

# MetaTrass: High-quality metagenome assembly on the human

## gut microbiome by co-barcoding sequencing reads

Yanwei Qi<sup>1,2,3#</sup>, Shengqiang Gu<sup>1,4#</sup>, Yue Zhang<sup>1</sup>, Lidong Guo<sup>1,4</sup>, Mengyang Xu<sup>1,2,3,5</sup>,  
Xiaofang Cheng<sup>5,6</sup>, Ou Wang<sup>5,6</sup>, Jianwei Chen<sup>1</sup>, Xiaodong Fang<sup>5,7</sup>, Xin Liu<sup>1,2,3</sup>, Li  
Deng<sup>1,2,3\*</sup>, Guangyi Fan<sup>1,2,3,5\*</sup>

<sup>1</sup> BGI-Qingdao, BGI-Shenzhen, Qingdao 266555, China

<sup>2</sup> State Key Laboratory of Agricultural Genomics, BGI-Shenzhen, Shenzhen 518083,  
China

<sup>3</sup> China National GeneBank, BGI-Shenzhen, Shenzhen 518120, China

<sup>4</sup> College of Life Sciences, University of Chinese Academy of Sciences, Beijing 100049,  
China

<sup>5</sup> BGI-Shenzhen, Shenzhen 518083, China

<sup>6</sup> MGI, BGI-Shenzhen, Shenzhen 518083, China

<sup>7</sup> BGI Genomics, BGI-Shenzhen, Shenzhen 518083, China

<sup>#</sup>These authors contributed equally to this work.

<sup>\*</sup>Corresponding authors (Guangyi Fan, email: fanguangyi@genomics.cn; Li Deng,  
email: dengli1@genomics.cn.)

**Running title:** Qi Y et al / Metagenome assembling of the human gut microbiome

## 30 Abstract

31 With the development of sequencing technologies and computational analysis in  
 32 metagenomics, the genetic diversity of non-conserved regions has been receiving  
 33 intensive attention to unravel the human gut microbial community. However, it  
 34 remains a challenge to obtain enough microbial draft genomes at a high resolution  
 35 from a single sample. In this work, we presented MetaTrass with a strategy of binning  
 36 first and assembling later to assemble high-quality draft genomes based on  
 37 metagenomics co-barcoding reads and the public reference genomes. We applied the  
 38 tool to the single tube long fragment reads datasets for four human faecal samples,  
 39 and generated more high-quality draft genomes with longer contiguity and higher  
 40 resolution than the common combination strategies of genome assembling and  
 41 binning. A total of 178 high-quality genomes was successfully assembled by  
 42 MetaTrass, but the maximum of 58 was generated by the optimal common  
 43 combination strategy in our tests. These high-quality genomes paved the way for  
 44 genetic diversity and lineage analysis among different samples. With the high  
 45 capability of assembling high-quality genomes of metagenomics datasets, MetaTrass  
 46 will facilitate the study of spatial characters and dynamics of complex microbial  
 47 communities at high resolution. The open-source code of MetaTrass is available at  
 48 <https://github.com/BGI-Qingdao/MetaTrass>.

49  
 50 **KEYWORDS:** Metagenome assembly, Synthetic long reads, Taxonomic binning,  
 51 Microbiome composition, Phylogenetic trees

## 54 Introduction

55 Through sequencing and analyzing the DNA of microbial communities directly from  
 56 the environment, metagenomics has showed important roles in advancing the study of  
 57 uncultured microbiomes [1, 2]. Comprehensive databases of metagenome-assembled  
 58 genomes, especially for the human gut microbiome, are massively expanded to

59 completely understand the genomic taxonomic structure of different microbiome  
60 communities according to genetic similarity [3, 4]. The progresses in metagenomics  
61 have shed new light on the study of spatial distribution and dynamics of complex  
62 microbial communities from the human gut [5, 6].

63 Based on the function mining of high-quality strain-resolved genomes, it is  
64 realized that genotypic differences among strains are strongly correlated with their  
65 phenotypic differences [7, 8]. The importance of intra-species non-homologous genes  
66 have been intensively studied in the field of pathogenicity, and many new species  
67 with both pathogenic and commensal strains have been found [9-11]. Indeed, the  
68 percentage of conserved intra-species homologous genes shared between strains is as  
69 low as 40% [12], and the large part of non-conservation region is thought of as the  
70 genetic origin of phenotypic diversity. Thus, complete draft genomes from a  
71 microbiome sample at the species level will enable a more comprehensive study of  
72 intra-species genome diversity, but it is still a challenge to generate sufficient high-  
73 quality genomes from metagenomic datasets.

74 Most of current approaches to analyze the microbiome communities are based on  
75 high-throughput and low-cost next-generation sequencing (NGS) reads. Many highly  
76 modularized computational tools have been developed such as genome assemblers,  
77 genome binners, taxonomic binners and taxonomic profilers [13-15]. The  
78 combinations of assembling first and binning later have been commonly used to  
79 generate metagenome-assembled genomes. In these strategies, a mass of short reads  
80 from a microbial community are firstly assembled to generate longer sequences by  
81 metagenomics assemblers with the consideration of uneven coverage depths of  
82 different microbial species [16-18]. Then the assembled sequences are grouped into  
83 individual genomes by genome binners based on similar *K-mer* composition and read  
84 coverage [19-21]. As a result, draft genomes with non-conserved genes are retrieved  
85 from various microbial communities. However, it is impossible to solve the  
86 assembling problem of the long inter-species repeats by the short NGS reads, so the  
87 contiguity of the draft genomes assembled by NGS reads is still not enough long for  
88 studying the long structural variations in metagenomics.

89 Various sequencing technologies with long-range information accompanied by  
 90 specialized computational tools are promised to overcome the problem of long repeats.  
 91 Third-generation single-molecule real-time sequencing (TGS) technologies developed  
 92 by Pacific Biosciences and Oxford Nanopore Technology (ONT) can produce  
 93 contiguous reads with lengths up to hundreds of kb, and show great potential to  
 94 generate complete genomes from both cultured and uncultured microbial communities  
 95 [22-24]. With using the chromatin-level contact probability information generated by  
 96 high-throughput chromosome conformation capture (Hi-C) technology, more high-  
 97 quality genome bins with improved contiguity can be retrieved [25]. Additionally, the  
 98 co-abundance of species in multiple samples with the common *K-mer* composition are  
 99 also used to improve the capability to retrieve high-quality genome bins for NGS  
 100 datasets [26]. However, there are limitations for these approaches. The high  
 101 sequencing error rate in TGS long reads hampers the distinction between true  
 102 variations and sequencing errors. An effective contact map with Hi-C library can only  
 103 be established for a draft genome with preferable contiguity. Constructing co-  
 104 abundance in multiple samples ignores the genome characteristics of a single sample  
 105 and increase the sequencing cost.

106 The co-barcoding sequencing library [27-31], an improved short-read sequencing  
 107 with long-range genomic information, can provide an alternative way to improve  
 108 metagenomics analyzing. In a co-barcoding library construction, long fragments  
 109 sheared from DNA samples are firstly distributed into different isolated partitions, and  
 110 then short-read fragments from the long fragment in the same partition are labeled  
 111 with a unique barcode sequence, finally the co-barcoded fragments are sequenced by  
 112 standard short-read sequencing platforms. For different co-barcoding libraries such as  
 113 BGI's single tube long fragment reads (stLFR) library [30], 10X Genomics' linked-  
 114 reads library [32] and Illumina's contiguity preserving transposase sequencing library  
 115 [28], different technical metrics in the total barcode number and the short-read  
 116 coverage of the long fragment have a great impact on their powers in the downstream  
 117 analysis [33-36]. The co-barcoding correlation on the draft sequences or the  
 118 assembled graph have been successfully applied to improve the contiguity of

119 assembled genomes for both large eukaryotic genomes [37-39] and metagenomes [29,  
120 40, 41]. All these methods are still the common combination strategy in principle,  
121 leaving the inherent problem of long repeats among species with uneven abundance  
122 unsolved in efficiently constructing high-quality draft genomes for complex microbial  
123 communities.

124 In this work, we introduced a pipeline of **Meta**genomics **Taxonomic Read**  
125 **Assembly of Single Species** (MetaTrass) based on co-barcoding sequencing data and  
126 references. Different from the common strategies, MetaTrass was a strategy of  
127 binning first and assembling later. The co-barcoding information was used not only to  
128 improve the assemblies by implementing co-barcoding assemblies, but also to  
129 simplify the dataset before assembling using microbial references with the help of  
130 taxonomic binning. We apply MetaTrass to stLFR datasets of a mock metagenome  
131 community and four real gut microbiome communities to evaluate its capability of  
132 producing high-quality draft genomes with high contiguity and high taxonomic  
133 resolution. The results were benchmarked by comparing to the common combinations  
134 of several mainstream tools. Meanwhile, the microbiome composition and genetic  
135 diversity in the four human gut samples were quantitatively analyzed with using the  
136 high-quality draft genomes assembled by MetaTrass. We expected that the high-  
137 quality draft genomes with taxonomic information at the species level assembled by  
138 our tools would be convenient to make more extensively use to investigate various  
139 microbial communities.

140

## 141 **Materials and methods**

### 142 **Datasets**

143 A mock microbial and four gut microbial communities were analyzed to evaluate the  
144 efficiency of MetaTrass. The mock microbial community (ZymoBIOMICS™  
145 Microbial Community DNA Standard) consists of 8 isolated bacteria with the  
146 abundance of about 12% and 2 fungi with the abundance of about 2%. The four gut  
147 microbial DNA samples include faeces of three healthy volunteers and one patient

with inflammatory bowel disease. The stLFR libraries were constructed according to the standard protocol [30]. The DNA samples were firstly sheared into long fragments, and then the long fragments were captured into a magnetic microbead with a unique barcode sequence. Finally, each long fragment was broken and hybridized with a unique barcode by the Tn5 transposase on the surface of the microbead. The stLFR libraries of the mock and the patient sample were sequenced on BGISEQ500 platform, and those of healthy samples were sequenced on MGISEQ2000 platform. The read length in the read pair was 100 bp for all datasets. The mock and three healthy sample libraries were individually allocated to a half lane, and a total of about 50 Gb raw reads were generated. The patient library was allocated to a full lane, generating about 100 Gb raw reads. Barcode sequences were extracted from the end of read2 and then replaced by numerical symbols in the read names in the FASTQ file with an in-house script. SOAPfilter\_v2.2 with parameters (-y -F CTGTCTCTTATACATCTTAGGAAGACAAGCACTGACGACATGA -R TCTGCTGAGTCGAGAACGTCTCTGTGAGCCAAGGAGTTGCTCTGG -p -M 2 -f -I -Q 10) was used to clean out low-quality raw reads with adaptors, excessive confused bases, and high duplications. Finally, 55.65 Gb clean data were retained for the mock microbiome, 34.48 Gb for the first healthy sample (H\_Gut\_Meta01), 35.33 Gb for the second (H\_Gut\_Meta02), 37.88 Gb for the third (H\_Gut\_Meta03), and 97.20 Gb for the patient sample (P\_Gut\_Meta01).

168

## 169 Taxonomic binning

We adopted Kraken2 (version 2.0.9-beta) [42] to classify stLFR reads into different species. Firstly, a customized reference databases were constructed according to the microbial community. Specially, references attached to the ZYMO product were used for the mock sample. The Kraken2 database of the Unified Human Gastrointestinal Genomes (UHGG) collection [3] was used to study the gut samples, and which included 4542 representative genomes at the species level. Then, the corresponding stLFR reads were classified with default parameters.

177

## 178 **Co-barcoding reads refining**

179 Since a taxonomic tree of references was established to reduce the number of multiple  
 180 hits of a *K-mer* from inter-species homologous sequences in Kraken2, the reads from  
 181 these regions were classified into the lowest common ancient (LCA) rank higher than  
 182 its corresponding species. Several works tried to reallocate these reads to species by  
 183 statistical inferences using the coverage depth or co-barcoding information of intra-  
 184 species homologous region of a species [43, 44]. In MetaTrass pipeline, the co-  
 185 barcoding correlation between reads classified into a species and those classified into  
 186 high LCA ranks was used to reduce the false negative of reads classified into high  
 187 LCA rank. Reads classified into a species level is defined as the taxonomic reads of  
 188 the species. In this step, we collected and refined reads for each barcode according to  
 189 the number of reads in the taxonomic reads (*Num\_T*) and the ratio of these reads to  
 190 the total reads (*Ratio\_T*). Barcodes appearing in the taxonomic reads were firstly  
 191 extracted as candidates. Then, we ranked candidates first in order of *Num\_T* from  
 192 large to small, and then *Ratio\_T* for those with the same *Num\_T*. Finally, reads with a  
 193 barcode of *Ratio\_T* larger than a threshold were chose based on the barcode rank.  
 194 Since sufficient read coverage is required for assembling a complete genome, only the  
 195 read sets of one species with abundance higher than 10× were refined by co-barcoding  
 196 information. The abundance of each species was roughly calculated according to the  
 197 coverage depth of the taxonomic reads on the reference. Meanwhile, we set a data size  
 198 threshold of the refined reads to reduce the computational consumption for species  
 199 with extremely high abundance (e.g., 300×). Paired-end reads were extracted by Seqtk  
 200 (version 1.3-r114-dirty) according to the barcode-related read names from the FASTQ  
 201 file of clean reads. Note that there were still some false positive reads, although  
 202 *Ratio\_T* was set to reduce them caused by the collision of long fragments from  
 203 different species in the same microbead. Sequences assembled by these reads would  
 204 be further filtered as following description in the section of sequences purifying.

205

## 206 **Co-barcoding reads assembling**

207 Reads of a single species with abundance higher than 10× were assembled by

Supernova (version 2.1.1), which is a co-barcoding *de novo* assembler for single large eukaryotic genomes with high performances. Supernova was designed for linked-reads of 10X Genomics, which have different barcode sequences and formats from stLFR reads. Thus, we converted the stLFR reads into linked-reads FASTQ files with an in-house script. Additionally, the parameter *--accept-extreme-coverage* was set to *yes* to adapt to large coverage depth differences.

## Sequences purifying

The similarity between whole genomes based on the alignment fraction (AF) and average nucleotide identity (ANI) have been commonly adopted to circumscribe species [3, 4]. MetaTrass also used the parameters of AF and ANI between assembled contigs and the reference to purify the sequences assembled by the refined co-barcoding reads. ANI was calculated independently for each alignment. AF was defined as the ratio of total alignment length to the total contig length. The alignments with ANI larger than a threshold were counted. In our practice, we set ANI threshold to 90%, and AF threshold to 50%. The alignments between contigs and references were generated by QUAST (version 5.0.2) [45] with default parameters, except that the identity threshold to obtain valid alignment was set to 90%.

## Combinations of assembling first and binning later

In a standard analysis of NGS metagenomics dataset, the combination of *de novo* genome assembling first and binning later was commonly adopted. We compared different combinations to MetaTrass by analyzing the mock and four gut samples. In our tests, the stLFR co-barcoding reads were assembled by NGS assemblers including IDBA-UD (version 1.1.3), MEGAHIT (version 1.1.3), and MetaSPAdes (version 3.10.1) or co-barcoding assemblers including Supernova [37], Athena (version 1.3.0) [29], and CloudSPAdes (version 3.13.1) [40]. Then, all these draft assemblies were binned by two genome bidders, MetaBAT2 (version 2.12.1)[21] and Maxbin2.0 (version 2.2.5) [20]. Since CloudSPAdes and Athena were not designed for stLFR reads, we made an appropriate format conversion with an in-house script where



LongRanger (version 2.2.2) [46] was used. In genome assembling, Supernova was run with the same parameters as those have been adopted in MetaTrass. IDBA-UD, MEGAHIT, MetaSPAdes, Athena, and CloudSPAdes were run with default parameters. All the assembling results were deposited into CNGB Sequence Archive (CNSA) [47] (<https://db.cngb.org/cnsa/>) of China National GeneBank DataBase (CNGBdb) [48] with accession number CNP0002163. In genome binning, MetaBAT2 and Maxbin2.0 were run with default parameters.

245

## 246 **Evaluations**

Both reference-based and reference-free assessments were used to evaluate the quality of assemblies obtained using different strategies. For the mock microbial community with definite references, the reference-based tool QUAST was used to evaluate contiguity and accuracy of metagenomics assemblies. Minimap2 was used to map assemblies to references and get valid alignments with the identity threshold of 95%. Then, the statistics such as genome fraction, NG50/NGA50, and number of misassemblies were assessed from the alignments with default parameters. For the real gut microbial communities, the reference-free tool CheckM (version 1.1.2) [49] were run with default parameters to evaluate the completeness and contamination of each genome from metagenomics assemblies in addition to QUAST. Following the guidance proposed in CheckM, we defined a high-quality assembly if it has >90% completeness and <5% contamination and a medium-quality assembly if it has >50% completeness and <10% contamination and does not meet the high-quality criterion. In addition, the statistics of each genome such as N50, genome size, and taxonomic rank were also obtained by CheckM, where the taxonomic rank was used to demonstrate the resolution of a genome bin.

263

## 264 **Variation and phylogenetic analysis**

All the high-quality genomes assembled by MetaTrass were used to call variations for the four gut samples. We aligned each genome to the corresponding reference using minimap2 (2.17-r974-dirty) with parameters (-x asm5) to prevent an alignment

267

268 extending to regions with diversity >5%. SAMtools (version 1.9) [50] and PAFtools  
 269 were used to convert the BAM file of initial unsorted alignments into a PAF file of  
 270 sorted alignments. We identified variations using the “call” module in PAFtools with  
 271 parameters (-L 10000) to filter out the alignments shorter than 10,000 bp. SNVs only  
 272 referred to single nucleotide substitutions, excluded single-base insertions or deletions.  
 273 Insertions or deletions with length shorter than 50 bp were defined as small indels,  
 274 and the others were large indels. In determination of shared variations among species  
 275 in different samples, the position and sequence information of a variation were used.  
 276 When variation information is the same for species genomes in different samples, the  
 277 variation was shared.

278 We used the “classify\_wf” function of GTDB-tk (version 0.3.1) [51] to conduct  
 279 taxonomic annotation of the genome bins obtained using the common strategies with  
 280 default parameters. Considering the procedure of UHGG database construction [4],  
 281 genome bins were assigned at the species level if the AF to the close species  
 282 representative genomes was higher than 30% and ANI was higher than 95%. We used  
 283 FastTree (version 2.1.10) [52] to build maximum-likelihood phylogenetic trees of the  
 284 high-quality genomes assembled by MetaTrass. The input of protein sequence  
 285 alignments was produced by GTDB-Tk using marker gene set of 120 bacteria and 122  
 286 archaea. Interactive Tree of Life (iTOL version 4.4.2) [53] was used to visualize and  
 287 annotate trees.

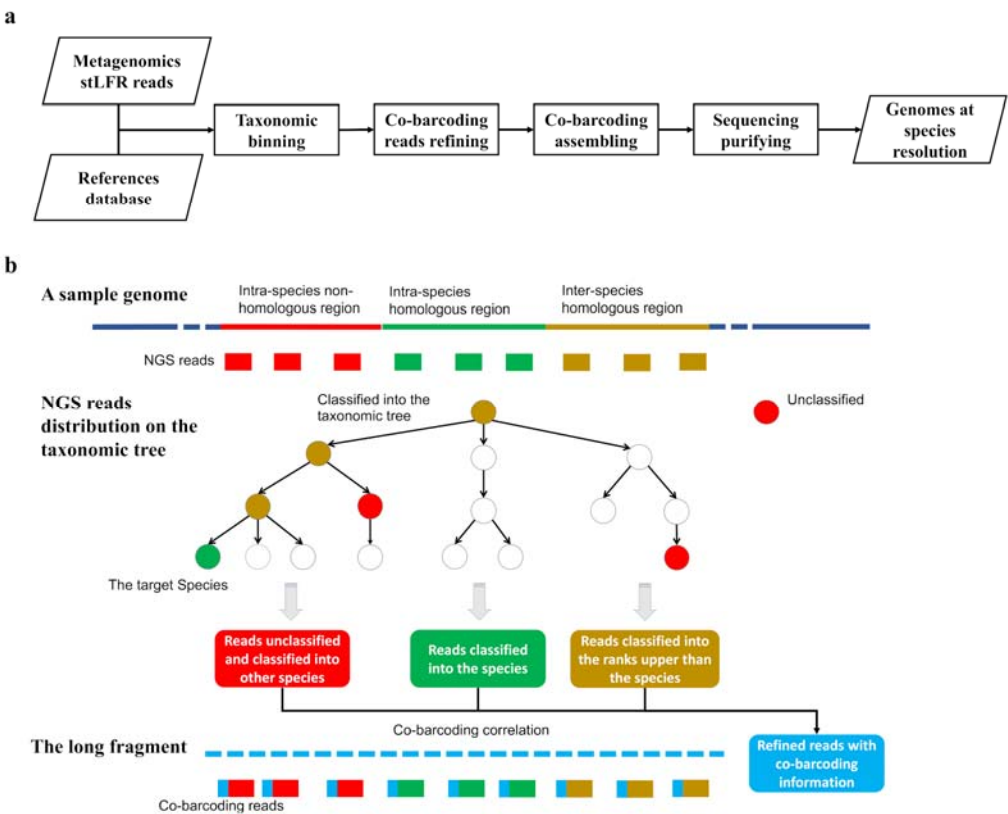
288

## 289 **Results and discussion**

### 290 **MetaTrass pipeline**

291 In this work, we developed a metagenomics assembling pipeline named MetaTrass to  
 292 combine the references and long-range co-barcoding information of stLFR library.  
 293 From the flowchart (**Figure 1a**), the taxonomic binning of stLFR reads were  
 294 processed before the genome assembling, and different from the previous common  
 295 combination strategies of assembling first and binning later. In taxonomic binning, the  
 296 metagenomics stLFR reads were classified into different taxonomic ranks by Kraken2  
 297 [42]. Since the phylogenetic relations among references were used in Kraken2, only

the reads from intra-species homologous region of a sample genome can be classified into the target species, but the reads from inter-species homologous and intra-species non-homologous regions were not classified effectively (Figure 1b). The reads from inter-species homologous regions were classified into the higher ranks of the target species and those from intra-species non-homologous region were unclassified or classified into irrelevant ranks. Totally, about 10% of the reads were classified into the high ranks for the four human gut datasets and about 9% of the reads were unclassified (Table S1). In co-barcoding refining, the co-barcoding correlation between the reads from intra-species homologous region and those from intra-species non-homologous and inter-species homologous region was used to refine the final reads set for a target species (Figure 1b). The barcodes of the intra-species homologous reads were firstly extracted as the candidate barcodes. Then, the final barcodes were collected by a constraint of data size and the quality of co-barcoding information. Finally, the reads with a barcode belong to the final barcodes were gathered to form the refined reads set for the target species. The constraint of data size was set to reduce computational consumption for the species with extremely high abundance. Since the barcodes with more reads classified into the target species are more possible to retain the long-range genomic information, the quality of co-barcoding information of a barcode was quantified by the number of reads classified into the species and the number ratio of these reads to total reads. In co-barcoding assembling, the refined reads of each species were independently assembled by Supernova. In practice, multiple long fragments from different species would share the same barcode in real stLFR libraries (Figure S1). Thus, the impure sequences assembled by the false positive reads from non-target species should be removed finally according to AF and ANI values of alignments between the assembly and references. Overall, the comprehensive use of co-barcoding information and references in our approach could reduce the false negative effects of taxonomic binning and the false positive effects of co-barcoding read refining.



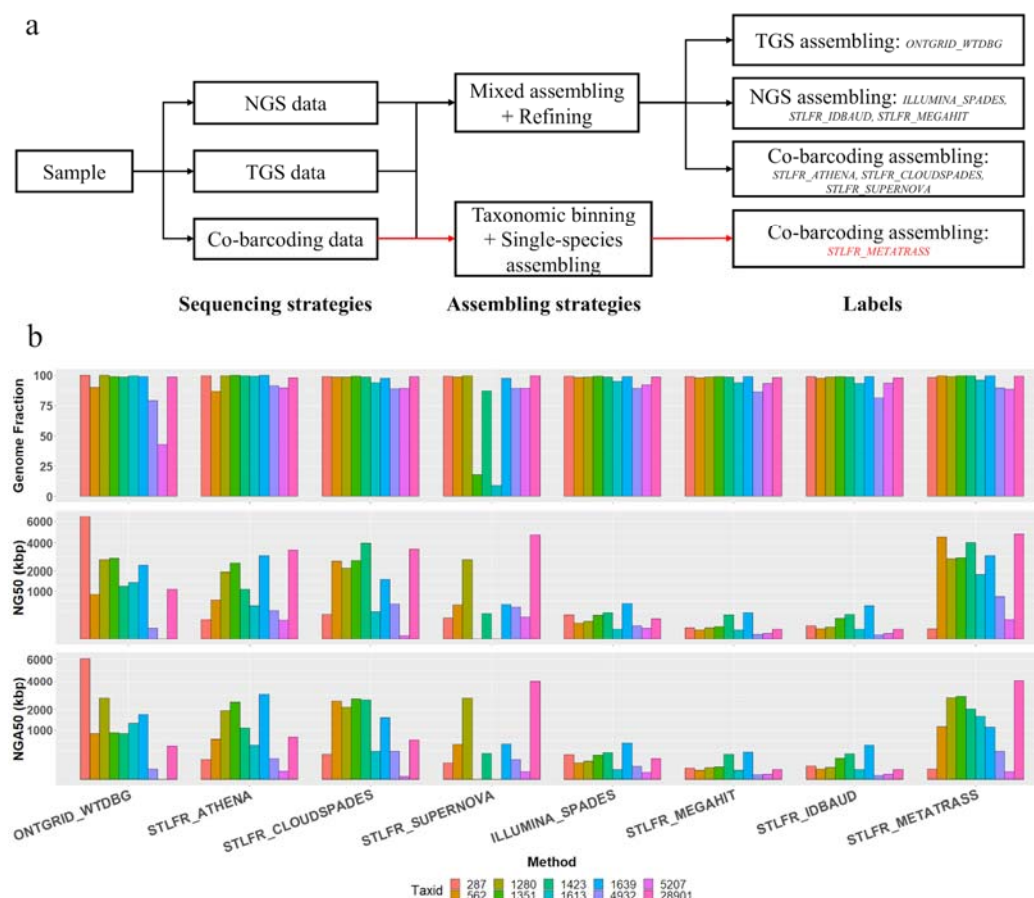
**Figure 1** Flowchart and scheme of MetaTrass. a) Flowchart of MetaTrass assembling pipeline. b) Scheme of the homologous relation and co-barcoding correlation of difference reads sets classified by taxonomic binning.

### Assembly of the mock microbiome

The strategy of binning first and assembling later have been widely adopted to assemble haplotype genomes for eukaryotes with large sizes [54, 34]. But it has been rarely used to assemble metagenomes. We firstly applied MetaTrass to assemble stLFR read sets of the mock microbial community. Totally, up to 99.4% of reads were assigned to different datasets of species due to the simplicity of the microbial community with low inter-species homology and intra-species non-homology (Table S2). To investigate the efficiency of our strategy, we compared it with the mainstream mixed assembling strategies (Figure 2a). Besides the MetaTrass, the stLFR reads were also directly assembled by IDBA-UD, MEGAHIT, Supernova, CloudSPAdes, and Athena in the mixed assembling. Additionally, the optimal mixed assemblies of ONT reads and Illumina NGS reads in Nicholls's work [55] were also

343 used to make a comparison, where the ONT result was assembled by WTDBG and  
344 the NGS result was by SPAdes. The draft genome of each species in a mixed  
345 assembly was extracted by our sequence purifying module.

346 Overall, our pipeline was superior in the production of draft genomes with high  
347 genome fractions and long contiguity (**Figure 2**). Two species *Enterococcus faecalis*  
348 and *Lactobacillus fermentum* were incompletely assembled by Supernova, and their  
349 genome fractions were only 17.7% and 8.9%. However, both species were properly  
350 recovered in MetaTrass, indicating that the assembling complexity caused by uneven  
351 abundances was reduced by taxonomic binning. All the assemblies by MetaTrass  
352 showed high genome fraction as those by NGS and co-barcoding assemblers designed  
353 for metagenome, which were higher than those of ONT assemblies. Compared to  
354 NGS assemblers, the co-barcoding and TGS assembler generated draft assemblies  
355 with significantly better contiguity, where Metatrass generated the best performance.  
356 MetaTrass produced seven draft genomes with NG50 around 2 Mb. Furthermore, the  
357 accuracy was guaranteed by MetaTrass, which obtained the most assemblies with  
358 NGA50 around 2 Mb. Meanwhile, assemblies by MetaTrass had less assembly errors  
359 compared to ONT assemblies (Figure S2). The average mismatch and indel numbers  
360 per 100 kb in assemblies with stLFR reads were 60 and 10, which were obviously  
361 smaller than that of the ONT assemblies.



**Figure 2** Scheme and evaluations for different strategies. a) Difference labels of the assemblies based on different sequencing and assembling strategies. b) Genome fraction, NG50 and NGA50 evaluated by QUASt for the assemblies.

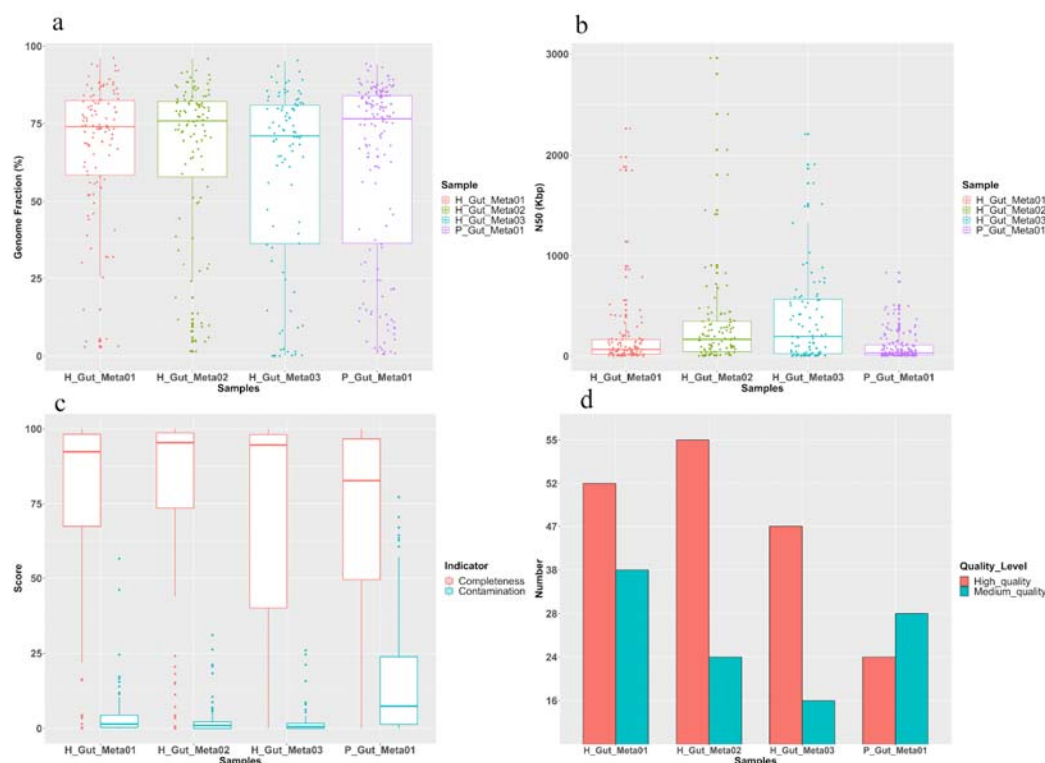
### Assembly of four human gut microbiomes

To evaluate the robustness of our approach, we applied MetaTrass to four human faecal samples. The comprehensive genome references of UHGG were used to classify NGS reads by Kraken2, and the community compositions were estimated by the classified reads at different taxonomic ranks (Figure S3-S6). The three healthy samples had a similar microbial community, where the major microbiomes were from *Firmicutes* A phylum. This microbial community was different from the patient microbial community dominated by *Proteobacteria* which is strongly correlated with the enteric diseases caused by dysbiosis in gut microbiota [56]. The total numbers of species with higher than 10× abundance were 113, 108, 93, and 158 in

377 H\_Gut\_Meta01, H\_Gut\_Meta02, H\_Gut\_Meta03, and P\_Gut\_Meta01 samples,  
378 respectively. The relations between these abbreviated notations and detailed sample  
379 information were described in the section of Materials and methods.

380 The genome fraction of an assembly to the reference is used to evaluate the  
381 completeness in single genome assembling. The genome fraction for all samples  
382 widely ranges from 0% to 90%, and the distributions of H\_Gut\_Meta01 and  
383 H\_Gut\_Meta02 were more concentrated than those of H\_Gut\_Meta03 and  
384 P\_Gut\_Meta01 (**Figure 3a**). However, more than half of the assembled genomes were  
385 with a genome fraction of at least 50%. Considering the large genetic diversity  
386 between sample genomes and the references [7], these results indicated that our  
387 pipeline could assemble complete genomes for species abundance of higher than 10×.  
388 The genetic diversity was also proved by the significant differences in genome  
389 fraction and the ratio of assembled length to the reference length among the four  
390 samples (Figure S7). The distributions of genomes N50 were generally dispersed, and  
391 the medians of H\_Gut\_Meta02 and H\_Gut\_Meta03 were obviously higher than those  
392 of H\_Gut\_Meta01 and P\_Gut\_Meta01 (Figure 3b). Nevertheless, the third quartiles in  
393 the box plots for the samples were larger than 100 kb, demonstrating that our pipeline  
394 had a strong capability to generate draft genomes with high contiguity. Note that for  
395 these three healthy samples plenty of ultra-long draft genomes (N50>1 Mb) was  
396 obtained, which provide possibilities to study the large genome difference in the  
397 microbiome.





**Figure 3** QUAST and CheckM evaluations of MetaTrass assemblies for the four human gut samples. a) Genome fraction. b) Scaffold N50. c) Box plot of completeness and contamination. d) Number of high- and medium-quality genomes.

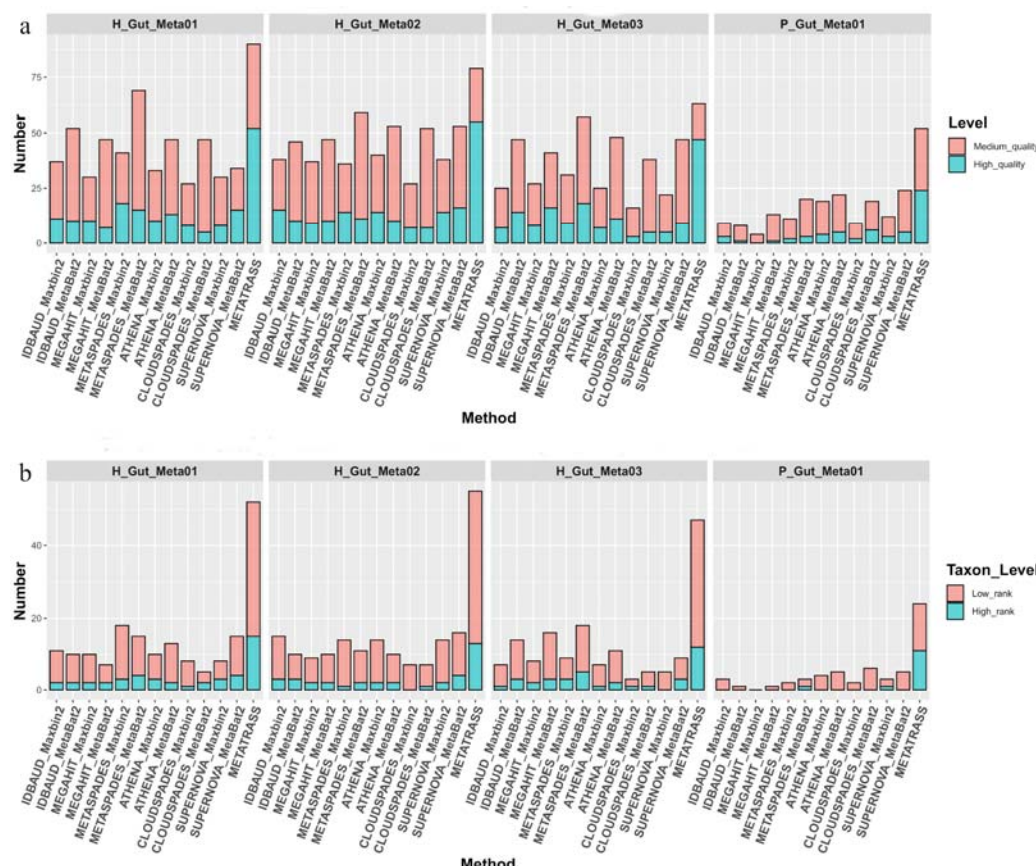
Considering the intra-species genetic diversity, we also evaluated the quality of metagenomics assemblies based on the conserved marker genes by CheckM. The completeness medians of three healthy samples were larger than 92%, and the contamination medians were smaller than 2% (Figure 3c). The completeness of the patient sample was about 83%, and the contamination median was about 7% (Figure S8). Meanwhile, a great number of high- and medium-quality genomes were assembled by MetaTrass for the four samples (Figure 3d). 52 high-quality and 37 medium-quality genomes were produced for H\_Gut\_Meta01, and 55 and 24 for H\_Gut\_Meta02, and 47 and 16 for H\_Gut\_Meta03, and 24 and 28 for P\_Gut\_Meta01, respectively.

### Comparison to the common combination strategy



415 To further evaluate our approach's efficiency, we compared it with common  
416 combinations of assembling tools and genome binning tools as listed in the section of  
417 Datasets and Methods. It should be noted that currently, there are still no genome  
418 binning tools to directly exploit the co-barcoding information. By counting the  
419 number of bins with completeness >50% and at least one conserved marker genes  
420 (Table S3), we observed that MetaTrass perform best of all these methods. Especially  
421 for P\_Gut\_Meta01, the optimal combination between Supernova and Maxbin2.0  
422 obtained 66 bins with completeness higher than 50%, but it was significantly less than  
423 117 obtained by MetaTrass.

424 By comprehensively analyzing the completeness, contamination and taxonomic  
425 rank of each bin, we assessed MetaTrass and common strategies in the ability to get  
426 high- and medium-quality genomes and resolution of taxonomic rank (**Figure 4**). For  
427 different samples, the best combination to produce the optimal results was different.  
428 The combinations of MetaSPAdes and Maxbin2.0, Supernova and MetaBAT2,  
429 MetaSPAdes and MetaBAT2, and Athena and MetaBAT2 is optimal for  
430 H\_Gut\_Meta01, H\_Gut\_Meta02, H\_Gut\_Meta03, and P\_Gut\_Meta01, respectively.  
431 For the four samples, the optimal results of the common strategies were still inferior  
432 to those of MetaTrass. For the example of H\_Gut\_meta01, the combination of  
433 MetaSPAdes and Maxbin2.0 produced 41 high- and medium-quality genomes, which  
434 was significantly less than 90 obtained by MetaTrass. There were only 3 out of totally  
435 18 high-quality genomes with a taxonomic rank lower than the order, but 15 out of 52  
436 for MetaTrass. Comparing the strategies only using NGS read information, the  
437 combination strategies of co-barcoding assembler and binner showed no obvious  
438 advantages in generating genomes with high quality and resolution, but MetaTrass  
439 was significantly superior to them. These results demonstrated that the usage of co-  
440 barcoding information in MetaTrass was more efficient and accurate than those in a  
441 mixed assembling.



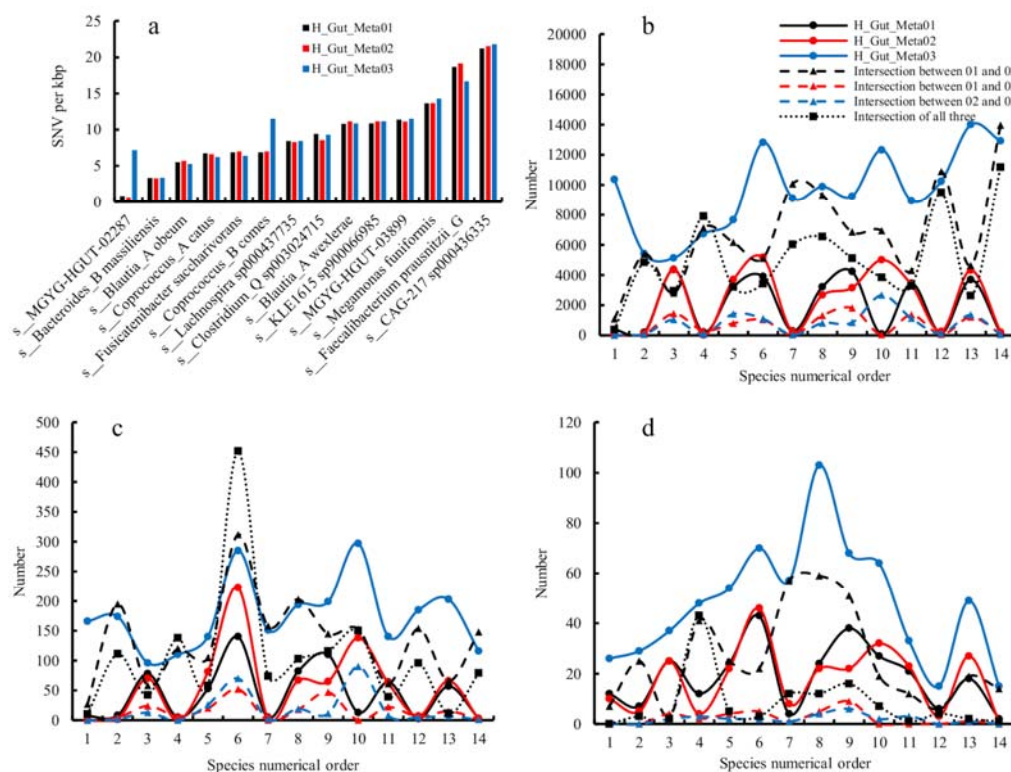
**Figure 4** Comparison of metagenome assembling for different methods. a) Number of high- and medium-quality genomes assembled with different methods. b) Number of high-quality genomes with high- and low-rank with different methods.

The human gut microbiome composition attracts much attention due to its strong correlation with personality traits [57]. To compare the microbiome composition structures of the high-quality genomes with different methods, we uniformly classified the high-quality genome bins into species using GTDB-tk. Using the large number of high-quality genomes obtained by MetaTrass, the phylogenetic trees of these genomes were constructed and the corresponding N50 were attached in the left histogram as shown in **Figure 5**. Meanwhile, the high-quality genome bins obtained by the common strategies were marked in red in the middle heat map (Figure 5), if the genome of the same species were also assembled by MetaTrass. The topology of the phylogenetic tree of genomes assembled by MetaTrass gave comprehensive insights

457 of the microbial composition structure. From the trees in Figure 5 and Figure S9-S11,  
458 the numbers of the order with high-quality genomes assembled by MetaTrass were 9,  
459 11, 7, and 7 for H\_Gut\_Meta01, H\_Gut\_Meta02, H\_Gut\_Meta03, and P\_Gut\_Meta01,  
460 respectively. Notably, some orders contained more than 5 high-quality genomes, and  
461 this provide convenience to study the microbiome structure at the genome-wide scale.  
462 For the sample of H\_Gut\_Meta01 (Figure 5), there were 27 high-quality genomes  
463 classified into *Lachnospirales* order and 14 into *Oscillospirales*. These two were  
464 exactly the dominating orders according to the taxonomic abundance distribution.  
465 Similar results were obtained for the other two healthy samples (Figure S9 and S10),  
466 indicating that the microbiome with higher sequencing coverage could be better  
467 assembled in MetaTrass. In contrast, the orders with more than 5 high-quality  
468 genomes were *Enterobacterales* and *Actinomycetales* for P\_Gut\_Meta01 (Figure S11).  
469 The obvious difference between the healthy and patient samples was consistent with  
470 the microbial compositions differences observed in the taxonomic binning results.  
471 MetaTrass successfully assemble most of the high-quality genomes of all common  
472 combinations in our tests. For instance, they generated 137 genome bins, while only  
473 25 genome bins were not assembled by MetaTrass (Figure 5). From the heat maps,  
474 most of the common strategies could assemble draft genomes for each order, but the  
475 total numbers in each order were relatively small. The maximal number of genomes in  
476 one order was 6 and obtained by the combination of Supernova and MetaBAT2 for  
477 *Lachnospirales*. Moreover, 146 of 179 high-quality genomes were with N50 values  
478 larger than 100 kb, demonstrating that MetaTrass had a strong ability to improve the  
479 contiguity of assemblies.



502       Based on the taxonomic information of high-quality genomes, we found 15  
503 species shared by three samples, where 14 species appeared in the three healthy  
504 samples but only one species of *Escherichia* appeared in the patient and two healthy  
505 samples. By analyzing the SNV density and intersection of variations between  
506 different samples for each species in three healthy samples, we further investigated  
507 the genetic diversity between species from different samples. The SNV densities were  
508 different for different species even in the same sample, but similar for the same  
509 species in different sample (**Figure 6a**). From Figure 6b to 6d, the number of unique  
510 and shared variations in different types significantly fluctuated for different species,  
511 but their difference among samples showed great consistency. The total shared  
512 numbers between H\_Gut\_Meta01 and H\_Gut\_meta02 were obviously more than  
513 those between H\_GutMeta03 and the other two samples for all variations.  
514 Furthermore, the ratio of large indels shared by all three samples to the total number  
515 was much smaller than those of SNVs and small indels. These results demonstrated  
516 that large variations were more specific than small variations in the huge genetic  
517 diversity between different samples, were consistent with the observation in the study  
518 of association between host health and structural variations in gut microbiome [58].



**Figure 6** SNV density and number of unique and shared variations for each species appearing in all three healthy samples. a) is the SNV density. b), c) and d) are the number of SNVs, small and large indels, respectively. The species numerical order in b), c) and d) corresponds to the appearance order of species from left to right in a).

## Computational performance

Runtime and used thread number of each assembler were recorded for all the human gut datasets (**Table 1**). Most of the assemblers were test on 24 Intel(R) Xeon(R) Silver 4116 CPU @ 2.10GHz, except for Athena and Supernova which were test on HPC Cluster for their large memory requirements. The thread number used in each assembler was the same for different samples. The time consumption of the format conversion from stLFR reads to 10X linked-reads were not included, and was about 500 minutes for dataset with 50 Gb with one thread. We found that MetaTrass was less time consuming than Athena but more than other assemblers. This may come from that both MetaTrass and Athena contained many sub-assembling, which took most of the time among all sub-processes in MetaTrass (Table S4). Since the sub-



assembling was independent, it could be run in parallel to further speed up the assembling by increasing the parallel number and the parallel number was 8 in default.

## Conclusion

High-quality genomes at species level are strongly demanded to investigate the genetic origins of diseases associated with the human gut, but how to get sufficient number of them in one sample is still a challenge due to the inter-species repeats and uneven abundance in metagenomics assembling. In this work, we developed a tool to get high-quality genomes with high taxonomic resolutions by combining the co-barcoding information with public references. Compared with the common combination strategies, our pipeline generated a large number of high-quality genomes for the human microbiome co-barcoding datasets. Meanwhile, plenty of draft genomes were also assembled with NG50 values of larger than 1 Mb, some of which were even longer than the references for both mock and human gut datasets. For all the four real gut samples, 178 draft genomes with high completeness and low contamination were generated, but their genome fractions relative to the references were low. The differences between the sample genomes assembled by MetaTrass and the reference genomes demonstrated that the co-barcoding information could be used to reduce the false negative reads in taxonomic binning. These reads retrieved from inter-species homologous and intra-species non-homologous regions by co-barcoding refining could significantly improve the assembly results. For the patient sample, the number of high-quality genomes with long contiguity assembled by MetaTrass was significantly larger than that without co-barcoding refining (Figure S13).

The efficiency of our pipeline depended on the co-barcoding information quality including the read coverage and length of long fragments. By aligning reads to the species reference, we calculated the genome fraction with different read coverage depths for different read sets including the taxonomic reads, the refined reads, and all reads. According to the genome fraction with high coverage depths, we evaluate the efficiency of the co-barcoding refining. From the results of species with medium abundance in P\_Gut\_Meta01 (Figure S5), We observed that the fraction with high depths of the refined reads was higher than those of the taxonomic reads, but still lower than those of all aligned reads. These results indicated that there were still some false negative reads introduced by the low coverage or short length of long fragments.

569 Thus, improvements on co-barcoding library and the co-barcoding refining would  
570 improve the performance of MetaTrass.

571 In summary, the application of MetaTrass in human gut samples showed great  
572 promise of generating high-quality genomes for real complex microbial community at  
573 a high resolution. With the increasing number of reference genomes from various  
574 microbial communities and the development of co-barcoding sequencing library, the  
575 combination strategy of binning first and assembling later in MetaTrass will be  
576 extended and facilitate the investigation of the association between host phenotypes  
577 and microbial genotypes for different microbial communities.

578  
579

## 580 **Acknowledgements**

581 This research was supported by the National Key Research and Development  
582 Program of China (2018YFD0900301-05), and Science Technology and Innovation  
583 Committee of Shenzhen Municipality of China (SGDX20190919142801722). We  
584 would thank Yufen Huang and many other BGI-Shenzhen employees for fruitful  
585 discussions in the development and performance test.  
586

## 587 **Conflicts of interest**

588 All authors are employees of the BGI group.  
589

## 590 **Authors' contributions**

591 Li Deng, Guangyi Fan and Yanwei Qi contributed to the software design. Yanwei Qi,  
592 Shengqiang Gu, Yue Zhang and Lidong Guo contributed to the software  
593 implementation. Li Deng, Yanwei Qi, Shengqiang Gu, Mengyang Xu and Jianwei  
594 Chen contributed to data analyses. Xiaofang Chen, Ou Wang and Xiaodong Fang  
595 contribute to the data curation, collection. Guangyi Fan, Li Deng and Xin Liu  
596 contributed to the benchmarking design. All authors contributed to the manuscript  
597 writing. Li Deng and Guangyi Fan supervised the project. All authors read and  
598 approved the final manuscript.



599

## 600 **Data availability statement**

601 MetaTrass is freely available at <https://github.com/BGI-Qingdao/MetaTrass>. The  
602 assembling results of the four human faecal samples were deposited into CNSA  
603 (<https://db.cngb.org/cnsa/>) of CNGBdb with accession number CNP0002163 and  
604 available from authors upon reasonable request and with permission of CNGBdb. The  
605 metagenomics stLFR datasets used in the study were available from the  
606 corresponding author on reasonable request.

607

## 608 **References**

- 609 1. Schloss, Patrick D, and Jo Handelsman. 2005. "Metagenomics for studying  
610 unculturable microorganisms: cutting the Gordian knot." *Genome Biology* 6: 1-4.
- 611 2. Qin, Junjie, Ruiqiang Li, Jeroen Raes, Manimozhiyan Arumugam, Kristoffer  
612 Solvsten Burgdorf, Chaysavanh Manichanh, Trine Nielsen, *et al.* 2010. "A human  
613 gut microbial gene catalogue established by metagenomic sequencing." *Nature*  
614 464: 59-65.
- 615 3. Parks, Donovan H, Maria Chuvochina, Pierre-Alain Chaumeil, Christian Rinke,  
616 Aaron J Mussig, and Philip Hugenholtz. 2020. "A complete domain-to-species  
617 taxonomy for Bacteria and Archaea." *Nature Biotechnology* 38: 1079-86.
- 618 4. Almeida, Alexandre, Stephen Nayfach, Miguel Boland, Francesco Strozzi,  
619 Martin Beracochea, Zhou Jason Shi, Katherine S. Pollard, *et al.* 2021. "A unified  
620 catalog of 204,938 reference genomes from the human gut microbiome." *Nature*  
621 *Biotechnology* 39: 105-14.
- 622 5. Sheth, Ravi U, Mingqiang Li, Weiqian Jiang, Peter A Sims, Kam W Leong, and  
623 Harris H Wang. 2019. "Spatial metagenomic characterization of microbial  
624 biogeography in the gut." *Nature Biotechnology* 37: 877-83.
- 625 6. Martino, Cameron, Liat Shenhav, Clarisse A. Marotz, George Armstrong, Daniel  
626 McDonald, Yoshiki Vázquez-Baeza, James T. Morton, *et al.* 2021. "Context-  
627 aware dimensionality reduction deconvolutes gut microbial community

- 628       dynamics." *Nature Biotechnology* 39: 165-8.
- 629   7.   Van Rossum, Thea, Pamela Ferretti, Oleksandr M Maistrenko, and Peer Bork.  
630       2020. "Diversity within species: interpreting strains in microbiomes." *Nature*  
631       *Reviews: Microbiology* 18: 491-506.
- 632   8.   Olm, Matthew R, Alexander Crits-Christoph, Keith Bouma-Gregson, Brian A  
633       Firek, Michael J Morowitz, and Jillian F Banfield. 2021. "inStrain profiles  
634       population microdiversity from metagenomic data and sensitively detects shared  
635       microbial strains." *Nature Biotechnology* 39: 727-36.
- 636   9.   Leimbach, Andreas, Jörg Hacker, and Ulrich Dobrindt. 2013. "E. coli as an all-  
637       rounder: the thin line between commensalism and pathogenicity." *Current Topics*  
638       *in Microbiology and Immunology* 358: 3-32.
- 639   10. Pierce, Jessica V, and Harris D Bernstein. 2016. "Genomic diversity of  
640       enterotoxigenic strains of *Bacteroides fragilis*." *PLoS One* 11: e0158171.
- 641   11. Yao, G., W. Zhang, M. Yang, H. Yang, J. Wang, H. Zhang, L. Wei, Z. Xie, and W.  
642       Li. 2020. "MicroPhenoDB Associates Metagenomic Data with Pathogenic  
643       Microbes, Microbial Core Genes, and Human Disease Phenotypes." *Genomics*  
644       *Proteomics Bioinformatics* 18: 760-72.
- 645   12. Welch, R. A., V. Burland, G. Plunkett, P. Redford, P. Roesch, D. Rasko, E. L.  
646       Buckles, *et al.* 2002. "Extensive mosaic structure revealed by the complete  
647       genome sequence of uropathogenic *Escherichia coli*." *Proceedings of the*  
648       *National Academy of Sciences* 99: 17020-4.
- 649   13. Sczyrba, Alexander, Peter Hofmann, Peter Belmann, David Koslicki, Stefan  
650       Janssen, Johannes Dröge, Ivan Gregor, *et al.* 2017. "Critical assessment of  
651       metagenome interpretation—a benchmark of metagenomics software." *Nature*  
652       *Methods* 14: 1063-71.
- 653   14. Breitwieser, Florian P, Jennifer Lu, and Steven L Salzberg. 2019. "A review of  
654       methods and databases for metagenomic classification and assembly." *Briefings*  
655       *in Bioinformatics* 20: 1125-36.
- 656   15. Eun Kang, J., A. Ciampi, and M. Hijri. 2020. "SeSaMe: Metagenome Sequence  
657       Classification of Arbuscular Mycorrhizal Fungi-associated Microorganisms."

658        *Genomics Proteomics Bioinformatics* 18: 601-12.

659    16. Li, Dinghua, Chi-Man Liu, Ruibang Luo, Kunihiko Sadakane, and Tak-Wah Lam.

660        2015. "MEGAHIT: an ultra-fast single-node solution for large and complex

661        metagenomics assembly via succinct de Bruijn graph." *Bioinformatics* 31: 1674-6.

662    17. Nurk, Sergey, Dmitry Meleshko, Anton Korobeynikov, and Pavel A Pevzner.

663        2017. "metaSPAdes: a new versatile metagenomic assembler." *Genome Research*

664        27: 824-34.

665    18. Peng, Yu, Henry CM Leung, Siu-Ming Yiu, and Francis YL Chin. 2012. "IDBA-

666        UD: a de novo assembler for single-cell and metagenomic sequencing data with

667        highly uneven depth." *Bioinformatics* 28: 1420-8.

668    19. Wu, Yu-Wei, and Yuzhen Ye. 2011. "A novel abundance-based algorithm for

669        binning metagenomic sequences using 1-tuples." *Journal of Computational*

670        *Biology* 18: 523-34.

671    20. Wu, Yu-Wei, Blake A Simmons, and Steven W Singer. 2016. "MaxBin 2.0: an

672        automated binning algorithm to recover genomes from multiple metagenomic

673        datasets." *Bioinformatics* 32: 605-7.

674    21. Kang, Dongwan D, Jeff Froula, Rob Egan, and Zhong Wang. 2015. "MetaBAT,

675        an efficient tool for accurately reconstructing single genomes from complex

676        microbial communities." *PeerJ* 3: e1165.

677    22. Bertrand, Denis, Jim Shaw, Manesh Kalathiyappan, Amanda Hui Qi Ng, M.

678        Senthil Kumar, Chenhao Li, Mirta Dvornicic, *et al.* 2019. "Hybrid metagenomic

679        assembly enables high-resolution analysis of resistance determinants and mobile

680        elements in human microbiomes." *Nature Biotechnology* 37: 937-44.

681    23. Chin, Chen-Shan, David H. Alexander, Patrick Marks, Aaron A. Klammer, James

682        Drake, Cheryl Heiner, Alicia Clum, *et al.* 2013. "Nonhybrid, finished microbial

683        genome assemblies from long-read SMRT sequencing data." *Nature Methods* 10:

684        563-9.

685    24. Kolmogorov, Mikhail, Derek M. Bickhart, Bahar Behsaz, Alexey Gurevich,

686        Mikhail Rayko, Sung Bong Shin, Kristen Kuhn, *et al.* 2020. "metaFlye: scalable

687        long-read metagenome assembly using repeat graphs." *Nature Methods* 17: 1103-

- 688 10.
- 689 25. DeMaere, Matthew Z, and Aaron E Darling. 2019. "bin3C: exploiting Hi-C  
690 sequencing data to accurately resolve metagenome-assembled genomes."  
691 *Genome Biology* 20: 1-16.
- 692 26. Cleary, Brian, Ilana Lauren Brito, Katherine Huang, Dirk Gevers, Terrance Shea,  
693 Sarah Young, and Eric J Alm. 2015. "Detection of low-abundance bacterial  
694 strains in metagenomic datasets by eigengenome partitioning." *Nature*  
695 *Biotechnology* 33: 1053-60.
- 696 27. Peters, Brock A., Bahram G. Kermani, Andrew B. Sparks, Oleg Alferov, Peter  
697 Hong, Andrei Alexeev, Yuan Jiang, *et al.* 2012. "Accurate whole-genome  
698 sequencing and haplotyping from 10 to 20 human cells." *Nature* 487: 190-5.
- 699 28. Adey, Andrew, Jacob O Kitzman, Joshua N Burton, Riza Daza, Akash Kumar,  
700 Lena Christiansen, Mostafa Ronaghi, *et al.* 2014. "In vitro, long-range sequence  
701 information for de novo genome assembly via transposase contiguity." *Genome*  
702 *Research* 24: 2041-9.
- 703 29. Bishara, Alex, Eli L Moss, Mikhail Kolmogorov, Alma E Parada, Ziming Weng,  
704 Arend Sidow, Anne E Dekas, Serafim Batzoglou, and Ami S Bhatt. 2018. "High-  
705 quality genome sequences of uncultured microbes by assembly of read clouds."  
706 *Nature Biotechnology* 36: 1067-75.
- 707 30. Wang, Ou, Robert Chin, Xiaofang Cheng, Michelle Wu, Qing Mao, Jingbo Tang,  
708 Yuhui Sun, *et al.* 2019. "Efficient and unique cobarcoding of second-generation  
709 sequencing reads from long DNA molecules enabling cost-effective and accurate  
710 sequencing, haplotyping, and de novo assembly." *Genome Research* 29: 798-808.
- 711 31. Chen, Zhoutao, Long Pham, Tsai-Chin Wu, Guoya Mo, Yu Xia, Peter L. Chang,  
712 Devin Porter, *et al.* 2020. "Ultralow-input single-tube linked-read library method  
713 enables short-read second-generation sequencing systems to routinely generate  
714 highly accurate and economical long-range sequencing information." *Genome*  
715 *Research* 30: 898-909.
- 716 32. Zheng, Grace X. Y., Billy T. Lau, Michael Schnall-Levin, Mirna Jarosz, John M.  
717 Bell, Christopher M. Hindson, Sofia Kyriazopoulou-Panagiotopoulou, *et al.* 2016.

- 718 "Haplotyping germline and cancer genomes with high-throughput linked-read  
719 sequencing." *Nature Biotechnology* 34: 303-11.
- 720 33. Danko, David C, Dmitry Meleshko, Daniela Bezdan, Christopher Mason, and  
721 Iman Hajirasouliha. 2019. "Minerva: an alignment-and reference-free approach  
722 to deconvolve Linked-Reads for metagenomics." *Genome Research* 29: 116-24.
- 723 34. Xu, Mengyang, Lidong Guo, Xiao Du, Lei Li, Brock A Peters, Li Deng, Ou  
724 Wang, *et al.* 2021. "Accurate haplotype-resolved assembly reveals the origin of  
725 structural variants for human trios." *Bioinformatics* 37: 2095-102.
- 726 35. Bishara, Alex, Yuling Liu, Ziming Weng, Dorna Kashef-Haghighi, Daniel E  
727 Newburger, Robert West, Arend Sidow, and Serafim Batzoglou. 2015. "Read  
728 clouds uncover variation in complex regions of the human genome." *Genome*  
729 *Research* 25: 1570-80.
- 730 36. Guo, Junfu, Chang Shi, Xi Chen, Ou Wang, Ping Liu, Huanming Yang, Xun Xu,  
731 Wenwei Zhang, and Hongmei Zhu. 2021. "stLFRsv: A Germline Structural  
732 Variant Analysis Pipeline Using Co-barcoded Reads." *Frontiers in Genetics* 12:  
733 222.
- 734 37. Weisenfeld, Neil I, Vijay Kumar, Preyas Shah, Deanna M Church, and David B  
735 Jaffe. 2017. "Direct determination of diploid genome sequences." *Genome*  
736 *Research* 27: 757-67.
- 737 38. Yeo, Sarah, Lauren Coombe, René L Warren, Justin Chu, and Inanç Birol. 2017.  
738 "ARCS: scaffolding genome drafts with linked reads." *Bioinformatics* 34: 725-31.
- 739 39. Guo, Lidong, Mengyang Xu, Wenchao Wang, Shengqiang Gu, Xia Zhao, Fang  
740 Chen, Ou Wang, *et al.* 2021. "SLR-superscaffolder: a de novo scaffolding tool for  
741 synthetic long reads using a top-to-bottom scheme." *BMC Bioinformatics* 22: 1-  
742 16.
- 743 40. Tolstoganov, Ivan, Anton Bankevich, Zhoutao Chen, and Pavel A Pevzner. 2019.  
744 "cloudSPAdes: assembly of synthetic long reads using de Bruijn graphs."  
745 *Bioinformatics* 35: i61-i70.
- 746 41. Kuleshov, Volodymyr, Michael P Snyder, and Serafim Batzoglou. 2016.  
747 "Genome assembly from synthetic long read clouds." *Bioinformatics* 32: i216-i24.

748 42. Wood, Derrick E, Jennifer Lu, and Ben Langmead. 2019. "Improved  
749 metagenomic analysis with Kraken 2." *Genome Biology* 20: 1-13.

750 43. Lu, Jennifer, Florian P Breitwieser, Peter Thielen, and Steven L Salzberg. 2017.  
751 "Bracken: estimating species abundance in metagenomics data." *PeerJ Computer*  
752 *Science* 3: e104.

753 44. Danko, David C, Dmitry Meleshko, Daniela Bezdan, Christopher Mason, and  
754 Iman Hajirasouliha. 2019. "Novel Algorithms for the Taxonomic Classification of  
755 Metagenomic Linked-Reads." *bioRxiv* 549667.

756 45. Gurevich, Alexey, Vladislav Saveliev, Nikolay Vyahhi, and Glenn Tesler. 2013.  
757 "QUAST: quality assessment tool for genome assemblies." *Bioinformatics* 29:  
758 1072-5.

759 46. Marks, Patrick, Sarah Garcia, Alvaro Martinez Barrio, Kamila Belhocine, Jorge  
760 Bernate, Rajiv Bharadwaj, Keith Bjornson, *et al.* 2019. "Resolving the full  
761 spectrum of human genome variation using Linked-Reads." *Genome Research* 29:  
762 635-45.

763 47. Guo, X., F. Chen, F. Gao, L. Li, K. Liu, L. You, C. Hua, *et al.* 2020. "CNSA: a  
764 data repository for archiving omics data." *Database* 2020: baaa055.

765 48. Chen, Feng Zhen, Li Jin You, Fan Yang, Li Na Wang, Xue Qin Guo, Fei Gao,  
766 Cong Hua, *et al.* 2020. "CNCBdb: China National GeneBank DataBase." *Hereditas* 42: 799-809.

768 49. Parks, Donovan H, Michael Imelfort, Connor T Skennerton, Philip Hugenholtz,  
769 and Gene W Tyson. 2015. "CheckM: assessing the quality of microbial genomes  
770 recovered from isolates, single cells, and metagenomes." *Genome Research* 25:  
771 1043-55.

772 50. Li, Heng, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer,  
773 Gabor Marth, *et al.* 2009. "The Sequence Alignment/Map format and SAMtools." *Bioinformatics* 25: 2078-9.

775 51. Chaumeil, Pierre-Alain, Aaron J Mussig, Philip Hugenholtz, and Donovan H  
776 Parks. 2019. "GTDB-Tk: a toolkit to classify genomes with the Genome  
777 Taxonomy Database." *Bioinformatics* 36: 1925-7.

- 778 52. Price, Morgan N., Paramvir S. Dehal, and Adam P. Arkin. 2010. "FastTree 2 –  
779 Approximately Maximum-Likelihood Trees for Large Alignments." *PLOS ONE* 5:  
780 e9490.
- 781 53. Letunic, Ivica, and Peer Bork. 2019. "Interactive Tree Of Life (iTOL) v4: recent  
782 updates and new developments." *Nucleic Acids Research* 47: W256-W9.
- 783 54. Koren, Sergey, Arang Rhie, Brian P Walenz, Alexander T Dilthey, Derek M  
784 Bickhart, Sarah B Kingan, Stefan Hiendleder, *et al.* 2018. "De novo assembly of  
785 haplotype-resolved genomes with trio binning." *Nature Biotechnology* 36: 1174-  
786 82.
- 787 55. Nicholls, Samuel M, Joshua C Quick, Shuiquan Tang, and Nicholas J Loman.  
788 2019. "Ultra-deep, long-read nanopore sequencing of mock microbial community  
789 standards." *GigaScience* 8: giz043.
- 790 56. Shin, Na-Ri, Tae Woong Whon, and Jin-Woo Bae. 2015. "Proteobacteria:  
791 microbial signature of dysbiosis in gut microbiota." *Trends in Biotechnology* 33:  
792 496-503.
- 793 57. Johnson, Katerina V-A. 2020. "Gut microbiome composition and diversity are  
794 related to human personality traits." *Human Microbiome Journal* 15: 100069.
- 795 58. Zeevi, David, Tal Korem, Anastasia Godneva, Noam Bar, Alexander Kurilshikov,  
796 Maya Lotan-Pompan, Adina Weinberger, *et al.* 2019. "Structural variation in the  
797 gut microbiome associates with host health." *Nature* 568: 43-8.
- 798 59. Chen, Lianmin, Daoming Wang, Sanzhima Garmaeva, Alexander Kurilshikov,  
799 Arnau Vich Vila, Ranko Gacesa, Trishla Sinha, *et al.* 2021. "The long-term  
800 genetic stability and individual specificity of the human gut microbiome." *Cell*  
801 184: 2302-15.

802  
803  
804  
805  
806  
807

808

809

810

811

812

## 813 **Tables**

814 **Table 1** Runtimes and thread number of each assembler for all the human gut  
815 datasets.

Assembler	Runtime (min)				
	Thread number				
	All samples	H_Gut_meta01	H_Gut_Meta02	H_Gut_Meta03	P_Gut_Meta01
IBDA-UD	6	863	884	911	2657
MEGAHIT	16	179	161	163	611
MetaSPAdes	16	1478	1289	1429	3459
CloudSPAdes	16	1024	1163	1039	2627
Supernova	8	1249	864	1098	6776
Athena	16	13813	8689	6361	--
MetaTrass	16	5145	2631	3147	8363

816 *Note:* The exact runtime of assembling P\_Gut\_Meta01 sample by Athena was not  
817 collected correctly due to several uncontrolled interrupts on HPC cluster.

818

819

## 820 **Supporting information**

821 **Table S1** Read number on different ranks classified by Kraken2 for the four gut  
822 samples.

823 **Table S2** Classified read information of the mock dataset.

824 **Table S3** The overall view of genome bins obtained by MetaTrass and all common  
825 strategies “Comp >50%” means the completeness higher than 50%.

826 **Table S4** The runtime of MetaTrass step by step for all human gut datasets.

827 **Table S5** Genome fraction with different coverage depths for different read sets



828 including the taxonomic read (TR), refined reads by co-barcoding (BR), and total  
829 reads (Total) for five species with medium abundances in P\_Gut\_Meta01.

830 **Figure S1** The probability of barcodes with long fragments from different species for  
831 four gut samples.

832 **Figure S2** Mismatches and Indels of different assemblies for the mock dataset.

833 **Figure S3** Distributions of classified reads at different phyla for four gut samples.

834 **Figure S4** Distributions of classified reads at different classes for four gut samples.

835 **Figure S5** Distributions of classified reads at different orders for four gut samples.

836 **Figure S6** Distributions of classified reads at different families for four gut samples.

837 **Figure S7** Genome faction and ratio of assembly length to reference length of all  
838 species assembled in MetaTrass for four gut samples, and the species are ordered by  
839 the completeness.

840 **Figure S8** Two-dimensional scatter plot of completeness and contamination  
841 evaluated by CheckM for four gut samples.

842 **Figure S9** Phylogenetic tree of the high-quality genomes assembled by MetaTrass for  
843 H\_Gut\_Meta02. The phylogenetic tree is on the left. Distribution of the high-quality  
844 genomes assembled by other methods are colored as red in the middle heat map.  
845 N50 of each high-quality genome is shown in the right histogram.

846 **Figure S10** Phylogenetic tree of the high-quality genomes assembled by MetaTrass  
847 for H\_Gut\_Meta03. The phylogenetic tree is on the left. Distribution of the high-  
848 quality genomes assembled by other methods are colored as red in the middle heat  
849 map. N50 of each high-quality genome is shown in the right histogram.

850 **Figure S11** Phylogenetic tree of the high-quality genomes assembled by MetaTrass  
851 for P\_Gut\_Meta01. The phylogenetic tree is on the left. N50 of each high-quality  
852 genome is shown in the right histogram. Because the genome bins obtained by the  
853 combination strategies cannot be classified into species by GTDB-tk, the heat map is  
854 not showed for this sample.

855 **Figure S12** Box plot of variations. Box plots of SNVs (a), small indels (b), large indels  
856 (c) and SNV density (d) called from the high-quality genomes for four gut samples.

857 **Figure S13** Number of genomes with different quality (a) and contiguity (b)

858 assembled by MetaTrass and MetaTrass\_TR for the patient gut sample. Since  
859 MetaTrass\_TR excluded the co-barcoding refining process compared to MetaTrass,  
860 the input dataset of co-barcoding assembling in MetaTrass\_TR is the taxonomic reads  
861 set. Only the high-quality genomes are considered to count the number of genomes  
862 with different contiguity.