

**Supplementary Table 1:** List of organisms used for constructing simulated test data sets. For each test organism, the representation status\* with respect to modified reference database (modified RefDB) and the phylogenetic similarity status\*\* with respect to the genomes of training data set are also presented.

Sanger data set				
S. No.	Organism	Number of sequences	Status with respect to modified RefDB*	Status with respect to training data set**
1.	<i>Acaryochloris marina</i> MBIC11017	1000	All Known	Superkingdom Shared
2.	<i>Agrobacterium tumefaciens</i> str. C58	1000	Species Unknown	Class Shared
3.	<i>Aquifex aeolicus</i> VF5	1000	All Known	Superkingdom Shared
4.	<i>Bacillus halodurans</i> C-125	1000	Species Unknown	Phylum Shared
5.	<i>Bdellovibrio bacteriovorus</i> HD100	1000	Order Unknown	Phylum Shared
6.	<i>Bordetella pertussis</i> Tohama I	1000	Family Unknown	Phylum Shared
7.	<i>Bradyrhizobium</i> sp. BTAi1	1000	Species Unknown	Class Shared
8.	<i>Burkholderia cenocepacia</i> AU 1054	1000	All Known	Phylum Shared
9.	<i>Campylobacter jejuni</i> subsp. jejuni NCTC 11168	1000	Genus Unknown	Phylum Shared
10.	<i>Chlamydomonada pneumoniae</i> AR39	1000	Genus Unknown	Superkingdom Shared
11.	<i>Chlorobium chlorochromatii</i> CaD3	1000	Genus Unknown	Superkingdom Shared
12.	<i>Clostridium tetani</i> E88	1000	Species Unknown	Class Shared
13.	<i>Corynebacterium diphtheriae</i> NCTC 13129	1000	Species Unknown	Class Shared
14.	<i>Francisella philomiragia</i> subsp. philomiragia ATCC 25017	1000	Species Unknown	Class Shared
15.	<i>Haemophilus influenzae</i> Rd KW20	1000	Species Unknown	Class Shared
16.	<i>Helicobacter pylori</i> J99	1000	Strain Unknown	Phylum Shared
17.	<i>Lactobacillus acidophilus</i> NCFM	1000	Family Unknown	Phylum Shared
18.	<i>Leptospira interrogans</i> serovar Lai str. 56601	1000	Family Unknown	Superkingdom Shared
19.	<i>Mycoplasma pneumoniae</i> M129	1000	Phylum Unknown	Superkingdom Shared
20.	<i>Neisseria meningitidis</i> MC58	1000	Genus Unknown	Phylum Shared
21.	<i>Nostoc</i> sp. PCC 7120	1000	Order Unknown	Superkingdom Shared
22.	<i>Polaromonas</i> sp. JS666	1000	Species Unknown	Phylum Shared
23.	<i>Pseudomonas aeruginosa</i> PAO1	1000	Genus Unknown	Class Shared
24.	<i>Pseudomonas mendocina</i> ymp	1000	Genus Unknown	Class Shared
25.	<i>Rhizobium etli</i> CFN 42	1000	Species Unknown	Class Shared
26.	<i>Rhodospseudomonas palustris</i> BisA53	1000	Genus Unknown	Class Shared
27.	<i>Rickettsia bellii</i> RML369-C	1000	Genus Unknown	Class Shared
28.	<i>Saccharophagus degradans</i> 2-40	1000	Genus Unknown	Class Shared
29.	<i>Shewanella baltica</i> OS185	1000	Species Unknown	Class Shared
30.	<i>Staphylococcus aureus</i> subsp. aureus	1000	Species Unknown	Phylum Shared
31.	<i>Streptomyces coelicolor</i> A3(2)	1000	Species Unknown	Class Shared
32.	<i>Thermotoga maritima</i> MSB8	1000	Genus Unknown	Superkingdom Shared
33.	<i>Trichodesmium erythraeum</i> IMS101	1000	Order Unknown	Superkingdom Shared
34.	<i>Wigglesworthia glossinidia</i> endosymbiont of <i>Glossina brevipalpis</i>	1000	Genus Unknown	Class Shared
35.	<i>Xanthomonas campestris</i> pv. <i>campestris</i> str. ATCC 33913	1000	Genus Unknown	Class Shared

\*All Known: Genome fragments of that particular microbe is present in the modified RefDB.

\*Strain Unknown: Genome fragments of this strain are absent in the modified RefDB.

\*Species Unknown: Genome fragments of all strains belonging to this species are absent in the modified RefDB.

\*Genus Unknown: Genome fragments of all the strains and species belonging to this genus are absent in modified RefDB.

\*Family Unknown: Genome fragments of all genera belonging to this family are absent in the modified RefDB.

\*Order Unknown: Genome fragments of all families belonging to this order are absent in the modified RefDB.

\*Phylum Unknown: Genome fragments of all classes belonging to this phylum are absent in modified RefDB.

\*\*Class Shared: The training data set is devoid of genomes that can share either strain or species or genus or family or order with that particular microbe.

\*\*Phylum Shared: The training data set does not consists of any genomes that belong to either strain or species or genus or family or order or class of that particular microbe.

\*\*Superkingdom Shared: The training data set does not consists of any genomes that belong to either strain or species or genus or family or order or class or phylum of that particular microbe.

Supplementary Table 1: (contd)

454-400 and 454-250 data sets				
S. No.	Organism	Number of sequences	Status with respect to modified RefDB*	Status with respect to training data set**
1.	<i>Acidovorax citrulli</i> AAC00-1	1000	All Known	Phylum Shared
2.	<i>Agrobacterium tumefaciens</i> str. C58	1000	Species Unknown	Class Shared
3.	<i>Aquifex aeolicus</i> VF5	1000	All Known	Superkingdom Shared
4.	<i>Bacillus halodurans</i> C-125	1000	Species Unknown	Phylum Shared
5.	<i>Bdellovibrio bacteriovorus</i> HD100	1000	Order Unknown	Phylum Shared
6.	<i>Bordetella pertussis</i> Tohama I	1000	Family Unknown	Phylum Shared
7.	<i>Bradyrhizobium</i> sp. BTAi1	1000	Species Unknown	Class Shared
8.	<i>Burkholderia cenocepacia</i> AU 1054	1000	All Known	Phylum Shared
9.	<i>Campylobacter jejuni</i> subsp. <i>jejuni</i> NCTC 11168	1000	Genus Unknown	Phylum Shared
10.	<i>Chlamydomonas pneumoniae</i> AR39	1000	Genus Unknown	Superkingdom Shared
11.	<i>Chlorobium chlorochromatii</i> CaD3	1000	Genus Unknown	Superkingdom Shared
12.	<i>Clostridium tetani</i> E88	1000	Species Unknown	Class Shared
13.	<i>Corynebacterium diphtheriae</i> NCTC 13129	1000	Species Unknown	Class Shared
14.	<i>Francisella philomiragia</i> subsp. <i>philomiragia</i> ATCC 25017	1000	Species Unknown	Class Shared
15.	<i>Haemophilus influenzae</i> Rd KW20	1000	Species Unknown	Class Shared
16.	<i>Helicobacter pylori</i> J99	1000	Strain Unknown	Phylum Shared
17.	<i>Lactobacillus acidophilus</i> NCFM	1000	Family Unknown	Phylum Shared
18.	<i>Leptospira interrogans</i> serovar <i>Lai</i> str. 56601	1000	Family Unknown	Superkingdom Shared
19.	<i>Mycoplasma pneumoniae</i> M129	1000	Phylum Unknown	Superkingdom Shared
20.	<i>Neisseria meningitidis</i> MC58	1000	Genus Unknown	Phylum Shared
21.	<i>Nostoc</i> sp. PCC 7120	1000	Order Unknown	Superkingdom Shared
22.	<i>Polaromonas</i> sp. JS666	1000	Species Unknown	Phylum Shared
23.	<i>Pseudomonas aeruginosa</i> PAO1	1000	Genus Unknown	Class Shared
24.	<i>Pseudomonas mendocina</i> ymp	1000	Genus Unknown	Class Shared
25.	<i>Rhizobium etli</i> CFN 42	1000	Species Unknown	Class Shared
26.	<i>Rhodopseudomonas palustris</i> BisA53	1000	Genus Unknown	Class Shared
27.	<i>Rickettsia bellii</i> RML369-C	1000	Genus Unknown	Class Shared
28.	<i>Saccharophagus degradans</i> 2-40	1000	Genus Unknown	Class Shared
29.	<i>Shewanella baltica</i> OS185	1000	Species Unknown	Class Shared
30.	<i>Staphylococcus aureus</i> subsp. <i>aureus</i>	1000	Species Unknown	Phylum Shared
31.	<i>Streptomyces coelicolor</i> A3(2)	1000	Species Unknown	Class Shared
32.	<i>Thermotoga maritima</i> MSB8	1000	Genus Unknown	Superkingdom Shared
33.	<i>Trichodesmium erythraeum</i> IMS101	1000	Order Unknown	Superkingdom Shared
34.	<i>Wigglesworthia glossinidia</i> endosymbiont of <i>Glossina brevipalpis</i>	1000	Genus Unknown	Class Shared
35.	<i>Xanthomonas campestris</i> pv. <i>campestris</i> str. ATCC 33913	1000	Genus Unknown	Class Shared

\*All Known: Genome fragments of that particular microbe is present in the modified RefDB.

\*Strain Unknown: Genome fragments of this strain are absent in the modified RefDB.

\*Species Unknown: Genome fragments of all strains belonging to this species are absent in the modified RefDB.

\*Genus Unknown: Genome fragments of all the strains and species belonging to this genus are absent in modified RefDB.

\*Family Unknown: Genome fragments of all genera belonging to this family are absent in the modified RefDB.

\*Order Unknown: Genome fragments of all families belonging to this order are absent in the modified RefDB.

\*Phylum Unknown: Genome fragments of all classes belonging to this phylum are absent in modified RefDB.

\*\*Class Shared: The training data set is devoid of genomes that can share either strain or species or genus or family or order with that particular microbe.

\*\*Phylum Shared: The training data set does not consists of any genomes that belong to either strain or species or genus or family or order or class of that particular microbe.

\*\*Superkingdom Shared: The training data set does not consists of any genomes that belong to either strain or species or genus or family or order or class or phylum of that particular microbe.

Supplementary Table 1: (contd)

454-100 data set				
S. No.	Organism	Number of sequences	Status with respect to modified RefDB*	Status with respect to training data set**
1.	<i>Acidovorax avenae</i>	1000	All Known	Phylum Shared
2.	<i>Agrobacterium tumefaciens</i> str. C58	1000	Species Unknown	Class Shared
3.	<i>Bacillus halodurans</i> C-125	1000	Species Unknown	Phylum Shared
4.	<i>Bdellovibrio bacteriovorus</i> HD100	1000	Order Unknown	Phylum Shared
5.	<i>Bordetella pertussis</i> Tohama I	1000	Family Unknown	Phylum Shared
6.	<i>Bradyrhizobium</i> sp. BTai1	1000	Species Unknown	Class Shared
7.	<i>Burkholderia cenocepacia</i> AU 1054	1000	All Known	Phylum Shared
8.	<i>Campylobacter jejuni</i> subsp. <i>jejuni</i> NCTC 11168	1000	Genus Unknown	Phylum Shared
9.	<i>Chlamydomonada pneumoniae</i> AR39	1000	Genus Unknown	Superkingdom Shared
10.	<i>Chlorobium chlorochromatii</i> CaD3	1000	Genus Unknown	Superkingdom Shared
11.	<i>Clostridium tetani</i> E88	1000	Species Unknown	Class Shared
12.	<i>Corynebacterium diphtheriae</i> NCTC 13129	1000	Species Unknown	Class Shared
13.	<i>Francisella philomiragia</i> subsp. <i>philomiragia</i> ATCC 25017	1000	Species Unknown	Class Shared
14.	<i>Haemophilus influenzae</i> Rd KW20	1000	Species Unknown	Class Shared
15.	<i>Helicobacter pylori</i> J99	1000	Strain Unknown	Phylum Shared
16.	<i>Lactobacillus acidophilus</i> NCFM	1000	Family Unknown	Phylum Shared
17.	<i>Leptospira interrogans</i> serovar <i>Lai</i> str. 56601	1000	Family Unknown	Superkingdom Shared
18.	<i>Magnetococcus</i> sp. MC-1	1000	All Known	Phylum Shared
19.	<i>Mycoplasma pneumoniae</i> M129	1000	Phylum Unknown	Superkingdom Shared
20.	<i>Neisseria meningitidis</i> MC58	1000	Genus Unknown	Phylum Shared
21.	<i>Nostoc</i> sp. PCC 7120	1000	Order Unknown	Superkingdom Shared
22.	<i>Polaromonas</i> sp. JS666	1000	Species Unknown	Phylum Shared
23.	<i>Pseudomonas aeruginosa</i> PAO1	1000	Genus Unknown	Class Shared
24.	<i>Pseudomonas mendocina</i> ymp	1000	Genus Unknown	Class Shared
25.	<i>Rhizobium etli</i> CFN 42	1000	Species Unknown	Class Shared
26.	<i>Rhodopseudomonas palustris</i> BisA53	1000	Genus Unknown	Class Shared
27.	<i>Rickettsia bellii</i> RML369-C	1000	Genus Unknown	Class Shared
28.	<i>Saccharophagus degradans</i> 2-40	1000	Genus Unknown	Class Shared
29.	<i>Shewanella baltica</i> OS185	1000	Species Unknown	Class Shared
30.	<i>Staphylococcus aureus</i> subsp. <i>aureus</i>	1000	Species Unknown	Phylum Shared
31.	<i>Streptomyces coelicolor</i> A3(2)	1000	Species Unknown	Class Shared
32.	<i>Thermotoga maritima</i> MSB8	1000	Genus Unknown	Superkingdom Shared
33.	<i>Trichodesmium erythraeum</i> IMS101	1000	Order Unknown	Superkingdom Shared
34.	<i>Wigglesworthia glossinidia</i> endosymbiont of <i>Glossina brevipalpis</i>	1000	Genus Unknown	Class Shared
35.	<i>Xanthomonas campestris</i> pv. <i>campestris</i> str. ATCC 33913	1000	Genus Unknown	Class Shared

\*All Known: Genome fragments of that particular microbe is present in the modified RefDB.

\*Strain Unknown: Genome fragments of this strain are absent in the modified RefDB.

\*Species Unknown: Genome fragments of all strains belonging to this species are absent in the modified RefDB.

\*Genus Unknown: Genome fragments of all the strains and species belonging to this genus are absent in modified RefDB.

\*Family Unknown: Genome fragments of all genera belonging to this family are absent in the modified RefDB.

\*Order Unknown: Genome fragments of all families belonging to this order are absent in the modified RefDB.

\*Phylum Unknown: Genome fragments of all classes belonging to this phylum are absent in modified RefDB.

\*\*Class Shared: The training data set is devoid of genomes that can share either strain or species or genus or family or order with that particular microbe.

\*\*Phylum Shared: The training data set does not consists of any genomes that belong to either strain or species or genus or family or order or class of that particular microbe.

\*\*Superkingdom Shared: The training data set does not consists of any genomes that belong to either strain or species or genus or family or order or class or phylum of that particular microbe.