

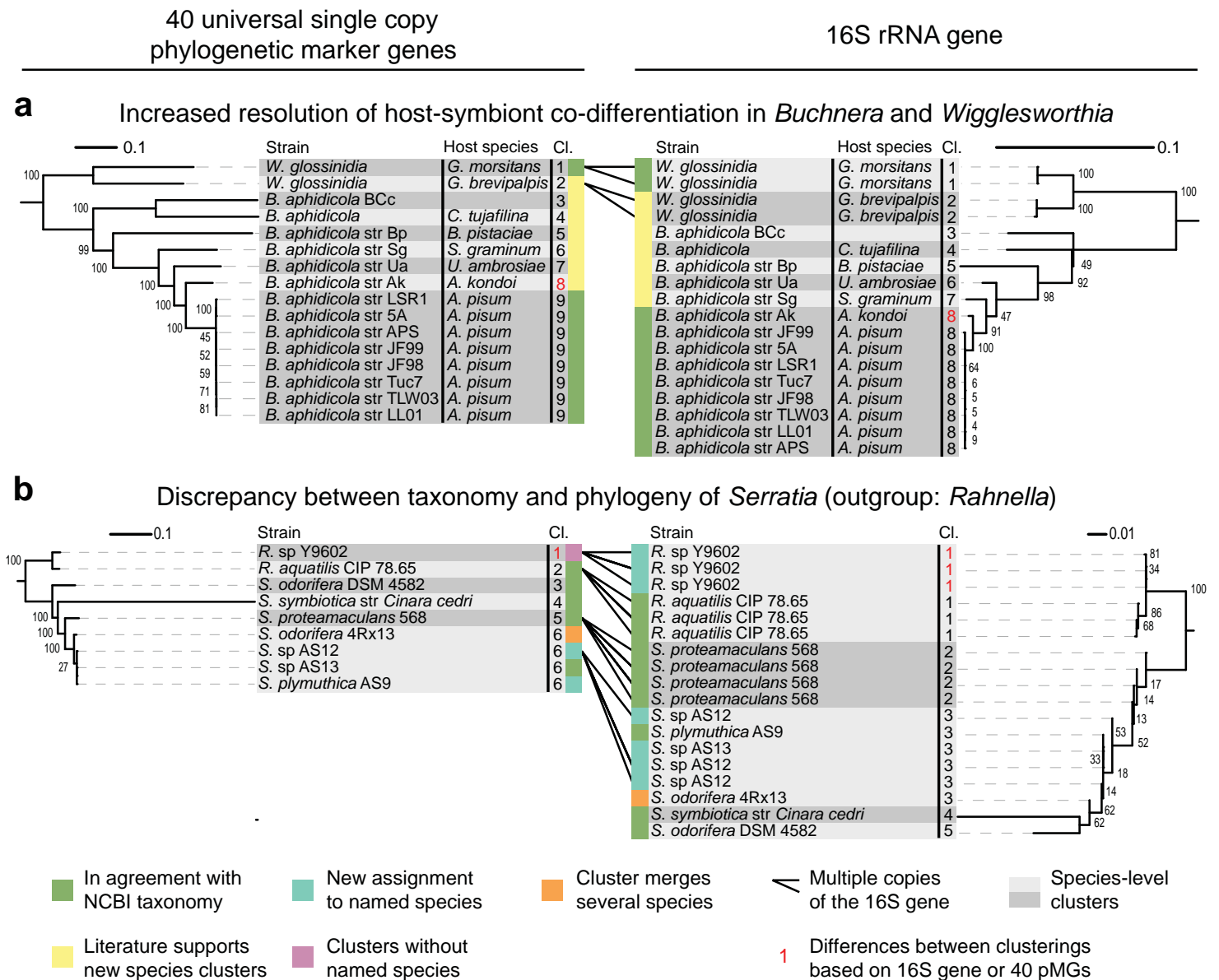








Supplementary Figure 5: Phylogenetic trees and species-level clustering of two additional clades displaying discrepancies to the NCBI taxonomy



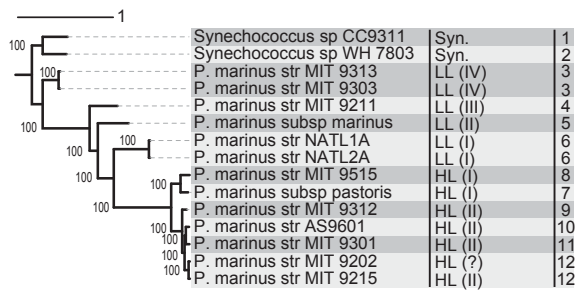
Phylogenetic trees and species-level clustering of two additional clades displaying discrepancies to the NCBI taxonomy. Trees were independently built either from a concatenation of 40 universal single-copy phylogenetic marker genes (pMGs) or from the 16S rRNA gene.

a) Host-symbiont co-speciation in *Buchnera* and *Wigglesworthia* is consistent with the highly resolved species clustering based on the 40 concatenated pMGs and in part also supported by the 16S rRNA gene.

b) Discrepancy between taxonomy and reconstructed phylogeny of *Serratia*. Both the combined MGs and the 16S rRNA gene support a reclassification of strain *S. odorifera* 4Rx13 forming one species together with *S. plymuthica* AS9 (and the two unnamed genomes sp AS12 and sp AS13). Please also note the long branch found for *S. symbiotica* str *Cinara cedri*, that is representative of a change of lifestyle from facultative to obligatory symbionts (Lamelas et al. PLoS Genetics, 2011). Color key: clustering in agreement with the NCBI taxonomy (blue); clustering discrepant to the NCBI taxonomy, but supported by literature information (yellow); assignment of an unnamed species to a cluster formed of a named species (pink); species clusters that do not contain any named species (dark violet); clusters which merge different species and thus have multiple species names (orange); genomes of species that are split among different clusters (red boxes); discrepancies between the clusterings of 40 pMGs and the 16S rRNA gene (red numbers, other are in black print); multiple copies of the 16S rRNA gene in one genome (black lines; identical copies were removed beforehand).

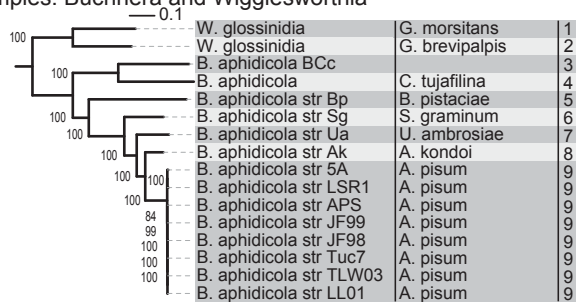
### One-to-one orthologs

#### a) Ecotype delineation in *Prochlorococcus*



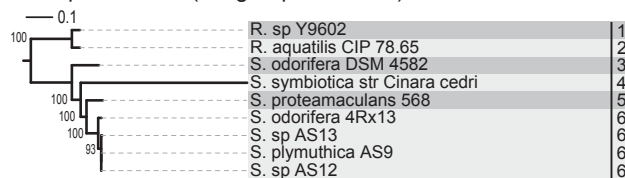
#### b) Increased resolution of host-symbiont co-differentiation.

Examples: *Buchnera* and *Wigglesworthia*



#### c) Incongruence between taxonomy and phylogeny.

Example: *Serratia*(Outgroup: *Rahnella*)

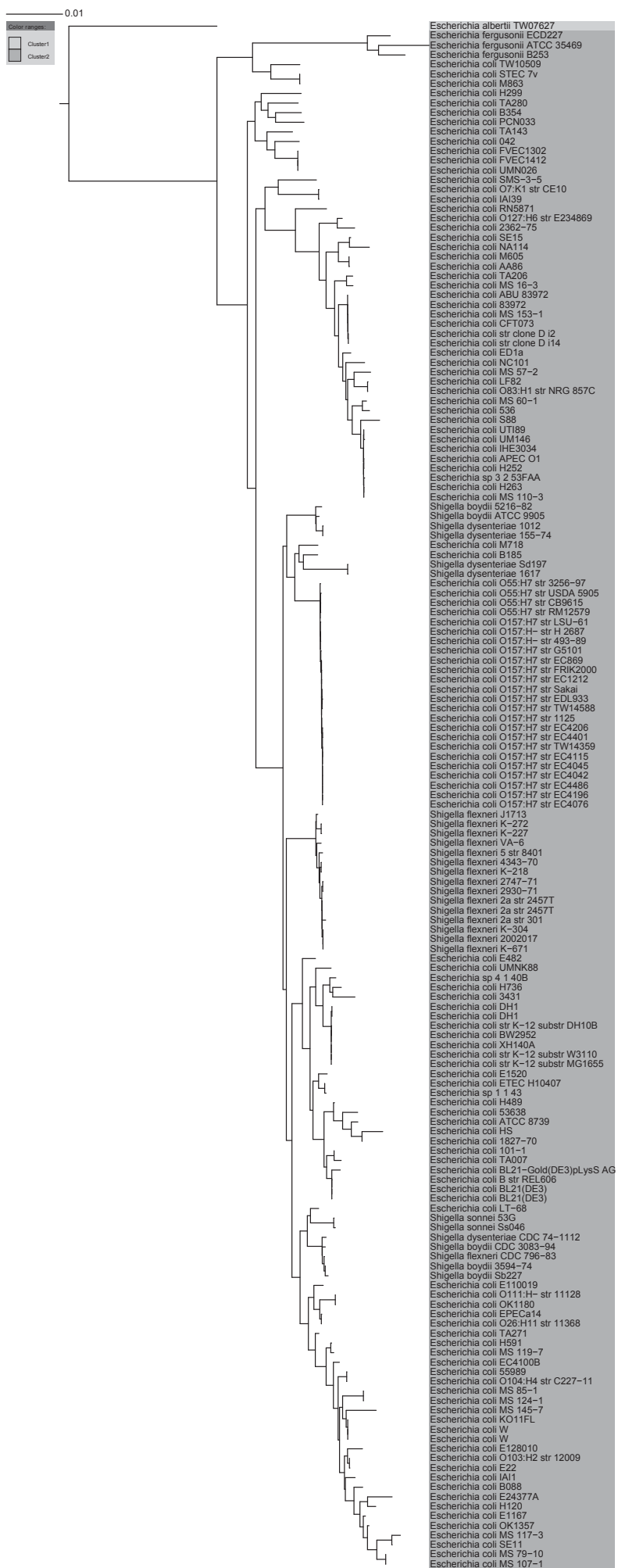


Phylogenetic trees and species level clustering of three clades displaying discrepancies to the NCBI taxonomy. Trees were independently built from a concatenation of all one-to-one orthologs found in the core-genomes used to build the trees.

a) Species-level clusterings and phylogenetic trees of the all one-to-one orthologs support speciation of so-called ecotypes in *Prochlorococcus marinus*, supporting the finer species resolution of the 40 combined pMGs resolving species in comparison to the 16S rRNA gene.

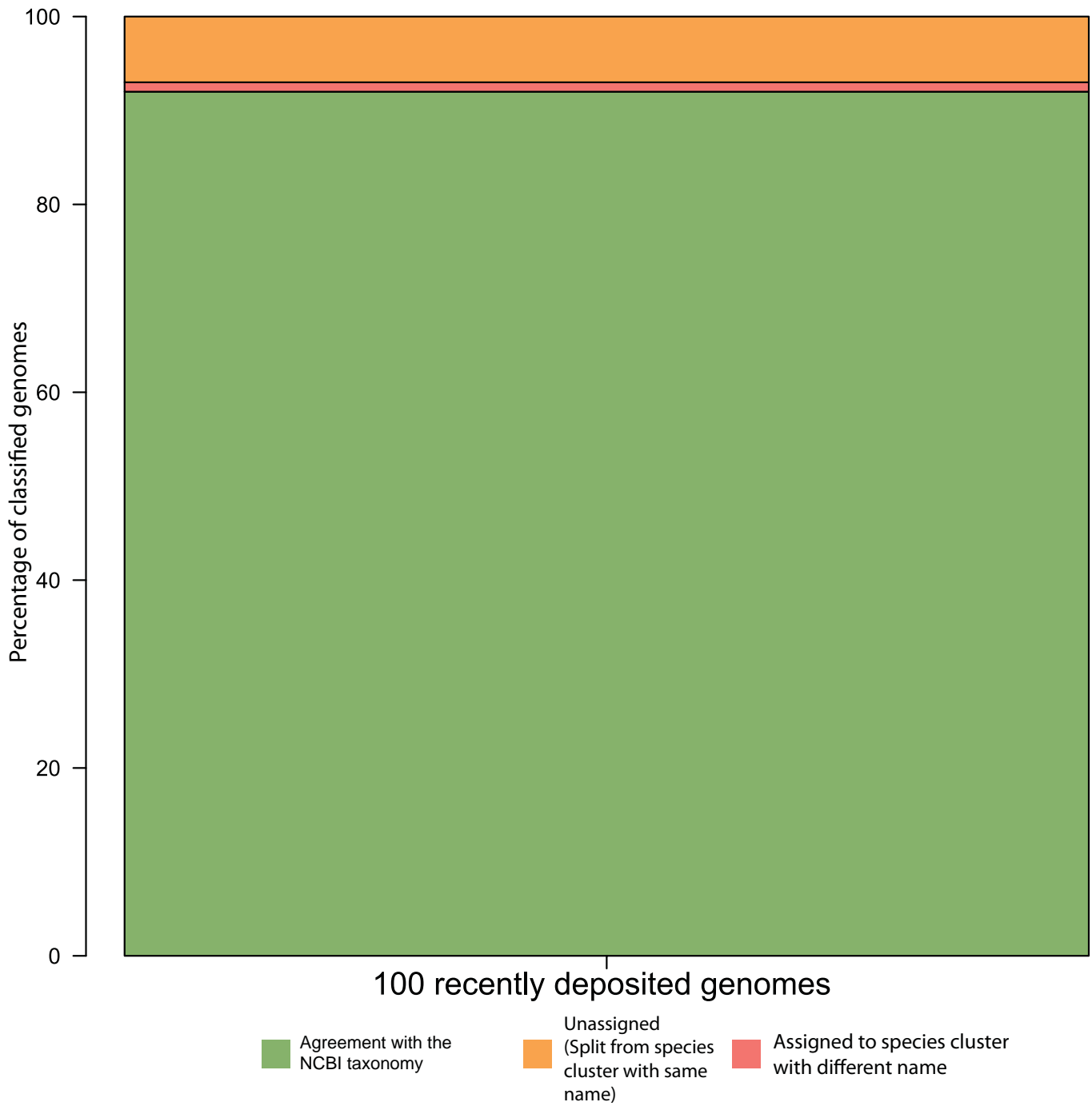
b) Discrepancy between taxonomy and reconstructed phylogeny of *Serratia*, especially strain *S. odorifera* 4Rx13 is also supported by trees built from all one-to-one orthologs

c) Host-symbiont co-differentiation in *Buchnera* and *Wigglesworthia* visible when using the 40 combined pMGs (Supplementary Fig. 5a) is supported by phylogenetic trees built from all one-to-one orthologs.



Phylogenetic tree of the *Escherichia/Shigella* clade based on the 40 universal single copy phylogenetic marker genes (pMGs). The polyphyly of *Shigella* reported earlier is clearly visible. Additionally, the tree shows that the phylogenetic distance between *Escherichia fergusonii* and the closest *Escherichia coli* genome are phylogenetically closer related than some of the *Escherichia coli* genomes, explaining why *Escherichia fergusonii* and *Escherichia coli* are found to be in the same cluster by our method.

Supplementary Figure 8: Evaluation of the accuracy of genomic placements for spec



Evaluation of the accuracy of genomic placements for spec based on 100 recently sequenced genomes (also see Supplementary Table 19 for details on the discrepancies).



Supplementary Table 1

An example of high intra-species 16S rRNA gene divergences. All pair-wise divergences between different 16S rRNA gene copies of *Desulfitobacterium hafniense* Y51. All values above 3% are marked in red as this represents the classical 97% sequence identity cutoff used to delineate species.

Nucleotide divergence between 16S rRNA gene copies	138119.DSY_16SrRNA1	138119.DSY_16SrRNA2	138119.DSY_16SrRNA3	138119.DSY_16SrRNA4	138119.DSY_16SrRNA5	138119.DSY_16SrRNA6
138119.DSY_16SrRNA1	0.00%	0.38%	0.06%	4.65%	4.76%	5.20%
138119.DSY_16SrRNA2	0.38%	0.00%	0.32%	4.78%	4.64%	4.95%
138119.DSY_16SrRNA3	0.06%	0.32%	0.00%	4.71%	4.70%	5.28%
138119.DSY_16SrRNA4	4.65%	4.78%	4.71%	0.00%	4.42%	0.56%
138119.DSY_16SrRNA5	4.76%	4.64%	4.70%	4.42%	0.00%	4.59%
138119.DSY_16SrRNA6	5.20%	4.95%	5.28%	0.56%	4.59%	0.00%

*Desulfitobacterium hafniense* Y51 (NCBI Taxonomy ID: 138119)

Supplementary Table 2

An example of extremely low inter-species 16S rRNA gene divergences.

The rows represent different copies of the 16S rRNA gene in *Aeromonas salmonicida* subsp. *salmonicida* A449 while the columns represent different copies of the 16S rRNA gene in *Aeromonas hydrophila* subsp. *hydrophila* ATCC 7966

Nucleotide divergence between 16S rRNA gene copies	380703.AHA_0077	380703.AHA_0145	380703.AHA_0198	380703.AHA_0335	380703.AHA_0752	380703.AHA_0865	380703.AHA_3516	380703.AHA_3944	380703.AHA_4016	380703.AHA_4246
382245.ASA_0079	1.23%	1.23%	1.30%	1.23%	1.23%	1.23%	1.30%	1.23%	1.23%	1.23%
382245.ASA_0295	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%
382245.ASA_0801	0.96%	0.96%	1.02%	0.96%	0.96%	0.96%	1.02%	0.96%	0.96%	0.96%
382245.ASA_3425	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%
382245.ASA_3599	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%
382245.ASA_4060	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%
382245.ASA_4192	1.23%	1.23%	1.30%	1.23%	1.23%	1.23%	1.30%	1.23%	1.23%	1.23%
382245.ASA_4244	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%	1.16%	1.09%	1.09%	1.09%
382245.ASA_4317	1.16%	1.16%	1.23%	1.16%	1.16%	1.16%	1.23%	1.16%	1.16%	1.16%

*Aeromonas salmonicida* subsp. *salmonicida* A449 (NCBI Taxonomy ID: 382245)

*Aeromonas hydrophila* subsp. *hydrophila* ATCC 7966 (NCBI Taxonomy ID: 380703)

Supplementary Table 3

List of the 40 universal single copy phylogenetic marker genes used in this study.

COG0012	Predicted GTPase, probable translation factor
COG0016	Phenylalanyl-tRNA synthetase alpha subunit
COG0018	Arginyl-tRNA synthetase
COG0048	Ribosomal protein S12
COG0049	Ribosomal protein S7
COG0052	Ribosomal protein S2
COG0080	Ribosomal protein L11
COG0081	Ribosomal protein L1
COG0085	DNA-directed RNA polymerase, beta subunit/140 kD subunit
COG0087	Ribosomal protein L3
COG0088	Ribosomal protein L4
COG0090	Ribosomal protein L2
COG0091	Ribosomal protein L22
COG0092	Ribosomal protein S3
COG0093	Ribosomal protein L14
COG0094	Ribosomal protein L5
COG0096	Ribosomal protein S8
COG0097	Ribosomal protein L6P/L9E
COG0098	Ribosomal protein S5
COG0099	Ribosomal protein S13
COG0100	Ribosomal protein S11
COG0102	Ribosomal protein L13
COG0103	Ribosomal protein S9
COG0124	Histidyl-tRNA synthetase
COG0172	Seryl-tRNA synthetase
COG0184	Ribosomal protein S15P/S13E
COG0185	Ribosomal protein S19
COG0186	Ribosomal protein S17
COG0197	Ribosomal protein L16/L10E
COG0200	Ribosomal protein L15
COG0201	Preprotein translocase subunit SecY
COG0202	DNA-directed RNA polymerase, alpha subunit/40 kD subunit
COG0215	Cysteinyl-tRNA synthetase
COG0256	Ribosomal protein L18
COG0495	Leucyl-tRNA synthetase
COG0522	Ribosomal protein S4 and related proteins
COG0525	Valyl-tRNA synthetase
COG0533	Metal-dependent proteases with possible chaperone activity
COG0541	Signal recognition particle GTPase
COG0552	Signal recognition particle GTPase

Supplementary Table 5

Discrepancies found between the species clustering and the type strain taxonomy. All species for which a sequenced type strain could be found were included. Please note that some of the discrepancies could later be curated using literature information (see Supplementary Table 16 for details)

Type of Discrepancy	NCBI Taxonomy ID	Genome Project ID	NCBI Taxonomy Name	Cluster Majority Taxonomy Name	Cluster ID
Split	491916	59115	Rhizobium etli	Rhizobium etli	682
Merged	590168	42777	Thermotoga naphthophila	Thermotoga petrophila	1224
Split	682634	42253	Serratia odorifera	Serratia odorifera	373
Split	525919	59219	Anaerococcus prevotii	Anaerococcus prevotii	1501
Split	1037409	67463	Bradyrhizobium japonicum	Bradyrhizobium japonicum	699
Merged	888063	67395	Enterobacter hormaechei	Enterobacter cloacae	384
Split	718254	45967	Enterobacter cloacae	Enterobacter cloacae	386
Split	701347	59969	Enterobacter cloacae	Enterobacter cloacae	386
Split	716541	48363	Enterobacter cloacae	Enterobacter cloacae	386
Split	717608	45855	Clostridium saccharolyticum	Clostridium saccharolyticum	1602
Merged	583358	49481	Thermoanaerobacter mathranii	Thermoanaerobacter italicus	1524
Merged	509193	55639	Thermoanaerobacter brockii	Thermoanaerobacter pseudethanolicus	1525
Split	866774	51527	Atopobium vaginae	Atopobium vaginae	810
Split	706437	61277	Streptococcus anginosus	Streptococcus anginosus	1393
Split	521097	59197	Capnocytophaga ochracea	Capnocytophaga ochracea	1121
Merged	997348	70617	Neisseria macacae	Neisseria sicca	437
Merged	435590	58253	Bacteroides vulgatus	Bacteroides dorei	1104
Merged	702446	47771	Bacteroides vulgatus	Bacteroides dorei	1104
Split	883	59089	Desulfovibrio vulgaris	Desulfovibrio vulgaris	754
Split	349102	58451	Rhodobacter sphaeroides	Rhodobacter sphaeroides	601
Split	391616	54699	Octadecabacter antarcticus	Octadecabacter antarcticus	635
Split	452638	58967	Polynucleobacter necessarius	Polynucleobacter necessarius	466
Merged	632516	60575	Caldicellulosiruptor lactoaceticus	Caldicellulosiruptor kristjanssonii	1511
Split	393480	54419	Fusobacterium nucleatum	Fusobacterium nucleatum	1479
Merged	566552	55369	Bifidobacterium catenulatum	Bifidobacterium pseudocatenulatum	973
Split	331678	58131	Chlorobium phaeobacteroides	Chlorobium phaeobacteroides	1046
Merged	868595	67317	Desulfotomaculum carboxydivorans	Desulfotomaculum nigrificans	1680
Split	637911	55947	Actinobacillus minor	Actinobacillus minor	310
Merged	526997	55207	Bacillus mycooides	Bacillus cereus	1337
Split	553190	43211	Gardnerella vaginalis	Gardnerella vaginalis	965
Split	682147	49423	Gardnerella vaginalis	Gardnerella vaginalis	965
Split	682148	49673	Gardnerella vaginalis	Gardnerella vaginalis	965
Split	457424	55575	Bacteroides fragilis	Bacteroides fragilis	1090
Split	927666	65449	Streptococcus oralis	Streptococcus oralis	1410
Split	871237	59469	Streptococcus infantis	Streptococcus infantis	1415
Split	1005705	67809	Streptococcus infantis	Streptococcus infantis	1415
Split	1035189	67191	Streptococcus infantis	Streptococcus infantis	1415
Merged	340047	58525	Mycoplasma capricolum	Mycoplasma leachii	1252
Merged	440085	58933	Methylobacterium chloromethanicum	Methylobacterium extorquens	710
Split	1048834	68273	Alicyclobacillus acidocaldarius	Alicyclobacillus acidocaldarius	1663
Merged	281309	58089	Bacillus thuringiensis	Bacillus cereus	1338
Merged	412694	58795	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527019	55239	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527021	55223	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527022	55215	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527024	55219	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527025	55225	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527028	55237	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527029	55233	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527030	55235	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527031	55229	Bacillus thuringiensis	Bacillus cereus	1338
Merged	527032	55231	Bacillus thuringiensis	Bacillus cereus	1338
Merged	541229	43737	Bacillus thuringiensis	Bacillus cereus	1338
Merged	714359	49135	Bacillus thuringiensis	Bacillus cereus	1338
Merged	930170	60447	Bacillus thuringiensis	Bacillus cereus	1338
Split	699184	50375	Bacillus thuringiensis	Bacillus cereus	1338
Split	526972	55161	Bacillus thuringiensis	Bacillus cereus	1338
Split	526990	55197	Bacillus thuringiensis	Bacillus cereus	1338
Split	526976	55169	Bacillus thuringiensis	Bacillus cereus	1338
Split	526988	55193	Bacillus thuringiensis	Bacillus cereus	1338
Split	579138	68445	Zymomonas mobilis	Zymomonas mobilis	580
Split	864569	52359	Streptococcus equinus	Streptococcus equinus	1381
Split	451756	54851	Clostridium perfringens	Clostridium perfringens	1541
Split	451757	54853	Clostridium perfringens	Clostridium perfringens	1541
Merged	632348	60491	Caldicellulosiruptor kronotskyensis	Caldicellulosiruptor bescii	1510
Split	246199	47357	Ruminococcus albus	Ruminococcus albus	1566
Merged	1032845	70839	Rickettsia heilongjiangensis	Rickettsia rickettsii	1747
Merged	272951	54113	Rickettsia sibirica	Rickettsia rickettsii	1747
Merged	652620	73963	Rickettsia japonica	Rickettsia rickettsii	1747
Merged	698757	42375	Pyrobaculum oguniense	Pyrobaculum arsenaticum	40
Split	562981	66131	Gemella haemolysans	Gemella haemolysans	1299

Supplementary Table 6

Genomes that were recently made available and are not part of the 3,496 genomes in the species clustering that were not classified by specl according to their NCBI species name (also see Supplementary Figure 8).

NCBI taxonomy ID	NCBI species name	specl assignment	Average pMG sequence identity to most closely related species cluster
1004785	<i>Alteromonas macleodii</i>	UNASSIGNED	92.4038405
1123519	<i>Pseudomonas stutzeri</i>	UNASSIGNED	89.36955966
1185652	<i>Sinorhizobium fredii</i>	UNASSIGNED	94.00208992
1202538	<i>Candidatus Carsonella ruddii</i>	UNASSIGNED	<80%
1224746	<i>Gluconobacter oxydans</i>	UNASSIGNED	84.39660304
245012	butyrate-producing bacterium SM4/1	<i>Clostridium saccharolyticum</i>	99.67464208
670307	<i>Hyphomicrobium denitrificans</i>	UNASSIGNED	88.82299904
710686	<i>Mycobacterium smegmatis</i>	UNASSIGNED	86.28364602

Supplementary Table 7

Type strain co-clustering: in this table all co-clustering type strains are listed together with their ANI values as well as their combined 40 pMG percent nucleotide identity, showing that the vast majority of incongruencies between the spec clustering and the type strain taxonomy are supported ANI.

Type of error	Cluster ID (also see table 14)	NCBI Taxonomy ID type strain 1	NCBI Taxonomy species name type strain 1	NCBI Taxonomy name type strain 1	NCBI Taxonomy ID type strain 2	NCBI Taxonomy species name type strain 2	NCBI Taxonomy name type strain 2	ANI	ANI	Combined 40 pMG % nucleotide identity
Merge	1224	390168	<i>Thermoplasma raphanophila</i>	<i>Thermoplasma raphanophila</i> RKU-10	390874	<i>Thermoplasma penophila</i>	<i>Thermoplasma penophila</i> RKU-1	98.39	97.13	0.987144705
Merge	1224	58333	<i>Thermoterrivibrio thalassius</i>	<i>Thermoterrivibrio thalassius</i> ATCC 33223	58333	<i>Thermoterrivibrio thalassius</i>	<i>Thermoterrivibrio thalassius</i> DSM 11979	97.34	96.14	0.982195036
Merge	1225	439293	<i>Thermoterrivibrio rufocellulosus</i>	<i>Thermoterrivibrio rufocellulosus</i> ATCC 33223	58333	<i>Thermoterrivibrio thalassius</i>	<i>Thermoterrivibrio thalassius</i> DSM 11979	95.65	96.71	0.992830079
Merge (supported by literature)	437	691748	<i>Nassaria macgregori</i>	<i>Nassaria macgregori</i> ATCC 33982	541243	<i>Nassaria sicca</i>	<i>Nassaria sicca</i> ATCC 29256	96.05	96.71	0.972498263
Merge	1104	439293	<i>Bacteroides vulgatus</i>	<i>Bacteroides vulgatus</i> ATCC 25462	483217	<i>Bacteroides dentis</i>	<i>Bacteroides dentis</i> DSM 17805	94.03	95.85	0.985169284
Merge	1811	632516	<i>Cardiobacterium lactoaceticus</i>	<i>Cardiobacterium lactoaceticus</i> EA	632338	<i>Cardiobacterium kristianssoni</i>	<i>Cardiobacterium kristianssoni</i> 177R18	98.13	98.22	0.989202005
Merge	1974	547434	<i>Haloproducium pseudocitricum</i>	<i>Haloproducium pseudocitricum</i> DSM 20438 – JCM 1201	666602	<i>Haloproducium californicum</i>	<i>Haloproducium californicum</i> DSM 11990 – JCM 1184	91.96	92.84	0.987420205
Merge	1680	698369	<i>Ostreobacterium nigrificans</i>	<i>Ostreobacterium nigrificans</i> DSM 674	688859	<i>Ostreobacterium carboxyliferans</i>	<i>Ostreobacterium carboxyliferans</i> CD-1-SRB	99.09	99.06	0.974308554
Merge (supported by literature)	1253	34054	<i>Methylobacterium capillatum</i>	<i>Methylobacterium capillatum</i> ATCC 23243 sub605	691447	<i>Methylobacterium lealii</i>	<i>Methylobacterium lealii</i> PCC 6806	93.25	94.53	0.986942317
Merge	710	44056	<i>Methylobacterium chthonotharsicum</i>	<i>Methylobacterium chthonotharsicum</i> DMM 1461617	661410	<i>Methylobacterium entosquans</i>	<i>Methylobacterium entosquans</i> CMA 1488933	96.49	96.91	0.980511344
Merge (supported by literature)	1338	228200	<i>Bacillus cereus</i>	<i>Bacillus cereus</i> ATCC 14879	671011	<i>Bacillus thuringiensis</i>	<i>Bacillus thuringiensis</i> serovar Berlin ATCC 10792	96.63	96.73	0.979302999
Merge	1910	614463	<i>Campylobacter jejuni</i>	<i>Campylobacter jejuni</i> DSM 6726	632344	<i>Campylobacter homocitricus</i>	<i>Campylobacter homocitricus</i> 2009	95.93	96.14	0.974984753
Merge (supported by literature)	40	340102	<i>Pyrobaculum armeniacum</i>	<i>Pyrobaculum armeniacum</i> DSM 13514	666797	<i>Pyrobaculum ogarensis</i>	<i>Pyrobaculum ogarensis</i> TE7	95.4	96.02	0.986339392
Merge (supported by literature)	1747	1032842	<i>Rickettsia helvetica</i>	<i>Rickettsia helvetica</i> DSM 654	650262	<i>Rickettsia agona</i>	<i>Rickettsia agona</i> strain YH	99.16	99.04	0.989341186
Merge (supported by literature)	1747	1032842	<i>Rickettsia helvetica</i>	<i>Rickettsia helvetica</i> DSM 654	27295	<i>Rickettsia sibirica</i>	<i>Rickettsia sibirica</i> 246	97.4	97.09	0.988646231
Merge (supported by literature)	1747	1032842	<i>Rickettsia helvetica</i>	<i>Rickettsia helvetica</i> DSM 654	650262	<i>Rickettsia rickettsii</i>	<i>Rickettsia rickettsii</i> str. Sheila Smith	97.77	96.88	0.988725217
Merge (supported by literature)	1747	652620	<i>Rickettsia japonica</i>	<i>Rickettsia japonica</i> strain YH	27295	<i>Rickettsia sibirica</i>	<i>Rickettsia sibirica</i> 246	97.54	97.18	0.988971202
Merge (supported by literature)	1747	652620	<i>Rickettsia japonica</i>	<i>Rickettsia japonica</i> strain YH	36202	<i>Rickettsia rickettsii</i>	<i>Rickettsia rickettsii</i> str. Sheila Smith	97.31	96.84	0.986705217
Merge (supported by literature)	1747	27295	<i>Rickettsia sibirica</i>	<i>Rickettsia sibirica</i> 246	36202	<i>Rickettsia rickettsii</i>	<i>Rickettsia rickettsii</i> str. Sheila Smith	96.18	97.89	0.989711344
Merge (supported by literature)	667	620466	<i>Buella racionae</i>	<i>Buella racionae</i> SK-33	204272	<i>Buella sus</i>	<i>Buella sus</i> 1330	96.75	99.71	0.998392072
Merge (supported by literature)	667	620466	<i>Buella racionae</i>	<i>Buella racionae</i> SK-33	444178	<i>Buella ovis</i>	<i>Buella ovis</i> ATCC 25840	99.03	99.4	0.997792305
Merge (supported by literature)	667	620466	<i>Buella racionae</i>	<i>Buella racionae</i> SK-33	483178	<i>Buella carsi</i>	<i>Buella carsi</i> ATCC 23365	99.73	99.69	0.998284046
Merge (supported by literature)	667	620466	<i>Buella racionae</i>	<i>Buella racionae</i> SK-33	233894	<i>Buella melitensis</i>	<i>Buella melitensis</i> bv. 1 str. 1684	99.86	99.82	0.997732763
Merge (supported by literature)	667	620466	<i>Buella racionae</i>	<i>Buella racionae</i> SK-33	568818	<i>Buella micropi</i>	<i>Buella micropi</i> CCM 4815	99.8	99.77	0.998870204
Merge (supported by literature)	667	204272	<i>Buella sus</i>	<i>Buella sus</i> 1330	444178	<i>Buella ovis</i>	<i>Buella ovis</i> ATCC 25840	99.63	99.4	0.998431515
Merge (supported by literature)	667	204272	<i>Buella sus</i>	<i>Buella sus</i> 1330	483178	<i>Buella carsi</i>	<i>Buella carsi</i> ATCC 23365	99.9	99.87	0.999077634
Merge (supported by literature)	667	204272	<i>Buella sus</i>	<i>Buella sus</i> 1330	224814	<i>Buella melitensis</i>	<i>Buella melitensis</i> bv. 1 str. 16M	99.88	99.82	0.998168548
Merge (supported by literature)	667	204272	<i>Buella sus</i>	<i>Buella sus</i> 1330	568818	<i>Buella micropi</i>	<i>Buella micropi</i> CCM 4815	99.8	99.77	0.999216098
Merge (supported by literature)	667	444178	<i>Buella ovis</i>	<i>Buella ovis</i> ATCC 25840	483178	<i>Buella carsi</i>	<i>Buella carsi</i> ATCC 23365	99.62	99.57	0.997922928
Merge (supported by literature)	667	444178	<i>Buella ovis</i>	<i>Buella ovis</i> ATCC 25840	233894	<i>Buella melitensis</i>	<i>Buella melitensis</i> bv. 1 str. 16M	99.56	99.51	0.997411783
Merge (supported by literature)	667	444178	<i>Buella ovis</i>	<i>Buella ovis</i> ATCC 25840	568818	<i>Buella micropi</i>	<i>Buella micropi</i> CCM 4815	99.87	99.85	0.998936398
Merge (supported by literature)	667	483178	<i>Buella carsi</i>	<i>Buella carsi</i> ATCC 23365	233894	<i>Buella melitensis</i>	<i>Buella melitensis</i> bv. 1 str. 16M	99.64	99.61	0.997777050
Merge (supported by literature)	667	483178	<i>Buella carsi</i>	<i>Buella carsi</i> ATCC 23365	568818	<i>Buella micropi</i>	<i>Buella micropi</i> CCM 4815	99.78	99.76	0.998262845
Merge (supported by literature)	667	224814	<i>Buella melitensis</i>	<i>Buella melitensis</i> bv. 1 str. 16M	568818	<i>Buella micropi</i>	<i>Buella micropi</i> CCM 4815	99.71	99.68	0.998279471

Supplementary Table 15

Summary of the number of pMGs and 16S rRNA gene copies found per genome.

Number of pMGs	Number of Genomes with this number of pMGs	Number of Genomes with at least this number of pMGs	Percentage of genomes with at least this number of pMGs
40	2860	2860	0.818077803
39	320	3180	0.909610984
38	118	3298	0.943363844
37	61	3359	0.960812357
36	28	3387	0.96882151
35	21	3408	0.974828375
34	18	3426	0.979977117
33	23	3449	0.986556064
32	14	3463	0.990560641
31	9	3472	0.993135011
30	10	3482	0.995995423
29	2	3484	0.996567506
27	2	3486	0.997139588
26	2	3488	0.99771167
24	1	3489	0.997997712
22	1	3490	0.998283753
21	1	3491	0.998569794
18	1	3492	0.998855835
17	1	3493	0.999141876
13	1	3494	0.999427918
11	1	3495	0.999713959
10	1	3496	1

Average Number of pMGs per genome

39.49141876

Genomes with full length 16S rRNA gene sequence

2899

Total full length 16S rRNA gene copies sequences found

9435

Average Number of 16S rRNA gene copies

3.254570542

Supplementary Table 17

Genome quality statistics and the number of found pMGs for all genomes from the clades used for the examples showing in Figure 2 and Supplementary Figures 5 and 6.

NCBI taxonomy name	NCBI Taxonomy ID	NCBI genome project ID	N50	Number of contigs	Number of pMGs found	Number of predicted genes	Comments
<i>Prochlorococcus marinus</i> str. MIT 9202	93058	54709	1691453	1	39	1927	
<i>Prochlorococcus marinus</i> subsp. <i>pastoris</i> str. CCMP1986	59919	57761	1657990	1	40	1760	
<i>Prochlorococcus marinus</i> str. MIT 9313	74547	57773	2410873	1	40	2324	
<i>Prochlorococcus marinus</i> subsp. <i>marinus</i> str. CCMP1375	167539	57995	1751080	1	40	1929	
<i>Prochlorococcus marinus</i> str. MIT 9303	59922	58305	2682675	1	40	3049	
<i>Prochlorococcus marinus</i> str. AS9601	146891	58307	1669886	1	40	1963	
<i>Prochlorococcus marinus</i> str. MIT 9211	93059	58309	1688963	1	40	1897	
<i>Prochlorococcus marinus</i> str. MIT 9515	167542	58313	1704176	1	40	1947	
<i>Prochlorococcus marinus</i> str. MIT 9312	74546	58357	1709204	1	40	1854	
<i>Prochlorococcus marinus</i> str. NATL2A	59920	58359	1842899	1	40	2205	
<i>Prochlorococcus marinus</i> str. NATL1A	167555	58423	1864731	1	40	2236	
<i>Prochlorococcus marinus</i> str. MIT 9301	167546	58437	1641879	1	40	1948	
<i>Prochlorococcus marinus</i> str. MIT 9215	93060	58819	1738790	1	40	2024	
<i>Synechococcus</i> sp. CC9311	64471	58123	2606748	1	40	2944	
<i>Synechococcus</i> sp. WH 7803	32051	61607	2366980	1	40	2586	
<i>Serratia odorifera</i> 4Rx13	682634	42253	457543	17	40	4741	
<i>Serratia odorifera</i> DSM 4582	667129	46991	170690	91	37	5194	
<i>Serratia proteamaculans</i> 568	399741	58725	5448853	1	40	5000	
<i>Serratia</i> sp. AS13	768493	60455	5442549	1	40	5060	
<i>Serratia symbiotica</i> str. Tucson	914128	62293	201972	389	0	2156	*
<i>Serratia</i> sp. AS9	768492	67313	5442880	1	40	5061	
<i>Serratia</i> sp. AS12	768490	67315	5443009	1	40	5061	
<i>Serratia symbiotica</i> str. 'Cinara cedri'	568817	82363	1762765	1	40	769	
<i>Rahnella aquatilis</i> CIP 78.65 = ATCC 33071	745277	60227	4861101	1	40	4442	
<i>Rahnella</i> sp. Y9602	741091	62715	4864217	1	40	4490	
<i>Wigglesworthia glossinidia</i> endosymbiont of <i>Glossina brevipalpa</i>	36870	57853	697724	1	40	651	
<i>Wigglesworthia glossinidia</i> endosymbiont of <i>Glossina morsitans</i>	1142511	82191	719535	1	40	676	
<i>Buchnera aphidicola</i> str. LL01 ( <i>Acyrtosiphon pisum</i> )	713603	43505	641799	1	39	612	
<i>Buchnera aphidicola</i> str. TLW03 ( <i>Acyrtosiphon pisum</i> )	713602	43507	641770	1	40	608	
<i>Buchnera aphidicola</i> str. JF99 ( <i>Acyrtosiphon pisum</i> )	713601	43509	641716	1	38	625	
<i>Buchnera aphidicola</i> str. JF98 ( <i>Acyrtosiphon pisum</i> )	713600	43511	641771	1	32	512	
<i>Buchnera aphidicola</i> str. LSR1 ( <i>Acyrtosiphon pisum</i> )	593275	55847	642011	1	40	629	
<i>Buchnera aphidicola</i> str. APS ( <i>Acyrtosiphon pisum</i> )	107806	57805	640681	1	40	600	
<i>Buchnera aphidicola</i> str. Bp ( <i>Baizongia pistaciae</i> )	224915	57827	615980	1	40	540	
<i>Buchnera aphidicola</i> str. Sg ( <i>Schizaphis graminum</i> )	198804	57913	641454	1	40	582	
<i>Buchnera aphidicola</i> BCC	372461	58579	416380	1	40	393	
<i>Buchnera aphidicola</i> str. Tuc7 ( <i>Acyrtosiphon pisum</i> )	561501	59283	641895	1	40	588	
<i>Buchnera aphidicola</i> str. 5A ( <i>Acyrtosiphon pisum</i> )	563178	59285	642122	1	40	590	
<i>Buchnera aphidicola</i> str. Ua ( <i>Uroleucon ambrosiae</i> )	1005057	65479	615380	1	40	564	
<i>Buchnera aphidicola</i> str. Ak ( <i>Acyrtosiphon kondoi</i> )	1005090	65481	641794	1	40	594	
<i>Buchnera aphidicola</i> ( <i>Cinara tujaefilina</i> )	261317	68101	444925	1	40	394	

\* = Was excluded before the number of marker genes was computed as it failed the initial genome quality control



Supplementary Table 19

Clusters of size two or more containing only not fully identified genomes and a list of the genomes falling into these clusters.

Cluster	NCBI Taxonomy ID	NCBI Taxonomy ID NCBI Project ID	NCBI Taxonomy Name of sequenced organi
1	395962	395962.5914	Cyanothece sp. PCC 8802
1	41431	41431.59027	Cyanothece sp. PCC 8801
2	1148	1148.57659	Synechocystis sp. PCC 6803
2	1148	1148.67081	Synechocystis sp. PCC 6803
3	60480	60480.58345	Shewanella sp. MR-4
3	60481	60481.58343	Shewanella sp. MR-7
3	94122	94122.58347	Shewanella sp. ANA-3
4	675810	675810.4111	Vibrio sp. RC341
4	675815	675815.4162	Vibrio sp. RC586
5	314267	314267.5426	Sulfitobacter sp. NAS-14.1
5	52598	52598.54191	Sulfitobacter sp. EE-36
6	439495	439495.5471	Pseudovibrio sp. JE062
6	911045	911045.8237	Pseudovibrio sp. FO-BEG1
7	164756	164756.5847	Mycobacterium sp. MCS
7	164757	164757.5849	Mycobacterium sp. JLS
7	189918	189918.5849	Mycobacterium sp. KMS
8	355249	355249.6692	Streptomyces sp. Tu6071
8	465543	465543.4842	Streptomyces sp. SPB74
8	591157	591157.5582	Streptomyces sp. SPB78
8	683219	683219.619	Streptomyces sp. SA3_actG
9	1000568	1000568.678	Megasphaera sp. UPII 199-6
9	699218	699218.4657	Megasphaera genomsp. type_1
10	481743	481743.4113	Paenibacillus sp. Y412MC10
10	908341	908341.6615	Paenibacillus sp. HGF5
11	457403	457403.682	Fusobacterium sp. 11_3_2
11	457405	457405.5561	Fusobacterium sp. 7_1
11	469601	469601.6819	Fusobacterium sp. 21_1A
11	469603	469603.4085	Fusobacterium sp. 3_1_33
11	556264	556264.5563	Fusobacterium sp. D11
12	469602	469602.4779	Fusobacterium sp. 3_1_27
12	469604	469604.56	Fusobacterium sp. 3_1_36A2
12	469606	469606.5561	Fusobacterium sp. 4_1_13
13	1029718	1029718.714	Candidatus Arthromitus sp. SFB-mouse
13	1041809	1041809.678	Candidatus Arthromitus sp. SFB-mouse
14	245018	245018.4596	butyrate-producing bacterium SSC/2
14	411484	411484.5455	Clostridium sp. SS2/1
14	658089	658089.6188	Lachnospiraceae bacterium 5_1_63FAA
15	552395	552395.6358	Lachnospiraceae bacterium 4_1_37FAA
15	658088	658088.6642	Lachnospiraceae bacterium 9_1_43BFAA

Supplementary Table 20

The 130 genomes used for the benchmarking shown in Figure 1b.

NCBI Taxonomy ID	Genome Project ID	NCBI Taxonomy Name
1009464	51067	<i>Gardnerella vaginalis</i> HMP9231
1010838	66075	<i>Mycobacterium avium</i> subsp. <i>paratuberculosis</i> S397
1027843	66485	<i>Chlamydomphila psittaci</i> 08DC60
1027844	66484	<i>Chlamydomphila psittaci</i> C18/98
1027845	66481	<i>Chlamydomphila psittaci</i> 02DC15
1027846	66479	<i>Chlamydomphila psittaci</i> 01DC11
1029822	66655	<i>Lactobacillus salivarius</i> NIAS840
1031710	66839	<i>Oligotropha carboxidovorans</i> OM4
1039817	66401	<i>Bifidobacterium longum</i> subsp. <i>longum</i> KACC 91563
1036743	67474	<i>Rhodospirillum rubrum</i> F11
1037409	67463	<i>Bradyrhizobium japonicum</i> USDA 6
1041521	67679	<i>Lactobacillus salivarius</i> GJ-24
1042403	67865	<i>Bifidobacterium animalis</i> subsp. <i>lactis</i> CNCM I-2494
1045856	80739	<i>Enterobacter cloacae</i> EcWSU1
1048834	68273	<i>Alicyclobacillus acidocaldarius</i> subsp. <i>acidocaldarius</i> Tc-4-1
1049565	67293	<i>Klebsiella pneumoniae</i> KC TC 2242
1051072	68327	<i>Streptococcus equi</i> subsp. <i>zooepidemicus</i> ATCC 35246
1075106	71815	<i>Bifidobacterium animalis</i> subsp. <i>lactis</i> BLC1
1107890	75083	<i>Leuconostoc mesenteroides</i> subsp. <i>mesenteroides</i> J18
1112856	85495	<i>Mycoplasma pneumoniae</i> 309
1125630	78789	<i>Klebsiella pneumoniae</i> subsp. <i>pneumoniae</i> HS11286
349123	54165	<i>Lactobacillus reuteri</i> 100-23
351581	54341	<i>Francisella tularensis</i> subsp. <i>holarctica</i> FSC200
353496	16120	<i>Lactobacillus delbrueckii</i> subsp. <i>bulgaricus</i> 2038
406984	17947	<i>Chlamydomphila pneumoniae</i> LPCoLN
436113	19245	<i>Mycoplasma mycoides</i> subsp. <i>capri</i> str. GM12
436113	39245	<i>Mycoplasma mycoides</i> subsp. <i>capri</i> str. GM12
445334	54845	<i>Clostridium perfringens</i> C str. JGS1495
445983	54817	<i>Borrelia burgdorferi</i> 156a
445984	54819	<i>Borrelia burgdorferi</i> B026
451754	54847	<i>Clostridium perfringens</i> B str. ATCC 3626
451755	54849	<i>Clostridium perfringens</i> E str. JGS1987
451757	54853	<i>Clostridium perfringens</i> NCTC 8239
476210	54835	<i>Borrelia burgdorferi</i> 118a
476211	54837	<i>Borrelia burgdorferi</i> 72a
476212	54833	<i>Borrelia burgdorferi</i> 94a
483214	54983	<i>Methanobrevibacter smithii</i> DSM 2375
486409	55085	<i>Bifidobacterium animalis</i> subsp. <i>lactis</i> HN019
488537	55051	<i>Clostridium perfringens</i> D str. JGS1721
492476	55013	<i>Clostridium thermocellum</i> JW20
491978	28333	<i>Acinetobacter baumannii</i> MDR-ZJ06
498736	55057	<i>Borrelia burgdorferi</i> 80a
498737	55061	<i>Borrelia burgdorferi</i> W019-23
498738	55055	<i>Borrelia burgdorferi</i> 29805
498739	55063	<i>Borrelia burgdorferi</i> CA-11.2A
498740	55067	<i>Borrelia burgdorferi</i> 64b
510831	38289	<i>Francisella tularensis</i> subsp. <i>tularensis</i> NE061598
515609	54573	<i>Ureaplasma parvum</i> serovar 14 str. ATCC 33697
515611	54877	<i>Ureaplasma urealyticum</i> serovar 13 str. ATCC 33698
518603	54749	<i>Ureaplasma urealyticum</i> serovar 5 str. ATCC 27817
519851	54747	<i>Ureaplasma parvum</i> serovar 6 str. ATCC 27818
521002	55123	<i>Methanobrevibacter smithii</i> DSM 2374
521007	29357	<i>Borrelia burgdorferi</i> N40
521008	29359	<i>Borrelia burgdorferi</i> JD1
525262	40867	<i>Cornebacterium jeikeium</i> ATCC 43734
525280	55483	<i>Erysipelothrix rhusiopathiae</i> ATCC 19414
525283	49043	<i>Fusobacterium nucleatum</i> subsp. <i>nucleatum</i> ATCC 23726
525325	55497	<i>Lactobacillus fermentum</i> ATCC 14931
525330	55505	<i>Lactobacillus johnsonii</i> ATCC 33200
525338	55521	<i>Lactobacillus plantarum</i> subsp. <i>plantarum</i> ATCC 14917
525341	55541	<i>Lactobacillus reuteri</i> CF48-3A
525369	55525	<i>Proteus mirabilis</i> ATCC 29906
526347	54751	<i>Ureaplasma urealyticum</i> serovar 11 str. ATCC 33695
543302	55305	<i>Alicyclobacillus acidocaldarius</i> LA1
548480	55465	<i>Bifidobacterium longum</i> subsp. <i>longum</i> ATCC 55813
548485	55517	<i>Lactobacillus reuteri</i> M14-1A
550694	38067	<i>Escherichia ferroussii</i> B253
550747	54753	<i>Ureaplasma urealyticum</i> serovar 12 str. ATCC 33696
550748	54647	<i>Ureaplasma urealyticum</i> serovar 4 str. ATCC 27816
550773	54651	<i>Ureaplasma urealyticum</i> serovar 9 str. ATCC 33175
553201	55453	<i>Rothia mucilaginosa</i> ATCC 25296
553481	55387	<i>Mycobacterium avium</i> subsp. <i>avium</i> ATCC 25291
555217	30987	<i>Zymomonas mobilis</i> subsp. <i>mobilis</i> ATCC 10988
557601	55405	<i>Acinetobacter baumannii</i> AB900
572545	55637	<i>Clostridium thermocellum</i> DSM 2360
573059	32603	<i>Desulfovibrio vulgaris</i> RCH1
573236	32515	<i>Bifidobacterium animalis</i> subsp. <i>lactis</i> V9
585198	51715	<i>Mobiluncus curtisi</i> subsp. <i>curtisi</i> ATCC 35241
585199	55885	<i>Lactobacillus reuteri</i> MM2-3
585520	55883	<i>Lactobacillus helveticus</i> DSM 20075
586220	55887	<i>Leuconostoc mesenteroides</i> subsp. <i>cremoris</i> ATCC 19254
592014	47501	<i>Acinetobacter baumannii</i> 6013113
637389	55945	<i>Acidithiobacillus caldus</i> ATCC 51756
637387	41553	<i>Clostridium thermocellum</i> DSM 1313
644651	37877	<i>Erwinia pyritifolia</i> DSM 12163
655813	47255	<i>Streptococcus oralis</i> ATCC 35037
657317	39161	<i>Eubacterium rectale</i> M104/1
657318	39159	<i>Eubacterium rectale</i> DSM 17629
663917	39745	<i>Geobacter sulfurreducens</i> KN400
667127	41361	<i>Klebsiella pneumoniae</i> subsp. <i>rhinoscleromatis</i> ATCC 13884
673196	42955	<i>Lactobacillus gasseri</i> 224-1
679936	40777	<i>Sulfobacillus acidophilus</i> DSM 10332
682147	49423	<i>Gardnerella vaginalis</i> AMD
682148	49673	<i>Gardnerella vaginalis</i> 5-1
686660	42953	<i>Veillonella parvula</i> ATCC 17745
696749	42153	<i>Acinetobacter baumannii</i> 1656-2
712861	43535	<i>Lactobacillus salivarius</i> CECT 5713
718254	45867	<i>Enterobacter cloacae</i> subsp. <i>cloacae</i> NCTC 9394
722911	45863	<i>Bifidobacterium longum</i> subsp. <i>longum</i> F8
762633	51619	<i>Thermus thermophilus</i> SGO.5JP17-16
767462	49153	<i>Lactobacillus helveticus</i> H10
768065	49335	<i>Halobacterium salinarum</i> R
768728	50761	<i>Lactobacillus salivarius</i> ACS-116-V-Col5a
862962	84217	<i>Bacteroides fragilis</i> 638R
863638	50455	<i>Clostridium acetobutylicum</i> EA 2018
869211	50823	<i>Sairochaeta thermophila</i> DSM 6578
873517	61279	<i>Capnocytophaga ochracea</i> F0287
876138	46683	<i>Streptococcus agalactiae</i> FSL S3-026
879307	52049	<i>Gardnerella vaginalis</i> 315-A
886886	52391	<i>Acinetobacter baumannii</i> AB210
887326	61495	<i>Mobiluncus curtisi</i> ATCC 51333
887899	61497	<i>Mobiluncus curtisi</i> subsp. <i>holmesii</i> ATCC 35242
888027	53001	<i>Lactobacillus delbrueckii</i> subsp. <i>lactis</i> DSM 20072
888048	62941	<i>Streptococcus parasanguinis</i> ATCC 903
888745	68679	<i>Streptococcus agalactiae</i> ATCC 13813
888826	63561	<i>Haemophilus parainfluenzae</i> ATCC 33392
905067	60563	<i>Streptococcus parasanguinis</i> F0405
907287	53881	<i>Mycoplasma hyopneumoniae</i> 168
909954	56053	<i>Lactobacillus johnsonii</i> DPC 6026
910312	61039	<i>Porphyromonas asaccharolytica</i> PR426713P-1
930944	62779	<i>Yersinia enterocolitica</i> subsp. <i>palearctica</i> Y11
980514	62779	<i>Acinetobacter baumannii</i> TCCD-AB0715
983328	62521	<i>Campylobacter fetus</i> subsp. <i>veneralis</i> NCTC 10354
984894	66151	<i>Chlamydomphila psittaci</i> Cal10
990315	63187	<i>Xanthomonas campestris</i> pv. <i>raphani</i> 756C
992402	63335	<i>Acinetobacter baumannii</i> ABNIH1
992403	63337	<i>Acinetobacter baumannii</i> ABNIH2
992404	63339	<i>Acinetobacter baumannii</i> ABNIH3
992405	63341	<i>Acinetobacter baumannii</i> ABNIH4
996285	63543	<i>Pseudomonas stutzeri</i> DSM 4166

Supplementary Note: Speciation of *Escherichia/Shigella*

*Escherichia/Shigella* complex is known to be polyphyletic<sup>19,38</sup> with *Shigella* genomes branching out in different regions of the phylogenetic tree of the *Escherichia/Shigella* complex (also see Supplementary Fig. 7) that forms a large cluster in the specI species clustering excluding *Escherichia albertii*, but including *Escherichia fergusonii*. *Escherichia coli* is one of the best studied organisms and its speciation has been studied thoroughly<sup>39</sup>. It seems clear that *Escherichia coli* has a conserved, ancient backbone<sup>39</sup>. This backbone was also recovered as the core genome in a recent analysis of the *Escherichia coli* pangenome<sup>40</sup>. In addition, to the core genome *Escherichia coli* seems to have a rather large dispensable genome<sup>40</sup>, supporting the notion that *Escherichia coli* seems to have an extensive mosaic structure<sup>39</sup>. More recently, there have also been claims that the *Escherichia/Shigella* complex is currently undergoing speciation into different environmental and host associated subgroups due to evolutionary events in the dispensable genome<sup>41</sup>. The same study detected a genetic continuum of the core genome within the *Escherichia* genus (more specifically between *E. coli* and *E. albertii*). In conclusion, our results coincide with the core genome hypothesis<sup>42</sup> as our clustering of this clade seems to coincide with the presence of a conserved core genome of the *Escherichia/Shigella* complex. The dispensable genome is not considered for species delineation by specI and if the dispensable genome has a strong impact on speciation, this might represent a limitation of the method. Analyses of the pan-genome of different species might help to explain the effect of the dispensable genome on speciation in the future.

38. Sims, G. & Kim, S.-H. Proceedings of the National Academy of Sciences of the United States of America 108, 8329–8334 (2011).
39. Welch, R. et al. Proceedings of the National Academy of Sciences of the United States of America 99, 17020–17024 (2002).
40. Rasko, D. et al. Journal of bacteriology 190, 6881–6893 (2008).
41. Luo, C. et al. Proceedings of the National Academy of Sciences of the United States of America 108, 7200–7205 (2011).
42. Riley, M. & Lizotte-Waniewski, M. Methods in molecular biology (Clifton, N.J.) 532, 367–377 (2009).