

MetaBAT, an efficient tool for accurately reconstructing single genomes from complex microbial communities

Dongwan D Kang, Jeff Froula, Rob Egan, Zhong Wang

Grouping large genomic fragments assembled from shotgun metagenomic sequences to deconvolute complex microbial communities, or metagenome binning, enables the study of individual organisms and their interactions. Because of the complex nature of these communities, existing metagenome binning methods often miss a large number of microbial species. In addition, most of the tools are not scalable to large datasets. Here we introduce automated software, called MetaBAT that integrates empirical probabilistic distances of genome abundance and tetranucleotide frequency for accurate metagenome binning. MetaBAT outperforms alternative methods in accuracy and computational efficiency on both synthetic and real metagenome datasets. It automatically forms hundreds of high quality genome bins on a very large assembly consisting millions of contigs in a matter of hours on a single node. MetaBAT is open source software and available at <https://bitbucket.org/berkeleylab/metabat> .

2 **MetaBAT, An Efficient Tool for Accurately Reconstructing Single** 3 **Genomes from Complex Microbial Communities**

4 Dongwan D. Kang^{1,2}, Jeff Froula^{1,2}, Rob Egan^{1,2}, and Zhong Wang^{1,2,3*}

5 ¹ Department of Energy, Joint Genome Institute, Walnut Creek, CA 94598, USA

6 ² Genomics Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

7 ³ School of Natural Sciences, University of California at Merced, Merced, CA, 95343, USA;

8 * To whom correspondence should be addressed. Tel: +1 925 296 5795; Email:

9 zhongwang@lbl.gov

10 **ABSTRACT**

11 Grouping large genomic fragments assembled from shotgun metagenomic sequences to
12 deconvolute complex microbial communities, or metagenome binning, enables the study of
13 individual organisms and their interactions. Because of the complex nature of these communities,
14 existing metagenome binning methods often miss a large number of microbial species. In
15 addition, most of the tools are not scalable to large datasets. Here we introduce automated
16 software, called MetaBAT that integrates empirical probabilistic distances of genome abundance
17 and tetranucleotide frequency (TNF) for accurate metagenome binning. MetaBAT outperforms
18 alternative methods in accuracy and computational efficiency on both synthetic and real
19 metagenome datasets. It automatically forms hundreds of high quality genome bins on a very
20 large assembly consisting millions of contigs in a matter of hours on a single node. MetaBAT is
21 open source software available at <https://bitbucket.org/berkeleylab/metabat>.

22 **INTRODUCTION**

23 High throughput metagenome shotgun sequencing is a powerful tool to study microbial
24 communities directly taken from their environment, thereby avoiding the requirement for
25 cultivation or the biases that may arise from it. Assembling short metagenome shotgun reads into
26 larger genomic fragments (contigs) by short read assemblers (Pevzner & Tang 2001; Pevzner et
27 al. 2001) often fail to produce full-length genomes. Predicting draft genomes from assembled
28 metagenomic contigs by metagenome binning provides a substitute for full-length genomes
29 (Mande et al. 2012; Mavromatis et al. 2007). Despite their fragmented nature, these draft
30 genomes are often derived from individual species (or population genomes representing

31 consensus sequences of different strains, Imelfort et al. 2014), and they approximate full
32 genomes as they can contain a near full set of genes.

33 Two metagenome binning approaches have been developed (reviewed in (Mande et al.
34 2012)). The supervised binning approach uses known genomes as references and relies on either
35 sequence homology or sequence composition similarity for binning (Krause et al. 2008; Wu &
36 Eisen 2008). This approach does not work well on environmental samples where many microbes
37 do not have closely related species with known genomes. In contrast, the unsupervised approach
38 relies on either discriminative sequence composition (Teeling et al. 2004b; Yang et al. 2010) or
39 species (or genomic fragments) co-abundance (Cotillard et al. 2013; Le Chatelier et al. 2013;
40 Nielsen et al. 2014; Qin et al. 2012; Wu & Ye 2011) or both (Albertsen et al. 2013; Alneberg et
41 al. 2014; Imelfort et al. 2014; Sharon et al. 2013; Wrighton et al. 2012; Wu et al. 2014) for
42 binning. Recent studies have shown that species co-abundance feature can be very effective to
43 deconvolute complex communities if there are many samples available (Albertsen et al. 2013;
44 Alneberg et al. 2014; Cotillard et al. 2013; Imelfort et al. 2014; Karlsson et al. 2013; Le Chatelier
45 et al. 2013; Nielsen et al. 2014; Sharon et al. 2013). A few recent methods, particularly
46 CONCOCT (Alneberg et al. 2014) and GroopM (Imelfort et al. 2014), are also fully automated
47 binning procedures.

48 Many of the above tools do not scale well to large metagenomic datasets. In this study, we
49 developed MetaBAT (Metagenome Binning with Abundance and Tetra-nucleotide frequencies),
50 as an efficient, fully automated software tool that is capable of binning millions of contigs from
51 thousands of samples. By using a novel statistical framework to combine tetra-nucleotide
52 frequency (TNF) and contig abundance probabilities, we demonstrated that MetaBAT produces
53 high quality genome bins.

54

55 **MATERIALS AND METHODS**

56 **An overview of MetaBAT software and its probabilistic models**

57 As a pre-requisite for binning, the user must create BAM files (Li et al. 2009) by aligning the
58 reads of each sample separately to the assembled metagenome (Figure 1. steps from 1 to 3).
59 MetaBAT takes an assembly file (fasta format, required) and sorted bam files (one per sample,
60 optional) as inputs. For each pair of contigs in a metagenome assembly, MetaBAT calculates

61 their probabilistic distance based on tetranucleotide frequency (TNF) and abundance (i.e. mean
62 base coverage), then the two distances are integrated into one composite distance. All the
63 pairwise distances form a matrix, which then is supplied to a modified k-medoid clustering
64 algorithm to bin contigs iteratively and exhaustively into genome bins (Figure 1).

65 We use tetranucleotide frequency as sequence composition signatures as it has been
66 previously shown that different microbial genomes have distinct TNF biases (Mrazek 2009;
67 Pride et al. 2003; Saeed et al. 2012; Teeling et al. 2004a). To empirically derive a distance to
68 discriminate TNFs of different genomes, we calculated the likelihood of inter- and intra-species
69 Euclidean distance by using 1,414 unique, complete genome references from NCBI (Figure 2a).
70 This empirically derived distance is termed tetranucleotide frequency probability distance (TDP).

71 To evaluate the effect of contig sizes on inter-species distance, we obtained posterior
72 probability distributions of inter-species distance with several fixed sizes and observed better
73 inter-species separation as contig size increases (Figure 2b). As contigs in real metagenome
74 assemblies have various sizes, we then modeled TDP between contigs of different sizes by fitting
75 a logistic function to reflect the dynamic nature of non-linear relationship between Euclidean
76 TNF distance and TDP in different contig sizes. The results (Figure 2, c and d) suggested that the
77 values of two parameters of the model, b and c , are unstable if the size of either contig is very
78 small ($< 2\text{kb}$) and one should be cautious to allow smaller contigs to be binned.

79 Although contigs originating from the same genome are expected to have similar sequence
80 coverage, i.e. genome abundance, the coverage of contigs can vary significantly within a library
81 due to biases originated from the current sequencing technology (Benjamini & Speed 2012;
82 Harismendy et al. 2009; Nakamura et al. 2011; Ross et al. 2013). As illustrated in Figure 2e, the
83 observed coverage variance derived from data consisting of isolate genome sequencing projects
84 (total 99 from IMG Database (Markowitz et al. 2012), henceforth referred as the IMG dataset)
85 significantly deviate from the theoretical Poisson distribution, consistent with the notion that
86 both the variance and the mean should be modelled (Clark et al. 2013). For computational
87 convenience, we chose the normal distribution as an approximation since it fits the observation
88 much better (Figure 2e). To compute the abundance distance of two contigs in one sample, we
89 use the area not shared by their inferred normal distributions with given coverage mean and
90 variance (Figure 2f). A geometric mean of the distances for all samples is used for the final
91 abundance distance probability (ADP) of two contigs. In addition, we applied a progressive

92 weighting mechanism to adjust the relative strength of the information from abundance distance,
 93 meaning that we put more weight on abundance distance when it was calculated from many
 94 samples (see below).

95 We then integrate TDP and ADP of each contig pair as the following:

$$96 \quad P(\mu_1, \sigma_1^2, \mu_2, \sigma_2^2) = \begin{cases} \max(TDP, ADP), & \text{if } TDP > 0.05 \\ ADP \cdot w + TDP \cdot (1 - w), & \text{otherwise} \end{cases}$$

97 , where $w = \min [\log(n+1) / \log(m+1), \alpha]$. n , m , and α represent the number of samples, a large
 98 number (100 as the default), and the maximum weight of ADP (0.9 as the default), respectively.
 99 For instance, in the default setting, the weight would be about 0.5 when the number of samples is
 100 10 and TDP is less than 0.05. The resulting distance matrix is used for binning (see below).

101 **Tetranucleotide Frequency Probability Distance (TDP)**

102 To establish empirical probabilities of intra- and inter- species for tetranucleotide frequency
 103 distance, we downloaded 1,414 unique, completed bacterial genomes from the NCBI database
 104 and shredded them into fragments ranging from 2.5kb to 500kb. Next, we obtained 1 billion
 105 random contig pairs from within or between genomes. The empirical posterior probability that
 106 two contigs are from different genomes is given as the following:

$$107 \quad P(T|D) = \frac{P(T)P(D|T)}{P(T)P(D|T) + P(R)P(D|R)}$$

109 , where T or R represent cases where two contigs are from different (inter) or the same (intra)
 110 species, respectively. D is the Euclidean TNF distance between two contigs. The same
 111 uninformative priors of T and R were chosen. In reality, P(T) is expected to be much bigger than
 112 P(R), thus we set $P(T) = 10 * P(R)$ as the default implementation to adjust the possible under-
 113 sampling issue in inter species distance.

114 The TDP of contig pairs with different sizes is approximated using logistic regression:

$$115 \quad P(D_{ij}; b_{ij}, c_{ij}) = \frac{1}{(1 + e^{-(b_{ij} + c_{ij} * D_{ij})})}$$

116 , where D_{ij} represents a Euclidean TNF distance between contig i and j . b and c , the two
 117 parameters for the logistic regression, are estimated from the empirical data.

118

119 **Abundance Distance Probability (ADP)**

120 The probabilistic abundance distance was calculated as follows: Suppose two contigs have the
 121 mean coverage of μ_1 and μ_2 and the variances of σ_1^2 and σ_2^2 , then we defined the abundance
 122 distance as the non-shared area of two normal distribution of $N(\mu_1, \sigma_1^2)$ and $N(\mu_2, \sigma_2^2)$:

$$123 \quad P(\mu_1, \sigma_1^2, \mu_2, \sigma_2^2) = \frac{1}{2} \int |\phi_{\mu_1, \sigma_1^2} - \phi_{\mu_2, \sigma_2^2}|$$

124 , where ϕ represents a normal distribution having two parameters μ and σ^2 . Numerically this can
 125 be simplified using cumulative distribution functions as follows assuming σ_2^2 is greater than or
 126 equal to σ_1^2 :

$$127 \quad P(\mu_1, \sigma_1^2, \mu_2, \sigma_2^2) = \begin{cases} \Phi_{\mu_1, \sigma_1^2}(k_0) - \Phi_{\mu_2, \sigma_2^2}(k_0), & \text{if } \sigma_1^2 = \sigma_2^2 \\ \Phi_{\mu_1, \sigma_1^2}(k_2) - \Phi_{\mu_1, \sigma_1^2}(k_1) + \Phi_{\mu_2, \sigma_2^2}(k_1) - \Phi_{\mu_2, \sigma_2^2}(k_2), & \text{otherwise} \end{cases}$$

128 , where Φ represents a cumulative normal distribution, and

$$k_0 = \frac{\mu_1 + \mu_2}{2},$$

$$k_1^* = \frac{\sqrt{\sigma_1^2 \cdot \sigma_2^2 \cdot ((\mu_1 - \mu_2)^2 - 2 \cdot (\sigma_1^2 - \sigma_2^2) \cdot \log(\sigma_2/\sigma_1))} - \mu_1 \cdot \sigma_2^2 + \mu_2 \cdot \sigma_1^2}{\sigma_1^2 - \sigma_2^2},$$

$$k_2^* = \frac{\sqrt{\sigma_1^2 \cdot \sigma_2^2 \cdot ((\mu_1 - \mu_2)^2 - 2 \cdot (\sigma_1^2 - \sigma_2^2) \cdot \log(\sigma_2/\sigma_1))} + \mu_1 \cdot \sigma_2^2 - \mu_2 \cdot \sigma_1^2}{\sigma_1^2 - \sigma_2^2},$$

$$129 \quad k_1 = \min(k_1^*, k_2^*) \text{ and } k_2 = \max(k_1^*, k_2^*)$$

130 To combine multiple abundance probabilities across different samples, we calculated the
 131 geometric mean of probabilities:

$$132 \quad P_{ij} = \sqrt{\prod_n^{\Sigma_n \mathbb{1}\{\mu_{in} > c \text{ OR } \mu_{jn} > c\}} P_{ijn}(\mu_{in}, \sigma_{in}^2, \mu_{jn}, \sigma_{jn}^2) * \mathbb{1}\{\mu_{in} > c \text{ OR } \mu_{jn} > c\}}$$

133 , where P_{ijn} represents the probability calculated from two abundances μ_{in} and μ_{jn} , and c
 134 represents a cut-off for reasonable minimum abundance for a contig.

135 As metagenome assemblies contain many small contigs, whether or not to include them is a
 136 dilemma; including small contigs will likely improve the genome completeness, at a cost in
 137 genome quality because their larger abundance variations make it harder to bin them correctly.
 138 We tried to empirically determine a reasonable contig size cut-off by plotting the ratio of mean
 139 and variance from the IMG single genome dataset. Although most genome variances are much
 140 larger than their means, their ratio becomes stabilized after contig size increases to 2.5kb

141 (Supplementary Figure S1). Therefore we used 2.5kb or larger contigs for the initial binning.
142 Smaller contigs can be recruited after binning based on their correlation to the bins (Imelfort et al.
143 2014).

144 **Iterative Binning**

145 We modified the k-medoid clustering algorithm (Kaufman & Rousseeuw 1987) to eliminate the
146 need to input a value for k and to reduce search space for efficient binning. Specifically, the
147 binning algorithm works as the following:

- 148 1. Find a seed contig (e.g. having the greatest coverage), and set it as the initial medoid.
- 149 2. Recruit all other contigs within a cutoff distance (i.e. parameters p_1 and p_2) to the seed.
- 150 3. Find a new medoid out of all member contigs.
- 151 4. Repeat 2-3 until there are no further updates to the medoid. These contigs form a bin.
- 152 5. For the rest of the contigs, repeat 1-4 to form more bins until no contigs are left.
- 153 6. Keep large bins (e.g. >200kb), and dissolve all other bins into free contigs.
- 154 7. (optional) For dataset with at least 10 samples, recruit additional free contigs to each bin
155 based on their abundance correlation.

156

157 **RESULTS**

158 **Binning performance on “error-free” metagenomic assemblies**

159 A metagenomic dataset (Accession #: ERP000108) from the MetaHIT consortium (henceforth
160 referred to as the MetaHIT dataset) (Qin et al. 2010) was chosen to benchmark MetaBAT
161 because it contains a large number of samples and the community contains many species with
162 reference genomes. To derive a reference genome set, we selected 290 known genomes from
163 NCBI that are present in MetaHIT at >5X mean coverage (Supplementary Table S1). These
164 reference genomes were then shredded into contigs of random sizes (> 2.5kb) following an
165 exponential distribution modeled to mimic real metagenome assemblies. The abundance of each
166 contig in every sample was also obtained using real data. These “error-free” metagenome contigs,
167 their abundance information, along with their parental reference genomes (as “true answers”),
168 were used in the following analysis to benchmark binning performance. For a full description of
169 the experiment, refer to MetaBAT wiki page:

170 <https://bitbucket.org/berkeleylab/metabat/wiki/Home>

171 For comparison, we ran several alternative binning tools on the same dataset described above.
172 These software include Canopy (Nielsen et al. 2014), CONCOCT v.0.4.0 (Alneberg et al. 2014),
173 GroopM v.0.3.0 (Imelfort et al. 2014), and MaxBin v.1.4.1 (Wu et al. 2014). Among them,
174 CONCOCT, GroopM, and MaxBin are also fully automated binning tools. An optional manual
175 step in GroopM for improving the quality of bins was excluded. Since MaxBin does not consider
176 multiple samples, we combined multiple samples into one.

177 We used >90% precision (lack of contamination) and >30% recall (completeness) as the
178 minimum criteria for a bin to be considered “good” which basically means the bin should be
179 composed of one or more strains of a single species (for results of other thresholds, refer to
180 Supplementary Figure S3 and S4). Formulas for this calculation are described in the
181 Supplementary Material. Among all binning tools, MetaBAT binned the greatest number of
182 genomes at almost every recall threshold (Figure 3A). CONCOCT is the only tool that produces
183 more genome bins with over 80% or 90% completeness than MetaBAT, but MetaBAT produces
184 many more bins at 70% completeness threshold. Interestingly, we found these tools complement
185 each other in forming genome bins (Figure 3B, GroopM was omitted since it detects only a few
186 genomes). Among the unique 133 genome bins collectively formed by all tools, MetaBAT
187 binned the most number of genomes (111, 83.5%), with 23 bins (17.2%) that were not found by
188 any other tool.

189 MetaBAT runs very efficiently in computation, as the entire binning process only took 14
190 minutes and 4GB of RAM (Table 1). Multiple simulations produced almost identical
191 performance results and thus were not shown.

192 **Binning performance on real metagenomic assemblies**

193 We next tested the performance of MetaBAT on real metagenomic assemblies. Using the same
194 raw sequence data from the above MetaHIT dataset, we pooled the sequences from all samples
195 and then assembled them using Ray Meta assembler (Boisvert et al. 2012). Because real
196 metagenomic assemblies often contain many small contigs, we lowered the minimum contig size
197 requirement from 2.5kb to 1.5kb to include more contigs into our binning experiment. As a result,
198 118,025 contigs were used for binning. We then ran the above 5 binning tools with their default
199 settings on this dataset. In contrast to the previous simulation experiment, in this experiment we

200 do not have a reference genome for every genome bin; we instead used CheckM (Parks et al.
201 2014) to calculate the approximate recall (percent of expected single-copy-genes that are binned)
202 and precision (the absence of genes from different genomes) rates. A full description of this
203 experiment is available of MetaBAT wiki page:

204 https://bitbucket.org/berkeleylab/metabat/wiki/Benchmark_MetaHIT.

205 Similar to the previous “error-free” experiment, MetaBAT again identified the greatest
206 number of unique genome bins having >90% precision (Figure 4A). In this experiment with real
207 metagenomic contigs the superior completeness we saw in CONCOCT during the “error-free”
208 experiment was lost. Moreover, the number of genome bins formed by MetaBAT was
209 consistently greater than the others at every completeness threshold. Similarly different tools
210 produced complementary binning results as before (Figure 4B). MetaBAT’s contribution appears
211 to be more pronounced this time. It missed 17 bins formed by all other tools combined, but
212 recovered 31 bins that no other tools produced. MetaBAT alone recovered 90.2% (133/144) of
213 genome bins from all tools. These results suggest MetaBAT is very robust against a real
214 metagenome assembly. Consistent with the simulation experiment, MetaBAT is computationally
215 very efficient and requires only 4 minutes to complete this experiment (Table 2).

216 To test the performance of MetaBAT on large-scale metagenomic data sets, we used a
217 dataset containing 1,704 (with replicates) human gut microbiome samples (Accession #:
218 ERP002061) (Nielsen et al. 2014). The entire dataset was assembled using Ray Meta assembler
219 (Boisvert et al. 2012) and Megahit (Li et al. 2015) resulting in 1,058,952 contigs (>1kb) that
220 were then used for binning. MetaBAT took less than 2 hours to generate 1,634 genome bins
221 (>200kb) using a single node with 16 CPU cores (32 hyper-threads), and the peak memory
222 consumption was at 17G. In comparison, Canopy took 18 hours using 36G memory in the same
223 setting. The other binning tools--CONCOCT, GroopM, and MaxBin--all failed to generate any
224 genome bins for this data set likely due to their inability to scale. For accuracy evaluation, we
225 used CheckM (Parks et al. 2014) and identified 610 high quality bins (>90% precision and >50%
226 completeness) among the bins predicted by MetaBAT, which is 35% more than the published
227 CAG bins (Nielsen et al. 2014) and 11% more than Canopy bins using our assembly. For details
228 to use MetaBAT on this large dataset please refer to:

229 https://bitbucket.org/berkeleylab/metabat/wiki/Example_Large_Data.

230

231 **Post-binning processing to further improve quality**

232 Assembly from pooled samples in the above experiment raises the possibility that similar
233 genomes (e.g., different strains) present in different samples are assembled into chimeric contigs.
234 This level of contamination may not be tolerated in some applications. Based on the assumption
235 that a single sample will contain fewer strains of the same species than all pooled samples, we
236 implemented an optional post-binning process to reduce the strain-level contamination. Briefly,
237 we first selected a single sample with the most reads mapped to a specific bin, and then
238 assembled these reads into a new set of contigs. If the new contigs produces better CheckM
239 results, we subsequently replaced the old contigs in this bin with the new ones.

240 This post-processing step significantly improved both completeness and precision (lack of
241 contamination) for 61% (992 out of 1,634) of the genome bins. Overall there were 571 bins with
242 >95% precision and >50% completeness, compared with 375 without post-processing. This
243 improvement was more obvious when we increased the precision threshold to >99%, as the
244 number high quality bins increases from 46 to 186.

245 By incorporating additional sequencing data and other post-binning optimizations, Nielson et
246 al. generated 373 high quality draft genomes (“MGS genomes”) (Nielson et al. 2014). We
247 therefore used these MGS draft genomes as reference for additional quality assessment of the
248 MetaBAT genome bins after post-processing. As shown in Figure 5A, 31 MGS draft genomes
249 were not well represented by MetaBAT bins, but MetaBAT recovered 55 additional genome bins
250 not reported by MGS draft genomes. For those overlapping bins, most MetaBAT bins closely
251 approximate the MGS draft genomes in accuracy—94% precision and 82% completeness
252 (Figure 5B).

253

254 **DISCUSSION**

255 In conclusion, we developed an efficient and fully automated metagenome binning tool,
256 MetaBAT, and evaluated its capability to reconstruct genomes using both synthetic and real
257 world metagenome datasets. Applying MetaBAT to a large-scale complex human microbiome
258 data recovers hundreds of high quality genome bins including many missed by alternative tools.
259 An optional post-processing step improves the overall binning quality.

260 The algorithm implemented in MetaBAT is different from existing methods in several ways.
261 For example, MetaBAT uses both sequence composition and co-abundance as features, while

262 Canopy (Nielsen et al. 2014) only uses co-abundance. Among the methods that use sequence
263 composition feature, CONCOCT and GroopM both rely on principal component analysis (PCA)
264 to discriminate contigs belong to different genomes. In contrast, MetaBAT and Maxbin calculate
265 posterior probabilities by modeling intra- and inter-genome TNF distances. While Maxbin uses a
266 universal model to calculate TNF distance probability for all contig sizes, MetaBAT uses a
267 dynamic model weighing on different contig sizes.

268 One of the noticeable improvements in MetaBAT over other automated tools is its
269 computational efficiency. In addition to the low memory requirement and fast computing speed,
270 if one runs several rounds of binning to fine-tune parameters on a large dataset (by default
271 MetaBAT does little parameter optimization), MetaBAT can be even faster as it saves
272 intermediate calculations. For example, binning with ~1M contigs and ~1K samples for a second
273 time only takes a few minutes.

274 There are a couple of considerations to keep in mind before applying MetaBAT. One
275 important consideration is the minimum number of samples required for a reasonable binning
276 performance. Although MetaBAT can run with only one sample or even in TNF-only mode for
277 binning, as shown in Supplementary Figure S6, our general advice is that more samples achieve
278 better binning results. The greater the abundance variation among samples of a target species, the
279 more likely MetaBAT will produce a good genome bin for this species. A second consideration
280 is the quality of the metagenome assembly. We do not expect binning to work well with poor
281 metagenome assemblies, e.g. assemblies including many small contigs less than 1kb since the
282 distance metrics computed for small contigs will not be very reliable.

283

284 **ACKNOWLEDGEMENT**

285 The authors thank Drs. Matt Blow, Rex Malmstrom and Tanja Woyke for their stimulating
286 discussions and critical comments.

287 **REFERENCES**

288 Albertsen M, Hugenholtz P, Skarshewski A, Nielsen KL, Tyson GW, and Nielsen PH. 2013. Genome
289 sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple metagenomes.
290 *Nat Biotechnol* 31:533-538.

- 291 Alneberg J, Bjarnason BS, de Bruijn I, Schirmer M, Quick J, Ijaz UZ, Lahti L, Loman NJ, Andersson AF,
292 and Quince C. 2014. Binning metagenomic contigs by coverage and composition. *Nat Methods* 11:1144-
293 1146.
- 294 Benjamini Y, and Speed TP. 2012. Summarizing and correcting the GC content bias in high-throughput
295 sequencing. *Nucleic Acids Res* 40:e72.
- 296 Boisvert S, Raymond F, Godzaridis E, Laviolette F, and Corbeil J. 2012. Ray Meta: scalable de novo
297 metagenome assembly and profiling. *Genome Biol* 13:R122.
- 298 Clark SC, Egan R, Frazier PI, and Wang Z. 2013. ALE: a generic assembly likelihood evaluation
299 framework for assessing the accuracy of genome and metagenome assemblies. *Bioinformatics* 29:435-
300 443.
- 301 Cotillard A, Kennedy SP, Kong LC, Prifti E, Pons N, Le Chatelier E, Almeida M, Quinquis B, Levenez F,
302 Galleron N, Gougis S, Rizkalla S, Batto JM, Renault P, consortium ANRM, Dore J, Zucker JD, Clement K,
303 and Ehrlich SD. 2013. Dietary intervention impact on gut microbial gene richness. *Nature* 500:585-588.
- 304 Daniel R. 2005. The metagenomics of soil. *Nat Rev Microbiol* 3:470-478.
- 305 Harismendy O, Ng PC, Strausberg RL, Wang X, Stockwell TB, Beeson KY, Schork NJ, Murray SS, Topol
306 EJ, Levy S, and Frazer KA. 2009. Evaluation of next generation sequencing platforms for population
307 targeted sequencing studies. *Genome Biol* 10:R32.
- 308 Imelfort M, Parks D, Woodcroft BJ, Dennis P, Hugenholtz P, and Tyson GW. 2014. GroopM: an
309 automated tool for the recovery of population genomes from related metagenomes. *PeerJ* 2:e603.
- 310 Kaufman L, and Rousseeuw P. 1987. *Clustering by means of medoids*: North-Holland.
- 311 Krause L, Diaz NN, Goesmann A, Kelley S, Nattkemper TW, Rohwer F, Edwards RA, and Stoye J. 2008.
312 Phylogenetic classification of short environmental DNA fragments. *Nucleic Acids Res* 36:2230-2239.
- 313 Le Chatelier E, Nielsen T, Qin J, Prifti E, Hildebrand F, Falony G, Almeida M, Arumugam M, Batto JM,
314 Kennedy S, Leonard P, Li J, Burgdorf K, Grarup N, Jorgensen T, Brandslund I, Nielsen HB, Juncker AS,
315 Bertalan M, Levenez F, Pons N, Rasmussen S, Sunagawa S, Tap J, Tims S, Zoetendal EG, Brunak S,
316 Clement K, Dore J, Kleerebezem M, Kristiansen K, Renault P, Sicheritz-Ponten T, de Vos WM, Zucker JD,
317 Raes J, Hansen T, Meta HITc, Bork P, Wang J, Ehrlich SD, and Pedersen O. 2013. Richness of human
318 gut microbiome correlates with metabolic markers. *Nature* 500:541-546.
- 319 Li D, Liu CM, Luo R, Sadakane K, and Lam TW. 2015. MEGAHIT: an ultra-fast single-node solution for
320 large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*.
- 321 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, and
322 Genome Project Data Processing S. 2009. The Sequence Alignment/Map format and SAMtools.
323 *Bioinformatics* 25:2078-2079.
- 324 Mande SS, Mohammed MH, and Ghosh TS. 2012. Classification of metagenomic sequences: methods
325 and challenges. *Brief Bioinform* 13:669-681.
- 326 Markowitz VM, Chen IM, Palaniappan K, Chu K, Szeto E, Grechkin Y, Ratner A, Jacob B, Huang J,
327 Williams P, Huntemann M, Anderson I, Mavromatis K, Ivanova NN, and Kyrpides NC. 2012. IMG: the

328 Integrated Microbial Genomes database and comparative analysis system. *Nucleic Acids Res* 40:D115-
329 122.

330 Mavromatis K, Ivanova N, Barry K, Shapiro H, Goltsman E, McHardy AC, Rigoutsos I, Salamov A,
331 Korzeniewski F, Land M, Lapidus A, Grigoriev I, Richardson P, Hugenholtz P, and Kyrpides NC. 2007.
332 Use of simulated data sets to evaluate the fidelity of metagenomic processing methods. *Nat Methods*
333 4:495-500.

334 Mrazek J. 2009. Phylogenetic Signals in DNA Composition: Limitations and Prospects. *Molecular Biology*
335 *and Evolution* 26:1163-1169.

336 Nakamura K, Oshima T, Morimoto T, Ikeda S, Yoshikawa H, Shiwa Y, Ishikawa S, Linak MC, Hirai A,
337 Takahashi H, Altaf-Ul-Amin M, Ogasawara N, and Kanaya S. 2011. Sequence-specific error profile of
338 Illumina sequencers. *Nucleic Acids Res* 39:e90.

339 Nielsen HB, Almeida M, Juncker AS, Rasmussen S, Li J, Sunagawa S, Plichta DR, Gautier L, Pedersen
340 AG, Le Chatelier E, Pelletier E, Bonde I, Nielsen T, Manichanh C, Arumugam M, Batto JM, Quintanilha
341 Dos Santos MB, Blom N, Borruel N, Burgdorf KS, Boumezbear F, Casellas F, Dore J, Dworzynski P,
342 Guarner F, Hansen T, Hildebrand F, Kaas RS, Kennedy S, Kristiansen K, Kultima JR, Leonard P,
343 Levenez F, Lund O, Moumen B, Le Paslier D, Pons N, Pedersen O, Prifti E, Qin J, Raes J, Sorensen S,
344 Tap J, Tims S, Ussery DW, Yamada T, Meta HITC, Renault P, Sicheritz-Ponten T, Bork P, Wang J,
345 Brunak S, and Ehrlich SD. 2014. Identification and assembly of genomes and genetic elements in
346 complex metagenomic samples without using reference genomes. *Nat Biotechnol*.

347 Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, and Tyson GW. 2014. CheckM: assessing the
348 quality of microbial genomes recovered from isolates, single cells, and metagenomes. *PeerJ PrePrints*.

349 Pevzner PA, and Tang H. 2001. Fragment assembly with double-barreled data. *Bioinformatics* 17 Suppl
350 1:S225-233.

351 Pevzner PA, Tang H, and Waterman MS. 2001. An Eulerian path approach to DNA fragment assembly.
352 *Proc Natl Acad Sci USA* 98:9748-9753.

353 Pride DT, Meinersmann RJ, Wassenaar TM, and Blaser MJ. 2003. Evolutionary implications of microbial
354 genome tetranucleotide frequency biases. *Genome Res* 13:145-158.

355 Qin J, Li Y, Cai Z, Li S, Zhu J, Zhang F, Liang S, Zhang W, Guan Y, Shen D, Peng Y, Zhang D, Jie Z, Wu
356 W, Qin Y, Xue W, Li J, Han L, Lu D, Wu P, Dai Y, Sun X, Li Z, Tang A, Zhong S, Li X, Chen W, Xu R,
357 Wang M, Feng Q, Gong M, Yu J, Zhang Y, Zhang M, Hansen T, Sanchez G, Raes J, Falony G, Okuda S,
358 Almeida M, LeChatelier E, Renault P, Pons N, Batto JM, Zhang Z, Chen H, Yang R, Zheng W, Li S, Yang
359 H, Wang J, Ehrlich SD, Nielsen R, Pedersen O, Kristiansen K, and Wang J. 2012. A metagenome-wide
360 association study of gut microbiota in type 2 diabetes. *Nature* 490:55-60.

361 Qin JJ, Li RQ, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons N, Levenez F, Yamada
362 T, Mende DR, Li JH, Xu JM, Li SC, Li DF, Cao JJ, Wang B, Liang HQ, Zheng HS, Xie YL, Tap J, Lepage
363 P, Bertalan M, Batto JM, Hansen T, Le Paslier D, Linneberg A, Nielsen HB, Pelletier E, Renault P,
364 Sicheritz-Ponten T, Turner K, Zhu HM, Yu C, Li ST, Jian M, Zhou Y, Li YR, Zhang XQ, Li SG, Qin N,

365 Yang HM, Wang J, Brunak S, Dore J, Guarner F, Kristiansen K, Pedersen O, Parkhill J, Weissenbach J,
366 Bork P, Ehrlich SD, Wang J, and Consortium M. 2010. A human gut microbial gene catalogue established
367 by metagenomic sequencing. *Nature* 464:59-U70.

368 Ross MG, Russ C, Costello M, Hollinger A, Lennon NJ, Hegarty R, Nusbaum C, and Jaffe DB. 2013.
369 Characterizing and measuring bias in sequence data. *Genome Biol* 14:R51.

370 Saeed I, Tang SL, and Halgamuge SK. 2012. Unsupervised discovery of microbial population structure
371 within metagenomes using nucleotide base composition. *Nucleic Acids Res* 40:e34.

372 Sharon I, Morowitz MJ, Thomas BC, Costello EK, Relman DA, and Banfield JF. 2013. Time series
373 community genomics analysis reveals rapid shifts in bacterial species, strains, and phage during infant
374 gut colonization. *Genome Res* 23:111-120.

375 Teeling H, Meyerdierks A, Bauer M, Amann R, and Glockner FO. 2004a. Application of tetranucleotide
376 frequencies for the assignment of genomic fragments. *Environ Microbiol* 6:938-947.

377 Teeling H, Waldmann J, Lombardot T, Bauer M, and Glockner FO. 2004b. TETRA: a web-service and a
378 stand-alone program for the analysis and comparison of tetranucleotide usage patterns in DNA
379 sequences. *BMC Bioinformatics* 5:163.

380 Wrighton KC, Thomas BC, Sharon I, Miller CS, Castelle CJ, VerBerkmoes NC, Wilkins MJ, Hettich RL,
381 Lipton MS, Williams KH, Long PE, and Banfield JF. 2012. Fermentation, hydrogen, and sulfur metabolism
382 in multiple uncultivated bacterial phyla. *Science* 337:1661-1665.

383 Wu M, and Eisen JA. 2008. A simple, fast, and accurate method of phylogenomic inference. *Genome Biol*
384 9:R151.

385 Wu Y-W, Tang Y-H, Tringe SG, Simmons BA, and Singer SW. 2014. MaxBin: an automated binning
386 method to recover individual genomes from metagenomes using an expectation-maximization algorithm.
387 *Microbiome* 2:26.

388 Wu YW, and Ye Y. 2011. A novel abundance-based algorithm for binning metagenomic sequences using
389 I-tuples. *J Comp Biol* 18:523-534.

390 Yang B, Peng Y, Leung HC, Yiu SM, Chen JC, and Chin FY. 2010. Unsupervised binning of
391 environmental genomic fragments based on an error robust selection of I-mers. *BMC Bioinformatics* 11
392 Suppl 2:S5.

393

394 **FIGURES LEGENDS**

395 Figure 1. Overview of the MetaBAT pipeline. Three preprocessing steps before MetaBAT are
396 applied: 1) A typical metagenome experiment may contain many spatial or time-series samples,
397 each consisting of many different genomes (different color circles). 2) Each sample is sequenced
398 by next-generation sequencing technology to form a sequencing library with many short reads. 3)
399 The libraries are combined for de novo assembly. After assembly, the reads from each sample
400 are aligned to form separate BAM files. MetaBAT then automatically performs the remaining

401 steps: 4) For each contig pair, a tetranucleotide frequency distance probability (TDP) is
 402 calculated. 5) For each contig pair, an abundance distance probability (ADP) across all the
 403 samples is calculated. 6) The TDP and ADP of each contig pair are then combined, and the
 404 resulting distance for all pairs form a distance matrix. 7) Each bin will be formed iteratively and
 405 exhaustively from the distance matrix.

406 Figure 2. Probabilistic modeling of TNF and Abundance distances. **a-d)** TNF distance modeling.
 407 **a)** Empirical probabilities of intra- (solid gray line) or inter- (dotted gray line) species Euclidean
 408 TNF distance are estimated from sequenced genomes. The posterior probability of two contigs
 409 originated from different genomes given a TNF distance is shown as a red solid line. All
 410 probabilities are calculated using a fixed contig size of 10kb. **b)** Different posterior inter-species
 411 probabilities for two equal-size contigs under various contig sizes. **c, d)** The estimation of
 412 parameters for a logistic curve with two contigs of different sizes. x and y axis represent the
 413 lengths of short and long contig, respectively, and z axis represents the estimates of each
 414 parameter *b* or *c* in a logistic curve, $TDP = 1/(1+\exp(-(b+c*TNF)))$, where TNF and TDP
 415 represents the Euclidean TNF distance and probabilistic TNF distance, respectively. **e-f)**
 416 Abundance distance modeling. **e)** The relationship between mean and variance of base depths
 417 (coverage) which were shown in x and y axis, respectively. Each dot represents this relationship
 418 in each genome, which calculated by median of mean and variance of the coverage. Theoretical
 419 Poisson model was shown as blue line and normal model was shown as red line. **f)** Probabilistic
 420 abundance distance between two contigs. The shaded area represents the abundance distance
 421 between two contigs in a given library.

422 Figure 3. Binning performance on synthetic metagenomic assemblies. A) The number of
 423 genomes (X-axis) identified by each binning method (Y-axis) in different recall (completeness)
 424 threshold and the same >90% precision threshold. B) A Venn diagram of identified genomes by
 425 top 4 binning methods.

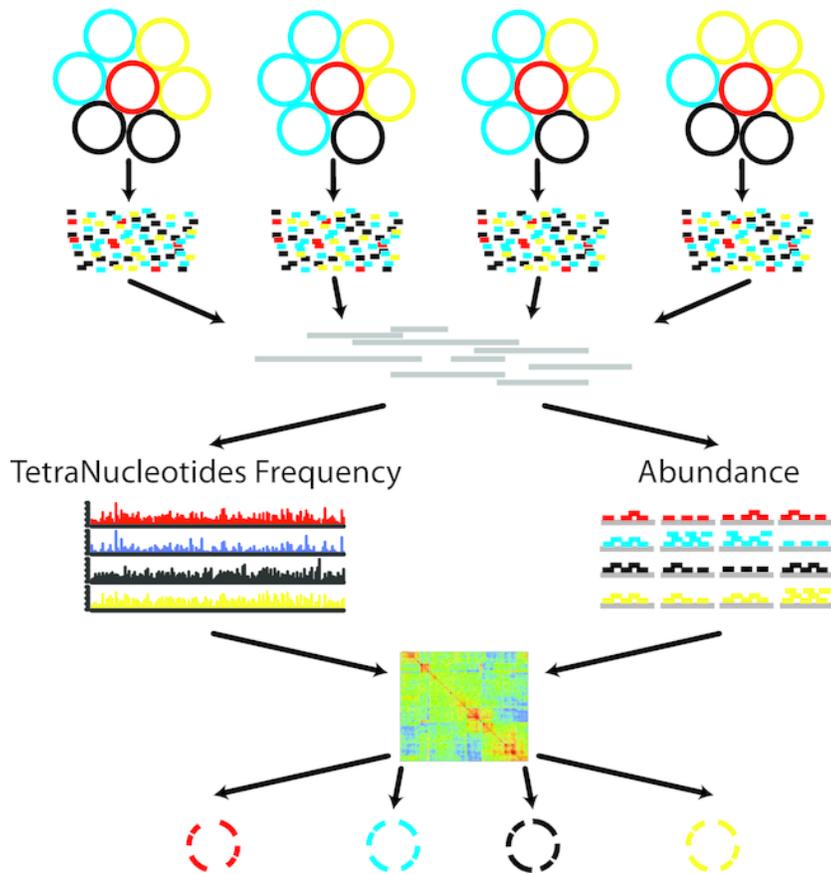
426 Figure 4. Binning performance on real metagenomic assemblies. A) The number of genomes (X-
 427 axis) identified by each binning method (Y-axis) in different recall (completeness) threshold and
 428 the same >90% precision threshold. B) A Venn diagram of identified genomes by top 4 binning
 429 methods.

430 Figure 5. A comparison between MetaBAT bins after post-processing and MGS draft genomes
431 from Nielsen et al. A) A Venn diagram of identified genome bins by MetaBAT and MGS draft
432 genomes. B) A scatterplot of completeness and precision for MetaBAT genome bins using MGS
433 draft genomes as reference. X-axis represents shared proportion of bases in terms of MetaBAT
434 bins (i.e. precision), and y-axis represents shared proportion of bases in terms of MGS genomes
435 (i.e. completeness). Each circle represents a unique MetaBAT bin having a corresponding MGS
436 genome (342 bins in total), and the size of the circle corresponds to the bin size.

1

Figure 1. Overview of the MetaBAT pipeline.

Three preprocessing steps before MetaBAT is applied: 1) A typical metagenome experiment may contain many spatial or time-series samples, each consisting of many different genomes (different color circles). 2) Each sample is sequenced by next-generation sequencing technology to form a sequencing library with many short reads. 3) The libraries may be combined before de novo assembly. After assembly, the reads from each sample must be aligned in separate BAM files. MetaBAT then automatically performs the remaining steps: 4) For each contig pair, a tetranucleotide frequency distance probability (TDP) is calculated from a distribution modelled from 1,414 reference genomes. 5) For each contig pair, an abundance distance probability (ADP) across all the samples is calculated. 6) The TDP and ADP of each contig pair are then combined, and the resulting distance for all pairs form a distance matrix. 7) Each bin will be formed iteratively and exhaustively from the distance matrix.



Preprocessing

- 1 Samples from multiple sites or times
- 2 Metagenome libraries
- 3 Initial de-novo assembly using the combined library

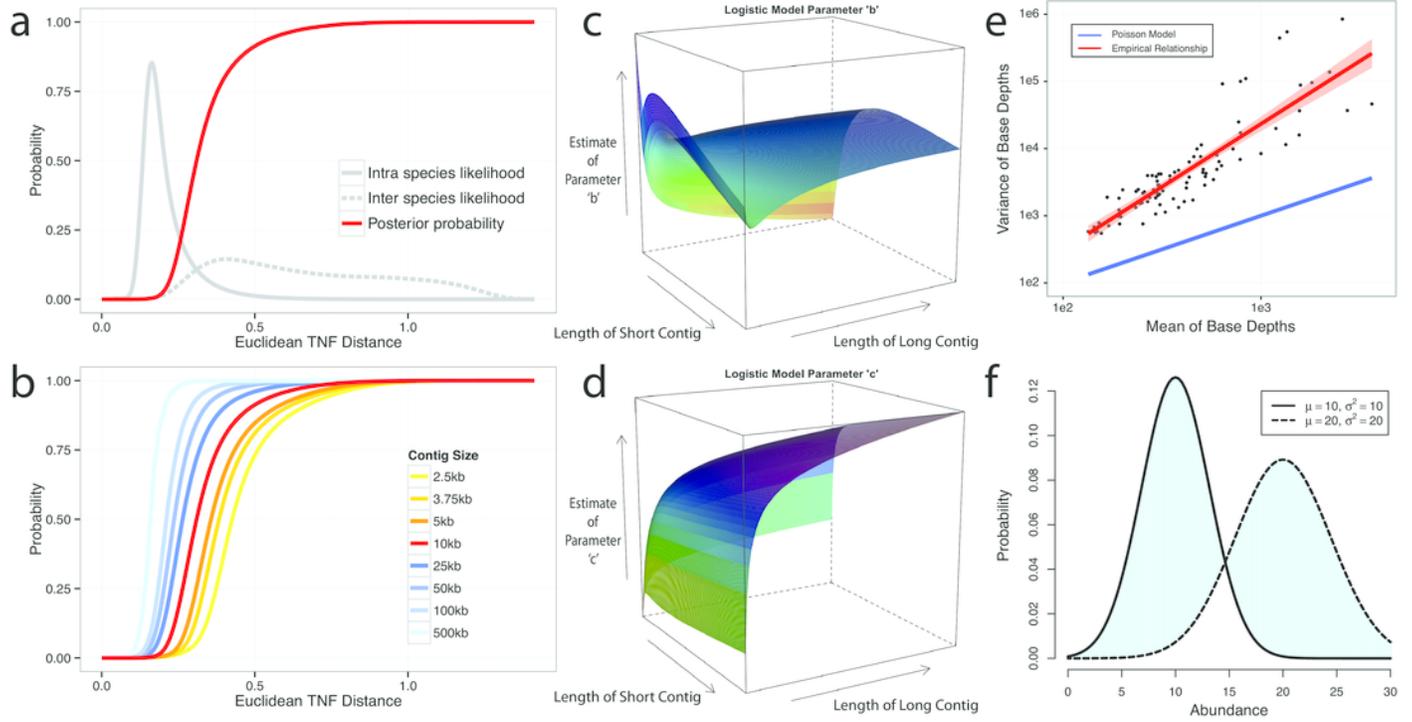
MetaBAT

- 4 Calculate TNF for each contig
- 5 Calculate Abundance per library for each contig
- 6 Calculate the pairwise distance matrix using pre-trained probabilistic models
- 7 Forming genome bins iteratively

2

Figure 2. Probabilistic modeling of TNF and Abundance distances.

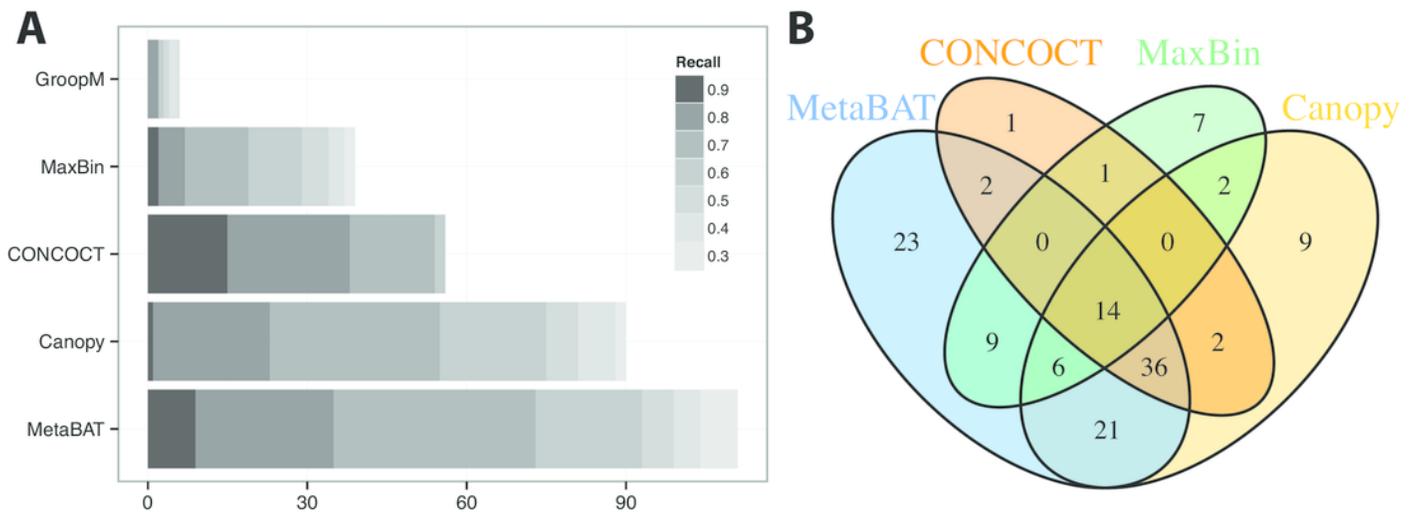
a-d) TNF distance modeling. **a)** Empirical probabilities of intra- (solid gray line) or inter- (dotted gray line) species Euclidean TNF distance are estimated from sequenced genomes. The posterior probability of two contigs originated from different genomes given a TNF distance is shown as a red solid line. All probabilities are calculated using a fixed contig size of 10kb. **b)** Different posterior inter-species probabilities for two equal-size contigs under various contig sizes. **c, d)** The estimation of parameters for a logistic curve with two contigs of different sizes. x and y axis represent the lengths of short and long contig, respectively, and z axis represents the estimates of each parameter b or c in a logistic curve, $TDP = 1/(1+\exp(-(b+c*TNF)))$, where TNF and TDP represents the Euclidean TNF distance and probabilistic TNF distance, respectively. **e-f)** Abundance distance modeling. **e)** The relationship between mean and variance of base depths (coverage) which were shown in x and y axis, respectively. Each dot represents this relationship in each genome, which calculated by median of mean and variance of the coverage. Theoretical Poisson model was shown as blue line and normal model was shown as red line. **f)** Probabilistic abundance distance between two contigs. The shaded area represents the abundance distance between two contigs in a given library.



3

Figure 3. Binning performance on synthetic metagenomic assemblies.

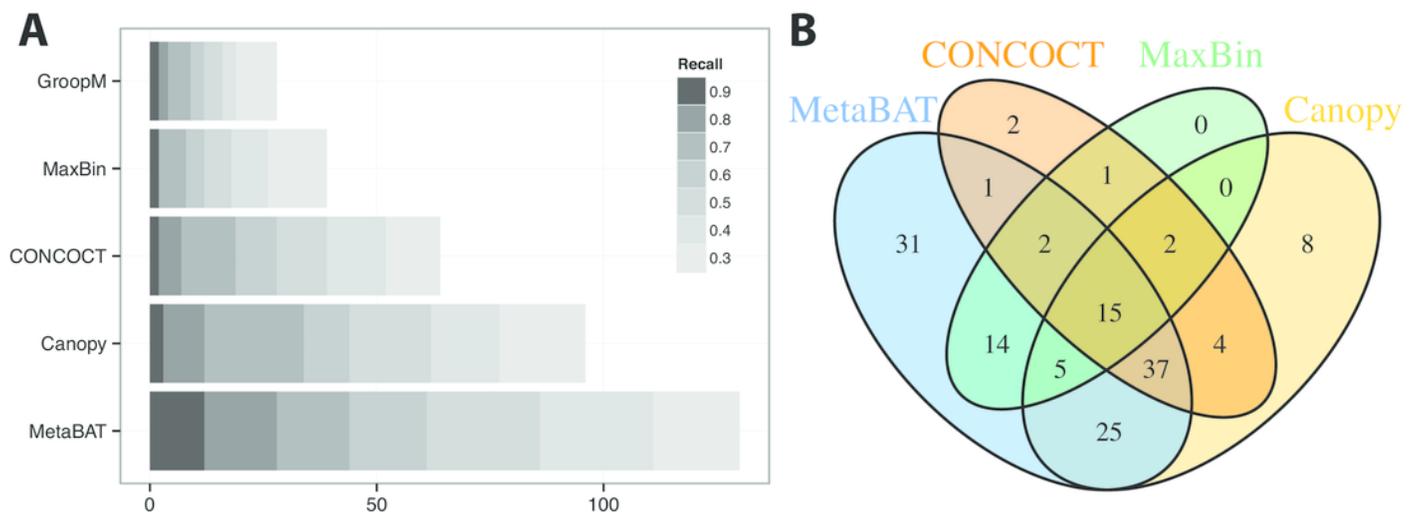
A) The number of genomes (X-axis) identified by each binning method (Y-axis) in different recall (completeness) threshold and >90% precision, which calculates the lack of contamination. B) Venn diagram of identified genomes by top 4 binning methods.



4

Figure 4. Binning performance on real metagenomic assemblies.

A) The number of genomes (X-axis) identified by each binning method (Y-axis) in different recall (completeness) threshold and >90% precision, which calculates the lack of contamination. B) Venn diagram of identified genomes by top 4 binning methods.



5

Figure 5. Comparison between MetaBAT bins after post-processing and MGS draft genomes from Nielsen et al.

A) Venn diagram of identified genome bins by MetaBAT having >90% precision and >30% completeness calculated by CheckM and one-to-one corresponding genomes in MGS draft genomes. B) Scatterplot of completeness and precision for MetaBAT genome bins when considered MGS draft genomes as the gold standard. X-axis represents shared proportion of bases in terms of MetaBAT bins (i.e. precision), and y-axis represents shared proportion of bases in terms of MGS genomes (i.e. completeness). Each circle represents a unique MetaBAT bins having uniquely corresponding MGS genomes (342 bins in total), and the size of it corresponds to bin size.

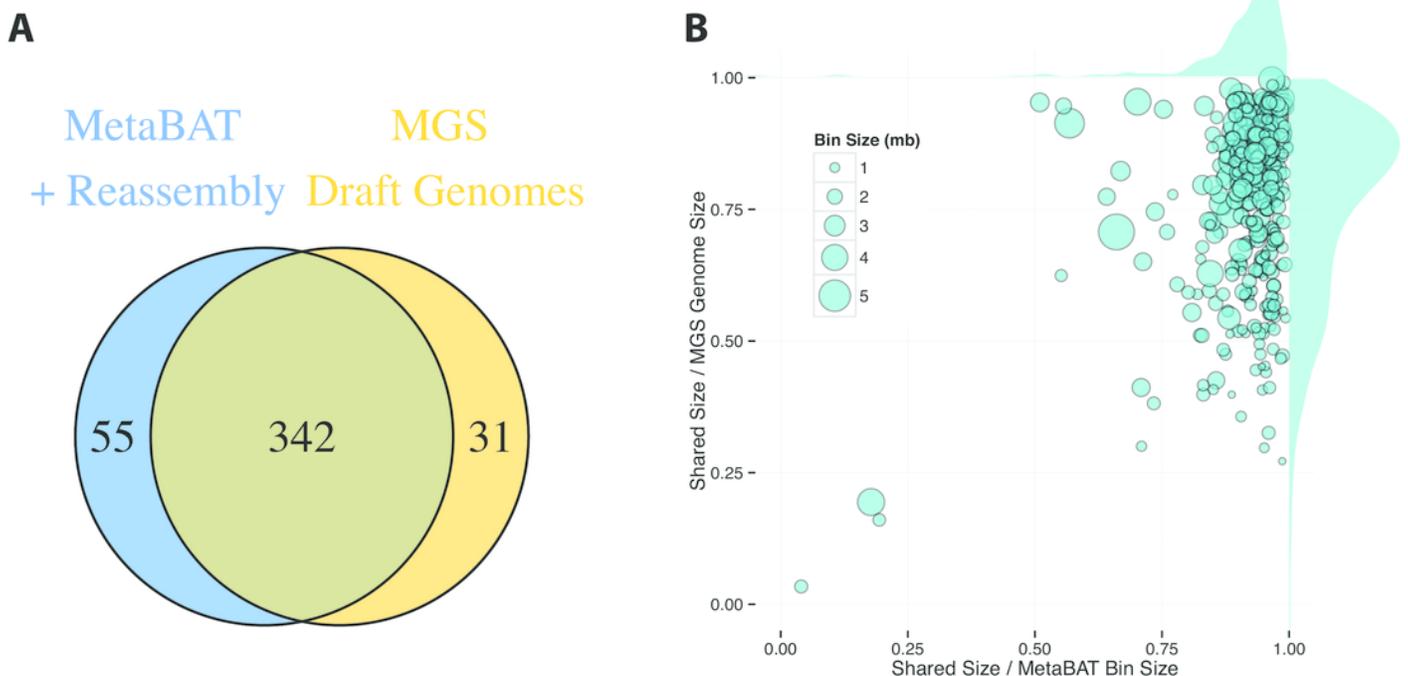


Table 1 (on next page)

Table 1

A summary of the binning performance on synthetic metagenomic assembly.

2 **Table 1. A summary of the binning performance on synthetic metagenomic assembly.**

	MetaBAT	Canopy	CONCOCT	MaxBin	GroopM**
Number of Bins Identified (>200kb)	340	230	235	260	445
Number of Genomes Detected (Precision > .9 & Recall > .3)	111	90	56	39	6
Wall Time (16 cores; 32 hyper-threads)	00:13:55	00:21:01*	104:58:01	20:51:19	45:29:39
Peak Memory Usage (for binning step)	3.9G	1.82G*	9.55G	7.7G	38G

*Canopy only use abundance table as input, so it should have taken more time and memory to read and write sequence data like the others

**Manual steps were not used

3

Table 2 (on next page)

Table 2

A summary of the binning performance on real metagenomic assembly.

2 **Table 2. A summary of the binning performance on real metagenomic assembly.**

	MetaBAT	Canopy	CONCOCT	MaxBin	GroopM**
Number of Bins Identified (>200kb)	234	223	260	168	335
Number of Genomes Detected (Precision > .9 & Recall > .3)	130	96	64	39	28
Wall Time (16 cores; 32 hyper-threads)	00:03:36	00:02:31*	82:19:53	06:49:39	12:19:12
Peak Memory Usage (for binning step)	3.0G	1.6G*	7G	5.8G	6.3G

*Canopy only use abundance table as input, so it should have taken more time and memory to read and write sequence data like the others

**Manual steps were not used

3