# High-resolution single-molecule long-fragment rRNA gene amplicon sequencing for uncultured bacterial and fungal communities

5   Chao Fang[1,2†], Xiaohuan Sun[1†], Fei Fan[1†], Xiaowei Zhang[1†], Ou Wang[1†], Haotian
6   Zheng[1,3], Zhuobing Peng[1,3], Xiaoqing Luo[1,4], Ao Chen[1], Wenwei Zhang[1], Radoje
7   Drmanac[5,6], Brock A. Peters[5,6*], Zewei Song[1*], Karsten Kristiansen[1,2*]

8   † These authors contributed equally to this work

9   * Corresponding authors.

10

11   1. BGI-Shenzhen, Shenzhen 518083, China

12   2. Laboratory of Genomics and Molecular Biomedicine, Department of Biology,
13      University of Copenhagen, 2100 Copenhagen, Denmark

14   3. BGI Education Center, University of Chinese Academy of Sciences, Shenzhen,
15      518083, China

16   4. State Key Laboratory of Biocontrol and Guangdong Provincial Key Laboratory of
17      Plant Resources, School of Life Sciences, Sun Yat-Sen University, Guangzhou,
18      510275, PR China

19   5. Advanced Genomics Technology Lab, Complete Genomics Inc., 2904 Orchard
20      Parkway, San Jose, California 95134, USA

21   6. MGI, BGI-Shenzhen, Shenzhen 518083, China

22

23   **Although several large-scale environmental microbial projects have been initiated in**
24   **the past two decades, understanding of the role of complex microbiotas is still**
25   **constrained by problems of detecting and identifying unknown microorganisms[1-6].**
26   **Currently, hypervariable regions of rRNA genes as well as internal transcribed spacer**
27   **regions are broadly used to identify bacteria and fungi within complex communities[7,**
28   **8], but taxonomic and phylogenetic resolution is hampered by insufficient sequencing**
29   **length[9-11]. Direct sequencing of full length rRNA genes is currently limited by read**
30   **length using second generation sequencing or sacrificed quality and throughput by**
31   **using single molecule sequencing. We developed a novel method to sequence and**
32   **assemble nearly full length rRNA genes using second generation sequencing.**
33   **Benchmarking was performed on mock bacterial and fungal communities as well as**
34   **two forest soil samples. The majority of rRNA gene sequences of all species in the**
35   **mock community samples were successfully recovered with identities above 99.5%**
36   **compared to the reference sequences. For soil samples we obtained exquisite**
37   **coverage with identification of a large number of putative new species, as well as high**
38   **abundance correlation between replicates. This approach provides a cost-effective**

39 **method for obtaining extensive and accurate information on complex environmental**
40 **microbial communities.**

41 Currently, there are two main approaches for sequencing rRNA genes in complex microbial
42 environmental samples: second generation technologies such as sequencing by synthesis
43 and DNA nanoball sequencing (DNBseq), and third-generation technologies such as
44 nanopore and single molecule real-time sequencing. Second-generation sequencing suffers
45 from relatively short read lengths (less than 300 bases) making it difficult to sequence the
46 entire rRNA gene. Recently, co-barcoding technologies combined with second generation
47 sequencing and *de novo* assembly have enabled nearly full-length rRNA genes, but their
48 quantification reproducibility using complex samples is unknown[9, 10]. Third generation
49 technologies have much longer read lengths enabling full coverage of the rRNA genes, but
50 they are still hampered by relatively high cost, low throughput, and high error rates[11-13]. As
51 such, there is still a great demand for a high-throughput and cost-effective approach to
52 support the deciphering of complex microbial communities. In this study, we developed a
53 strategy whereby a novel combination of rolling circle replication and DNA co-barcoding [14]
54 techniques allows sequence information of long amplicons to be obtained by a high-
55 throughput, single-molecule level *de novo* assembly approach to achieve species-resolution
56 and highly accurate profiles in an efficient and cost-effective manner using a short read
57 sequencing platform. To demonstrate the applicability of this method, we benchmarked it
58 using two mock communities of various bacterial and fungal species. We also applied it to
59 field soil samples as a demonstration for discovering unknown species and were able to
60 recover rRNA sequences even longer than those found in current reference databases.

61 The basic idea is to assemble a single DNA molecule by using DNA co-barcoding
62 technology enabled by the single-tube long fragment read (stLFR) method. Using stLFR, it is
63 possible to label DNA sub-fragments from a single long DNA molecule with the same
64 barcode. Importantly, in a single stLFR library there are approximately 3.6 billion different
65 barcodes allowing each long DNA molecule in a sample to be labeled by a unique
66 barcode[15]. One shortcoming of the current stLFR method is that sequence coverage of each
67 long molecule is less than 1X. As a result, using the De Bruijn graph (DBG) algorithm, it is
68 usually impossible to achieve full assembly of DNA fragments from a single barcode. To
69 overcome this lack of co-barcoded sequence coverage, we applied a modified version of the
70 stLFR technology, named stcLFR (single tube complete-coverage Long Fragment Read).
71 This approach involves an initial process of rolling-circle replication (RCR), which performs a
72 linear amplification process whereby tandem copies of a single DNA molecule are joined
73 head to tail in one long single stranded concatemer. This method also avoids cumulative
74 replication errors since the same DNA fragment is used as the template for replication. After
75 conversion to double stranded DNA, the RCR amplified products are disrupted at regular
76 intervals of approximately 500 base pairs by a transposon containing a capture sequence.
77 The transposon inserted molecules are then captured by barcode containing
78 oligonucleotides anchored to the surface of a micron sized magnetic bead, whereupon the
79 DNA is fragmented and ligated to these oligonucleotides (Figure 1a). After PCR on the
80 beads, the DNA fragments with unique barcodes, corresponding to multiple copies of each
81 single molecule, are subjected to DNB based sequencing. By decoding barcoded sequences
82 from the sequencing dataset, copy depth can be estimated in each barcoded binning group.
83 Bins with depth > 5X are then used for parallel *de novo* assembly (Figure 1b). Successfully

84  assembled contigs are termed pre-binning assemblies (PBAs) and used for subsequent
85  analysis.

86  The ZymoBIOMICS™ mock community which contains 8 bacterial species (See Methods)
87  was used to test the ability of stcLFR to properly assemble the rDNA amplicons. A region
88  covering 4.5 kb from the SSU (515Fng forward primer) to the LSU (TW13 reverse primer)
89  was amplified for sequencing. From 186,580,349 clean reads generated from 3 replicates,
90  we decoded 26,512,829 valid barcoded read bins of which 148,814 having a coverage > 5X
91  were used for *de novo* assembly. With a relatively small number of reads and simple
92  genomic content, the assembly proceeded faster and was more complete compared to
93  general metagenomic assembly. We finally retrieved 146,509 assembled PBAs with a
94  centroid length distribution close to 800 bp (Supplementary figure 1a). Most of them covered
95  either the SSU (34.15%) or the LSU (38.58%). Only 25.23 % of PBAs covered both the SSU
96  and the LSU (Supplementary figure 1b). The assembly quality was high, with 9.72% of the
97  aligned PBAs exhibiting 100% identity to the reference genomes, and more than 60% of the
98  PBAs achieved 99.5% or higher identity with the reference genomes. By contrast, only
99  1.78% of the PBAs could not align to the reference genomes (Figure 2a).

100  For tests on fungi, a mock community comprising 7 common fungal species was used to
101  assess the classification resolution at the species level (See Methods). Primers spanning a
102  region ~ 2.5 kb covering ITS1 and ITS2, and parts of the flanking SSU and LSU rRNA genes
103  were used. Following the DNA library preparation and sequencing protocols used for
104  bacterial species, we generated 190,208,124 high-quality reads, decoded 21,215,624
105  barcoded bins, with 270,794 having a coverage > 5X. In total, 263,416 PBAs were
106  generated with a size centered on 2.3 kb, which is close to the value expected from the
107  design (Supplementary figure 1c). Exceeding the results obtained for the bacterial samples,
108  12.54% of the fungal PBAs exhibited a 100% alignment to the reference genomes and more
109  than 72.4% of the PBAs achieved 99.5% or higher identity to the reference genomes. Only
110  0.43% of the PBAs were unable to be aligned to the reference genomes (Figure 2b). 76.8%
111  of the PBAs covered the entire region from the SSU to the LSU, and the ITS region was
112  included in 83.7% of the PBAs (Supplementary figure 1d).

113  For assessing the accuracy of quantification, we mapped all decoded reads to the
114  assembled complete rRNA gene sequences. Mapped reads sharing the same barcode
115  represent a single DNA molecule and as such were counted once. We determined the
116  difference between the observed and the calculated relative abundance of each species to
117  estimate the error. For analyses of the mock bacterial community, we calculated the error to
118  range from -9.3% to 15.6%, with an error standard deviation of 7.9% (Figure 2c). For
119  analyses of the fungal mock community, we estimated the error to range from -9.7% to
120  8.2%, with an error standard deviation of 5.6% (Figure 2d). We estimated that the variation
121  in the relative abundance of each species, as defined by the variance of error, was as low as
122  0.32% for bacterial species and 0.15% for fungi (Figure 2c). Though we only obtained a low
123  percentage of amplicons covering the entire bacterial rRNA gene region as designed, we
124  successfully recovered long fragments of bacterial 16S and 23S rRNA genes, and fungal
125  fragments covering the entire ITS region and parts of the flanking 18S and 28S rRNA genes.

126  From both bacterial and fungal mock samples, the identity between PBAs and the
127  corresponding reference sequences exceeded in most cases 99%, a percentage higher than
128  97% which often is used for amplicon-based species identification. In addition, the observed

129  low variation in the relative abundance of each species also demonstrated the consistency
130  of the protocol, an important requirement for comparative analyses using high-resolution
131  profiles.

132  We had limited success in retrieving the entire ~4 kb length of bacterial rDNA regions and
133  observed a gap in the coverage of the bacterial rDNA genes, but still successfully recovered
134  long fragments of bacterial 16S and 23S rRNA genes. The gap harbors transfer RNA (tRNA)
135  genes (Supplementary figure 2), which may be a target for cleavage by the Tn5
136  transposase, the enzyme used for fragmentation during sequencing library preparation[16-19].
137  Although cleavage by the Tn5 transposase exhibits limited sequence bias, target
138  preferences might still exist and cause the failure of effectively retrieving fragments covering
139  the entire rDNA gene region[22-24]. In eukaryotes, tRNA genes are generally located outside
140  the SSU and LSU regions, not in-between, which may at least in part explain why the
141  eukaryotic amplicons survived, but the bacterial amplicons were split.

142  From the reference alignment benchmarking, we observed that when a PBA exceeded its
143  designed size, part of the sequences could be aligned to reference sequences belonging to
144  different species. We consider such PBAs as chimeric PBAs, which were subsequently
145  filtered out (Supplementary figure 3b, red lines). Since we successfully enriched for the
146  designed size of fungal rDNA sequences, this trend was prominent when the length of PBAs
147  exceeded 2.3 kb for fungal samples. In addition, for bacterial PBAs with length < 500 bp we
148  also observed a higher probability for achieving the exact same identity and bit score by
149  multiple reference sequences representing different species, leading to an ambiguous
150  annotation (Supplementary figure 3a, solid blue line). This trend became more severe in
151  fungal PBAs with size < 2000 bp (Supplementary figure 3b, the solid blue line). These
152  findings might point to one limitation of using short sequences to distinguish taxonomies at
153  the species level, especially for fungi. However, coverage of multiple regions greatly
154  eliminated the ambiguous annotation (Supplementary figure 4c).

155  To further explore the abundance and reproducibility of taxonomies in real environmental
156  samples, we sequenced 3 replicates of two natural soil samples to identify bacterial rRNA
157  genes and fungal ITS regions using the same primer sets and protocols used for the mock
158  samples. We obtained 632,574,076 bacterial reads and 1,006,826,680 fungal reads with
159  valid barcodes. For bacterial amplicons, we successfully generated 544,433 PBAs from
160  724,038 candidate bins (Figure 2e). As we observed using mock communities, only a small
161  fraction (5.86%) of the bacterial PBAs covered both the SSU and the LSU regions. Most
162  bacterial PBAs only covered the SSU (49.79%) or the LSU (30.97), whereas 13.38% did not
163  match any currently known rDNA region (Figure 2f). For fungal amplicons, 1,807,779 PBAs
164  from 1,845,429 high-coverage bins were generated with a size distribution that peaked
165  around 2.3kb and 1.8kb (Figure 2g). A PBA size of 2.3kb was consistent with the primer
166  design showing full coverage of the flanking region of the SSU region, the ITSs, and the
167  flanking region of the LSU region. By contrast, half of the PBAs with sizes that peaked
168  around 1.8kb lacked the LSU region. We detected ITS regions from 65.05% of the PBAs
169  (Figure 2h). This percentage increased when the size exceeded 1 kb.

170  We also tried to profile the bacterial and fungal community at different rank levels. As soil
171  samples contain very complex microbial communities, we chose Kraken2 to handle the
172  profiling task. To build a Kraken2 database we first used PBAs to generate cluster trees for
173  bacteria and eukaryotes (mostly fungi), respectively. Since we retrieved individual rDNA

174    subunit assemblies for bacteria, we selected PBAs harboring SSU sequences for clustering.
175    For eukaryotes, and especially fungi, the ITS region provided the highest level of
176    discrimination, and accordingly, we selected PBAs covering ITS1 and ITS2 for clustering.
177    The operational taxonomic unit (OTU) clusters were constructed using a set of identity
178    thresholds enabling classification of taxa at multiple taxonomy levels, from domain to
179    species. These thresholds were determined by pairwise global alignments of ribosomal gene
180    subunit sequences including SSU and LSU from the SILVA database and ITS from the
181    UNITE database (Supplementary figure 4), setting 97% identity as the threshold for genus
182    association and 99% for species association for bacteria, and 95% identity and 97% identity
183    as the thresholds for genus and species association, respectively, for fungi. Additional
184    thresholds were the same for bacteria and fungi (See Methods). The trees generated using
185    PBAs were then merged with taxonomy trees based on the SILVA and the UNITE databases
186    (Figure 3a). In this process, clades were merged with taxonomies when the representing
187    PBAs achieved sequence identities greater than the rank's threshold. The combined tree
188    contained three categories of branches. One type was represented by the public SILVA and
189    UNITE sequences (PUBs, green in Figure 3, defined as exclusively PUBs), one represented
190    by sequences only present in PBAs (blue in Figure 3, exclusively PBAs), and one
191    represented by both the PBA and the PUB sequences (red in Figure 3, the shared PBAs and
192    PUBs). The merged taxonomy tree and combined sequences were then used to build the
193    database.

194    By this approach, we mapped all barcoded reads to the combined database retrieving
195    60,942 bacterial OTUs at the species level of which 1,078 OTUs were supported by both the
196    PBA sequences and PUB sequences, 18,403 could not be annotated, representing putative
197    novel species, and the remaining 41,461 were annotated by publicly available sequences. At
198    the genus level, 6,765 OTUs represented PBAs with no annotation. Using 72% identity as
199    the threshold for association with a phylum, we identified 241 OTUs of which 19 were novel.
200    For the eukaryotic communities, we retrieved 37,355 fungal OTUs at the species level of
201    which 1592 OTUs were supported by both the PBA sequences and PUB sequences, 1,817
202    could not be annotated, representing putative novel species, and the remaining 33,946
203    OTUs were annotated by PUB sequences. At the genus level, 1,266 OTUs represented
204    PBAs with no annotation. Using 72% identity as the threshold for association with a phylum,
205    we identified 244 OTUs, all of which were annotated. (data of sample S are shown, Figure
206    3b, left panel).

207    We found that species level OTUs supported by both PBA and PUB sequences represented
208    only 1.8% of all OTUs; the relative abundance of these OTUs combined corresponded to
209    about 12% of the relative abundances of all OTUs. At the genus level, the relative
210    abundance of OTUs supported by both the PBA and PUB sequences was even higher
211    (62%), and at the phylum level the relative abundance of the OTUs supported by both PBA
212    and PUB sequences reached close to 100%. This trend was more obvious for fungal OTUs,
213    where 4.2% OTUs supported by both the PBAs and PUBs contributed 79% of relative
214    abundance at the species level (figure 3b, right panel). This indicated that OTUs supported
215    by PBAs (whether annotated or novel) were abundant and dominated the microbial
216    community. Compared with the fungal community, more dominant taxonomies may
217    represent novel species.

218    Further analysis showed that the three categories of branches in the combined tree
219    exhibited distinct differences in terms of reproducibility between replicates. OTUs supported

220    by both PBA and PUB sequences exhibited high correlation between replicates (r=0.93 for
221    bacterial communities and r=0.88 for eukaryotic communities; Spearman's correlation, figure
222    3c). Higher correlations were observed at the genus and the phylum levels. By contrast, the
223    correlation between replicates of OTUs supported exclusively by PBA was lower and the
224    correlation decreased significantly when the abundance became lower than $10^{-4}$, both for
225    bacteria and eukaryotes. Some of the unannotated bacterial phyla taxa still varied, indicating
226    that they might not represent an independent phylum. For the taxa supported by public
227    references, though well organized by taxonomic information, the species taxa performed
228    worse than those supported by unannotated PBAs. However, the performance was
229    improved dramatically at the phylum levels. These dominant annotated OTUs supported
230    both by PBAs and PUBs and highly abundant ($>10^{-4}$) novel OTUs exhibited high
231    reproducibility for quantification comparison tasks at the species and the genus level.

232    The rarefaction curves also showed a better performance in terms of number of reads
233    needed for sufficient coverage of OTUs from the shared PBAs and PUBs sequences. To
234    retrieve the majority of taxa (here defined as 95% of maximum OTUs), the publicly available
235    database required the highest number of read-pairs, 410 million read-pairs for bacteria and
236    425 million read-pairs for fungi. To assess exclusively PBA OTUs, this requirement was
237    reduced to 110 million and 150 million for bacteria and eukaryotes, respectively. The shared
238    PBAs and PUBs representing OTUs required the lowest number of reads, only 40 million
239    bacterial read-pairs and 55 million eukaryotic read-pairs, respectively.

240    An alternative approach for long-read amplicon sequencing was recently published by Karst
241    and co-workers[20]. In this work, re-designed UMIs and single-molecular sequencing were
242    used for obtaining long-read high-accuracy amplicon sequences using Nanopore or PacBio
243    sequencing. While this is an elegant approach, our approach relying on second-generation
244    sequencing enables a more cost-effective very high sequencing throughput coupled with
245    robust and reproducible quantification.  In summary, our approach provided a 99% identity of
246    species level, high-throughput strategy to expand current rRNA gene databases by including
247    long marker sequences and potential novel taxonomies, as well as a comparable accurate
248    quantification profiling strategy. This approach provides a cost-effective method for obtaining
249    extensive and accurate information on environmental complex microbial communities.

# Methods

250

### Sample collection and DNA extraction

251
252    The mock bacterial DNA standard, ZymoBIOMICS Microbial Community Standard (D6305),
253    was purchased from ZymoBIOMICS™. Fungi strains were purchased from China General
254    Microbiological Culture Collection Center. The fungal mock DNA sample consisted of equally
255    mixed genomic DNA of 7 fungi, including *Aspergillus ustus* (No. BNCC144426), *Trichoderma*
256    *koningii* (No. BNCC144774), *Penicillium expansum* (No. BNCC146144), *Aspergillus nidulans*
257    (No. BNCC336164), *Penicillium chrysogenum* (No. BNCC336234), *Trichoderma reesei* (No.
258    BNCC341839) and *Trichoderma longibrachiatum* (No. BNCC336352). The sequence of the
259    rDNA region for each strain was rechecked by Sanger sequencing. Soil samples were
260    collected from Nanbanhe tropical rainforest plot (21.612 N, 101.574 E) in Yunnan, China in
261    2017. Soil samples were immediately stored at -80℃ after collection, and DNA extraction
262    was performed within two months. DNA extraction of all samples was performed using the
263    PowerSoil®DNA Isolation Kit (Mobio) according to manufacturer's instructions. DNA
264    concentration was measured by Qubit flex fluorometer (Invitrogen).

## rDNA long fragments amplification

265     We initially amplified the region of bacterial 16S-23S rRNA gene and fungi full length ITS
266 region with Kapa Hifi DNA polymerase (Roche) and EX Taq DNA polymerase (Takara Bio).
267 No correct PCR products (examined by Sanger sequencing) were amplified by Kapa,
268 probably due to its high fidelity and the fact that the conserved regions where primers locate
269 can result in mismatches of several nucleotides. By contrast, EX Taq DNA polymerase
270 (Takara Bio), which also has proofreading activity, successfully amplified the rDNA targets, it
271 was therefore used for full length rDNA amplification. For PCR amplification we used 5'
272 phosphorylated primers. The primer sequences are listed below. For bacteria, the forward
273 primer was 27F (AGAGTTTGATCATGGCTCAG) and the reverse primer was our in-house
274 designed 23S-2850R (CTTAGATGCCTTCAGCRVTTATC). For fungi, highly universal
275 eukaryotic primers for high eukaryotic and fungal taxonomic coverage were selected based
276 on a previous study (Tedersoo et al., 2018). The forward primer was SSU515Fngs
277 (GCCAGCAACCGCGGTAA) and the reverse primer was TW13
278 (GGTCCGTGTTTCAAGACG). The conditions for PCR were as follows: 1 µL of diluted
279 template DNA, 1 µL of forward primer (10 M), 1 µL of reverse primer (10 M), 15.5 µL of
280 nuclease-free water, 5 µL of 10X Ex Taq Buffer, 1 µL of dNTPs Mixture (2.5 mM each), and
281 0.5 µL of EX Taq Polymerase. We amplified samples using the following cycling conditions:
282 95 °C for 5 min; 30 cycles of 95°C for 30s, 55°C for 30s, and 72°C for 3 min; and then a final
283 extension at 72 °C for 10 min. The amplified long amplicons were purified using QIAquick
284 PCR Purification Kit (Qiagen).

## Circularization and Rolling Circle Replication

286     To enable a highly efficient circularization of long molecules, a double stranded
287 circularization method was applied. 100 ng amplicon were first subjected to end-repair and
288 A-tailing according to the protocol described previously[21]. Then the product was incubated
289 with 8 pmol of adapter and 3000 units T4 DNA ligase (MGI, 1000004279) in 80 µL of 1X
290 PNK buffer (NEB, B0201S) with extra 1 mM ATP and 7.5% PEG-8000 at room temperature
291 for 1 hour, followed by a 0.5X SPRI beads purification (Beckman, A63882) and eluted with
292 20 µL of TE buffer (Thermo Fisher, AM9849). Next polymerase extension with 1 pmol of
293 primer containing uracil was carried out, by adding prior heat-activated 200 units Pfu Turbo
294 Cx (Agilent Technologies, Inc., 600414) in 50 µL of 1X PfuCx buffer at 72℃ for 20 minutes,
295 followed by a 1.5X SPRI beads purification and elution with 30 µL of TE buffer. Sticky ends
296 were created by injecting 20 units of USER enzyme (NEB, M5505S) and incubating with 50
297 µL of 1X TA buffer (Teknova, T0380) at 37℃ for 1 hour. T4 DNA ligase-mediated
298 circularization was performed in 150 µL of 1X TA buffer with extra 1 mM ATP. All linearized
299 DNA were removed with 0.4 units Plasmid-Safe™ DNase (Lucigen, E3110K) followed by a
300 1X SPRI beads purification.

302     The double stranded circle from the last step was designed with a 3nt gap on one strand
303 acting as the initial extending site for rolling circle replication. The RCR reaction was carried
304 out by incubation with 5 units Phi29 (MGI, 1000007887) in 21 µL of 1X Phi buffer at 30℃ for
305 1 hour. In this step the original long amplicons were transformed into long concatemers with
306 multiple copies of each long molecule. Next, unlike other studies[22, 23], we applied 60 units
307 warm-started Bst 2.0 polymerase (NEB, M0538M) and primer extension with specific
308 sequences (CGCTGATAAGGTCGCCATGCCTCTCAGTAC) to generate the second strand
309 of the RCR product ready for standard stLFR library preparation.

MANUSCRIPT                    7

### stLFR barcoded DNA library construction

310    Extended amplicons produced after RCR were labeled by unique barcodes with MGIEasy
312    stLFR Library Prep Kit (MGI, 1000005622). Briefly, indexed transposons are inserted into
313    1ng of double strand rolling circle replication products from different samples, followed by
314    hybridization of the transposon integrated DNA onto clonally barcoded beads. After capture,
315    the sub-fragments of each transposon inserted DNA molecule are ligated to the barcode
316    oligo. A few additional library processing steps are performed followed by PCR. At this point
317    the co-barcoded sub-fragments are ready for sequencing. This method and a detailed
318    protocol were previously described by Wang et. al.[24] and Cheng et. al.[25], respectively.

### Sequencing and decoding long fragments

320    DNA libraries were sequenced using the pair-end 100 bp mode on the BGISEQ-500
321    platform[26]. During sequencing, the barcode part was sequenced first and attached to the tail
322    of read2. The barcode detection and sequence read quality control were managed by the
323    fastp software[27].This tool was initially designed for traditional NGS data QC, with a parallel
324    function to speed up the process. We added barcode detection module to it. In the detection
325    process, each 10-base barcode sequence (3 10-base sequences make up a full barcode) is
326    scanned against our available barcode list, both by forward and reverse strand, within 1
327    base mismatch tolerance. In addition, a module to perform the BGISEQ platform specific
328    quality filtering and trimming was also implemented as described previously[26].  Reads were
329    sorted by barcode and stored in a fastq format. More than 85% of reads were associated
330    with a valid barcode. Less than 1% of detected barcodes were recovered after barcode error
331    correction (data not shown).

### Single-molecular pre-binning and *de novo* assembly strategy

333    An estimation of kmer coverage was initially performed to ensure that each bead had
334    enough coverage for an independent de novo assembly. The kmer coverage was
335    determined by the formula:

$$cov_{kmer} = \frac{n_{reads} \times L_{read} - k + 1}{n_{kmer}}$$

337    Where $n_{reads}$ means reads number belonging to a bead; $L_{read}$ means the length of each
338    read, which is 100bp here; $k$ means the length of kmer, which is 31; $n_{kmer}$ means the unique
339    number of kmers from $n_{reads}$ , here calculated by mash version 2.1.1[28].

340    According to the coverage distribution, we only allowed beads with $cov_{kmer} \geq 5$ for the
341    assembly process. Each bead's assembly was performed by megahit version 1.1.2 [29], with
342    parameters *--k-min 21 --k-step 20 --prune-level 0 --min-count 1*.

343    Since amplicons were head to tail adjoined by the RCR adaptor, it is possibile to assemble a
344    circular sequence with the RCR adaptor. In this case, a script to clip the RCR adaptor from
345    contigs and linearize the sequence was used.

### Taxonomic annotation for PBAs

347    For taxonomic assignment, PBAs were initially aligned to SILVA version 138 SSU and LSU
348    refseq[30] by blastn[31], separately. For fungal PBAs, UNITE[32, 33] (released at Apr 2, 2020) was
349    used as reference of ITS region. Alignment results were then summarized to provide the
350    most possible taxonomy for each PBA according to the summary of alignments of the SSU,

351   ITS and LSU regions, measured by sequence identity, bit score, and length coverage. PBAs
352   with ambiguous annotation pointing to multiple different taxonomies with the same score
353   were eliminated. Chimeric PBAs with several fragments uniquely assigned to different
354   taxonomies were also discarded.

### Default taxonomic rank threshold determination

356   The pairwise global alignments were performed for the SSU, LSU and ITS region
357   sequences, separately. The SSU sequences were collected from SILVA version 138.1 SSU
358   refseq[30]. The LSU sequences were collected from SILVA version 138.1 LSU refseq[30]. The
359   ITS sequences were collected from UNITE[32, 33] (released on Apr 2, 2020). For each region,
360   the alignments were initiated from its top rank (species or subspecies). For each taxonomy,
361   if multiple associated sequences existed, a pairwise global alignment was performed by
362   vsearch v2.14.1[34] to calculate sequences identities between each two sequences, with the
363   following command:

364   vsearch --allpairs_global <belonging sequence file> --acceptall --uc -

365   To minimize computing, no more than 100 sequences from a given taxa were used. After
366   completion of assignment to all ranks, the sequence identity distribution of the median value
367   of each taxon was used for visualization. For each rank, the peak position was selected
368   manually as the estimated threshold for this rank.

### Cluster tree generation

370   For each rank, taxonomy assignment of a given taxon was determined and organized for
371   Kraken version 2.0.8-beta[35] database generation.

372   Bacterial PBAs ranging from 500 bp to 1500 bp and fungal PBAs ranged from 1700 bp to
373   2500bp were picked

374   PBAs used to process clusters satisfied the following criteria: 1) target rDNA region(s)
375   detected; 2) size of length in defined range; 3) no ambiguous annotation; 4) no chimera
376   detected. For bacterial PBAs, barrnap version 0.9 (https://github.com/tseemann/barrnap)
377   was used for rDNA regions detection, and the picked length ranged from 500 bp to 1500bp.
378   For fungal PBAs, barrnap version 0.9 and ITSx version 1.0.11[36] were both used for rDNA
379   regions detection, and picked length ranged from 1700 bp to 2500 bp. Any PBA with partial
380   segments aligned with conflicting annotations was discarded. Retained PBAs were then
381   clustered by vsearch with a series of identities at 0.4, 0.5, 0.66, 0.72 (for phylum), 0.77, 0.83,
382   0.89, 0.92, 0.95, 0.97(for genus), 0.98, 0.99 (for species), 0.995, 0.999 and 1, according to
383   the taxonomies similarity centroids calculated from all annotated sequences from the SILVA
384   and the UNITE databases (See Supplemental table 3). Each cluster with a certain identity
385   was regard as a clade of a taxonomy tree if the member PBAs could not provide a unified
386   rank annotation. Note that the cutoff of rank above species varied greatly among different
387   groups. The cutoff values were only assigned for clades without any useful information,
388   otherwise they were determined by the clade members' classifications. Finally, at species
389   rank, singleton OTUs without annotation were discarded in order to reduce the false
390   discovery rate.

### Relative abundance computation of each rank

392   The rank abundance was calculated by Kraken2, which uses kmer alignments to determine
393   the position of each read. We then rescaled it to barcode unit whereby reads sharing the

394      same barcode represent a single DNA molecule. The generated results are compatible with
395      Kraken2 so each rank's profile can be classified directly.

396      **Phylogenetic tree building**
397      For bacterial SSU rRNA gene sequences, barrnap version 0.9
398      (https://github.com/tseemann/barrnap) was used to secure quality of the PBAs sequences
399      with no chimeras. Fungal ITS regions were predicted by ITSx version 1.0.11[36] which also
400      enabled determination of the boundary of the SSU and the LSU. The SSU and LSU
401      sequence were then joined together for heterogeneity rate computation. Both bacterial SSU
402      and fungal joint SSU+LSU sequences were multiple aligned by mafft v7.407[37] and terminal
403      gaps trimmed by trimAl v1.4.rev22[38]. RAxML version 8.2.1[39] was then employed to perform
404      100 boostraps General Time Reversible model of nucleotide substitution under the Gamma
405      model of rate heterogeneity, with accommodated searches incorporated (GTRCAT). The
406      processes were executed on a 40-core node of the Danish National Supercomputer for Life
407      Sciences (Computerome 2.0) with following parameters:

408      raxmlHPC-HYBRID-AVX -f a -p 12345 -x 12345 -T 40 -# 100 -m GTRCAT -s
409      <input.aligned.trimed.fasta>

410      The results of the best-scoring ML tree with support values were imported and visualized by
411      ggtree package[40] in Rstudio v1.3.1073[41] IDE for R 4.0.3[42].

412      **Data availability**
413      The data that support the findings of this study have been deposited into CNGB Sequence
414      Archive (CNSA)[43] of China National GeneBank DataBase (CNGBdb)[44] with accession
415      number CNP0001509. All data processing scripts are integrated in a toolbox pipeline
416      available on github (https://github.com/Scelta/metaSeq).

417 # Acknowledgments

421 # Contributions
422      F.F, O.W. and B.P. conceived the wet lab method. X.S., F.F. and O.W performed wet lab
423      experiments and BGISEQ500 sequencing. C.F. and Z.S. conceived the bioinformatics
424      method. C.F developed the software pipeline and performed data analysis as well as
425      visualization. X.S. and X.Z. interpreted the microbial characteristics of data. X.S. performed
426      the phylogenetic analysis. H.Z., Z.P. and X.L. assisted to perform the data. C.F., X.S., X.Z.,
427      O.W., B.P., Z.S. and K.K. wrote the manuscript. All authors participated in discussions and
428      contributed to the revision of the manuscript. All authors read and approved the final
429      manuscript.

430 # Competing interests
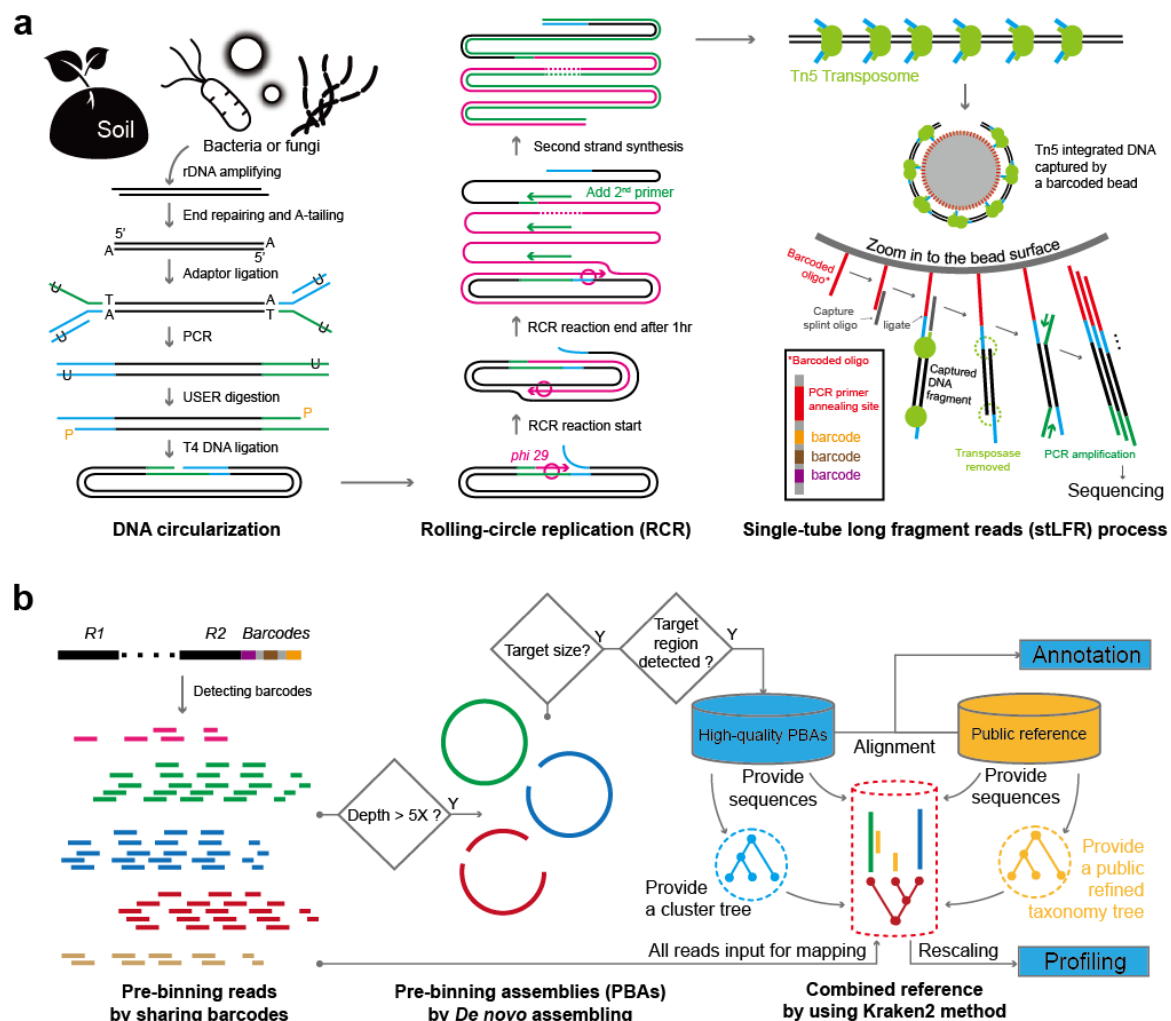431      All authors declare no competing financial interests.

# Figures



**Figure 1 Sequencing framework (a).** *In vitro* processes of rolling-circle replication and stLFR. DNA from bacteria and fungi was extracted and amplified using designed rDNA primers. Amplicons were turned into long concatemers with circularization and rolling-circle replication (RCR). The products of RCR were labeled with unique barcodes by integrating transposons and hybridizing onto 30 million different clonal barcoded beads in a single tube. After PCR, these sheared short sub-fragments were subjected to DNA sequencing. **(b).** *In silico* processes of assembly, classification, and quantification. Reads with attached barcodes were decoded during the quality-control process. Pre-binning was done by grouping reads sharing the same barcode. Subsequently, *de novo* assembly was performed for each bin independently and in parallel. The pre-binning assemblies (PBAs) were then selected by target size range and were detected rDNA subunits or ITS regions by Barnnap and ITSx softweare were used for. The open reference databases SILVA and UNITE were used for taxonomic classification. Kraken2 was used to generate a database for profiling by mapping all barcode-detected reads.
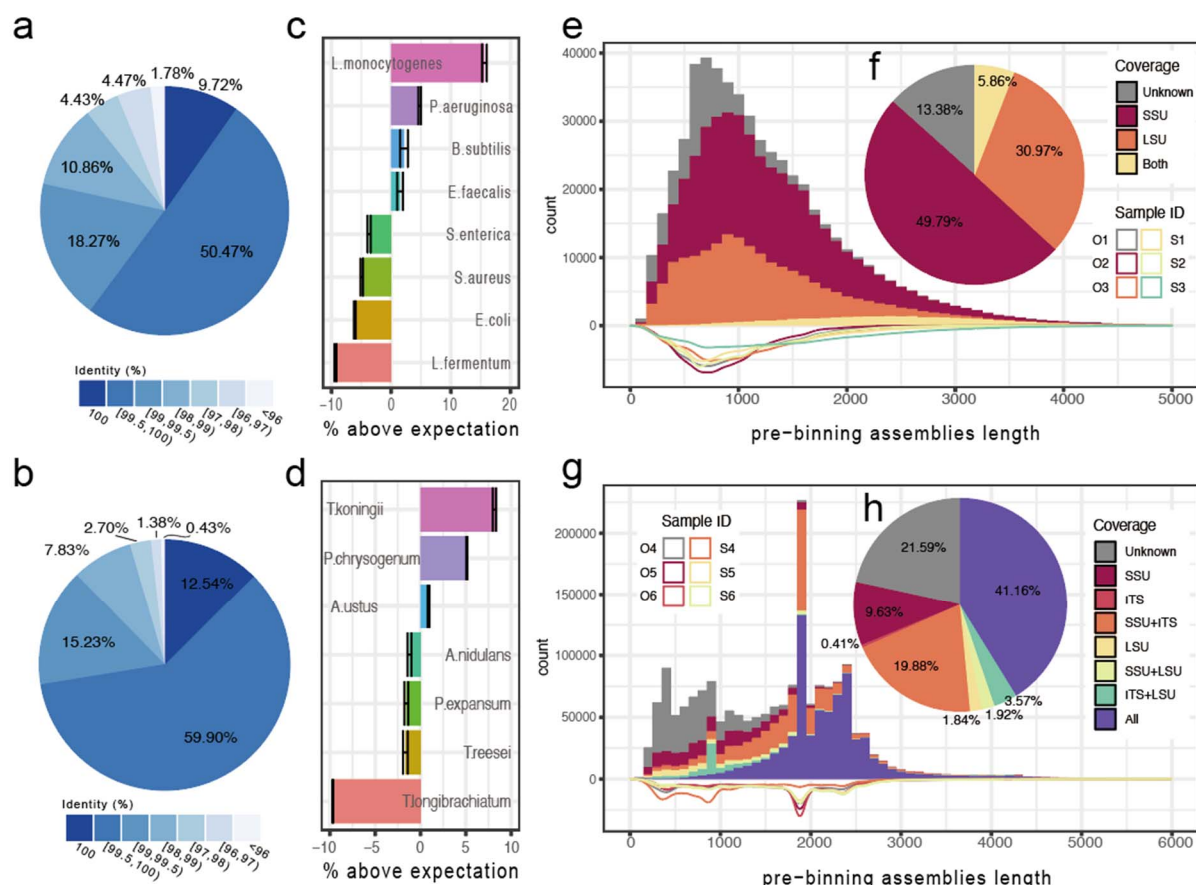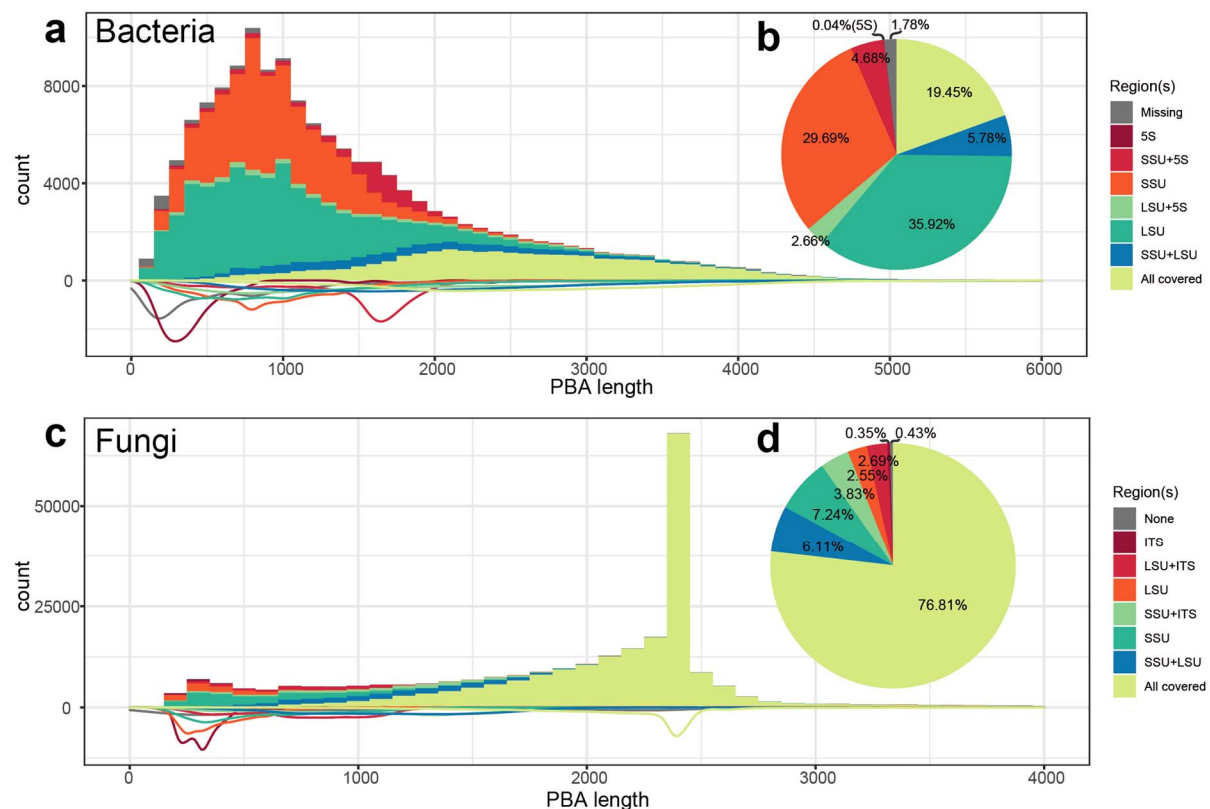
**Figure 2 Long fragments restore performance**. (a) Identity of alignments of mock bacterial rDNA PBAs to reference sequences. (b) Identity of alignments of mock fungal rDNA PBAs to reference sequences. (c) Difference between the observed and the calculated abundances of bacterial species in the mock community. (d) Difference between the observed and the calculated abundances fungal species in the mock community. (e) rDNA sequences assembled from soil bacterial DNA. Subunits (16S/SSU and 23S/LSU) were detected by alignment to the SILVA and UNITE database or predicting by barrnap. The length distribution of each sample is plotted on the negative part of the y-axis. Insert, pie plot of percentages of subunits detected. (f) rDNA sequences assembled from soil fungal DNA. Percentage of covered 18S/SSU, 28S/LSU and ITS is indicated by different colors. Insert, pie plot of percentages of subunits and ITS region detected.

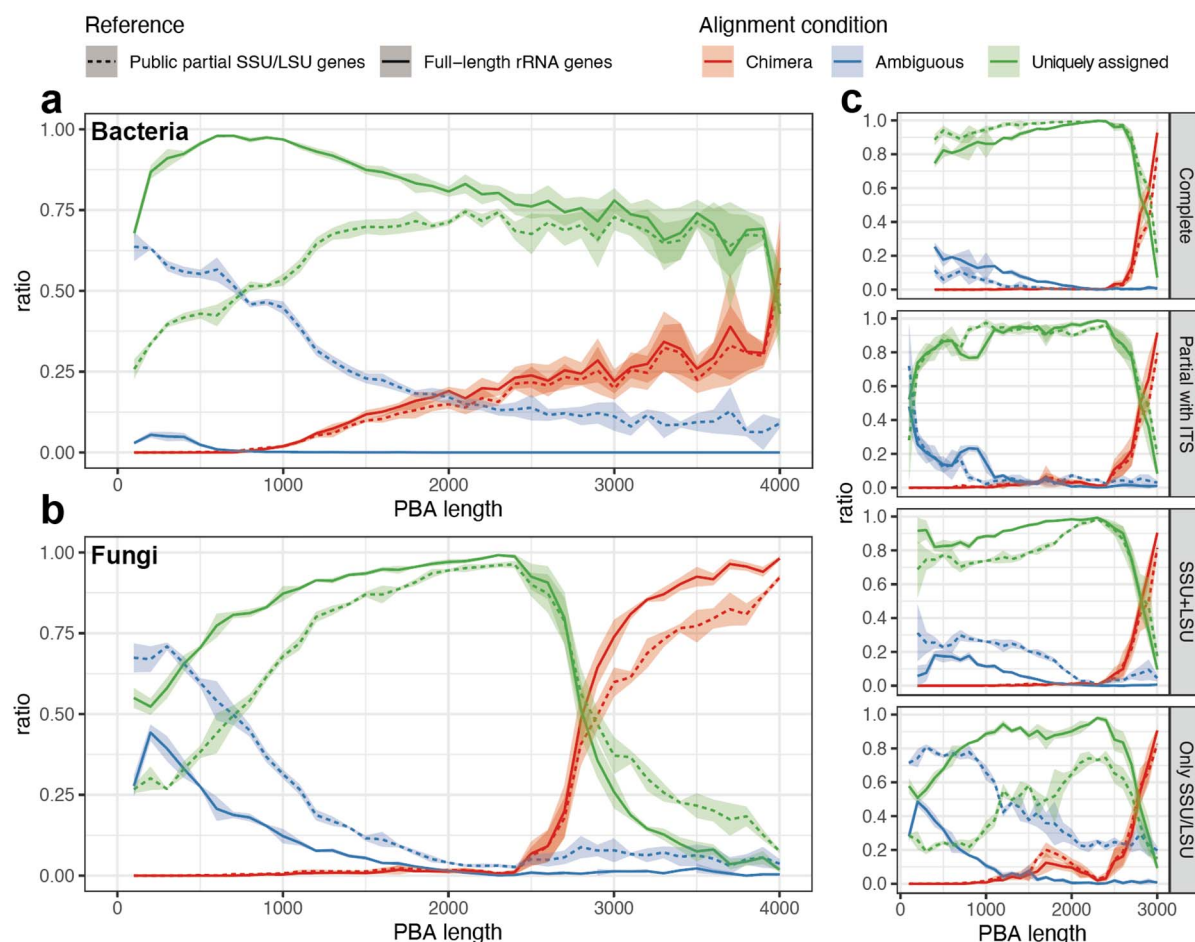**Figure 3 Identification and performance of a combined Kraken2 profiling strategy.** (a) The pre binning assembly (PBA) cluster tree was generated by a series of clustering of PBAs at default identity cut-off for each level of taxonomy. The identity cut-offs were estimated by pairwise alignments of annotated SSU, LSU and ITS sequences from public databases. The public taxonomy tree was mainly based on the SILVA SSU taxonomy tree and supplemented with new taxonomies from the SILVA LSU and UNITE databases. A combination of PBAs and public reference sequences (PUBs) was established based on the identities between the PBA and PUB sequences. The combined taxonomy tree contains three branch categories of sequences: 1) Sequences shared by PBAs and the PUBs (red); 2) sequences unique to PBAs with no taxonomy association (blue); 3) Reference sequences from public taxonomy tree not shared by PBA (green). (b) The performance of the three branches of the taxonomic trees. Left: The number of OTUs shown at the species, genus and phylum levels, visualized as bar plots. Right, relative abundances. The colors indicate OTUs belonging to PUBs (green) or PBAs with (red) or without (blue) annotation. (c) Correlation of relative abundance between duplicates. At each level of taxonomy, the relative abundance observed in each duplicate is colored light red for bacteria and light green for fungi. The Spearman correlation values for bacteria and fungi are indicated (d) Rarefaction curves of observed bacterial and fungal OTUs at the species level using the different databases. One bacterial sample (solid line) and one fungal sample (dash line) with 3 replicates and sequenced to more than 500 million read pairs were analyzed. On each curve, a cross marks the number of read needed to cover 95% of the total counts.
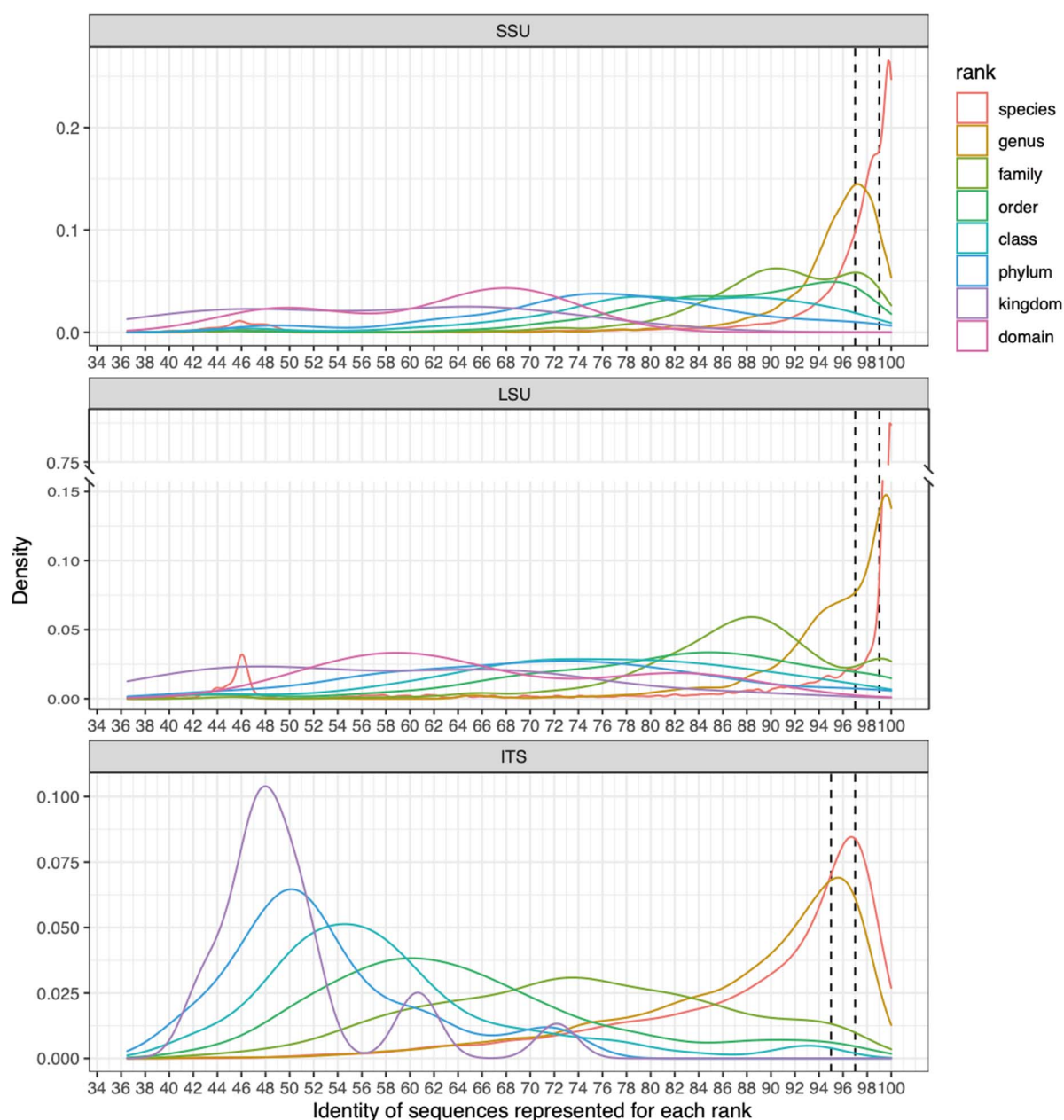
# Supplementary Material



**Supplementary figure 1 Length and rDNA subunits coverage of bacterial and fungal PBAs**. (a) Length distribution of bacterial mock PBAs. Percentages of SSU, LSU and 5S coverage are shown in colors. Density lines of subunit coverage are shown on the negative side of the y-axis. (b) Pie plot of percentage of subunits detected. (c) Length distribution of fungal mock PBAs. Percentages of SSU, LSU and ITS coverage are shown in colors. Density lines of subunit coverage conditions shown on the negative side of the y-axis. (d) Pie plot of percentages of subunits and ITS region detected.

**Supplementary figure 2 PBA coverage of 8 bacterial mock species' rRNA gene sequences.** All PBAs were mapped to their corresponding reference sequences. 16S, ITS and 23S regions were detected by barrnap and colored. Several positions where the coverage immediately dropped were observed in some species (red arrows).

**Supplementary figure 3 Taxonomic classification accuracy estimation.** Fraction of ambiguous (blue) and chimera PBAs (red) were summarized in each 100 bp window of PBA length in bacterial (a) and fungi (b) samples, respectively. Mock specific complete reference sequences (solid line) and public references from PUB (dashed line) were used for comparison. (c) Fungal PBAs were divided according to covered regions for estimation of the ratio of ambiguous and chimeric sequences relative to the length of the PBAs. .

**Supplementary figure 4 The sequence identity distribution of each rank based on rRNA gene sub-region sequences.** The median value of sequence identities between each two members of a given taxonomy was used. At each rank, the curve peak position was selected manually as the estimated threshold to determine a default rank of an unknown clade.

# Reference

1.  Vogel, T.M. et al. TerraGenome: a consortium for the sequencing of a soil metagenome. *Nature Reviews Microbiology* **7**, 252-252 (2009).
2.  Ehrlich, S.D. in Metagenomics of the Human Body. (ed. K.E. Nelson) 307-316 (Springer New York, New York, NY; 2011).
3.  Gevers, D. et al. The Human Microbiome Project: a community resource for the healthy human microbiome. *PLoS Biol* **10**, e1001377 (2012).
4.  Gilbert, J.A., Jansson, J.K. & Knight, R. The Earth Microbiome project: successes and aspirations. *BMC biology* **12**, 69 (2014).
5.  Sunagawa, S. et al. Structure and function of the global ocean microbiome. *Science* **348**, 1261359 (2015).
6.  Bahram, M. et al. Structure and function of the global topsoil microbiome. *Nature* **560**, 233-237 (2018).
7.  Hamady, M., Lozupone, C. & Knight, R. Fast UniFrac: facilitating high-throughput phylogenetic analyses of microbial communities including analysis of pyrosequencing and PhyloChip data. *The ISME Journal* **4**, 17-27 (2010).
8.  Schoch, C.L. et al. Nuclear ribosomal internal transcribed spacer (ITS) region as a universal DNA barcode marker for <em>Fungi</em>. *Proceedings of the National Academy of Sciences* **109**, 6241-6246 (2012).
9.  Bishara, A. et al. High-quality genome sequences of uncultured microbes by assembly of read clouds. *Nature biotechnology* (2018).
10. Karst, S.M. et al. Retrieval of a million high-quality, full-length microbial 16S and 18S rRNA gene sequences without primer bias. *Nature biotechnology* **36**, 190-195 (2018).
11. Deamer, D., Akeson, M. & Branton, D. Three decades of nanopore sequencing. *Nature biotechnology* **34**, 518-524 (2016).
12. Wagner, J. et al. Evaluation of PacBio sequencing for full-length bacterial 16S rRNA gene classification. *BMC Microbiology* **16**, 274 (2016).
13. Benítez-Páez, A., Portune, K.J. & Sanz, Y. Species-level resolution of 16S rRNA gene amplicons sequenced through the MinION™ portable nanopore sequencer. *GigaScience* **5** (2016).
14. Peters, B.A., Liu, J. & Drmanac, R. Co-barcoded sequence reads from long DNA fragments: a cost-effective solution for "perfect genome" sequencing. *Frontiers in Genetics* **5** (2015).
15. Wang, O. et al. Efficient and unique cobarcoding of second-generation sequencing reads from long DNA molecules enabling cost-effective and accurate sequencing, haplotyping, and de novo assembly. *Genome research* **29**, 798-808 (2019).
16. Adey, W. Algal Turf Scrubber (ATS), Algae to Energy Project: Cleaning Rivers while Producing Biofuels and Agricultural and Health Products. Progress Report to the Lewis Foundation. *Smithsonian Institution* (2010).
17. Picelli, S. et al. Tn5 transposase and tagmentation procedures for massively scaled sequencing projects. *Genome research* **24**, 2033-2040 (2014).
18. Hennig, B.P. et al. Large-scale low-cost NGS library preparation using a robust Tn5 purification and tagmentation protocol. *G3: Genes, Genomes, Genetics* **8**, 79-89 (2018).
19. Wang, Y. et al. A practical random mutagenesis system for Ralstonia solanacearum strains causing bacterial wilt of Pogostemon cablin using Tn5 transposon. *World Journal of Microbiology and Biotechnology* **35**, 7 (2019).
20. Karst, S.M. et al. High-accuracy long-read amplicon sequences using unique molecular identifiers with Nanopore or PacBio sequencing. *Nature methods* **18**, 165-169 (2021).
21. Dong, Z. et al. Development of coupling controlled polymerizations by adapter-ligation in mate-pair sequencing for detection of various genomic variants in one single assay. *DNA Research* **26**, 313-325 (2019).

564  22.  Volden, R. et al. Improving nanopore read accuracy with the R2C2 method enables
565       the sequencing of highly multiplexed full-length single-cell cDNA. *Proceedings of the*
566       *National Academy of Sciences* **115**, 9726-9731 (2018).
567  23.  Adams, M. et al. One fly–one genome: chromosome-scale genome assembly of a
568       single outbred Drosophila melanogaster. *Nucleic acids research* **48**, e75-e75 (2020).
569  24.  Wang, O. et al. Efficient and unique cobarcoding of second-generation sequencing
570       reads from long DNA molecules enabling cost-effective and accurate sequencing,
571       haplotyping, and de novo assembly. *Genome research* **29**, 798-808 (2019).
572  25.  Peters, B. et al. A simple bead-based method for generating cost-effective co-
573       barcoded sequence reads. *Protocol Exchange* (2018).
574  26.  Fang, C. et al. Assessment of the cPAS-based BGISEQ-500 platform for
575       metagenomic sequencing. *Gigascience* **7**, 1-8 (2018).
576  27.  Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ
577       preprocessor. *Bioinformatics* **34**, i884-i890 (2018).
578  28.  Ondov, B.D. et al. Mash: fast genome and metagenome distance estimation using
579       MinHash. *Genome biology* **17**, 132 (2016).
580  29.  Li, D., Liu, C.M., Luo, R., Sadakane, K. & Lam, T.W. MEGAHIT: an ultra-fast single-
581       node solution for large and complex metagenomics assembly via succinct de Bruijn
582       graph. *Bioinformatics* **31**, 1674-1676 (2015).
583  30.  Quast, C. et al. The SILVA ribosomal RNA gene database project: improved data
584       processing and web-based tools. *Nucleic acids research* **41**, D590-D596 (2012).
585  31.  Altschul, S.F., Gish, W., Miller, W., Myers, E.W. & Lipman, D.J. Basic local alignment
586       search tool. *J Mol Biol* **215**, 403-410 (1990).
587  32.  Nilsson, R.H. et al. The UNITE database for molecular identification of fungi:
588       handling dark taxa and parallel taxonomic classifications. *Nucleic acids research* **47**,
589       D259-D264 (2019).
590  33.  Abarenkov, K.Z., Allan; Piirmann, Timo; Pöhönen, Raivo; Ivanov, Filipp; Nilsson, R.
591       Henrik; Kõljalg UNITE general FASTA release for eukaryotes 2. Version
592       04.02.2020. . *UNITE Community* (2020).
593  34.  Rognes, T., Flouri, T., Nichols, B., Quince, C. & Mahe, F. VSEARCH: a versatile
594       open source tool for metagenomics. *PeerJ* **4**, e2584 (2016).
595  35.  Wood, D.E., Lu, J. & Langmead, B. Improved metagenomic analysis with Kraken 2.
596       *Genome biology* **20**, 257 (2019).
597  36.  Bengtsson-Palme, J. et al. Improved software detection and extraction of ITS1 and
598       ITS2 from ribosomal ITS sequences of fungi and other eukaryotes for analysis of
599       environmental sequencing data. *Methods in Ecology and Evolution*, n/a-n/a (2013).
600  37.  Nakamura, T., Yamada, K.D., Tomii, K. & Katoh, K. Parallelization of MAFFT for
601       large-scale multiple sequence alignments. *Bioinformatics* **34**, 2490-2492 (2018).
602  38.  Capella-Gutiérrez, S., Silla-Martínez, J.M. & Gabaldón, T. trimAl: a tool for
603       automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*
604       **25**, 1972-1973 (2009).
605  39.  Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis
606       of large phylogenies. *Bioinformatics* **30**, 1312-1313 (2014).
607  40.  Yu, G., Smith, D.K., Zhu, H., Guan, Y. & Lam, T.T.-Y. ggtree: an r package for
608       visualization and annotation of phylogenetic trees with their covariates and other
609       associated data. *Methods in Ecology and Evolution* **8**, 28-36 (2017).
610  41.  Team, R. RStudio: Integrated Development Environment for R. *RStudio, PBC* (2020).
611  42.  Team, R.C. R: A Language and Environment for Statistical Computing. *R Foundation*
612       *for Statistical Computing* (2020).
613  43.  Guo, X. et al. CNSA: a data repository for archiving omics data. *Database* **2020**
614       (2020).
615  44.  Chen, F.Z. et al. CNGBdb: China National GeneBank DataBase. *Yi Chuan* **42**, 799-
616       809 (2020).

617