

Supplementary Table 2: The organisms constituting the simLC, simMC and the simHC data sets along with their representation status* with respect to modified reference database (modified RefDB) and their phylogenetic similarity status** with respect to the genomes of training data set.

S. No.	Organism	Number of sequences			Status with respect to modified RefDB*	Status with respect to training data set**
		simLC	simMC	simHC		
1	<i>Actinobacillus succinogenes</i> 130Z	252	290	483	All Known	Class Shared
2	<i>Alkalilimnicola ehrlichei</i> MLHE-1	373	384	829	All Known	Family Shared
3	<i>Alkaliphillus metalliredigenes</i> UNDEF	489	504	1091	Genus Unknown	Class Shared
4	<i>Anabaena variabilis</i> ATCC 29413	855	881	1703	Class Unknown	Superkingdom Shared
5	<i>Anaeromyxobacter dehalogenans</i> 2CP-C	584	578	1273	All Known	Phylum Shared
6	<i>Arthrobacter</i> sp. FB24	570	613	1211	All Known	Class Shared
7	<i>Azotobacter vinelandii</i> AvOP	650	652	1311	All Known	Class Shared
8	<i>Bacillus cereus</i> NVH391-98	520	454	1000	All Known	Phylum Shared
9	<i>Bifidobacterium longum</i> DJO10A	288	235	610	All Known	Species Shared
10	<i>Bradyrhizobium</i> sp. BTAi1	9277	22956	2127	All Known	Class Shared
11	<i>Brevibacterium linens</i> BL2	542	506	1088	Family Unknown	Class Shared
12	<i>Burkholderia ambifaria</i> AMMD	955	926	1937	All Known	Phylum Shared
13	<i>Burkholderia cenocepacia</i> AU 1054	879	917	1791	All Known	Phylum Shared
14	<i>Burkholderia cenocepacia</i> HI2424	956	1010	2045	All Known	Phylum Shared
15	<i>Burkholderia</i> sp. sp.strain 383	1074	1024	2191	All Known	Phylum Shared
16	<i>Burkholderia vietnamiensis</i> G4	992	985	2083	All Known	Phylum Shared
17	<i>Burkholderia xenovorans</i> LB400	1149	1146	2384	All Known	Phylum Shared
18	<i>Caldicellulosiruptor saccharolyticus</i> UNDEF	367	316	658	Order Unknown	Order Shared
19	<i>Chlorobium limicola</i> DSMZ 245(T)	381	353	671	Genus Unknown	Superkingdom Shared
20	<i>Chlorobium phaeobacteroides</i> BS1	483	515	1082	Genus Unknown	Superkingdom Shared
21	<i>Chlorobium phaeobacteroides</i> DSM 266	359	390	719	Genus Unknown	Superkingdom Shared
22	<i>Chlorobium vibrioforme</i> f. thiosulfatophilum DSMZ 265(T)	199	286	534	Genus Unknown	Superkingdom Shared
23	<i>Chloroflexus aurantiacus</i> J-10-fl	679	668	1277	All Known	Family Shared
24	<i>Chromohalobacter salexigens</i> DSM3043	454	481	888	All Known	Order Shared
25	<i>Clostridium beijerincki</i> NCIMB 8052	737	702	1411	Species Unknown	Class Shared
26	<i>Clostridium thermocellum</i> ATCC 27405	461	494	932	Species Unknown	Class Shared
27	<i>Crocospaera watsonii</i> WH 8501	812	776	1593	Family Unknown	Superkingdom Shared
28	<i>Cytophaga hutchinsonii</i> ATCC 33406	5168	580	1161	All Known	Phylum Shared
29	<i>Dechloromonas aromatica</i> RCB	537	559	1132	All Known	Phylum Shared
30	<i>Deinococcus geothermalis</i> DSM11300	415	404	809	All Known	All Shared
31	<i>Desulfitobacterium hafniense</i> DCB-2	769	725	1486	Family Unknown	Class Shared
32	<i>Desulfovibrio desulfuricans</i> G20	484	424	919	All Known	Phylum Shared
33	<i>Ehrlichia canis</i> Jake	196	197	283	All Known	Class Shared
34	<i>Ehrlichia chaffeensis</i> sapulpa	150	155	255	All Known	Class Shared
35	<i>Enterococcus faecium</i> DO	359	360	676	Family Unknown	Phylum Shared
36	<i>Exiguobacterium</i> UNDEF 255-15	377	392	788	All Known	Phylum Shared

S. No.	Organism	Number of sequences			Status with respect to modified RefDB*	Status with respect to training data set**
		simLC	simMC	simHC		
37	<i>Ferroplasma acidarmanus fer1</i>	238	264	471	Family Unknown	Order Shared
38	<i>Frankia sp. Ccl3</i>	645	643	1334	All Known	Class Shared
39	<i>Frankia sp. EAN1pec</i>	1109	1111	2248	All Known	Class Shared
40	<i>Geobacter metallireducens GS-15</i>	515	520	1025	All Known	Phylum Shared
41	<i>Haemophilus somnus 129PT</i>	232	265	513	All Known	Class Shared
42	<i>Jannaschia sp. CCS1</i>	543	543	1148	All Known	Family Shared
43	<i>Kineococcus radiotolerans SRS30216</i>	566	580	1187	All Known	Class Shared
44	<i>Lactobacillus brevis ATCC 367</i>	177	247	445	Family Unknown	Phylum Shared
45	<i>Lactobacillus casei ATCC 334</i>	362	344	648	Family Unknown	Phylum Shared
46	<i>Lactobacillus delbrueckii bulgaricus ATCC BAA-365</i>	195	203	391	Family Unknown	Phylum Shared
47	<i>Lactobacillus gasseri ATCC 33323</i>	244	258	551	Family Unknown	Phylum Shared
48	<i>Lactococcus lactis cremoris SK11</i>	301	345	584	Genus Unknown	Phylum Shared
49	<i>Leuconostoc mesenteroides mesenteroides ATCC 8293</i>	235	241	473	Family Unknown	Phylum Shared
50	<i>Magnetococcus sp. MC-1</i>	504	520	1153	All Known	Phylum Shared
51	<i>Marinobacter aquaeolei VT8</i>	547	586	1164	All Known	Class Shared
52	<i>Mesorhizobium sp. BNC1</i>	567	605	1289	All Known	Class Shared
53	<i>Methanococcoides burtonii DSM6242</i>	268	318	663	All Known	Family Shared
54	<i>Methanosarcina barkeri Fusaro</i>	545	608	1213	All Known	Genus Shared
55	<i>Methanospirillum hungatei JF-1</i>	429	459	919	All Known	Order Shared
56	<i>Methylobacillus flagellatus strain KT</i>	365	321	687	All Known	Phylum Shared
57	<i>Moorella thermoacetica ATCC 39073</i>	377	765	740	Order Unknown	Family Shared
58	<i>Nitrobacter hamburgensis UNDEF</i>	630	568	1272	All Known	Class Shared
59	<i>Nitrobacter winogradskyi Nb-255</i>	427	358	857	All Known	Class Shared
60	<i>Nitrosococcus oceani UNDEF</i>	409	406	868	All Known	Order Shared
61	<i>Nitrosomonas eutropha C71</i>	314	320	649	All Known	Phylum Shared
62	<i>Nitrospira multiformis ATCC 25196</i>	378	420	814	All Known	Phylum Shared
63	<i>Nocardioides sp. JS614</i>	636	691	1337	All Known	Class Shared
64	<i>Novosphingobium aromaticivorans DSM 12444 (F199)</i>	520	536	1093	All Known	Family Shared
65	<i>Oenococcus oeni PSU-1</i>	182	221	422	Family Unknown	Phylum Shared
66	<i>Paracoccus denitrificans PD1222</i>	585	620	1362	All Known	Family Shared
67	<i>Pediococcus pentosaceus ATCC 25745</i>	217	173	456	Family Unknown	Phylum Shared
68	<i>Pelobacter carbinolicus DSM 2380</i>	489	472	896	All Known	Phylum Shared
69	<i>Pelobacter propionicus DSM 2379</i>	508	550	1145	All Known	Phylum Shared
70	<i>Pelodictyon luteolum UNDEF</i>	250	314	581	Genus Unknown	Superkingdom Shared
71	<i>Pelodictyon phaeoclathratiforme BU-1 (DSMZ 5477(T))</i>	402	359	703	Genus Unknown	Superkingdom Shared
72	<i>Polaromonas sp. JS666</i>	733	769	1489	Species Unknown	Phylum Shared
73	<i>Prochlorococcus marinus str. MIT 9312</i>	183	218	404	All Known	Superkingdom Shared
74	<i>Prochlorococcus sp. NATL2A</i>	253	191	480	All Known	Superkingdom Shared
75	<i>Prosthecochloris aestuarii SK413/DSMZ 271(t)</i>	282	269	692	Species Unknown	Superkingdom Shared
76	<i>Pseudoalteromonas atlantica T6c</i>	588	657	1301	All Known	Class Shared

S. No.	Organism	Number of sequences			Status with respect to modified RefDB*	Status with respect to training data set**
		simLC	simMC	simHC		
77	<i>Pseudomonas fluorescens</i> PfO-1	730	791	1587	Genus Unknown	Class Shared
78	<i>Pseudomonas putida</i> F1	675	745	1528	Genus Unknown	Class Shared
79	<i>Pseudomonas syringae</i> B728a	746	738	1545	Genus Unknown	Class Shared
80	<i>Psychrobacter arcticum</i> 273-4	327	329	623	All Known	Class Shared
81	<i>Psychrobacter cryopegella</i> UNDEF	422	366	793	All Known	Class Shared
82	<i>Rhodobacter sphaeroides</i> 2.4.1	514	528	1119	All Known	Species Shared
83	<i>Rhodoferrax ferrireducens</i> UNDEF	599	588	1276	All Known	Phylum Shared
84	<i>Rhodopseudomonas palustris</i> BisA53	636	680	1392	Genus Unknown	Class Shared
85	<i>Rhodopseudomonas palustris</i> BisB18	699	6107	1348	Genus Unknown	Class Shared
86	<i>Rhodopseudomonas palustris</i> BisB5	575	16577	1200	Genus Unknown	Class Shared
87	<i>Rhodopseudomonas palustris</i> HaA2	28861	671	1339	Genus Unknown	Class Shared
88	<i>Rhodospirillum rubrum</i> ATCC 11170	559	4868	1062	All Known	All Shared
89	<i>Rubrobacter xylanophilus</i> DSM 9941	409	444	799	All Known	All Shared
90	<i>Ruegeria</i> sp. TM1040	469	470	1065	Species Unknown	Family Shared
91	<i>Saccharophagus degradans</i> 2-40	582	642	1324	Genus Unknown	Class Shared
92	<i>Shewanella amazonensis</i> SB2B	536	482	1055	All Known	Class Shared
93	<i>Shewanella baltica</i> OS155	621	649	1313	Species Unknown	Class Shared
94	<i>Shewanella frigidimarina</i> NCMB400	551	559	1257	All Known	Class Shared
95	<i>Shewanella putefaciens</i> UNDEF	565	530	1153	All Known	Class Shared
96	<i>Shewanella</i> sp. ANA-3	664	586	1279	All Known	Class Shared
97	<i>Shewanella</i> sp. MR-7	568	515	1177	All Known	Class Shared
98	<i>Shewanella</i> sp. PV-4	524	537	1165	All Known	Class Shared
99	<i>Shewanella</i> sp. W3-18-1	533	532	1214	All Known	Class Shared
100	<i>Sphingopyxis alaskensis</i> RB2256	438	402	846	All Known	Family Shared
101	<i>Streptococcus suis</i> 89/1591	263	276	490	Species Unknown	Phylum Shared
102	<i>Streptococcus thermophilus</i> LMD-9	178	187	501	Species Unknown	Phylum Shared
103	<i>Sulfurimonas denitrificans</i> DSM 1251	277	312	490	All Known	Phylum Shared
104	<i>Synechococcus</i> sp. PCC 7942 (<i>elongatus</i>)	316	309	646	All Known	Superkingdom Shared
105	<i>Syntrophobacter fumaroxidans</i> MPOB	606	552	1181	All Known	Phylum Shared
106	<i>Syntrophomonas wolfei</i> Goettingen	314	363	708	Family Unknown	Class Shared
107	<i>Thermoanaerobacter ethanolicus</i> 39E	315	265	570	Order Unknown	All Shared
108	<i>Thermobifida fusca</i> YX	434	431	930	All Known	Class Shared
109	<i>Thiobacillus denitrificans</i> ATCC 25259	395	400	741	All Known	Phylum Shared
110	<i>Thiomicrospira crunogena</i> XCL-2	274	296	603	Family Unknown	Class Shared
111	<i>Trichodesmium erythraeum</i> IMS101	977	927	2051	Class Unknown	Superkingdom Shared
112	<i>Xylella fastidiosa</i> Dixon	601	10484	1303	All Known	Class Shared
	Total Number of Reads	97495	114834	116771		
	Total Number of Reads Considered after exclusion of low length reads***	96732	113373	115592		

- *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.

- **All Shared: The training data set includes the sequences from the genome of that particular microbe.
- **Species Shared: The training data set is devoid of genomes belonging to the strain of that particular microbe
- **Genus Shared: The training data set is devoid of genomes that can share either strain or species with that particular microbe.
- **Family Shared: The training data set does not consists of any genomes that belong to either strain or species or genus of that particular microbe.
- **Order Shared: The training data set does not consists of any genomes that belong to either strain or species or genus or family of that particular microbe.
- **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.
- **Super-kingdom 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.

***Sequences having a length less than 400bp were removed from all three data sets.