

Additional File 1

Supporting Information

Reference-free Learning with Multiple Metagenomic Samples

1 Simulated Data Sets

In our simulation studies, we generated reads using MetaSim [1] with Illumina 80-bp paired-end reads error model. To model the high variation of the bacteria distribution across different samples, we simulated the relative abundances of species in each sample from a Dirichlet distribution. The 150 microbial species used in all simulation settings are listed in Table **S1**. MetaGen, MaxBin and CLARK are based on contigs longer than 1000bp. MetaBAT is based on contigs longer than 1500bp, which is the algorithm’s minimum length cut-off. In simulation settings 1-3, we considered the sequencing depth, the number of samples, and the number of species as three variables. In each setting, we fixed two variables and change the other one.

Simulation Settings 1-3:

In the first setting, we fixed the number of samples at 80 and the number of specie at 100, and changed the pooled sequencing depth as 80x, 120x and 160x, which is equivalent to 1x, 1.5x and 2x per sample. In the second setting, we fixed the pooled sequencing depth at 120x and the number of species at 100, and changed the number of sample as “20”, “40” and “80”. In the third setting, we fixed the pooled sequencing depth at 120x and the number of samples at 80, and changed the number of species as “50”, “100” and “150”. All the detailed binning results of MetaGen, CLARK, MaxBin, MetaBAT, and CONCOCT are shown in Figures **5** and **S2-S8**. The computing times are shown in Figure **S1**. The correlations between the estimated relative abundances and the true relative abundances is shown in Figure **6**.

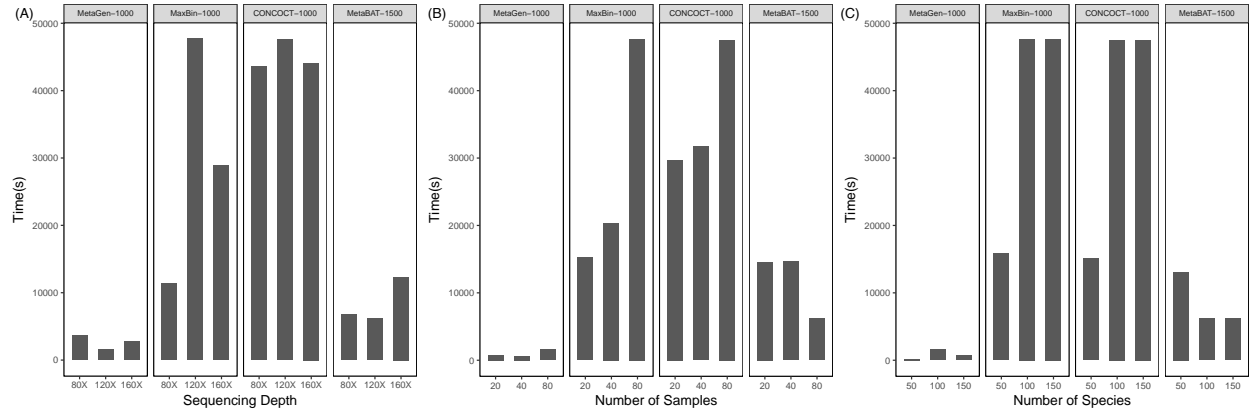


Figure S1: Computing times for MetaGen, MaxBin, CONCOCT and MetaBAT under different sequencing depths, different numbers of samples and different numbers of species.

80x-80sample-100species

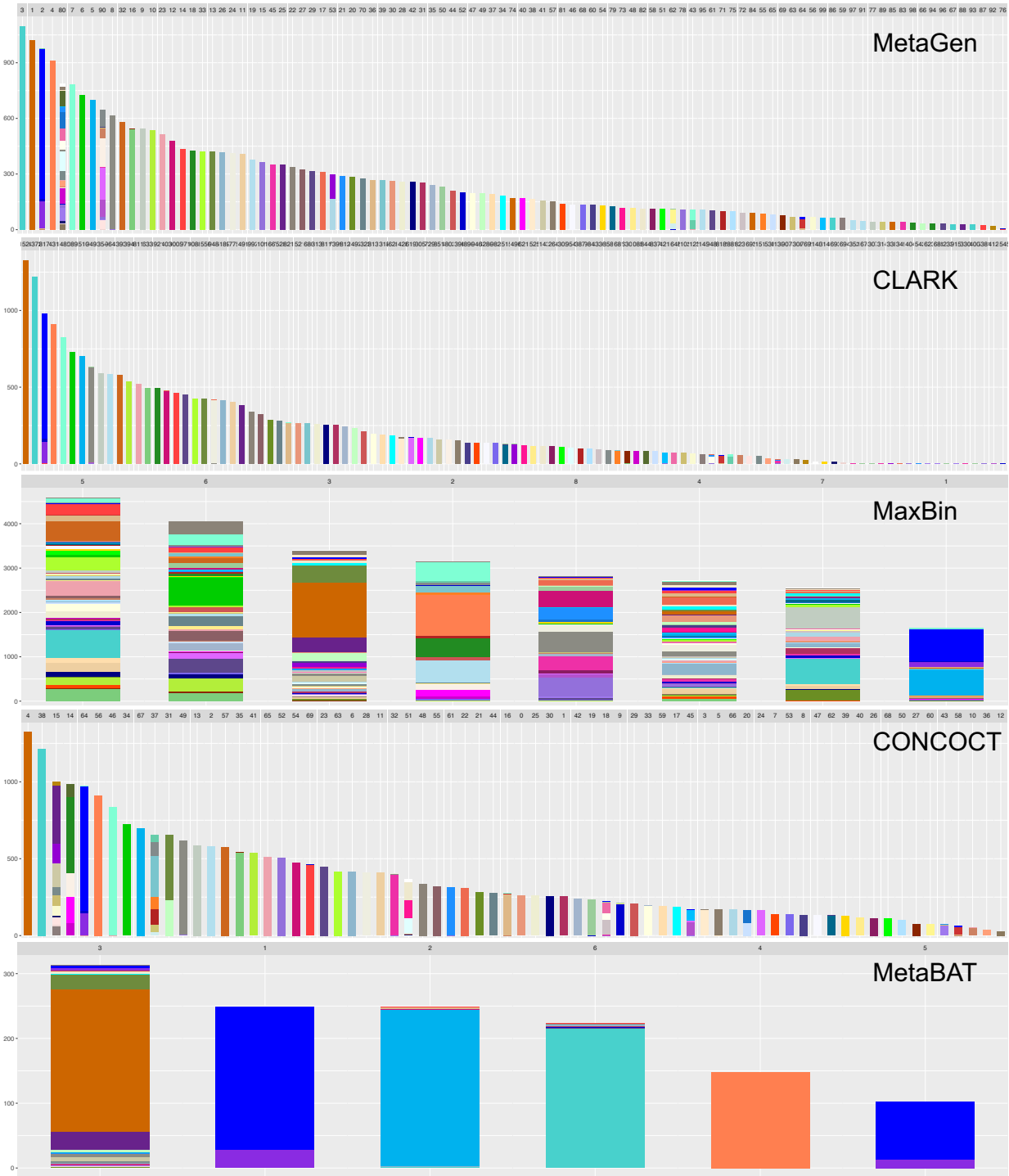


Figure S2: Binning results for MetaGen, CLARK, MaxBin, CONCOCT and MetaBAT with the pooled sequencing depth fixed at 80x, the number of samples at 80 and the number of species at 100. Color legends are given in **Figure S9**.

120x-80sample-100species

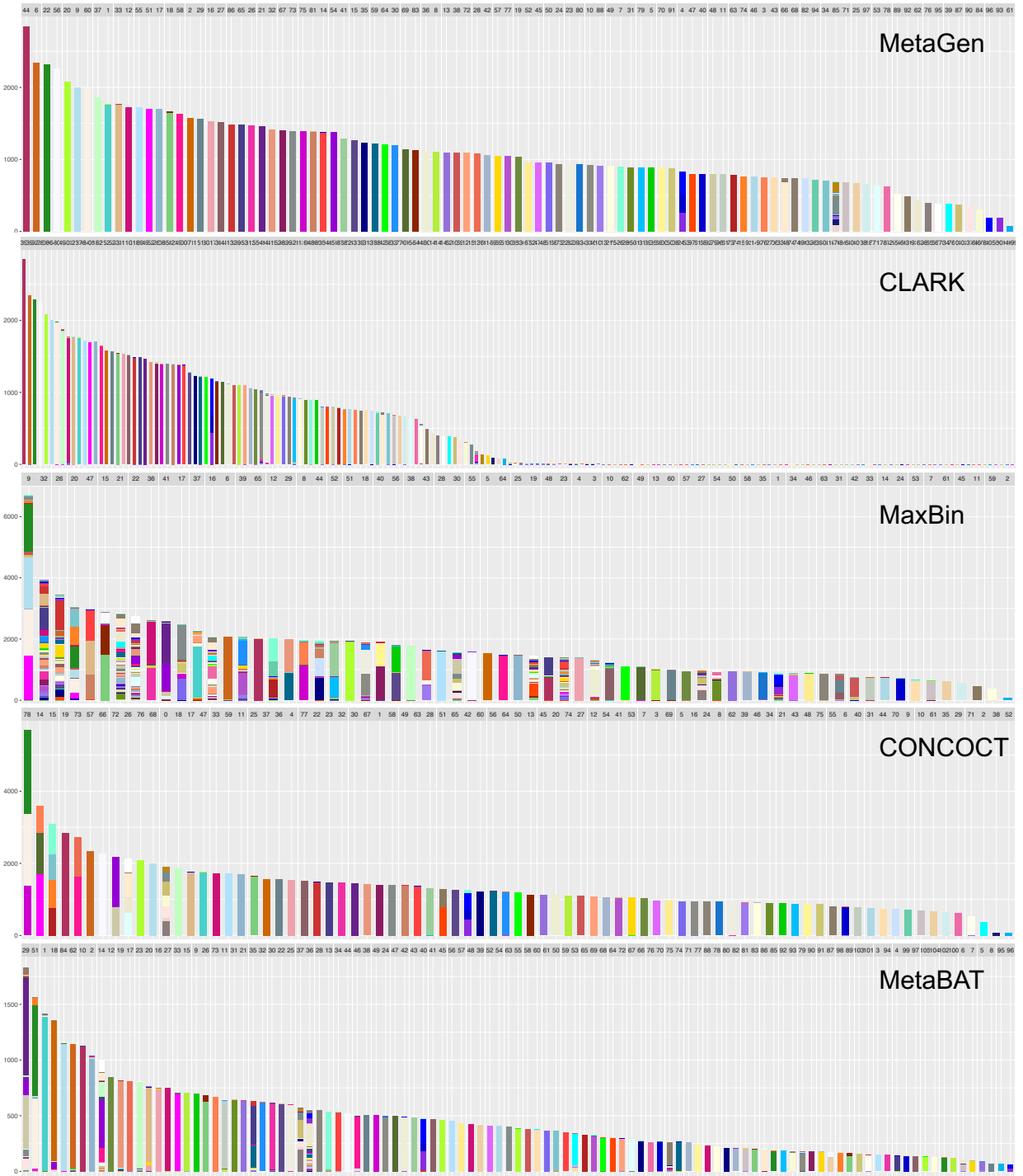


Figure S3: Binning results for MetaGen, CLARK, MaxBin, CONCOCT and MetaBAT with the pooled sequencing depth at 120x, the number of samples at 80, and the number of species at 100. Color legends are given in **Figure S9**.

160x-80sample-100species

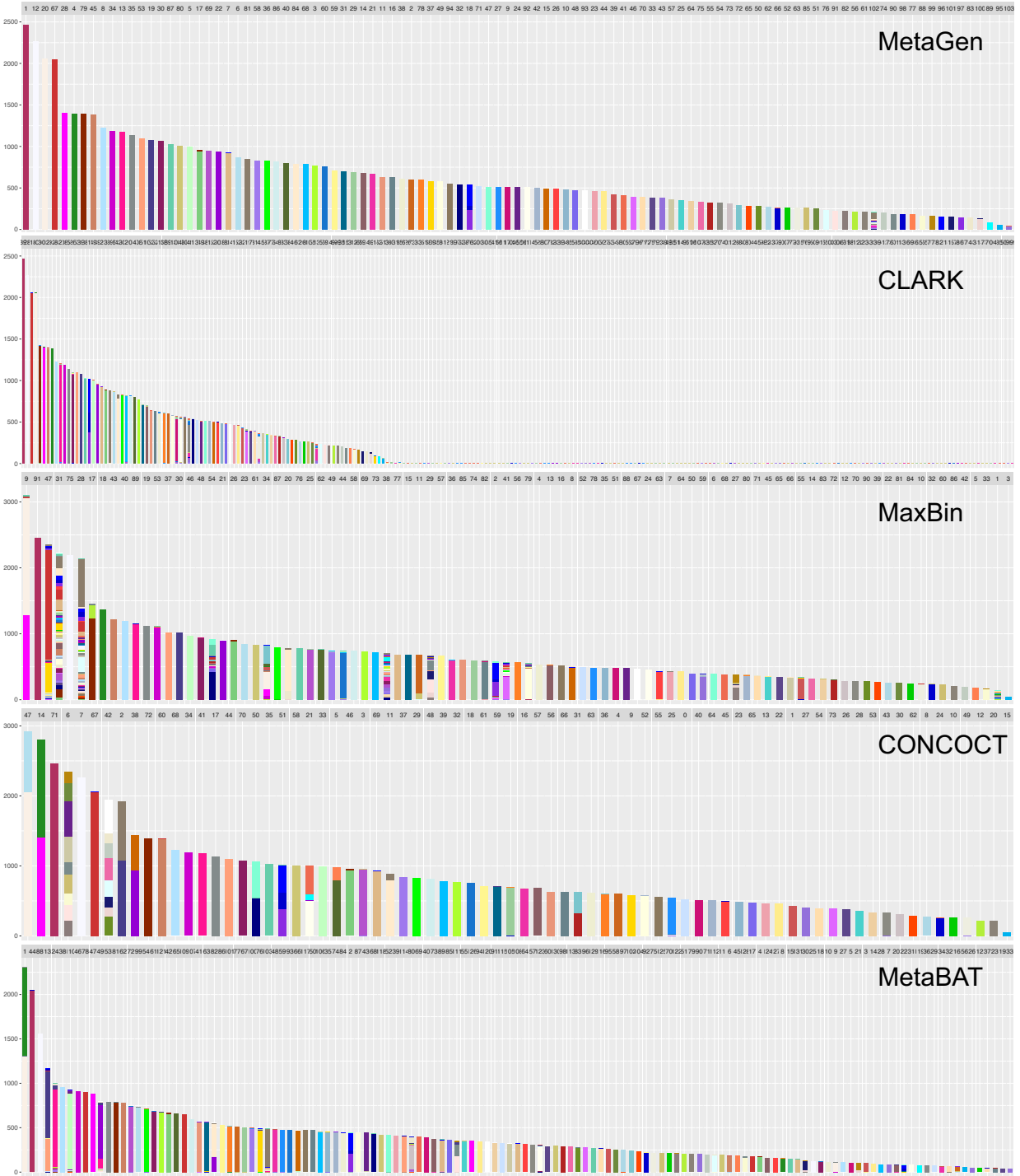


Figure S4: Binning results for MetaGen, CLARK, MaxBin, CONCOCT, and MetaBAT with the pooled sequencing depth at 160x, the number of samples at 80, and the number of species at 100. Color legends are given in **Figure S9**.

120x-20sample-100species

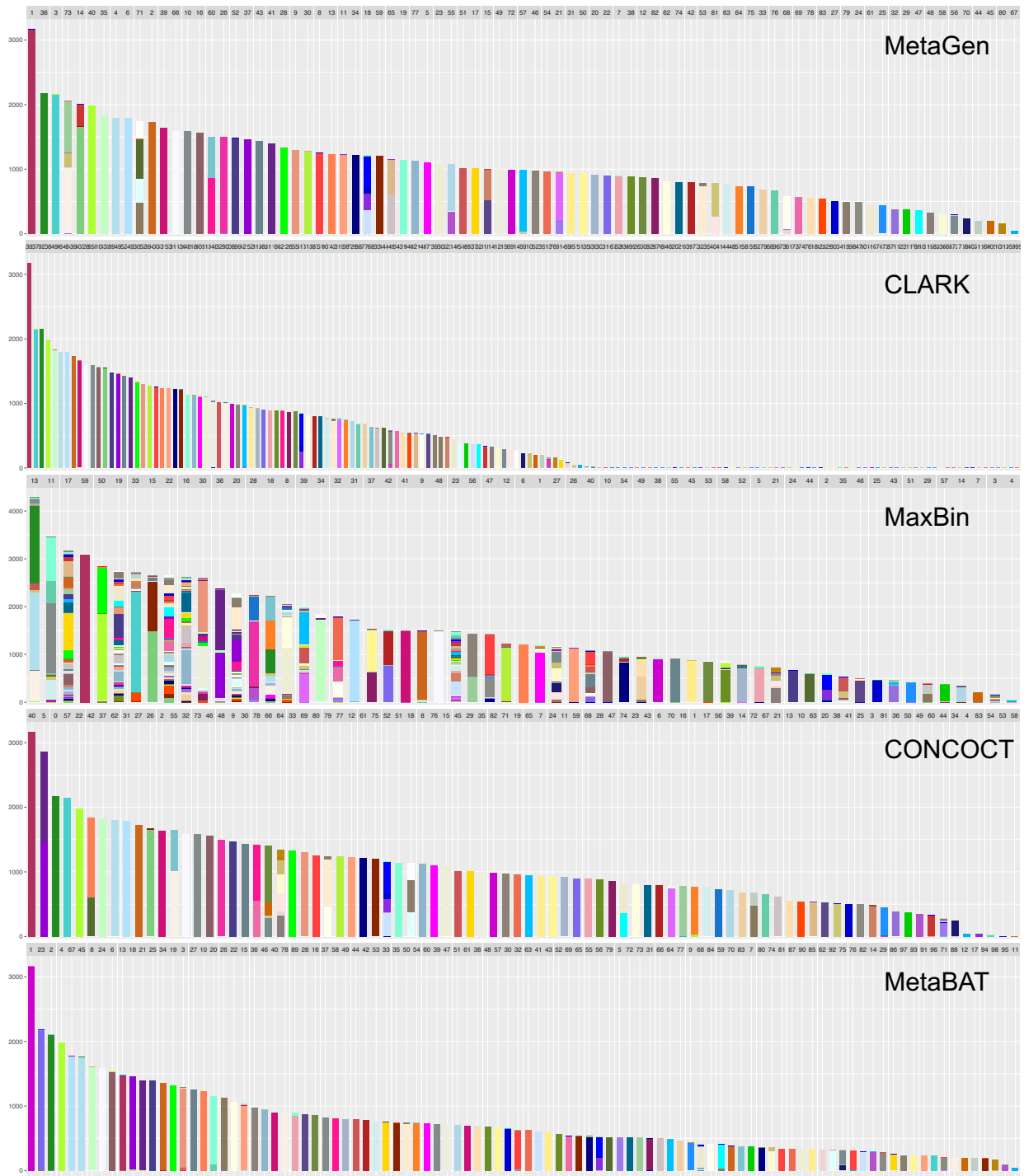


Figure S5: Binning results for MetaGen, CLARK, MaxBin, CONCOCT, and MetaBAT with the pooled sequencing depth at 120x, the number of samples at 20 and the number of species at 100. Color legends are given in **Figure S9**.

120x-40sample-100species

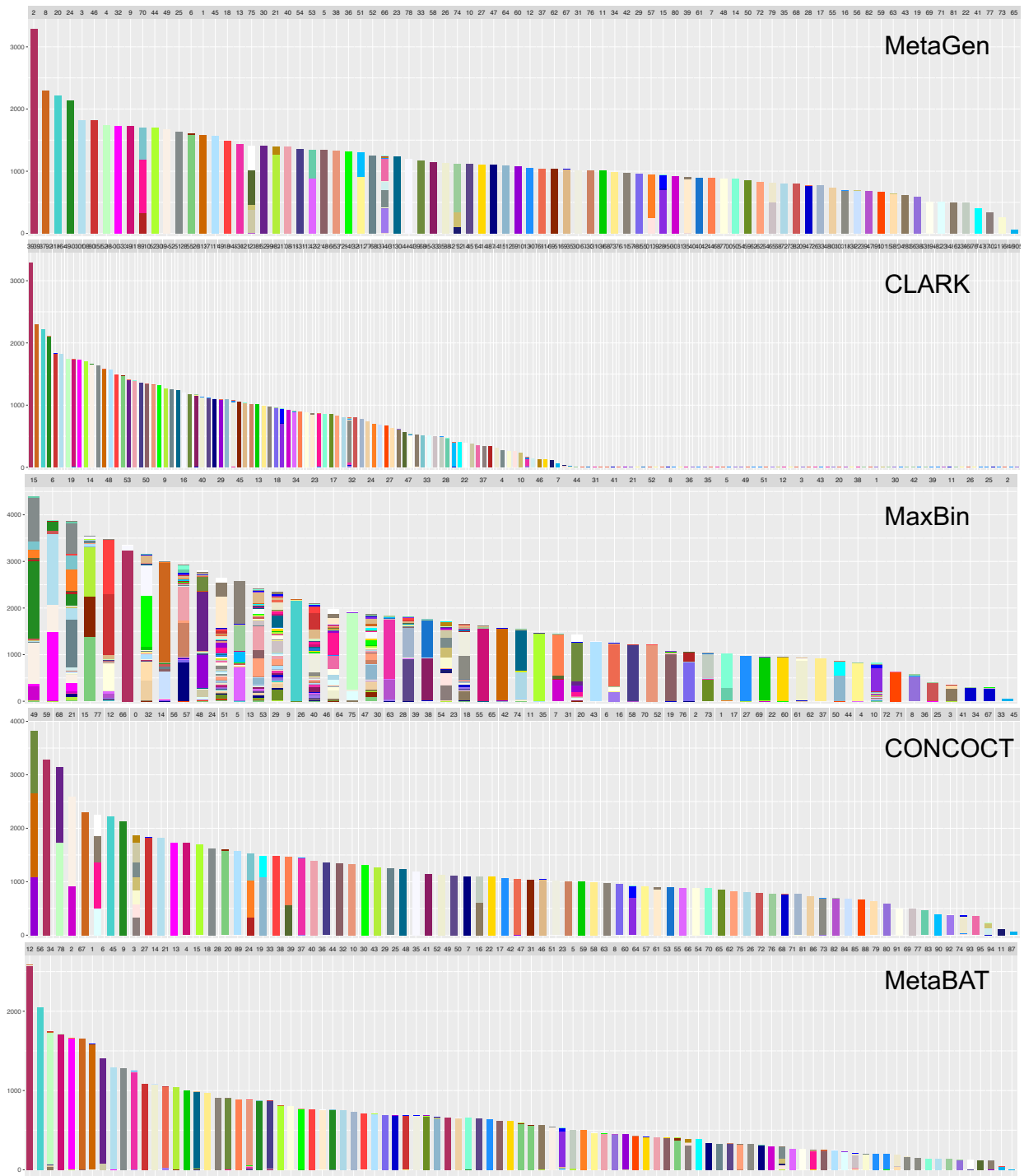


Figure S6: Binning results for MetaGen, CLARK, MaxBin, CONCOCT, and MetaBAT with the pooled sequencing depth at 120x, the number of samples at 40, and the number of species at 100. Color legends are given in **Figure S9**.

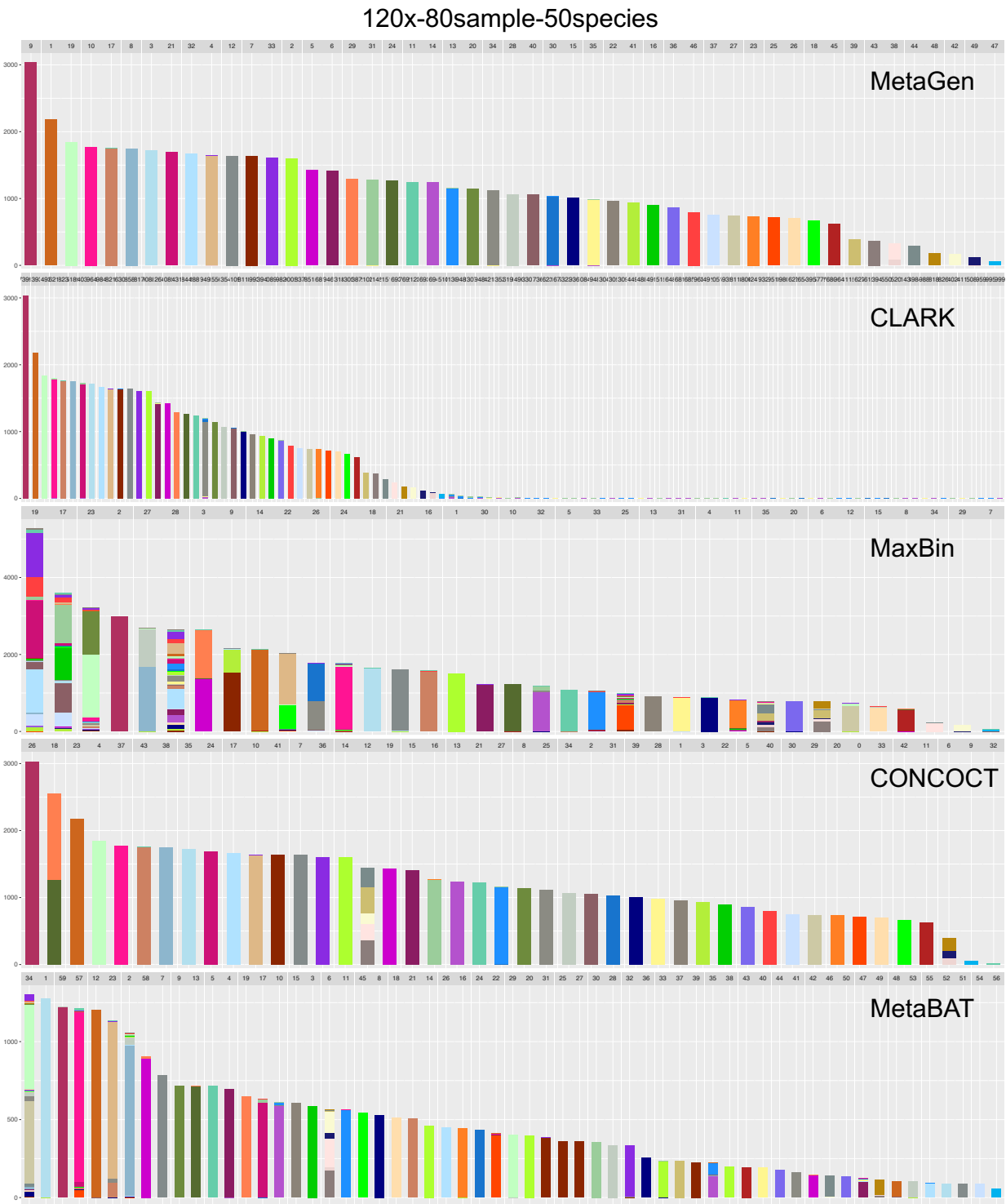


Figure S7: Binning results for MetaGen, CLARK, MaxBin, CONCOCT, and MetaBAT with the pooled sequencing depth at 120x, the number of samples at 80, and the number of species at 50. Color legends are given in **Figure S9**.

120x-80sample-150species

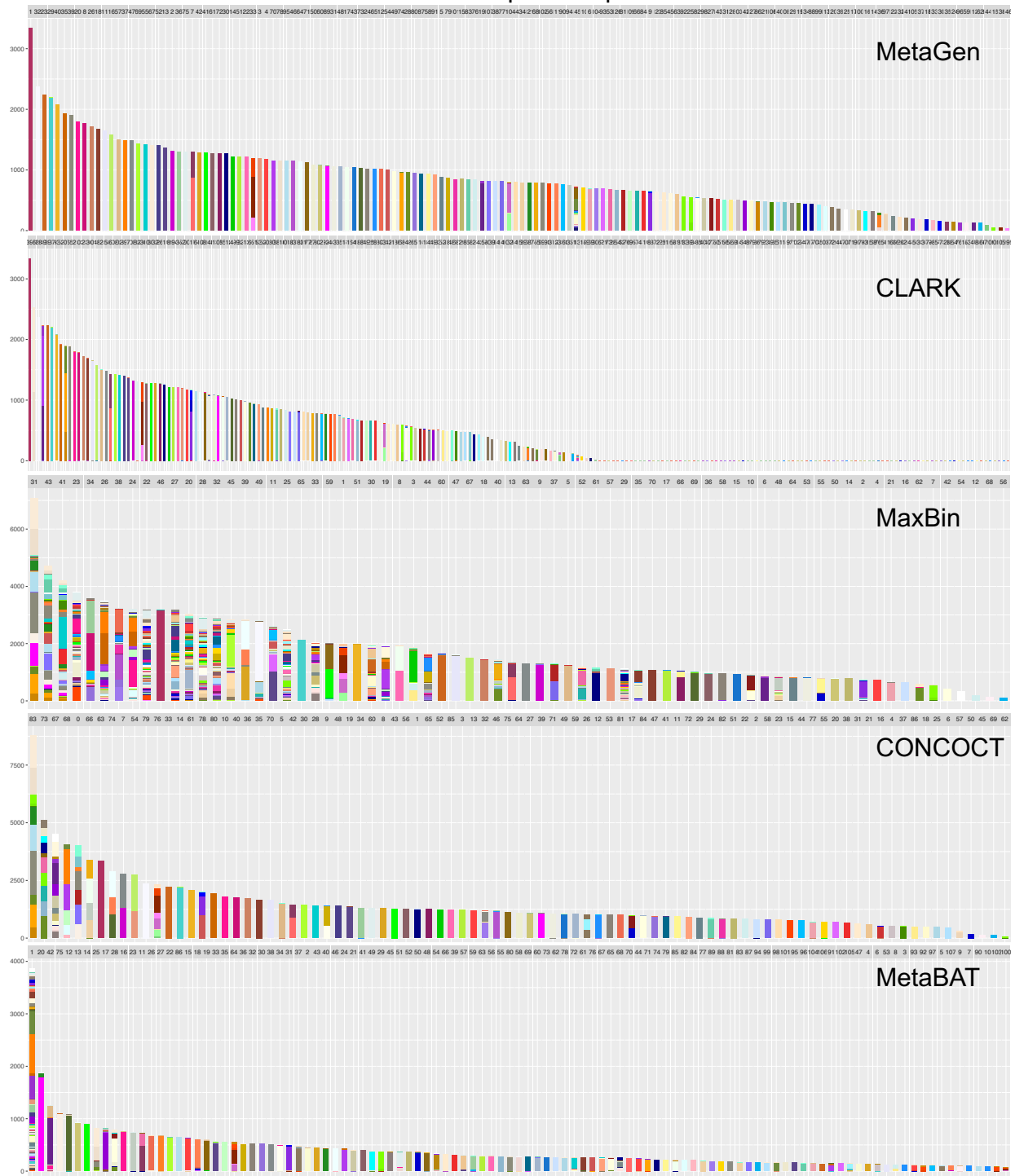


Figure S8: Binning results for MetaGen, CLARK, MaxBin, CONCOCT, and MetaBAT with the pooled sequencing depth at 120x, the number of samples at 80, and the number of species at 150. Color legends are given in **Figure S9**.

Acholeplasma laidlawii PG-8A	Butyrivibrio proteoclasticus B316	Coprococcus sp. ART55/1	Flavobacteriaceae bacterium 3519-10	Ruminococcus albus 7
Achromobacter xylosoxidans A8	Candidatus Arthromitus sp.	Coriobacterium glomerans PW2	Fusobacterium nucleatum subsp.	Ruminococcus bromii L2-63
Achromobacter xylosoxidans NH44784-1996	Candidatus Azobacteroides pseudotrichonymphae	Corynebacterium argentoratense DSM	Haemophilus parainfluenzae T3T1	Ruminococcus champanellensis 18P13
Akkermansia muciniphila ATCC	Candidatus Chloracidobacterium thermophilum	Cryptobacterium curtum DSM	Herbaspirillum seropedicae SmR1	Ruminococcus obeum A2-162
Alistipes finegoldii DSM	Candidatus Liberibacter solanacearum	Cupriavidus metallururgans CH34	Janthinobacterium sp. Marseille	Ruminococcus sp. SR1/5
Alistipes shahii WAL	Candidatus Nitrospira defluvi	Cylinndrospermum stagnale PCC	Lawsonia intracellulularis PHE/MN1-00	Ruminococcus torques L2-14
Alkalimnicola ehrlichii MLHE-1	Candidatus Pelagibacter sp.	Dehalobacter sp. CF	Megasphaera elsdenii DSM	Selenomonas ruminantium subsp.
Alkaliphilus metalliredigens QYMF	Candidatus Phytoplasma mali	Denitrovibrio acetophilus DSM	Mesoplasma florum L1	Selenomonas sputigena ATCC
Asticcacaulis excentricus CB	Candidatus Protochlamydia amoebophila	Desulfarculus baarsii DSM	Moorella thermoacetica ATCC	Sphingobacterium sp. 21
Bacillus amyloliquefaciens XH7	Candidatus Rickettsia amblyommii	Desulfatibacillum aikenivorans AK-01	Mycoplasma mycoides subsp.	Spiroplasma chrysopicola DF-1
Bacillus amyloliquefaciens Y2	Capnocytophaga canimorsus Cc5	Desulfotobacterium dehalogenans ATCC	Mycoplasma putrefaciens Mput9231	Spiroplasma diminutum CUAS-1
Bacillus cereus E33L	Clavibacter michiganensis subsp.	Desulfotobacterium dichloroeliminans LMG	Mycoplasma suis K13806	Spiroplasma syphidicola EA-1
Bacillus thuringiensis serovar	Clostridium acetobutylicum EA	Desulfobacca acetoxidans DSM	Odoribacter splanchnicus DSM	Spiroplasma taiwanense CT-1
Bacteroides fragilis 638R	Clostridium acidurici 9a	Echinicola vietnamensis DSM	Oscillibacter valericigenes Sjm18-20	Streptococcus parasanguinis FW213
Bacteroides fragilis NCTC	Clostridium autoethanogenum DSM	Enterococcus faecium Aus0085	Paludibacter propionicigenes WB4	Streptococcus pseudopneumoniae IS7493
Bacteroides fragilis YCH46	Clostridium beijerinckii NCIMB	Enterococcus hirae ATCC	Parabacteroides distans ATCC	Streptococcus thermophilus ND03
Bacteroides helcogenes P	Clostridium botulinum A	Enterococcus mundtii QU	Pedococcus pentosaceus SL4	Tannerella forsythia ATCC
Bacteroides thetaiotaomicron VPI-5482	Clostridium botulinum B	Erysipelothrix rhusiopathiae SY1027	Pedobacter heparinus DSM	Tetragenococcus halophilus NBRC
Bacteroides vulgatus ATCC	Clostridium botulinum H04402	Escherichia coli LF82	Pedobacter saltans DSM	Tolomonas auensis DSM
Bartonella australis Aust/NH1	Clostridium cellulolyticum H10	Escherichia coli NA114	Pelagibacterium halotolerans B2	Variovorax paradoxus B4
Bartonella bacilliformis KCS83	Clostridium cellulovorans 743B	Ethanoligenens harbinense YUAN-3	Porphyromonas gingivalis ATCC	Variovorax paradoxus EPS
Bartonella clarridgeiae 73	Clostridium cf. saccharolyticum	Eubacterium cylindroides T2-87	Prevotella denticola F0289	Variovorax paradoxus S110
Bifidobacterium adolescentis ATCC	Clostridium clariflavum DSM	Eubacterium eligens ATCC	Prevotella intermedia 17	Veillonella parvula DSM
Bifidobacterium bifidum BGN4	Clostridium difficile R20291	Eubacterium rectale ATCC	Prevotella melaninogenica ATCC	Yersinia enterocolitica subsp. enterocolitica 8081
Bifidobacterium bifidum PRL2010	Clostridium phytofermentans ISDg	Eubacterium siraeum V10Sc8a	Ralstonia eutropha H16	Yersinia enterocolitica subsp. palearctica 105.5R(r)
Bifidobacterium breve UCC2003	Clostridium saccharolyticum WM1	Faecalibacterium prausnitzii L2-6	Ralstonia eutropha JMP134	Yersinia enterocolitica subsp. palearctica Y11
Bifidobacterium longum DJO10A	Clostridium sp. SY8519	Ferrimonas balearica DSM	Ralstonia solanacearum CMR15	Yersinia pestis A1122
Brevibacillus brevis NBRC	Clostridium stercorarium subsp.	Fibrobacter succinogenes subsp.	Rhodospirillum rubrum F11	Yersinia pestis Angola
Brevundimonas subvibrioides ATCC	Clostridium thermocellum DSM	Filifactor alocis ATCC	Rivularia sp. PCC	Yersinia pestis Antiqua
Butyrate-producing bacterium SSC/2	Coprococcus catus GD/7	Finexgoldia magna ATCC	Roseburia intestinalis XB6B4	Yersinia pseudotuberculosis IP

Figure S9: Color codes for all the 150 species used in Figures S1-S8.

Simulation Setting 4:

In this simulation setting, we first simulated a metagenomic community including 54 strains of *E. coli* (listed in Table S2). We generated 40 samples each with 2 million pair-end reads using MetaSim [1] with Illumina 80-bp paired-end reads error model.

Second, we simulated a 100x coverage of 100-bp paired-end Illumina reads of 6 *ecoli* strains (see Table S2) by ART simulator [2] with default settings for Illumina and library settings as -m 350 -s 50. These samples were further grouped together to simulate 6 samples using the relative abundance provided in (supplementary Table 1 in [3].)

Simulation Setting 5:

In this simulation study, we generated reads using MetaSim [1] with Illumina 80-bp paired-end reads error model. we generated 10 samples each with 25 million reads. To model the high variation of the bacteria distribution across different samples, the relative abundance of species in each sample was based on the 10 randomly selected real samples in our real examples. The histogram of the log-transformed relative abundance is shown in Figure S10. The 545 microbial species used in this simulation setting are listed in Table S3.

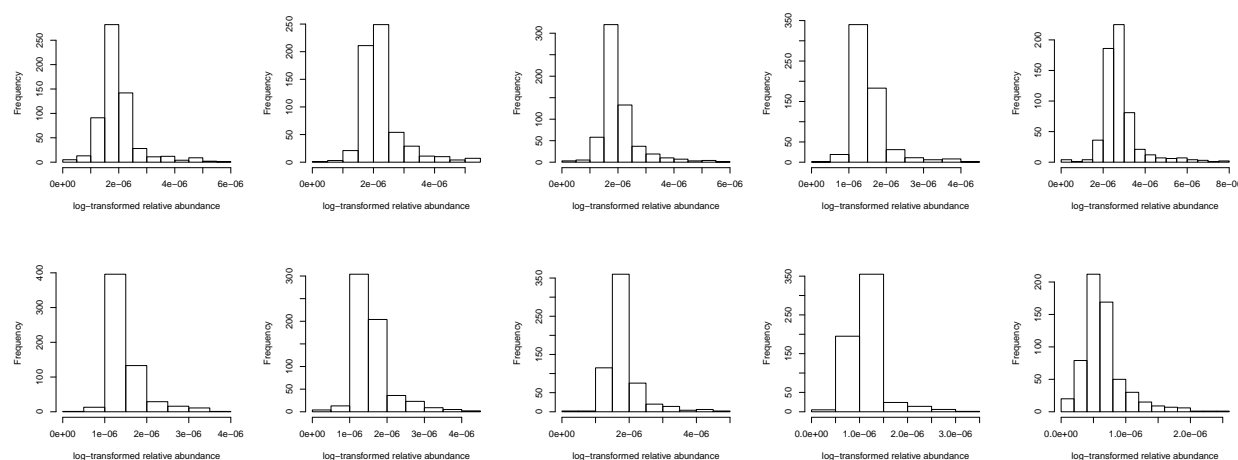


Figure S10: Histogram of the log-transformed relative abundance of different species for the 10 samples for the complex metagenomic data sets with 545 genomes and corresponding 439 circular elements.

Simulation Setting 6:

We fixed the pooled sequencing depth at 120x, the number of samples at 80, and the number of specie at 100, and changed the proportion of species in each sample from 50% to 100%. Performances of MetaGen, CLARK, MaxBin, MetaBAT and CONCOCT are shown in Figure S11. The detailed binning results of MetaGen, CLARK, MaxBin, MetaBAT, and CONCOCT are shown in

Figures S12, S13, 5.

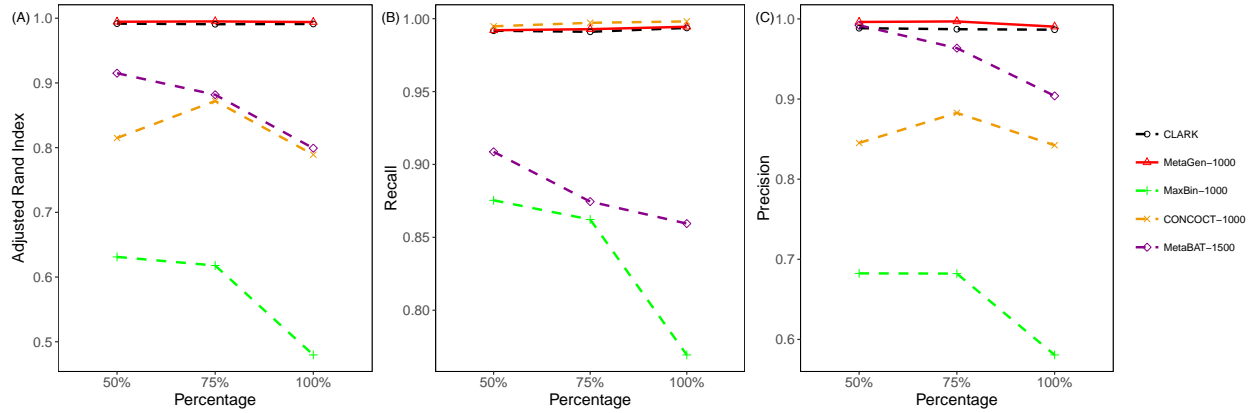


Figure S11: Adjusted Rand Index (A), Precision (B), and Recall (C) of CLARK, MetaGen, MaxBin, CONCOCT, and MetaBAT are evaluated on the simulated data sets with 50%, 75%, and 100% species, respectively, in each sample. All the three simulated data sets have 80 samples, 100 species, and the sequencing depth of 120x.

120x-80sample-100species-50

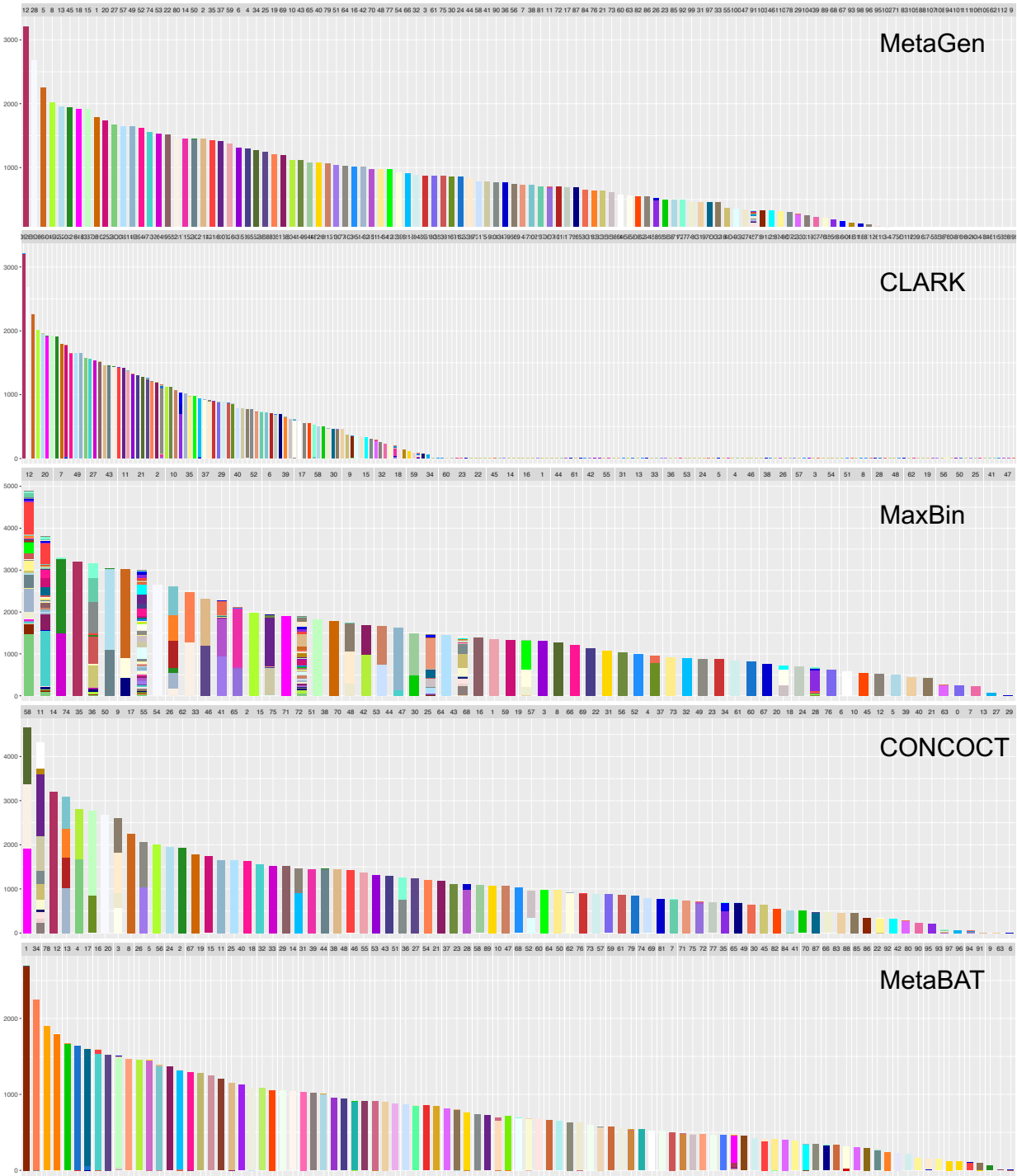


Figure S12: Binning results for MetaGen, CLARK, MaxBin, CONCOCT, and MetaBAT with the pooled sequencing depth at 120x, the number of samples at 80, and the number of species at 150, and with 50% of species in each sample. Color legends are given in **Figure S16**.

120x-80sample-100species-75

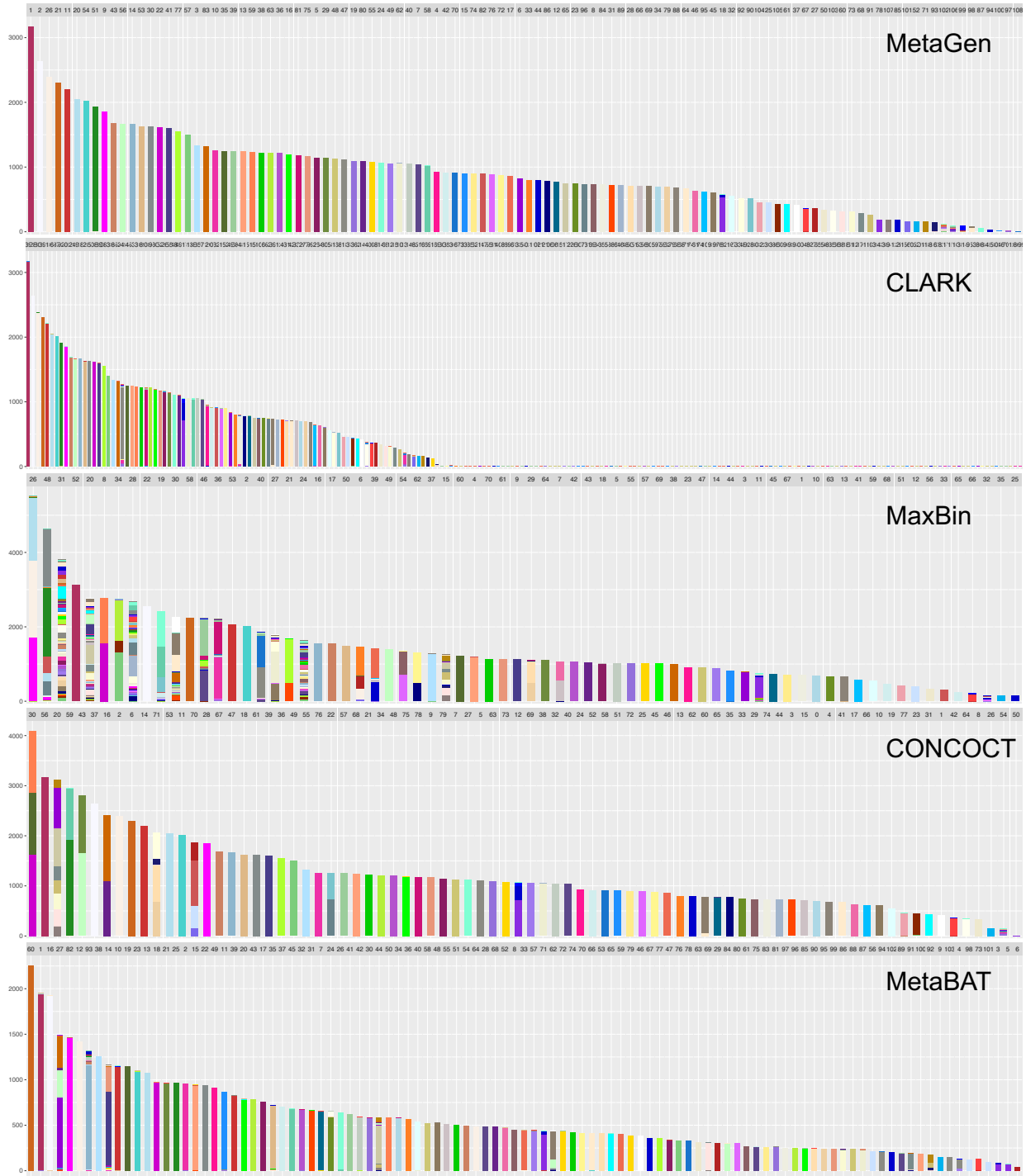


Figure S13: Binning results for MetaGen, CLARK, MaxBin, CONCOCT and MetaBAT with the pooled sequencing depth at 120x, the number of samples at 80, the number of species at 150, and with 75% of species in each sample. Color legends are given in **Figure S16**.

Simulation Setting 7:

We fixed the pooled sequencing depth at 120x, the number of samples at 80, and the number of specie at 100, and used two different assemblers, MegaHIT and Ray. Performances of MetaGen, CLARK, MaxBin, MetaBAT and CONCOCT are shown in Figure S14. The detailed binning results of MetaGen, CLARK, MaxBin, MetaBAT and CONCOCT are shown in Figures S15, 5.

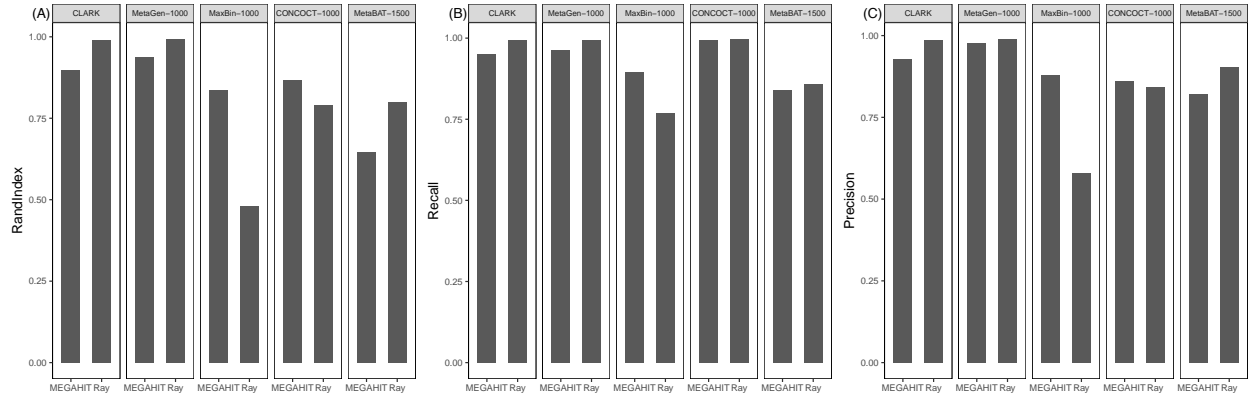


Figure S14: Adjusted Rand Index (A), Precision (B), and Recall (C) of CLARK, MetaGen, MaxBin, CONCOCT, and MetaBAT are evaluated using two different assemblers, MegaHIT and Ray. The simulated data sets have 80 samples, 100 species and the sequencing depth of 120x.

120x-80sample-100species-MegaHIT

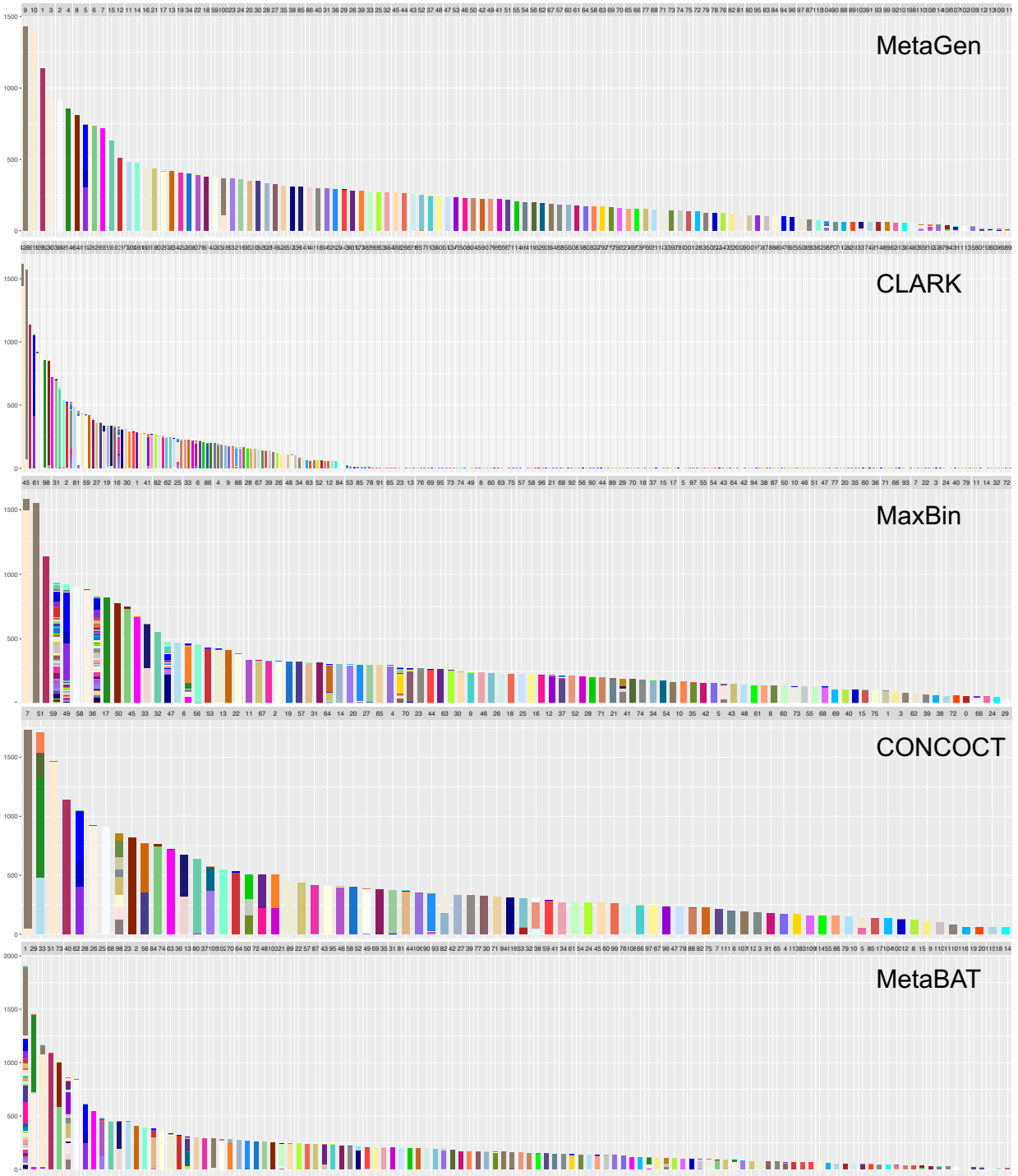


Figure S15: Binning results for MetaGen, CLARK, MaxBin, CONCOCT, and MetaBAT with the pooled sequencing depth at 120x, the number of samples at 80, and the number of species at 150, using the MegaHIT assembler. Color legends are given in **Figure S16**.

Acholeplasma laidlawii PG-8A	Butyrivibrio proteoclasticus B316	Coprococcus sp. ART55/1	Flavobacteriaceae bacterium 3519-10	Ruminococcus albus 7
Achromobacter xylosoxidans A8	Candidatus Arthromitus sp.	Coriobacterium glomerans PW2	Fusobacterium nucleatum subsp.	Ruminococcus bromii L2-63
Achromobacter xylosoxidans NH44784-1996	Candidatus Azobacteroides pseudotrichonymphae	Corynebacterium argentoratense DSM	Haemophilus parainfluenzae T3T1	Ruminococcus champanellensis 18P13
Akkermansia muciniphila ATCC	Candidatus Chloracidobacterium thermophilum	Cryptobacterium curtum DSM	Herbaspirillum seropedicae SmR1	Ruminococcus obeum A2-162
Alistipes finegoldii DSM	Candidatus Liberibacter solanacearum	Cupriavidus metallurians CH34	Janthinobacterium sp. Marseille	Ruminococcus sp. SR1/5
Alistipes shahii WAL	Candidatus Nitrospira defluvi	Cylinndrospermum stagnale PCC	Lawsonia intracellularis PHE/MN1-00	Ruminococcus torques L2-14
Alkalimnicola ehrlichii MLHE-1	Candidatus Pelagibacter sp.	Dehalobacter sp. CF	Megasphaera elsdenii DSM	Selenomonas ruminantium subsp.
Alkaliphilus metalliredigens QYMF	Candidatus Phytoplasma mali	Denitrovibrio acetophilus DSM	Mesoplasma florum L1	Selenomonas sputigena ATCC
Asticcacaulis excentricus CB	Candidatus Protochlamydia amoebophila	Desulfarculus baarsii DSM	Moorella thermoacetica ATCC	Sphingobacterium sp. 21
Bacillus amyloliquefaciens XH7	Candidatus Rickettsia amblyommii	Desulfatibacillum aikenivorans AK-01	Mycoplasma mycoides subsp.	Spiroplasma chrysopticola DF-1
Bacillus amyloliquefaciens Y2	Capnocytophaga canimorsus Cc5	Desulfotobacterium dehalogenans ATCC	Mycoplasma putrefaciens Mput9231	Spiroplasma diminutum CUAS-1
Bacillus cereus E33L	Clavibacter michiganensis subsp.	Desulfotobacterium dichloroeliminans LMG	Mycoplasma suis K13806	Spiroplasma syrrhidicola EA-1
Bacillus thuringiensis serovar	Clostridium acetobutylicum EA	Desulfobacca acetoxidans DSM	Odoribacter splanchnicus DSM	Spiroplasma taiwanense CT-1
Bacteroides fragilis 638R	Clostridium acidurici 9a	Echinicola vietnamensis DSM	Oscillibacter valericigenes Sjm18-20	Streptococcus parasanguinis FW213
Bacteroides fragilis NCTC	Clostridium autoethanogenum DSM	Enterococcus faecium Aus0085	Paludibacter propionigenes WB4	Streptococcus pseudopneumoniae IS7493
Bacteroides fragilis YCH46	Clostridium beijerinckii NCIMB	Enterococcus hirae ATCC	Parabacteroides distans ATCC	Streptococcus thermophilus ND03
Bacteroides helcogenes P	Clostridium botulinum A	Enterococcus mundtii QU	Pedococcus pentosaceus SL4	Tannerella forsythia ATCC
Bacteroides thetaiotaomicron VPI-5482	Clostridium botulinum B	Erysipelothrix rhusiopathiae SY1027	Pedobacter heparinus DSM	Tetragenococcus halophilus NBRC
Bacteroides vulgatus ATCC	Clostridium botulinum H04402	Escherichia coli LF82	Pedobacter saltans DSM	Tolomonas auensis DSM
Bartonella australis Aust/NH1	Clostridium cellulolyticum H10	Escherichia coli NA114	Pelagibacterium halotolerans B2	Variovorax paradoxus B4
Bartonella bacilliformis KCS83	Clostridium cellulovorans 743B	Ethanoligenens harbinense YUAN-3	Porphyromonas gingivalis ATCC	Variovorax paradoxus EPS
Bartonella clarridgeiae 73	Clostridium cf. saccharolyticum	Eubacterium cylindroides T2-87	Prevotella denticola F0289	Variovorax paradoxus S110
Bifidobacterium adolescentis ATCC	Clostridium clariflavum DSM	Eubacterium eligens ATCC	Prevotella intermedia 17	Veillonella parvula DSM
Bifidobacterium bifidum BGN4	Clostridium difficile R20291	Eubacterium rectale ATCC	Prevotella melaninogenica ATCC	Yersinia enterocolitica subsp. enterocolitica 8081
Bifidobacterium bifidum PRL2010	Clostridium phytofermentans ISDg	Eubacterium siraeum V10Sc8a	Ralstonia eutropha H16	Yersinia enterocolitica subsp. palearctica 105.5R(r)
Bifidobacterium breve UCC2003	Clostridium saccharolyticum WM1	Faecalibacterium prausnitzii L2-6	Ralstonia eutropha JMP134	Yersinia enterocolitica subsp. palearctica Y11
Bifidobacterium longum DJO10A	Clostridium sp. SY8519	Ferrimonas balearica DSM	Ralstonia solanacearum CMR15	Yersinia pestis A1122
Brevibacillus brevis NBRC	Clostridium stercorarium subsp.	Fibrobacter succinogenes subsp.	Rhodospirillum rubrum F11	Yersinia pestis Angola
Brevundimonas subvibrioides ATCC	Clostridium thermocellum DSM	Filifactor alocis ATCC	Rivularia sp. PCC	Yersinia pestis Antiqua
Butyrate-producing bacterium SSC/2	Coprococcus catus GD/7	Finexgoldia magna ATCC	Roseburia intestinalis XB6B4	Yersinia pseudotuberculosis IP

Figure S16: Color codes for all the 150 species used in simulation settings 4-5.

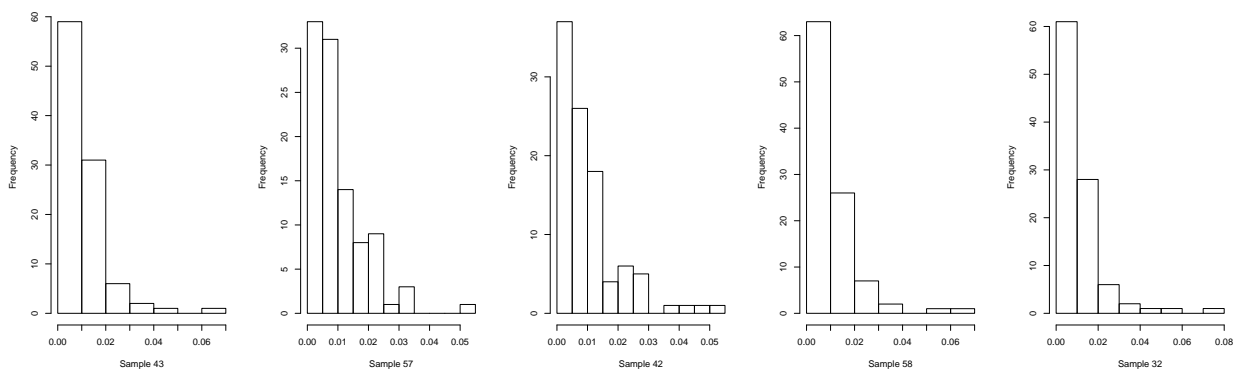


Figure S17: Histogram of the relative abundance of different species for the 5 randomly selected samples for the metagenomic data sets with the pooled sequencing depth at 80x, the number of samples at 80, and the number of species at 100.

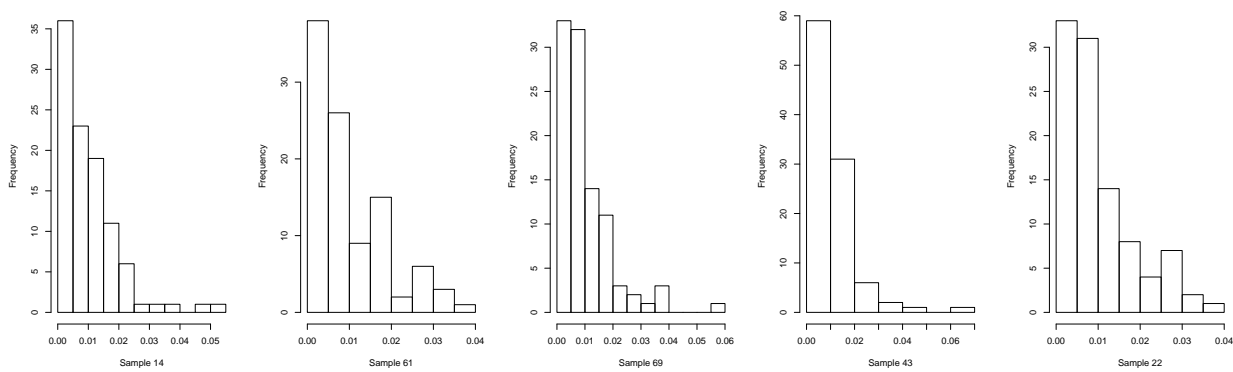


Figure S18: Histogram of the relative abundance of different species for the 5 randomly selected samples for the metagenomic data sets with the pooled sequencing depth at 120x, the number of samples at 80, and the number of species at 100.

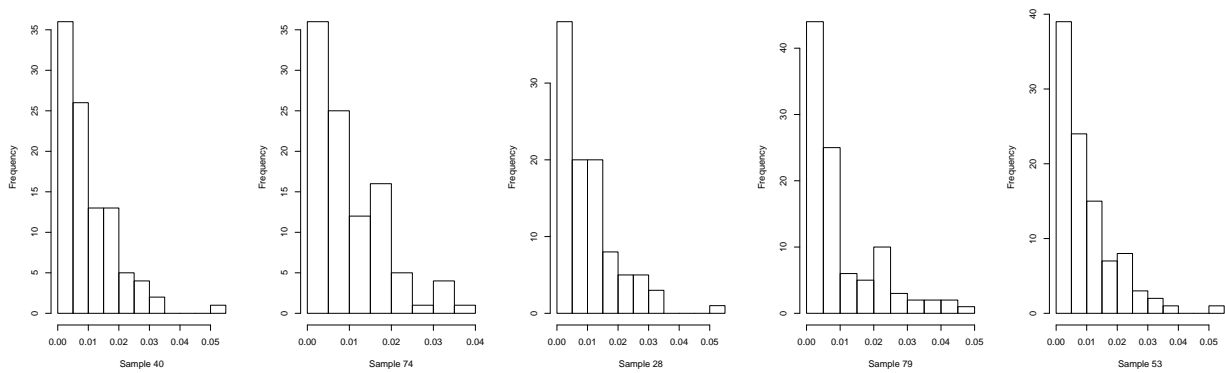


Figure S19: Histogram of the relative abundance of different species for the 5 randomly selected samples for the metagenomic data sets with the pooled sequencing depth at 160x, the number of samples at 80, and the number of species at 100.

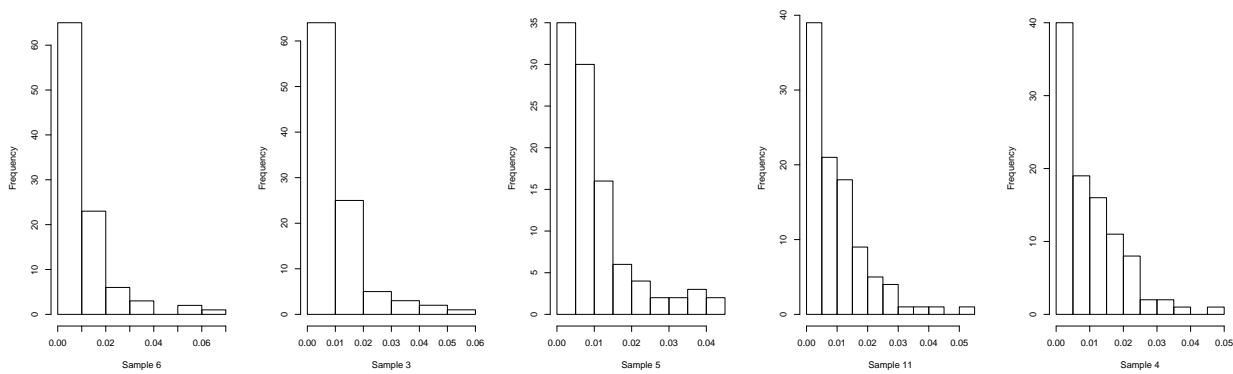


Figure S20: Histogram of the relative abundance of different species for the 5 randomly selected samples for the metagenomic data sets with the pooled sequencing depth at 120x, the number of samples at 20, and the number of species at 100.

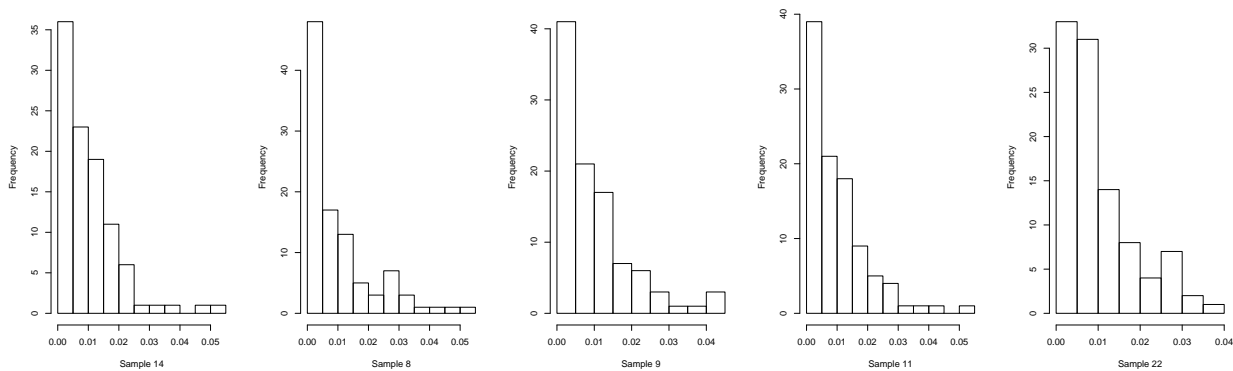


Figure S21: Histogram of the relative abundance of different species for the 5 randomly selected samples for the metagenomic data sets with the pooled sequencing depth at 120x, the number of samples at 40, and the number of species at 100.

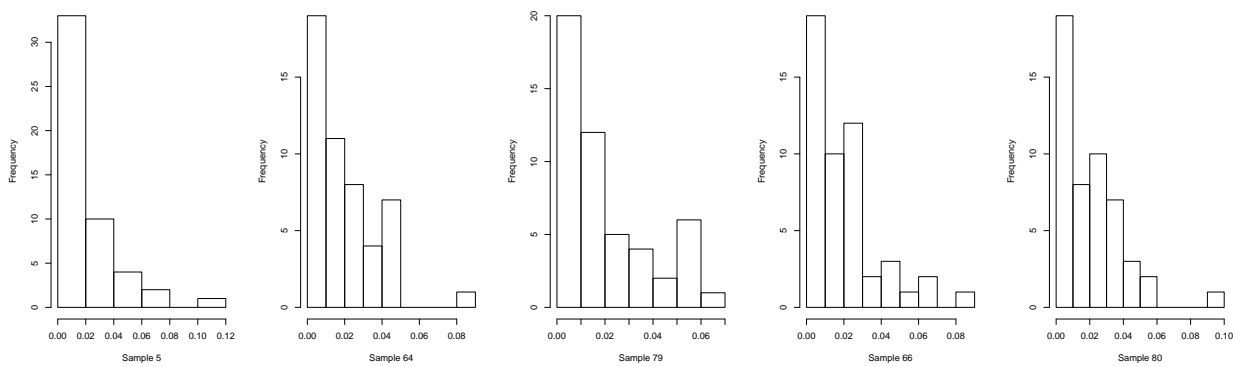


Figure S22: Histogram of the relative abundance of different species for the 5 randomly selected samples for the metagenomic data sets with the pooled sequencing depth at 120x, the number of samples at 80, and the number of species at 50.

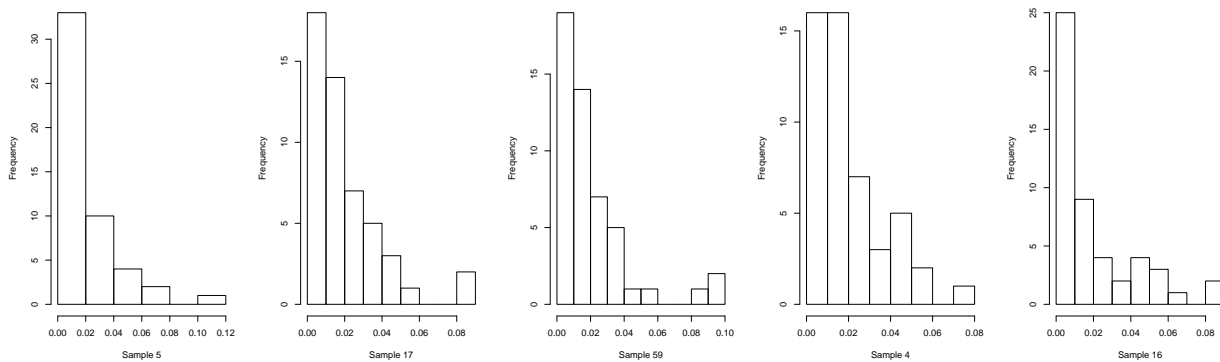


Figure S23: Histogram of the relative abundance of different species for the 5 randomly selected samples for the metagenomic data sets with the pooled sequencing depth at 120x, the number of samples at 80, and the number of species at 150.

2 Software parameter settings

We used the default parameter settings for MetaGen, MaxBin and MetaBAT for automatically selecting the number of species. For CONCOCT, we set the true number of species as the maximum number of species and kept the other settings as default. For the reference-based method CLARK, we selected the “bacteria” database and set the taxonomy rank as “species”. To map the contigs assembled in our real data examples, we searched the NCBI nucleotide database and used TAXAassign (<https://github.com/umerijaz/TAXAassign>) to assign them to taxonomic groups. We set the thresholds for the minimum percentage identity as 60, 70, 80, 95, 95, 97 for phylum, class, order, family, genus, species, respectively.

In the simulation study, we ran MetaGen with default pipeline in our software user manual to do the pooled assembly (with Ray or MEGAHIT assembler) and extract the read counting mapping matrix. For simulation 1-4 and 6-7, we use parameter settings `-s 10 -i 10 -o 2` (with `-t 0.05` for the case with number of sample less or equal than 10). For simulation setting 5, we use the parameter settings `-t 0.01 -s 10 -i 10 -o 2`. As documented in our user manual, the option `-t` is used to specify the threshold for the initial values. It is recommended to set this number smaller(0.01 -0.1) when the number of samples is less than 10 and (0.1-0.2) when the number of samples is larger than 10. The default value is set as 0.1. The option `-s` sets the minimum number of bins. The option `-i` sets the step size to search the optimal number of species using BIC criterion.

3 Real Data Sets

Metagenomic Analysis of IBD

In this study, we assembled the paired-end reads published in Qin et al. [4] using Ray assembler [5], which generated a total of 3,476,781 contigs. We applied MetaGen on the 71,4582 contigs that are longer than 500bp and identified 2,150 clusters according to our BIC scores (Figure S24). We also aligned all the contigs to the NCBI nucleotide database using BLAST, and found that 47,565 contigs have a close match to the species in the database.

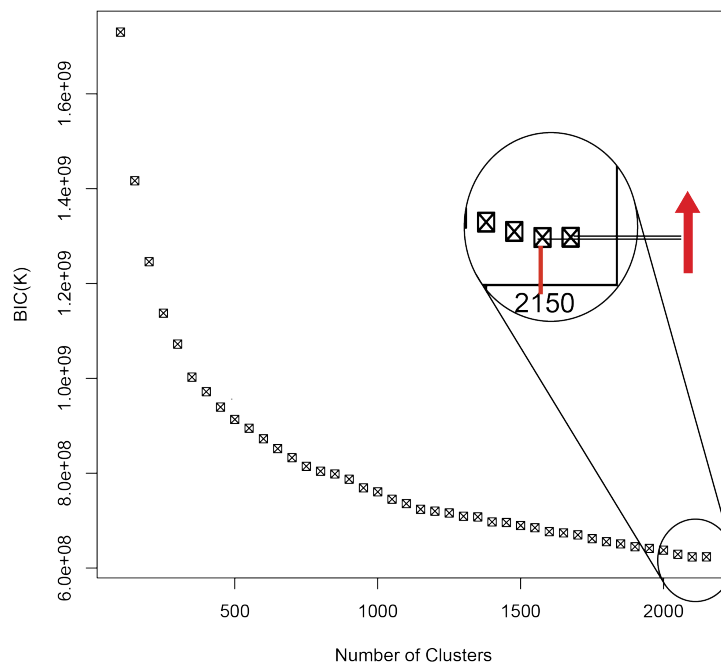


Figure S24: BIC scores of MetaGen for the IBD dataset containing 124 metagenomic samples.

Metagenomic Analysis for T2D

In this study, we assembled the paired-end reads published in Qin et al. [6] using Ray assembler [5], which generated a total of 465,496 contigs that are longer than 500bp. In Figure **S25**, the BIC score is minimized when the number of clusters reaches 2450. Among all the contigs, 44,297 can be aligned to reference genomes in the NCBI nucleotide database using BLAST.

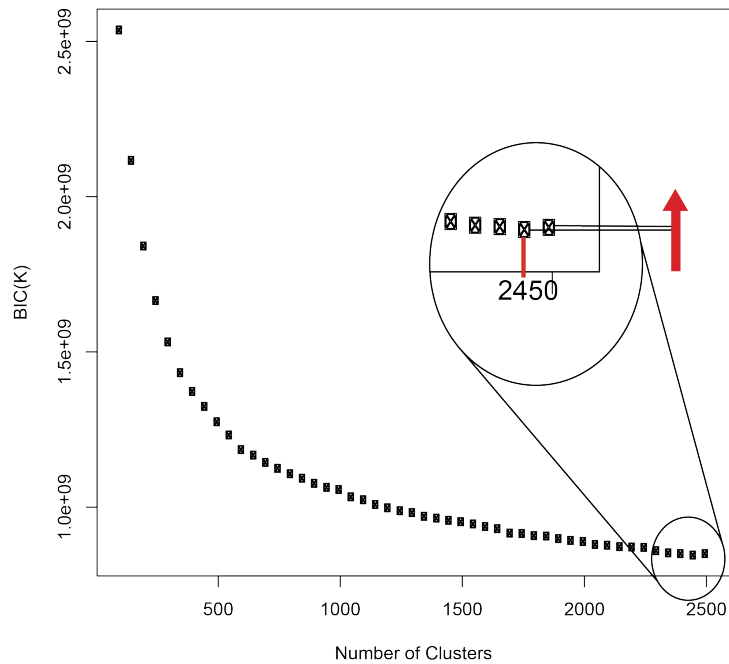


Figure S25: BIC scores of MetaGen for the T2D dataset containing 145 metagenomic samples.

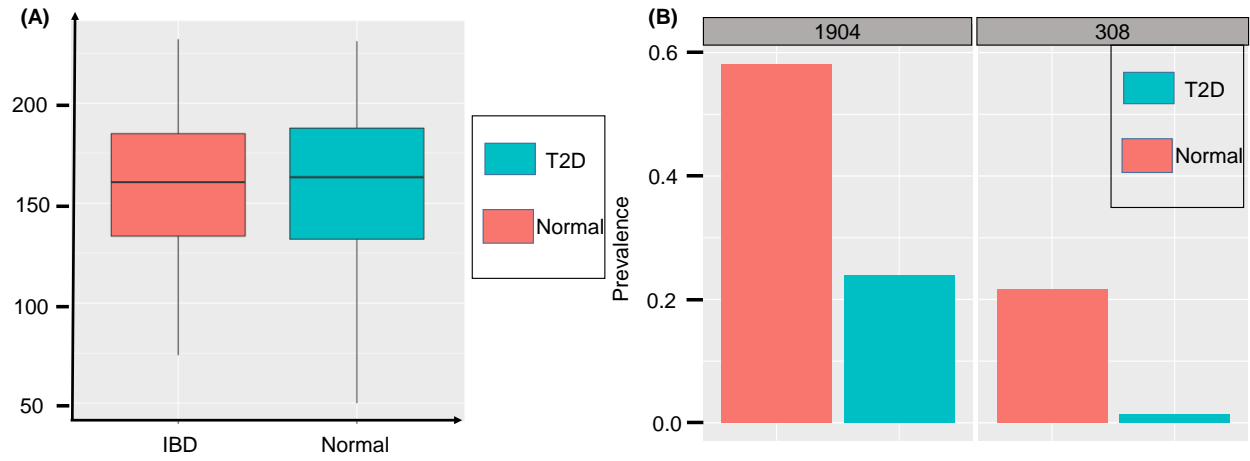


Figure S26: The boxes in (A) represent the distributions of the number of bacterial species in individuals of the normal control and T2D groups, respectively. Plotted in (B) are the prevalences of the 2 clusters (1904 and 308) that are less commonly seen in T2D patients compared to normal controls.

Metagenomic Analysis for Obesity

The DNA samples were sequenced using Pyrosequencing 454 with 9,395,811 total reads for all samples. The pooled short reads from 18 samples were assembled into contigs by the Ray Assembler [5]. A total of 712,751 contigs were constructed, among which 33,107 are longer than 500bp. We further removed 7,724 contigs with fewer than 10 total mapped reads.

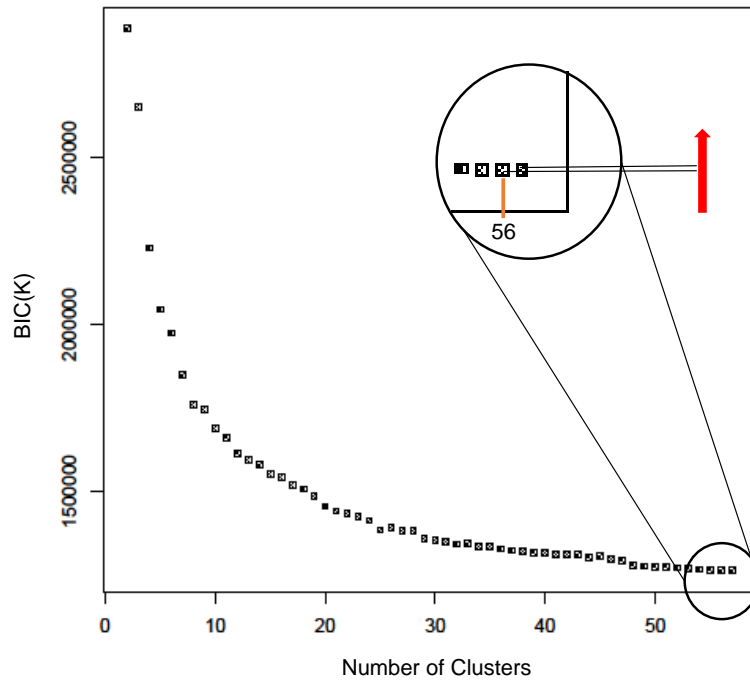


Figure S27: BIC scores of MetaGen for the obesity dataset containing 18 metagenomic samples.

4 Supplementary Note

LASSO-logistic regression for real data

Let y_j denote the group membership of the j th individual, where $y_j = 0$ indicates the control group and $y_j = 1$, the patient group. The logistic regression model assumes that

$$\log \frac{Pr(y_j = 0)}{Pr(y_j = 1)} = \beta_0 + \beta_1 \log(\hat{b}_{j1}) + \cdots + \beta_K \log(\hat{b}_{jK}) \quad (1)$$

where \hat{b}_{jk} is the relative species abundance (see **Materials and Methods** (7)). To alleviate the overfitting, we employed LASSO penalty on β_0, \dots, β_K by assuming $\sum_{k=1}^K |\beta_k| \leq \tau$, where τ is selected by leave-one-out cross validation [7].

Calculating the phylum level relative abundance

Bins are pooled into one phylum if more than 50% of contigs in these bins are mapped to that phylum. For the obesity data set, there are 25 bins that can be mapped to one of the four given phyla. We calculated the relative abundance (see **Materials and Methods**) for each phylum by aggregating the relative abundance of all identified species in this phylum.

Prediction result with CLARK

We applied the same prediction procedure to predict the disease state using the relative abundance outputted by CLARK. The leave-one-out of cross validation error rate is 0.137 and 0.290 for the IBD and T2D data set, which is higher than the error rate predicted by MetaGen.

Supplementary Tables for species name and taxon ID

Species Name	Taxon ID	Species Name	Taxon ID
<i>Acholeplasma laidlawii</i> PG-8A	441768	<i>Achromobacter xylosoxidans</i> A8	762376
<i>Achromobacter xylosoxidans</i> NH44784-1996	1167634	<i>Akkermansia muciniphila</i> ATCC BAA 835	349741
<i>Alistipes finegoldii</i> DSM 17242	679935	<i>Alistipes shahii</i> WAL 8301	717959
<i>Alkalilimnicola ehrlichii</i> MLHE-1	187272	<i>Alkaliphilus metalliredigens</i> QYMF	293826
<i>Asticcacaulis excentricus</i> CB 48	573065	<i>Bacillus amyloliquefaciens</i> XH7	1034836
<i>Bacillus amyloliquefaciens</i> Y2	1126211	<i>Bacillus cereus</i> E33L	288681
<i>Bacillus thuringiensis</i> serovar konkukian 97 27	1279365	<i>Bacteroides fragilis</i> 638R	862962
<i>Bacteroides fragilis</i> NCTC 9343	272559	<i>Bacteroides fragilis</i> YCH46	295405
<i>Bacteroides helcogenes</i> P 36 108	693979	<i>Bacteroides thetaiotaomicron</i> VPI 5482	226186
<i>Bacteroides vulgatus</i> ATCC 8482	435590	<i>Bartonella australis</i> Aust/NH1	1094489
<i>Bartonella bacilliformis</i> KC583	360095	<i>Bartonella clarridgeiae</i> 73	696125
<i>Bifidobacterium adolescentis</i> ATCC 15703	367928	<i>Bifidobacterium bifidum</i> BGN4	484020
<i>Bifidobacterium bifidum</i> PRL2010 chromosome	702459	<i>Bifidobacterium breve</i> UCC2003	326426
<i>Bifidobacterium longum</i>	205913	<i>Brevibacillus brevis</i> NBRC 100599	358681
<i>Brevundimonas subvibrioides</i> ATCC 15264	633149	Butyrate producing bacterium SSC 2	245018
<i>Butyrivibrio proteoclasticus</i> B316	515622	Candidatus <i>Accumulibacter phosphatis</i> clade IIA UW 1	511995
Candidatus <i>Arthromitus</i> sp. SFB-mouse-Yit	1041809	Candidatus <i>Chloracidobacterium thermophilum</i> B	981222
Candidatus <i>Liberibacter solanacearum</i> CLso ZC1	658172	Candidatus <i>Nitrospira defluvii</i>	330214
Candidatus <i>Pelagibacter</i> sp. IMCC9063	1002672	Candidatus <i>Phytoplasma mali</i>	37692
Candidatus <i>Protochlamydia amoebophila</i> UWE25	264201	Candidatus <i>Rickettsia amblyommii</i> str. GAT-30V	1105111
<i>Capnocytophaga canimorsus</i> Cc5	860228	<i>Clavibacter michiganensis</i> sepedonicus	31964
<i>Clostridium acetobutylicum</i> EA 2018	863638	<i>Clostridium acidurici</i> 9a	1128398
<i>Clostridium autoethanogenum</i> DSM 10061	1341692	<i>Clostridium beijerinckii</i>	290402
<i>Clostridium botulinum</i> A str. ATCC 3502	413999	<i>Clostridium botulinum</i> B str. Eklund 17B	935198
<i>Clostridium botulinum</i> H04402 065	941968	<i>Clostridium cellulolyticum</i> H10	394503
<i>Clostridium cellulovorans</i> 743B	573061	<i>Clostridium cf saccharolyticum</i> K10	717608
<i>Clostridium clariflavum</i> DSM 19732	720554	<i>Clostridium difficile</i> R20291	645463
<i>Clostridium phytofermentans</i> ISDg	357809	<i>Clostridium saccharolyticum</i> WM1	610130
<i>Clostridium stercoarium</i> DSM 8532	1121335	<i>Clostridium</i> SY8519	1042156
<i>Clostridium thermocellum</i> DSM 1313	637887	<i>Coprococcus</i> ART55 1	751585

<i>Coprococcus catus</i> GD 7	717962	<i>Coriobacterium glomerans</i> PW2	700015
<i>Corynebacterium argentoratense</i> DSM 44202	1348662	<i>Cryptobacterium curtum</i> DSM 15641	469378
<i>Cupriavidus metallidurans</i> CH34	266264	<i>Cylindrospermum stagnale</i> PCC 7417	56107
<i>Dehalobacter</i> CF	1131462	<i>Denitrovibrio acetiphilus</i> DSM 12809	522772
<i>Desulfarculus baarsii</i> DSM 2075	644282	<i>Desulfatibacillum alkenivorans</i> AK-01	439235
<i>Desulfitobacterium dehalogenans</i> ATCC 51507	756499	<i>Desulfitobacterium dichloroeliminans</i> LMG P 21439	871963
<i>Desulfotomaculum acetoxidans</i> DSM 11109	880072	<i>Echinicola vietnamensis</i> DSM 17526	926556
<i>Enterococcus faecium</i> Aus0085	1305849	<i>Enterococcus hirae</i> ATCC 9790	768486
<i>Enterococcus mundtii</i> QU 25	1300150	<i>Erysipelothrix rhusiopathiae</i> SY1027	1313290
<i>Escherichia coli</i> LF82	591946	<i>Escherichia coli</i> NA114	1033813
<i>Ethanoligenens harbinense</i> YUAN 3	663278	<i>Eubacterium cylindroides</i> T2 87	717960
<i>Eubacterium eligens</i> ATCC 27750	515620	<i>Eubacterium rectale</i> ATCC 33656	515619
<i>Eubacterium siraeum</i> V10Sc8a	717961	<i>Faecalibacterium prausnitzii</i> L2-6	718252
<i>Ferrimonas balearica</i> DSM 9799	550540	<i>Fibrobacter succinogenes</i> subsp. <i>succinogenes</i> S85	59374
<i>Filifactor alocis</i> ATCC 35896	546269	<i>Finegoldia magna</i> ATCC 29328	334413
<i>Flavobacteriaceae bacterium</i> 3519-10	531844	<i>Fusobacterium nucleatum</i> ATCC 25586	190304
<i>Haemophilus parainfluenzae</i> T3T1	862965	<i>Herbaspirillum seropedicae</i> SmR1	757424
<i>Janthinobacterium</i> Marseille	375286	<i>Lawsonia intracellularis</i> PHE MN1 00	363253
<i>Megasphaera elsdenii</i> DSM 20460	1064535	<i>Mesoplasma florum</i> L1	265311
<i>Moorella thermoacetica</i> ATCC 39073	264732	<i>Mycoplasma mycoides capri</i> LC 95010	862259
<i>Mycoplasma putrefaciens</i> Mput9231	1292033	<i>Mycoplasma suis</i> KI3806	708248
<i>Odoribacter splanchnicus</i> DSM 20712	709991	<i>Oscillibacter valericigenes</i> Sjm18-20	693746
<i>Paludibacter propionicigenes</i> WB4	694427	<i>Parabacteroides distasonis</i> ATCC 8503	435591
<i>Pediococcus pentosaceus</i> SL4	1408206	<i>Pedobacter heparinus</i> DSM 2366	485917
<i>Pedobacter saltans</i> DSM 12145	762903	<i>Pelagibacterium halotolerans</i> B2	1082931
<i>Porphyromonas gingivalis</i> ATCC 33277	431947	<i>Prevotella denticola</i> F0289	767031
<i>Prevotella intermedia</i> 17	246198	<i>Prevotella melaninogenica</i> ATCC 25845	553174
<i>Ralstonia eutropha</i> H16	381666	<i>Ralstonia eutropha</i> JMP134	264198
<i>Ralstonia solanacearum</i> CMR15	859655	<i>Rhodospirillum rubrum</i> F11	1036743
<i>Rivularia</i> PCC 7116	373994	<i>Roseburia intestinalis</i> XB6B4	718255
<i>Ruminococcus albus</i> 7	697329	<i>Ruminococcus bromii</i> L2-63	657321
<i>Ruminococcus champanellensis</i> 18P13	213810	<i>Ruminococcus obeum</i>	657314
<i>Ruminococcus</i> sp. SR1/5	657323	<i>Ruminococcus torques</i> L2-14	657313
<i>Selenomonas ruminantium lactilytica</i> TAM6421	927704	<i>Selenomonas sputigena</i> ATCC 35185	546271
<i>Sphingobacterium</i> 21	743722	<i>Spiroplasma chrysopicola</i> DF-1	1276227
<i>Spiroplasma diminutum</i> CUAS 1	1276221	<i>Spiroplasma syrphidicola</i> EA 1	1276229
<i>Spiroplasma taiwanense</i> CT 1	1276220	<i>Streptococcus parasangnis</i> FW213	1114965

Streptococcus pseudopneumoniae IS7493	1054460	Streptococcus thermophilus ND03	767463
Tannerella forsythia ATCC 43037	203275	Tetragenococcus halophilus NBRC 12172	945021
Tolomonas auensis DSM 9187	595494	Variovorax paradoxus B4	1246301
Variovorax paradoxus EPS	595537	Variovorax paradoxus S110	543728
Veillonella parvula DSM 2008	479436	Yersinia enterocolitica palearctica 105 5R r	994476
Yersinia enterocolitica subsp. enterocolitica 8081	393305	Yersinia enterocolitica subsp. palearctica Y11	930944
Yersinia pestis A1122	1035377	Yersinia pestis Angola	349746
Yersinia pestis Antiqua	360102	Yersinia pseudotuberculosis IP 32953	273123

Table 1: 150 genomes used in the simulation study with species names and NCBI taxon id.

Species Name	Taxon ID	Species Name	Taxon ID
E coli 536	362663	E coli APEC O1	405955
E coli IAI1	585034	E coli SE11	409438
E coli DH1	536056	E coli HS	331112

Table 2: 6 *E coli* strains used in the simulation setting 4

Species Name	Taxon ID	Species Name	Taxon ID
<i>Eubacterium rectale</i>	39491	<i>Propionibacterium freudenreichii</i>	1744
<i>Bacteroides vulgatus</i>	821	<i>Conexibacter woesei</i>	191495
<i>Roseburia hominis</i>	301301	<i>Alkalilimnicola ehrlichii</i>	351052
<i>Akkermansia muciniphila</i>	239935	<i>Geobacter uraniireducens</i>	351604
<i>Alistipes finegoldii</i>	214856	<i>Laribacter hongkongensis</i>	168471
<i>Alistipes shahii</i>	328814	<i>Streptococcus anginosus</i>	1328
<i>Eubacterium eligens</i>	39485	<i>Syntrophomonas wolfei</i>	863
<i>Bacteroides thetaiotaomicron</i>	818	<i>Aeromonas hydrophila</i>	644
<i>Odoribacter splanchnicus</i>	28118	<i>Bacillus megaterium</i>	1404
<i>Bacteroides fragilis</i>	817	<i>Corynebacterium aurimucosum</i>	169292
<i>Parabacteroides distasonis</i>	823	<i>Salinibacter ruber</i>	146919
<i>Bacteroides helcogenes</i>	290053	<i>Weissella koreensis</i>	165096
<i>Bifidobacterium adolescentis</i>	1680	<i>Alicyclophilus denitrificans</i>	179636
<i>Prevotella melaninogenica</i>	28132	<i>Rhodospirillum rubrum</i>	1085
<i>Escherichia coli</i>	562	<i>Shigella flexneri</i>	623
<i>Prevotella ruminicola</i>	839	<i>Zobellia galactanivorans</i>	63186
<i>Clostridium</i> sp. SY8519	1042156	<i>Asticcacaulis excentricus</i>	78587
<i>Haemophilus parainfluenzae</i>	729	<i>Herbaspirillum seropedicae</i>	964
<i>Clostridium saccharolyticum</i>	84030	<i>Polaribacter</i> sp. MED152	313598
<i>Oscillibacter valericigenes</i>	351091	<i>Thermoanaerobacterium</i> xy- <i>lanolyticum</i>	29329
<i>Prevotella dentalis</i>	52227	<i>Marivirga tractuosa</i>	1006
<i>Bacteroides salanitronis</i>	376805	<i>Phycisphaera mikurensis</i>	547188
<i>Bifidobacterium longum</i>	216816	<i>Deinococcus proteolyticus</i>	55148
<i>Prevotella intermedia</i>	28131	<i>Geobacter sulfurreducens</i>	35554
<i>Streptococcus suis</i>	1307	<i>Spirochaeta thermophila</i>	154
<i>Prevotella denticola</i>	28129	<i>Streptococcus</i> sp. I-G2	1156431
<i>Prevotella</i> sp. oral taxon 299	652716	<i>Streptococcus intermedius</i>	1338
<i>Ruminococcus champanellensis</i>	1161942	<i>Alkaliphilus metalliredigens</i>	208226
<i>Methanobrevibacter smithii</i>	2173	<i>Brachyspira pilosicoli</i>	52584
<i>Streptococcus thermophilus</i>	1308	<i>Desulfovibrio salexigens</i>	880
<i>Adlercreutzia equolifaciens</i>	446660	<i>Synechococcus</i> sp. WH 7803	32051
<i>Ethanoligenens harbinense</i>	253239	<i>Advenella kashmirensis</i>	310575
<i>Bifidobacterium bifidum</i>	1681	<i>Sinorhizobium medicae</i>	110321
<i>Butyrivibrio proteoclasticus</i>	43305	<i>Streptococcus oralis</i>	1303
<i>Ruminococcus albus</i>	1264	<i>Acetohalobium arabaticum</i>	28187
<i>Enterococcus faecium</i>	1352	<i>Anaeromyxobacter</i> sp. Fw109-5	404589
<i>Eubacterium limosum</i>	1736	<i>Cupriavidus taiwanensis</i>	164546
<i>Lactobacillus delbrueckii</i>	1584	<i>Micavibrio aeruginosavorus</i>	349221
<i>Streptococcus salivarius</i>	1304	<i>Streptococcus agalactiae</i>	1311

<i>Streptococcus lutetiensis</i>	150055	<i>Corynebacterium variabile</i>	1727
<i>Tannerella forsythia</i>	28112	<i>Xanthobacter autotrophicus</i>	280
<i>Gordonibacter pamelaee</i>	471189	<i>Phenylobacterium zucineum</i>	284016
<i>Veillonella parvula</i>	29466	<i>Ramlibacter tataouinensis</i>	94132
<i>Bifidobacterium breve</i>	1685	<i>Paenibacillus</i> sp. JDR-2	324057
<i>Desulfovibrio desulfuricans</i>	876	<i>Spirochaeta africana</i>	46355
<i>Campylobacter jejuni</i>	197	<i>Thermosediminibacter oceani</i>	291990
<i>Eggerthella lenta</i>	84112	<i>Pseudomonas mendocina</i>	300
<i>Streptococcus parasanguinis</i>	1318	<i>Caldicellulosiruptor kronotskyensis</i>	413889
<i>Clostridium botulinum</i>	1491	<i>Deinococcus maricopensis</i>	309887
<i>Desulfovibrio vulgaris</i>	881	<i>Pseudomonas protegens</i>	380021
<i>Acidaminococcus fermentans</i>	905	<i>Anaeromyxobacter dehalogenans</i>	161493
<i>Enterococcus faecalis</i>	1351	<i>Rhizobium</i> sp. IRBG74	424182
<i>Acidaminococcus intestini</i>	187327	<i>Thermanaerovibrio acidaminovorans</i>	81462
<i>Porphyromonas gingivalis</i>	837	<i>Pseudoxanthomonas suwonensis</i>	314722
<i>Haemophilus influenzae</i>	727	<i>Verminephrobacter eiseniae</i>	364317
<i>Treponema succinifaciens</i>	167	<i>Syntrophobacter fumaroxidans</i>	119484
<i>Eggerthella</i> sp. YY7918	502558	<i>Treponema azotonutricium</i>	150829
<i>Streptococcus pyogenes</i>	1314	<i>Treponema brennaboreense</i>	81028
<i>Lactobacillus ruminis</i>	1623	<i>Beijerinckia indica</i>	533
<i>Slackia heliotrinireducens</i>	84110	<i>Acholeplasma brassicae</i>	61635
<i>Olsenella uli</i>	133926	<i>Actinobacillus pleuropneumoniae</i>	715
<i>Lactococcus lactis</i>	1358	<i>Oligotropha carboxidovorans</i>	40137
<i>Salmonella enterica</i>	28901	<i>Wolinella succinogenes</i>	844
<i>Streptococcus pasteurianus</i>	197614	<i>Desulfitobacterium dehalogenans</i>	36854
<i>Porphyromonas asaccharolytica</i>	28123	<i>Rhodococcus equi</i>	43767
<i>Solitaea canadensis</i>	995	<i>Arthrobacter</i> sp. FB24	290399
<i>Bifidobacterium animalis</i>	28025	<i>Magnetococcus marinus</i>	1124597
<i>Selenomonas ruminantium</i>	971	[<i>Bacillus</i>] <i>selenitireducens</i>	85683
<i>Coriobacterium glomerans</i>	33871	<i>Chromohalobacter salexigens</i>	158080
<i>Megasphaera elsdenii</i>	907	<i>Desulfosporosinus acidiphilus</i>	885581
<i>Sorangium cellulosum</i>	56	<i>Methylocystis</i> sp. SC2	187303
<i>Selenomonas sputigena</i>	69823	<i>Corynebacterium glutamicum</i>	1718
<i>Treponema denticola</i>	158	<i>Acidovorax</i> sp. KKS102	358220
<i>Paenibacillus mucilaginosus</i>	61624	<i>Chlorobium limicola</i>	1092
<i>Heliobacterium modesticaldum</i>	35701	<i>Rhizobium etli</i>	29449
<i>Paludibacter propionicigenes</i>	185300	<i>Xanthomonas campestris</i>	339
<i>Bifidobacterium dentium</i>	1689	<i>Brachybacterium faecium</i>	43669
<i>Saprospira grandis</i>	1008	<i>Gloeobacter violaceus</i>	33072
<i>Riemerella anatipestifer</i>	34085	<i>Kocuria rhizophila</i>	72000
<i>Lawsonia intracellularis</i>	29546	<i>Acholeplasma palmae</i>	38986
<i>Streptococcus mitis</i>	28037	<i>Leifsonia xyli</i>	1575

Desulfovibrio gigas	879	Stackebrandtia nassauensis	283811
Syntrophobotulus glycolicus	51197	Thermus thermophilus	274
Gardnerella vaginalis	2702	Actinoplanes missouriensis	1866
Desulfotomaculum acetoxidans	58138	Streptobacillus moniliformis	34105
Anaerococcus prevotii	33034	Dichelobacter nodosus	870
Symbiobacterium thermophilum	2734	Pseudomonas syringae	317
Chitinophaga pinensis	79329	Corynebacterium halotolerans	225326
Clostridium tetani	1513	Klebsiella pneumoniae	573
Mycobacterium tuberculosis	1773	Macrococcus caseolyticus	69966
Clostridium cellulovorans	1493	Streptococcus pseudopneumoniae	257758
Pseudomonas stutzeri	316	Thermacetogenium phaeum	85874
Clostridium perfringens	1502	Tistrella mobilis	171437
Candidatus Azobacteroides pseudotriconymphae	511435	Wolbachia sp. wRi	66084
Clostridium stercoarium	1510	Calditerrivibrio nitroreducens	477976
Filifactor alocis	143361	Candidatus Puniceispirillum marinum	767892
Cytophaga hutchinsonii	985	Frateuria aurantia	81475
Clostridium kluveri	1534	Micrococcus luteus	1270
Desulfovibrio magneticus	184917	Neisseria meningitidis	487
Clostridium saccharobutylicum	169679	Azorhizobium caulinodans	7
Clostridium clariflavum	288965	Azotobacter vinelandii	354
Aggregatibacter aphrophilus	732	Belliella baltica	232259
Rhodopseudomonas palustris	1076	Marinithermus hydrothermalis	186192
Desulfarculus baarsii	453230	Paenibacillus terrae	159743
Fusobacterium nucleatum	851	Thiobacillus denitrificans	36861
Pseudomonas putida	303	Candidatus Arthromitus sp. SFB-mouse	49118
Fibrella aestuarina	651143	Myxococcus stipitatus	83455
Mycoplasma hyopneumoniae	2099	Pseudovibrio sp. FO-BEG1	911045
Clostridium saccharoperbutylacetonicum	36745	Candidatus Accumulibacter phosphatis	327160
Finegoldia magna	1260	Catenulispora acidiphila	304895
Flavobacterium johnsoniae	986	Geobacter bemidjiensis	225194
Desulfomicrobium baculatum	899	Oceanimonas sp. GK1	511062
Bifidobacterium thermophilum	33905	Spiroplasma taiwanense	2145
Clostridium sp. BNL1100	755731	Elusimicrobium minutum	423605
Brevibacillus brevis	1393	Novosphingobium sp. PP1Y	702113
Thermaerobacter marianensis	73919	Clavibacter michiganensis	28447
Desulfotomaculum gibsoniae	102134	Leptotrichia buccalis	40542
Clostridium cellulolyticum	1521	Myxococcus xanthus	34
Niastella koreensis	354356	Saccharothrix espanaensis	103731
Halothiobacillus neapolitanus	927	Flavobacterium branchiophilum	55197

<i>Clostridium acetobutylicum</i>	1488	<i>Gemmatimonas aurantiaca</i>	173480
<i>Thermobacillus composti</i>	377615	<i>Polymorphum gilvum</i>	991904
<i>Acetobacterium woodii</i>	33952	<i>Bacillus cereus</i>	1396
<i>Desulfovibrio africanus</i>	873	<i>Caldilinea aerophila</i>	133453
<i>Capnocytophaga ochracea</i>	1018	<i>Cyanothece</i> sp. ATCC 51142	43989
<i>Pyrolobus fumarii</i>	54252	<i>Ruegeria pomeroyi</i>	89184
<i>Fibrobacter succinogenes</i>	833	<i>Arthrospira platensis</i>	118562
<i>Mahella australiensis</i>	252966	<i>Desulfohalobium retbaense</i>	45663
<i>Pedobacter heparinus</i>	984	<i>Geobacillus</i> sp. C56-T3	691437
<i>Aeromonas salmonicida</i>	645	<i>Myxococcus fulvus</i>	33
<i>Desulfurispirillum indicum</i>	936456	<i>Azoarcus</i> sp. KH32C	748247
<i>Dyadobacter fermentans</i>	94254	<i>Caldisericum exile</i>	693075
<i>Lactobacillus acidophilus</i>	1579	<i>Azospirillum brasilense</i>	192
<i>Melioribacter roseus</i>	1134405	<i>Brevundimonas subvibrioides</i>	74313
<i>Dehalococcoides mccartyi</i>	61435	<i>Desulfosporosinus orientis</i>	1563
<i>Enterobacter cloacae</i>	550	<i>Marinobacter hydrocarbonoclasticus</i>	2743
<i>Pseudomonas aeruginosa</i>	287	<i>Streptomyces cattleya</i>	29303
<i>Serratia marcescens</i>	615	<i>Agrobacterium fabrum</i>	1176649
<i>Achromobacter xylosoxidans</i>	85698	<i>Polynucleobacter necessarius</i>	576610
<i>Azospirillum lipoferum</i>	193	<i>Thermobifida fusca</i>	2021
<i>Clostridium novyi</i>	1542	<i>Burkholderia pseudomallei</i>	28450
<i>Desulfotobacterium hafniense</i>	49338	<i>Cupriavidus metallidurans</i>	119219
<i>Desulfotomaculum ruminis</i>	1564	<i>Ferrimonas balearica</i>	44012
<i>Streptococcus pneumoniae</i>	1313	<i>Cellvibrio japonicus</i>	155077
<i>Bacillus coagulans</i>	1398	<i>Sideroxydans lithotrophicus</i>	63745
<i>Desulfovibrio aespoensis</i>	182210	<i>Burkholderia</i> sp. YI23	1097668
<i>Rhizobium leguminosarum</i>	384	<i>Deinococcus peraridilitoris</i>	432329
<i>Ralstonia solanacearum</i>	305	<i>Geitlerinema</i> sp. PCC 7407	1173025
<i>Haliangium ochraceum</i>	80816	<i>Mycoplasma agalactiae</i>	2110
<i>Desulfotobacterium dichloroeliminans</i>	233055	<i>Caldicellulosiruptor saccharolyticus</i>	44001
<i>Escherichia fergusonii</i>	564	<i>Mycobacterium abscessus</i>	36809
<i>Variovorax paradoxus</i>	34073	<i>Streptomyces bingchenggensis</i>	379067
<i>Cyanobium gracile</i>	59930	<i>Candidatus Amoebophilus asiaticus</i>	281120
<i>Atopobium parvulum</i>	1382	<i>Carnobacterium maltaromaticum</i>	2751
<i>Erysipelothrix rhusiopathiae</i>	1648	<i>Paracoccus denitrificans</i>	266
<i>Chromobacterium violaceum</i>	536	<i>Sphingomonas</i> sp. MM-1	745310
<i>Sphaerobacter thermophilus</i>	2057	<i>Zymomonas mobilis</i>	542
<i>Desulfovibrio alaskensis</i>	58180	<i>Chamaesiphon minutus</i>	1173032
<i>Alicyclobacillus acidocaldarius</i>	405212	<i>Pelobacter carbinolicus</i>	19
<i>Clostridium beijerinckii</i>	1520	<i>Sebaldella termitidis</i>	826
<i>Staphylococcus aureus</i>	1280	<i>Singulisphaera acidiphila</i>	466153
<i>Geobacter</i> sp. M21	443144	<i>Sphingobium</i> sp. SYK-6	627192

Rhodothermus marinus	29549	Burkholderia thailandensis	57975
Sphingobacterium sp. 21	743722	Gluconacetobacter diazotrophicus	33996
Mycoplasma pulmonis	2107	Kytococcus sedentarius	1276
Clostridium pasteurianum	1501	Nitratifactor salsuginis	269261
Rubrivivax gelatinosus	28068	Staphylococcus saprophyticus	29385
Desulfatibacillum alkenivorans	259354	Synechococcus sp. JA-3-3Ab	321327
Francisella tularensis	263	Arthrobacter sp. Rue61a	1118963
Aromatoleum aromaticum	551760	Coprothermobacter proteolyticus	35786
Mobiluncus curtisii	2051	Flavobacteriaceae bacterium 3519-10	531844
Rhodospirillum centenum	34018	Shewanella amazonensis	60478
Desulfobacterium autotrophicum	2296	Syntrophus aciditrophicus	316277
Desulfobulbus propionicus	894	Thermoanaerobacter wiegeli	46354
Candidatus Desulforudis audaxviator	471827	Deinococcus radiodurans	1299
Candidatus Solibacter usitatus	332163	Rhizobium tropici	398
Runella slithyformis	106	Yersinia enterocolitica	630
Geobacter metallireducens	28232	Candidatus Koribacter versatilis	658062
Alkaliphilus oremlandii	461876	Corynebacterium maris	575200
Leptothrix cholodnii	34029	Hyphomicrobium denitrificans	53399
Moorella thermoacetica	1525	Streptomyces coelicolor	1902
Pseudogulbenkiania sp. NH8B	748280	Bacillus infantis	324767
Beutenbergia cavernae	84757	Desulfomonile tiedjei	2358
Owenweeksia hongkongensis	253245	Granulicella tundricola	940615
Pediococcus pentosaceus	1255	Ilyobacter polytropus	167642
Carboxydotherrnus hydrogenoformans	129958	Rhodobacter capsulatus	1061
Paenibacillus sp. Y412MC10	481743	Bacillus clausii	79880
Tepidanaerobacter acetatoxydans	499229	Gallionella capsiferriformans	370405
Stenotrophomonas maltophilia	40324	Pectobacterium carotovorum	554
Comamonas testosteroni	285	Aequorivita sublithincola	101385
Robiginitalea biformata	252307	Anaerolinea thermophila	167964
Treponema primitia	88058	Halanaerobium praevalens	2331
Paenibacillus larvae	1464	Sphaerochaeta pleomorpha	1131707
Geobacter sp. M18	443143	Erwinia billingiae	182337
Thiomicrospira crunogena	39765	Pelagibacterium halotolerans	531813
Campylobacter coli	195	Anoxybacillus flavithermus	33934
Ignavibacterium album	591197	Collimonas fungivorans	158899
Pelobacter propionicus	29543	Desulfovibrio piezophilus	879567
Thermoanaerobacterium thermosaccharolyticum	1517	Methylobacterium sp. 4-46	426117
Gluconobacter oxydans	442	Rothia dentocariosa	2047
Acholeplasma laidlawii	2148	Sulfuricella denitrificans	649841
Streptococcus sanguinis	1305	Thermomonospora curvata	2020
Streptomyces violaceusniger	68280	Bdellovibrio bacteriovorus	959

Magnetospirillum gryphiswaldense	55518	Burkholderia mallei	13373
Lactobacillus reuteri	1598	Frankia alni	1859
Allochromatium vinosum	1049	Methylomonas methanica	421
Jonesia denitrificans	43674	Pasteurella multocida	747
Thauera sp. MZ1T	85643	Prosthecochloris aestuarii	1102
Bacillus cellulosityticus	1413	Starkeya novella	921
Magnetospirillum magneticum	84159	Terriglobus roseus	392734
Stigmatella aurantiaca	41	Thermus scotoductus	37636
Flavobacterium columnare	996	Chloracidobacterium thermophilum	458033
Methylomicrobium alcaliphilum	271065	Chlorobium phaeovibrioides	1094
Pseudomonas fluorescens	294	Corynebacterium efficiens	152794
Bifidobacterium asteroides	1684	Geobacter daltonii	1203471
Chlorobaculum parvum	274539	Nakamurella multipartita	53461
Geobacter lovleyi	313985	Streptococcus constellatus	76860
Listeria monocytogenes	1639	Ammonifex degensii	42838
Mycoplasma hyorhinitis	2100	Polaromonas sp. JS666	296591
Cupriavidus necator	106590	Rhodococcus erythropolis	1833
Rothia mucilaginosa	43675	Amphibacillus xylanus	1449
Paenibacillus polymyxa	1406	Halobacillus halophilus	1570
Streptococcus gordonii	1302	Methylococcus capsulatus	414
Lactobacillus salivarius	1624	Mycoplasma penetrans	28227
Zunongwangia profunda	398743	Photorhabdus asymbiotica	291112
Streptococcus equi	1336	Hirschia baltica	2724
Thermincola potens	863643	Ochrobactrum anthropi	529
Buchnera aphidicola	9	Blastococcus saxobsidens	138336
Cellulophaga lytica	979	Cryptobacterium curtum	84163
Desulfotomaculum reducens	59610	Natranaerobius thermophilus	375929
Streptosporangium roseum	2001	Nocardia cyriacigeorgica	135487
Coralloccoccus coralloides	184914	Thiomonas intermedia	926
Solibacillus silvestris	76853	Burkholderia sp. RPE64	758793
Desulfurivibrio alkaliphilus	427923	Deinococcus deserti	310783
Gallibacterium anatis	750	Kineococcus radiotolerans	131568
Ralstonia pickettii	329	Mesorhizobium ciceri	39645
Shigella sonnei	624	Vibrio furnissii	29494
Streptococcus sp. I-P16	1156433	Bacillus cytotoxicus	580165
Thiocystis violascens	73141	Deinococcus geothermalis	68909
Citrobacter rodentium	67825	Pantoea sp. At-9b	592316
Prochlorococcus marinus	1219	Sulfurimonas autotrophica	202747
Weeksella virosa	1014	Corynebacterium jeikeium	38289
Shigella boydii	621	Nitrospira defluvi	330214
Deinococcus gobiensis	502394	Echinicola vietnamensis	390884
Haliscomenobacter hydroxiss	2350	Lactococcus garvieae	1363

Leadbetterella byssophila	316068	Nocardiopsis dassonvillei	2014
Oceanobacillus iheyensis	182710	Turneriella parva	29510
Parvibaculum lavamentivorans	256618	Amycolatopsis mediterranei	33910
Desulfococcus oleovorans	181663	Bacillus sp. 1NLA3E	666686
Hyphomonas neptunium	81032	Pseudomonas resinovorans	53412
Micrococcus phosphovorans	29405	Pseudomonas sp. T7-7	1007105
Lactobacillus fermentum	1613	Thermus sp. CCB.US3.UF1	1111069
Oceanithermus profundus	187137	Candidatus Phytoplasma australiense	59748
Rhodobacter sphaeroides	1063	Edwardsiella ictaluri	67780
Rubrobacter xylanophilus	49319	Methylobium petroleiphilum	105560
Staphylococcus haemolyticus	1283	Xenorhabdus bovienii	40576
Thioflavococcus mobilis	80679	Caulobacter segnis	88688
Acetobacter pasteurianus	438	Novosphingobium aromaticivorans	48935
Bacillus subtilis	1423	Ornithobacterium rhinotracheale	28251
Agrobacterium sp. H13-3	861208	Polaromonas naphthalenivorans	216465
Bacillus thuringiensis	1428	Roseiflexus sp. RS-1	357808
Bordetella petrii	94624	Truepera radiovictrix	332249
Nautilia profundicola	244787	Bordetella avium	521
Opitutus terrae	107709	Burkholderia gladioli	28095
Sinorhizobium fredii	380	Erwinia amylovora	552
Meiothermus ruber	277	Exiguobacterium sibiricum	332410
Shigella dysenteriae	622	Frankia sp. EAN1pec	298653
Acidobacterium capsulatum	33075	Melissococcus plutonius	33970
Brachyspira hyodysenteriae	159	Thermobispora bispora	2006
Dechloromonas aromatica	259537	Methylobacterium nodulans	114616
Gordonia sp. KTR9	337191	Sphingomonas wittichii	160791
Methylobacterium extorquens	408	Rhodococcus vanniellii	1069
Rivularia sp. PCC 7116	373994		

Table 3: 545 species used in the simulation setting 5

References

- [1] Daniel C Richter, Felix Ott, Alexander F Auch, Ramona Schmid, and Daniel H Huson. Metasim—a sequencing simulator for genomics and metagenomics. *PloS One*, 3(10):e3373, 2008.
- [2] Weichun Huang, Leping Li, Jason R Myers, and Gabor T Marth. Art: a next-generation sequencing read simulator. *Bioinformatics*, 28(4):593–594, 2012.
- [3] Chengwei Luo, Rob Knight, Heli Siljander, Mikael Knip, Ramnik J Xavier, and Dirk Gevers.

- Constrains identifies microbial strains in metagenomic datasets. *Nature biotechnology*, 33(10): 1045–1052, 2015.
- [4] Junjie Qin, Ruiqiang Li, Jeroen Raes, Manimozhiyan Arumugam, Kristoffer Solvsten Burgdorf, Chaysavanh Manichanh, Trine Nielsen, Nicolas Pons, Florence Levenez, Takuji Yamada, et al. A human gut microbial gene catalogue established by metagenomic sequencing. *nature*, 464 (7285):59–65, 2010.
- [5] Sébastien Boisvert, Frédéric Raymond, Élénie Godzaridis, François Laviolette, Jacques Corbeil, et al. Ray meta: scalable de novo metagenome assembly and profiling. *Genome Biol*, 13(12): R122, 2012.
- [6] Junjie Qin, Yingrui Li, Zhiming Cai, Shenghui Li, Jianfeng Zhu, Fan Zhang, Suisha Liang, Wenwei Zhang, Yuanlin Guan, Dongqian Shen, et al. A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature*, 490(7418):55–60, 2012.
- [7] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.