

# 1 Characterizing enterotypes in human metagenomics: a viral perspective

2 Li Song<sup>1</sup>, Lu Zhang<sup>2\*</sup>, Xiaodong Fang<sup>1\*</sup>

3 <sup>1</sup> College of Life Sciences, University of Chinese Academy of Sciences, Beijing 100049, China

4 <sup>2</sup> Department of Computer Science, Hong Kong Baptist University, Hong Kong

5 **\*Correspondence:**

6 Corresponding Authors

7 [fangxd@bgi.com](mailto:fangxd@bgi.com), [ericluzhang@hkbu.edu.hk](mailto:ericluzhang@hkbu.edu.hk)

8 **Keywords: Gut metagenomics, Virome, Enterotype, Biomarker, Perspective**

## 9 Abstract

10 The diversity and high genomic mutation rates of viral species hinder our understanding of viruses  
11 and their contributions to human health. Here we investigated the human fecal virome using  
12 previously published sequencing data of 2,690 metagenomes from seven countries. We found that the  
13 virome was dominated by double-stranded DNA viruses, and young children and adults showed  
14 dramatic differences in their fecal enterovirus composition. Beta diversity showed there were  
15 significantly higher distances to centroids in individuals with severe phenotypes, such as cirrhosis. In  
16 contrast, there were no significant differences in lengths to centroids or viral components between  
17 patients with mild phenotypes, such as hypertension. Enterotypes showed the same specific viruses  
18 and enrichment direction after independent determination of enterotypes in various projects.  
19 Confounding factors, such as different sequencing platforms and library construction, did not result  
20 in a batch effect to confuse enterotype assignment. The gut virome composition pattern could be  
21 described by two viral enterotypes, which supported a discrete, rather than a gradient, distribution.  
22 Compared with enterotype 2, enterotype 1 had a higher viral count and Shannon index, but a lower  
23 beta diversity, indicating more resistance to the external environment's harmful effects. Disease was  
24 usually accompanied by a viral enterotype disorder. However, a sample outside of the enterotyping  
25 mathematical space of enterotype database did not necessarily indicate sickness. Therefore, the  
26 background context must be carefully considered when using a viral enterotype as a biomarker for  
27 disease prediction. The disease, second only to the enterotype, explains significant variation in viral  
28 community composition, implying that double-stranded DNA is relevant to human health. Our results  
29 of investigating a baseline viral database highlight important insights into the virome composition of  
30 human ecosystems, and provide an alternate biomarker for early disease screening.

## 31 1 Background

32 In recent years, many studies have shown that viral colonization in the human body is highly related  
33 to human health and life. Cross-species virus transmission poses an extraordinary threat to human  
34 and animal health (Daszak et al., 2000). With advanced sequencing technology, the primary material  
35 for viral research has become viral genomes (virome), which enable viral identification and  
36 classification at the molecular level (Fujimoto et al., 2020; Gregory et al., 2020). The success of  
37 virome studies greatly relies on high-quality viral genomes (Minot et al., 2011). However, viruses are  
38 highly diverse and individual-specific (ref) and traditional purification strategies, culture, and

39 sequencing are labor-intensive and inefficient (Reddy et al., 2015), thus severely preventing the  
40 comprehensive and intensive study of viruses.

41 The strategy of assembling the viral genome involves a comprehensive and in-depth analysis of the  
42 virome. David and colleagues launched the “Uncovering Earth’s virome” project to build The  
43 Integrated Microbial Genome/Virus (IMG/VR) database in 2016(Paez-Espino et al., 2016, 2017,  
44 2019; Roux et al., 2021). Recently, data of 28,060 metagenomes were used to mine 142,809 human  
45 gut viruses, and Gubaphage was found to be the second common virus branch in the human  
46 gut(Camarillo-Guerrero et al., 2021). These projects opened the prelude to the construction of viral  
47 genome database and laid the foundation for a comprehensive analysis of the human gut  
48 virome(Gregory et al., 2020). Disease is associated with the gut virome, but studies have ignored the  
49 importance of viral sequencing information in massive metagenome sequencing data. The  
50 construction of the viral genome database has enabled detailed research on the human gut virome.

51 An enterotype is a cluster of microbes in the human gut and it describes the distributional of the  
52 human gut microbial community(Arumugam et al., 2011). Multiple studies have reported that there  
53 are two dominant enterotypes, which correspond to the individuals’ preference for digesting plant  
54 fiber or animal meat (Costea et al., 2017). The gut is an ecosystem, and the enterotype summarizes its  
55 microbial characteristics using mathematical methods(Arumugam et al., 2011; Holmes et al., 2012),  
56 but such knowledge is insufficient (Jeffery et al., 2012). Research on the composition patterns and  
57 function of the gut microbiome will significantly improve our understanding of its relationship with  
58 health and disease(Knights et al., 2014). Enterotypes can be used for gut microbial analysis, to  
59 inform disease treatment and prevention strategies, and may also provide a theoretical basis for diet  
60 therapy. The relationship between viral enterotypes and the human disease status is still largely  
61 unknown. Whether enterotypes can be used as biomarkers for predicting the disease status requires  
62 further research.

63 In this study, we collected previously published human metagenomic sequencing data, conducted  
64 sample quality control through a fast pipeline, identified virus species, and determined viral  
65 abundance. Furthermore, we established a baseline database of the human gut virome based on 2,690  
66 metagenomes. We demonstrated the relationship between virus species and abundance in various  
67 ethnicities, countries, and diseases using different DNA library construction methods and sequencing  
68 platforms, and analyzed the association between viral community diversity and disease. Viral  
69 enterotypes were assigned by the Dirichlet multinomial mixture model (DMM). We independently  
70 identified enterotype-specific viral operational taxonomic units (vOTUs) for each dataset and  
71 resolved the inter-relationships among enterotypes from different projects by comparing the  
72 abundance of enterotype-specific viruses. Further, we compared the ecological diversity of viruses  
73 between different enterotypes, and evaluated the correlation of viral enterotype disorders and their  
74 diversity with diseases. The results of study elucidate the relationship between enteroviruses and  
75 human health in a large population and highlight the decisive role of viruses as molecular markers in  
76 identifying high-risk individuals. Viral research is likely to make an indispensable contribution to  
77 improving human health.

## 78 **2 Materials and methods**

### 79 **2.1 Choosing an alignment method**

80 Alignment and assembly methods are used to detect viruses and estimate their abundance.  
81 MetaPhlAn(Segata et al., 2012)and its upgraded version, MetaPhlAn2(Truong et al., 2015), are  
82 alignment tools that use marker genes for alignment and have achieved great success in bacterial

83 genome alignment. However, many viral genomes do not have marker genes, and therefore, this  
84 strategy is not useful for viral classification and abundance estimation. Virome(Eric Wommack et al.,  
85 2012), VirSorter(Roux et al., 2015), and VirFinder(Ren et al., 2017) use assembly methods to  
86 classify viruses and calculate abundance, but these tools require a large amount of computing  
87 resources and time and therefore cannot be applied to large projects. Some recently developed  
88 alignment tools, such as ViromeScan(Rampelli et al., 2016), VIP(Li et al., 2016), and HoloVir(Laffy  
89 et al., 2016) have been shown to perform well for bacterial genomes. However, they are impractical  
90 for aligning viral reads to genomes. Moreover, many software tools are for online use, which means  
91 that they are unsuitable for large-scale projects. VirMap(Ajami et al., 2018) software developed for  
92 processing protein and genome data can provide good results. It can be accurately identify the virus  
93 species regardless of the sequencing depth. However, this software involves substantial computing  
94 resources. After comparing the advantages and disadvantages of different software(Ajami et al.,  
95 2018), we finally chose FastViromeExplorer(Tithi et al., 2018), a software based on k-mer alignment  
96 used by Kallisto(Bray et al., 2016). This software maps all reads to the reference and then uses the  
97 expectation-maximization algorithm to estimate the virus species and their corresponding abundance.

## 98 **2.2 Data collection and processing**

99 We downloaded all data from the National Center for Biotechnology Information (NCBI) sequence  
100 read archive (SRA). The SRA numbers for each project are shown in Supplementary Table 1. We  
101 only chose pair-end data from projects sequenced by the Illumina HiSeq 2000 or 2500 platforms.  
102 After processing the original data sample (Supplementary Figure 1) using Trimmomatics(Bolger et  
103 al., 2014)to remove the raw data and adapters of low-quality reads, we detected and removed  
104 contamination from the host's DNA and RNA data, and discarded the unpaired reads. Finally, we  
105 used FastViromeExplorer software to align reads to IMG/VR v2.

## 106 **2.3 Viral contig taxonomic annotation**

107 We used Glimmer3 toolkit Version 3.02b(Delcher et al., 2007) to predict and extract the open  
108 reading frame of viral contigs with a minimum length threshold of 100 amino acids. The protein  
109 sequences were aligned to the UniProt TrEMBL database as of February 2021(Bateman et al., 2021)  
110 using BLASTX(Boratyn et al., 2012). The major voting system was then used as described  
111 previously to ascertain the family of a viral contig. A contig needed to be supported by five proteins  
112 to be considered as successful assignment; otherwise, the assignment was considered a failure. When  
113 a virus sequence was annotated to multiple families in taxonomic assignment, we choose the family  
114 with the largest proteins. When multiple families have the same number of proteins, the size of the  
115 accumulated E-value (BLASTX alignment) of all proteins was compared.

## 116 **2.4 Calculation of ecological diversity**

117 We first used Tximport(Soneson et al., 2015) R package to read the original abundance information  
118 of the virus (the output of Kallisto) from each project. The “betadiver” in the Vegan R package was  
119 used for calculating alpha and beta diversity. For alpha diversity, we first transformed the abundance  
120 information into integers and then used the “rrarefy” function to normalize abundance. We then used  
121 the “estimate,” “diversity,” and “specnumber” functions to obtain various measurement values of  
122 alpha diversity. We used the “RLE” method embedded in “calcNormFactors” to normalize raw  
123 abundance for beta diversity. We used the “Hellinger” method in “Decostand” to transform the data  
124 and eliminate false similarities caused by many viruses whose abundance was 0. When the  
125 abundance of many viruses in the two samples were 0, some algorithms might consider them to have  
126 similar abundance distribution and conclude that they were close to each other. The real reason may

127 be that many viruses have not been detected. We used the “betadiver” and “betadisper” functions to  
128 obtain beta diversity, and then used Adonis2 to analyze the viral ecological differences between cases  
129 and controls, and corrected them with raw data size. The Kruskal test was used to determine whether  
130 there was a significant difference in the distance from the centroid between cases and controls.  
131 Tukey’s honestly significant difference test was used to determine differences in variance within and  
132 between groups.

## 133 **2.5 Enterotyping and MaAsLin2 analysis**

134 We used the DMM method to determine viral enterotypes in each project independently. Enterotypes  
135 were assigned using the “DirichletMultinomial” R package, with predetermined parameters of 1 to 10  
136 enterotypes, and enterotype data from each project were run 10 times. The smallest Laplace value  
137 corresponding to the number of enterotypes was considered as the optimal result. MaAsLin2  
138 ([huttenhower.sph.harvard.edu/maaslin2](http://huttenhower.sph.harvard.edu/maaslin2)) analysis was used to determine the specific vOTUs  
139 associated with enterotypes, with correlations considered significant at the 5% level (after multiple  
140 testing correction). We also applied the envfit function in Vegan to estimate the effect size of the  
141 structural variance explained by factors such as enterotype and disease.

## 142 **3 Results**

### 143 **3.1 Sequencing data and summarization**

144 We collected 12.36 TB of metagenomic sequencing data from 18 previously published projects  
145 (Supplementary Tables 1 and 2). We selected data from 2,690 metagenome samples of high quality  
146 for the subsequent analysis (Supplementary Figure 1 and Supplementary Table 1), of which 1,092  
147 were samples were from women, 859 were from men, and 739 were from unknown sex. The length  
148 of sequencing reads from each sample were 2.26 to 8.55 G (Supplementary Table 1), and  
149 approximately 10% of strictly filtered reads were aligned against IMG/VR v2 viral sequences  
150 (Supplementary Table 3). We obtained 2,690 metagenome samples by choosing paired-end  
151 sequencing data from the Illumina HiSeq 2000 and 2500 platforms and excluding projects with a  
152 small data size (< 1 G).

153 [insert figure 1 here]

154 We annotated the geographic locations of the included projects on the basis of their predominant  
155 samples (Figure 1A). Because there were no specific sampling coordinates, each project was located  
156 by country. We annotated the viral taxonomy at the family level based on the protein sequence  
157 similarities (Minot et al., 2013; Hannigan et al., 2015). Approximately 50% of the viral genomes  
158 failed taxonomic assignment (Figure 1B), and double-stranded (ds) DNA viruses, such as  
159 Siphoviridae, Myoviridae, and Podoviridae, were the dominant enteroviruses as previously reported  
160 (Zuo et al., 2020). The density peak was close to zero, which indicated that the viruses were rarely  
161 shared among individuals (Supplementary Figure 2). The samples from Finland were outliers in the  
162 PCoA and tSNE plots (Figure 1C and D) because of the low viral diversity (Supplementary Figure  
163 3). This finding might be explained by age. The average age of individuals in the Finland project was  
164 1.5, and their gut communities did not reach stable states. Although the samples from the other six  
165 countries showed substantial variability in the PCoA and tSNE plots (Figure 1C and D), they  
166 belonged to the same cluster, especially the samples from the studies conducted in China. The studies  
167 from China had the most individuals, and the samples were spread over almost the entire plot. In the  
168 tSNE plot, we found that the samples from the USA and Peru were clustered in a local region, which  
169 indicated that the gut virome showed characteristics of geographical distribution.

170 [insert figure 2 here]

171 To study the distribution characteristics of the viral species in samples with different phenotypes, we  
172 divided all samples from studies with a case–control design into three categories. These categories of  
173 controls, cases, and all represented healthy people, patients with various diseases, and all individuals,  
174 respectively. As more samples were included, the number of viral species showed exponential  
175 growth, with no significant difference between cases and controls until samples from ~100  
176 individuals were included (Figure 2A). After including ~100 individuals, the “case” curve showed a  
177 steep increased viral count. As expected, a significant increment in the number of viral species was  
178 observed when the number of samples was increased in the “all” curve. However, the three growth  
179 curves were essentially parallel (Figure 2A), which suggested that the overall number of viruses in  
180 the patient population after viral community disruption was limited. More interestingly, the “case”  
181 and “all” curves overlapped with each other after ~1000 samples. The reason for this finding could be  
182 that the case population contained all species of viruses in the control population. When we  
183 compared the growth curves of different projects, we found that the curves for Finland, Peru, and  
184 Chinese populations with cirrhosis had significant differences (Figure 2B). The samples from the  
185 Finland project were obtained from only 1.5-year-old children, at which age the enterovirus  
186 community is not well established. It is unclear why the number of viral species in Peru samples were  
187 small at the beginning of the curve. The dramatic increase in the number of viruses in the Chinese  
188 population with cirrhosis may be due to severe disruption of the enterovirus community. We used  
189 unique species in cases and controls to define group-specific viruses and compared the change in the  
190 proportion of unique viral species between cases and controls (Supplementary Table 4). We found  
191 that the mean proportion of viruses in case samples was 26% and in control samples was 14%.  
192 Among all samples, the proportion of viruses that were unique to cases was 23%. Each case  
193 individual had an average of 10.99 viruses, and the ratio of viruses that were unique to controls was  
194 4%, and each control individual had an average of 2.43 viruses. Overall, there was an enrichment of  
195 viruses in cases.

Table 1. Beta diversity for measuring the sample distance in projects with a case–control design.

| Project                    | Adonis2 for disease | Adonis2 for raw data | ANOVA    | Kruskal test |
|----------------------------|---------------------|----------------------|----------|--------------|
| Sweden T2D*                | 5.00E-03            | 1.00E-03             | 6.34E-03 | 1.62E-03     |
| China cirrhosis            | 1.00E-03            | 1.00E-03             | 3.91E-09 | 8.41E-14     |
| China rheumatoid arthritis | 2.40E-02            | 1.50E-02             | 0.86     | 0.48         |
| Austria carcinoma          | 1.00E-03            | 0.55                 | 5.29E-04 | 7.89E-05     |
| China colorectal cancer    | 4.00E-03            | 3.10E-02             | 0.12     | 9.71E-03     |
| China hypertension         | 0.08                | 0.25                 | 0.13     | 0.08         |



|                              |          |          |          |          |
|------------------------------|----------|----------|----------|----------|
| China coronary heart disease | 1.00E-03 | 0.81     | 4.15E-02 | 1.79E-02 |
| China T2D discovery          | 2.60E-02 | 3.50E-02 | 1.75E-02 | 0.15     |
| China T2D validation         | 1.00E-03 | 2.90E-02 | 0.38     | 0.45     |
| China obesity                | 0.06     | 0.37     | 0.78     | 0.78     |

---

196 \*Type 2 diabetes (T2D).

### 197 3.2 Relationship of ecological diversity of viruses and disease

198 The beta diversity of a microbial community is usually used to evaluate dynamic changes in an  
199 ecosystem (Koleff et al., 2003). A comparison of the results of projects with a case–control design  
200 revealed that the degree of imbalance in the viral community composition was related to the severity  
201 of the disease phenotype. An example of this finding is that the viral community in patients with  
202 cirrhosis (Figure 3A) was significantly different from that in healthy people (Adonis2,  $p = 0.001$ ,  
203 adjusted for raw data size). Comparison of the distance to the centroid between patients and healthy  
204 individuals by the Mann–Whitney  $U$  test showed a significant dissimilarity (Figure 3B). Specifically,  
205 patients had a significantly larger distance than healthy individuals, which indicated that patients had  
206 a considerably disordered viral community. In contrast, we did not detect a significant difference  
207 between patients and healthy individuals in the hypertension project (Adonis2,  $p = 0.08$ , Figure 3C).  
208 We also compared the distance to the centroid for each pair of three cohorts (Figure 3D), and a  
209 significant difference was found only between patients with hypertension and healthy individuals  
210 (Wilcoxon,  $p = 0.036$ ).

211 We further investigated statistical differences in gut viral composition between case and control  
212 samples from various aspects to investigate changes in the viral community across different  
213 phenotypes. Using Adonis2, we found a significant difference in enteroviruses between cases and  
214 controls expect hypertension and obesity (Table 1), which suggested that their gut viral community  
215 was less affected by the disease state. Consistently, in cases with relatively mild phenotypes, such as  
216 hypertension or obesity, there was no noticeable differences in body metabolism compared with the  
217 controls. An analysis of variance (ANOVA) was used to determine whether there was a significant  
218 difference between two centroids (to test the component of viruses) between cases and controls. We  
219 found that the cirrhosis and cancer cohorts showed a substantial difference between two centroids  
220 (Table 1). The Kruskal-Wallis test was performed to determine whether the distance to the centroid  
221 in principal coordinates analysis was significantly different between the case and control groups, and  
222 the results were consistent with those of ANOVA. Compared with the controls, cases with more  
223 severe phenotypes, such as cirrhosis and cancer, showed substantial differences in gut viral  
224 composition (Table 1), whereas cases with relatively mild phenotypes, such as hypertension, showed  
225 no significant differences.

226 [insert figure 3 here]

### 227 3.3 Characterizing viral enterotypes

228 The characteristics of enterotypes of the gut virome were the focus of this study. Data on enterotypes  
229 are generally used to help adjust population stratification in Metagenome-wide association studies

230 (MWAS) analysis (Wang et al., 2012). The correlation between enterotypes and disease phenotypes  
231 has received much attention in this field. The DMM method is commonly used for determining  
232 enterotypes of the gut microbiome and is more effective than the partitioning around medoids (Ding  
233 and Schloss, 2014). Different library construction methods, sequencing platforms, and other factors  
234 may lead to false-positive assignment of enterotypes. To avoid this situation, we adopted a project-  
235 independent strategy for determining enterotypes. There were two or three enterotypes in most  
236 projects, while some projects only had one enterotype (Figure 1A, Table 2, Supplementary Figure 4).  
237 Enterotypes with the same intrinsic composition pattern were considered as the same. We used  
238 Maaslin2 to discover enterotype-specific vOTUs and then determined their enrichment direction on  
239 the basis of mean abundance. Similar enterotypes had the same specific vOTUs and the same  
240 enrichment trend. We manually classified enterotypes in all of the projects into three groups (Table 2,  
241 Supplementary Table 5). Enterotypes 1 and 2, which are the two major types, were widely distributed  
242 in all projects, which indicated that these two types of enterotypes were common across the project  
243 populations. However, enterotype 3 was rare. Unclassified individuals were not able to be  
244 confidently assigned to enterotype 1 or 2.

245 A permutation test was performed to demonstrate the validity of manual classification, which  
246 involved randomly paired enterotypes from different projects. We assumed that paired enterotypes  
247 had the same specific vOTUs and enrichment directions. We assigned a lower error rate to paired  
248 enterotypes if they had more identical vOTUs and similar enrichment trends. We repeated pairing 5  
249 million times to obtain the distribution of pairing scores. These scores showed that our manually  
250 classified enterotypes had the lowest error rate (Figure 4A). Moreover, random pairing supported the  
251 three major enterotypes. Enterotypes 1- and 2-specific vOTUs were dominant (Figure 4B). The same  
252 enterotype-specific vOTUs with highly consistent enrichment trends indicated that the enterotypes  
253 from different projects had a similar pattern of virome composition (Figure 4B). Different DNA  
254 processing methods, sequencing platforms, ethics, age, and other confounding factors did not affect  
255 the identification of viral enterotypes. The vOTUs that were specific to unclassified enterotypes  
256 appeared complex. They intersected with either enterotype 1 or 2. Enterotype 3-specific vOTUs in  
257 different projects were less concordant than enterotypes 1- and 2-specific vOTUs.

258 [insert figure 4 here]

259 The microbiome is an ecosystem, the stability of which is reflected by the diversity of species in the  
260 system. As a species becomes more prosperous and uniform, the system's diversity increases and it  
261 becomes more resistant to the effects of the external environment (Keesing et al., 2010). There are  
262 two dominant enterotypes in the viral community (Zuo et al., 2020), one of which has a high alpha  
263 diversity. The results of our study are remarkably close to expected results. Although the viral count  
264 varied among samples from different projects, enterotype 1 across the samples had more viruses than  
265 enterotype 2 (Figure 4C). A higher value of the Shannon index and a smaller sample distance in  
266 enterotype 1, compared with enterotype 2, indicated its more stable composition pattern. We found  
267 that more individuals were categorized as enterotype 1 than enterotype 2 (1204 vs. 716). By  
268 comparing the proportion of healthy samples with the two enterotypes, we found that individuals  
269 who were categorized as enterotype 2 had a higher risk of being sick than those who were  
270 categorized as enterotype 1 (odds ratio: 1.38, Fisher's exact test,  $p = 0.01$ ). We observed an  
271 interesting finding when we compared samples from the cirrhosis project and the Sweden mother-  
272 child project. The third enterotype had the most discrete sample distribution in the cirrhosis project,  
273 and a higher viral count and Shannon index compared with the Sweden mother-child project (Figure  
274 3C). In contrast, the third enterotype had a large sample distance and the lowest viral count and  
275 Shannon Index in the Sweden mother-child project. The cases in these two projects had

276 diverse medical conditions. Specifically, the case cohort in the cirrhosis project had disordered gut  
277 virome due to the disease, which explains why the number of viruses in the samples did not decrease.  
278 In contrast, children in the Swedish mother-child project lacked a stable gut virome and had a lower  
279 viral count, which suggested that enterotype 3 in the samples of this project was not caused by any  
280 disease.

281 The viral enterotype may play a dominant role in influencing the structural variance of the gut virome  
282 via a variety of factors. The Adonis test was used to determine the significance of viral enterotypes.  
283 The results were significant in all projects. Our results explain most of the structural variance in the  
284 gut virome (Figure 4D). In the Peru and cirrhosis projects, the Adonis R squared values were 0.62  
285 and 0.57, respectively. Age, disease, BMI, raw data, and sex were not significant factors affecting  
286 viral enterotypes in most projects, but Adonis p values reached significance in several projects.  
287 Disease was the second most significant factor in the projects, which suggested that illness had a  
288 higher ability to reshape the gut microbiome than other factors. Characterizing the interaction  
289 between the gut virome and external stimuli was complex. Whether a single factor has a particular  
290 contribution requires consideration of the context of this factor. An example of this situation is that,  
291 in liver cirrhosis, the association between the gut virome and age was strong, but it was not  
292 significant for diabetes.

293 Table 2: Manually categorized results for each project.

| Enterotype                   | Enterotype 1 | Enterotype 2 | Enterotype 3 |
|------------------------------|--------------|--------------|--------------|
| Denmark no phenotype         | GP1          | GP2          | -            |
| China cirrhotic              | GP2          | GP1          | GP3          |
| Sweden mother-offspring pair | GP3          | GP2          | GP1          |
| China rheumatoid arthritis   | GP1          | GP2          | -            |
| Austria carcinoma            | GP1          | -            | GP2          |
| UK no phenotype              | GP1          | GP2          | -            |
| China colorectal cancer      | GP1          | GP2          | -            |
| China hypertension           | GP2          | GP3          | -            |
| China coronary heart disease | GP1          | GP2          | -            |
| China T2D discovery          | GP1          | GP2          | -            |
| China T2D validation         | GP1          | GP2          | -            |
| China healthy Mongolian      | GP1          | GP2          | -            |
| China ankylosing spondylitis | GP1          | GP2          | -            |



294 Groups in the same column were considered to belong to one enterotype.

295 [insert figure 5 here]

296 Enterotypes are useful for describing the gut microbial community, and determining the association  
297 between diseases and enterotype is important to detect high risk individual in population. In the liver  
298 cirrhosis project, individuals could be broadly divided into three categories (Figure 5A). Enterotypes  
299 1 and 3 were enriched in healthy individuals and patients, respectively (69 controls/16 cases vs. 2  
300 controls/64 cases, Supplementary Table 6), and enterotype 2 accounted for half of them (43 controls,  
301 43 cases, Supplementary Table 6). We found that the viral enterotype was significantly related to  
302 liver cirrhosis (Fisher's exact test,  $p = 5.99E-24$ , Supplementary Table 7). Enterotype 3 was loosely  
303 distributed in individuals (Figure 5A). However, enterotypes 1 and 2 showed a closer relationship.  
304 These three groups did not have discrete clustering boundaries and demonstrated some overlap with  
305 one another in the PCoA plot. There was no apparent clustering of samples enriched locally due to  
306 the viral count or the Shannon index (Figure 5B). In the hypertension project, the clustering  
307 boundaries of enterotypes 1 and 3 were more pronounced than those for enterotype 2 (Figure 5C),  
308 and there was no overlapping area between the two clusters. This finding was surprising because  
309 individuals in enterotype 2 had a smaller viral count and a lower Shannon index (Figure 5D). Some  
310 of them were close to enterotype 1, while others had clusters of enterotype 3. However, the specific  
311 vOTUs and enrichment direction of individuals in enterotype 2 showed a high consistency (Figure  
312 4B), indicating that enterotype 2 was real. We found no significant association between the viral  
313 enterotype and hypertension (Fisher's exact test,  $p = 0.3$ , Supplementary Table 7). Gut virome  
314 community disorders showed significant differences in the cirrhosis and hypertension projects, which  
315 indicated that not all diseases caused evident ecological perturbation in the human gut. Thus,  
316 applying viral enterotypes as biomarkers for predicting clinical disease requires specific  
317 consideration.

#### 318 **4 Discussion**

319 Recent investigations have shown that enterotypes of the human gut can be divided into two  
320 categories based on their predominant flora (Bacteroidetes/Prevotella). Their functions correspond to  
321 the digestion of meat and vegetarian food (Arumugam et al., 2011; Costea et al., 2017). However,  
322 some researchers consider that the distribution of enterotypes is not discrete, but rather gradient. This  
323 viewpoint suggests that those two enterotypes are the two endpoints of the gradient distribution of  
324 Bacteroidetes/Prevotella (Jeffery et al., 2012). This study used the DMM method to assign viral  
325 enterotypes and showed that there were two enterotypes in most projects. Viral enterotypes did not  
326 have an apparent dominant virus. An explanation for this finding may be that most human  
327 enteroviruses are dsDNA viruses. As previously reported, dsDNA viruses are less harmful than RNA  
328 virus to the human body (Dutilh et al., 2014; Camarillo-Guerrero et al., 2021). These viruses do not  
329 undergo strong selection when colonizing the human gut, and there is no dominant viral strain that  
330 can occupy the whole human intestine. Recent studies have shown two common and harmless  
331 dsDNA virus branches in the human intestine, namely crAssphage (Dutilh et al., 2014) and  
332 Gubaphage (Camarillo-Guerrero et al., 2021). These two dominant virus branches may correlate with  
333 the two main viral enterotypes observed in this study. We analyzed the abundance of vOTUs  
334 corresponding to each enterotype and found that in different projects, the OTUs and enrichment  
335 direction of a specific virus in the same viral enterotype were consistent. Therefore, existing evidence  
336 and our findings support the view of two discrete viral enterotypes.

337 There was a third enterotype in several projects, but because of limited evidence, we could not  
338 conclude that it is ubiquitous in the human gut. In the hypertension project, researchers found that all  
339 specific vOTUs in the enterotype were shared with enterotypes 1 and 2. PCoA analysis also showed  
340 that individuals were located in the interconnection area, which is likely to explain a gradient  
341 distribution. In the liver cirrhosis project, 64 of the 66 samples were from patients, and the third  
342 enterotype was significantly related to patients. Unlike the hypertension project, the third enterotype  
343 in patients with liver cirrhosis was not related to the number of viruses. In other projects, viral  
344 enterotypes 1 and 2 had more specific viruses and a higher consistent enrichment direction in contrast  
345 to the rare specific viruses of viral enterotype 3 and an inconsistent enrichment direction. This result  
346 suggests that a third viral enterotype in various projects may not have belonged to the same cluster.  
347 Therefore, we cannot conclude that there was a stable presence of viral enterotype 3 in the  
348 population. We speculate that interaction between emergence of a disease and disorder of the gut  
349 virome may contribute to emergence of viral enterotype 3. Our results are in agreement with existing  
350 studies on the disruption of the human gut community that accompanies disease(Wang and Jia, 2016;  
351 Yu et al., 2017; Nakatsu et al., 2018).

352 In this study, enterotype 1 had a higher viral count and Shannon index compared with enterotype 2.  
353 In addition to having a smaller sample distance, we speculate that enterotype 1 might have more  
354 stable viral ecological communities than enterotype 2. We found enterotype 2 had 1.38 times more  
355 patients than the one in enterotype 1. This result suggests that a stable microbial community has a  
356 higher ability to resist the influence of external stimulations. In the cirrhosis project, enterotype 3 was  
357 characterized as being enriched in patients and having an extremely disordered gut virome.  
358 Enterotype 3 had the largest sample distance and highest number of virus species. A possible reason  
359 for this finding is that bile acid secretion in patients with liver cirrhosis is obstructed, which leads to  
360 drastic changes in the gut microbiome of the patients. This may have resulted in large-scale  
361 replacement of the virome and reduced similarity of virus species in this patient population. It is also  
362 possible that the microenvironment of viral evolution in the human body is disturbed owing to  
363 disease progression or the similarity of virus species is decreased due to a shift in the distribution of  
364 the ecological gradient. We also found a large distance in samples from the Sweden mother-child pair  
365 project, with the viral count being significantly lower than the global average level. A possible  
366 explanation for this finding is that the gut microbiome of young children is developing and has yet to  
367 reach a stable state(Derrien et al., 2019).

368 We determined enterotypes at the bacterial and viral levels in the China diabetes(Wang et al., 2012),  
369 and found a strong correlation between enterotypes at these two levels (China T2D discovery:  $p =$   
370  $1.70E-07$ ; China T2D validation:  $p = 1.58E-11$ , Fisher's exact test, Supplementary Table 8). We  
371 found that enterotypes (bacterial and viral levels) were not randomly distributed and that the bacterial  
372 community had a strong selection effect on the viral community. However, bacterial- and viral-level  
373 enterotypes were not correlated or were weakly correlated with sex, age, BMI, and disease  
374 (Supplementary Table 9). This finding may be explained by the use of high-abundance bacterial and  
375 viral species for determining enterotypes. A high abundance of bacteria or viruses in the intestine is  
376 significantly related to disease. Such a high abundance directly and severely affects the human body,  
377 That is a microbial infection, not a harmonious symbiosis, which is in contrast to the current  
378 understanding of the gut community and health. Wang et al. used enterotype as a covariate in their  
379 study and proposed a useful method to stratify human gut microbiomes in MWAS, which effectively  
380 improved the power of hypothesis testing(Wang et al., 2012). Although we found a strong correlation  
381 between bacterial and viral enterotypes, they were not proven to be equivalent. Therefore, we suggest  
382 using bacterial and viral enterotypes as independent covariates in MWAS. More in-depth

383 investigations are warranted to determine whether this strategy can efficiently reduce false-positive  
384 and false-negative rates of investigating pathogenic microbiomes.

385 We used MaAsLin2 to identify viruses that were specific to enterotypes and disease. We found that  
386 the number of vOTUs that were specific to disease was significantly lower when we simultaneously  
387 entered these two factors into the software program than when we entered only disease. In the  
388 cirrhosis project, 56 vOTUs were specific to the disease state ( $q$ -value  $\leq 0.05$ ), when the enterotype  
389 was excluded. In contrast, we found 241 and 7 vOTUs were specific to the enterotype and disease,  
390 respectively, when we included these two factors simultaneously. There were 21 vOTUs tested as  
391 disease-related became enterotype-associated. These results suggest that viral enterotypes need to be  
392 taken into account is MWAS. Although much of the literature suggests that most dsDNA viruses are  
393 not strongly associated with disease, we cannot rule out the contribution of dsDNA viruses to illness.

394 To better determine the efficacy of using enterotypes to identify the disease state, a standardized  
395 method is first required for determining enterotypes. One such standardized pipeline for enterotypes  
396 was previously reported by Costea et al. (2017). Researchers have also established an enterotype  
397 database of the gut microbiome community on the basis of the MetaHIT dataset (Qin et al., 2010; le  
398 Chatelier et al., 2013). They then built a machine learning model and trained it by applying the  
399 MetaHIT dataset and the corresponding enterotypes. Finally, this model was used to predict the  
400 enterotypes of testing samples on the basis of their bacterial abundance matrix. We independently  
401 assigned enterotypes in different projects and found that the two manually adjusted categories shared  
402 the most specific viruses and similar enrichment directions. This consistency masked the batch  
403 effects among different datasets, and demonstrates the ability of viral enterotypes to identify  
404 individuals with disease. The construction of a large-scale viral enterotype database to define the  
405 enterotyping mathematical space of healthy individuals might be helpful to detect individuals with  
406 disease outside the mathematical space. Therefore, we believe that using viral enterotypes of the gut  
407 virome community as a feature for disease prediction will significantly improve the accuracy of  
408 disease prediction.

## 409 **5 Author Contributions**

410 XF conceived this study. LS, LZ, and XF analyzed data, prepared the figures, and drafted the  
411 manuscript.

## 412 **6 Funding**

413 This work was financially supported by the Science Technology and Innovation Committee of  
414 Shenzhen Municipality, China (SGDX20190919142801722).

## 415 **7 Conflict of Interest Statement**

416 The authors declare that the research was conducted in the absence of any commercial or financial  
417 relationships that could be construed as a potential conflict of interest.

## 418 **8 Acknowledgments**

419 The authors thank many interns and former colleagues for collecting data, and their colleague Yufen  
420 Huang for discussing the analysis strategy.

## 421 **9 References**

- 422 Ajami, N. J., Wong, M. C., Ross, M. C., Lloyd, R. E., and Petrosino, J. F. (2018). Maximal viral  
423 information recovery from sequence data using VirMAP. *Nature Communications* 9.  
424 doi:10.1038/s41467-018-05658-8.
- 425 Arumugam, M., Raes, J., Pelletier, E., Paslier, D. le, Yamada, T., Mende, D. R., et al. (2011).  
426 Enterotypes of the human gut microbiome. *Nature* 473, 174–180. doi:10.1038/nature09944.
- 427 Bateman, A., Martin, M. J., Orchard, S., Magrane, M., Agivetova, R., Ahmad, S., et al. (2021).  
428 UniProt: The universal protein knowledgebase in 2021. *Nucleic Acids Research* 49, D480–  
429 D489. doi:10.1093/nar/gkaa1100.
- 430 Bolger, A. M., Lohse, M., and Usadel, B. (2014). Trimmomatic: A flexible trimmer for Illumina  
431 sequence data. *Bioinformatics* 30, 2114–2120. doi:10.1093/bioinformatics/btu170.
- 432 Boratyn, G. M., Schäffer, A. A., Agarwala, R., Altschul, S. F., Lipman, D. J., and Madden, T. L.  
433 (2012). Domain enhanced lookup time accelerated BLAST. *Biology Direct* 7. doi:10.1186/1745-  
434 6150-7-12.
- 435 Bray, N. L., Pimentel, H., Melsted, P., and Pachter, L. (2016). Near-optimal probabilistic RNA-seq  
436 quantification. *Nature Biotechnology* 34, 525–527. doi:10.1038/nbt.3519.
- 437 Camarillo-Guerrero, L. F., Almeida, A., Rangel-Pineros, G., Finn, R. D., and Lawley, T. D. (2021).  
438 Massive expansion of human gut bacteriophage diversity. *Cell* 184, 1098-1109.e9.  
439 doi:10.1016/j.cell.2021.01.029.
- 440 Costea, P. I., Hildebrand, F., Manimozhiyan, A., Bäckhed, F., Blaser, M. J., Bushman, F. D., et al.  
441 (2017). Enterotypes in the landscape of gut microbial community composition. *Nature*  
442 *Microbiology* 3, 8–16. doi:10.1038/s41564-017-0072-8.
- 443 Daszak, P., Cunningham, A. A., and Hyatt, A. D. Emerging Infectious Diseases of Wildlife-  
444 Threats to Biodiversity and Human Health. Available at: [www.sciencemag.org](http://www.sciencemag.org).
- 445 Delcher, A. L., Bratke, K. A., Powers, E. C., and Salzberg, S. L. (2007). Identifying bacterial genes  
446 and endosymbiont DNA with Glimmer. *Bioinformatics* 23, 673–679.  
447 doi:10.1093/bioinformatics/btm009.
- 448 Derrien, M., Alvarez, A. S., and de Vos, W. M. (2019). The Gut Microbiota in the First Decade of  
449 Life. *Trends in Microbiology* 27, 997–1010. doi:10.1016/j.tim.2019.08.001.
- 450 Ding, T., and Schloss, P. D. (2014). Dynamics and associations of microbial community types across  
451 the human body. *Nature* 509, 357–360. doi:10.1038/nature13178.
- 452 Dutilh, B. E., Cassman, N., McNair, K., Sanchez, S. E., Silva, G. G. Z., Boling, L., et al. (2014). A  
453 highly abundant bacteriophage discovered in the unknown sequences of human faecal  
454 metagenomes. *Nature Communications* 5. doi:10.1038/ncomms5498.
- 455 Eric Wommack, K., Bhavsar, J., Polson, S. W., Chen, J., Dumas, M., Srinivasiah, S., et al. (2012).  
456 VIROME: A standard operating procedure for analysis of viral metagenome sequences.  
457 *Standards in Genomic Sciences* 6, 427–439. doi:10.4056/sigs.2945050.

- 458 Fujimoto, K., Kimura, Y., Shimohigoshi, M., Satoh, T., Sato, S., Tremmel, G., et al. (2020).  
459 Metagenome Data on Intestinal Phage-Bacteria Associations Aids the Development of Phage  
460 Therapy against Pathobionts. *Cell Host and Microbe* 28, 380-389.e9.  
461 doi:10.1016/j.chom.2020.06.005.
- 462 Gregory, A. C., Zablocki, O., Zayed, A. A., Howell, A., Bolduc, B., and Sullivan, M. B. (2020). The  
463 Gut Virome Database Reveals Age-Dependent Patterns of Virome Diversity in the Human Gut.  
464 *Cell Host and Microbe* 28, 724-740.e8. doi:10.1016/j.chom.2020.08.003.
- 465 Hannigan, G. D., Meisel, J. S., Tyldsley, A. S., Zheng, Q., Hodkinson, B. P., Sanmiguel, A. J., et al.  
466 (2015). The human skin double-stranded DNA virome: Topographical and temporal diversity,  
467 genetic enrichment, and dynamic associations with the host microbiome. *mBio* 6.  
468 doi:10.1128/mBio.01578-15.
- 469 Holmes, I., Harris, K., and Quince, C. (2012). Dirichlet multinomial mixtures: Generative models for  
470 microbial metagenomics. *PLoS ONE* 7. doi:10.1371/journal.pone.0030126.
- 471 Jeffery, I. B., Claesson, M. J., O'Toole, P. W., and Shanahan, F. (2012). Categorization of the gut  
472 microbiota: Enterotypes or gradients? *Nature Reviews Microbiology* 10, 591–592.  
473 doi:10.1038/nrmicro2859.
- 474 Keesing, F., Belden, L. K., Daszak, P., Dobson, A., Harvell, C. D., Holt, R. D., et al. (2010). Impacts  
475 of biodiversity on the emergence and transmission of infectious diseases. *Nature* 468, 647–652.  
476 doi:10.1038/nature09575.
- 477 Knights, D., Ward, T. L., McKinlay, C. E., Miller, H., Gonzalez, A., McDonald, D., et al. (2014).  
478 Rethinking enterotypes. *Cell Host and Microbe* 16, 433–437. doi:10.1016/j.chom.2014.09.013.
- 479 Koleff, P., Gaston, K. J., and Lennon, J. J. (2003). Measuring beta diversity for presence-absence  
480 data.
- 481 Laffy, P. W., Wood-Charlson, E. M., Turaev, D., Weynberg, K. D., Botté, E. S., van Oppen, M. J. H.,  
482 et al. (2016). HoloVir: A workflow for investigating the diversity and function of viruses in  
483 invertebrate holobionts. *Frontiers in Microbiology* 7. doi:10.3389/fmicb.2016.00822.
- 484 le Chatelier, E., Nielsen, T., Qin, J., Prifti, E., Hildebrand, F., Falony, G., et al. (2013). Richness of  
485 human gut microbiome correlates with metabolic markers. *Nature* 500, 541–546.  
486 doi:10.1038/nature12506.
- 487 Li, Y., Wang, H., Nie, K., Zhang, C., Zhang, Y., Wang, J., et al. (2016). VIP: An integrated pipeline  
488 for metagenomics of virus identification and discovery. *Scientific Reports* 6.  
489 doi:10.1038/srep23774.
- 490 Minot, S., Bryson, A., Chehoud, C., Wu, G. D., Lewis, J. D., and Bushman, F. D. (2013). Rapid  
491 evolution of the human gut virome. *Proceedings of the National Academy of Sciences of the*  
492 *United States of America* 110, 12450–12455. doi:10.1073/pnas.1300833110.
- 493 Minot, S., Sinha, R., Chen, J., Li, H., Keilbaugh, S. A., Wu, G. D., et al. (2011). The human gut  
494 virome: Inter-individual variation and dynamic response to diet. *Genome Research* 21, 1616–  
495 1625. doi:10.1101/gr.122705.111.

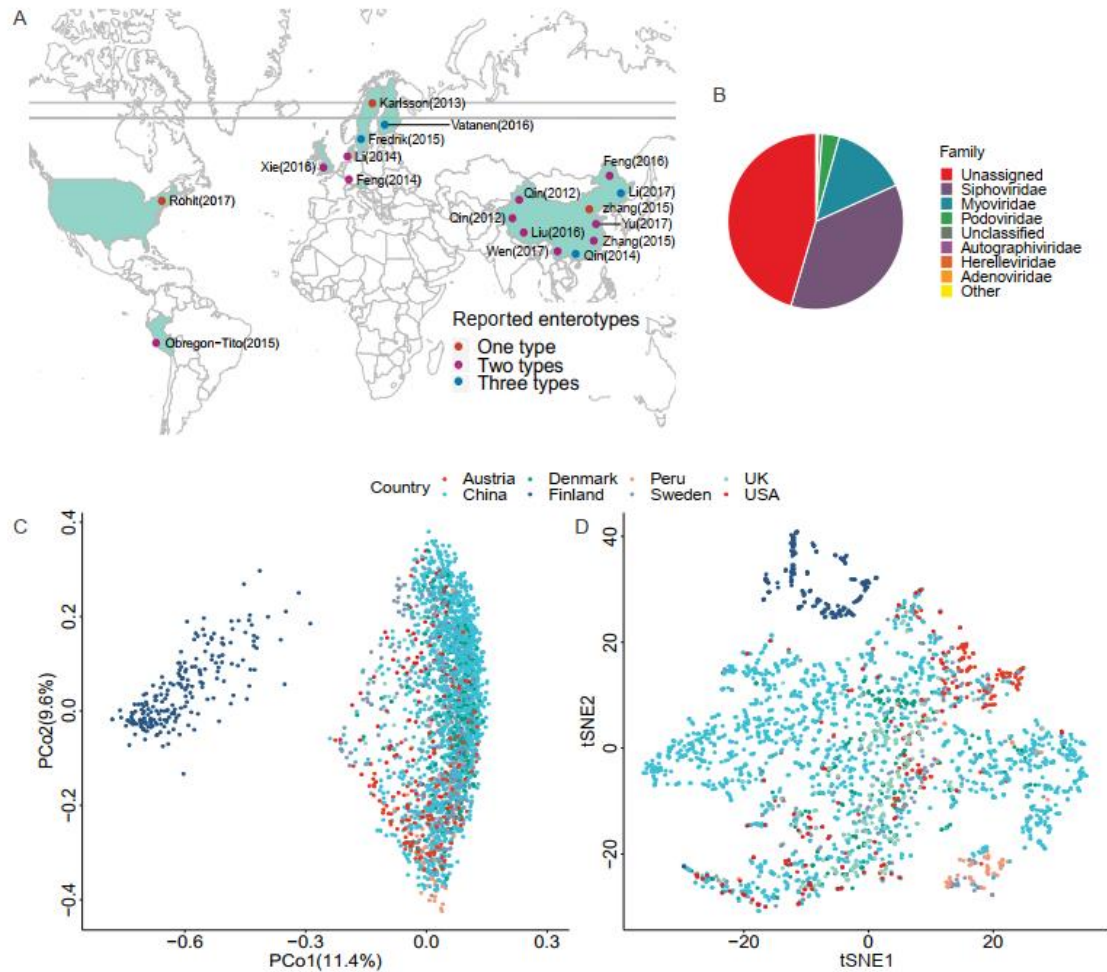


- 496 Nakatsu, G., Zhou, H., Wu, W. K. K., Wong, S. H., Coker, O. O., Dai, Z., et al. (2018). Alterations in  
497 Enteric Virome Are Associated With Colorectal Cancer and Survival Outcomes.  
498 *Gastroenterology* 155, 529-541.e5. doi:10.1053/j.gastro.2018.04.018.
- 499 Paez-Espino, D., Chen, I. M. A., Palaniappan, K., Ratner, A., Chu, K., Szeto, E., et al. (2017).  
500 IMG/VR: A database of cultured and uncultured DNA viruses and retroviruses. *Nucleic Acids*  
501 *Research* 45, D457–D465. doi:10.1093/nar/gkw1030.
- 502 Paez-Espino, D., Eloë-Fadrosch, E. A., Pavlopoulos, G. A., Thomas, A. D., Huntemann, M.,  
503 Mikhailova, N., et al. (2016). Uncovering Earth’s virome. *Nature* 536, 425–430.  
504 doi:10.1038/nature19094.
- 505 Paez-Espino, D., Roux, S., Chen, I. M. A., Palaniappan, K., Ratner, A., Chu, K., et al. (2019).  
506 IMG/VR v.2.0: An integrated data management and analysis system for cultivated and  
507 environmental viral genomes. *Nucleic Acids Research* 47, D678–D686.  
508 doi:10.1093/nar/gky1127.
- 509 Qin, J., Li, R., Raes, J., Arumugam, M., Burgdorf, K. S., Manichanh, C., et al. (2010). A human gut  
510 microbial gene catalogue established by metagenomic sequencing. *Nature* 464, 59–65.  
511 doi:10.1038/nature08821.
- 512 Rampelli, S., Soverini, M., Turrone, S., Quercia, S., Biagi, E., Brigidi, P., et al. (2016). ViromeScan:  
513 A new tool for metagenomic viral community profiling. *BMC Genomics* 17.  
514 doi:10.1186/s12864-016-2446-3.
- 515 Reddy, T. B. K., Thomas, A. D., Stamatis, D., Bertsch, J., Isbandi, M., Jansson, J., et al. (2015). The  
516 Genomes OnLine Database (GOLD) v.5: A metadata management system based on a four level  
517 (meta)genome project classification. *Nucleic Acids Research* 43, D1099–D1106.  
518 doi:10.1093/nar/gku950.
- 519 Ren, J., Ahlgren, N. A., Lu, Y. Y., Fuhrman, J. A., and Sun, F. (2017). VirFinder: a novel k-mer  
520 based tool for identifying viral sequences from assembled metagenomic data. *Microbiome* 5, 69.  
521 doi:10.1186/s40168-017-0283-5.
- 522 Roux, S., Enault, F., Hurwitz, B. L., and Sullivan, M. B. (2015). VirSorter: Mining viral signal from  
523 microbial genomic data. *PeerJ* 2015. doi:10.7717/peerj.985.
- 524 Roux, S., Páez-Espino, D., Chen, I. M. A., Palaniappan, K., Ratner, A., Chu, K., et al. (2021).  
525 IMG/VR v3: An integrated ecological and evolutionary framework for interrogating genomes of  
526 uncultivated viruses. *Nucleic Acids Research* 49, D764–D775. doi:10.1093/nar/gkaa946.
- 527 Segata, N., Waldron, L., Ballarini, A., Narasimhan, V., Jousson, O., and Huttenhower, C. (2012).  
528 Metagenomic microbial community profiling using unique clade-specific marker genes. *Nature*  
529 *Methods* 9, 811–814. doi:10.1038/nmeth.2066.
- 530 Soneson, C., Love, M. I., and Robinson, M. D. (2015). Differential analyses for RNA-seq: transcript-  
531 level estimates improve gene-level inferences. *F1000Research* 4, 1521.  
532 doi:10.12688/f1000research.7563.1.

- 533 Tithi, S. S., Aylward, F. O., Jensen, R. v., and Zhang, L. (2018). FastViromeExplorer: A pipeline for  
534 virus and phage identification and abundance profiling in metagenomics data. *PeerJ* 2018.  
535 doi:10.7717/peerj.4227.
- 536 Truong, D. T., Franzosa, E. A., Tickle, T. L., Scholz, M., Weingart, G., Pasolli, E., et al. (2015).  
537 MetaPhlan2 for enhanced metagenomic taxonomic profiling. *Nature Methods* 12, 902–903.  
538 doi:10.1038/nmeth.3589.
- 539 Wang, J., and Jia, H. (2016). Metagenome-wide association studies: Fine-mining the microbiome.  
540 *Nature Reviews Microbiology* 14, 508–522. doi:10.1038/nrmicro.2016.83.
- 541 Wang, J., Qin, J., Li, Y., Cai, Z., Li, S., Zhu, J., et al. (2012). A metagenome-wide association study  
542 of gut microbiota in type 2 diabetes. *Nature* 490, 55–60. doi:10.1038/nature11450.
- 543 Yu, T. C., Guo, F., Yu, Y., Sun, T., Ma, D., Han, J., et al. (2017). *Fusobacterium nucleatum*  
544 Promotes Chemoresistance to Colorectal Cancer by Modulating Autophagy. *Cell* 170, 548-  
545 563.e16. doi:10.1016/j.cell.2017.07.008.
- 546 Zuo, T., Sun, Y., Wan, Y., Yeoh, Y. K., Zhang, F., Cheung, C. P., et al. (2020). Human-Gut-DNA  
547 Virome Variations across Geography, Ethnicity, and Urbanization. *Cell Host and Microbe* 28,  
548 741-751.e4. doi:10.1016/j.chom.2020.08.005.

549

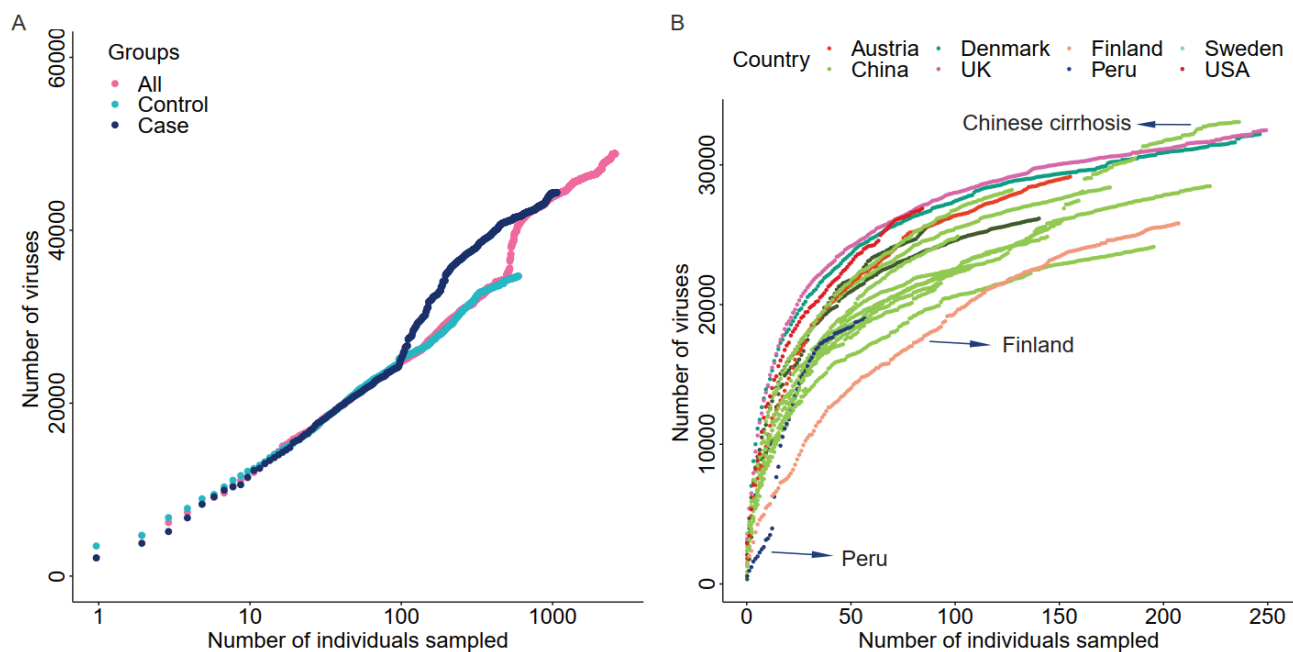
550



551

552 Figure 1: Location, taxonomic assignment, and abundance of the 2,690 samples. A: Geographic  
553 locations of the 18 projects, with classification by the number of enterotypes. B: Pie chart shows viral  
554 taxonomic assignment at the family level by protein alignment. C: Principal Coordinates Analysis  
555 (PCoA) plot based on the Bray–Curtis distance and the relative abundance of viruses. D: t-  
556 Distributed Stochastic Neighbor Embedding (tSNE) plot based on the relative abundance of viruses.

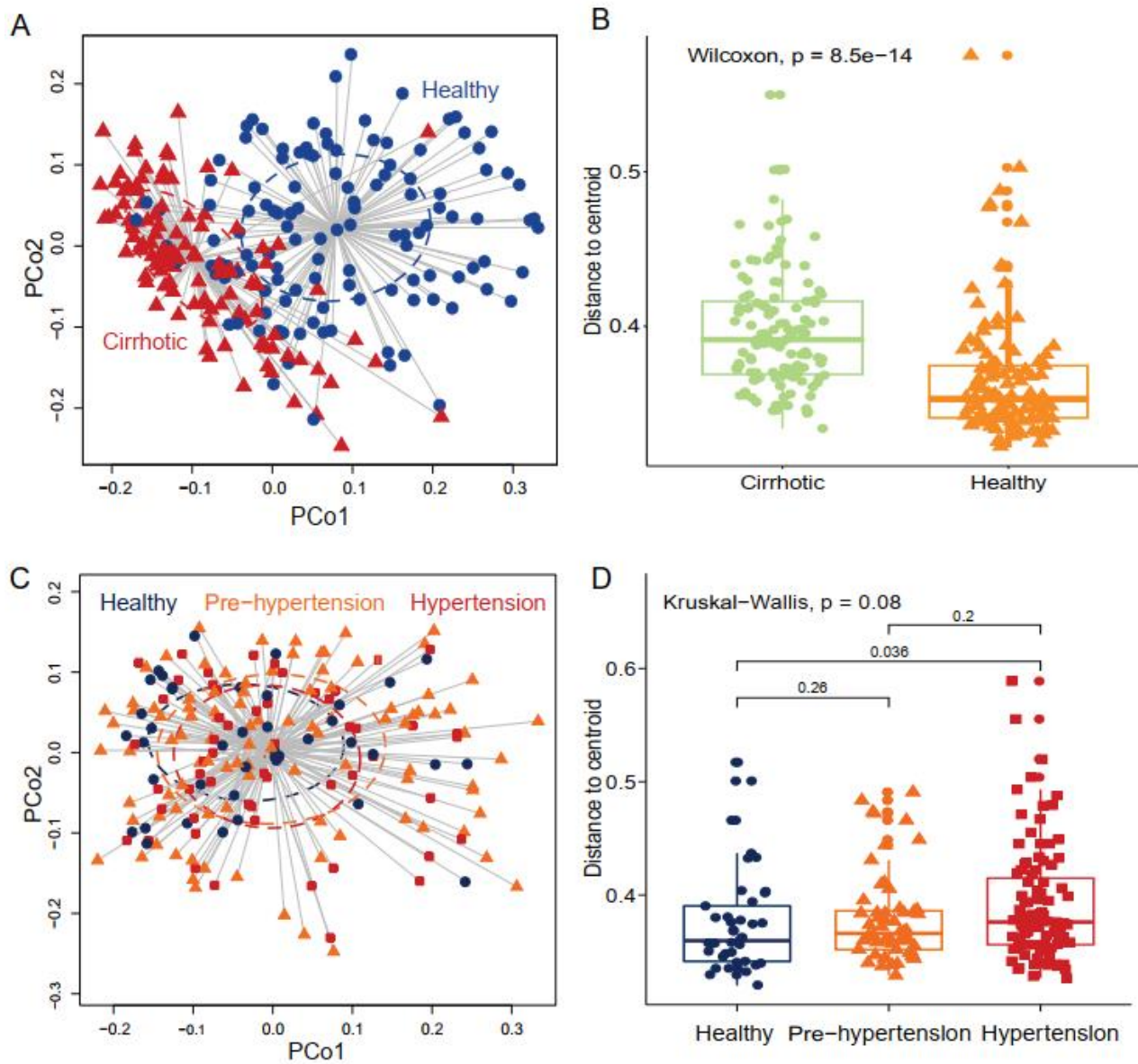
557



558

559 Figure 2: Cumulative curves of the number of virus species against the number of samples. A:  
560 Cumulative curves of cases, controls, and all samples. Only samples from studies with a case-control  
561 design were included. B: Cumulative curves of sample data divided into seven countries.

562

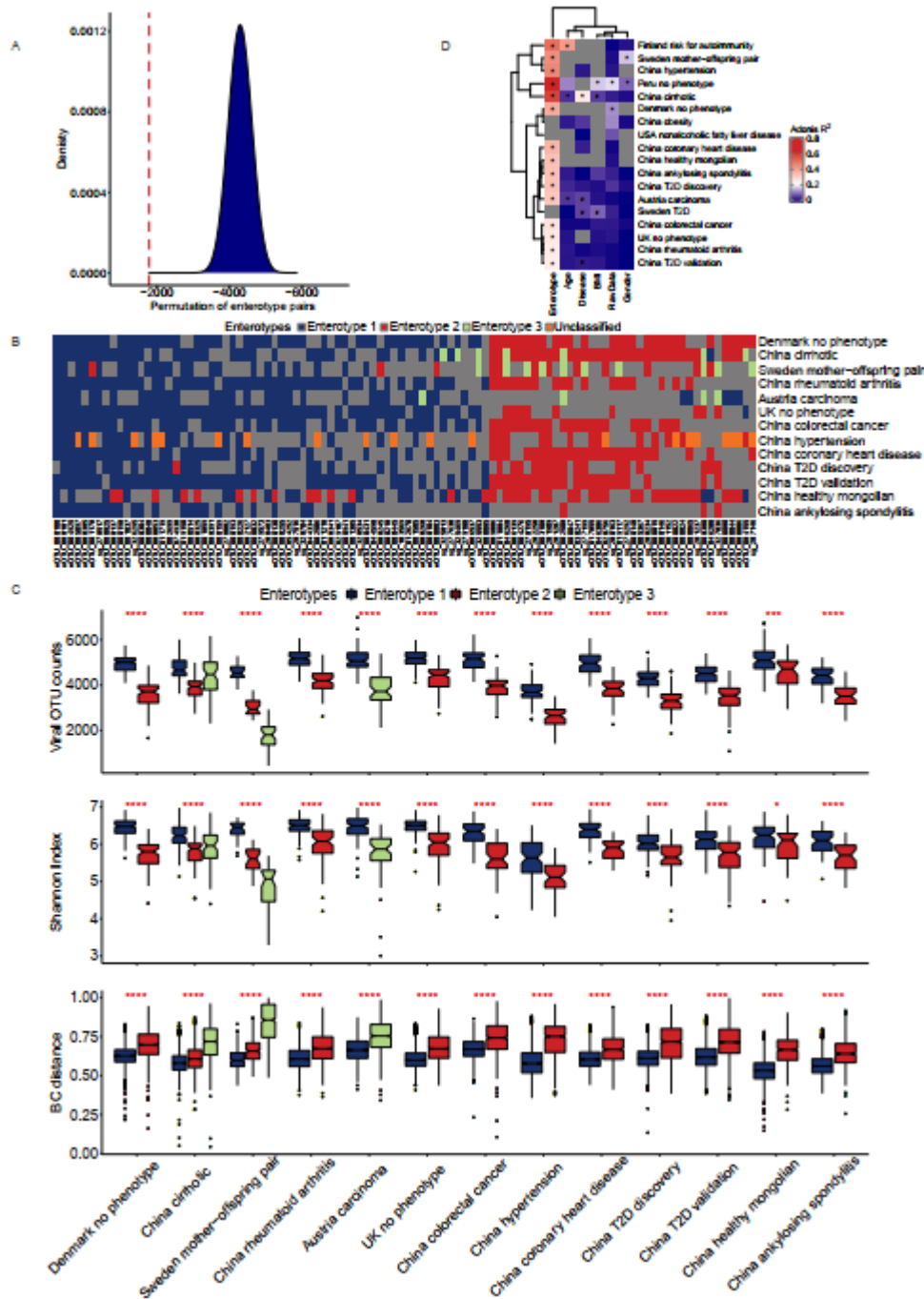


563

564 Figure 3: Gut virome characterized by beta diversity in the included projects. (A) Principal  
565 coordinates analysis plot of the cirrhosis project. Each ellipse represents a cohort, and the point  
566 connected by the straight gray lines represents the centroid. (B) Boxplot of the distance to the  
567 centroid. A significant difference in the distance to the centroid was found between the two groups.  
568 (C) Principal coordinates analysis plot of the hypertension project. (D) Boxplot of the hypertension  
569 project with comparison for each pair of the three groups.

570

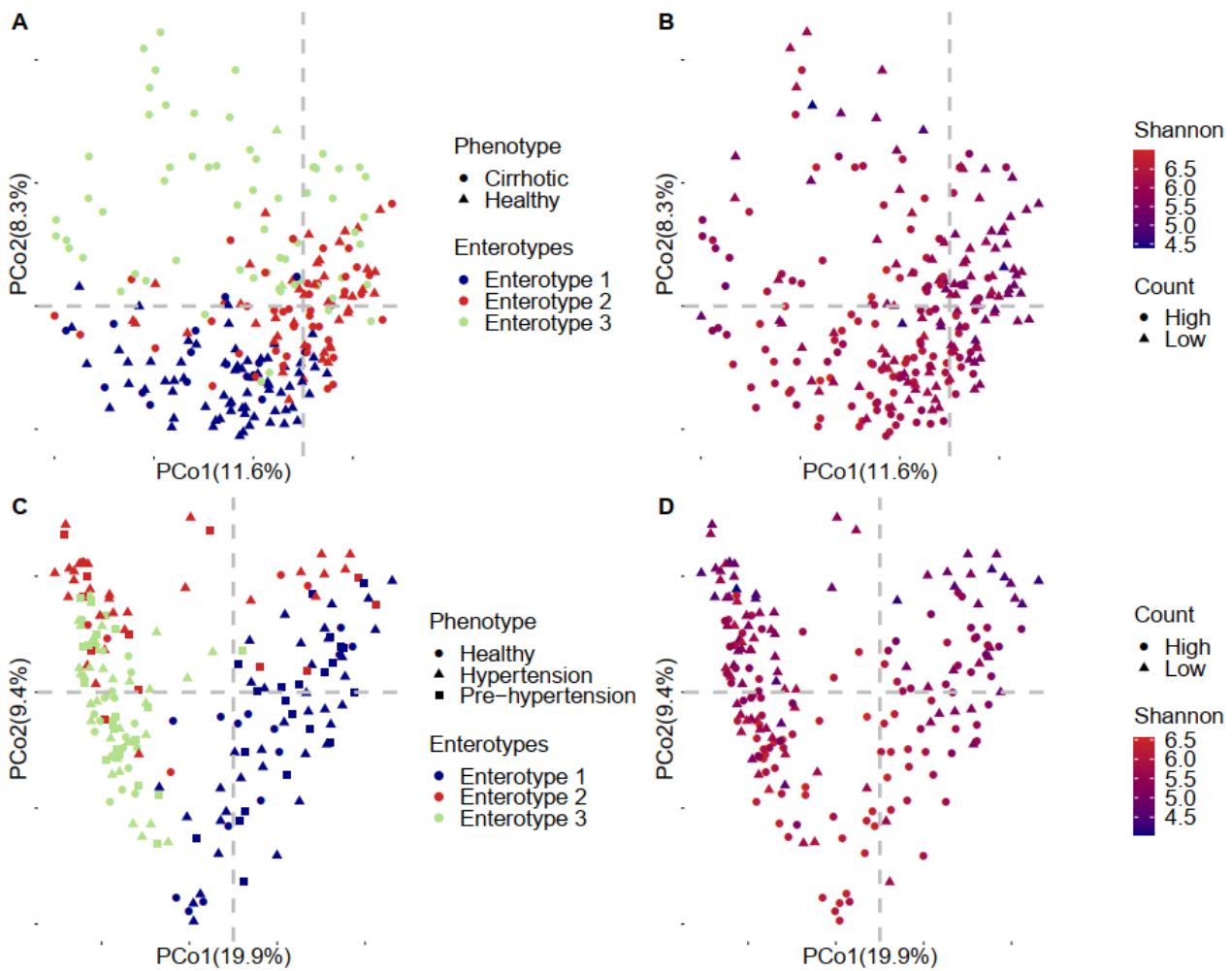




571

572 Figure 4: Characterization of viral enterotypes in all projects. A: We used the random pairing method  
 573 to confirm the accuracy of artificial enterotype classification. The density map shows the score  
 574 distribution of 5 million permutations, and the red line indicates the score of the manual category. B:  
 575 The categories of manual enterotypes in different projects show a high concordance of their specific  
 576 vOTUs and enrichment direction. C: Ecological diversity of different viral enterotype populations. D:  
 577 Effect of different covariates on the structural variance of the gut virome community.

578



579

580 Figure 5: Detailed PCoA map of liver cirrhosis and hypertension. A: Samples of liver cirrhosis were  
581 plotted in relation to their phenotype and enterotypes. B: Samples of liver cirrhosis were plotted in  
582 relation to their viral count and Shannon index. C: Samples of hypertension were plotted in relation  
583 to their phenotype and enterotype. D: Samples of hypertension were plotted in relation to their viral  
584 count and Shannon index.