

# 1 **Profiling of Human Gut Virome with Oxford Nanopore Technology**

2 Jiabao Cao<sup>1,2,#</sup>, Yuqing Zhang<sup>1,2,#</sup>, Min Dai<sup>3</sup>, Jiayue Xu<sup>1</sup>, Liang Chen<sup>1</sup>, Faming

3 Zhang<sup>3,4</sup>, Na Zhao<sup>1,\*</sup>, Jun Wang<sup>1,\*</sup>

## 4 **Author Affiliations:**

5 1. CAS Key Laboratory of Pathogenic Microbiology and Immunology, Institute of  
6 Microbiology, Chinese Academy of Sciences, Beijing, China.

7 2. University of Chinese Academy of Sciences, Beijing 100049, China.

8 3. Medical Center for Digestive Diseases, the Second Affiliated Hospital of Nanjing  
9 Medical University, Nanjing 210011, China.

10 4. Key Lab of Holistic Integrative Enterology, Nanjing Medical University, Nanjing  
11 210011, China.

12 <sup>#</sup> Jiabao Cao, and Yuqing Zhang contributed equally to this manuscript.

13 \* Correspondence to:

14 Professor Jun Wang, CAS Key Laboratory of Pathogenic Microbiology and  
15 Immunology, Institute of Microbiology, Chinese Academy of Sciences, No. 1-3  
16 Beichenxi Road, Chaoyang District, Beijing, China; [junwang@im.ac.cn](mailto:junwang@im.ac.cn); Dr. Na Zhao,

17 CAS Key Laboratory of Pathogenic Microbiology and Immunology, Institute of  
18 Microbiology, Chinese Academy of Sciences, No. 1-3 Beichenxi Road, Chaoyang  
19 District, Beijing, China; [zhaona@im.ac.cn](mailto:zhaona@im.ac.cn).

20

21 ***Abstract***

22 Human gut virome play critical roles in maintaining gut microbial composition and  
23 functionality, as well as host physiology and immunology. Yet, there are insufficient  
24 amount of studies on this topic mainly due to methodological limitations, including  
25 enrichment of viruses (phages and host viruses) as well as short read-length from  
26 current sequencing technology. Here we developed a full working protocol for  
27 analyzing human gut virome using physical enrichment, reverse transcription and  
28 random amplification, and eventually the state-of-art single-molecule real-time  
29 sequencing (SMRT) platform of Oxford Nanopore Technology (ONT). We  
30 demonstrate that sequencing viral DNA directly, or viral DNA/cDNA after  
31 amplification using ONT achieves much longer reads and provides more information  
32 regarding virome diversity, many of the virome sequences do not have match in  
33 current databases. Moreover, direct DNA sequencing of virome provides first insights  
34 into the epigenetic modifications on phages, where signals of methylations can be  
35 directly detected. Our study demonstrates that progressing sequencing technology and  
36 bioinformatic improvements will bring more knowledge into virome composition,  
37 diversity and potentially their important functions.

38 ***Highlights:***

- 39 1. Virus-like particles were enriched from human stool samples;
- 40 2. Viral nucleotides were sequenced with Oxford Nanopore Technology with and  
41 without amplification;
- 42 3. Gut virome in humans showed highly individualized composition;

43 4. Novel sequences and contigs were found to be the majority in the resulted

44 sequences;

45 5. Epigenetic modifications were detected directly on virus genomes.

46 ***Keywords:***

47 Human gut virome; Enrichment; Amplification; Oxford Nanopore Technology;

48 Epigenomics

## 49 *Introduction*

50       The human gut is home to tremendous amount of microbes [1]. They inhabit  
51 different ecological niches in the gut, forming complex interaction networks between  
52 themselves and with the human cells, and the dynamic balance between gut  
53 microbiome and host is required for human health [2-6]. Studies in human cohort and  
54 mouse models, among others, have confirmed that gut microbial communities are  
55 associated with increasing number of some of metabolism diseases and infectious  
56 diseases, providing insights as well as potential targets for future monitoring and  
57 therapies [3, 7-14].

58       The gut microbiome contains bacteria, archaea, fungi, protozoa, and, lastly but  
59 not leastly virus. The most abundant cellular members of the microbiome are bacteria  
60 and archaea (account for more than 99% of biomass), and have received most  
61 attention in human microbiome studies over the years [15-17]. Yet, the advances in  
62 next-generation sequencing (NGS) technology and bioinformatic tools have also  
63 facilitated the development of human virome studies. Metagenomic analysis suggests  
64 that the gut of healthy humans harbors commensal virus, including phages, DNA  
65 virus and RNA virus [18-22]. Virome (phages and other host viruses) play roles in  
66 intestinal physiology, enteric immune system, host health and disease [23, 24]. The  
67 dynamic balance between the virome and the intestinal immune system is finely  
68 regulated by cytokines secreted by immune cells [20, 25]. For instance, virome  
69 changes in inflammatory bowel disease (IBD) (Crohn's disease and ulcerative colitis)  
70 are disease specific [26]. Phages residing in mucosal surfaces can influence the host

71 by providing nonhost-derived immunity against bacterial infections [26, 27]. By  
72 inducing interferons (IFNs), commensal virus can protect from gut inflammation  
73 during tissue damage [28, 29]. However, using current short read sequencing  
74 technologies, such as Illumina, can only offer knowledge on gut virome that is both  
75 biased and fragmentary.

76 Oxford Nanopore Technologies (ONT) as one of the emerging single-molecule  
77 real-time sequencing technology (SMRT) has the advantage of rapid library  
78 preparation, ultra-long reads and real-time data acquisition [30-32]. For virome, ONT  
79 sequencing has the potential to acquire virus genome by producing genome-length  
80 reads that cover all of mutation within a single virus particle. In addition, biological  
81 nanopores are able to discriminate not only the genome but also single base  
82 modifications such as 5-methylation of cytosine (5mC for DNA and m5C for RNA)  
83 and 6-methylation of adenine (6mA for DNA and m6A for RNA) in native DNA/RNA  
84 [32]. Increasing evidence in last years suggests that DNA/RNA methylation can  
85 influence biological function, including regulation of DNA/RNA replication and  
86 repair, and gene expression [33-36]. Recently, *Oliveira* et al. reported that DNA  
87 methyltransferase in *Clostridioides difficile* has able to mediate sporulation, *C.*  
88 *difficile* disease transmission and pathogenesis [37]. *Xue* et al. reported viral  
89 N<sup>6</sup>-methyladenosine could upregulate replication and pathogenesis of human  
90 respiratory syncytial virus [38]. These findings suggest that epigenetic regulation is  
91 also important for the pathogenesis of important pathogens.

92 To profile the gut virome in healthy adults, including identity as well as potential

93 epigenetic information in the virome, we developed a protocol combining physical  
94 enrichment, optional reverse transcription and amplification of nucleotides, and  
95 bioinformatic analytical pipelines, and firstly characterized the virome in five healthy  
96 humans using the ONT PromethION platform. We were able to generate long reads  
97 for virome up to tens of kilobases, resulting in many novel contigs that do not have  
98 matches in the available databases, and also for abundant virus we could detect  
99 epigenetic signals. These discoveries are instructive to future investigations into the  
100 genomics, epigenomics and potential function of human gut virome.

101

## 102 *Material and Methods*

### 103 1. Enrichment and purification of virus-like particles (VLPs)

104 Each frozen faecal samples (approximately 1.5 g) from five individuals who  
105 provided written informed consent were resuspended in 15 ml sterile Phosphate  
106 Buffered Saline and homogenized thoroughly. The suspension was centrifuged at  
107 4,500 rpm for 10 min at the 4 °C to remove large food residues (Beckman Coulter  
108 Allegra™ X-22R). Transfer the supernatant to fresh tubes and centrifuged at 4,500  
109 rpm for 10 min at the 4 °C again. The supernatant was filtered through 0.45 µm  
110 PVDF membrane (Millipore) to remove eukaryotic and bacterial cell-sized particles  
111 before ultracentrifugating at 180,000× g for 3 hours at the 4 °C (Beckman Coulter  
112 XP-100). The pellets were resuspended in 400 µl sterile Phosphate Buffered Saline  
113 and treated with 8 U of TURBO DNase I (Ambion) and 20 U of RNase A (Fermentas)  
114 at 37 °C for 30 min. The viral nucleic acids (DNA and RNA) were extracted by using

115 QIAamp MinElute Virus Spin Kit (Qiagen) following the manufacturer's instructions

116 and eluted into RNase-free water. [39, 40]

## 117 2. Reverse transcription and random amplification

118 Viral first strand cDNA was synthesized in a 20  $\mu$ l reaction mixture with 13  $\mu$ l of

119 purified viral nucleic acids from each sample and 100 pmol of primer Rrm

120 (5'-GACCATCTAGCGACCTCCAC - NNNNNN-3'), as previously described

121 [39-41]. For the double-strand cDNA synthesis, 100 pmol of primer Rrm and Klenow

122 fragment (3.5 U/ $\mu$ l; Takara) were added. Random amplification was conducted with

123 8 $\mu$ l of the double-strand cDNA template in a final reaction volume of 200 $\mu$ l, which

124 contained 4 $\mu$ M primer Rm (5'-GCCGGAGCTCTGCAGAATTC-3'), 90  $\mu$ M dNTPs

125 each, 80  $\mu$ M Mg<sup>2+</sup>, 10x Buffer and 1 U of KOD-Plus DNA polymerase (Toyobo). The

126 amplification product was purified by agarose gel electrophoresis and QIAquick Gel

127 Extraction Kit (Qiagen).

## 128 3. PromethION library preparation and sequencing

129 PromethION library preparation was performed according to the manufacturer's

130 instructions for the barcoding cDNA/DNA and native DNA (SQK-LSK109 and

131 EXP-NBD104). When multiplexing, all the samples were pooled together. ONT

132 MinKNOW software (v.19.10.1) was used to collect raw sequencing data, and Guppy

133 (v.3.2.4) was used for local base-calling of the raw data after sequencing runs were

134 completed. The PromethION was run for up to 96 h.

## 135 4. ONT sequence analysis and assembly

136 Qcat (Oxford Nanopore Technology), python command-line tool for

137 demultiplexing ONT reads from FASTQ files, was used to trim adaptor and barcode  
138 sequences. With genomeSize = 2k and default parameters, trimmed raw reads were  
139 analysis using Canu v1.9 [42] for virome genome de novo assembly, which includes  
140 read correction, read trimming and contig construction.

#### 141 5. Matching to current database

142 Raw reads qcat trimmed were analysis using mimimap2 [43] to identify gut  
143 virome composition, which aligned reads to the reference genome in the National  
144 Center for Biotechnology Information (NCBI) virus genome database, including all  
145 known viruse genome sequences. To improve the accuracy of the viral taxonomy, two  
146 following criteria were adapted: (1) the depth of coverage of reference viral  
147 genome  $\geq 5X$ ; (2) the breadth of coverage of the reference viral genome  $\geq 50\%$ .

148 To assign the taxonomy of assembled contigs by Canu, two approaches and three  
149 databases were applied, including minimap2, blastn, NCBI virus genome database,  
150 The human gut virome database (GVD) and NCBI nucleotide database. The filter  
151 criteria of alignment results of contigs by minimap2 was the same with raw reads,  
152 while contig matched length  $\geq 1000$  bp with nucleic similarity  $\geq 98\%$  and e-value  
153  $\leq 10^{-5}$  was adapted as the identified criteria of viruses by blastn.

#### 154 6. Identification and annotation of bacteriophage ORFs

155 Seeker [44], a new prediction tool via deep learning framework, was used to  
156 identify putative phages from contigs in amplified cDNA/DNA group with default  
157 parameters. According to the multiPhATE pipeline [45], ORFs in phages were  
158 identified by PHANOTATE [46], a tool to annotate phage genomes. Consequently,



159 amino sequences of ORFs were aligned to Phantome (<http://www.phantome.org>) and  
160 pVOGs [47] databases by blastp with parameters “percent of identity  $\geq$  60, e-value  
161  $\leq$  0.01” and hmm searched to pVOGs by jackhmmmer [48] with default parameters,  
162 respectively. ORFs repeated in two or more samples were clustered by usearch [49] with  
163 percent of identity  $\geq$  99. To validate these ORFs, we mapped our in-house NGS  
164 data using bowtie2 [50] with default parameters. Coverage of ORFs was calculated by  
165 weeSAM (<https://github.com/centre-for-virus-research/weeSAM>) and only ORFs  
166 with mapped reads  $\geq$  10 were counted.

## 167 7. Methylation analysis

168 Tombo v1.5 was used to detect the methylation states of nucleic acids from raw  
169 DNA samples [51]. A log-likelihood threshold of 2.5 was used to call methylation and  
170 the filter cutoff of methylation sites was defined by the estimated fraction of  
171 significantly modified reads  $\geq$  0.7 and coverage depth  $\geq$  10X. Different  
172 methylation sites of 5-methylcytosine (5mC) and N6-Methyladenine DNA  
173 Modification (6mA) were visualized by Integrative Genomics Viewer (IGV Version  
174 2.5.3) [52] with default parameters and the putative methyltransferase recognition  
175 motifs were identified by MEME (Version 5.0.5) [53] with the following parameters:  
176 “-dna -mod zoops”, and webLogo (<https://weblogo.berkeley.edu/logo.cgi>) was used  
177 to plot the logo of the motifs we identified.

## 178 8. Accession number

179 The sequencing data were deposited at GSA (Genome Sequence Archive) under  
180 BioProject accession no. PRJCA002499. The full protocol is available at

181 <https://github.com/caojiabao/VirPipeline>.

182

## 183 *Results*

### 184 1. Virome separation, enrichment and sequencing

185       Since metagenomic sequencing using fecal DNA usually results in only minor  
186 fraction of virome sequences, and most of reads will be either from bacteria or  
187 archaea, enrichment of viruses are necessary for studying virome in human gut  
188 samples. Thus, we have combined a series of enrichment methods including filtration  
189 and super centrifugation, to enrich for virus-like-particles (VLPs) in the fecal samples  
190 (Figure 1). After VLPs were isolated, additional DNase/RNase treatment were used to  
191 remove any potential free-DNA that were not virus-originated. The left-over  
192 DNA/RNA were quantified, and half of the nucleotides were directly subjected to  
193 ONT DNA library preparation, to profile the DNA virus abundances as well as  
194 methylation; and the other half were first subjected to RNA reverse transcription and  
195 then amplificaiton with short random primers, to have RNA virus sequenced as cDNA,  
196 and improve the chance of low-abundance DNA viruses to be detected in the  
197 sequencing results.

198       In our study as a primary investigation, we have first profiled five healthy  
199 volunteers' fecal samples with our protocol. Virus-like particle (VLP) fractions of five  
200 individuals were enriched, and raw DNA, as well as enriched cDNA/DNA were  
201 sequenced using ONT PromthION platform. With one flowcell, the ONT PromthION  
202 yielded a total of 8.2 Gb raw data, with a median of 1.7 Gb per sample in amplified

203 group; and 452 Mb raw data, with a median of 67 Mb per sample in raw DNA group  
204 (Table S1).

## 205 2. Virome composition in healthy individuals revealed by ONT sequencing

206 With sequencing reads from amplified cDNA/DNA results, we have mapped the  
207 ONT reads onto NCBI virus genomes and then analyzed the composition across five  
208 individuals. Consistent with other studies, our result from amplified cDNA/DNA of  
209 virome showed that bacteriophage families were the most frequently detected and  
210 accounted for the majority of the intestinal virome in number. The final catalogue of  
211 bacteriophages included the *Caudovirals* order (families *siphoviridae*, *podoviridae*),  
212 family *Inoviridae* and family *Microviridae*. Meanwhile, the eukaryotic CRESS-DNA  
213 viruses (family *genomoviridae*) was also detected (Figure 2; Table S2). Of special  
214 note was the presence of the RNA plant viruses including family *Virgaviridae* and  
215 *Alphaflexiviridae*, which showed good agreement with the findings of Shkoporov *et*  
216 *al* [54]. We observed that individual specificity is probably a feature of the faecal viral  
217 communities (Figure 2), which had been already demonstrated by several previous  
218 studies [55-57]. However, uncultured phage WW-nAnB strain 3 belonging to family  
219 *inoviridae* was detected to be presented in amplified cDNA/DNA among five  
220 individuals.

221 To characterize the potential biases regarding virome composition caused by  
222 amplification, we have compared the results from amplified cDNA/DNA and raw  
223 virome DNA. We could not compare raw RNA due to the fact that it is still not yet  
224 possible to multiplex RNA samples on ONT platforms. The relative abundance of

225 virus each virome was defined using relative proportion of each virus in terms of  
226 breadth of coverage on the assembled genome, similar to the definition of bacterial  
227 abundances in metagenomic studies. As expected, numbers of viruses are higher in  
228 amplified cDNA/DNA results, except for individual 2 and individual 5 who remained  
229 the same in terms of virus diversity. Further, abundances of the common existing  
230 viruses between raw DNA and amplified cDNA/DNA showed essential shifts in all  
231 five individuals, demonstrated detectable bias of reverse transcription as well as  
232 random amplification approach we adapted.

### 233 3. Virome assembly using ONT sequences

234 We next focused on the assembly of nanopore sequencing reads in five  
235 individuals. Assembler Canu was used to assemble the virome sequences separately  
236 from raw DNA and amplified cDNA/DNA groups into contigs, which yielded a total  
237 of 1564 contigs, with a median of 15 and 347 contigs per sample for raw DNA and  
238 amplified cDNA/DNA group. Average length of contigs from raw DNA group was  
239 longer than those from amplified cDNA/DNA group (Figure S1). The contigs vary  
240 largely in length, ranging from 1kb to 53kb (Figure S2). Consequently, we obtained  
241 the identity of certain contigs by matching with NCBI virus genomes, human gut  
242 virus database (GVD) [58] and NCBI nucleotide databases, however all with very low  
243 matching rate, suggesting a large collection of potentially novel genomes in our  
244 results (Table S3). Thus, Seeker was used to identify bacteriophages in amplified  
245 cDNA/DNA group. As a result, more than 50% of contigs per sample were identified  
246 as phages in amplified cDNA/DNA group. To characterize these bacteriophages, we

247 performed phages ORFs prediction and functional annotation, which yielded average  
248 7 ORFs per contig, ranging from 6 to 10, and average 3 ORFs per kb in length. In  
249 addition, ORFs repeated in two or more samples were validated by mapping to  
250 additional Illumina Hiseq sequencing data of metagenomics from same samples. We  
251 founded that over half ORFs (59.3%, 35/59) can be matched, indicating that these  
252 phages stably exist in our data (Table S4).

253 To examine the potential of covering full virus genome of long sequence  
254 produced from the oxford nanopore technology sequencing, we mapped raw reads to  
255 the contigs by canu separated from amplified cDNA/DNA and raw DNA. We  
256 aggregated and counted the length of raw reads aligned to the longest contigs of our  
257 choice from these two groups, and calculated the proportion of raw reads to contigs in  
258 length (Figure 3). In amplified VLP cDNA/DNA group, the max value of reads length  
259 in each sample was all more than 15% of contigs in length, the highest achieving 40%  
260 of contigs' full length in individual 4 (Figure S3). In raw DNA group, this proportion  
261 was higher in general, understandably due to the fact random amplification usually  
262 can not reach full length available in the raw DNA. A small number of reads were  
263 longer than the final contigs (Figure 3), for which the possible reason for this result  
264 was that some long reads were trimmed by canu during assembly, due to the lower  
265 quality of part of the sequences being abandoned by Canu.

#### 266 4. Viral epigenomic detection using ONT

267 In addition, methylation states of bases in DNA from reads can be detected  
268 directly by the Oxford Nanopore sequencing without extra laboratory techniques. In

269 this study, we have analyzed methylation signals on a few contigs that reached >10X  
270 coverage in raw DNA data in any of individuals. Methylation detection can not be  
271 carried out on cDNA and amplified DNA samples, for they will lose all the  
272 methylation signals. For the only one of the contigs (contig00000015) with known  
273 identity (Uncultured crAssphage), we detected in total 17 5mC and 120 6mA  
274 methylation sites covered in the 8kb genome (Figure 4). For 5mC and 6mA  
275 methylation, the nucleotide motifs YCHYTTACTWMRECT (e-value =  $1.2 \times 10^{-2}$ )  
276 and motif MADWDTWANADYYWW (e-value =  $2.5 \times 10^{-4}$ ) were identified  
277 respectively, with the methylated nucleotide highlighted in bold italic.

278 We have also found another 4 contigs with >10X coverage, but they do not have  
279 detectable relative in the databases we have searched against. They are potentially  
280 novel viruses and very unlikely to be bacteria or archaea, as our protocol has removed  
281 most of the none VLPs, and bacteria/archaea contigs would most likely have matches  
282 in the databases we have searched. They have 5mC methylation sites ranging from 0.3  
283 to 1.8 (% of genome) and 6mA methylation sites ranging from 0.7 to 2.5 (% of  
284 genome). The motifs for methylation are also extremely diverse, including  
285 NHHYYKGC DHNNN, WDWADDWCDWYNDDW, MNNNNTRCGBNNNND and  
286 HDNBYDVCVVVVNNH for 5mC, and WHWHNYDAHNNYYHTT,  
287 VNWWDWHAYBYNNNT, DRNVRKKABBNDNNN and  
288 NRNARNDA SYAHHNH for 6mA (Figure S4). Since most of the methylation studies  
289 are performed in eukaryotes, and only starting in bacteria with limited information  
290 available, it is yet difficult to compare the methylation profiles to understand its

291 underlying mechanisms.

292

### 293 *Discussion*

294 Our study combined the physical enrichment of VLPs in fecal samples,  
295 nucleotide amplification with the latest sequencing technology to establish a complete  
296 workflow of human gut virome profiling. With longer reads as well as richer  
297 information of additional epigenetic modifications, developments in sequencing  
298 technology could bring another round of revolution in metagenome as well as virome  
299 investigations. More importantly, despite the complex steps before sequencing were  
300 designed for maximum enrichment of VLPs, as well as removal of any DNA/RNA  
301 that were not of virus origin, and thus inevitably makes it relatively time-consuming,  
302 ONT could carry out sequencing and produce reads nearly simultaneously, making it  
303 possible to finish data generation from samples with five days of working time, and  
304 potentially even shorter if PromethION was run for less time, or data were analyzed  
305 during the course of being generated (real-time); comparatively, Illumina based  
306 platforms usually takes longer to generate enough read length for downstream  
307 analyses. There are several studies who have already utilized this property for fast  
308 pathogen detection in infectious diseases [59-61], there are cases of virome analysis  
309 that might require such time efficiency as well, and our protocol provides a feasible  
310 choice for virome studies that also requires short turn-around time.

311 We also compared the effect of reverse transcription and consequent  
312 amplifications on virome analysis with ONT. Such steps were used for several reasons.

313 Firstly, despite ONT is capable of sequencing DNA or RNA directly, multiplexing is  
314 still only possible for DNA libraries, and directly sequencing RNA is not yet  
315 cost-effective for virome analysis, while cDNA is a better alternative. Also, to utilize  
316 the capacity of sequencing on ONT, very high concentrations of DNA or RNA  
317 libraries are required to generate enough reads, yet this is also very difficult for  
318 virome DNA/RNA, who usually only reaches 10% of the required input from 1.5  
319 grams of fecal samples. Amplification does lead to higher amount of reads and  
320 enables detection of low abundance viruses to be detected, as revealed in our study;  
321 but it also leads to on average shorter reads and certain biases in the estimations of  
322 virus abundances, resulting from both affinity to random primers as well as PCR  
323 produced artefacts. Lastly, amplified cDNA/DNA loses all the methylation  
324 modifications on the viral nucleotides and prevents investigations into this potentially  
325 vital epigenetic information. Thus, future investigator needed to balance the pros and  
326 cons of amplification processes, and could take advantage of our approach of using  
327 both raw DNA (and/or RNA) and amplified cDNA/DNA for sequencing, gaining  
328 complementary information within the same sequencing run.

329 The profiles of virome in our studied individuals suggest a highly diverse virome,  
330 with only small amount of “core” viruses shared in between. This core could shrink  
331 even further with increasing number of individuals while the total diversity of viruses  
332 increase, which we plan to investigate in the future. We found phages making the  
333 majority part of the virome in the human gut, plus a few host viruses; while the  
334 mystery of plant RNA viruses is again present with ONT data, many reads achieving >



335 1kb in length, whether they are left-overs from our plant-based food, or rather human  
336 viruses with their closest relatives in the plant-associated viruses, still call for more  
337 investigations, especially within functional experiments. Our data suggest that there  
338 are potentially high number of unknown, novel viruses in the human gut, as our  
339 assembled contigs have very low rate of matches in the current databases; we consider  
340 those not likely to be contaminations due to our vigorous depletion of any non-VLPs  
341 and non-viral nucleotides, plus the fact that they do not have any match in nucleotide  
342 collection of NCBI either.

343 It's also the first time we demonstrate that the phage genomes are methylated via  
344 direct sequencing. RSV viruses and influenza viruses are known to have m6A  
345 methylation on their RNA genome [38, 62], detected with more complex methods  
346 with low throughput, while DNA phages (or other DNA viruses) are not yet studied to  
347 our knowledge. In *E. coli* and *C. difficile* it is shown that 6mA is the main form of  
348 DNA methylation, while eukaryotes usually lack this form, and in our results phages  
349 also have 6mA as the main form of DNA methylation [37, 63, 64]. Since DNA  
350 methylations play an important role in bacterial defense against phages, how phage  
351 genome becomes methylated, and the consequent impact on phage life cycle and  
352 interactions with bacterial hosts remain to be explored with dedicated studies. Besides,  
353 it remains possible the motifs of methylation between phage and bacteria are  
354 intrinsically linked, and provide additional information to determine the host range of  
355 phages; this would require increasing the current knowledge on epigenetics of both  
356 bacteria and phages in the future.

357

## 358 *Conclusions*

359 To summarize, we developed and pilot-tested a thorough protocol for human gut  
360 virome analysis using the latest ONT sequencing platform, and generated novel  
361 insights into the individuality, diversity of gut virome with new sequencing data. Our  
362 approach of course can be applied for other virome studies, including animal gut, soil  
363 and water virome etc., and accumulating both sequences as well as epigenetic  
364 information on those samples, have the long potential of opening up new directions in  
365 metagenomic, microbiological and medical researches.

366

## 367 *Acknowledgements*

368 This work was supported by the National Key Research and Development  
369 Program of China (grant number 2018YFC2000500), the Strategic Priority Research  
370 Program of the Chinese Academy of Sciences (grant number XDB29020000), and the  
371 National Natural Science Foundation of China (grant number 31771481 and  
372 91857101).

373

## 374 *References*

- 375 [1] Sender R, Fuchs S, Milo R. Are We Really Vastly Outnumbered? Revisiting the Ratio of  
376 Bacterial to Host Cells in Humans. *Cell* 2016;164(3):337-40.  
377 <https://doi.org/10.1016/j.cell.2016.01.013>.
- 378 [2] Falony G, Joossens M, Vieira-Silva S, Wang J, Darzi Y, Faust K, et al. Population-level

379 analysis of gut microbiome variation. *Science* 2016;352(6285):560-4. [https://](https://doi.org/10.1126/science.aad3503)  
380 [doi.org/10.1126/science.aad3503](https://doi.org/10.1126/science.aad3503).

381 [3] Moschen AR, Gerner RR, Wang J, Klepsch V, Adolph TE, Reider SJ, et al. Lipocalin 2  
382 Protects from Inflammation and Tumorigenesis Associated with Gut Microbiota Alterations. *Cell*  
383 *Host Microbe* 2016;19(4):455-69. <http://doi.org/10.1016/j.chom.2016.03.007>.

384 [4] Wang J, Thingholm LB, Skieceviciene J, Rausch P, Kummén M, Hov JR, et al. Genome-wide  
385 association analysis identifies variation in vitamin D receptor and other host factors influencing  
386 the gut microbiota. *Nat Genet* 2016;48(11):1396-406. [https:// doi.org/10.1038/ng.3695](https://doi.org/10.1038/ng.3695).

387 [5] Tschurtschenthaler M, Wang J, Fricke C, Fritz TMJ, Niederreiter L, Adolph TE, et al. Type I  
388 interferon signalling in the intestinal epithelium affects Paneth cells, microbial ecology and  
389 epithelial regeneration. *Gut* 2014;63(12):1921-31. [https:// doi.org/10.1136/gutjnl-2013-305863](https://doi.org/10.1136/gutjnl-2013-305863).

390 [6] Wang J, Chen L, Zhao N, Xu XZ, Xu YK, Zhu BL. Of genes and microbes: solving the  
391 intricacies in host genomes. *Protein Cell* 2018;9(5):446-61. [https://](https://doi.org/10.1007/s13238-018-0532-9)  
392 [doi.org/10.1007/s13238-018-0532-9](https://doi.org/10.1007/s13238-018-0532-9).

393 [7] Belkaid Y, Hand TW. Role of the Microbiota in Immunity and Inflammation. *Cell*  
394 2014;157(1):121-41. <https://doi.org/10.1016/j.cell.2014.03.011>.

395 [8] Wang ZN, Klipfell E, Bennett BJ, Koeth R, Levison BS, Dugar B, et al. Gut flora metabolism  
396 of phosphatidylcholine promotes cardiovascular disease. *Nature* 2011;472(7341):57-U82.  
397 <https://doi.org/10.1038/nature09922>.

398 [9] Ridaura VK, Faith JJ, Rey FE, Cheng JY, Duncan AE, Kau AL, et al. Gut Microbiota from  
399 Twins Discordant for Obesity Modulate Metabolism in Mice. *Science* 2013;341(6150):1079-U49.  
400 <https://doi.org/10.1126/science.1241214>.

- 401 [10] Jostins L, Ripke S, Weersma RK, Duerr RH, McGovern DP, Hui KY, et al. Host-microbe  
402 interactions have shaped the genetic architecture of inflammatory bowel disease. *Nature*  
403 2012;491(7422):119-24. <https://doi.org/10.1038/nature11582>.
- 404 [11] Giongo A, Gano KA, Crabb DB, Mukherjee N, Novelo LL, Casella G, et al. Toward defining  
405 the autoimmune microbiome for type 1 diabetes. *Isme J* 2011;5(1):82-91.  
406 <https://doi.org/10.1038/ismej.2010.92>.
- 407 [12] Cox LM, Blaser MJ. Pathways in Microbe-Induced Obesity. *Cell Metab* 2013;17(6):883-94.  
408 <https://doi.org/10.1016/j.cmet.2013.05.004>.
- 409 [13] Zhu WF, Gregory JC, Org E, Buffa JA, Gupta N, Wang ZN, et al. Gut Microbial Metabolite  
410 TMAO Enhances Platelet Hyperreactivity and Thrombosis Risk. *Cell* 2016;165(1):111-24.  
411 <https://doi.org/10.1016/j.cell.2016.02.011>.
- 412 [14] Wang ZN, Roberts AB, Buffa JA, Levison BS, Zhu WF, Org E, et al. Non-lethal Inhibition of  
413 Gut Microbial Trimethylamine Production for the Treatment of Atherosclerosis. *Cell*  
414 2015;163(7):1585-95. <https://doi.org/10.1016/j.cell.2015.11.055>.
- 415 [15] Zarate S, Taboada B, Yocupicio-Monroy M, Arias CF. Human Virome. *Arch Med Res*  
416 2017;48(8):701-16. <http://doi.org/10.1016/j.arcmed.2018.01.005>.
- 417 [16] Zou SM, Caler L, Colombini-Hatch S, Glynn S, Srinivas P. Research on the human virome:  
418 where are we and what is next. *Microbiome* 2016;4. <http://doi.org/10.1186/s40168-016-0177-y>.
- 419 [17] Scarpellini E, Ianiro G, Attili F, Bassanelli C, De Santis A, Gasbarrini A. The human gut  
420 microbiota and virome: Potential therapeutic implications. *Digest Liver Dis* 2015;47(12):1007-12.  
421 <http://doi.org/10.1016/j.dld.2015.07.008>.
- 422 [18] Liu L, Gong T, Tao WY, Lin BL, Li C, Zheng XS, et al. Commensal viruses maintain

423 intestinal intraepithelial lymphocytes via noncanonical RIG-I signaling. *Nat Immunol*  
424 2019;20(12):1681-1691. <https://doi.org/10.1038/s41590-019-0513-z>.

425 [19] Reyes A, Haynes M, Hanson N, Angly FE, Heath AC, Rohwer F, et al. Viruses in the faecal  
426 microbiota of monozygotic twins and their mothers. *Nature* 2010;466(7304):334-U81.  
427 <https://doi.org/10.1038/nature09199>.

428 [20] Virgin HW. The Virome in Mammalian Physiology and Disease. *Cell* 2014;157(1):142-50.  
429 <http://doi.org/10.1016/j.cell.2014.02.032>.

430 [21] Zhang T, Breitbart M, Lee WH, Run JQ, Wei CL, Soh SWL, et al. RNA viral community in  
431 human feces: Prevalence of plant pathogenic viruses. *Plos Biol* 2006;4(1):108-18.  
432 <https://doi.org/ARTN e310.1371/journal.pbio.0040003>.

433 [22] Shi Y, Mu LB. An expanding stage for commensal microbes in host immune regulation. *Cell*  
434 *Mol Immunol* 2017;14(4):339-48. <https://doi.org/10.1038/cmi.2016.64>.

435 [23] Handley SA. The virome: a missing component of biological interaction networks in health  
436 and disease. *Genome Med* 2016;8. <http://doi.org/10.1186/s13073-016-0287-y>.

437 [24] Mukhopadhy I, Segal JP, Carding SR, Hart AL, Hold GL. The gut virome: the 'missing link'  
438 between gut bacteria and host immunity? *Ther Adv Gastroenter* 2019;12.  
439 <http://doi.org/10.1177/1756284819836620>.

440 [25] Foca A, Liberto MC, Quirino A, Marascio N, Zicca E, Pavia G. Gut Inflammation and  
441 Immunity: What Is the Role of the Human Gut Virome? *Mediat Inflamm* 2015.  
442 <http://doi.org/10.1155/2015/326032>.

443 [26] Norman JM, Handley SA, Baldrige MT, Droit L, Liu CY, Keller BC, et al. Disease-Specific  
444 Alterations in the Enteric Virome in Inflammatory Bowel Disease. *Cell* 2015;160(3):447-60.

445 <http://doi.org/10.1016/j.cell.2015.01.002>.

446 [27] Manrique P, Dills M. Young MJ. The Human Gut Phage Community and Its Implications for  
447 Health and Disease. *Viruses-Basel* 2017;9(6). <http://doi.org/10.3390/v9060141>.

448 [28] Yang JY, Kim MS, Kim E, Cheon JH, Lee YS, Kim Y, et al. Enteric Viruses Ameliorate Gut  
449 Inflammation via Toll-like Receptor 3 and Toll-like Receptor 7-Mediated Interferon-beta  
450 Production. *Immunity* 2016;44(4):889-900. <https://doi.org/10.1016/j.immuni.2016.03.009>.

451 [29] Broggi A, Tan Y, Granucci F, Zanoni I. IFN-lambda suppresses intestinal inflammation by  
452 non-translational regulation of neutrophil function. *Nat Immunol* 2017;18(10):1084-1093.  
453 <https://doi.org/10.1038/ni.3821>.

454 [30] Deamer D, Akeson M, Branton D. Three decades of nanopore sequencing. *Nat Biotechnol*  
455 2016;34(5):518-24. <https://doi.org/10.1038/nbt.3423>.

456 [31] Laszlo AH, Derrington IM, Ross BC, Brinkerhoff H, Adey A, Nova IC, et al. Decoding long  
457 nanopore sequencing reads of natural DNA. *Nat Biotechnol* 2014;32(8):829-33.  
458 <https://doi.org/10.1038/nbt.2950>.

459 [32] Schatz MC. Nanopore sequencing meets epigenetics. *Nat Methods* 2017;14(4):347-8.  
460 <https://doi.org/10.1038/nmeth.4240>.

461 [33] Low DA, Weyand NJ, Mahan MJ. Roles of DNA adenine methylation in regulating bacterial  
462 gene expression and virulence. *Infect Immun* 2001;69(12):7197-204. <https://doi.org/10.1128/iai.69.12.7197-7204.2001>.

464 [34] Casadesus J, Low D. Epigenetic gene regulation in the bacterial world. *Microbiol Mol Biol R*  
465 2006;70(3):830-56. <https://doi.org/10.1128/Mmbr.00016-06>.

466 [35] Oliveira PH, Touchon M, Rocha EPC. Regulation of genetic flux between bacteria by

467 restriction-modification systems. *P Natl Acad Sci USA* 2016;113(20):5658-63.

468 <https://doi.org/10.1073/pnas.1603257113>.

469 [36] Cohen NR, Ross CA, Jain S, Shapiro RS, Gutierrez A, Belenky P, et al. A role for the

470 bacterial GATC methylome in antibiotic stress survival. *Nat Genet* 2016;48(5):581-6.

471 <https://doi.org/10.1038/ng.3530>.

472 [37] Oliveira PH, Ribis JW, Garrett EM, Trzilova D, Kim A, Sekulovic O, et al. Epigenomic

473 characterization of *Clostridioides difficile* finds a conserved DNA methyltransferase that mediates

474 sporulation and pathogenesis. *Nat Microbiol* 2019; <https://doi.org/10.1038/s41564-019-0613-4>.

475 [38] Xue MG, Zhao BS, Zhang ZJ, Lu MJ, Harder O, Chen P, et al. Viral N-6-methyladenosine

476 upregulates replication and pathogenesis of human respiratory syncytial virus. *Nat Commun*

477 2019;10 <https://doi.org/ARTN.459510.1038/s41467-019-12504-y>.

478 [39] Thurber RV, Haynes M, Breitbart M, Wegley L, Rohwer F. Laboratory procedures to generate

479 viral metagenomes. *Nat Protoc* 2009;4(4):470-83. <https://doi.org/10.1038/nprot.2009.10>.

480 [40] Ge XY, Li Y, Yang XL, Zhang HJ, Zhou P, Zhang YZ, et al. Metagenomic Analysis of Viruses

481 from Bat Fecal Samples Reveals Many Novel Viruses in Insectivorous Bats in China. *J Virol*

482 2012;86(8):4620-30. <https://doi.org/10.1128/Jvi.06671-11>.

483 [41] Froussard P. rPCR: a powerful tool for random amplification of whole RNA sequences.

484 *Genome Research* 1993;2:185-90. <https://doi.org/10.1101/gr.2.3.185>.

485 [42] Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and

486 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res*

487 2017;27(5):722-36. <https://doi.org/10.1101/gr.215087.116>.

488 [43] Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*

- 489 2018;34(18):3094-100. <https://doi.org/10.1093/bioinformatics/bty191>.
- 490 [44] Auslander N, Gussow AB, Benler S, Wolf YI, Koonin EV. Seeker: Alignment-free  
491 identification of bacteriophage genomes by deep learning. *bioRxiv* 2020;2020.04.04.025783.  
492 <https://doi.org/10.1101/2020.04.04.025783>.
- 493 [45] Ecale Zhou CL, Malfatti S, Kimbrel J, Philipson C, McNair K, Hamilton T, et al.  
494 multiPhATE: bioinformatics pipeline for functional annotation of phage isolates. *Bioinformatics*  
495 (Oxford, England) 2019;35(21):4402-4. <https://doi.org/10.1093/bioinformatics/btz258>.
- 496 [46] McNair K, Zhou C, Dinsdale EA, Souza B, Edwards RA. PHANOTATE: a novel approach to  
497 gene identification in phage genomes. *Bioinformatics* (Oxford, England) 2019;35(22):4537-42.  
498 <https://doi.org/10.1093/bioinformatics/btz265>.
- 499 [47] Graziotin AL, Koonin EV, Kristensen DM. Prokaryotic Virus Orthologous Groups (pVOGs):  
500 a resource for comparative genomics and protein family annotation. *Nucleic acids research*  
501 2017;45(D1):D491-d8. <https://doi.org/10.1093/nar/gkw975>.
- 502 [48] Johnson LS, Eddy SR, Portugaly E. Hidden Markov model speed heuristic and iterative  
503 HMM search procedure. *BMC Bioinformatics* 2010;11(1):431.  
504 <https://doi.org/10.1186/1471-2105-11-431>.
- 505 [49] Edgar R. Taxonomy annotation and guide tree errors in 16S rRNA databases. *PeerJ*  
506 2018;6e5030. <https://doi.org/10.7717/peerj.5030>.
- 507 [50] Langmead B, Wilks C, Antonescu V, Charles R. Scaling read aligners to hundreds of threads  
508 on general-purpose processors. *Bioinformatics* (Oxford, England) 2019;35(3):421-32.  
509 <https://doi.org/10.1093/bioinformatics/bty648>.
- 510 [51] Marcus Stoiber JQ, Rob Egan, Ji Eun Lee, Susan Celniker, Robert K. Neely, Nicholas Loman,



511 Len A Pennacchio, James Brown. De novo Identification of DNA Modifications Enabled by  
512 Genome-Guided Nanopore Signal Processing. bioRxiv 2017; <https://doi.org/10.1101/094672>.

513 [52] Robinson JT, Thorvaldsdottir H, Wenger AM, Zehir A, Mesirov JP. Variant Review with the  
514 Integrative Genomics Viewer. *Cancer Res* 2017;77(21):e31-e4.  
515 <https://doi.org/10.1158/0008-5472.CAN-17-0337>.

516 [53] Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, et al. MEME SUITE: tools  
517 for motif discovery and searching. *Nucleic Acids Res* 2009;37(Web Server issue):W202-8.  
518 <https://doi.org/10.1093/nar/gkp335>.

519 [54] Shkoporov AN, Clooney AG, Sutton TDS, Ryan FJ, Daly KM, Nolan JA, et al. The human  
520 gut virome is highly diverse, stable and individual-specific. *Cell Host & Microbe* 2019;  
521 <https://doi.org/10.1101/657528>.

522 [55] Moreno-Gallego JL, Chou S-P, Di Rienzi SC, Goodrich JK, Spector TD, Bell JT, et al.  
523 Virome Diversity Correlates with Intestinal Microbiome Diversity in Adult Monozygotic Twins.  
524 *Cell Host & Microbe* 2019;25(2):261-72.e5. <https://doi.org/10.1016/j.chom.2019.01.019>.

525 [56] Garmaeva S, Sinha T, Kurilshikov A, Fu J, Wijmenga C, Zhernakova A. Studying the gut  
526 virome in the metagenomic era: challenges and perspectives. *BMC Biol* 2019;17(1):84.  
527 <https://doi.org/10.1186/s12915-019-0704-y>.

528 [57] Minot S, Sinha R, Chen J, Li H, Keilbaugh SA, Wu GD, et al. The human gut virome:  
529 inter-individual variation and dynamic response to diet. *Genome Res* 2011;21(10):1616-25.  
530 <https://doi.org/10.1101/gr.122705.111>.

531 [58] Gregory AC ZO, Howell A, Bolduc B, Sullivan MB. The human gut virome database.  
532 bioRxiv 2019; <https://doi.org/10.1101/655910>.

- 533 [59] Gardy JL, Loman NJ. Towards a genomics-informed, real-time, global pathogen surveillance  
534 system. *Nat Rev Genet* 2018;19(1):9-20. <https://doi.org/10.1038/nrg.2017.88>.
- 535 [60] Depledge DP, Srinivas KP, Sadaoka T, Bready D, Mori Y, Placantonakis DG, et al. Direct  
536 RNA sequencing on nanopore arrays redefines the transcriptional complexity of a viral pathogen.  
537 *Nat Commun* 2019;10 <https://doi.org/10.1038/s41467-019-08734-9>.
- 538 [61] Charalampous T, Kay GL, Richardson H, Aydin A, Baldan R, Jeanes C, et al. Nanopore  
539 metagenomics enables rapid clinical diagnosis of bacterial lower respiratory infection. *Nat*  
540 *Biotechnol* 2019;37(7):783-792. <https://doi.org/10.1038/s41587-019-0156-5>.
- 541 [62] Courtney DG, Kennedy EM, Dumm RE, Bogerd HP, Tsai K, Heaton NS, et al.  
542 Epitranscriptomic Enhancement of Influenza A Virus Gene Expression and Replication. *Cell Host*  
543 *Microbe* 2017;22(3):377-386. <https://doi.org/10.1016/j.chom.2017.08.004>.
- 544 [63] Beaulaurier J, Schadt EE, Fang G. Deciphering bacterial epigenomes using modern  
545 sequencing technologies. *Nat Rev Genet* 2019;20(3):157-72.  
546 <https://doi.org/10.1038/s41576-018-0081-3>.
- 547 [64] Fang G, Munera D, Friedman DI, Mandlik A, Chao MC, Banerjee O, et al. Genome-wide  
548 mapping of methylated adenine residues in pathogenic *Escherichia coli* using single-molecule  
549 real-time sequencing (vol 30, pg 1232, 2012). *Nat Biotechnol* 2013;31(6):566-9.  
550 <https://doi.org/10.1038/nbt0613-565f>.

551 *Figure Legends:*

552 **Figure 1. An integrated novel workflow for enrichment of virus-like particles**

553 **(VLPs), extraction of nucleic acids and ONT sequencing.** The complete workflow

554 consists of four fragments: (1) Washing and filtration of faecal samples using

555 PBS/Sterile and PVDF membrane including step 1-3; (2) Precipitation of VLPs

556 including step 4-5; (3) Extraction, amplification and purification of viral nucleic acids

557 including step 6-8; (4) Construction of library and ONT sequencing including step

558 9-10.

559 **Figure 2. Composition and relative abundance of viruses in each individual.**

560 Different color of the bar represents different viral species or strains. The line between

561 unamplified and amplified group in each individual represents common virus species

562 between two groups.

563 **Figure 3. Proportion of raw reads to contigs in length.** The proportion of raw reads

564 to contigs in length is shown by combination of scatter diagram, boxplot and violin

565 chart. Different color represents five different individuals. (A) raw DNA group, (B)

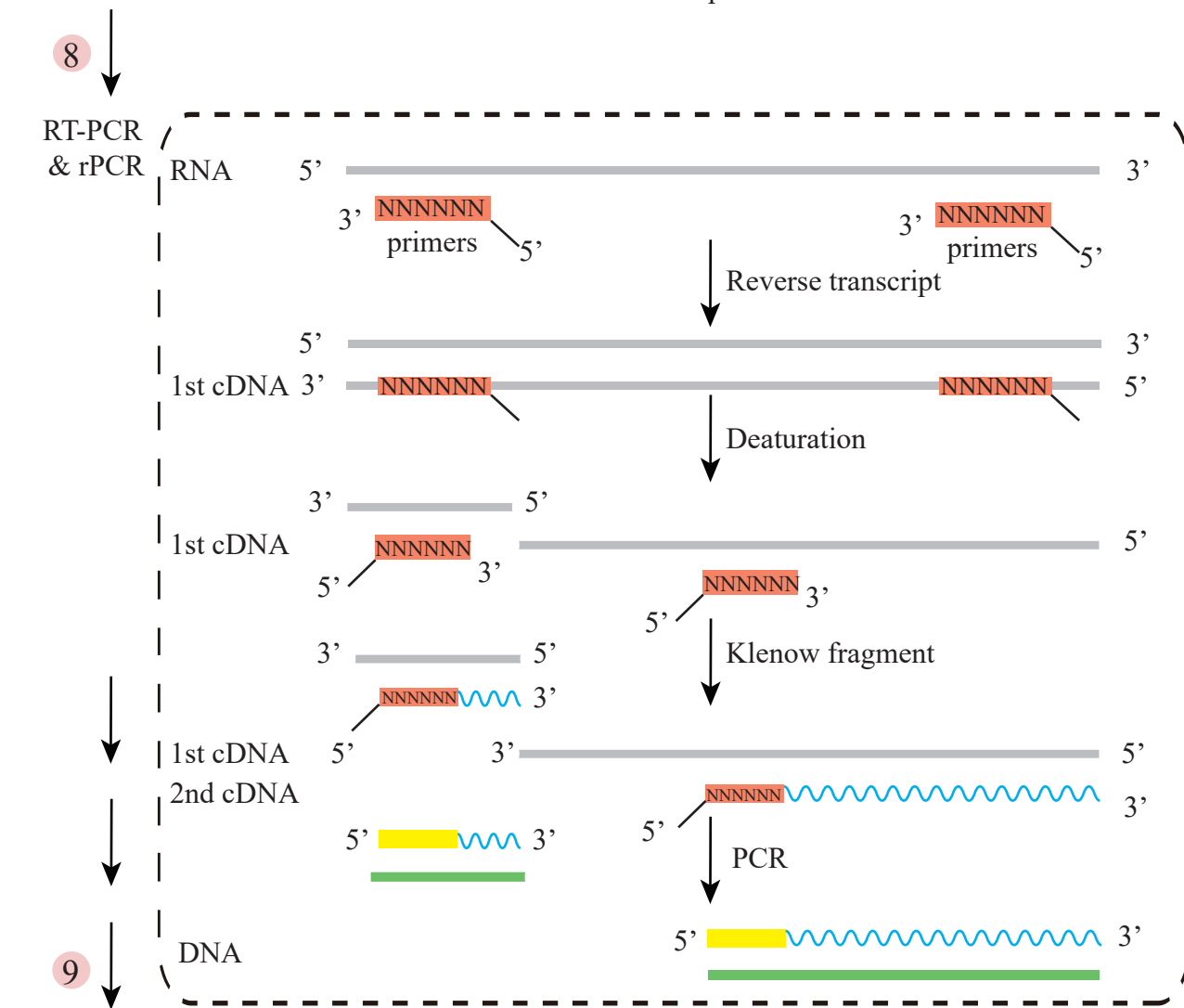
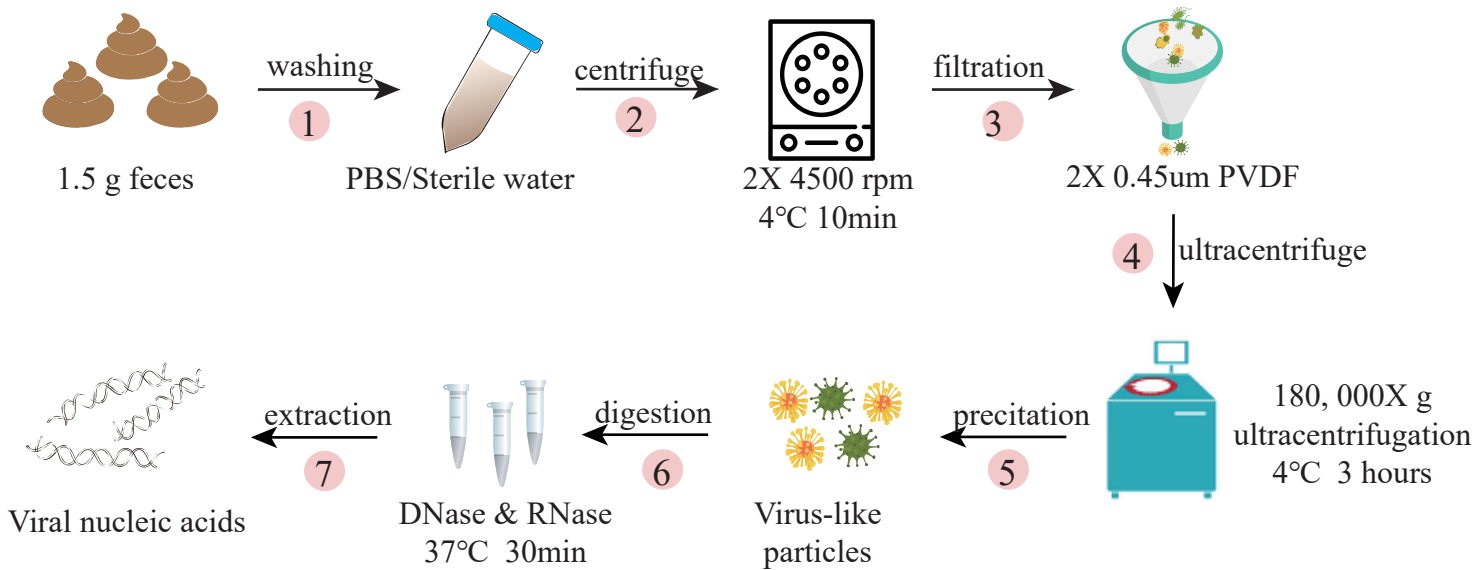
566 amplified DNA/cDNA group.

567 **Figure 4. Different methylation sites identification and motif recognition of viral**

568 **contig.** (A) the distribution of methylation sites (5mC and 6mA) in the

569 contig00000015; (B) 5mC motif; (C) 6mA motif.

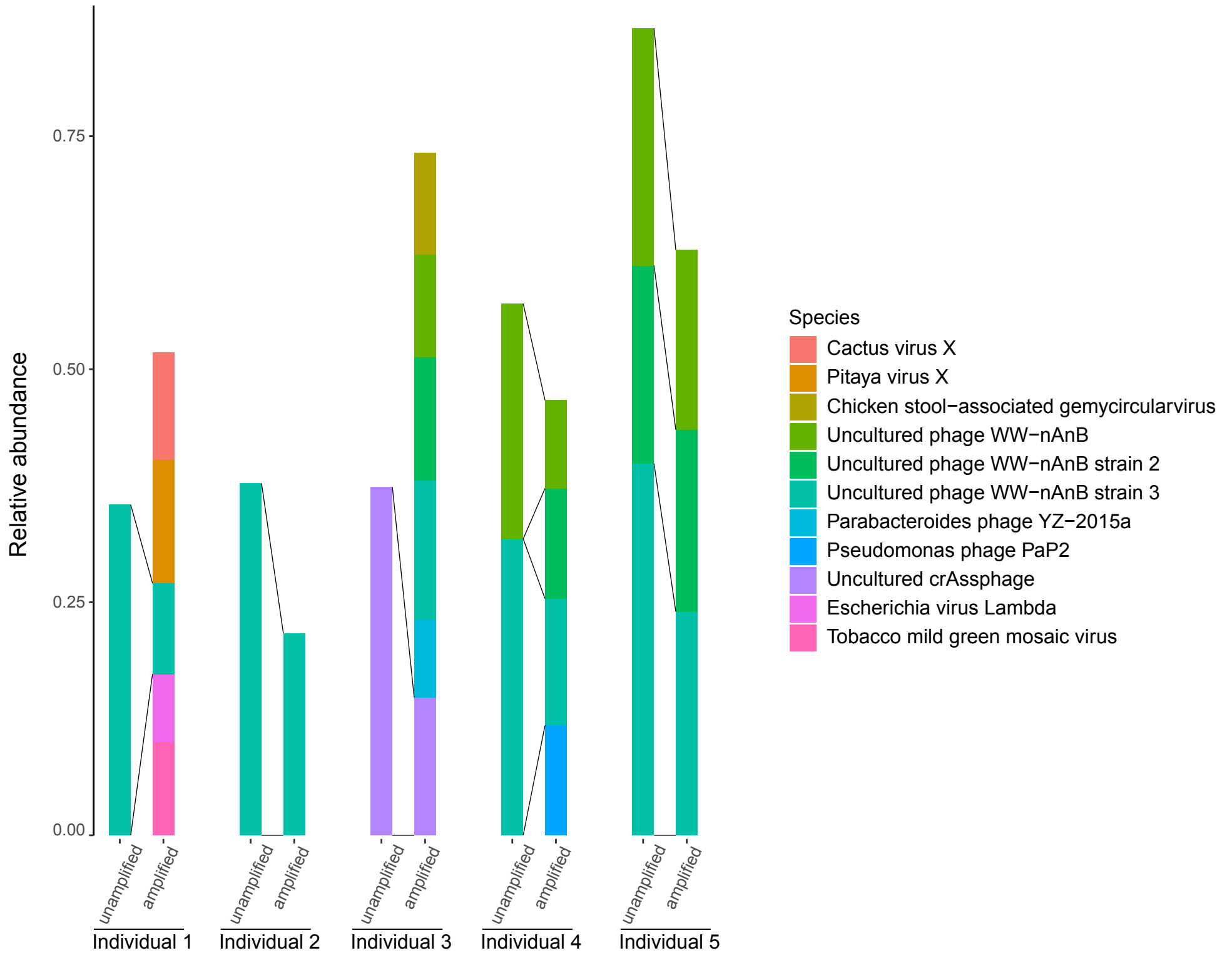
570



Library  
preparation

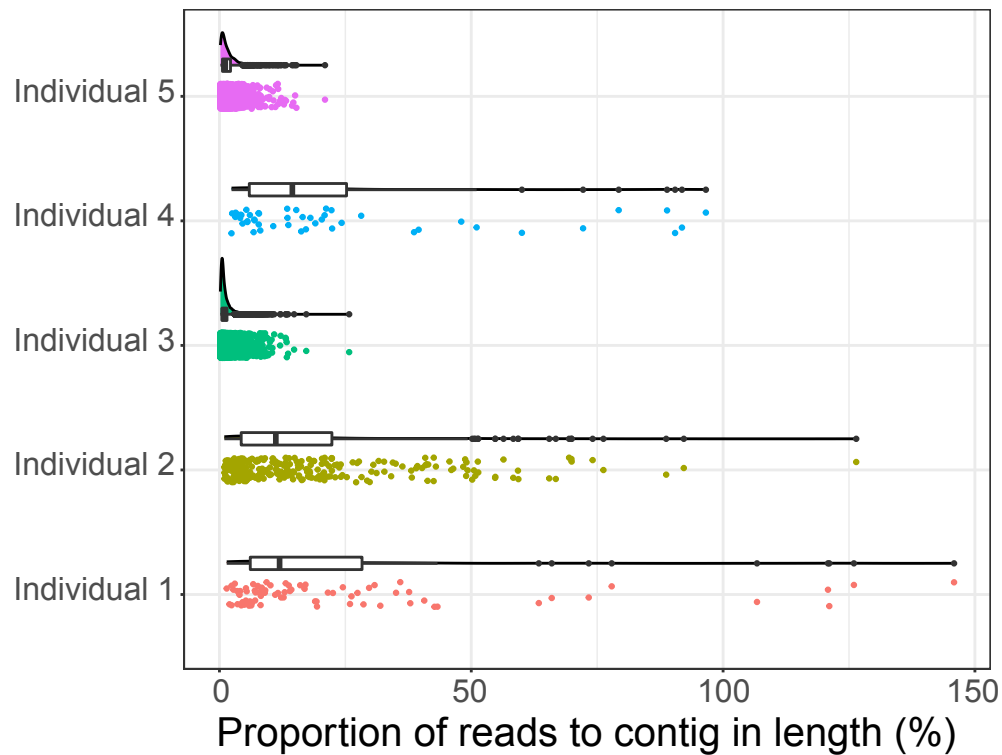
sequencing  
10

PromethIon



A

## Raw DNA group



B

## Amplified DNA/cDNA group

