# Genome Enrichment of Rare, Unknown Species from Complicated Microbiome by Nanopore Selective Sequencing

Yuhong Sun [a, #], Xiang Li [a, b, c, #], Qing Yang [a], Bixi Zhao [a], Ziqi Wu [a], Yu Xia [a, b, c] *

[a] School of Environmental Science and Engineering, College of Engineering, Southern University of Science and Technology, Shenzhen 518055, China
[b] State Environmental Protection Key Laboratory of Integrated Surface Water-Groundwater Pollution Control, School of Environmental Science and Engineering, Southern University of Science and Technology, Shenzhen 518055, China
[c] Guangdong Provincial Key Laboratory of Soil and Groundwater Pollution Control, School of Environmental Science and Engineering, Southern University of Science and Technology, Shenzhen, 518055, China

**Keywords:**

Nanopore sequencing; selective sequencing; Read Until; Readfish; rare species; thermophilic anaerobic digestion

# These authors contributed equally to this work
*Corresponding author:
Yu Xia
Address: School of Environmental Science and Engineering, College of Engineering, Southern University of Science and Technology, Shenzhen 518055, China
E-mail: xiay@sustech.edu.cn

# Abstract

Rare species are vital members of a microbial community, but retrieving their genomes is difficult due to their low abundance. The ReadUntil (RU) approach allows nanopore devices to sequence specific DNA molecules selectively in real-time, which provides an opportunity for enriching rare species. However, there is still a gap in RU-based enriching of rare and unknown species in environmental samples whose community composition is unclear, and many species lack corresponding reference in public databases. Here we present metaRUpore to overcome this challenge. We applied metaRUpore to a thermophilic anaerobic digester (TAD) community, it successfully redirected the sequencing throughput from high-abundance populations to rare species while facilitating the recovery of 41 high-quality metagenome-assembled genomes (MAGs) at low sequencing effort. The simplicity and robustness of the approach make it accessible for labs with moderate computational resources and hold the potential to become the standard practice in future metagenomic sequencing of complicated microbiomes.

# 1 Introduction

Microbial communities are composed of a high number of rare species[1]. Rare species play a vital role in ecosystem health and stability[2]. For example, the slow-growing autotrophic microbes of ammonia-oxidizing bacteria or archaea (AOB/AOA) and anammox enable the rate-limiting step for natural nitrogen turnover[3,4]. Therefore, identifying the functional capacities of these rare species is essential to understanding the community dynamics and ecological function of a natural microbiome[2,3].

The recovery of draft genomes (referred to as metagenome-assembled genomes, MAGs)

25   from high-throughput metagenomic whole-genome sequencing (thereafter short as

26   metagenomic) datasets ushered in a new era for understanding the ecological and

27   evolutionary traits of the unculturable majority of natural microbiomes. However, high-

28   quality (HQ, usually defined as >90% completeness with <5% contamination and the

29   intact rRNA operon[44]) MAGs recovery for low abundant species is always difficult. In

30   metagenomic sequencing, the low-abundance microorganisms are often missed or

31   simply neglected due to low sequencing coverage. To get sufficient genome coverage of

32   low-abundance species, extremely deep sequencing will be required. It would be a great

33   waste if the study aims were to focus on rare species. Things can become more

34   intractable during the data analyses that recovering the unknown genomes from

35   hundreds of gigabytes to terabytes of data is a massive computational challenge[4].

36

37   To raise coverage of rare taxa from a high-abundance background, molecular biology-

38   based methods including hybrid capture or CRISPR-Cas9 enrichment are adapted in

39   library preparation to enrich target[5,6]. On the other hand, depletion of high abundance

40   species may serve the same purpose. Saponin-based host DNA depletion in human

41   metagenomic communities is used for rapid clinical diagnosis of relatively low

42   abundance pathogenic bacteria[7]. What is evident, however, is that these approaches

43   require the use of extra reagents and preparatory procedures. This is compounded by the

44   fact that they require known information about the enrichment or depletion targets in

45   order to design the experiment, which does not appear to work for enriching low

46   abundance species in metagenomic communities with unknown compositions.

47

48   Unlike the endeavors made prior to sequencing, Nanopore sequencing (Oxford

49   Nanopore Technology, ONT) users can program their system to reverse the voltage

50   polarity of the sequencing pore to eject reads identified as not of interest, which provides

51   a potential solution to enrich for rare species in metagenomic samples. This 'selective

52   sequencing' or Read Until (RU) strategy was first implemented by Loose and colleagues

53   in 2016[8]. The earliest adopted dynamic time warp (DTW) algorithm-based approach

54   could not scale to references larger than millions of bases, which limits its widespread

55 usage[8]. With the similar goal of mapping streaming raw signal to DNA reference,

56 UNCALLED has a lighter computational footprint than DTW[9]. Still, it requires abundant

57 computational resources. The newly designed Readfish toolkit eliminates the need for

58 complex signal mapping algorithms, and exploits existing ONT tools to provide a robust

59 toolkit for designing and controlling selective sequencing experiments[10]. Until now, the

60 application of RU is principally limited to the elimination of known host species[9, 10, 11]

61 or the enrichment of known targets such as mitogenomes of blood-feeding insects[12, 13].

62

63 By ejecting dominant species while accepting low-abundance species, selective

64 sequencing provides a potential solution to enrich rare species in metagenomic samples.

65 Nonetheless, enrichment for low abundance species in real metagenomic samples by

66 selective sequencing remains challenging because the community composition is never

67 known, and a large proportion of the species lacks a corresponding reference in public

68 databases. To specifically address such metagenomic-issue and to realize effective

69 targeted enrichment of rare species within a complicated environment microbiome, here

70 we introduced metaRUpore, a protocol consisting of know-how for configuring selective

71 nanopore sequencing and necessary bioinformatic scripts to achieve efficient enrichment

72 of rare species within a complicated environment microbiome. We initially assessed the

73 efficacy of enriching low abundance species in a mock community. Based on this

74 evaluation, we elaborated the principles and processes of metaRUpore and applied it to

75 a thermophilic anaerobic digester (TAD) community that was treating waste sludge of a

76 domestic wastewater treatment plant (WWTP). Meanwhile, we demonstrate a robust and

77 effective procedure for assembling and binning HQ-MAGs from RU-based nanopore

78 datasets. And an archaeal HQ-MAG retrieved from the TAD community revealed a giant

79 (112Kbp) function-related genomic island, extending the evolutionary traits of the

80 important *Bathyarchaeota* phylum.

## 2 Results

## *H. mediterrane* enrichment in a mock community

To evaluate nanopore performance on enriching low abundance species with RU, we firstly constructed a mock community. The *Haloferax mediterranei* strain which accounts for 1% of the mock community, was the target of our enrichment, while the other seven bacteria species were targets to be depleted during the RU run. In the mock run, a MinION flow cell was configured into two parts, where the first half of the channels did selective sequencing, and the other half did normal sequencing as a control. In the RU channels, the reads were basecalled and then mapped to a 33-M reference which contained all these eight microorganisms when they are being sequenced. A DNA molecule would be firstly sequenced for 0.4s before the obtained sequence was aligned to decide it should be sequenced continually or ejected. The average length of rejected reads was 537 bases, it demonstrated that the entire process of basecalling, mapping, and rejection decision could be completed in about 1.3s, based on the average nanopore sequencing seed of 400bp/s with R9.4.1 chemistry[10]. In the RU-delivered dataset, >99.9% of archaeal reads were kept while >99% of bacterial reads were ejected. *H. mediterranei* got enriched to the absolute dominant population within the community with a relative abundance of 62% in kept reads (Fig. 1a) with the coverage increased twice to 21.19× in RU data (Fig. 1b).

Despite the high rejection precision and fairly ideal enrichment result, it must be noted that the total yield of selective sequencing was approximately 60% lower than that of normal sequencing (Fig. 1c). This reduction in throughput can be partly attributed to the increased idle time of each nanopore caused by a large number of ejections[9]. At an enriched target prevalence of 1% within a community, each nanopore ejected an average of 2,430 short fragments while 267 continuous long fragments were sequenced in a 7-hour run. In addition, a rapid drop in active channels happed after 1-hour sequencing in RU channels (Fig. 1 d and Supplementary Fig. 1) and the effective pore got depleted

109  after 6-hour runtime which was 4 times shorter than normal run whose pores could

110  normally last for 24 hours (Fig. 1d). Consequently, it's critical to establish an appropriate

111  target proportion for selective sequencing to achieve the best tradeoff between

112  enrichment effectiveness and throughput loss. Fortunately, increasing sequencing effort

113  could easily compensate for the RU-induced per flow cell throughput loss.

## In situ Metagenomic selective sequencing protocol and performance

116  We introduced a pipeline, MetaRUpore ( https://github.com/sustc-xylab/metaRUpore),

117  to selectively sequence rare populations in complex microbiome samples. The protocol

118  consists of three consecutive steps (Fig. 2a): (1) 1h normal sequencing to obtain an

119  overall picture of the community structure and the genomic profile of the dominant

120  populations, (2) bioinformatics analysis to determine the reference and target dataset for

121  optimized RU configuration, and (3) finally a 40h selective sequencing for enriching

122  rare populations in the sample. The pore control of the nanopore device was

123  implemented by Readfish[10] which combines Guppy with minimap2[14] to determine the

124  eject/keep action for a pore.

125

126  Here we show our results in applying the metaRUpore protocol to facilitate the genome

127  recovery of rare populations within the TAD community, which consists of 2,977 OTUs

128  with a Shannon index of 8.74, representing a typical diversity level of bioreactor systems

129  (Supplementary Fig. 2). Rarefaction analysis demonstrated that the reads sequenced in

130  the first 1 h normal sequencing already cover 90% of the overall diversity in the TAD

131  community (Supplementary Fig. 5). Among the 125,606 reads sequenced, 66% of them

132  could be assigned to a known reference by Centrifuge[15]. All of these classified reads

133  obtained in the first 1 h run were set as the target for ejection in subsequent RU run as it

134  mostly consisted of the known and abundant populations within the community. Notably,

135  using whole-genome sequences from close species (same family or genus) as the

136  reference for RU run will result in poor performance in ejecting the dominant

137  populations because environmental microbiomes typically contain a high proportion of

138  genetic fragments that are distinct from all the sequences deposited in whole-genome

139  collections. In fact, even with the entire bacterial whole genome collection set as the

140  ejection target, only an ejection efficiency of 22% was achieved in RU sequencing of

141  the TAD community, leaving the community profile largely unchanged after selective

142  sequencing. Another thing to note is that the classified reads obtained in the firstly 1h

143  normal sequencing, inevitably contain genomic fragments from the rare and unknown

144  populations we intend to enrich, which will result in incomplete genome coverage of

145  rare populations in the sequences obtained in the RU channels. Therefore, a small

146  fraction of the channels still needed to be set to normal sequencing in the subsequent 40h

147  RU run and the delivered dataset needs to be assembled together with the RU-derived

148  datasets. For our RU-sequencing of the TAD community, we set 1/8 channels to normal

149  sequencing (--channels 1 448) (Fig. 2b). Our subsequent data analysis revealed that 29

150  HQ-MAGs would be missed if reads derived from selective sequencing were assembled

151  alone. To further manipulate the selection, the users can manually select which taxa to

152  keep during subsequent RU run; reads belonging to these taxa will be subtracted from

153  the target dataset based on their taxonomic affiliations determined by ARGpore2[16]. For

154  example, in our TAD community, we intended to keep all the archaea reads, so we

155  eliminated them from the ejection target datasets. The entire aforementioned

156  bioinformatic analysis can be completed in less than 30 min, such short suspension will

157  not affect the flow cell chemistry and the subsequent RU run may directly start without

158  refreshing the sequencing library.

159

160  The 40h RU run on one flow cell delivered 6.84 Gbp of effective long reads with an

161  average read length of 3.46 kbp, while the normal sequencing channels produced 1.71

162  Gbp reads with an average read length of 3.60 kbp (Supplementary Fig. 3). To ensure

163  adequate genome coverage, we sequenced the TAD community following metaRUpore

164  protocol using three flow cells one by one on GridION X5. Given the concern to exhaust

165  computation capacity on GridION X5, we did not test RU run with multiple flow cells

166  sequenced simultaneously. RU sequencing using metaRUpore protocol resulted in a

167   marked change in the community structure. As shown in the 3D density plot of

168   phylogeny distribution of the overall TAD community (Fig. 3a), several density peaks

169   of the original TAD community were depleted in the RU-run delivered datasets,

170   indicating DNA of the high abundance populations of the TAD community was

171   effectively ejected during RU-sequencing and the community got homogeneous with

172   coverage of different populations become much more unified. Such unified coverage of

173   different populations will help to minify the disparity of kmer frequency in the dataset,

174   preventing kmers of the rare species from being filtered out as error-containing kmers

175   due to coverage drop during the kmer-counting step of a *de novo* assembly algorithm[17,

176   18].

# Bioinformatics pipeline for *de novo* metagenomic assembly and genome recovery

179   As illustrated in the assembly pipeline (Fig. 2c), the 31G data consisting of RU and

180   normal sequencing were assembled together respectively using three different

181   assemblers, namely Canu[19], Unicycler[20], and metaFlye[21]. The basic statistics of

182   assembled contigs were summarized in Supplementary Table 1. To improve the

183   robustness of the binning, 139 > 1Mbp contigs were firstly picked, as the candidate of

184   HQ genome[22]. The rest shorter contigs derived by the three assemblers were respectively

185   binned by MetaBAT2[23]. Only contigs longer than 100 kbp were kept for subsequent

186   binning. The MAGs retrieved above were subject to consensus correction by Medaka

187   with nanopore data and polished by Pilon[24] with Illumina short reads (SRs). Next,

188   polished MAGs were further corrected for frame-shift errors using MEGAN-LR[22] based

189   on DIAMOND alignment against the *nr* database. Finally, MAGs obtained by the

190   different assemblers were de-duplicated using dRep[25] with a relatedness threshold of

191   ANI > 0.95 to obtain species-level representative MAGs. Totally, we obtained 46 draft-

192   quality MAGs after dereplication. Among them, 41 MAGs including 6 complete circular

193   genomes were high-quality (HQ) (Supplementary Fig. 8 and Supplementary Table 2).

194   32 of these HQ MAGs were firstly picked single >1Mbp contigs, while the remaining

195  15 HQ MAGs were obtained by binning. All of these MAGs contained less than 13

196  contigs with an average N50 > 2 Mbp, demonstrating that they are highly continuous. In

197  comparison, the normal nanopore sequencing dataset yielded 29 draft-quality MAGs,

198  including 16 HQ MAGs. 15 of them were included in the 41 HQ MAGs retrieved by

199  metaRUpore strategy (Supplementary Fig. 8). Worth noting is that the 26 HQ MAGs

200  that are additionally obtained by RU-based selective sequencing were mainly from the

201  rare populations of the TAD community (Fig. 3b). Additionally, evident coverage

202  reduction was observed in the dominant populations that the coverage of MAG17,

203  MAG4, and MAG30, which together accounted for 21% of the TAD community,

204  dramatically reduced by 78% after RU-based selective sequencing (Fig. 3b and

205  Supplementary Table 3), demonstrating the effectiveness of metaRUpore protocol in

206  eliminating dominant populations during sequencing. Despite the lowered overall

207  throughput, coverage of the rare species MAG33, MAG35, MAG57, and MAG56 was

208  doubled at the current sequencing effort and the application of the metaRUpore protocol

209  has reduced the abundance limit for HQ-MAG recovery in the TAD community to 0.7%.

210  It could be expected that by using additional flow cells, HQ-MAGs could be obtained

211  for populations with even lower prevalence.

212

## 3 Discussion

## Complete genomes recovered from TAD community

215  The 41 HQ MAGs introduce 5 new phyla, namely *WOR-3*, *OLB16*, *Omnitrophota*,

216  *Gemmatimonadota*, and *Deferribacterota,* into the global HQ genome collection of AD

217  microbiome[26] (Fig. 4). Furthermore, our MAGs show much better integrity and

218  continuity than those in the previous collection assembled with SRs in terms of N50,

219  number of contigs as well as intact rRNA operon. Additionally, evolutional traits

220  analysis reveals a much more conservative scale of gene flow based on HQ genomes we

221  assembled than that based on fragmented MAGs[27] (Supplementary Fig. 9) .

222

## Versatile metabolic capacities of *Bathyarchaeota* phylum in TAD community

*Bathyarchaeota* was recently recognized as a methanogenesis contributor[28] that may play active roles in global biogeochemical cycles[31]. However, the absence of pure cultures of the phyla has hampered our understanding of their ecological functions and evolutionary positions from a genome-centric perspective[29,30]. Genomes reported for this phylum so far are highly fragmented (Fig 5a). In this work, MetaRUpore has boosted the abundance of *Bathyarchaeota* in the TAD community from 0.19% to 0.32%, facilitating its genome recovery as MAG56, which to the best of our knowledge, is the first complete genome for this phylum. MAG56 represented a novel *Bathyarchaeota* lineage with the closest neighbor being Bathy-5 (Fig 5b). The genome size of MAG56 is 1.9Mbp, notably larger than the average size of previously assembled genomes of *Bathyarchaeota* phylum (1.23Mbp)[29,30,31]. *Bathyarchaeota* was previously proposed to have methyl-dependent hydrogenotrophic methanogenic potential[28,32] as MAGs recovered from deep aquifers[34] possess an MCR-like complex. However, no MCR homology could be detected in MAG56. Given the complete nature of the genome obtained in this study, a functioning methanogenic pathway in the TAD community lineage of *Bathyarchaeota* seemed implausible.

Remarkably, we found three genomic islands (GIs) (Fig 5a) in MAG56 with the largest being 36 kbp in length. These GIs were always missing in previously genomes assembled by short reads due to the defective resolving of repetitive fragments flanking the exogenous genetic island[33,34]. In the largest GIs of 36 Kbp, we identified six copies of Tyrosine recombinase (*xerA*, *xerC*, or *xerD*), which had previously been reported to facilitate the insertion of gene islands into the host chromosome by catalyzing site-specific, energy-independent DNA recombination[34,36]. Additionally, we identified a heat shock protein, *HtpX*, that may contribute to the heat shock response facilitating the cell's

250    survival in a thermophilic environment. Collectively, this GI represents a highly mobile

251    fitness island[33] that offers selective advantages for the archaeal population within the

252    thermophilic digester community. And the recovery of complete MAGs by metaRUpore

253    undoubtedly enabled the discovery of the role of large GIs in shaping *Bathyarchaeota*'s

254    evolution.

255

256    Overall, we proposed metaRUpore, a method for enriching low-abundance and

257    undiscovered microorganisms in complex microbial communities based on nanopore

258    selective sequencing. The heuristic ejecting targets determined through initial short-term

259    *de novo* sequencing of the dominant populations, overcome the constraints imposed by

260    the absence of reference genomes for selective sequencing of complex communities.

261    metaRUpore unifies the sequenced community structure and increases the genome

262    coverage of low-abundance species, facilitating the assembly of additional HQ genomes

263    of rare species within the microbiota. HQ MAGs retrieved from the TAD community by

264    metaRUpore contribute to the building of a more comprehensive database of AD-

265    associated microbes, which will ultimately allow for an in-depth understanding of their

266    biological characteristics. More importantly, metaRUpore protocol is robust and requires

267    minimal modification to the experimental procedure of nanopore library construction

268    and sequencing, making it easily applicable to metagenomic investigations of other

269    environmental microbiomes. Even though selective sequencing for the rare sphere is

270    inevitably associated with a reduction in per-flow cell data yield. Future implementation

271    of the RU API on PromethION will easily provide a throughput boost, overcoming the

272    coverage barrier and enabling complete genome recovery of rare species with even lower

273    abundance from complex microbiomes using the metaRUpore protocol.

## 274  4 Methods

275    **Sampling and DNA extraction**

276    Genomic DNA of the eight microorganisms of the mock community was extracted by QIAamp DNA

277    Micro Kit (50). Samples for TAD community were taken when the methanogenic bacteria were at

278    their highest activity. Genomic DNA of the TAD community samples was extracted by QIAGEN

279    DNeasyR PowerSoilR Kit (100). DNA concentration was determined using the Life Technologies

280    Qubit high sensitivity assay kits. The quality of the DNA was measured by Thermo Scientific™
281    NanoDrop™ to assure that it all met the requirements for library construction.
282
283    **Construction of the synthetic mocks**
284    We synthesized a mock community of eight microorganisms, of which Archaea accounted for 1%
285    and the other seven bacteria species shared the rest equally based on DNA concentration determined
286    from qubit average measurements. The archaeal species is *Haloferax mediterranei* and these seven
287    bacteria are *Acinetobacter baumannii*, *Enterococcus faecalis*, *Escherichia coli*, *Klebsiella*
288    *pneumoniae*, *Pseudomonas aeruginosa*, *Serratia marcescens*, and Staphylococcus aureus.
289
290    **Library construction and Sequencing**
291    All sequencing libraries were constructed using the ONT Ligation Sequencing Kit (no. SQK-LSK109)
292    according to the manufacturer's instructions. When preparing the reactor sample libraries, in order to
293    remove as many very short DNA fragments as possible, 0.4X beads was used for each step of the
294    cleanup, and therefore the initial amount of genomic DNA was increased to 2ug to ensure a sufficient
295    amount of DNA of the final library. ONT MinION flowcells v.R9.4.1 were used for all sequencing
296    on an ONT GridION.
297
298    **Selective sequencing via metaRUpore**
299    The execution of metaRUpore to enrich for unknown low abundance taxa is divided into the
300    following three steps: firstly, a period (in this case 60 min) of normal sequencing is performed to
301    generate reference file for selective sequncing using Readfish[10] which should contain the vast
302    majority of taxa in the community. Next, the sequenced data is fed into metaRUpore to obtain the
303    reference and target needed to configure Readfish TOML for selective sequencing. During this time,
304    it is advisable to keep the MinION flowcell with the DNA library in a 4℃ refrigerator to avoid the
305    loss of activity of the nanopores affecting the subsequent sequencing. We put the reference and target
306    paths into the TOML file and set config_name = "dna_r9.4.1_450bps_fast", single_on = unblock,
307    multi_on = unblock, single_off = stop_receiving, multi_off = stop_receiving, no_seq = proceed,
308    no_map = proceed. As recommended by the author of Readfish, we deactivated adapter scaling by
309    editing the config files (dna_r9.4.1_450bps_fast.cfg) in the guppy data directory. Next, selective
310    sequencing was started. the configuration on MinKNOW was the same as for normal sequencing.
311    Readfish runs at the same time as the sequencing starts.
312
313    **Analysis of long-read sequence data**
314    Sequencing-derived fastq reads were performed adaptor trimming using Porechop (GitHub -
315    rrwick/Porechop) (version 0.2.2) with default settings. These reads were subsequently assembled by
316    the three tools: Canu[19] (version 2.2, default setting except -nanopore, genomeSize=3m,
317    maxInputCoverage=10000, corOutCoverage=10000, corMhapSensitivity=high, corMinCoverage=0,
318    redMemory=32, oeaMemory=32, batMemory=200 useGrid=false), Unicycler[20] (version 0.4.9b,
319    default setting except -t 40, --keep 3) and Flye[17] (version 2.8.3, default setting except –nano-raw, --
320    threads 50, --plasmids, --meta, --debug). Generated contigs that was at least 1Mbp in length were
321    regarded as potential whole-chromosome sequence. Among the remaining contigs that are less than
322    1Mbp, we did metagenomic binning for the contigs that are greater than 100kbp in length. Metabat2[21]

323    (version 2.12.1 with default setting) is used to respectively binning the contigs assembled by above
324    three assemblers.
325    Next, we took multiple steps to correct the >1Mbp potential chromosome and bins we obtained.
326    Firstly, we used nanopore data to perform consensus correction on them using Medaka (GitHub -
327    nanoporetech/medaka)(version 1.4.3, default setting except -t 20, -m r941_min_high_g360). They
328    were then further corrected with the short reads data using Pilon[23] (version 1.24 with default setting
329    except --fix all, --vcf). We used DIAMOND[35] (version 0.9.24) to align the Pilon polished potential
330    chromosome (with default settings except -f 100 -p 40 -v --log --long-reads -c1 -b12) against the
331    NBCI–NR database[38] (July 2021). We used daa-meganizer in MEGAN Community Edition suite[39]
332    (version 6.21.7, run with default settings except --longReads, --lcaAlgorithm longReads, --
333    lcaCoveragePercent 51, --readAssignmentMode alignedBases) to format the .daa output file and
334    receive frame-shift corrected sequence with 'Export Frame-Shift Corrected Reads' option.
335    We checked the completeness and contamination of these potential genomes with CheckM[40] (version
336    v1.0.12, run with default setting except lineage_wf, -t 20). All the putative genomes were de-
337    replicated using the dRep[25] (version 3.2.2, run with default setting except -p 40 -sa 0.95 –genomeInfo)
338    to get species-level unique MAGs. Next, gene annotations were obtained using Prokka[41] (version
339    1.13). Microbial taxonomic classifications were assigned using GTDB-Tk[42] (version 1.3.0, GTDB-
340    Tk reference data version r89).
341
342    **Calculation of the abundance and assessment of the quality of MAG**
343    Abundance was calculated from both selective sequencing data and normal sequencing data, by
344    mapping these data to the MAGs using minimap2[14] (version 2.17) separately using the following
345    flags -ax map-ont -t 40. We used samtools[43] (version 1.11) to extract .sam file that matched each
346    MAG individually. The abundance of each MAG is calculated by dividing the number of bases in all
347    reads in this .sam file by the total number of bases selectively sequenced or normally sequenced.
348    Analogously, sorted .bam files were used in the calculation of coverage of the MAGs.
349    We defined high-quality (HQ) MAGs as encoding multiple rRNA genes (23S/16S/5S), SCG-
350    completeness > 90% and contamination < 5%[44]. Draft-quality (DQ) MAGs means MAGs having >
351    70% SCG-completeness, < 10% contamination, and the presence of 16S rRNA. While if a MAG
352    meets all of the DQ criteria but misses 16S rRNA were regarded as low-quality (LQ) genomes.
353
354
355    **Code availability**
356    The metaRUpore workflow is available on the GitHub page: https://github.com/sustc-
357    xylab/metaRUpore.
358
359    **Availability of data and materials**
360    The raw nucleotide sequence data (both Illumina and Nanopore) used in the present study has
361    been deposited in the NCBI database under project ID PRJNA794848.
362

371  **Conflict of interests**
372  The authors claim no conflict of interests.
373

# Reference

375  1.  Jousset, Alexandre, et al. Where less may be more: how the rare biosphere
376      pulls ecosystems strings. *ISME J.* **11**, 853-862 (2017).
377  2.  Shade, Ashley, et al. Conditionally rare taxa disproportionately contribute to.
378      temporal changes in microbial diversity. *MBio* **5**, e01371-14 (2014).
379  3.  Xiong, Chao, et al. Rare taxa maintain the stability of crop mycobiomes and
380      ecosystem functions. *Environmental Microbiology* **23,** 907-1924 (2021).
381  4.  Pop, Mihai. Genome assembly reborn: recent computational
382      challenges. *Briefings in bioinformatics* **10**, 354-366 (2009).
383  5.  Gilpatrick, T. et al. Targeted nanopore sequencing with Cas9-guided adapter
384      ligation. *Nat. Biotechnol.* **38**, 433-438 (2020).
385  6.  Gu, Wei, et al. Depletion of Abundant Sequences by Hybridization (DASH):
386      using Cas9 to remove unwanted high-abundance species in sequencing libraries
387      and molecular counting applications. *Genome biology* **17**, 1-13 (2016).
388  7.  Charalampous, Themoula, et al. Nanopore metagenomics enables rapid clinical
389      diagnosis of bacterial lower respiratory infection. *Nat. Biotechnol.* **37**, 783-792
390      (2019).
391  8.  Loose, M., Malla, S. & Stout, M. Real-time selective sequencing using nanopore
392      technology. *Nat. Methods* **13**, 751–754 (2016).
393  9.  Kovaka, Sam, et al. Targeted nanopore sequencing by real-time mapping of raw
394      electrical signal with UNCALLED. *Nat. Biotechnol.* **39**, 431-441 (2020).
395  10. Payne, A. et al. Readfish enables targeted nanopore sequencing of gigabase-
396      sized genomes. *Nat. Biotechnol.* **39**, 442-450 (2020).
397  11. Gan, M. et al. Combined nanopore adaptive sequencing and enzyme-based host
398      depletion efficiently enriched microbial sequences and identified missing
399      respiratory pathogens. *BMC genomics* **22**, 1-11 (2021).
400  12. Kipp, E. J. et al. Nanopore adaptive sampling for mitogenome sequencing and
401      bloodmeal identification in hematophagous insects. *bioRxiv* (2021).
402  13. Martin, S. et al. Nanopore adaptive sampling: a tool for enrichment of low
403      abundance species in metagenomic samples. *bioRxiv* (2021).
404  14. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*
405      **34,** 3094–3100 (2018).
406  15. Kim, D., Song, L., Breitwieser, F. P. & Salzberg, S. L. Centrifuge: Rapid and
407      sensitive classification of metagenomic sequences. *Genome Res.* **26**, 1721–1729
408      (2016).

409   16.   Xia, Y. et al. MinION Nanopore sequencing enables correlation between
410         resistome phenotype and genotype of coliform bacteria in municipal sewage.
411         *Front. Microbiol.* **8**, 1–13 (2017).
412   17.   Kolmogorov, M. et al. metaFlye: scalable long-read metagenome assembly
413         using repeat graphs. *Nat. Methods* **17**, 1103-1110 (2020).
414   18.   Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes : a
415         new versatile metagenomic assembler. *Genome research* **27**, 824-834 (2017).
416   19.   Koren, S. et al. Canu : scalable and accurate long-read assembly via adaptive k -
417         mer weighting and repeat separation. *Genome research* **27**, 722–736 (2017).
418   20.   Wick, Ryan R., et al. Unicycler: resolving bacterial genome assemblies from
419         short and long sequencing reads. *PLoS computational biology* **13,** e1005595
420         (2017).
421   21.   Kang, D. D. et al. MetaBAT 2 : an adaptive binning algorithm for robust and ef
422         fi cient genome reconstruction from metagenome assemblies. *PeerJ* **7**, e7359
423         (2019).
424   22.   Arumugam, K. et al. Recovery of complete genomes and non-chromosomal
425         replicons from activated sludge enrichment microbial communities with long
426         read metagenome sequencing. *NPJ Biofilms Microbiomes* **7**, 1-13 (2021).
427   23.   Walker, B. J. et al. Pilon : An Integrated Tool for Comprehensive Microbial
428         Variant Detection and Genome Assembly Improvement. *PloS one* **9**, e112963
429         (2014).
430   24.   Huson, D. H. et al. MEGAN-LR: New algorithms allow accurate binning and
431         easy interactive exploration of metagenomic long reads and contigs. *Biol. Direct*
432         **13**, 1–17 (2018).
433   25.   Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. DRep: A tool for fast
434         and accurate genomic comparisons that enables improved genome recovery
435         from metagenomes through de-replication. *ISME J.* **11**, 2864–2868 (2017).
436   26.   Campanaro, S. et al. New insights from the biogas microbiome by
437         comprehensive genome-resolved metagenomics of nearly 1600 species
438         originating from multiple anaerobic digesters. *Biotechnol. Biofuels* **13**, 1–18
439         (2020).
440   27.   Alvarez-Ponce, David, et al. Gene similarity networks provide tools for
441         understanding eukaryote origins and evolution. *Proceedings of the National*
442         *Academy of Sciences* **110**, E1594-E1603 (2013).
443   28.   Evans, P. N. et al. Methane metabolism in the archaeal phylum Bathyarchaeota
444         revealed by genome-centric metagenomics. *Science* **350**, 434–438 (2015).
445   29.   Yu, Tiantian, et al. Growth of sedimentary Bathyarchaeota on lignin as an
446         energy source. *Proceedings of the National Academy of Sciences* **115**, 6022-
447         6027 (2018).
448   30.   Gagen, Emma J., et al. Novel cultivation-based approach to understanding the
449         miscellaneous crenarchaeotic group (MCG) archaea from sedimentary
450         ecosystems. *Applied and Environmental Microbiology* **79**, 6400-6406 (2013).
451   31.   Feng, X., Wang, Y., Zubin, R. & Wang, F. Core Metabolic Features and Hot
452         Origin of Bathyarchaeota. *Engineering* **5**, 498–504 (2019).

453 32. Borrel, G. et al. Wide diversity of methane and short-chain alkane metabolisms
454 in uncultured archaea. *Nat. Microbiol.* **4**, 603–613 (2019).

455 33. Juhas, M. et al. Genomic islands : tools of bacterial horizontal gene transfer and
456 evolution. *FEMS microbiology reviews* **33**, 376–393 (2009).

457 34. Nicholls, Samuel M., et al. Ultra-deep, long-read nanopore sequencing of mock
458 microbial community standards. *Gigascience* **8**, giz043 (2019).

459 35. Dorman, C. J. & Bogue, M. M. The interplay between DNA topology and
460 accessory factors in site-specific recombination in bacteria and their
461 bacteriophages. *Science progress* **99**, 420-437 (2016).

462 36. Badel, Catherine, Violette Da Cunha, and Jacques Oberto. "Archaeal tyrosine
463 recombinases." *FEMS Microbiology Reviews* **45**,1-27(2021).

464 37. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
465 DIAMOND. *Nat. methods* **12**, 59-60 (2015).

466 38. Leary, N. A. O. et al. Reference sequence ( RefSeq ) database at NCBI : current
467 status , taxonomic expansion , and functional annotation. *Nucleic acids research*
468 **44**, 733–745 (2016).

469 39. Huson, Daniel H., et al. MEGAN community edition-interactive exploration and
470 analysis of large-scale microbiome sequencing data. *PLoS computational*
471 *biology* **12**, e1004957 (2016).

472 40. Parks, Donovan H., et al. CheckM: assessing the quality of microbial genomes
473 recovered from isolates, single cells, and metagenomes. *Genome research* **25**
474 1043-1055 (2015).

475 41. Seemann, Torsten. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*
476 **30**, 2068-2069 (2014).

477 42. Chaumeil, Pierre-Alain, et al. GTDB-Tk: a toolkit to classify genomes with the
478 Genome Taxonomy Database. *Bioinformatics* **36**, 1925-1927 (2020).

479 43. Li, Heng, et al. The sequence alignment/map format and SAMtools.
480 *Bioinformatics* **25**, 2078-2079 (2009).

481 44. owers, R. M. et al. perspective Minimum information about a single amplified
482 genome ( MISAG ) and a metagenome-assembled genome ( MIMAG ) of
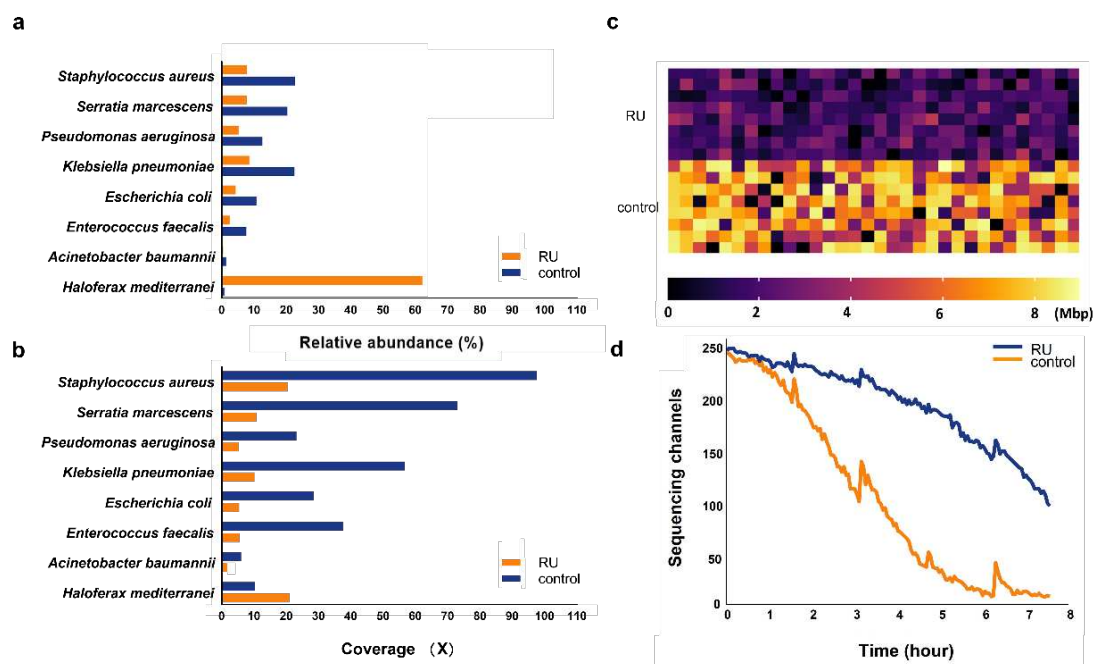483 bacteria and archaea. *Nat. Biotechnol.* **35**, (2017).

**Fig. 1 Enriching low abundance species in mock community with RU. a**, Bar plot of the abundance of the seven microbial species in RU and control runs. **b**, Bar plot of the coverage of the seven microbial species' genome in RU and control runs. **c**, heatmap of data yield per channel in RU and control runs, and **d**, plot of the number of sequencing channels over the course of the sequencing run.
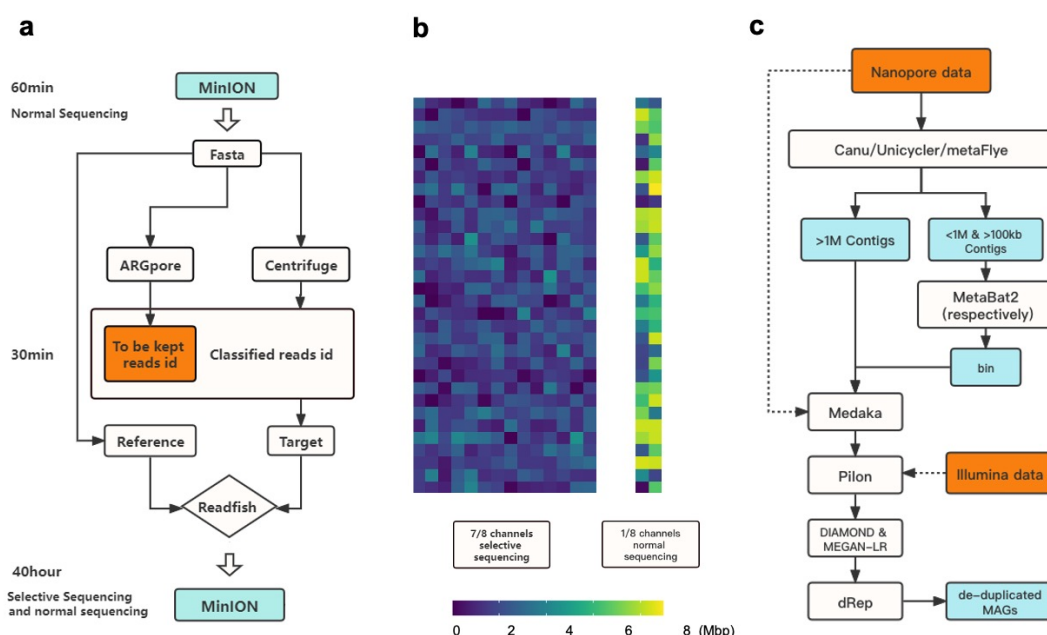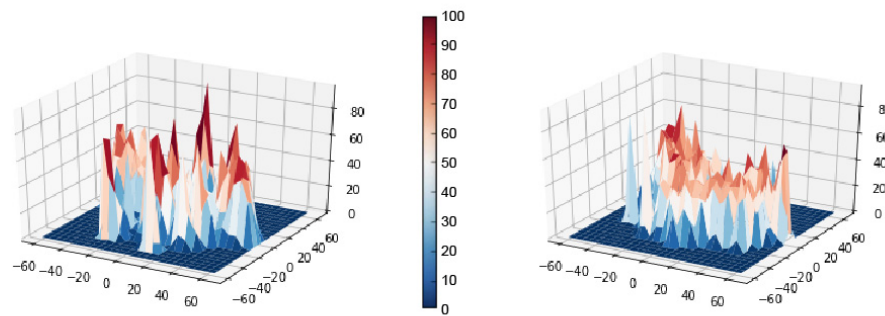


**Fig. 2 a,** The workflow of metaRUpore. **b,** A MinION flow cell in metaRUpore is configured into two parts, 1/8th of the channels for normal sequencing and the remaining channels for selective sequencing. **c,** The bioinformatic workflow for HQ-MAGs retrieval based on datasets derived from nanopore selective sequencing and Illumina sequencing.
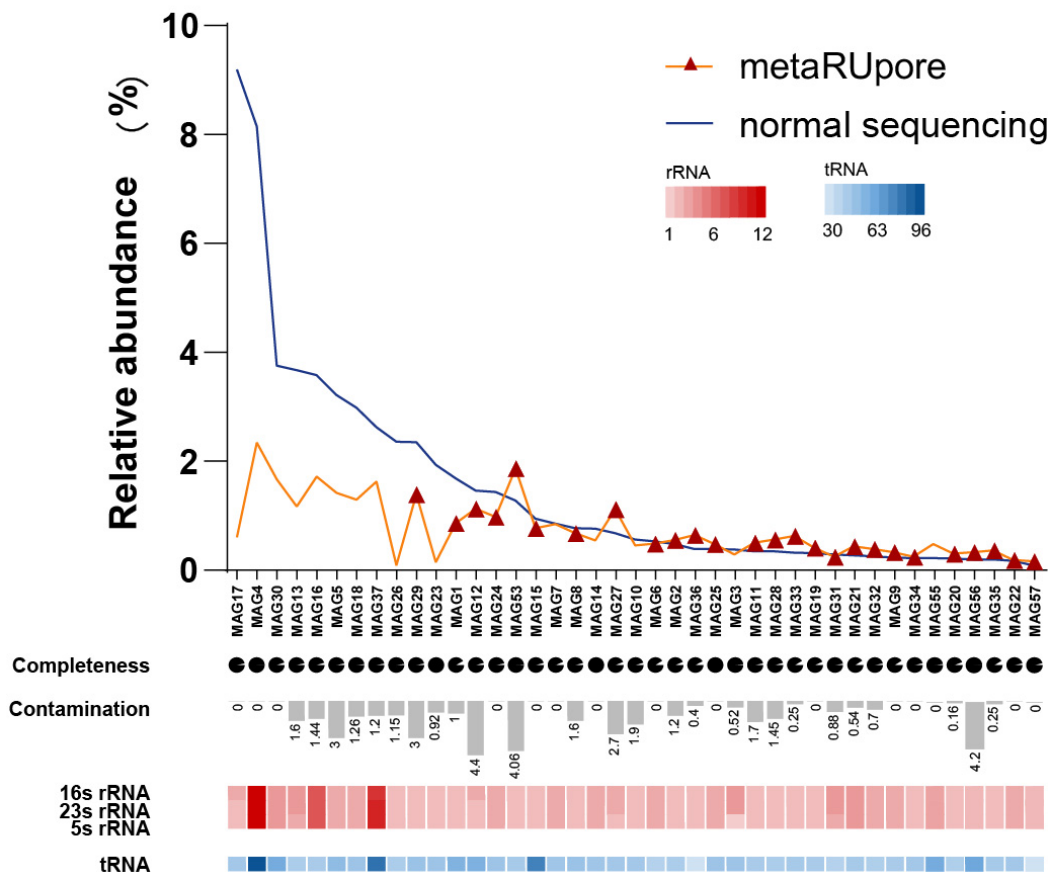
**Fig. 3 Performance of metaRUpore on recovery of high-quality MAGs in TAD community. a**, 3D density plots of t-SNE downscaling results for normal sequencing datasets and selective sequencing datasets by metaRUpore at four base frequencies, showing that metaRUpore renders the community structure homogenous. **b**, The distribution of 41 retrieved HQ MAGs in normal and RU sequencing dataset. The red triangles indicate MAGs that were could only be assembled in the metaRUpore dataset. The pie chart and bar chart represent the level of genomic completeness and contamination by CheckM. The copy number of 16S rRNA, 23S rRNA, and 5S rRNA is represented by the red heatmap, while the copy number of tRNA is represented by the blue heatmap.
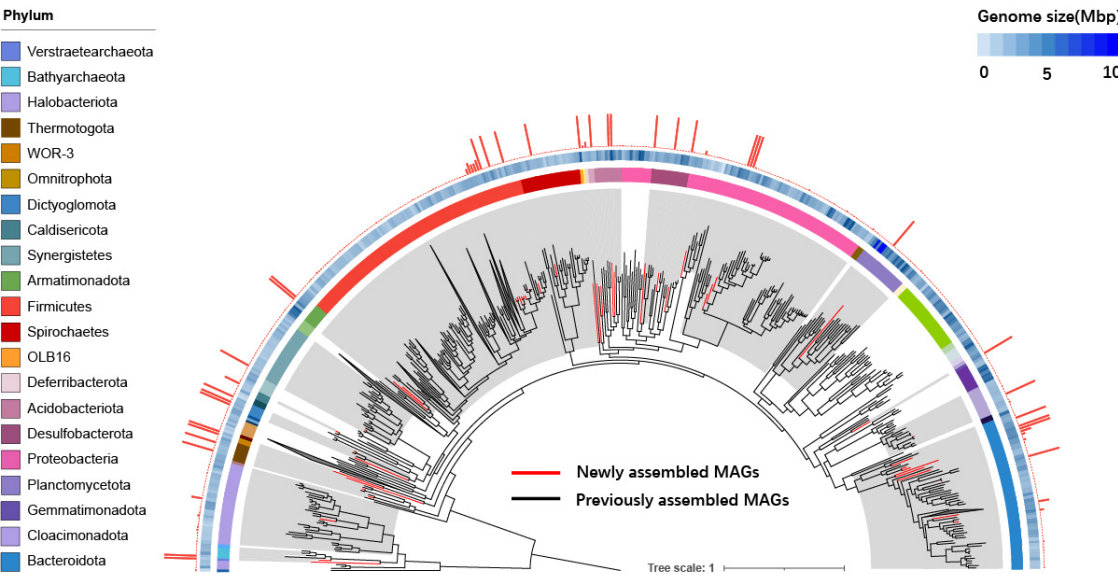
505



506

**Fig. 4 Phylogenomics of MAGs in anerobic reactor.** A phylogenetic tree was constructed from 41 HQ-MAGs derived by metaRUpore (red branches) and 1,108 HQ-MAGs collection derived from other AD systems (black branches). External circles represent, respectively: (1) taxonomic assignment at phylum level, (2) genome size (heatmap), (3) bar plot representing the genome continuity, which is calculated as the reciprocal of the number of contigs. The grey shaded areas indicate phyla with near-complete genomes obtained by metaRUpore，and the name of each phylum is in the legend on the left.
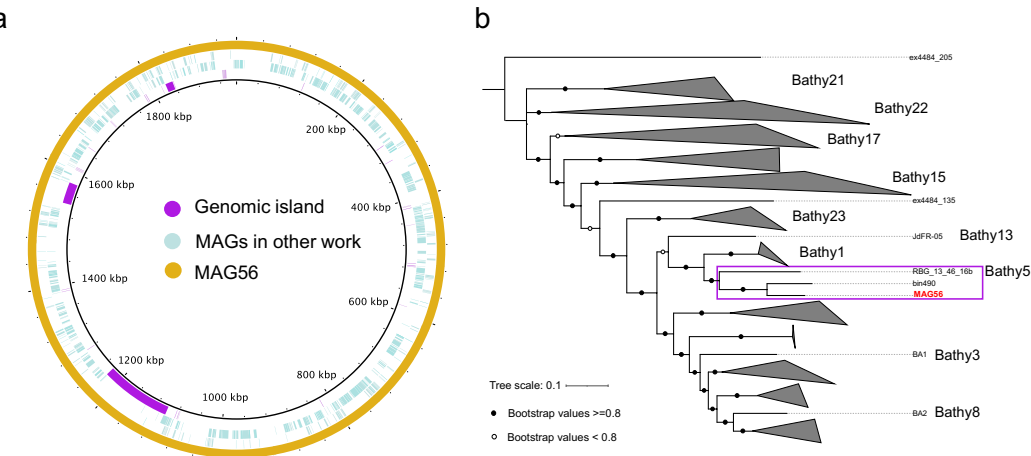


514

**Fig. 5** a, Genomes comparison of MAG56 and other MAGs of Bathyarchaeota from prior research. The outermost ring stands for the circular genome of MAG56 reconstructed by metaRUpore. The second to third circles from the outside represent the MAGs of phylum Bathyarchaeota reconstructed by short reads-only assembly method (MAGs covered by purple boxes in Figure 5c). The innermost purple circle represents the genomic island. b, A Maximum Likelihood Tree showing the phylogeny of Bathyarchaeota based on the MAGs from the current study (MAG56) and prior research[29]. Bootstrap values for these phylogenies are shown with open ( < 80%) and filled ( ≥ 80%) circles.

523