1  **Phylogenomics of expanding uncultured environmental Tenericutes**

2  **provides insights into their pathogenicity and evolutionary**

3  **relationship with Bacilli**

4  Yong Wang[1,*], Jiao-Mei Huang[1,2], Ying-Li Zhou[1,2], Alexandre Almeida[3,4], Robert D.

5  Finn[3], Antoine Danchin[5,6], Li-Sheng He[1]

6  [1]Institute of Deep Sea Science and Engineering, Chinese Academy of Sciences, Sanya,

7  Hai Nan, China

8  [2] University of Chinese Academy of Sciences, Beijing, China

9  [3]European Molecular Biology Laboratory, European Bioinformatics Institute

10  (EMBL-EBI), Wellcome Genome Campus, Hinxton, UK

11  [4]Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, UK.

12  [5]Department of Infection, Immunity and Inflammation, Institut Cochin INSERM

13  U1016 - CNRS UMR8104 - Université Paris Descartes, 24 rue du Faubourg

14  Saint-Jacques, 75014 Paris, France

15  [6]School of Biomedical Sciences, Li Kashing Faculty of Medicine, University of Hong

16  Kong, 21 Sassoon Road, SAR Hong Kong, China

17

18  [*]**Corresponding author:**

19  Yong Wang, PhD

20  Institute of Deep Sea Science and Engineering, Chinese Academy of Sciences

21  No. 28, Luhuitou Road, Sanya, Hai Nan, P.R. of China

22  **Phone:** 086-898-88381062

23  **E-mail:** wangy@idsse.ac.cn

24  **Running title:** *Genomics of environmental Tenericutes*

25  **Keywords:** Bacilli; autotrophy; pathogen; gut microbiome; environmental

26  Tenericutes

**ABSTRACT**

The metabolic capacity, stress response and evolution of uncultured environmental Tenericutes have remained elusive, since previous studies have been largely focused on pathogenic species. In this study, we expanded analyses on Tenericutes lineages that inhabit various environments using a collection of 840 genomes. Several novel environmental lineages were discovered inhabiting the human gut, ground water, bioreactors and hypersaline lake and spanning the Haloplasmatales and Mycoplasmatales orders. A phylogenomics analysis of Bacilli and Tenericutes genomes revealed that some uncultured Tenericutes are affiliated with novel clades in Bacilli, such as RF39, RFN20 and ML615. Erysipelotrichales and two major gut lineages, RF39 and RFN20, were found to be neighboring clades of Mycoplasmatales. We detected habitat-specific functional patterns between the pathogenic, gut and the environmental Tenericutes, where genes involved in carbohydrate storage, carbon fixation, mutation repair, environmental response and amino acid cleavage are overrepresented in the genomes of environmental lineages. We hypothesize that the two major gut lineages, namely RF39 and RFN20, are probably acetate and hydrogen producers. Furthermore, deteriorating capacity of bactoprenol synthesis for cell wall peptidoglycan precursors secretion is a potential adaptive strategy employed by these lineages in response to the gut environment. This study uncovers the characteristic functions of environmental Tenericutes and their relationships with Bacilli, which sheds new light onto the pathogenicity and evolutionary processes of Mycoplasmatales.

**IMPORTANCE**

Environmental Tenericutes bacteria were recently discovered in numerous environments. However, our current collection of Tenericutes genomes was overrepresented by those for pathogens. Our phylogenomics study displays the relationships between all the available Tenericutes, as well as those between Tenericutes and the clades in Bacilli, which casts lights into the uncertain boundary between the environmental lineages of Tenericutes and Bacilli. By comparing the

56 genomes of the environmental and pathogenic Tenericutes, we revealed the metabolic

57 pathways and adaptive strategies of the Tenericutes in the different environments and

58 hosts. We also predicted the metabolism of the two major gut lineages RF39 and

59 RFN20 of Tenericutes, indicating their potential importance in stabilization of the gut

60 microbiome and contribution to human health.

61

62 **INTRODUCTION**

63 The phylum Tenericutes is composed of bacteria lacking a peptidoglycan cell wall.

64 The most well-studied clade belonging to this phylum is Mollicutes, which contains

65 medically relevant genera, including *Mycoplasma*, *Ureaplasma* and *Acholeplasma.*

66 All reported mollicutes are commensals or obligate parasites of humans, domestic

67 animals, plants and insects (1). Most studies so far have focused on pathogenic strains

68 in the Mycoplasmatales order (which encompasses the genera such as *Mycoplasma*,

69 *Ureaplasma*, *Mesoplasma* and *Spiroplasma*), resulting in their overrepresentation in

70 current genome databases. However, Tenericutes can also be found across a wide and

71 diverse range of environments. Recently, free-living *Izemoplasma* and *Haloplasma*

72 were reported in a deep-sea cold seep and brine pool, respectively (2, 3). Based on

73 their genomic features, the cell wall-lacking *Izemoplasma* were predicted to be

74 hydrogen producers and DNA degraders. The *Haloplasma contractile* genome

75 encodes actin and tubulin homologues, which might be required for its specific

76 motility in deep-sea hypersaline lake (4). These marine environmental Tenericutes

77 exhibit metabolic versatility and adaptive flexibility. This points out the unwanted

78 limitation that we must take into account at present when working on isolates of

79 marine Tenericutes representatives. The paucity of marine isolates currently available

80 has limited further mechanistic insights.

81

82 Environmental Tenericutes might be pathogens and/or mutualistic symbionts in the

83 gut of their host species. For example, mycoplasmas and hepatoplasmas affiliated

84 with Mycoplasmatales play a role in degrading recalcitrant carbon sources in the

85 stomach and pancreas of isopods (5, 6). *Spiroplasma* symbionts discovered in sea

3

cucumber guts possibly protect the host intestine from invading viruses (7). Tenericutes were also found in the intestinal tract of healthy shallow-water fish, mussels and 305 insect specimens (8-10). Recently, over 100 uncultured Tenericutes displaying high phylogenetic diversity were discovered in human gut metagenomes (11), irrespective of age and health status. It remains to be determined whether these novel lineages found in the human gut are linked to the maintenance of gut homeostasis and microbiome function. As a consequence of the host cell-associated lifestyle, the Tenericutes bacteria show extreme reduction in their genomes as well as reduced metabolic capacities, eliminating genes related to regulatory elements, biosynthesis of amino acids and intermediate metabolic compounds that must be imported from the host cytoplasm or tissue (12). Beyond genome reduction, evolution of pathogenic Mycoplasmatales species has also been accompanied by acquisition of new core metabolic and virulence factors (13, 14). Therefore, a comparison of the genetic profiles between environmental lineages and pathogens is needed to obtain insights into the adaptation of beneficial symbionts and the emergence of new diseases.

Since Tenericutes were recently reclassified into a Bacilli clade of Firmicutes (15), the discovery of environmental Tenericutes renovates the question regarding the boundary between Tenericutes and other clades of Bacilli. RF39 and RFN20 are two novel lineages of Bacilli, reported in the gut of the humans and domestic animals (16, 17). The environmental lineages of Bacilli and Tenericutes are expected to consist in close relatives but their genetic relationship has not been studied. This is important to address, as uncultured environmental Tenericutes and Bacilli may potentially emerge as pathogens. In this study, we compiled the genomes of 840 Tenericutes and determined their phylogenomic relationships with Bacilli. By analyzing the functional capacity encoded in these genomes, we deciphered the major differences in metabolic spectra and adaptive strategies between the major lineages of Tenericutes, including the two dominant gut lineages RF39 and RFN20.

116

## RESULTS AND DISCUSSION

**Phylogenetic tree of 16S rRNA genes and phylogenomics of Tenericutes**

119 We retrieved all available Tenericutes genomes from the NCBI database (April, 2019).

120 A total of 840 genomes with ≥50% completeness and ≤10% contamination by foreign

121 DNA were selected (Supplementary file 1). From these, 685 16S rRNA genes were

122 extracted and clustered together when displaying at >99% identity, resulting in 227

123 representative sequences. Approximately 70% of the non-redundant sequences were

124 derived from the order Mycoplasmatales (highly represented by the hominis group),

125 which was largely composed of pathogens isolated from plants, humans and animals.

126 Together with 33 reference sequences from marine samples, a total of 260 16S rRNA

127 genes were used to build a maximum-likelihood (ML) tree. Using *Bacillus subtilis* as

128 an outgroup, Tenericutes 16S rRNA sequences were divided into several clades (Fig.

129 1A). *Acholeplasma* and *Phytoplasma* were grouped into one clade, while

130 *Izemoplasma* and *Haloplasma* were closer to the basal group. Tenericutes species

131 were detected across a range of environments, including mud, bioreactors, hypersaline

132 lake sediment, and ground water. The non-human hosts of Tenericutes included

133 marine animals, domestic animals and fungi. Sequences isolated from fungi and

134 mycoplasma-infected animal blood samples were associated with longer branches,

135 indicating the occurrence of a niche-specific evolution. *Hepatoplasma* identified as a

136 novel genus in Mycoplasmatales is also exclusively present in the gut microbiome of

137 amphipods and isopods (5, 18). *Spiroplasma* detected in a sea cucumber gut has been

138 described as a mutualistic endosymbiont (7), rather than a pathogen. These isolates

139 from environmental hosts were distantly related to others in the tree, indicating a high

140 diversity of Mycoplasmatales across a wide range of hosts and their essential role in

141 adaptation and health of marine invertebrates. Analyses of 135 16S rRNA amplicon

142 datasets and 141 Tara Ocean metagenomes (19) from marine waters revealed the

143 presence of mycoplasmas from the hominis group and other sequences from the basal

144 groups of the tree in more than 21.7% of the samples. Four of the five representative

145 16S rRNA sequences from the hominis group were similar (95.9%-99.3%) to that of

5

146    halophilic *Mycoplasma todarodis* isolated from squids collected near an Atlantic

147    island. The finding of the Tenericutes isolated from humans and other animal hosts in

148    the marine samples indicates that they may be spreading possibly through sewage.

149    The relative abundance of the twelve representative 16S rRNA genes from the marine

150    waters was extremely low (<0.1%) in the microbial communities of the oceans.

151    However, considering the tremendous body of marine water, the oceans harbor a

152    massive Tenericutes population composed of undetected novel lineages. We detected

153    two major clades of human gut lineages (hereafter referred to as HG1 and HG2) that

154    were placed between Mycoplasmatales and Acholeplasmatales (Fig. 1A). These two

155    lineages have been revealed recently as encompassing many previously unknown

156    species in the human gut (11). However, their contribution to human health and the

157    core gut microbiome stability remains unclear.

158

159    A phylogenomics analysis of Tenericutes was performed using concatenated

160    conserved proteins from 840 Tenericutes genomes and three Firmicutes genomes.

161    Interestingly, the topology of the phylogenomic tree coincides with that of the

162    phylogenetic tree based on 16S rRNA genes. However, 67.6% of the genomes were

163    derived from Mycoplasmatales, indicating a strong bias of Tenericutes genomes

164    towards pathogens and disease-inducing isolates. The human gut lineages HG1 ($n$=87)

165    and HG2 ($n$=21) were found to be neighboring clades of Mycoplasmatales as well

166    (Fig. 1B). The genetic distance between the genomes of the gut lineages was much

167    higher than that between the species in Mycoplasmatales, except for those in

168    mycoplasma-infected blood and fungi. *Acholeplasma* and *Phytoplasma* were within a

169    clade composed of uncultured environmental Tenericutes lineages from ground waters,

170    hypersaline sediments and mud, suggesting an environmental origin for the two

171    genera.

172

173    By calculating the relative evolutionary divergence (RED) of the genomes of several

174    Tenericutes lineages (15), the average RED values for HG1 and HG2 were 0.94±0.03

175    and 0.91±0.07, respectively. Considering an expected RED value of 0.92 at the genus

6

176   level, these two lineages can be considered new genera in Tenericutes. The RED value

177   for the sequences from hypersaline lake sediments was 0.70, which supports the

178   presence of a new order or family in Tenericutes.

179

**Phylogenomic position of Tenericutes in Bacilli**

181   Tenericutes were recently integrated into the Bacilli clade within the Firmicutes

182   phylum (15). To examine the phylogenetic positions of the new Tenericutes lineages

183   and Bacilli, we used representative genomes of the orders within Bacilli and those in

184   Tenericutes available on NCBI. The topology of the phylogenomic relationships was

185   supported by two ML methods. In the phylogenomic tree, four Bacilli orders, namely

186   Staphylococcales, Exiguobacterales, Bacillales, and Lactobacillales, were clearly split

187   from those of Tenericutes. Newly defined orders RF39, RFN20 and ML615 in Bacilli

188   clustered with HG1, HG2, and uncultured Tenericutes from bioreactors, respectively.

189   This suggests that most of uncultured environmental Tenericutes are probably novel

190   Bacilli orders, and that the boundary between Tenericutes and Bacilli is uncertain.

191   RF39, RFN20 and ML615 were also affiliated with Tenericutes if the boundary of

192   Tenericutes on the tree was set at Haloplasmatales. Although RF39 and RFN20 are

193   part of the HG1 and HG2 lineages, they have also been detected in domestic animals

194   (20). Interestingly, the Erysipelotrichales order was phylogenetically placed between

195   both human gut lineages (Fig. 2). Since all Erysipelotrichales species described in the

196   literature so far possess a cell wall (21), their phylogenomic affinity to cell

197   wall-lacking Tenericutes is unexpected.

198

199   We investigated the genome structure of Tenericutes and Erysipelotrichales species by

200   calculating genome completeness, size and GC content (Fig. S1). Most of the

201   high-quality genomes (>90% completeness and <5% contamination) were assigned to

202   Mycoplasmatales and Acholeplasmatales. In contrast to the rather stable genomes of

203   the pathogenic species, the genome sizes of the uncultured Tenericutes species

204   differed from each other and almost all were smaller than 2 Mb. Haloplasmatales

205   genomes were the largest on average. Most of the Tenericutes genomes have a low

7

206    GC content (<30%), whereas the average GC content of those from a hypersaline lake

207    was about 50%, consistent with a selection pressure exerted by ionic strength on the

208    DNA double helix (22, 23). Notably, GC contents calculated on 1 kb intervals in

209    Tenericutes genomes from ground water and HG1 (specifically RF39) varied from 20%

210    to 70%, suggesting great plasticity and frequent gene transfers. However, these results

211    were dependent on the number of genomes considered from different sources and may

212    be influenced by the quality of genome binning.

213

214    **Genomic and functional divergence between environmental Tenericutes and**

215    **pathogens**

216

217    Erysipelotrichales and Tenericutes genomes were functionally annotated to

218    characterize their metabolic pathways and stress responses that might determine the

219    versatility and niche-specific evolution of different orders and lineages in Tenericutes.

220    The annotation results against the Kyoto Encyclopedia of Genes and Genomes

221    (KEGG) (24) and the clusters of orthologous groups (COGs) databases were used to

222    calculate the percentages of the genes in the genomes (supplementary file 2). Based

223    on the frequency of all the COGs, Erysipelotrichales and Tenericutes were split into

224    two major agglomerative hierarchical clustering (AHC) clusters. Mycoplasmatales

225    and *Phytoplasma* formed AHC cluster 1, while the remaining formed cluster 2.

226

227    Using Mann-Whitney test, 203 KEGG genes and 420 COGs showed a significant

228    difference ($p < 0.01$) in frequency between the two AHC clusters (supplementary file 2).

229    We selected 62 of the genes to represent those for 16 functional categories that were

230    distinct in environmental adaptation and carbon metabolism between the two clusters

231    (Table S1 and Fig. 3). Sugars such as xylose, galactose and fructose might be

232    fermented to L-lactate, formate and acetate by Tenericutes. The sugar sources and

233    fermentation products differed between the groups (Fig. 3). Phosphotransferase (PTS)

234    systems responsible for sugar cross-membrane transport were encoded by most of the

235    genomes    of    *Spiroplasma*,    *Mesoplasma*,    *Entomoplasma*,    Haloplasmatales,

8

236    Erysipelotrichales, mycoides, and pneumoniae groups. Although most of the

237    environmental Tenericutes genomes did not maintain PTS systems, sugar uptake

238    might be carried out by ABC transporters. Almost all of the Tenericutes groups in the

239    AHC cluster 2 (containing all the environmental lineages) were found to encode genes

240    involved in starch synthesis (*glgABP*) and carbon storage, except for HG1. These

241    Tenericutes groups also encoded the pullulanase gene PulA involved in starch

242    degradation. Autotrophic pathways were present almost exclusively in environmental

243    Tenericutes genomes. $CO_2$ is fixed by two autotrophic steps mediated by the citrate

244    lyase genes that function in reductive citric acid cycle (rTCA) and the

245    2-oxoglutarate/2-oxoacid ferredoxin oxidoreductase genes (*korABCD*) that encode

246    enzymes for reductive acetyl-CoA pathway. The resulting pyruvate might be further

247    stored as glucose and glycan via reversible Embden–Meyerhof–Parnas (EMP)

248    pathway. PPDK is the key enzyme that controls the interconversion of

249    phosphoenolpyruvate and pyruvate in prokaryotes (25). Among all the environmental

250    lineages and Erysipelotrichales, *ppdK* gene was frequently identified (73.8%-100%)

251    except for Haloplasmatales and Acholeplasmatales.

252

253    Aromatic biosynthesis pathway was lost in Mycoplasmatales, indicating their

254    complete dependence on hosts for aromatic amino acids. Acquisition of amino acids

255    by some environmental Tenericutes was likely conducted by peptidases (*pepD2*) and

256    cross-membrane oligopeptide transporters. Glycine was also probably an important

257    carbon and nitrogen source for the environmental Tenericutes, as a high percentage of

258    their genomes (76.3%-100%) contained the glycine cleavage genes *gcvT* and *gcvH*.

259

260    Glycerol is a key intermediate between sugar and lipid metabolisms and is imported

261    by a facilitation factor GlpF. Phosphorylation of glycerol by a glycerol kinase (GK) is

262    followed by oxidation to dihydroxyacetone phosphate (DHAP) by

263    glycerol-3-phosphate (G3P) dehydrogenase (GlpD), which is further metabolized in

264    the glycolysis pathway (26). More than 95% of the genomes of *Mesoplasma*,

265    pneumoniae, mycoides and wastewater groups contained the *glpD* gene; in contrast,

266   *Phytoplasma* and *Ureaplasma* genomes lacked a *glpD* gene. 62% of RFN20 genomes

267   harbored the *glpD* gene, while it was only found in 2% of RF39. RF39 genomes also

268   lacked the GK-encoding gene, which suggests that RF39 cannot utilize glycerol from

269   diet or the gut membrane. Hydrogen peroxide ($H_2O_2$) is a by-product of G3P

270   oxidation, and has deleterious effects on epithelial surfaces in humans and animals

271   (27). On the other hand, these $H_2O_2$ catabolism genes were more frequently identified

272   in uncultured environmental Tenericutes (Fig. 3).

273

274   The DNA mismatch repair machinery components MutS and MutL were almost

275   entirely absent from Mycoplasmatales and *Phytoplasma* genomes. RFN20 genomes

276   also had a low percentage of the DNA repairing genes (33.3% for *mutS* and 57.1% for

277   *mutL)*. This lack of DNA repairing genes might have generated more mutants in small

278   asexual microbial populations capable of adapting to new environments due to

279   Muller's ratchet effect (28).

280

281   In *Mycoplasma* species as in mitochondria, tRNA anticodon base U34 can pair with

282   any of the four bases in codon family boxes (29). To makes this ability more efficient

283   U34 is modified in some organisms by enzymes using a carboxylated

284   S-adenosylmethionine. The SmtA enzyme, also known as CmoM, is a

285   methyltransferase that adds a further methyl group to U34 modified tRNA for precise

286   decoding of mRNA and rapid growth (30, 31). The high frequency of *smtA* gene in the

287   environmental Tenericutes genomes indicates a capacity to regulate their growth

288   under various conditions. OmpR is a two-component regulator tightly associated with

289   a histidine kinase/phosphatase EnvZ for regulatory response to environmental

290   osmolarity changes(32). Its presence in most of the environmental Tenericutes

291   genomes (>70.4%) suggests its involvement in regulating stress responses in these

292   organisms. The genomes of two gut lineages RFN20 and RF39 also contained a high

293   percentage of the *ompR* gene. In contrast, almost all Mycoplasmatales and

294   *Phytoplasma* genomes lacked the *ompR* gene.

295

296    The cell division/cell wall cluster transcriptional repressor MraZ can negatively

297    regulate cell division of Tenericutes (33). The *mraZ* gene that is thus responsible for

298    dormancy of bacteria is conserved in *Erysipelotrichales* and *Mycoplasmatales*.

299    Further studies are needed to examine whether this gene can be targeted to control

300    pathogenicity of the bacteria in the two orders.

301

302    The Rnf proton pump system evolved in anoxic condition and is employed by

303    anaerobes to generate proton gradients for energy conservation (34). In

304    single-membrane Tenericutes, proton gradients can hardly be established by the Rnf

305    system due to the leakage of protons directly to the environment. However, this

306    system was well preserved in genomes from Izemoplasmatales and the wastewater

307    group. The Rnf system in these species was likely used for pumping protons out of the

308    cell to balance cytoplasmic pH.

309

310    **Metabolic model of gut lineages RFN20 and RF39**

311    A recent study reported the genome features of RFN20 and RF39, the two main clades

312    comprising uncultured Tenericutes (16). The major findings on these two lineages

313    were their small genomes and the lack of several amino acid biosynthesis pathways.

314    After correction for genome completeness in this study, we found that the RF39

315    genomes were indeed significantly smaller than those of RFN20 genomes (t-test;

316    *p*=0.0012). We selected four nearly complete genomes of RFN20 and RF39 for

317    annotation and elaborated their metabolic potentials (Table 1). The genome sizes were

318    between 1.5 Mb-1.9 Mb, smaller than those from *Sharpea azabuensis* belonging to

319    the order Erysipelotrichales. We built a schematic metabolic map for the

320    representative RFN20 and RF39 species on the basis of the KEGG and COG

321    annotation results. The two lineages were predicted to be acetogens since the four

322    genomes encoded genes for acetate production (Fig. 4). We hypothesize that sugars

323    are imported from the environment by ABC sugar transporters, while autotrophic $CO_2$

324    fixation might occur via carboxylation of acetyl-CoA to pyruvate by the

325    pyruvate:ferredoxin oxidoreductase (PFOR). Glycerol is imported and enters

11

326 glycerophospholipid metabolism, which results in cardiolipin biosynthesis instead of

327 fermentation through the EMP pathway. In some pathogenic mycoplasmas, glycerol

328 can be taken into central carbon metabolism (26), as mentioned above.

329

330 RFN20 and RF39 are probably mixotrophic since $CO_2$ can be fixed to pyruvate and

331 stored as starch, while central carbon metabolism is also connected with amino acid

332 metabolism. After uptake of oligopeptides by the App ABC transporter system, an

333 endo-oligopeptidase encoded by *pepF* yields amino acids for protein synthesis.

334 Glycine and serine might feed into pyruvate metabolism. The peptidoglycan

335 biosynthesis pathway was found to be complete in all four RFN20 and RF39 genomes

336 here considered, but two genomes, namely HG1.1 and HG2.1 (Table 1), lacked the

337 genes encoding the enzymes for UDP-N-acetylglucosamine (UDP-NAG) synthesis.

338 Instead, these genomes harbored all the genes required for the subsequent synthesis

339 steps to generate extracellular peptidoglycan. *murG* and *mraY* genes, which are

340 involved in integration of UDP-NAG and UDP-N-acetylmuramate (UDP-NAM) into

341 the peptidoglycan unit, respectively, were identified in the four genomes. With the

342 addition of an oligopeptide, the peptidoglycan unit is secreted into the cell surface

343 with the assistance of bactoprenol (C55 isoprenoid alcohol) (35, 36), which is formed

344 by condensation of eight isopentenyl-diphosphate (IPP) units and one

345 farnesyl-diphosphate (FPP). The *uppS* gene responsible for the bactoprenol formation

346 was identified in the four RFN20 and RF39 genomes (37). In bacteria, IPP can be

347 synthesized by several metabolic steps. All the genomes contained the genes encoding

348 the respective enzymes involved in the intermediate steps of IPP and dimethylallyl

349 diphosphate (DPP) synthesis through MEP/DOXP pathway, except for *ispD* gene in

350 one genome (Fig. 4). However, the polyprenyl synthetase gene (*ispA*), which is

351 essential in the formation of FPP, was missing in three of the genomes. Given the loss

352 of the *ispA* gene, the source of FPP for bactoprenol synthesis is unclear. Overall, 86.9%

353 and 14.3% of the RF39 and RFN20 genomes contained the *mraY* gene, respectively,

354 while 68.7% and 5.2% of the RF39 and RFN20 genomes had the *murG* gene,

355 respectively. Therefore, most of the RFN20 genomes collected in this study lacked the

356    complete pathway for peptidoglycan synthesis. The two essential genes for

357    peptidoglycan synthesis were only frequently detected in Tenericutes genomes from

358    the bioreactor group (75.0% for both genes) and Erysipelotrichales genomes (80.0%

359    and 60.0% for *mraY* and *murG*, respectively). Therefore, the capacity of

360    peptidoglycan synthesis is possibly deteriorating in the gut lineages, as a potential

361    adaptive strategy to the gut environment. Similarly, the *H. contractile* was reported to

362    possess the peptidoglycan synthesis genes in its genome (4), although it also lacks a

363    cell wall. Our further examination of the genome found that the *murEF* genes

364    involved in extending the oligopeptide attached on UDP-NAM were absent. Hence,

365    the synthesis of aminosugars NAG and NAM probably served as a mechanism of

366    carbon and nitrogen storage for *H. contractile*.

367

368    RFN20 and RF39 are probably hydrogen producers, as the four genomes of HG1 and

369    HG2 had [FeFe]-hydrogenase encoding genes. All the genomes carried the *feo* and *fhu*

370    genes for ferrous iron uptake. Ferrous irons are taken by ABC transporters Feo into

371    the cells when ferrous iron concentration is high in the environment. The Fhu receptor

372    for ferrichrome absorption is required in iron-limiting condition such as the human

373    gut (38). The oxygen-sensitive [FeFe]-hydrogenases contain 4Fe-4S cluster and an

374    H-cluster consisting of several conserved catalytic motifs involved in hydrogen

375    production. Three distinct binding motifs of H-cluster in [FeFe]-hydrogenases,

376    TSCXP, $PCX_2KX_2E$ and $EXMXCXGGCX_2C$ (39), were present in the five

377    hydrogenases encoded by all the four genomes (Fig. S2). However, three of the

378    hydrogenases from HG1 and HG2 harbor specific sites that differ from the others in

379    some of the active sites. We have identified several orthologs with these distinct

380    amino acids in the conserved motifs. These [FeFe]-hydrogenases formed a novel

381    cluster in the phylogenetic tree. HG2.1 genome harbored two copies of the

382    [FeFe]-hydrogenase genes, which were diversified as shown by their positions in the

383    phylogenetic tree and the differences in conserved catalytic sites (Fig. S2). In the

384    human gut, three groups of [FeFe]-hydrogenases have been detected, and were

385    proposed to be involved in methanogenesis, acetogenesis and sulfate reduction (40).

386  Lignocellulose-feeding termites also produce a high concentration of hydrogen in
387  their guts, probably for degradation of wood (41). Therefore, the HG1 and HG2 gut
388  lineages are probably important for maintenance of a healthy gut microbial ecosystem
389  and degradation of recalcitrant carbon.

390

391  As indicated by the phylogenomics tree, there is a high genomic variation within the
392  RFN20 and RF39 lineages. Therefore, the predicted lifestyle of RFN20 and RF39
393  may vary among human populations. For example, 68.7% and 76.2% of RF39 and
394  RFN20 genomes, respectively, harbored the *uppS* gene for bactoprenol synthesis.
395  However, the lack of high-quality, isolate genomes representing these lineages hinders
396  the evaluation of their dynamics and evolutionary processes in the human gut.

397

398  In this study, the genomic features of RFN20 and RF39 were shown to be highly
399  dynamic among genomes from different sources. RF39 genomes lacked most of the
400  genes for carbohydrate storage but maintained *mutSL* genes involved in DNA repair
401  (Fig. 3). Except for this, there were no major differences between the two lineages,
402  although a previous study claimed that RF39 were prone to be autotrophic (16). In
403  deep-sea isopod gut, we also identified two types of Tenericutes bacteria, *Mycoplasma*
404  sp. Bg1 and Bg2 (6). M. sp. Bg1 was able to degrade sialic acids probably by
405  attachment to the host gut surface. The co-existence of two Tenericutes lineages in
406  human and animal intestinal tracts is still enigmatic and warrants further
407  investigations using microscopy and transcriptomics methods.

408

409  In conclusion, our study revealed phylogenetic diversity of the Tenericutes groups and
410  their phylogenomic relationships with Bacilli. In the environmental groups of
411  Tenericutes, we uncovered novel lineages in human guts and marine environments,
412  indicating the lack of environmental representatives for studies on their adaptive
413  strategies and pathogenicity. Our finding of the gut lineages and their metabolic
414  characteristics casts lights into unknown diversified mutualistic Tenericutes in gut
415  microbiome.

14

416

## MATERIAL AND METHODS

### Genome collection and quality check

419    A total of 857 Tenericutes genomes were downloaded from the NCBI database. Three

420    genomes of deep-sea symbiotic Tenericutes were collected from the previous studies

421    (6, 7). Completeness and contamination of the genomes were evaluated by CheckM

422    (v1.0.5) (42). Those with >10% contaminants and <50% complete were removed. To

423    explore variations of GC content in these genomes, GC content within 1-kb genome

424    intervals were calculated. 16S rRNA genes were identified from these genomes using

425    rRNA_HMM with default settings (43), and only those longer than 300 bp were

426    extracted. If there was more than one 16S rRNA gene in a genome, the longest one

427    was selected. The sequences were grouped with an identity cutoff of 99% using

428    CD-HIT (44) and only the longest was retained as the representative. From each order

429    of Bacilli, five genomes (see supplementary file 1) were obtained from the Genome

430    Tree Database (GTDB) (15). They were selected from different families if possible.

431

### Genome annotation and comparison

433    The protein coding sequences in the genomes were predicted by Prodigal (v2.6.2) (45)

434    (proteins from Tenericutes in particular were predicted with parameter –g 4). The

435    proteins were then searched against the eggNOG database by eggNOG-mapper (v2)

436    (46) (with parameters --seed_orthorlog_evalue 1e-10), KEGG   (24) and COGs (47)

437    databases by Blastp with E-value cutoff of 1e-05 and similarity threshold of 40%. The

438    functions of essential COGs belonging to Tenericutes were referred to those for a

439    synthetic bacterium JCVI-Syn3.0 with a minimal genome (48).

440

441    The collected Tenericutes genomes were grouped by taxonomy and source

442    (supplementary file 1). The percentage of the KEGG genes and COGs in the genomes

443    of each group was calculated. This was also accomplished for *Erysipelotrichales*

444    genomes. To filter low-frequency genes, at least one of the groups had a target gene

445    in >30% of the genomes. The percentages of the genes used for Bray-Curtis

15

446    dissimilarity estimates were calculated using the COG frequency table. AHC analysis

447    was conducted using the pairwise dissimilarities between groups. A Mann-Whitney

448    test was performed using the percentages of COGs and KEGG genes between the

449    AHC clusters. The KEGG genes with *p* value <0.01 were clustered into functional

450    modules on the KEGG website (www.kegg.jp).

451

452    **Phylogenetic and phylogenomic analyses**

453    The analyses on the datasets of 16S rRNA gene amplicons from marine samples were

454    described in our previous study (49). The representative reads of Tenericutes OTUs

455    were recruited for this study. Raw metagenomic data from Tara Ocean project were

456    checked by FastQC (version 0.11.4). Reads with low quality bases (PHRED quality

457    score < 20 over 70% of the reads) were removed using the NGS QC Toolkit (50). The

458    quality-filtered reads were merged using PEAR (v0.9.5) (51) and those 16S rRNA

459    fragments >140 bp were identified and extracted with rRNA_HMM (43). After

460    taxonomic classification of the fragments using the Ribosomal Database Project (RDP)

461    classifier version 2.2 against the SILVA 128 database (52, 53), those belonging to

462    Tenericutes were collected for the following phylogenetic study.

463

464    The 16S rRNA genes from the genomes, the amplicons and the Tara project were first

465    clustered by MUSCLE (v3.8) (54) and then trimmed by trimAl v1.4 (automated1)

466    (55). The ML phylogenetic tree of 16S rRNA genes was built by IQ-TREE (v1.6.10)

467    (56, 57) (with parameters -m GTR+F+R10 -alrt 1000 -bb 1000). Conserved proteins

468    of the Tenericutes genomes were identified by AMPHORA2 (58). A total of 31

469    conserved proteins were used to construct the phylogenomic tree for Tenericutes. The

470    conserved proteins were aligned with MUSCLE (v3.8)(54), concatenated and then

471    trimmed with trimAl (v1.4) (automated1) (55). The conserved proteins from

472    *Syntrophomonas wolfei* (NC_008346), *Thermacetogenium phaeum* (NC_018870) and

473    *Desulfallas geothermicus* (NZ_FOYM01000001) were combined with the dataset of

474    Tenericutes as an outgroup. The phylogenomics tree for Tenericutes was built by

475    IQ-TREE (v1.6.10) (56, 57) (with parameters -m LG+F+R10 -alrt 1000 -bb 1000).

16

476  The phylogenomic tree for Bacilli and Tenericutes was constructed first with

477  IQ-TREE (v1.6.10) using the same settings as that for the phylogenomics tree of

478  Tenericutes and then with RAxML 8.1.21 using PROTGAMMA+BLOSUM62 model

479  with 100 bootstrap replicates.

480

481  **Prediction of metabolic models of RFN20 and RF39**

482  Four genomes were selected from the downloaded genomes of Tenericutes to

483  represent RFN20 and RF39 with respect to their high genome completeness. The

484  protein-coding sequences were predicted by Prodigal (v2.6.2) (45) as mentioned

485  above. The proteins were then searched against COG database (47) by Blastp (59)

486  with an E-value cutoff of 1e-05. KEGG annotation was conducted using the online

487  BlastKOALA tool (24).

488

489  **ACKNOWLEDGEMENTS:**

494  Y.W., A.D. and L.S.H. designed the study; Y.W., J.M.H., and Y.L.Z. performed the

495  bulk of the phylogenomic analyses; A.A. and R.D.F. contributed data for analysis;

496  Y.W. wrote the manuscript. All of us contributed to manuscript revisions.

497  The authors declare that there is no conflict of interest.

498  **REFERENCES:**

499  1.  **Razin S, Herrmann R.** 2002. Molecular biology and pathogenicity of mycoplasmas. Springer,
500      Boston, MA.
501  2.  **Antunes A, Rainey FA, Wanner G, Taborda M, Pätzold J, Nobre MF, da Costa MS,**
502      **Huber R.** 2008. A new lineage of halophilic, wall-less, contractile bacteria from a brine-filled
503      deep of the Red Sea. J Bacteriol **190:**3580-3587.
504  3.  **Skennerton CT, Haroon MF, Briegel A, Shi J, Jensen GJ, Tyson GW, Orphan VJ.** 2016.
505      Phylogenomic analysis of Candidatus 'Izimaplasma' species: free-living representatives from
506      a Tenericutes clade found in methane seeps. ISME J **10:**2679-2692.
507  4.  **Antunes A, Alam I, El Dorry H, Siam R, Robertson A, Bajic VB, Stingl U.** 2011. Genome

17

508          sequence of Haloplasma contractile, an unusual contractile bacterium from a deep-sea anoxic
509          brine lake. J Bacteriol **193:**4551-4552.

510    5.    **Wang YJ, Stingl U, Anton-Erxleben F, Geisler S, Brune A, Zimmer M.** 2004. "*Candidatus*
511          Hepatoplasma crinochetorum," a new, stalk-forming lineage of Mollicutes colonizing the
512          midgut glands of a terrestrial isopod. Appl Environ Microb **70:**6166-6172.

513    6.    **Wang Y, Huang JM, Wang SL, Gao ZM, Zhang AQ, Danchin A, He LS.** 2016. Genomic
514          characterization of symbiotic mycoplasmas from the stomach of deep-sea isopod bathynomus
515          sp. Environ Microbiol **18:**2646-2659.

516    7.    **He L-S, Zhang P-W, Huang J-M, Zhu F-C, Danchin A, Wang Y.** 2018. The enigmatic
517          genome of an obligate ancient Spiroplasma symbiont in a hadal holothurian. Appl   Environ
518          Microbiol **84:**e01965-01917.

519    8.    **Sullam KE, Essinger SD, Lozupone CA, O'Connor MP, Rosen GL, Knight R, Kilham SS,**
520          **Russell JA.** 2012. Environmental and ecological factors that shape the gut bacterial
521          communities of fish: a meta-analysis. Mol Ecol **21:**3363-3378.

522    9.    **Yun JH, Roh SW, Whon TW, Jung MJ, Kim MS, Park DS, Yoon C, Nam YD, Kim YJ,**
523          **Choi JH, Kim JY, Shin NR, Kim SH, Lee WJ, Bae JW.** 2014. Insect gut bacterial diversity
524          determined by environmental habitat, diet, developmental stage, and phylogeny of host. Appl
525          Environ Microbiol **80:**5254-5264.

526   10.    **Aceves AK, Johnson P, Bullard SA, Lafrentz S, Arias CR.** 2018. Description and
527          characterization of the digestive gland microbiome in the freshwater mussel Villosa nebulosa
528          (Bivalvia: Unionidae). J Molluscan Studies **84:**240-246.

529   11.    **Almeida A, Mitchell AL, Boland M, Forster SC, Gloor GB, Tarkowska A, Lawley TD,**
530          **Finn RD.** 2019. A new genomic blueprint of the human gut microbiota. Nature **568:**499-504.

531   12.    **Moran NA.** 2002. Microbial minimalism: Genome reduction in bacterial pathogens. Cell
532          **108:**583-586.

533   13.    **Lo W-S, Gasparich GE, Kuo C-H.** 2018. Convergent evolution among ruminant-pathogenic
534          mycoplasma involved extensive gene content changes. Genome Biol Evol **10:**2130-2139.

535   14.    **Chernov VM, Chernova OA, Mouzykantov AA, Medvedeva ES, Baranova NB, Malygina**
536          **TY, Aminov RI, Trushin MV.** 2018. Antimicrobial resistance in mollicutes: known and
537          newly emerging mechanisms. FEMS Microbiol Lett **365**.

538   15.    **Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil PA,**
539          **Hugenholtz P.** 2018. A standardized bacterial taxonomy based on genome phylogeny
540          substantially revises the tree of life. Nature Biotechnol **36:**996-1004.

541   16.    **Nayfach S, Shi ZJ, Seshadri R, Pollard KS, Kyrpides NC.** 2019. New insights from
542          uncultivated genomes of the global human gut microbiome. Nature **568:**505-510.

543   17.    **Zhang LT, Huang XF, Xue B, Peng QH, Wang ZS, Yan TH, Wang LZ.** 2015.
544          Immunization against rumen methanogenesis by vaccination with a new recombinant protein.
545          PLoS ONE **10:**e0140086.

546   18.    **Cheng X-Y, Wang Y, Li J-Y, Yan G-Y, He L-S.** 2019. Comparative analysis of the gut
547          microbial communities between two dominant amphipods from the Challenger Deep, Mariana
548          Trench. Deep Sea Res I **151:**103081.

549   19.    **Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, Djahanschiri B,**
550          **Zeller G, Mende DR, Alberti A, Cornejo-Castillo FM, Costea PI, Cruaud C, d'Ovidio F,**
551          **Engelen S, Ferrera I, Gasol JM, Guidi L, Hildebrand F, Kokoszka F, Lepoivre C,**

552    **Lima-Mendez G, Poulain J, Poulos BT, Royo-Llonch M, Sarmento H, Vieira-Silva S,**
553    **Dimier C, Picheral M, Searson S, Kandels-Lewis S, Bowler C, de Vargas C, Gorsky G,**
554    **Grimsley N, Hingamp P, Iudicone D, Jaillon O, Not F, Ogata H, Pesant S, Speich S,**
555    **Stemmann L, Sullivan MB, Weissenbach J, Wincker P, Karsenti E, Raes J, Acinas SG,**
556    **Bork P.** 2015. Structure and function of the global ocean microbiome. Science **348:**1261359.

557  20.  **Pitta DW, Parmar N, Patel AK, Indugu N, Kumar S, Prajapathi KB, Patel AB, Reddy B,**
558    **Joshi C.** 2014. Bacterial diversity dynamics associated with different diets and different
559    primer pairs in the rumen of kankrej cattle. PLoS ONE **9:**e111710.

560  21.  **Shimoji Y, Yokomizo Y, Sekizaki T, Mori Y, Kubo M.** 1994. Presence of a capsule in
561    Erysipelothrix-Rhusiopathiae and its relationship to virulence for mice. Infect Imm
562    **62:**2806-2810.

563  22.  **Soppa J.** 2006. From genomes to function: haloarchaea as model organisms. Microbiology
564    **152:**585-590.

565  23.  **Lyubchenko YL, Shlyakhtenko LS.** 1997. Visualization of supercoiled DNA with atomic
566    force microscopy in situ. Proc Natl Acad Sci U S A **94:**496-501.

567  24.  **Kanehisa M, Goto S.** 2000. KEGG: kyoto encyclopedia of genes and genomes. Nucl Acids
568    Res **28:**27-30.

569  25.  **Tjaden B, Plagens A, Dorr C, Siebers B, Hensel R.** 2006. Phosphoenolpyruvate synthetase
570    and pyruvate, phosphate dikinase of *Thermoproteus tenax*: key pieces in the puzzle of archaeal
571    carbohydrate metabolism. Mol Microbiol **60:**287-298.

572  26.  **Yeh JI, Chinte U, Du S.** 2008. Structure of glycerol-3-phosphate dehydrogenase, an essential
573    monotopic membrane enzyme involved in respiration and metabolism. Proc Natl Acad Sci U S
574    A **105:**3280-3285.

575  27.  **Blotz C, Stulke J.** 2017. Glycerol metabolism and its implication in virulence in *Mycoplasma*.
576    FEMS Microbiol Rev **41:**640-652.

577  28.  **Andersson DI, Hughes D.** 1996. Muller's ratchet decreases fitness of a DNA-based microbe.
578    Proc Natl Acad Sci U S A **93:**906-907.

579  29.  **Grosjean H, Westhof E.** 2016. An integrated, structure- and energy-based view of the genetic
580    code. Nucl Acids Res **44:**8020-8040.

581  30.  **Sakai Y, Miyauchi K, Kimura S, Suzuki T.** 2016. Biogenesis and growth phase-dependent
582    alteration of 5-methoxycarbonylmethoxyuridine in tRNA anticodons. Nucl Acids Res
583    **44:**509-523.

584  31.  **Yamanaka K, Ogura T, Niki H, Hiraga S.** 1995. Characterization of the smtA gene encoding
585    an S-adenosylmethionine-dependent methyltransferase of Escherichia coli. FEMS Microbiol
586    Lett **133:**59-63.

587  32.  **Cai SJ, Inouye M.** 2002. EnvZ-OmpR interaction and osmoregulation in Escherichia coli. J
588    Biol Chem **277:**24155-24161.

589  33.  **Eraso JM, Markillie LM, Mitchell HD, Taylor RC, Orr G, Margolin W.** 2014. The highly
590    conserved MraZ protein is a transcriptional regulator in Escherichia coli. J Bacteriol
591    **196:**2053-2066.

592  34.  **Schuchmann K, Muller V.** 2014. Autotrophy at the thermodynamic limit of life: a model for
593    energy conservation in acetogenic bacteria. Nat Rev Microbiol **12:**809-821.

594  35.  **Thorne KJ, Kodicek E.** 1966. The structure of bactoprenol, a lipid formed by lactobacilli
595    from mevalonic acid. Biochem J **99:**123-127.

596   36.   **Manat G, Roure S, Auger R, Bouhss A, Barreteau H, Mengin-Lecreulx D, Touzé T.** 2014.
597         Deciphering the metabolism of undecaprenyl-phosphate: the bacterial cell-wall unit carrier at
598         the membrane frontier. Microb Drug Ris **20:**199-214.

599   37.   **Mostafavi AZ, Lujan DK, Erickson KM, Martinez CD, Troutman JM.** 2013. Fluorescent
600         probes for investigation of isoprenoid configuration and size discrimination by
601         bactoprenol-utilizing enzymes. Bioorganic Med Chem **21:**5428-5435.

602   38.   **Wooldridge KG, Williams PH.** 1993. Iron uptake mechanisms of pathogenic bacteria. FEMS
603         Microbiol Rev **12:**325-348.

604   39.   **Mulder David W, Shepard Eric M, Meuser Jonathan E, Joshi N, King Paul W, Posewitz**
605         **Matthew C, Broderick Joan B, Peters John W.** 2011. Insights into [FeFe]-hydrogenase
606         structure, mechanism, and maturation. Structure **19:**1038-1052.

607   40.   **Wolf PG, Biswas A, Morales SE, Greening C, Gaskins HR.** 2016. H-2 metabolism is
608         widespread and diverse among human colonic microbes. Gut Microbes **7:**235-245.

609   41.   **Ballor NR, Leadbetter JR.** 2012. Patterns of [FeFe] hydrogenase diversity in the gut
610         microbial communities of lignocellulose-feeding higher termites. Appl Environ Microb
611         **78:**5368-5374.

612   42.   **Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW.** 2015. CheckM:
613         assessing the quality of microbial genomes recovered from isolates, single cells, and
614         metagenomes. Genome Res **25:**1043-1055.

615   43.   **Huang Y, Gilna P, Li W.** 2009. Identification of ribosomal RNA genes in metagenomic
616         fragments. Bioinformatics **25:**1338-1340.

617   44.   **Fu LM, Niu BF, Zhu ZW, Wu ST, Li WZ.** 2012. CD-HIT: accelerated for clustering the
618         next-generation sequencing data. Bioinformatics **28:**3150-3152.

619   45.   **Hyatt D, Locascio PF, Hauser LJ, Uberbacher EC.** 2012. Gene and translation initiation
620         site prediction in metagenomic sequences. Bioinformatics **28:**2223-2230.

621   46.   **Huerta-Cepas J, Forslund K, Coelho LP, Szklarczyk D, Jensen LJ, Von MC, Bork P.**
622         2016. Fast genome-wide functional annotation through orthology assignment by
623         eggNOG-mapper. Mol Biol Evol **34:**2115.

624   47.   **Galperin MY, Makarova KS, Wolf YI, Koonin EV.** 2015. Expanded microbial genome
625         coverage and improved protein family annotation in the COG database. Nucl Acids Res
626         **43:**261-269.

627   48.   **Hutchison CA, 3rd, Chuang RY, Noskov VN, Assad-Garcia N, Deerinck TJ, Ellisman**
628         **MH, Gill J, Kannan K, Karas BJ, Ma L, Pelletier JF, Qi ZQ, Richter RA, Strychalski EA,**
629         **Sun L, Suzuki Y, Tsvetanova B, Wise KS, Smith HO, Glass JI, Merryman C, Gibson DG,**
630         **Venter JC.** 2016. Design and synthesis of a minimal bacterial genome. Science **351:**aad6253.

631   49.   **Li W-L, Huang J-M, Zhang P-W, Cui G-J, Wei Z-F, Wu Y-Z, Gao Z-M, Han Z, Wang Y.**
632         2019. Periodic and spatial spreading of alkanes and Alcanivorax bacteria in deep waters of the
633         Mariana Trench. Appl Environ Microbiol **85:**e02089-02018.

634   50.   **Patel RK, Jain M.** 2012. NGS QC toolkit: A toolkit for quality control of next generation
635         sequencing data. PLoS ONE **7:**e30619.

636   51.   **Zhang J, Kobert K, Flouri T, Stamatakis A.** 2014. PEAR: a fast and accurate Illumina
637         Paired-End reAd mergeR. Bioinformatics **30:**614.

638   52.   **Caporaso JG, Bittinger K, Bushman FD, Desantis TZ, Andersen GL, Knight R.** 2010.
639         PyNAST: a flexible tool for aligning sequences to a template alignment. Bioinformatics

640       **26:**266-267.

641   53.   **Wang Q, Garrity GM, Tiedje JM, Cole JR.** 2007. Naïve Bayesian Classifier for Rapid
642       Assignment of rRNA Sequences into the New Bacterial Taxonomy. Appl Environ Microbiol
643       **73:**5261.

644   54.   **Edgar RC.** 2004. MUSCLE: multiple sequence alignment with high accuracy and high
645       throughput. Nucl Acids Res **32:**1792-1797.

646   55.   **Capellagutiérrez S, Sillamartínez JM, Gabaldón T.** 2009. trimAl: a tool for automated
647       alignment trimming in large-scale phylogenetic analyses. Bioinformatics **25:**1972-1973.

648   56.   **Lam-Tung N, Schmidt HA, Arndt VH, Bui Quang M.** 2015. IQ-TREE: a fast and effective
649       stochastic algorithm for estimating maximum-likelihood phylogenies. Mol Biol Evol
650       **32:**268-274.

651   57.   **Kalyaanamoorthy S, Minh BQ, Wong TKF, Haeseler AV, Jermin LS.** 2017. ModelFinder:
652       fast model selection for accurate phylogenetic estimates. Nature Meth **14:** 587-589.

653   58.   **Wu M, Scott AJ.** 2012. Phylogenomic analysis of bacterial and archaeal sequences with
654       AMPHORA2. Bioinformatics **28:**1033-1034.

655   59.   **Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ.** 1990. Basic local alignment search
656       tool. J Mol Biol **215:**403-410.

657

658   Table 1. Representative genomes of RFN20 and RF39.

659   RF39 (HG1) was represented by HG1.1 and HG1.2 from the Tenericutes downloaded

660   from NCBI; RFN20 (HG2) was represented by HG2.1 and HG2.2. *S. azabuensis* was a

661   species in *Erysipetrichales*.

662

| ID | HG1.1 | HG1.2 | HG2.1 | HG2.2 | *Sharpea azabuensis* |
|---|---|---|---|---|---|
| Accession | UQAI01000000 | UQAG01000000 | UPZX01000000 | UQBB01000000 | JNKU00000000 |
| Genome size (bp) | 1,690,546 | 1,911,898 | 1,525,481 | 1,699,832 | 2,411,783 |
| %GC | 30 | 29.5 | 30.1 | 30.4 | 37.1 |
| No.contigs | 109 | 71 | 31 | 16 | 94 |
| %Complete | 98.7 | 98.7 | 98.9 | 98.5 | 99.1 |
| %Contaminant | 0 | 0 | 0 | 0 | 0.9 |
| No. tRNA | 38 | 35 | 34 | 45 | 57 |
| No. rRNA | 0 | 2 | 1 | 0 | 10 |
| %Coding density | 92 | 90.8 | 92.5 | 91.6 | 89 |
| No. CDSs | 1,548 | 1,834 | 1,488 | 1,570 | 2,424 |

663

664   Figure 1. Phylogenetic trees of Tenericutes

665   The maximum-likelihood phylogenetic trees were constructed by concatenated

666   conserved proteins (A) and 16S rRNA genes (B). The bootstrap values (>50) are

21

667    denoted by the dots on the branches.

668

669    Figure 2. Phylogenetic positions of Tenericutes families in Bacilli.

670    Representative genomes from orders of Bacilli were used to construct the

671    phylogenomics tree using concatenated conserved proteins by IQ-TREE and RAxML.

672    The bootstrap values were shown as triangles (50-90) and dots (>90) with a red color

673    for the results of RAxML and deep blue for those of IQ-TREE, respectively. The

674    purple clades represent the orders of Bacilli and the red ones denote Tenericutes.

675

676    Figure 3. Distribution of genes and pathways in the Tenericutes lineages.

677    Tenericutes lineages were grouped using an agglomerative hierarchical clustering on

678    the basis of the distribution of COGs within each group. The color and size of each

679    dot represent the percentage of genomes within each lineage that carries the gene. The

680    functions of these genes are shown in Table S1.

681

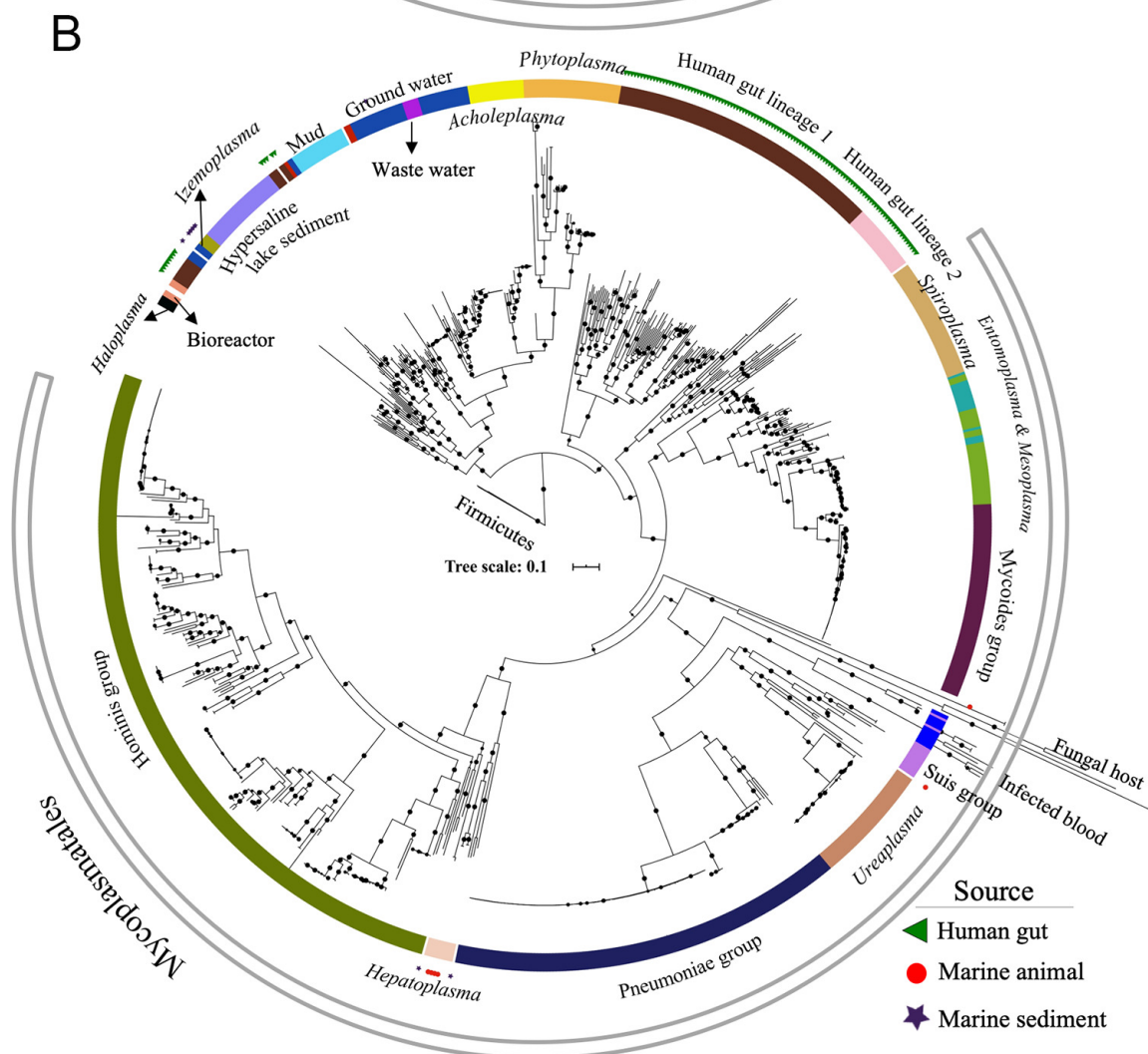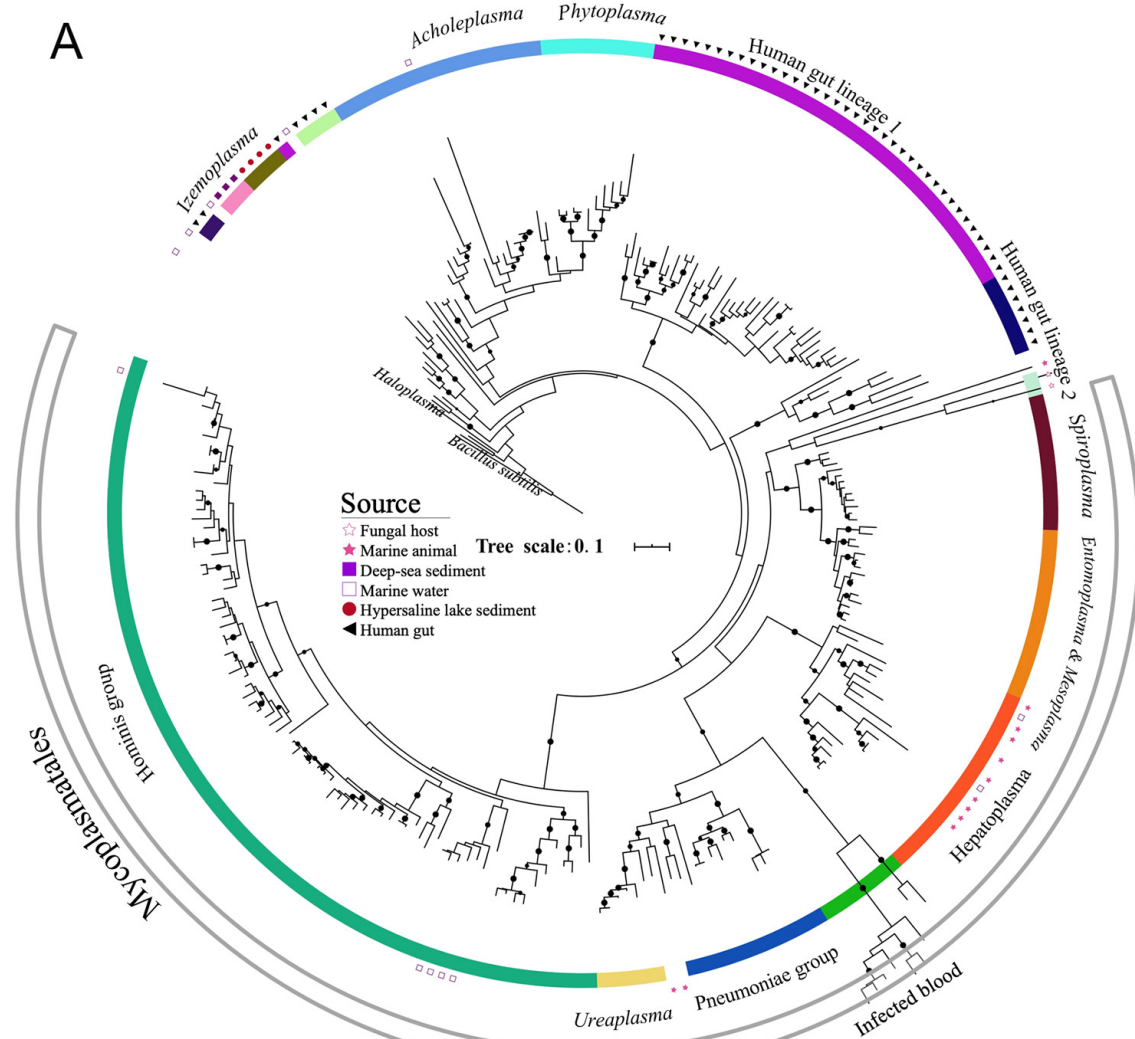682    Figure 4. Schematic metabolism of RFN20 and RF39

683    Metabolic models predicted by using gene annotation results of four representative

684    genomes of RFN20 and RF39 (see Table 1). Solid squares indicate presence of the

685    genes responsible for a step or a pathway. The products depicted in the MEP/DOXP

686    pathway are 1-deoxy-xylulose 5-P, 2-C-methyl-D-erythritol 4-P, 4-(Cytidine

687    5'-PP)-2-C-methyl-erythritol, 2-P-4-(cytidine 5'-PP)-2-C-methyl-erythritol,

688    2-C-methyl-erythritol 2,4-PP, 1-hydroxy-2-methyl-2-butenyl 4-PP, dimethylallyl-PP,
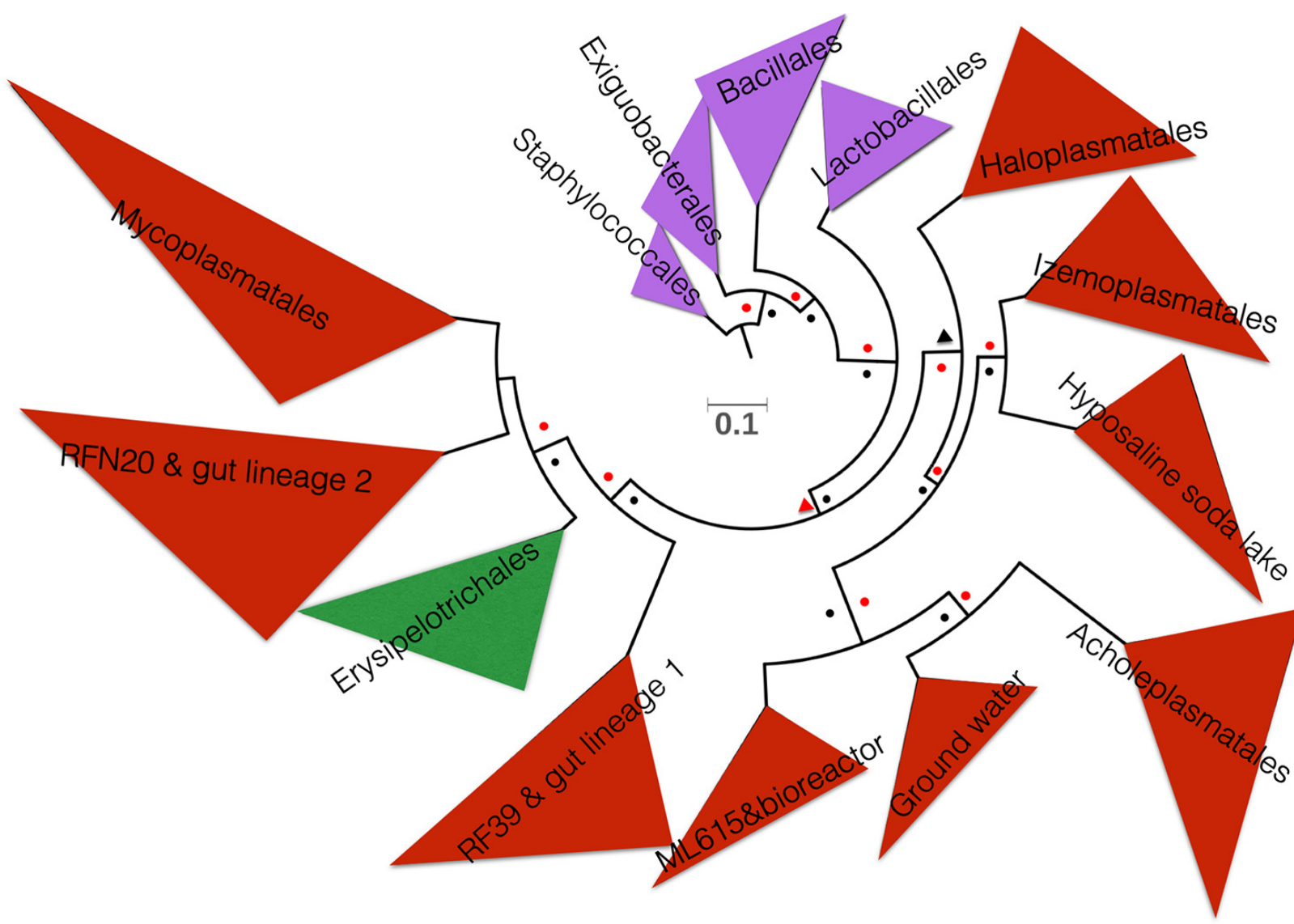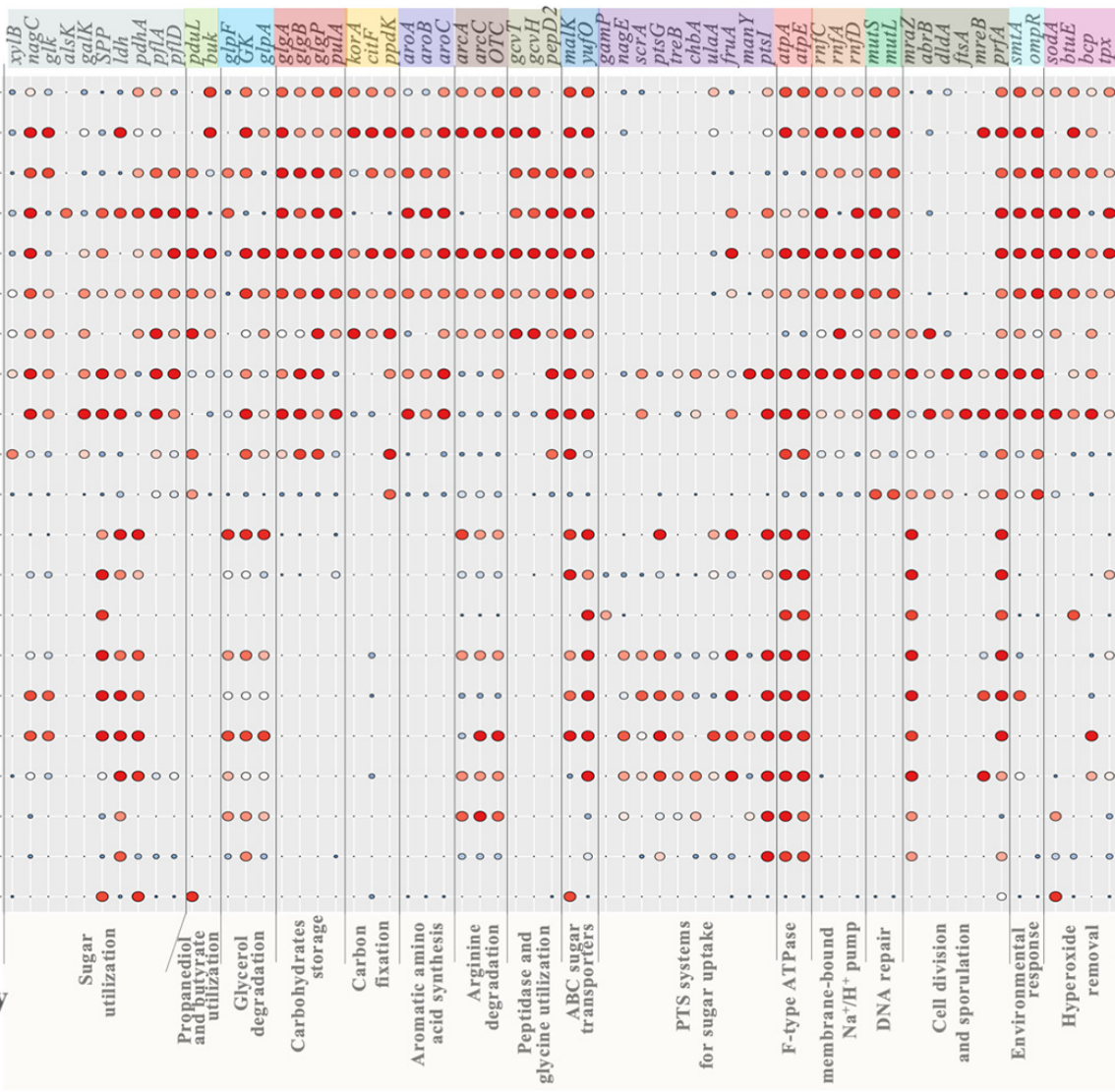
689    isopentenyl-PP, and farnesyl-PP.

690

691

692

693

A

*Acholeplasma*  *Phytoplasma*

*Izemoplasma*

*Haloplasma*

*Bacillus subtilis*

Source
☆ Fungal host
★ Marine animal
■ Deep-sea sediment
□ Marine water
● Hypersaline lake sediment
◄ Human gut

Tree scale: 0.1

Human gut lineage 1

Human gut lineage 2

*Spiroplasma*

*Entomoplasma & Mesoplasma*

*Hepatoplasma*

Infected blood

Pneumoniae group

*Ureaplasma*

Hominis group

*Mycoplasmatales*

B

*Izemoplasma*  Mud  Ground water  *Phytoplasma*  Human gut lineage 1

Hypersaline
lake sediment

*Acholeplasma*

Waste water

*Haloplasma*  Bioreactor

Firmicutes

Tree scale: 0.1

Human gut lineage 2

*Spiroplasma*

*Entomoplasma & Mesoplasma*

Mycoides group

Fungal host

Infected blood

Suis group

*Ureaplasma*

Pneumoniae group

*Hepatoplasma*

Hominis group

*Mycoplasmatales*

Source
▲ Human gut
● Marine animal
★ Marine sediment

Mycoplasmatales

RFN20 & gut lineage 2

Erysipelotrichales

Staphylococcales

Exiguobacterales

Bacillales

Lactobacillales

Haloplasmatales

Izemoplasmatales

Hyposaline soda lake

Acholeplasmatales

Ground water

ML615&bioreactor

RF39 & gut lineage 1

0.1