**FIG S7.  The evolutionary trajectory of six *Rickettsia* biosynthetic pathways that contain holes, or "missing enzymes".**

Our *Rickettsia* metabolic reconstruction identified six holes occurring in pathways for the synthesis of CMP-3-deoxy-β-D-manno-octulosonate (CMP-Kdo), Diaminopimelate (DAP), CDP-diacylglycerol (CDP-DG), terpenoid backbones, ubiquinone ($CoQ_8$), and queuosine (see manuscript for further details).  These pathways were analyzed in the genomes of Rickettsiales and Holosporales to determine the presence/absence of the missing enzymes.  If present, phylogenies were estimated to determine if the genes encoding these enzymes were present in the Rickettsiales/Holosporales common ancestor.

(A) *E. coli* enzymes corresponding to the holes in the *Rickettsia* CMP-Kdo (KdsC), DAP (DapC), CDP-DG (PlsX, PlsY or PlsB), terpenoid backbones (IspA), $CoQ_8$ (UbiC, UbiI), and queuosine (QueG) pathways were used as queries in blastp searches against the Rickettsiales and Holosporales databases at NCBI.  For significant matches, only top hits with coverage to the query greater than 85% were considered.  The four obtained rickettsial homologs were then used as queries in blastp searches against the *Alphaproteobacteria* (excluding Rickettsiales and Holosporales) and the NR (excluding *Alphaproteobacteria*) databases at NCBI.  Again, for significant matches, only top hits with coverage to the query greater than 85% were considered.

(B) Phylogenomics analysis of CMP-Kdo, DAP, CDP-DG, terpenoid backbones, $CoQ_8$, and queuosine pathways across selected Holosporales and Rickettsiales genomes.  The presence/absence of enzymes corresponding to the holes in the *Rickettsia* pathways are shown at right.  Genome-based phylogeny was estimated on 105 orthologous groups (OGs) of proteins (single copy CDS in over 95% of genomes).  We sampled 53 Rickettsiales/Holosporales and three outgroup *Alphaproteobacteria* genomes: IDs for sequences retrieved from PATRIC (P) and NCBI (N) are shown for each taxon.  Several genomes were specifically annotated for this project: *Caedibacter varicaedens* (NCBI, NZ_BBVC00000000.1); "*Candidatus* Arcanobacter lacustris" str. SCGC (NCBI, JYHA00000000.1); "*Candidatus* Paracaedibacter symbiosus" (NCBI,

NZ_JQAK00000000.1); "*Candidatus* Xenolissoclinum pacificiensis" str. L6 (NCBI, AXCJ00000000.1); *Rickettsia tamurae* str. AT-1 (NCBI, CCMG01000000.1); "*Candidatus* Hepatobacter penaei" str. NHPB (NCBI, JQAJ00000000.1); *Rickettsia* endosymbiont of *Adalia bipunctata* (unpublished); *Rickettsia* endosymbiont of *Ichthyophthirius multifiliis* (unpublished); *Rickettsia* sp. MEAM1 str. *Bemisia tabaci* (unpublished). Data for Rickettsiales endosymbiont of *Trichoplax adhaerens* is described elsewhere (1). OGs were generated using FastOrtho, a modified version of OrthoMCL (2). Multiple sequence alignment of each OG was performed using MUSCLE (default parameters) (3), with regions of poor alignment (length heterogeneous regions) masked using Gblocks (4). All modified alignments were concatenated into a single dataset for phylogeny estimation. Using PhyloBayes MPI (5), we analyzed the dataset with the CAT model of substitution, a nonparametric method for modeling site-specific features of sequence evolution (6, 7). Given the strong base compositional bias of rickettsial genomes (~30 %GC), the ability of the CAT model to accommodate saturation due to convergences and reversions (8) is of substantial importance for estimating rickettsial phylogeny, as demonstrated by us and others (1, 9–11). Two independent Markov chains were run in parallel using PhyloBayes MPI v.1.2e under the CAT-GTR model, with the bipartition frequencies analyzed at various time points using the bpcomp program. For tree-building, appropriate burn-in values were determined by plotting the log likelihoods for each chain over sampled generations (time). Analyses were considered complete when the maximum difference in bipartition frequencies between the two chains was less than 0.1. Ultimately, a burn-in value of 1000, with sampling every 2 trees, was used to build a consensus tree.

(C-F) Phylogeny estimations of DapC (C), PlsX (D), PlsY (E), and IspA (F) proteins. Datasets for each protein were constructed as follows: the rickettsial protein (panel A) was used in blastp queries against several taxon-specific databases: 1) "Rickettsiales", 2) "Holosporales", 3) "*Alphaproteobacteria* (minus Rickettsiales and Holosporales)", 4) "*Proteobacteria* (minus *Alphaproteobacteria*)", 5) "Bacteria (minus *Proteobacteria*)", and 6) "minus Bacteria". The top 5-

10 (query-dependent) subjects from each search resulting in significant (> 40 bits) alignments were all compiled and aligned using MUSCLE v3.8.31 (default parameters). Protein phylogenies were estimated under maximum likelihood with RAxML v8.2.4 (12), using a gamma model of rate heterogeneity and estimation of the proportion of invariant sites. Both the Lee and Gascuel (LG) and Blocks of Amino Acid Substitution Matrix (BLOSUM62) models of amino acid substitution were used, and branch support was assessed with 1,000 pseudo-replications.

## REFERENCES

1.  **Driscoll T**, **Gillespie JJ**, **Nordberg EK**, **Azad AF**, **Sobral BW**. 2013. Bacterial DNA sifted from the Trichoplax adhaerens (Animalia: Placozoa) genome project reveals a putative rickettsial endosymbiont. Genome Biol Evol **5**:621–45.

2.  **Li L**, **Stoeckert CJ**, **Roos DS**. 2003. OrthoMCL: identification of ortholog groups for eukaryotic genomes. Genome Res **13**:2178–89.

3.  **Edgar RC**. 2004. MUSCLE: Multiple sequence alignment with high accuracy and high throughput. Nucleic Acids Res **32**:1792–1797.

4.  **Talavera G**, **Castresana J**. 2007. Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. Syst Biol **56**:564–77.

5.  **Lartillot N**, **Rodrigue N**, **Stubbs D**, **Richer J**. 2013. PhyloBayes MPI: phylogenetic reconstruction with infinite mixtures of profiles in a parallel environment. Syst Biol **62**:611–5.

6.  **Lartillot N**, **Philippe H**. 2004. A Bayesian Mixture Model for Across-Site Heterogeneities in the Amino-Acid Replacement Process. Mol Biol Evol **21**:1095–1109.

7.  **Lartillot N**, **Philippe H**. 2006. Computing Bayes factors using thermodynamic integration. Syst Biol **55**:195–207.

8.  **Lartillot N**, **Brinkmann H**, **Philippe H**. 2007. Suppression of long-branch attraction artefacts in the animal phylogeny using a site-heterogeneous model. BMC Evol Biol **7**

**Suppl 1**:S4.

9.      **Rodríguez-Ezpeleta N**, **Embley TM**. 2012. The SAR11 group of alpha-proteobacteria is

not related to the origin of mitochondria. PLoS One **7**:e30520.

10.     **Viklund J**, **Ettema TJG**, **Andersson SGE**. 2012. Independent genome reduction and

phylogenetic reclassification of the oceanic SAR11 clade. Mol Biol Evol **29**:599–615.

11.     **Gillespie JJ**, **Driscoll TP**, **Verhoeve VI**, **Utsuki T**, **Husseneder C**, **Chouljenko VN**,

**Azad AF**, **Macaluso KR**. 2014. Genomic diversification in strains of Rickettsia felis

isolated from different arthropods. Genome Biol Evol.

12.     **Stamatakis A**. 2014. RAxML version 8: A tool for phylogenetic analysis and post-analysis

of large phylogenies. Bioinformatics **30**:1312–1313.

**FIG. S7**

**A**

| Hole | Reference protein | Top Rick (% ID) | | Top *Alpha* E value; % ID | | Top non-*Alpha* E value; % ID | |
|---|---|---|---|---|---|---|---|
| KdsC | P0ABZ4: *E. coli* str. K12 | --- | | --- | | --- | |
| DapC | P18335: *E. coli* str. K12 | WP_070064655: | *Wolbachia pipientis* (42%) | 4e-130; | 50% | 1e-122; | 48% |
| PlsX | P27247: *E. coli* str. K12 | WP_011452016: | *Neorickettsia sennetsu* (45%) | 9e-94; | 47% | 4e-77; | 45% |
| PlsY | P60782: *E. coli* str. K12 | KKB96285: | "*Cand.* Arcanobacter lacustris" (40%) | 5e-50; | 49% | 2e-36; | 49% |
| PlsB | P0A7A7: *E. coli* str. K12 | --- | | --- | | --- | |
| IspA | P22939: *E. coli* str. K12 | WP_039459403: | "*Cand.* Jidaibacter acanthamoeba" (44%) | 7e-109; | 59% | 2e-72; | 55% |
| UbiC | P26602: *E. coli* str. K12 | --- | | --- | | --- | |
| UbiI | P25535: *E. coli* str. K12 | --- | | --- | | --- | |
| QueG | P39288: *E. coli* str. K12 | --- | | --- | | --- | |

**B**



Columns (top): KdsC, DapC, PlsX, PlsY, PlsB, IspA, UbiC, UbiI, QueG

*Azospirillum* sp. B510 [P, 137722.3]
*Rhodospirillum rubrum* str. ATCC 11170 [P,269796.9]
alphaproteobacterium str. BAL199 [P, 331869.3] } Other *Alphaproteobacteria*

"*Candidatus* Paracaedibacter symbiosus"
"*Candidatus* Odyssella thessalonicensis" str. L13 [P, 985867.6]
"*Candidatus* Paracaedibacter acanthamoebae" [N, 32950]
"*Candidatus* Caedibacter acanthamoebae" [N, 32965]
*Caedibacter varicaedens*
"*Candidatus* Hepatobacter penaei" str. NHPB
*Holospora obtusa* str. F1 [N, 22931]
*Holospora elegans* str. E1 [N, 31836]
*Holospora undulata* str. HU1 [N, 17492]

**HOLOSPORALES**

Rickettsiales bacterium str. Ac37b [N, 33266]
"*Candidatus* Arcanobacter lacustris" str. SCGC
"*Candidatus* Occidentia massiliensis" str. Os18 [N, 30237]
*Orientia tsutsugamushi* str. Ikeda [P,334380.4]
*Orientia chuto* str. Fuller [N, 36736]

**RICKETTSIACEAE**

*Rickettsia* sp. (*Ichthyophthirius multifiliis*)
*Rickettsia bellii* str. RML369-C [P, 336407.7]
*Rickettsia* sp. MEAM1 str. *Bemisia tabaci*
*Rickettsia* sp. (*Adalia bipunctata*)
*Rickettsia canadensis* str. McKiel [P, 293613.4]
*Rickettsia helvetica* str. C9P9 [P, 1144888.4]  — basal lineages
*Rickettsia tamurae* str. AT-1
*Rickettsia conorii* str. Malish 7 [P, 272944.4]  — SFG
*Rickettsia massiliae* str. MTU5 [P, 416276.5]
*Rickettsia prowazekii* str. Madrid E [P, 272947.5]  — TG
*Rickettsia typhi* str. Wilmington [P, 257363.4]
*Rickettsia akari* str. Hartford [P, 293614.5]  — TRG
*Rickettsia felis* str. URRWXCal2 [P, 315456.7]

*Rickettsia* pathway holes

Branches increased horizontally by 400%

Rickettsiales endosymbiont of *Trichoplax adhaerens*
"*Candidatus* Midichloria mitochondrii" str. IricVA [P, 696127.4]
Rickettsial endosymbiont of Acanthamoeba sp. UWC8 [N, 32519]
"*Candidatus* Jidaibacter acanthamoeba" [N, 45806]

**MIDICHLORIACEAE**

"*Candidatus* Xenolissoclinum pacificiensis" str. L6
*Neorickettsia helminthoeca* str. Oregon [N, 16400]
*Neorickettsia sennetsu* str. Miyayama [P, 222891.8]
*Neorickettsia risticii* str. Illinois [P, 434131.3]
*Wolbachia* endosymbiont of *Onchocerca volvulus* str. Cameroon [N, 11053]
*Wolbachia* endosymbiont of *Cimex lectularius* [N, 35952]
Wolbachia endosymbiont strain TRS of *Brugia malayi* [P, 292805.6]
*Wolbachia* endosymbiont of *Wuchereria_bancrofti* [P, 96496.5]
*Wolbachia* endosymbiont of *Drosophila melanogaster* [P, 163164.16]
*Wolbachia* endosymbiont of *Drosophila simulans* [P, 77038.12]
*Wolbachia* endosymbiont of *Culex quinquefasciatus* JHB [P, 569881.3]
Wolbachia pipientis wAlbB [N, 11990]
"*Candidatus* Neoehrlichia lotoris" str. RAC413 [N, 36737]
*Anaplasma phagocytophilum* str. HZ [P, 212042.8]
*Anaplasma centrale* str. Israel [P, 574556.4]
*Anaplasma marginale* str. St. Maries [P, 234826.6]
*Ehrlichia ruminantium* str. Gardel [P, 302409.5]
*Ehrlichia canis* str. Jake [P, 269484.6]
*Ehrlichia minasensis* [N, 35818]
*Ehrlichia chaffeensis* str. Arkansas [P, 205920.11]
*Ehrlichia muris* str. AS145 [N, 23922]
*Ehrlichia* sp. HF [N, 31756]

**ANAPLASMATACEAE**

\* some genes may never have been present after the divergence from the other *Alpha-proteobacteria*

0.5 sub./site

Legend:
- Loss of KdsC
- Loss of LPS entirely (*X*)
- Loss of DapC
- DapC, but no PGN
- Loss of PlsB
- Loss of PlsX/Y
- Loss of IspA
- Loss of UbiC
- Loss of UbiI
- Loss of QueDECFAG (*X*)
- Gain of QueDECFA

**C**

**DapC: *N*-succinyldiaminopimelate aminotransferase**

Archaeoglobus_fulgidus__WP_010877594.1

Coccomyxa_subellipsoidea_C-169___XP_005651785.1

Caldanaerobacter_subterraneus__WP_022588199.1

Piscirickettsia_salmonis__WP_016209980.1

100 — Candidatus_Odyssella_thessalonicensis__WP_033444530.1

Candidatus_Paracaedibacter_symbiosus__WP_052046308.1

Anaerovibrio_lipolyticus__KHM53015.1

100 — Methanococcus_maripaludis__WP_012194296.1

Methanotorris_igneus___WP_013799319.1

100 — Smithella_sp._D17__KFZ44386.1

50 — Smithella_sp._SCADC__KFO68501.1

Moorella_glycerini__WP_054936021.1

Desulfococcus_multivorans__WP_040413854.1

100 — Geobacter_bremensis__WP_026841615.1

Malonomonas_rubra__WP_072909607.1

Idiomarina_abyssalis__WP_054490255.1

Mariprofundus_micogutta__WP_072659455.1

52 — Numidum_massiliense__WP_074011201.1

87 — Bacillus_mannanilyticus__WP_025028266.1

97 — Aneurinibacillus_migulanus__WP_043067692.1

Brevibacillus_brevis__WP_017252253.1

98 — Candidatus_Xenolissoclinum_pacificiensis_L6__ETO91444.1

100 — Neorickettsia_helminthoeca__WP_038560077.1

99 — Neorickettsia_risticii__WP_015816652.1

Neorickettsia_sennetsu__WP_011452223.1

93 — Rickettsiales_bacterium_Ac37b__WP_038603157.1

100 — Wolbachia_endosymbiont_of_Brugia_malayi___WP_011256330.1

85 — Ehrlichia_minasensis__WP_045170781.1

95 — Anaplasma_centrale__WP_012880909.1

Candidatus_Neoehrlichia_lotoris__WP_045809287.1

98 — Magnetospirillum_caucaseum__EME67964.1

100 — Rhodovibrio_salinarum__WP_037257124.1

Phormidium_willei__WP_068791131.1

86 — 100 — Paraburkholderia_terrae_BS001__EIM93111.1

69 — Burkholderia_sp._JPY251__WP_018437429.1

100 — Bradyrhizobium_neotropicale__WP_063680178.1

Rhodopseudomonas_palustris__WP_047307002.1

Cucumibacter_marinus__WP_035872516.1

69 — Lutibaculum_baratangense__WP_040484875.1

100 — Streptomyces_purpurogeneiscleroticus___KOX60391.1

Asanoa_ferruginea__WP_053620795.1

Aquamicrobium_defluvii__WP_035026380.1

72 — Ahrensia_marina__WP_053999841.1

Hoeflea_olei__WP_066182554.1

Martelella_endophytica__WP_045680504.1

51 — Agrobacterium_tumefaciens__WP_035219093.1

0.1 sub./site

FIG. S7

**D**

**PlsX: Phosphate acyltransferase**

Legend:
- Rickettsiales
- Holosporales
- other *Alphaproteobacteria*
- other Proteobacteria
- other bacteria

100 Holophaga_foetida___WP_005036190.1
Geothrix_fermentans__WP_026853009.1
Carboxydothermus_ferrireducens__WP_034542157.1
Thermovibrio_ammonificans__WP_013537666.1
Clostridium_tetani__WP_023438126.1
95 Brevibacillus_borstelensis___WP_024982984.1
100 Bacillus_flexus__WP_061784595.1
Parageobacillus_toebii__WP_062753429.1
100 Morganella_morganii____WP_071592728.1
Morganella_sp._HMSC11D09__WP_004240652.1
Buttiauxella_ferragutiae__WP_074388586.1
79 Kluyvera_georgiana__WP_074398851.1
Kluyvera_ascorbata__WP_072011380.1
Salmonella_enterica__WP_074177156.1
Citrobacter_farmeri__WP_072015706.1
Escherichia_fergusonii__WP_002431498.1
Lelliottia_amnigena__WP_064325693.1
Leclercia_adecarboxylata__WP_071886466.1
Enterobacter_cloacae__WP_028019051.1
62 Enterobacter_hormaechei_WP_044489107.1

**Branches increased horizontally by 300%**

98 Candidatus_Xenolissoclinum_pacificiensis_L6__ETO91046.1
100 Neorickettsia_helminthoeca__WP_038559727.1
100 Neorickettsia_sennetsu___WP_011452016.1
Neorickettsia_risticii__WP_041351503.1
94 Anaplasma_centrale___WP_012880566.1
92 Ehrlichia_minasensis__WP_045171340.1
Wolbachia_endosymbiont_of_Brugia_malayi__WP_011256813.1
Candidatus_Neoehrlichia_lotoris_str._RAC413___KJV69026.1
Rickettsiales_bacterium_Ac37b___WP_038602985.1
100 Candidatus_Jidaibacter_acanthamoeba__KIE06001.1
Candidatus_Midichloria_mitochondrii__WP_013950765.1
99 Candidatus_Paracaedibacter_symbiosus__WP_032113756.1
Candidatus_Odyssella_thessalonicensis__WP_010297662.1
77 Caedibacter_varicaedens__GAO97392.1
98 Candidatus_Hepatobacter_penaei__WP_052545554.1
Holospora_elegans___WP_035544044.1
Thalassospira_australica___WP_052065952.1
Tistrella_mobilis__WP_062765116.1
58 Kiloniella_spongiae__WP_047765600.1
Oceanibaculum_indicum_P24__EKE77690.1
Azospirillum_halopraeferens__WP_029011194.1
Magnetospirillum_marisnigri__WP_068491985.1
Caulobacter_vibrioides__WP_010919245.1
Rhodospira_trueperi__SDD69148.1
81 Rhodospirillum_rubrum__WP_011389356.1
99 Pararhodospirillum_photometricum__WP_041794325.1

100, 100, 53, 82, 100, 73, 94, 94, 64, 100, 53

0.1 sub./site

**FIG. S7**

**E**

**PlsY: glycerol-3-phosphate acyltransferase**

Legend:
- Rickettsiales
- Holosporales
- other *Alphaproteobacteria*
- other Proteobacteria
- other bacteria
- non-bacteria

Halolamina_sediminis__WP_071946839.1
Alcanivorax_jadensis_T9__KGD62212.1
Gallaecimonas_pentaromativorans__WP_050657144.1
Paraglaciecola_agarilytica__WP_039995076.1
Shewanella_waksmanii__WP_028774292.1
Suttonella_ornithocola__WP_072577108.1
Methylomicrobium_agile__WP_031429895.1
Methylomonas_koyamae__WP_064025399.1
Piscirickettsia_salmonis__WP_017376031.1
Taylorella_asinigenitalis__WP_014111403.1
Caedibacter_varicaedens___WP_062141522.1
Loktanella_rosea__WP_076659499.1
Pseudovibrio_denitrificans__WP_054784807.1
Bartonella_henselae__WP_011180713.1
Acetobacter_aceti__AQS86387.1
Candidatus_Odyssella_thessalonicensis__WP_010298129.1
Candidatus_Paracaedibacter_symbiosus__WP_032112393.1
Kordiimonas_lipolytica__WP_068144749.1
Dinoroseobacter_shibae___WP_012179756.1
Paracoccus_aminovorans__WP_074967348.1
Kiloniella_laminariae__WP_020593536.1
Rhodovibrio_salinarum__WP_027288463.1
Tistrella_mobilis__WP_062763456.1
Wolbachia_endosymbiont_of_Brugia_malayi__WP_011256621.1
Anaplasma_centrale__WP_012880297.1
Candidatus_Neoehrlichia_lotoris__WP_045809101.1
Ehrlichia_minasensis___WP_045171213.1
Candidatus_Jidaibacter_acanthamoeba__WP_039458599.1
Candidatus_Midichloria_mitochondrii__WP_013951329.1
Candidatus_Xenolissoclinum_pacificiensis_L6__ETO91536.1
Neorickettsia_helminthoeca__WP_038559902.1
Holospora__WP_006290506.1
Rickettsiales_bacterium_Ac37b__WP_038602759.1
Candidatus_Arcanobacter_lacustris__KKB96285.1
Hippea_jasoniae__WP_035587756.1
Leptospirillum_ferrooxidans__WP_041774774.1
Thermosipho_africanus___WP_004100835.1
Emergencia_timonensis__WP_067541539.1
Candidatus_Hepatobacter_penaei__WP_031933713.1
Candidatus_Heimdallarchaeota_archaeon_LC_2__OLS29070.1
Anaerococcus_prevotii__WP_015777848.1
Chlamydia_trachomatis__CRH66047.1
Paraclostridium_bifermentans__WP_025162696.1
Streptobacillus_ratti__WP_072592637.1
Staphylococcus_hyicus__WP_039645874.1
Bacillus_vietnamensis___WP_060672032.1
Lactobacillus_aviarius__WP_064208208.1

Support values: 89, 78, 93, 67, 74, 99, 71, 62, 79, 91, 90

0.1 sub./site

FIG. S7

**F**

**IspA: farnesyl diphosphate synthase**

Ostreococcus_lucimarinus_CCE9901__XP_001420781.1
Syntrophus_aciditrophicus__WP_041585559.1
100 ┌ Selenomonas_artemidis_F0399___EFW30636.1
    └ Mitsuokella_multacida__WP_040635779.1
Desulfovirgula_thermocuniculi__WP_027719263.1
82 Desulfotomaculum_ferrireducens__WP_077713465.1
Butyricicoccus_pullicaecorum__WP_016146829.1
55 Acetivibrio_cellulolyticus__WP_010252615.1
97 Thermoanaerobacter_italicus__WP_012995156.1
99 Caldanaerobacter_subterraneus__KUK34508.1
Geobacter_bemidjiensis__WP_012529688.1
Desulfococcus_multivorans__WP_020876973.1
Chloracidobacterium_thermophilum__ABV27206.1
uncultured_marine_thaumarchaeote_KM3_54_G03__AIF12168.1
74 Pseudoalteromonas_piscicida__WP_045964627.1
Acidithiobacillus_caldus__WP_004873145.1
Hahella_chejuensis__WP_011399574.1
51 Alkalilimnicola_ehrlichii__WP_011628694.1
Beggiatoa_leptomitiformis__WP_062148244.1
Ectothiorhodospira_haloalkaliphila__WP_025280698.1
74 Crenothrix_polyspora__SJM89924.1
98 Porphyrobacter_sanguineus__WP_072672911.1
Mastigocladus_laminosus__WP_072046828.1
65 Zymomonas_mobilis__WP_012817111.1
94 Insolitispirillum_peregrinum___WP_076400376.1
Novispirillum_itersonii___WP_019644662.1
77 Elstera_litoralis___WP_045775587.1
Azospirillum_halopraeferens__WP_051340468.1
50 Rhodospirillum_centenum__WP_041785096.1
91 Niveispirillum_irakense__WP_029014580.1
95 Candidatus_Paracaedibacter_symbiosus__WP_052046184.1
Candidatus_Odyssella_thessalonicensis___WP_010298776.1
Reyranella_massiliensis__WP_020696399.1
Micavibrio_aeruginosavorus__WP_015467093.1
Candidatus_Jidaibacter_acanthamoeba__WP_039459403.1
Candidatus_Arcanobacter_lacustris__KKB96385.1
74 Rickettsiales_bacterium_Ac37b__WP_038602013.1
98 Caedibacter_varicaedens__WP_062140753.1
Holospora__WP_006295417.1
87 Ehrlichia_minasensis__CEI85257.1
Candidatus_Neoehrlichia_lotoris_str._RAC413__KJV69019.1
91 Anaplasma_centrale__WP_041651145.1
Wolbachia_endosymbiont_of_Brugia_malayi__WP_011256934.1
Candidatus_Midichloria_mitochondrii__WP_013950557.1
74 100 Neorickettsia_sennetsu__WP_011452077.1
Neorickettsia_risticii__WP_015816528.1

— 0.1 sub./site

FIG. S7