

Supplementary Appendix

1

2 **Methods**

3 **Sample size estimation**

4 The main objective of this exploratory study was to investigate the difference in
5 the mucosal virome between healthy subjects and patients with ulcerative
6 colitis (UC). In a previous fecal virome study on healthy individuals and patients
7 with inflammatory bowel disease [1], they found a significant higher abundance
8 of *Caudovirales* in the feces of UC subjects than controls (0.6 ± 0.3 versus
9 0.3 ± 0.2 , mean \pm SD, in relative abundance). We herein hypothesize that the
10 mucosal *Caudovirales* levels in UC is also higher than in controls, and the
11 relative abundance in UC and control mucosa are 0.6 ± 0.2 and 0.3 ± 0.2
12 respectively. Based on these assumptions, we estimated that 16 number of
13 subjects in each group would achieve 80 percent power at 5 percent significant
14 level. It follows that our current sample size is of sufficient power to detect this
15 difference.

16

17 **Mucosal virus-like particles (VLPs) enrichment**

18 Rectal biopsy was digested in 1 mL digest buffer (1 mg/mL Collagenase D, 1U
19 Baseline DNase I, PBS pH7.5) at 37°C for 1 hour, with intermittent intensive
20 vortex every 20 minutes. Biopsy suspension was then cleared by
21 centrifugation at 5,000 x g for 5 minutes to remove debris and cells. Clarified
22 suspension was passed through one 0.22 μ m filter to remove residual host
23 and bacteria cells. Samples were treated with lysozyme (1 mg/ml at 37 °C for
24 30 min) followed by chloroform (0.2x volume at RT for 10 min) to degrade any
25 remaining bacterial and host cell membranes. Non-virus protected DNA was
26 degraded by treatment with a DNase cocktail (10U Tubro DNaseI (Ambion),
27 1U Baseline zero DNase (Epicenter)) followed by heat inactivation of DNases
28 at 65 °C for 10 min. VLPs were lysed (4% SDS plus 38 mg/ml Proteinase K at

29 56 °C for 20 min), treated with CTAB (2.5% CTAB plus 0.5 M NaCl at 65°C for
30 10 min), and nucleic acid was extracted with phenol:chloroform pH 8.0
31 (Invitrogen). The aqueous fraction was washed once with an equal volume of
32 chloroform, purified and concentrated on a column (DNA Clean &
33 Concentrator™-5, Zymo Research). VLP DNA was amplified for 4 hr using
34 Phi29 polymerase (GenomiPhi V2 kit, GE Healthcare) prior to sequencing. To
35 reduce amplification bias, four independent reactions were performed for each
36 sample and pooled together afterwards.

37

38 **Sequence Processing and Quality Control**

39 Raw reads were filtered by SOAPnuke (v 1.5.3) (<http://soap.genomics.org.cn/>)
40 developed by BGI as follows: (i) adaptors removed, (ii) read removed if N base
41 is more than 3% of the read, (iii) read removed if bases with quality low than 20
42 were more than 40% of read, (iv) all duplicates removed. Human sequences
43 were removed from the quality-trimmed dataset by DeconSeq (v 0.4.3) with
44 default parameters and the human reference GRCh38 [2].

45

46 **De Novo Contig Assembly and Taxonomy Annotation**

47 In view of that only part of the reads can be assembled in metagenomic
48 dataset, all other potential viral reads not able to assemble into contigs are
49 considered to present in very low abundances [3] and not taken into analysis.
50 Hence, we assembled those abundant metagenomic viral contigs and
51 classified viral reads. Contigs were assembled using the IDBA (v 1.1.1) [4],
52 using maximum kmer length 120, with a minimum contig length of 1,000 bp.
53 The assembled contigs were clustered at a 95% identity level using CD-HIT [5]
54 to generate a unique contig consortium. In total, 178,155 unique contigs were
55 assembled.

56 Open Reading Frame (ORF) were predicted and extracted from contigs using
57 the Glimmer3 toolkit (v 3.02) [6] and a minimum length threshold of 100 amino
58 acids. The translated amino acid sequences of predicted ORFs from the VLP
59 contigs were matched against the standard subset of the standalone entire
60 UniProt TrEMBL database as of June 15, 2018, that contained only virus and
61 phage reference proteins, using blastx ($e < 10^{-5}$) as described previously [7].
62 Each contig was assigned taxonomy based on the most abundant taxa
63 contained within that contig using a voting system as described previously for
64 virus taxonomic assignment at different taxon levels [8, 9]. The voting system
65 first annotated each ORF of a contig of interest with the best-hit virus taxonomy.
66 It then compared all of the taxonomic assignments of the ORFs within the
67 contig of interest, and annotated the contig with the majority ORF assignment.
68 Contigs with less than one ORF per 10kb were not assigned taxonomy as this
69 suggests a contig of only limited similarity [8]. Contigs without a majority ORF
70 taxonomic assignment due to ties of multiple major taxa were assigned as
71 having multiple possible taxonomic annotations. Among the assembled all
72 contigs, 51,241 contigs were considered as viral contig based on ORF binning.
73 To preclude false positive contigs classified as viral contig, all these viral
74 contigs were further queried against NCBI nonredundant nucleotide database
75 using BBSketch with the following settings: ambiguous=radom, qtrim=lr,
76 minid=0.97. Viral contigs mapped to organisms from other kingdoms with
77 *ANI*>0.97, most of which were bacteria, were excluded. Eventually, 30,853
78 contigs were remained and regarded as *bona fide* viral, which were then used
79 as our curated viral contig database. Some contigs shared the same
80 taxonomic identity, the contig table was therefore collapsed by taxonomic
81 identity, where the contig abundances were summed if they shared identity.
82 Mucosal virome metacommunity clustering analysis was performed according
83 to the previously reported method [10, 11], based on the virome species-level

84 community structures. Virus core species, common species and unique
85 species were defined when the viral species meet the below criteria: 1) core
86 species was deemed when it exists in >50% of all study subjects with
87 concordant RPKM>10 in any individual; 2) common species was deemed
88 when it exists in 20-50% of all study subjects with concordant RPKM>10 in any
89 individual; and 3) unique species was deemed when it exists in <20% of all
90 study subjects. The presence and percentage ratio were evaluated in each
91 individual and compared between healthy controls, UC Metacommunity 1 and
92 UC metacommunity 2 subjects.

93

94 **Non-metric multidimensional scaling analysis**

95 The difference in mucosal viral community structures between controls and UC
96 was performed via NMDS (Non-metric multidimensional scaling analysis) plot
97 based upon Bray-Curtis dissimilarities among all subjects.

98

99 **DESeq, Random Forest and LEfSe linear discriminant analysis**

100 To compare differences in the configurations of mucosal virome and bacterial
101 microbiome as well as the functions of virome between UC patients and
102 healthy household controls, between mucosal virome metacommunities,
103 Differential analyses were performed. *DESeq* and *Random Forest* were
104 performed in R via *DEseq* and *randomForest* package respectively. *Lefse*
105 analyses were performed on the Huttenhower lab Galaxy server
106 (<http://huttenhower.sph.harvard.edu/galaxy/>). *DB-RDA* analysis was also
107 performed in R to delineate the effect of medication and health on mucosal
108 virome configuration.

109

110 **Virome function analysis**

111 Virome functions were classified through annotating all viral-contig derived
112 reads via *HUMANN2 v0.9.4*. Predicted functions were collapsed by Gene
113 ontology terms and Pfam protein family identities, with abundance values
114 expressed in RPK (reads per kilobase). To establish the presence or absence
115 a function within a sample, a stringent RPK threshold value > 10 was used to
116 define as present.

117

118 **Mucosal bacterial DNA Extraction**

119 Bacterial DNA was extracted from rectal biopsies using Maxwell RSC Tissue
120 DNA kit (Promega, Madison, Wisconsin) according to the protocol. Briefly, 1
121 rectal biopsy was added to the Lysis Tube with 0.5mm plus 0.1 mm beads
122 filled with 300 μ l Lysis solution (500mM NaCl, 50mM Tris-HCl pH 8.0, 50mM
123 EDTA, 4% sodium dodecyl sulphate SDS), and then homogenized at
124 maximum speed for ≥ 3 minutes. The lysates were incubated at 70°C for 10
125 minutes followed by centrifugation at $\geq 16,000 \times g$ for 5 minutes. The
126 supernatant was transferred to Maxwell RSC Tissue tools and extracted
127 through Maxwell RSC Instrument (Promega, Madison, Wisconsin).

128

129 **16S rRNA sequencing and quality control**

130 The final fecal DNA samples were sequenced on the Illumina HiSeq 2500
131 platform (V4 region, 2 X 250 bp). Quality control and data analysis were
132 implemented in mothur (v 1.38.0) as previously described [12]. Any sequences
133 with ambiguous bases and anything longer than 275 bp were removed, and
134 aligned against the non-redundant Greengenes database (v 13.8) [13] using
135 the NAST algorithm. Any sequences that failed to align with the V3-4 region
136 were discarded. The remaining sequences were trimmed to the same
137 alignment coordinates over which they fully overlapped, followed by removal of
138 homopolymers and detection for the presence of chimeras by UChime.

139

140 **16S rRNA sequencing data analysis**

141 The resulting sequences were classified against the Greengenes database
142 and annotated with deepest level taxa represented by pseudo-bootstrap
143 confidence scores of at least 80% averaged over 1,000 iterations of the naive
144 Bayesian classifier. Any sequences that were classified as either being
145 originated from archaea, eukarya, chloroplasts, mitochondria, or unknown
146 kingdoms, were removed. The annotated sequences were assigned to
147 phylotypes according to their consensus taxonomy with which at least 80% of
148 the sequences agreed. Closed reference operational taxonomic units (OTUs)
149 sharing 97% identity were clustered as well and assigned taxonomy according
150 to the Greengenes database. Bacterial diversity and richness were calculated
151 in R via *phyloseq* package. Bacteria abundance tables at the phylum, family
152 and genus levels were also generated and compared between health and UC
153 mucosa. Lefse analysis was performed to define bacterial taxa associated with
154 UC and healthy controls.

155

156

157

158 Reference

- 159 1 Norman JM, Handley SA, Baldrige MT, Droit L, Liu CY, Keller BC, *et al.*
160 Disease-Specific Alterations in the Enteric Virome in Inflammatory Bowel Disease.
161 *Cell* 2015;**160**:447-60.
- 162 2 Schmieder R, Edwards R. Fast Identification and Removal of Sequence
163 Contamination from Genomic and Metagenomic Datasets. *PloS one* 2011;**6**.
- 164 3 Paez-Espino D, Pavlopoulos GA, Ivanova NN, Kyrpides NC. Nontargeted virus
165 sequence discovery pipeline and virus clustering for metagenomic data. *Nat Protoc*
166 2017;**12**:1673.
- 167 4 Peng Y, Leung HCM, Yiu SM, Chin FYL. IDBA - A Practical Iterative de Bruijn Graph
168 De Novo Assembler. *Lect N Bioinform* 2010;**6044**:426-40.
- 169 5 Fu LM, Niu BF, Zhu ZW, Wu ST, Li WZ. CD-HIT: accelerated for clustering the
170 next-generation sequencing data. *Bioinformatics* 2012;**28**:3150-2.
- 171 6 Delcher AL, Bratke KA, Powers EC, Salzberg SL. Identifying bacterial genes and
172 endosymbiont DNA with Glimmer. *Bioinformatics* 2007;**23**:673-9.

173 7 Hannigan GD, Meisel JS, Tyldsley AS, Zheng Q, Hodkinson BP, SanMiguel AJ, *et al.*
174 The Human Skin Double-Stranded DNA Virome: Topographical and Temporal
175 Diversity, Genetic Enrichment, and Dynamic Associations with the Host Microbiome.
176 *mBio* 2015;**6**.

177 8 Minot S, Bryson A, Chehoud C, Wu GD, Lewis JD, Bushman FD. Rapid evolution
178 of the human gut virome. *Proceedings of the National Academy of Sciences of the*
179 *United States of America* 2013;**110**:12450-5.

180 9 Hannigan GD, Meisel JS, Tyldsley AS, Zheng Q, Hodkinson BP, SanMiguel AJ, *et al.*
181 The human skin double-stranded DNA virome: topographical and temporal diversity,
182 genetic enrichment, and dynamic associations with the host microbiome. *mBio*
183 2015;**6**:e01578-15.

184 10 Costea PI, Hildebrand F, Manimozhiyan A, Bäckhed F, Blaser MJ, Bushman FD, *et*
185 *al.* Enterotypes in the landscape of gut microbial community composition. *Nat*
186 *Microbiol* 2018;**3**:8.

187 11 Arumugam M, Raes J, Pelletier E, Le Paslier D, Yamada T, Mende DR, *et al.*
188 Enterotypes of the human gut microbiome. *Nature* 2011;**473**:174.

189 12 Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, *et al.*
190 Introducing mothur: Open-Source, Platform-Independent, Community-Supported
191 Software for Describing and Comparing Microbial Communities. *Appl Environ Microb*
192 2009;**75**:7537-41.

193 13 McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A, *et al.* An
194 improved Greengenes taxonomy with explicit ranks for ecological and evolutionary
195 analyses of bacteria and archaea. *Isme J* 2012;**6**:610-8.
196
197

198 **Supplementary Figure 1. Mucosal virobiota composition at the order,**
199 **genus and species levels in healthy individuals and UC subjects. a.**
200 Mucosal viral abundance variation boxplot in healthy individuals. **b.** Mucosal
201 viral abundance variation boxplot in UC subjects. Only the most abundant viral
202 taxa are shown for clarity. For all box plots, the boxes extend from the 1st to
203 3rd quartile (25th to 75th percentile), with the median depicted by a horizontal
204 line.

205

206 **Supplementary Figure 2. Mucosal virobiota diversity alteration at the**
207 **species level in UC.** Comparison of the mucosal virobiota α diversity based
208 on Shannon diversity (a) and Chao1 richness (b) index in the mucosa of
209 controls and UC subjects. Statistical significance was determined by *t* test, $*P$
210 < 0.05 , $**P < 0.01$. For box plots, the boxes extend from the 1st to 3rd quartile
211 (25th to 75th percentile), with the median depicted by a horizontal line.

212

213 **Supplementary Figure 3. Mucosal inflammation was linked to virome**
214 **alterations in UC. a,** Comparison of the virus abundance in the rectal mucosa
215 of controls and UC subjects with and without inflammation, at the viral contig
216 level. Total virus abundance was calculated as RPKM sum of all viral contigs
217 recruited reads normalized by sequence depth in each subject.
218 Between-group difference was determined by Mann-Whitney test, $*P < 0.05$. **b,**
219 Comparison of the mucosal virome α diversity based on Shannon diversity,
220 evenness, Chao1 richness index in the mucosa of controls and UC subjects
221 with and without inflammation. Statistical significance was determined by
222 one-way ANOVA, $*P < 0.05$, $**P < 0.01$. **c,** Comparison of the *Caudovirales*
223 order abundance in the rectal mucosa of controls and UC subjects with and
224 without inflammation. *Caudovirales* abundance was calculated as RPKM sum
225 of *Caudovirales* contigs recruited reads normalized by sequence depth in each
226 subject. Between-group difference was determined by Mann-Whitney test, $*P$
227 < 0.05 . **d,** Comparison of the mucosal *Caudovirales* α diversity based on
228 Shannon diversity, evenness, Chao1 richness index in the mucosa of controls
229 and UC subjects with and without inflammation. Statistical significance was
230 determined by one-way ANOVA, $*P < 0.05$. For box plots, the boxes extend
231 from the 1st to 3rd quartile (25th to 75th percentile), with the median depicted
232 by a vertical line. Dots indicate individual values of the studied subjects.

233

234 **Supplementary Figure 4. Phylogentic dendrogram of mucosa viral taxa in**
235 **health and UC.** Differentially enriched viral taxa between healthy and UC
236 mucosa are color-coated, at the order, family and genus levels. Red colored
237 taxa correspond to UC mucosa enriched viruses, while blue colored taxa
238 correspond to healthy individual enriched viruses.

239

240 **Supplementary Figure 5. Validation of differential viral species in Beijing**
241 **and Xiangshan cohorts. a**, Geographic plot of study sites in China. 20
242 healthy controls versus 20 UC subjects and 8 healthy controls versus 8 UC
243 subjects were enrolled from study sites Beijing and Xiangshan, Zhejiang
244 province, respectively. Abundances of the species *Chrysochromulina ericina*
245 *virus* (**b**) and *Mimivirus* (**c**) at study sites Beijing and Xiangshan were plotted in
246 mean \pm s.e.m. Between-group difference was determined by Mann-Whitney
247 test, * $P < 0.05$, ** $P < 0.01$, * $P < 0.001$. Dots indicate individual values of the
248 studied subjects.

249
250 **Supplementary Figure 6. Differential viral species between mucosal**
251 **virome Metacommunity 1 subjects in health and UC. a**, Differentially
252 enriched viral species between the mucosa of healthy Metacommunity 1
253 subjects and UC Metacommunity 1 subjects, determined by *DESeq* analysis
254 with FDR correction. Only those taxa with adjusted p values < 0.05 and
255 $|\text{Log}_2(\text{between-group fold-change})| > 2$ are shown. Blue denotes viral species
256 enriched in healthy metacommunity 1 subjects, red denotes viral species
257 enriched in UC metacommunity 1 subjects. The abundances of two replicated
258 viral species, *Feldmannia species virus* (**b**) and *Pseudomonas virus* (**c**), in
259 Beijing and Xiangshan cohort, were plotted with respect to each study site.
260 Between-group statistical difference was determined by Mann-Whitney test, * P
261 < 0.05 . Dots indicate individual values of the studied subjects.

262
263 **Supplementary Figure 7. Mucosal virome functions in healthy individuals**
264 **and UC subjects, predicted at the Gene ontology and Pfam protein family**
265 **levels. a**. The abundance boxplot of mucosal viral function in healthy
266 individuals. **b**. The abundance boxplot of mucosal viral function in UC subjects.
267 Only the most abundant 20 viral functions, GO and Pfam protein functions
268 respectively, are shown for clarity. For box plots, the boxes extend from the 1st
269 to 3rd quartile (25th to 75th percentile), with the median depicted by a horizontal
270 line.

271
272 **Supplementary Figure 8. Altered bacterial microbiota in UC mucosa. a**,
273 Comparison of bacteria α diversities based on Simpson diversity and Chao1
274 richness in the mucosa of controls and UC subjects. Statistical significance
275 was determined by t test, * $P < 0.05$, ** $P < 0.01$. For box plots, the boxes
276 extend from the 1st to 3rd quartile (25th to 75th percentile), with the median
277 depicted by a horizontal line. The bacteria composition in health and UC
278 mucosa was plotted in relative abundance, at the phylum (**b**), family (**c**) and
279 genus (**d**) levels. **e**, Differentially enriched bacterial taxa between the mucosa
280 of health individuals and UC subjects were determined by *Lefse* analysis with

281 FDR correction. Only those taxa with adjusted P value < 0.05 and LDA effect
282 size >2 are shown.
283
284
285