore it

118  was excluded in our analysis (Table S2). The metagenomes were analyzed via the software

119  Anvi’o version 5.5[31] to sort the different bacterial families composing each metagenome into

120  distinct bins. In this analysis, the fasta sequence of a metagenome is used to create a contig

121  database and a profile database. Then, this contig database is visualized and bins are manually

122   created to maximize completeness while minimizing redundancy. Finally, the software

123   CheckM version 1.1[32] was used to identify the taxonomy of each bin (Figure S1).

124

125   **Bacterial biosynthetic gene clusters analysis**

126   The bacterial biosynthetic gene clusters (BGCs) of each genome and each metagenomic bin

127   were analyzed with antiSMASH 5.0[33] (Figure S1) with the following analysis options: strict

128   detection, and activation of search for KnownClusterBlast, ClusterBlast, SubClusterBlast and

129   Active site finder. BGCs smaller than 5kb were then filtered out of the data and were not

130   included. The taxonomic classification of each cluster was verified to the genus level using the

131   software Blast+[34]. NaPDos[35] was used to classify the ketosynthase (KS) domain and

132   condensation (C) domain sequences required in the biosynthesis of PK and NRP respectively

133   and to infer the KS and C phylogenies.

134

135   **Distance matrix calculations**

136   The published *Cephalotes* phylogeny was retrieved[36] and the packages ape[37] and phytools[38] of

137   the R software version 3.6.1[39] were used to exclude from this phylogeny the *Cephalotes* species

138   from which no genome or metagenome was available. The cophenetic distance option in the R

139   package stats[39] was used on this pruned phylogeny to create the *Cephalotes* phylogenetic

140   distance matrix. The cophenetic distance replaces original pairwise distances between the

141   *Cephalotes* species by the computed distances between their clusters. To calculate the

142   metagenomic BGCs distance matrix, a BGC matrix containing the numbers and types of each

143   BGC found in each metagenome was created. Then the distance function (Euclidean method)

144   of the R package stats was used on the BGC matrix to calculate the metagenomic BGCs distance

145   matrix. The Euclidean method was chosen because it calculates the absolute distance between

146   two samples with continuous numerical variables, without removing redundancies.

147

## Genetic networking construction

149 The genetic networking of the genomic and metagenomic bacterial BGCs was generated using

150 BiG-SCAPE version 20191011[40] (Figure S1). The network was constructed with the three

151 following options: "--include_singletons", "--mix", and "--cutoffs 1.0". The resulting similarity

152 matrix was filtered with different thresholds between 0.6 and 0.8[41–43], and the threshold 0.65

153 was chosen because it filtered out enough to form different groups while maximizing the

154 number of bacterial BGCs in each network[41,43]. The filtered similarity matrix was visualized

155 with Cytoscape version 3.7.2[44], using the MCL clusterization algorithm from clusterMaker2[45]

156 version 1.3.1 with an index value I = 2.0. To create a genetic networking of the bacterial BGCs

157 we the method of Lin et al.[46] and implemented in the software BiG-SCAPE by Navarro-Munoz

158 et al.[40,47]

159

## BGCs genomic core gene analysis

161 Although BGCs are made of many domains, some have been deemed the "genomic core". The

162 genomic core is a group of one or a few domains within a single BGC, which contain all the

163 required modules needed for the BGC to be functional and are highly conserved across bacteria.

164 The genomic core of four types of BGCs (arylpolyene, NRP, siderophore and T1PK) were

165 analyzed using the CORASON software version 1.0[40] (Figure S1). First, a database for each

166 type of BGCs was created from all the BGCs of the same type recovered in the genomes and

167 metagenomes, which were obtained from our previous AntiSMASH analysis. Then a BGC from

168 each of our database was chosen as the reference BGC for its database, and a gene from the

169 biosynthetic core of this reference BGC was chosen as the query protein. The choice of the

170 reference BGC and query protein for a database were made taking different variables into

171 account: (1) the reference BGC must be one of the longest BGCs in the database, (2) the query

172   protein must come from a biosynthetic core gene or an additional biosynthetic gene close to the

173   core, and (3) the query protein must be present in at least half of the BGCs in the database.

174   Once a reference BGC and a protein query were selected for a database, the CORASON

175   software was used to determine the genomic core of this type of BGCs, and to infer the

176   phylogenies. If less than half of the BGCs present in a database appear in the phylogeny, or if

177   many BGCs appear more than once in the phylogeny (which may happen if the selected query

178   protein belongs to a biosynthetic gene that can be found multiple times in the same BGC), the

179   reference BGC and query protein previously selected were deemed to be inadequate for the

180   analysis. In this case, the process to choose a reference BGC and query protein were done once

181   more, following the same procedure as before. The reference BGC and query protein selected

182   for the four types of BGCs studied are presented in Table S3.

183

184   **Statistical analysis**

185   To test for correlations between the phylogeny and geography of the *Cephalotes* ants and the

186   BGCs found in the bacterial metagenomes the Chi-squared function of the R package stats[48]

187   was implemented. The *Cephalotes* phylogenetic distance matrix and the bacterial BGCs

188   distance matrix from the metagenome analysis were compared using the Mantel test of the R

189   package ape[37]. The PCoA between the number and types of bacterial BGCs from the

190   metagenome analysis and *Cephalotes* geography were performed using the *prcomp* function of

191   the R package stats and the R package ggplot[49]. The correlations between the number of each

192   type of bacterial BGCs from the metagenome analysis and *Cephalotes* geography were

193   calculated using the one-way ANOVA function with the Tukey's pairwise of the Past software

194   version 3.25[50].

195

196

## Results

### Mining bacterial genomes

The whole bacterial genomes (N=14) studied herein originated from cultured gut bacteria associated with *Cephalotes varians* and *Cephalotes rohweri*. Among these 14 bacterial genomes, only three genomes did not possess any BGC, while in the others a total of 31 biosynthetic gene clusters (BGCs) were found and the type and length of the 31 BGCs are listed in Table S4. In bacteria associated with *C. varians*, seven types of BGCs were found for a total of 10 BGCs: arylpolyene, beta-lactone, ectoine, nonribosomal peptide (NRP), resorcinol, polyketide type I (T1PK) and terpene. In bacteria associated with *C. rohweri*, six types of BGCs were found for a total of 21 BGCs: arylpolyene, beta-lactone, ladderane, NRP, siderophore and T1PK. In *C. varians* the majority of BGCs originate from Burkholderiales with only one BGC coming from Pseudomonadales symbionts whereas in *C. rohweri* they originate from different bacterial orders (Burkholderiales, Opitutales, Pseudomonadales and Xanthomonadales). The number of bacterial BGCs from *C. varians* and *C. rohweri* were compared by an ANOVA test (Figure S2). Results show that only in the case of NRPs the number of BGCs are statically different between the two ants host (p=0.0036).

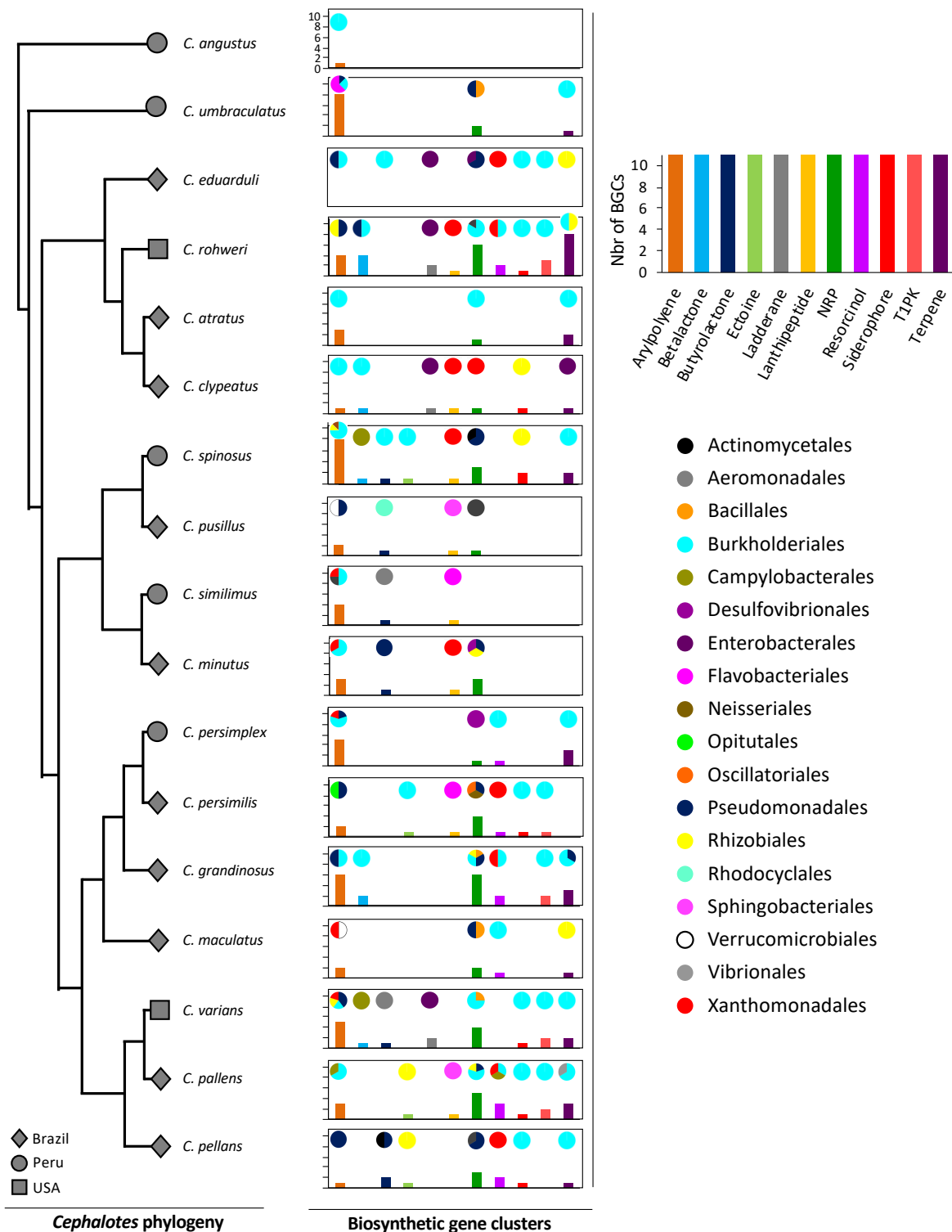### Assessing the quality of metagenomic binning

The analysis of the 17 *Cephalotes* metagenomes resulted in 168 bins and using AntiSMASH a total of 233 BGCs were found in 102 bins out of 168 (Table S5). Completeness of the BGCs ranges from 2.47% to 100%. More than half the bins (87 out of 168) have a predicted completeness level higher than 85%, 66 bins have a completeness level between 50 and 85%, and only 15 bins have a completeness level lower than 50%. The completeness is strongly correlated with the redundancy and number of contigs and genes in each bin (Figure S3A; N=168; PCA analysis, $1^{st}$ dimension=92.19%, $2^{nd}$ dimension=7.74%). The length of each

9

222    bacterial BGC is highly linked to the number of genes in the corresponding BGC (Figure S3B;

223    N=233; Spearman correlation, S=354810, p<2.2e-16, rho=0.832), and moderately linked to the

224    completeness of the bin in which the BGC is found (Figure S3C; N=233; Spearman correlation,

225    S=1579300, p=9.252e-13, rho=0.251). Among the 168 bins, 31 different bacterial orders were

226    identified, and among them 21 bacterial orders possess at least one bacterial genus predicted to

227    have a BGC (Figure S4A). The bacterial families having the highest number of BGCs are

228    Burkholderiales (41 BGCs), Pseudomonadales (41 BGCs), Rhizobiales (41 BGCs), and

229    Xanthomonadales (24 BGCs). The bacterial specificity of each bin has been assessed by

230    analyzing the number of different bacterial orders present in each bin (Figure S4B). A bin is

231    considered specific if it is associated with only one bacterial order. More than half of the bins

232    are specific, and 75% of the bins contain less than two bacterial orders.

233

234    **Mining bacterial metagenomes across the *Cephalotes* phylogeny**

235    From the 233 BGCs retrieved in the genomes and metagenomes we identified 20 different types

236    of BGCs: acyl amino acids, arylpolyene, bacteriocin, beta-lactone, butyrolactone,

237    cyclodipeptide (CDP), ectoine, furan, hserlactone, ladderane, lanthipeptide, linear azol(in)e-

238    containing peptides (LAP), non-ribosomal peptide (NRP), phenazine, resorcinol, siderophore,

239    polyketide type I (T1PK), terpene, and thiopeptide. The arylpolyene and NRP BGCs were the

240    more numerous with 66 and 56 respectively retrieved in the metagenomes. The type and length

241    of all the BGCs are presented in Table S6. Some bacterial BGCs found in the metagenomic

242    analysis are also hybrids made of different types of BGCs: arylpolyene-beta-lactone-resorcinol,

243    arylpolyene-ladderane, arylpolyene-resorcinol, and NRP-T1PK (Table S6). Although we

244    cannot rule misassembling during the genome construction, we are confident in the hybrid

245    NRP-T1PK because the sequence of the KS domain correspond to hybrid KS sequence in

246    NaPDoS database. In Figure 1, the most abundant types of BGCs across the different

247    *Cephalotes* species are investigated, representing a total of 202 BGCs originating from 19

248    different bacterial orders.



249
250
251    **Figure 1.** Diversity of biosynthetic gene clusters of bacteria associated with *Cephalotes* turtle

252    ants. The grey symbols represent the country from which an ant originated. The graphs

253    represent the bacterial BGCs found in the metagenomes of each *Cephalotes* species in this

254    study. The bars of the graphs represent the number of BGC of each type found in the

255    metagenomes. The pies represent the proportion of bacterial orders from which each BGC

256    originated.

257

258        The 17 species of *Cephalotes* studied here were collected in three different countries,

259    10 species from Brazil, five species from Peru and two species from the USA. Across these

260    species, a high number of BGCs are found in each species with a mean of ~14 (Figure 1; Table

261    S6). The lowest number of BGCs was found in *C. angustus* (N=1), the highest number of BGCs

262    was recorded for *C. rohweri* (N=31) and the median number of BGCs across species is 14. The

263    correlations between the diversity of bacterial BGCs and the phylogeny or geography of the

264    *Cephalotes* hosts were assessed (Figure 2). After a Chi-squared test confirmed that the

265    phylogeny and geography of the ants are not correlated (p = 0.558), a Mantel test indicates no

266    correlation between the *Cephalotes* phylogeny and the type and number of bacterial BGCs

267    (Figure 2A; N=17; p=0.205; Mantel test). On the other hand the type and number of bacterial

268    BGCs are statistically different in ants collected in Brazil, Peru and USA (Figure 2B; N=17;

269    p=0.004; ANOVA). A principal coordinates analysis shows that the bacterial BGCs associated

270    with *Cephalotes* coming from the same country group together (Figure 2B; PCoA, $1^{st}$

271    axis=51.89%, $2^{nd}$ axis=25.42%).

272        The number of each type of bacterial BGCs associated with *Cephalotes* from different

273    countries were compared using ANOVA (Figure 2C; N=17). Arylpolyene BGCs are more

274    numerous in bacteria of ants from Peru than from Brazil, but not in other comparisons (Peru-

275    Brazil: p=0.035; USA-Peru: p=0.604; USA-Brazil: p=0.409). Beta-lactone, ladderane, T1PK,

276    and terpene BGCs are more numerous in bacteria of ants from the USA compare to ants from

277    Brazil (p=0.047, p=0.014, p=0.049, p=0.048 respectively) and Peru (p=0.037, p=0.010,

278    p=0.038, p=0.040 respectively), while there is no difference between ants from Brazil and Peru

279    for these BGCs (p=0.315, p=0.417, p=0.536, p=1.000 respectively). NRP BGCs are more

280    numerous in bacteria of ants from the USA than in bacteria of ants from Peru (Peru-Brazil:

281    p=0.163; USA-Peru: p=0.046; USA-Brazil: p=0.241). No statistical differences were present in

282    the number of butyrolactone, ectoine, lanthipeptide, resorcinol, and siderophore BGCs in

283    bacteria of ants between the different geographic areas (p=0.720, p=0.560, p=0.760, p=0.880,
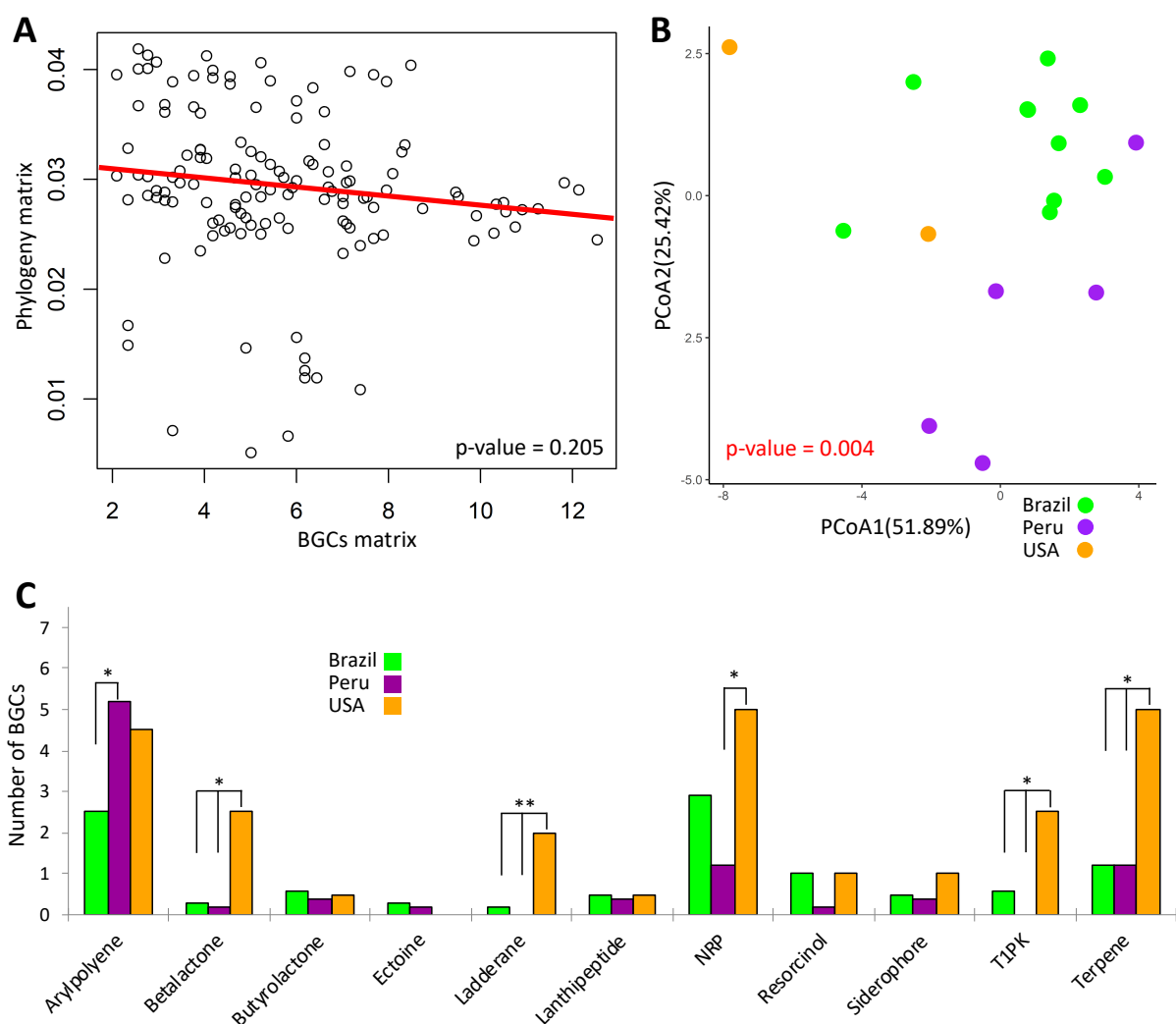
284    p=0.620 respectively).



285
286

287    **Figure 2.** (A) Correlation between the *Cephalotes* phylogeny distance matrix and the bacterial

288    BGCs distance matrix. (B) PCoA of the bacterial BGCs across *Cephalotes* geography. (C) One-

13

289    way ANOVA test showing the statistical differences between the most abundant bacterial BGCs

290    across *Cephalotes* geography. The symbol * represents p-value < 0.05.

291

292    **Genetic networking and BGCs genomic core genes phylogenies**

293    To assess the genetic architectural diversity of BGCs (number and type of gene, gene order,

294    gene domain, and BGC length) we created a genetic network of the BGCs found in the genomes

295    and metagenomes of *Cephalotes*-associated bacteria (Figure 3). We first used the type of BGCs

296    and the bacterial order to color the outer and inner circle of nodes respectively (Figure 3A). In

297    this analysis, 31 networks of BGCs are formed, each containing between 2 to 24 nodes with 90

298    nodes as singletons. Each network is mainly formed by the same type of bacterial BGCs, despite

299    the fact that a third of the networks contains one or two BGC nodes of a different type than the

300    rest of the nodes composing this network. Additionally, in half of the networks, BGCs

301    originated from the same bacterial order, while in the other half the BGCs originated from two

302    or three different bacterial orders. In Figure 3B we used the species of ant and the host

303    geography to color the inner and outer nodes respectively. However, neither the host phylogeny

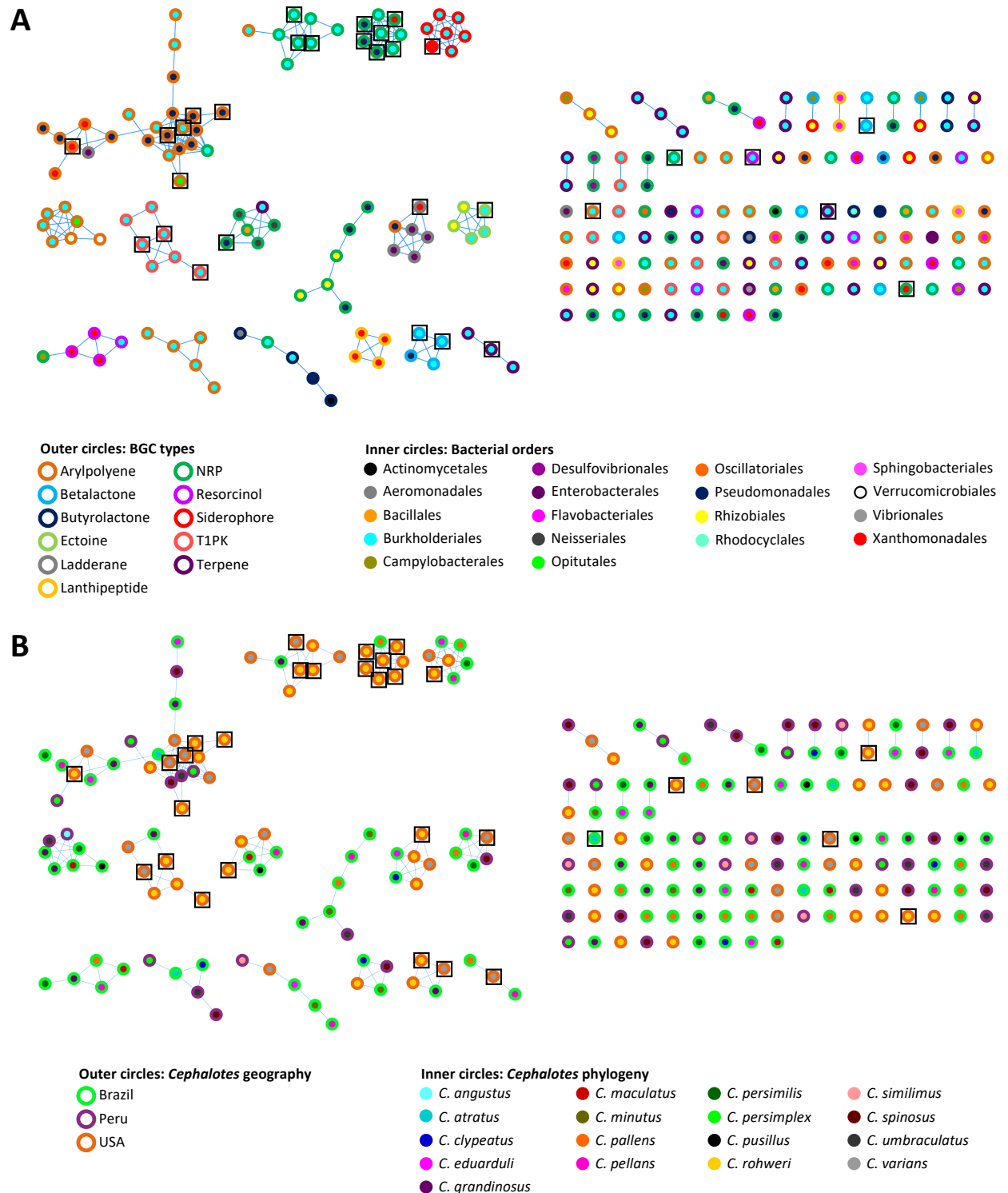304    nor geography structures the grouping of the bacterial BGCs.

**Figure. 3** Genetic networks of the bacterial BGCs in *Cephalotes* genomes and metagenomes, with a distance threshold of 0.65 and an MCL index value of 2.0. Color codes are respective of the bacterial BGCs types and bacterial order in (A) and the *Cephalotes* geography and phylogeny in (B). Black squares represent the 31 BGCs found in bacterial genomes.

311

312          The phylogenies of the genomic core genes of the BGCs were inferred using

313     CORASON (Figure 4, Figure S5-S6, Table S5). The four types of BGCs, arylpolyene, NRP,

314     siderophore, and T1PK, were chosen for this analysis because they are thought to play an

315     important role in bacteria colonization, interaction and competitiveness[51–54]. Using an

316     arylpolyene KS gene as a query 38 different BGCs were retrieved in the arylpolyene BGC

317     phylogeny out of the 66 detected in the *Cephalotes* bacterial genomes and metagenomes (Figure

318     S5). These 38 arylpolyene BGCs form several distinct clades. One clade is composed of nine

319     BGCs which are grouped in the largest arylpolyene network cluster (Figure 4A). These BGCs

320     originated from the genomes and metagenomes of Burkholderiales and Pseudomonadales

321     isolated from six different ant species (Figure 4A). In this clade, the arylpolyene BGCs have a

322     length of between 24 000 to 43 000 nucleotides (nt) and between 26 to 44 genes. In the NRP

323     BGC phylogeny based on an AMP-binding gene, 33 different NRP were retrieved out of the 57

324     detected in the *Cephalotes* bacterial genomes and metagenomes (Figure S6). These 33 NRP

325     BGCs form several distinct clades including a polytomy composed of seven BGCs found in the

326     network cluster (Figure 4B). In this clade, the NRP BGCs were identified in the genomes and

327     the metagenome from three bacterial orders (Burkholderiales, Pseudomonadales, and

328     Xanthomonadales) isolated from a single ant host (*C. rohweri*). The NRP BGCs have a length

329     of between 23 000 to 53 000 nt and between 25 to 53 genes. The PKS-KS gene results in the

330     T1PK BGC phylogeny of ten BGCs out of the 15 detected in the *Cephalotes* bacterial genomes

331     and metagenomes (Figure 4C). All BGCs originated from the genomes and metagenomes of

332     Burkholderiales isolated from six different ant species. The T1PK BGCs have a length of

333     between 7 000 to 10 000 nt and between 4 to 11 genes. Seven T1PK BGCs were found a single

334     network cluster and the three T1PK BGCs left all fall as singletons in the genetic network

335     analysis (Figure 3A). The IucA-IucC gene was used to infer the siderophore BGC phylogeny.

336    Eight siderophores were retrieved out of the 10 detected in the *Cephalotes* bacterial genomes

337    and metagenomes (Figure 4D). The siderophore BGCs have a length of between 7 000 to 10

338    000 nt and between 4 to 11 genes. A This network matches seven siderophore BGCs from the

339    genomes and metagenomes of Burkholderiales and Xanthomonadales each associated with a
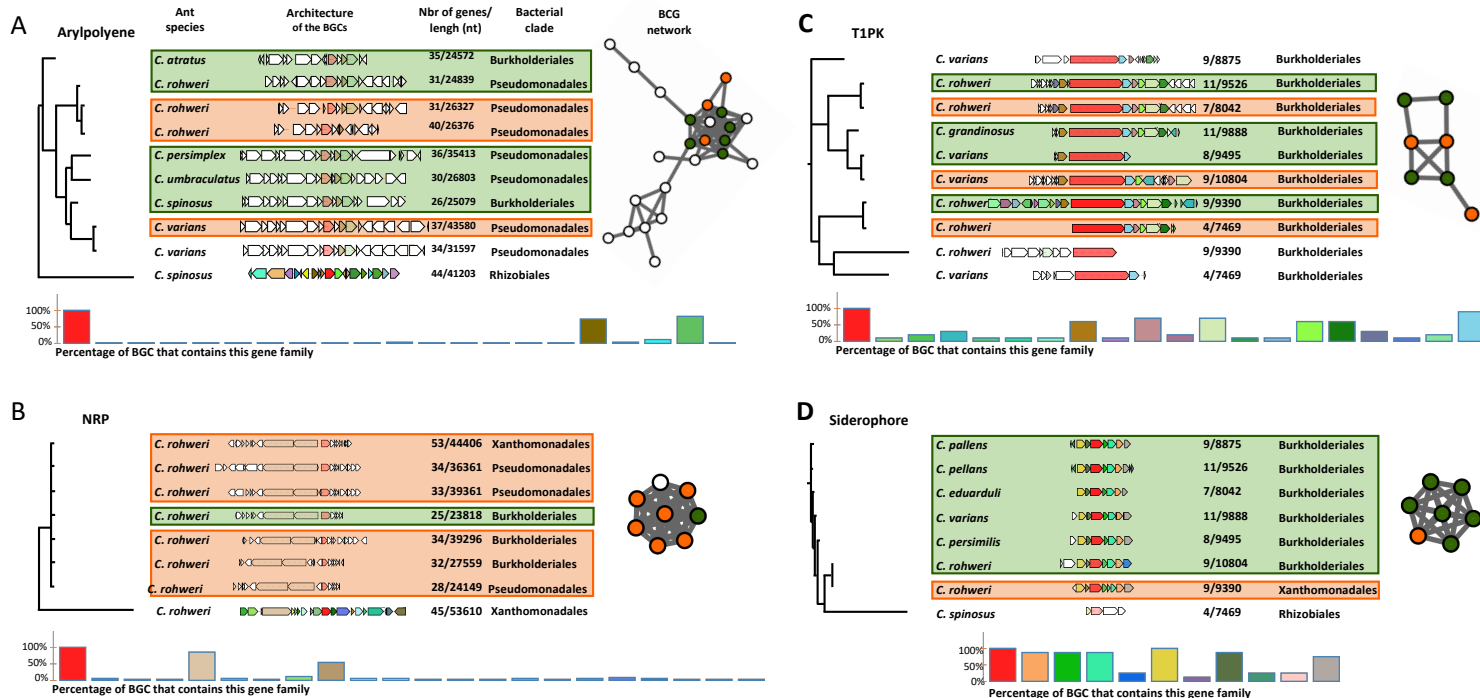
340    different ant species.



341

**Figure 4.** Representation of selected phylogenetic clades of arylpolyene (A), NRP (B), T1PK

(C) and siderophore (D) BGCs. The bacterial origin, genomic composition, size and number of

genes of each BGC, and the ant species hosting this bacterium are indicated next to the

phylogeny. The histograms represent the percentage of BGC that contains this gene family. The

BGCs network were retrieved from the previous analysis. Orange and green colors indicate if

the BGCs originate from a bacterial genome or metagenome respectively. The white nodes in

the network represent BGC which were not found in the selected phylogenetic clade. Full

phylogenies of arylpolyene and NRP are found in Figures S5-S6 respectively.

350

351       The phylogenies of the elongation domains of NRP synthase and PK synthase,

352   respectively C and KS domains, were inferred using NaPDos[35] to check if the sequences of

353   known metabolites were matching our data (Figure 5). The majority of the identified C domain

354   sequences (66 out of 89) from genomes and metagenomes belong to the cyclization domain

355   class catalyzing both peptide bond formation and subsequent cyclization. The retrieved KS

356   domain sequences (N=13) from genomes and metagenomes fall into a single clade

357   corresponding to *Cis*-AT modular class or the iterative domain class. Unfortunately, no close

358   similarity with the identified C and KS sequences of known non-ribosomal peptide and

359   polyketide were found when comparing with NaPDos and the NCBI database, thus no

360   prediction could be made about the chemical structure of bacterial metabolites.
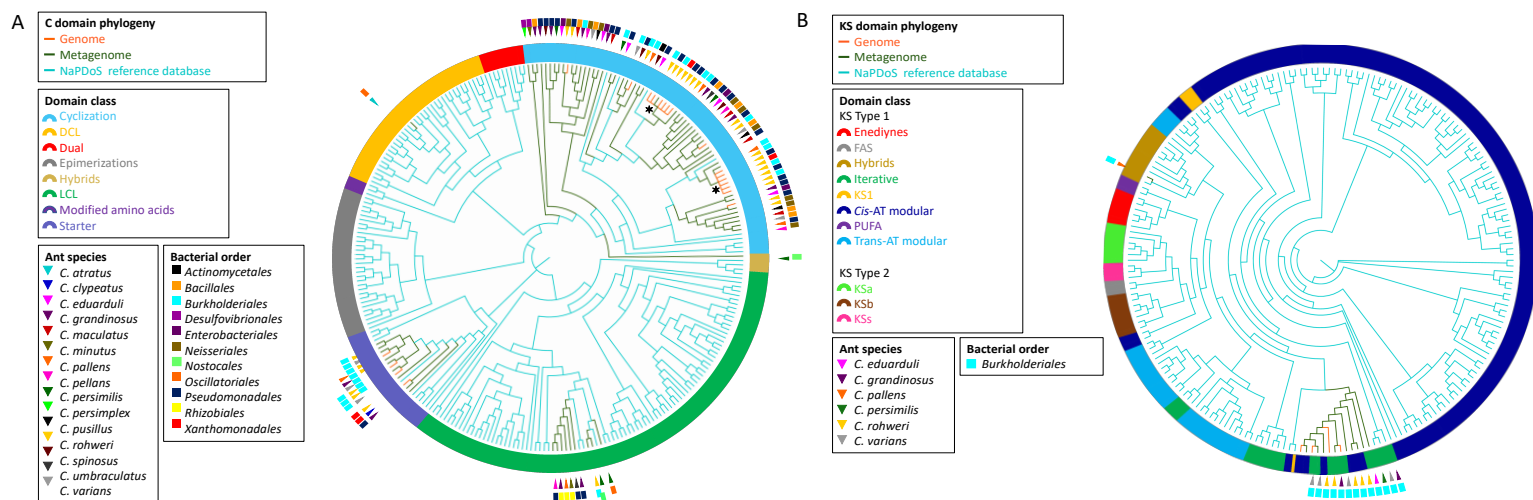


361

**Figure 5.** Phylogeny of retrieved C and KS elongation domain sequences in the NRP and PK

363   assemblies respectively. The phylogenies were inferred by implementing the C domain

364   sequences (A) and the KS domain sequences (B) from the *Cephalotes* genomes (orange

365   branches) and metagenomes (green branches) with the NaPDoS reference database (blue

366   branches). The ring surrounding the tips represent C domain class or KS domain class. The

367   bacterial order from which each domain sequence was retrieved is indicated by the colored

368   rectangle. The ant species associated with each sequence is indicated by the colored triangle.

369    The symbols * designate the clade of the two C domain sequences find in the NRP genomic

370    core analysis (see Fig. 4).

371

372

373    **Discussion**

374    Since the multiplication of computational tools for mining microbial genomes, the discovery of

375    BGCs of interest for pharmaceutical applications has become a growing field[19]. This approach

376    is complementary to natural products chemistry and metabolomics in helping with the

377    prioritization of strains synthesizing potential drugs having new chemical structures. Genome

378    mining studies on large scale prokaryote genomic data have revealed that BGCs are distributed

379    in large gene cluster families, with the vast majority of them still uncharacterized[41]. Some of

380    these BGCs are distributed widely across the entire bacterial domain, with the most prominent

381    being arylpolyenes, polyketides, saccharides and siderophores, having known and putative

382    functions in microbe-host and microbe-microbe interactions[41]. The rise of advanced genomic

383    tools to mine genomes and predict the chemical structure of metabolites from BGCs will

384    certainly benefit the field of host-microbe symbiosis. Herein we present the first genome and

385    metagenome mining study of gut bacteria across different species of herbivorous turtle ants to

386    assess the diversity of BGCs in a nutritional symbiosis. When applied across the host phylogeny

387    a genome mining approach provides crucial information about how the evolutionary history of

388    the host may impact the diversity of BGCs in the symbiont genome. For instance BGCs can

389    undergo a vast array of evolutionary process through *de novo* assembly, gene

390    duplication/deletion, horizontal gene transfer which drive the microbial chemodiversity[55]. The

391    systematic study of BGCs diversity of symbionts across host species from different geographic

392    areas may provide a list of BGCs likely to play a crucial functional role for the symbiont and

393    for the host.

394    Although the use of bacterial metagenomes can provide accurate genomic

395    information[56], one can argue that mining metagenomes using AntiSMASH[33] may generate

396    unreal BGCs by incorrectly assembled bins. Recently mining the metagenome for BGCs has

397    been reported for the microbiome involved in a defensive symbiosis associated with the eggs

398    of the beetle *Lagria villosa*[26] and with the marine sponge *Mycale hentscheli*[57]. Our results show

399    high similarities in the type and size of BGCs detected both in the bacterial genomes and

400    metagenomes, with nine prominent types of BGCs detected: beta-lactone, ectoine, ladderane,

401    NRP, resorcinol, siderophore, T1PK and terpene (Table S4 and S6). The genomic core gene

402    analysis for arylpolyene, NRP, T1PK and siderophore show strong resemblance in the

403    architecture and genomic core structure between bacterial genomes and metagenomes. Indeed,

404    bacterial genomes and metagenomes of the same bacterial order and same host ant species fall

405    together in the genetic network (Figure 3), the genomic core gene analysis (Figure 4), and the

406    C and KS domain phylogenies (Figure 5). Together, these results show that the BGCs analysis

407    of metagenomes are a good representation of the BGCs diversity in samples indicating that the

408    retrieved BGCs from metagenomic data are not chimeras created by misassembling

409    metagenomic bins.

410    As the 17 species of *Cephalotes* were collected in Brazil, Peru and the USA we

411    questioned whether the host phylogeny and/or geography have an impact on symbiont BGC

412    diversity. Our results show that the number and type of bacterial BGCs are not correlated with

413    *Cephalotes* evolutionary history (p=0.205; Mantel test) but are more similar in *Cephalotes*

414    samples collected in the same geographic area (p=0.004; ANOVA) (Figure 2). A significantly

415    higher number of beta-lactone, ladderane, NRP, T1PK, and terpene BGCs associated with

416    *Cephalotes* from the USA than associated with *Cephalotes* from South America. These results

417    are in accordance with the geographic mosaic theory of coevolution which stipules that in

418    strongly interacting species, geography and community ecology can shape coevolution through

419  local adaptation[58,59]. Implications of this theory in the context of metabolites and biosynthetic

420  gene clusters is starting to be investigated. Recently, it was shown that in soil bacteria the NRP

421  and PK main biosynthetic domains, respectively adenylation and ketosynthase, are different

422  across the Australian continent[60]. According to their results, the main factor for these

423  differences is the sample latitude which can be one geographical factor explaining the

424  differences between results from North and South America. As for the BGC architectures they

425  are correlated with host geography (Figure 3B). This is in accordance with previous work on

426  antibiotic-producing bacterial symbiont and fungus-growing ants[61]. In this conserved symbiosis

427  it has been demonstrated that the geographical isolation of populations in Central America has

428  an impact on antibiotic potency of the locally adapted symbiont but not on the BGC

429  architecture.

430    Four bacterial orders possess more than 80% of all the 233 BGCs detected in the

431  *Cephalotes* bacterial genomes and metagenomes: Burkholderiales (44%), Pseudomonadales

432  (15%), Xanthomonadales (15%) and Rhizobiales (9%) (Figure 1, Table S4 and S6).

433  Interestingly, these four bacterial orders belong to the *Cephalotes* core bacterial community[7].

434  The *Cephalotes* core bacterial community is maintained across the *Cephalotes* phylogeny and

435  benefit these herbivorous ants with a low-nitrogen diet by recycling nitrogen from urea into

436  essential and non-essential amino acids. In many insect nutritional symbioses the endosymbiont

437  resides in host cells or bacteriocytes and have a reduced genome with several copies of the

438  functional genes responsible for the biosynthesis of amino acids and vitamins to enrich in

439  nutrients the host diet. This is contrasting with the multipartite mutualism in the gut of

440  *Cephalotes* turtle ants. The maintenance of several bacteria is more complex compare to other

441  symbioses involving a single gut symbiont (bean bugs[70]), or bacteriocytes endosymbionts

442  (*Camponotus* ants[71], beetles[9]). The conserved bacterial symbionts in *Cephalotes* have

443  redundancy in N-metabolism and our results show these bacterial symbionts possess the vast

444    majority of the BGCs detected. This support the idea that metabolic functions of a complex

445    microbiome are often selected not solely for direct benefit of hosts, but for sustaining the gut

446    community[72]. Since the *Cephalotes* core symbionts are not transmitted maternally but are newly

447    acquired by *Cephalotes* ants in each generation, the maintenance of genes promoting gut

448    colonization to outcompete the transient bacteria is likely high selected to be maintained.

449          This hypothesis is in accordance with recent findings in a bean bug which needs to

450    acquire a specific *Burkholderia* gut symbiont from the soil environment in every new

451    generation[70]. They demonstrated that a large number of *Burkholderia* species were able to

452    access the bean bug gut, but in experiments of co-inoculations, the specific *Burkholderia*

453    symbiont always outcompete all the other *Burkholderia* species. In addition in bees it has been

454    shown that gut symbionts produce amino acids and siderophores, which are required to colonize

455    the insect gut[72]. In the *Cephalotes* gut the five main symbionts from different bacterial orders

456    need to colonize gut tissue, tolerate other community members and outcompete invading

457    bacteria. The means used by co-occurring gut symbionts to survive may be through the

458    production of specific metabolites mediating these functions. If this is true in this system, we

459    need to identify which types of BGCs would be necessary for gut symbionts.

460          To test the similarities of BGC structure across different hosts, we performed

461    CORASON genomic core gene analysis on four types of BGCs: arylpolyene, NRP, T1PK and

462    siderophore (Figure 4, Figure S5-S6). The selected clades presented in Figure 4A and 4B show

463    similarities in the sequence of the core genes for arylpolyene BGCs originating from

464    Burkholderiales and Pseudomonadales, and for NRP BGCs originating from Burkholderiales,

465    Pseudomonadales, and Xanthomonadales. These arylpolyene and NRP BGCs are each

466    clustered in the networking analysis (Figure 3 and 4) demonstrating the architectural

467    similarities. In the NRP BGC clade of Figure 4B the query gene AMP-binding and the two

468    genes coding for a C domain have 100% similarity and 100% coverage across the BGCs

22

469    retrieved from the six genomes and metagenomes of Burkholderiales, Pseudomonadales, and

470    Xanthomonadales from the same host *C. rohweri*. This high similarity between the six C

471    domain sequences from three bacterial orders is also confirmed in the phylogenetic analysis

472    (Figure 5). This pattern suggests a possible convergence of the core genes in the gut of *C.*

473    *rohweri* possibly via horizontal gene transfer. In insect-bacteria symbiosis horizontal gene

474    transfer has been shown to occur even in bacteria going through genome reduction[26]. However

475    the three sequences did not match any known sequence in the database limiting the

476    identification of the bacterial origin of these genes. On the other hand, the T1PK and

477    siderophore core gene phylogenies (Figure 4B and 4C) show that BGCs from Burkholderiales

478    across different host species share a strong likely genetic conservation. Several identical genes

479    of T1PK and siderophore BGCs are present across the core gene phylogenies supporting a

480    conserved architecture which suggests a common origin or convergent evolution of

481    Burkholderiales BGCs across *Cephalotes* species. Symbiotic convergent evolution has been

482    reported in insect-bacteria symbiosis[74,75], showing insect symbionts retain the ability to produce

483    a certain number of proteins and functions which have an importance for both bacteria and host

484    survival.

485        A final remark concerns the potential division of labor in a multipartite mutualism. In

486    honey bees the gut microbiota are dominated by five coevolved bacterial clusters. However the

487    gut bacterial species have different abilities in digesting polysaccharides which can be

488    explained by a divergence into different ecological niches inside the gut of their hosts[8]. In the

489    *Cephalotes* gut four of the five conserved symbionts possess the majority of the retrieved BGCs

490    but Opitutales a major player in the *Cephalotes* nutritional symbiosis is lacking BGCs.

491    Although multiple bins of Opitutales were retrieved only two BGCs (arylpolyene) were

492    detected (Figure S7). In the *Cephalotes* gut only Opitutales possess urease genes able to

493    transform urea into ammonia readily used by the four other symbionts. Opitutales are localized

494 in the midgut whereas Burkholderiales, Pseudomonadales, Rhizobiales, and Xanthomonadales

495 are found in the ileum of the hindgut (Flynn et al. in prep.). This spatial isolation prevents direct

496 interaction of Opitutales with other symbionts and their lack of BCGs supports the idea that

497 BGCs have an essential role for symbionts living in close proximity.

498      Our study sheds light on the question of why many obligate host-associated, non-

499 endosymbiotic bacteria do not have highly reduced genomes, like their endosymbiotic

500 counterparts. In this multipartite nutritional symbiosis with herbivorous turtle ants the core

501 bacterial community members not only maintain genes to supplement the host's nitrogen-poor

502 diet[7], but we show they also retain BGCs that likely contribute to their colonization of the gut

503 and mediate bacterial-bacterial communication and interactions. Future work on host-

504 associated bacterial communities should begin to explore the specific metabolites expressed by

505 these microbes and their functional outcome in these complex ecosystems.

506

507 **Author Contributions**

508 AC, CSM, and CD conceived the study. AC performed the analyses. AC, CSM, and CD

509 interpreted the results. AC drafted the manuscript. AC, CSM, and CD revised the manuscript.

510

516

517 **Conflict of Interest**

518 The authors declare no competing financial and non-financial competing interests.

519

526

527 **References**

528 1.    Hosokawa, T., Kikuchi, Y., Nikoh, N., Shimada, M. & Fukatsu, T. Strict host-symbiont

529        cospeciation and reductive genome evolution in insect gut bacteria. *PLoS Biol.* **4**,

530        1841–1851 (2006).

531 2.    Jiggins, F. M. & Hurst, G. D. D. Rapid insect evolution by symbiont transfer. *Science*

532        *(80-. ).* **332**, 185–186 (2011).

533 3.    Duron, O. *et al.* Tick-Bacteria Mutualism Depends on B Vitamin Synthesis Pathways.

534        *Curr. Biol.* **28**, 1896-1902.e5 (2018).

535 4.    Binetruy, F. *et al.*  Microbial community structure reveals instability of nutritional

536        symbiosis during evolutionary radiation of Amblyomma ticks . *Mol. Ecol.* 1016–1029

537        (2020) doi:10.1111/mec.15373.

538 5.    Boucias, D. G. *et al.* The hindgut lumen prokaryotic microbiota of the termite

539        Reticulitermes flavipes and its responses to dietary lignocellulose composition. *Mol.*

540        *Ecol.* **22**, 1836–1853 (2013).

541 6.    Hansen, A. K. & Moran, N. A. Aphid genome expression reveals host-symbiont

542        cooperation in the production of amino acids. *Proc. Natl. Acad. Sci. U. S. A.* **108**,

543        2849–2854 (2011).

544  7.  Hu, Y. *et al.* Herbivorous turtle ants obtain essential nutrients from a highly conserved

545     nitrogen-recycling gut microbiome. *Nat. Commun.* (2018).

546  8.  Zheng, H. *et al.* Division of labor in honey bee gut microbiota for plant polysaccharide

547     digestion. *PNAS* **116**, 25909–25916 (2019).

548  9.  Anbutsu, H. *et al.* Small genome symbiont underlies cuticle hardness in beetles. *Proc.*

549     *Natl. Acad. Sci.* 201712857 (2017) doi:10.1073/pnas.1712857114.

550  10. Scarborough, C. L., Ferrari, J. & Godfray, H. C. J. Ecology: Aphid protected from

551     pathogen by endosymbiont. *Science (80-. ).* **310**, 1781 (2005).

552  11. Hedges, L. M., Brownlie, J. C., O'Neill, S. L. & Johnson, K. N. Wolbachia and virus

553     protection in insects. *Science (80-. ).* **322**, 702 (2008).

554  12. Hibbing, M. E., Fuqua, C., Parsek, M. R. & Peterson, S. B. Bacterial competition:

555     surviving and thriving in the microbial jungle. *Nat. Rev. Microbiol.* **8**, 15–25 (2010).

556  13. Netzker, T. *et al.* Microbial interactions trigger the production of antibiotics. *Curr.*

557     *Opin. Microbiol.* **45**, 117–123 (2018).

558  14. Engl, T. *et al.* Evolutionary stability of antibiotic protection in a defensive symbiosis.

559     *Proc. Natl. Acad. Sci.* 201719797 (2018) doi:10.1073/pnas.1719797115.

560  15. Matarrita-Carranza, B. *et al.* Evidence for widespread associations between neotropical

561     hymenopteran insects and Actinobacteria. *Front. Microbiol.* **8**, 1–17 (2017).

562  16. Chevrette, M. G. *et al.* The antimicrobial potential of Streptomyces from insect

563     microbiomes. *Nat. Commun.* **10**, (2019).

564  17. Flórez, L. V. *et al.* An antifungal polyketide associated with horizontally acquired

565     genes supports symbiont-mediated defense in Lagria villosa beetles. *Nat. Commun.* **9**,

566     (2018).

567  18. Barke, J. *et al.* A mixed community of actinomycetes produce multiple antibiotics for

568     the fungus farming ant Acromyrmex octospinosus. *BMC Biol.* **8**, 109 (2010).

569    19.    Ziemert, N., Alanjary, M. & Weber, T. The evolution of genome mining in microbes –

570            a review. *Nat. Prod. Rep.* **33**, 988–1005 (2016).

571    20.    Adamek, M., Spohn, M., Stegmann, E. & Ziemert, N. Mining Bacterial Genomes for

572            Secondary Metabolite Gene Clusters. in *Antibiotics: Methods and Protocols* (ed. Sass,

573            P.) vol. 1520 23–47 (Springer Science+Business Media, 2017).

574    21.    Crits-Christoph, A., Diamond, S., Butterfield, C. N., Thomas, B. C. & Banfield, J. F.

575            Novel soil bacteria possess diverse genes for secondary metabolite biosynthesis. *Nature*

576            **558**, 440–444 (2018).

577    22.    Paulus, C. *et al.* New natural products identified by combined genomics-metabolomics

578            profiling of marine Streptomyces sp. MP131-18. *Sci. Rep.* **7**, 1–11 (2017).

579    23.    Beemelmanns, C., Guo, H., Rischer, M. & Poulsen, M. Natural products from microbes

580            associated with insects. *Beilstein J. Org. Chem.* **12**, 314–327 (2016).

581    24.    Kroiss, J. *et al.* Symbiotic streptomycetes provide antibiotic combination prophylaxis

582            for wasp offspring. *Nat. Chem. Biol.* **6**, 261–263 (2010).

583    25.    Cuadrat, R. R. C., Ionescu, D., Dávila, A. M. R. & Grossart, H. P. Recovering

584            genomics clusters of secondary metabolites from lakes using genome-resolved

585            metagenomics. *Front. Microbiol.* **9**, 1–13 (2018).

586    26.    Waterworth, S. C. *et al.* Horizontal gene transfer to a defensive symbiont with a

587            reduced genome amongst a multipartite beetle microbiome. *bioRxiv* **11**, 780619 (2019).

588    27.    Ashen, J. B. & Goff, L. J. Molecular and ecological evidence for species specificity

589            and coevolution in a group of marine algal-bacterial symbioses. *Appl. Environ.*

590            *Microbiol.* **66**, 3024–3030 (2000).

591    28.    Edlund, A., Loesgen, S., Fenical, W. & Jensen, P. R. Geographic distribution of

592            secondary metabolite genes in the marine actinomycete Salinispora arenicola. *Appl.*

593            *Environ. Microbiol.* **77**, 5916–5925 (2011).

594   29.   Russell, J. A. *et al.* Bacterial gut symbionts are tightly linked with the evolution of

595         herbivory in ants. *Proc. Natl. Acad. Sci.* **106**, 21236–21241 (2009).

596   30.   Chen, I. M. A. *et al.* IMG/M v.5.0: An integrated data management and comparative

597         analysis system for microbial genomes and microbiomes. *Nucleic Acids Res.* **47**,

598         D666–D677 (2019).

599   31.   Eren, A. M. *et al.* Anvi'o: An advanced analysis and visualization platformfor 'omics

600         data. *PeerJ* **2015**, 1–29 (2015).

601   32.   Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:

602         Assessing the quality of microbial genomes recovered from isolates, single cells, and

603         metagenomes. *Genome Res.* **25**, 1043–1055 (2015).

604   33.   Blin, K. *et al.* antiSMASH 5.0: updates to the secondary metabolite genome mining

605         pipeline. *Nucleic Acids Res.* **47**, W81–W87 (2019).

606   34.   Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic Local

607         Alignment Search Tool. *J. Mol. Biol.* **215**, 403–410 (1990).

608   35.   Ziemert, N. *et al.* The natural product domain seeker NaPDoS: A phylogeny based

609         bioinformatic tool to classify secondary metabolite gene diversity. *PLoS One* **7**, 1–9

610         (2012).

611   36.   Price, S. L. *et al.* Renewed diversification is associated with new ecological

612         opportunity in the Neotropical turtle ants. *J. Evol. Biol.* **27**, 242–258 (2014).

613   37.   Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of phylogenetics and evolution

614         in R language. *Bioinformatics* **20**, 289–290 (2004).

615   38.   Revell, L. J. phytools: An R package for phylogenetic comparative biology (and other

616         things). *Methods Ecol. Evol.* **3**, 217–223 (2012).

617   39.   Team, R. C. R: A language and environment for statistical computing. *R Found. Stat.*

618         *Comput. Vienna, Austria* (2019).

619    40.    Navarro-Munoz, J. C. *et al.* A computational framework to explore large-scale

620            biosynthetic diversity. *Nat. Chem. Biol.* **16**, 60–68 (2020).

621    41.    Cimermancic, P. *et al.* Insights into secondary metabolism from a global analysis of

622            prokaryotic biosynthetic gene clusters. *Cell* **158**, 412–421 (2014).

623    42.    Schorn, M. A. *et al.* Sequencing rare marine actinomycete genomes reveals high

624            density of unique natural product biosynthetic gene clusters. *Microbiol. (United*

625            *Kingdom)* **162**, 2075–2086 (2016).

626    43.    Adamek, M. *et al.* Comparative genomics reveals phylogenetic distribution patterns of

627            secondary metabolites in Amycolatopsis species. *BMC Genomics* **19**, 1–15 (2018).

628    44.    Shannon, P. *et al.* Cytoscape: A software Environment for integrated models of

629            biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).

630    45.    Morris, J. H. *et al.* ClusterMaker: A multi-algorithm clustering plugin for Cytoscape.

631            *BMC Bioinformatics* **12**, 1–14 (2011).

632    46.    Lin, K., Zhu, L. & Zhang, D. Y. An initial strategy for comparing proteins at the

633            domain architecture level. *Bioinformatics* **22**, 2081–2086 (2006).

634    47.    Yeong, B. M. *BiG-SCAPE*. https://git.wageningenur.nl/yeong001/BGC_networks

635            (2016).

636    48.    R Core Team. An introduction to dplR. *R Found. Stat. Comput. Vienna, Austria* (2019)

637            doi:10.1108/eb003648.

638    49.    Gómez-Rubio, V. ggplot2 - Elegant Graphics for Data Analysis (2nd Edition). *J. Stat.*

639            *Softw.* **77**, 3–5 (2017).

640    50.    Hammer, Ø., Harper, D. A. T. & Ryan, P. D. Past: Paleontological statistics software

641            package for education and data analysis. *Palaeontol. Electron.* **4**, (2001).

642    51.    Schöner, T. A. *et al.* Aryl Polyenes, a Highly Abundant Class of Bacterial Aryl

643            Polyenes, aHighly Abundant Class of Bacterial Natural Products, Are

644       FunctionallyRelated to Antioxidative Carotenoids. *ChemBioChem* **17**, 247–253 (2016).

645   52.   Zhang, W., Li, Z., Miao, X. & Zhang, F. The screening of antimicrobial bacteria with

646       diverse novel nonribosomal peptide synthetase (NRPS) genes from South China sea

647       sponges. *Mar. Biotechnol.* **11**, 346–355 (2009).

648   53.   Shevchuk, O. *et al.* Polyketide synthase (PKS) reduces fusion of Legionella

649       pneumophila-containing vacuoles with lysosomes and contributes to bacterial

650       competitiveness during infection. *Int. J. Med. Microbiol.* **304**, 1169–1181 (2014).

651   54.   Chaturvedi, K. S., Hung, C. S., Crowley, J. R., Stapleton, A. E. & Henderson, J. P. The

652       siderophore yersiniabactin binds copper to protect pathogens during infection. *Nat.*

653       *Chem. Biol.* **8**, 731–736 (2012).

654   55.   Rokas, A., Mead, M. E., Steenwyk, J. L., Raja, H. A. & Oberlies, N. H. Biosynthetic

655       gene clusters and the evolution of fungal chemodiversity. *Nat. Prod. Rep.* (2020)

656       doi:10.1039/c9np00045c.

657   56.   Alneberg, J. *et al.* Genomes from uncultivated prokaryotes: A comparison of

658       metagenome-assembled and single-amplified genomes 06 Biological Sciences 0604

659       Genetics. *Microbiome* **6**, 1–14 (2018).

660   57.   Storey, M. A. *et al.* Metagenomic Exploration of the Marine Sponge Mycale hentscheli

661       Uncovers Multiple Polyketide-Producing Bacterial Symbionts. *MBio* **11**, 1–16 (2020).

662   58.   Thompson, J. N. Coevolution: The Geographic Mosaic of Coevolutionary Arms Races.

663       *Curr. Biol.* **15**, 992–994 (2005).

664   59.   Nuismer, S. L. Parasite Local Adaptation in a Geographic Mosaic. *Evolution (N. Y).* **60**,

665       24 (2006).

666   60.   Lemetre, C. *et al.* Bacterial natural product biosynthetic domain composition in soil

667       correlates with changes in latitude on a continent-wide scale. *Proc. Natl. Acad. Sci. U.*

668       *S. A.* **114**, 11615–11620 (2017).

669  61.  Caldera, E. J., Chevrette, M. G., McDonald, B. R. & Currie, C. R. Local Adaptation of

670       Bacterial Symbionts within a Geographic Mosaic of Antibiotic Coevolution. *Appl.*

671       *Environ. Microbiol.* **85**, (2019).

672  62.  Fenical, W. Natural products chemistry in the marine environment. *Science (80-. ).*

673       **215**, 923–928 (1982).

674  63.  Gross, H. & Konig, G. Terpenoids from marine organisms: unique structures and their

675       pharmacological potential. *Phytochem. Rev.* **5**, 115–141 (2006).

676  64.  De Carvalho, C. C. C. R. & Fernandes, P. Production of metabolites as bacterial

677       responses to the marine environment. *Mar. Drugs* **8**, 705–727 (2010).

678  65.  Miller, D. L., Smith, E. A. & Newton, I. L. G. A bacterial symbiont protects honey

679       bees from fungal disease. *bioRxiv* 2020.01.21.914325 (2020)

680       doi:10.1101/2020.01.21.914325.

681  66.  Tanaka, A., Tapper, B. A., Popay, A., Parker, E. J. & Scott, B. A symbiosis expressed

682       non-ribosomal peptide synthetase from a mutualistic fungal endophyte of perennial

683       ryegrass confers protection to the symbiotum from insect herbivory. *Mol. Microbiol.*

684       **57**, 1036–1050 (2005).

685  67.  Piel, J. Metabolites from symbiotic bacteria. *Nat. Prod. Rep.* **26**, 338–362 (2009).

686  68.  Lesueur, D., del Carro Rio, M. & Diem, H. G. Modification of the growth and the

687       competitiveness of a Bradyrhizobium strain obtained through affecting its siderophore-

688       producing ability. in *Developments in Plant and Soil Sciences* (ed. Abadía, J.) 59–66

689       (1995).

690  69.  Johnson, L. J. *et al.* Biosynthesis of an extracellular siderophore is essential for

691       maintenance of mutualistic endophyte-grass symbioses. *Proc. sixth Int. Symp. fungal*

692       *endophytes grasses Grassl. Res. Pract. Ser.* 177–179 (2007).

693  70.  Itoh, H. *et al.* Host–symbiont specificity determined by microbe–microbe competition

694    in an insect gut. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 22673–22682 (2019).

695  71.  Gil, R. *et al.* The genome sequence of Blochmannia floridanus: Comparative analysis

696    of reduced genomes. *Proc. Natl. Acad. Sci. U. S. A.* **100**, 9388–9393 (2003).

697  72.  Powell, J. E., Leonard, S. P., Kwong, W. K., Engel, P. & Moran, N. A. Genome-wide

698    screen identifies host colonization determinants in a bacterial gut symbiont. *Proc. Natl.*

699    *Acad. Sci. U. S. A.* **113**, 13887–13892 (2016).

700  73.  Andrews, M. *et al.* Horizontal transfer of symbiosis genes within and between rhizobial

701    genera: Occurrence and importance. *Genes (Basel).* **9**, (2018).

702  74.  McCutcheon, J. P. & Moran, N. A. Functional convergence in reduced genomes of

703    bacterial symbionts spanning 200 my of evolution. *Genome Biol. Evol.* **2**, 708–718

704    (2010).

705  75.  López-Sánchez, M. J. *et al.* Evolutionary convergence and nitrogen metabolism in

706    Blattabacterium strain Bge, primary endosymbiont of the cockroach Blattella

707    germanica. *PLoS Genet.* **5**, (2009).

708